

Variance-Gated Ensembles: An Epistemic-Aware Framework for Uncertainty Estimation

Anonymous authors
Paper under double-blind review

Abstract

Machine learning applications require fast and reliable per-sample uncertainty estimation. A common approach is to use predictive distributions from Bayesian or approximation methods and additively decompose uncertainty into aleatoric (data-related) and epistemic (model-related) components. However, additive decomposition has recently been questioned, with evidence that it breaks down when using finite-ensemble sampling and/or mismatched predictive distributions. This paper introduces variance-gated ensembles (VGE), a differentiable framework that injects epistemic sensitivity *via* a signal-to-noise gate computed from ensemble statistics. VGE provides: (i) a variance-gated margin uncertainty (VGMU) score that couples decision margins with ensemble predictive variance; and (ii) a variance-gated normalization (VGN) layer that generalizes the variance-gated uncertainty mechanism to training *via* per-class, learnable normalization of ensemble member probabilities. We derive closed-form vector-Jacobian products enabling end-to-end training through ensemble sample mean and variance. VGE is positioned as a compute-efficient alternative to pairwise-divergence methods. It delivers competitive uncertainty quality relative to state-of-the-art information-theoretic baselines, while reducing per-sample cost from quadratic to linear in ensemble size. As a result, VGE provides a practical and scalable approach to epistemic-aware uncertainty estimation in ensemble models.

1 Introduction

Machine learning models achieve high accuracy on average yet still fail on individual predictions in ways that are obvious on inspection. Flagging such failures requires per-sample uncertainty estimates that indicate when a prediction should not be trusted. A common approach to estimate uncertainty is Bayesian model averaging (BMA) in Bayesian neural networks (BNNs) or through approximations such as Monte Carlo dropout (MCD) (Gal & Ghahramani, 2016), deep ensembles (DE) (Lakshminarayanan et al., 2017), and last-layer ensembles (LLEs) (Harrison et al., 2024; Sensoy et al., 2018). These methods approximate the predictive distribution by averaging predictions from an ensemble of models, and the standard approach decomposes the resulting predictive entropy into data- and model-related components (Houlsby et al., 2011).

However, existing approaches face a trade-off between efficiency and reliability. Entropy-based decompositions are efficient but unreliable, and finite ensembles introduce sampling error into the decomposition. Methods such as MCD and DE sample from approximate posteriors (*i.e.*, the dropout distribution or an implicit distribution over independently-trained networks) rather than the true posterior (Wimmer et al., 2023; Schweighofer et al., 2023). Pairwise alternatives such as the expected pairwise Kullback–Leibler (EPKL) divergence (Schweighofer et al., 2023) recover epistemic sensitivity under finite ensembles but incur an $O(M^2C)$ cost that scales poorly with ensemble size and class count.

Here, we introduce variance-gated ensembles (VGE), a framework that recovers epistemic sensitivity at linear cost through an exponential gating mechanism. We propose a variance-gated margin uncertainty (VGMU) score using a margin-based signal-to-noise of the top-2 ranked classes. This concept is extended to an end-to-end variance-gated normalization (VGN) layer using a per-class-based signal-to-noise, where parameters are used to adaptively scale a gating mechanism during training. The VGN layer suppresses high-

FIX:
VoEx::1

FIX:
pDhd::8

39 variance, low-consensus predictions and re-shapes member distributions. VGE achieves $O(MC)$ complexity
 40 for uncertainty decomposition and $O(C)$ for VGMU evaluation at inference, enabling real-time deployment
 41 without sacrificing epistemic sensitivity.

42 In summary, this research provides the following contributions:

- 43 1. **Variance-Gated Ensembles.** A framework that introduces epistemic sensitivity *via* a signal-to-
 44 noise gate computed from ensemble statistics.
- 45 2. **Variance-Gated Margin Uncertainty.** An inference-time score combining decision margin with
 46 ensemble predictive variance to quantify class separability.
- 47 3. **Variance-Gated Normalization.** A differentiable, epistemic-aware normalization layer that mod-
 48 ulates member probabilities through a per-class learnable parameter, enabling efficient end-to-end
 49 optimization with $O(MC)$ complexity.

50 Apart from this section, the remainder of the paper is organized as follows: [Section 2](#) provides the
 51 information-theoretic background for uncertainty decomposition and positions this work within the existing
 52 literature. [Section 3](#) introduces the variance-gated ensemble framework, including the VGMU scoring metric
 53 and the formulation of VGN. This section also presents a geometric-based interpretation of ensemble predic-
 54 tions. [Section 4](#) introduces the analytical gradients required for end-to-end training using VGN. [Section 5](#)
 55 presents experimental set-up and results, comparing the proposed approach with common uncertainty es-
 56 timation methods in the machine learning literature. [Section 6](#) provides an axiomatic comparison to prior
 57 uncertainty decomposition frameworks, and clarifies the intended scope and limitations of VGN and VGMU.
 58 Finally, [Section 7](#) concludes the paper with a summary of findings and recommendations.

59 2 Background and Related Work

60 2.1 Predictive Uncertainty

61 Predictive uncertainty under a probabilistic model is typically formalized through the Bayesian predictive
 62 distribution, which averages model predictions over a posterior on parameters. Let \mathbf{w} denote the parameters
 63 of a probabilistic model and $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ the observed data. Bayesian inference maintains a posterior
 64 $p(\mathbf{w} | \mathcal{D})$ over parameters. Predictions for a new input \mathbf{x} are obtained *via* BMA, which integrates over the
 65 posterior to yield the predictive distribution over labels y

$$p(y | \mathbf{x}, \mathcal{D}) = \int_{\mathbf{w}} p(y | \mathbf{x}, \mathbf{w}) p(\mathbf{w} | \mathcal{D}) d\mathbf{w}. \quad (1)$$

66 This integral relies on the assumption that y is *conditionally independent of \mathcal{D} given \mathbf{w}*

$$p(y | \mathbf{x}, \mathbf{w}, \mathcal{D}) = p(y | \mathbf{x}, \mathbf{w}). \quad (2)$$

67 This means that once \mathbf{w} is known, the training data provides no additional information for predicting y . To
 68 quantify this uncertainty, we measure the Shannon entropy of the predictive distribution

$$H(y | \mathbf{x}, \mathcal{D}) = - \sum_y p(y | \mathbf{x}, \mathcal{D}) \log p(y | \mathbf{x}, \mathcal{D}) \quad (3)$$

69 which captures the predictive uncertainty for a new input \mathbf{x} .

70 2.2 Additive Decomposition

71 Total uncertainty (TU) as measured by entropy conflates two distinct sources, aleatoric uncertainty (AU) and
 72 epistemic uncertainty (EU). By definition, the mutual information (I) between the label (y) and parameters \mathbf{w}
 73 is the reduction in predictive entropy from knowing \mathbf{w} , where $I(y; \mathbf{w} | \mathbf{x}, \mathcal{D}) = H(y | \mathbf{x}, \mathcal{D}) - H(y | \mathbf{w}, \mathbf{x}, \mathcal{D})$.

FIX:
pDhd::9

74 Applying the conditional independence assumption (Equation 2), we can decompose $H(y \mid \mathbf{x}, \mathcal{D})$. Since
 75 $p(y \mid \mathbf{x}, \mathbf{w}, \mathcal{D}) = p(y \mid \mathbf{x}, \mathbf{w})$, the conditional entropy of y given \mathbf{w} can be estimated from model sampling
 76 (*i.e.*, an ensemble)

$$H(y \mid \mathbf{w}, \mathbf{x}, \mathcal{D}) = \mathbb{E}_{\mathbf{w} \sim p(\mathbf{w} \mid \mathcal{D})} [H(y \mid \mathbf{x}, \mathbf{w})]. \quad (4)$$

77 Replacing the second term with Equation 4 results in the additive decomposition.

$$\underbrace{H(y \mid \mathbf{x}, \mathcal{D})}_{\text{TU}} = \underbrace{\mathbb{E}_{\mathbf{w} \sim p(\mathbf{w} \mid \mathcal{D})} [H(y \mid \mathbf{x}, \mathbf{w})]}_{\text{AU}} + \underbrace{I(y; \mathbf{w} \mid \mathbf{x}, \mathcal{D})}_{\text{EU}}. \quad (5)$$

78 The two terms have qualitatively different implications for downstream tasks such as active learning, selec-
 79 tive prediction, out-of-distribution (OOD) detection, and human-in-the-loop systems. The aleatoric term
 80 $\mathbb{E}_{\mathbf{w} \sim p(\mathbf{w} \mid \mathcal{D})} [H(y \mid \mathbf{x}, \mathbf{w})]$ averages the entropy of the individual model predictions over the posterior. It
 81 captures noise intrinsic to the data-generating process, uncertainty that persists even if the posterior were
 82 known exactly and therefore cannot be reduced by collecting more data. The epistemic term $I(y; \mathbf{w} \mid \mathbf{x}, \mathcal{D})$
 83 measures how much the predictions of the model would change if we knew which parameters \mathbf{w} were cor-
 84 rect. When ensemble members disagree, this term is large; when they agree, it vanishes. Unlike aleatoric
 85 uncertainty, epistemic uncertainty is in principle reducible. Observing additional data updates the posterior
 86 $p(\mathbf{w} \mid \mathcal{D})$ and reduces the disagreement among model hypotheses. The additive decomposition (Equation 5)
 87 enables downstream decisions for active learning (Houlsby et al., 2011), selective prediction (Geifman & El-
 88 Yaniv, 2017), and OOD detection (Lakshminarayanan et al., 2017) that selectively target only the reducible
 89 component of predictive uncertainty (Gawlikowski et al., 2023).

90 2.3 Expected Pairwise Kullback–Leibler Divergence

91 In practice, the posterior $p(\mathbf{w} \mid \mathcal{D})$ is approximated by a finite ensemble $\{\mathbf{w}_m\}_{m=1}^M$, with the posterior
 92 predictive distribution approximated as the mixture

$$p(y \mid \mathbf{x}, \mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^M p(y \mid \mathbf{x}, \mathbf{w}_m). \quad (6)$$

93 When a finite ensemble produces individual member distributions that are non-Gaussian, multimodal, or
 94 heavily skewed, the averaged prediction may not represent the output of any individual member, and its
 95 entropy fails to capture the true spread of the ensemble (Wimmer et al., 2023). Given this finite ensemble ap-
 96 proximation (Equation 6), an alternative for the epistemic component is the EPKL divergence (Schweighofer
 97 et al., 2023)

$$\text{EPKL}(\mathbf{x}) = \frac{1}{M(M-1)} \sum_{i=1}^M \sum_{\substack{j=1 \\ j \neq i}}^M D_{\text{KL}}(p(y \mid \mathbf{x}, \mathbf{w}_i) \parallel p(y \mid \mathbf{x}, \mathbf{w}_j)). \quad (7)$$

98 EPKL is non-negative, equals zero if and only if all ensemble members agree, and captures pairwise dis-
 99 agreement directly in predictive distribution space without requiring estimation of the mixture entropy.
 100 However, a practical limitation of EPKL is quadratic scaling. Evaluating Equation 7 requires $O(M^2C)$ pair-
 101 wise divergence computations, compared with the $O(MC)$ for entropy-based estimators. This cost becomes
 102 prohibitive as ensemble size and number of classes increases, motivating more efficient alternatives that still
 103 capture ensemble disagreement.

104 2.4 Ensemble Uncertainty Estimation

105 **Ensemble approximations.** Uncertainty estimation (Gawlikowski et al., 2023; Hüllermeier & Waegeman,
 106 2021) in classification tasks is often framed through Bayesian predictive distributions; however, an exact
 107 Bayesian solution is intractable for modern neural networks. Practical approximations include variational
 108 Bayesian neural networks (Radford, 1995), MCD (Gal & Ghahramani, 2016; Gal et al., 2017), DE (Lak-
 109 shminarayanan et al., 2017), and LLE approaches such as evidential deep learning (Sensoy et al., 2018),
 110 variational inference (Harrison et al., 2024; Steger et al., 2024), multihead (Lee et al., 2015), and other en-
 111 semble strategies (Huang et al., 2017; Kushibar et al., 2022; Wen et al., 2020). Among these approaches, DE

FIX:
pDhd::10

112 have demonstrated strong performance in calibration and OOD detection, often outperforming approximate
 113 Bayesian methods under dataset shift (Ovadia et al., 2019). As a result, many recent uncertainty estimation
 114 methods operate on ensemble predictive distributions.

115 **Moment-based methods.** An alternative line of work derives uncertainty signals from ensemble statistics,
 116 such as the mean and variance of predicted class probabilities (Depeweg et al., 2018; Smith & Gal, 2018).
 117 Moment-based approaches are attractive in large-scale or real-time settings. VGE follows this approach by
 118 deriving measures from ensemble variance, enabling linear-time computation while retaining sensitivity to
 119 predictive disagreement.

120 **Calibration methods.** Calibration methods, such as temperature scaling, operate post hoc and adjust
 121 predictive confidence without modifying uncertainty structure during training (Guo et al., 2017; Kumar
 122 et al., 2022). In contrast, recent work explores integrating uncertainty-aware mechanisms into the training
 123 process itself, including diversity-promoting losses and uncertainty-regularized objectives (Fort et al., 2020;
 124 Ashukha et al., 2021). VGN contributes to this area by introducing a differentiable normalization layer that
 125 modulates ensemble predictions based on epistemic variance, allowing uncertainty re-shaping to be learned
 126 end-to-end.

127 **Margin-based criteria.** In many applications, risk-based decisions depend on ambiguity between top-
 128 ranked classes rather than uncertainty of the full-simplex distribution. Margin-based criteria, such as Best-
 129 versus-Second Best (BvSB) scores, are used in active learning and selective prediction (Joshi et al., 2009;
 130 Geifman & El-Yaniv, 2017). VGMU extends this method by incorporating ensemble variance into the
 131 decision margin, offering a decision-focused epistemic uncertainty measure that emphasizes disagreement
 132 among top-ranked classes while remaining computationally efficient.

133 The variance-gated ensemble framework, introduced in the next section, derives uncertainty directly from
 134 ensemble moments, achieving linear-time epistemic sensitivity without requiring pairwise comparisons.

135 3 Framework Definition and Setup

136 This section describes the variance-gated ensemble framework for epistemic-aware ensemble modeling, defin-
 137 ing normalization and analytical properties. We first formalize variance-gated normalization, where member
 138 probabilities are modulated by a signal-to-noise gate, and analyze its sensitivity to sample mean confidence,
 139 predictive spread, and per-class learnable parameter \mathbf{k} . We then provide a geometric-based interpretation,
 140 followed by a variance-gated uncertainty decomposition for total uncertainty into aleatoric and epistemic com-
 141 ponents, and introduce the variance-gated margin uncertainty score to capture class separability. Figure 1
 142 provides a high-level overview of the VGE framework, illustrating how ensemble predictions flow through
 143 variance-gated normalization to produce epistemic-aware distributions and uncertainty scores. See Table S1
 144 in the supporting information for a listing of symbols and abbreviations used throughout the variance-gated
 145 ensemble framework.

NEW:
cddg::C1

NEW:
pDhd::11

146 **Notation convention.** Throughout this paper, boldface lowercase letters denote C -dimensional vectors
 147 (*e.g.*, \mathbf{p}_m , $\bar{\mathbf{p}}$, \mathbf{s} , $\mathbf{\Gamma}$) and all operations on such vectors, including squaring, division, exponentiation, multipli-
 148 cation, and inequalities are applied element-wise unless explicitly stated otherwise.

NEW:
pDhd::1

149 **Problem setup.** Consider a C -class classification task in which an input \mathbf{x} is mapped to a label $y \in$
 150 $\{1, \dots, C\}$. An ensemble of M models with parameters $\{\mathbf{w}_m\}_{m=1}^M$ is available, each producing a predictive
 151 categorical distribution over classes. The goal is per-sample uncertainty estimation. Given a new input
 152 \mathbf{x} , quantify both the predictive confidence of the ensemble and the degree of inter-member disagreement,
 153 decomposing predictive uncertainty into aleatoric and epistemic components.

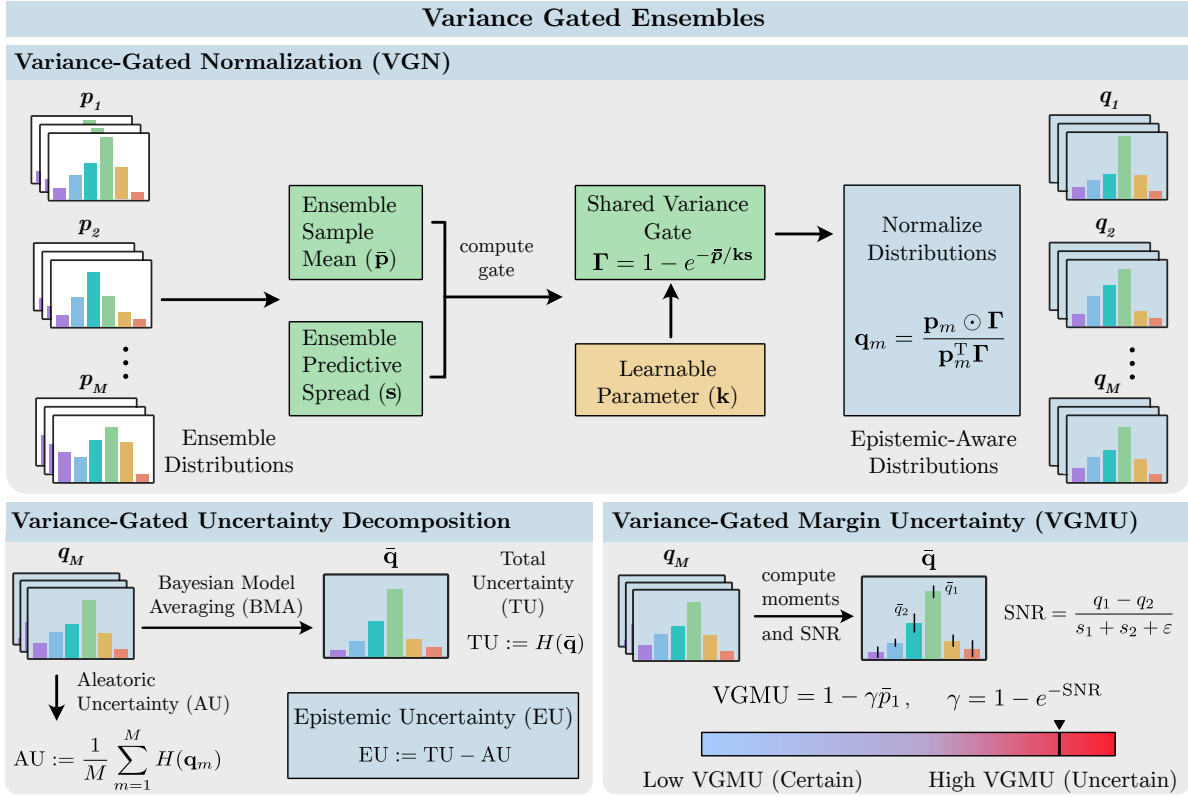


Figure 1: Overview of the variance-gated ensemble framework. Top: Variance-gated normalization computes a shared gate Γ from the ensemble mean $\bar{\mathbf{p}}$, predictive spread \mathbf{s} , and learnable parameter \mathbf{k} , producing epistemic-aware distributions \mathbf{q}_m . Bottom (left): Variance-gated additive decomposition of total uncertainty into aleatoric and epistemic components. Bottom (right): Variance-gated margin uncertainty, a decision-focused score computed from the top-2 class margins and their signal-to-noise ratio. VGMU does not require the variance-gated prediction q ; it can be evaluated on any ensemble output p given access to per-member means and standard deviations.

154 3.1 Ensemble Statistics and Variance Gate Definition

155 We first define the variance gate and analyze its local sensitivities. Let each ensemble member produce a
 156 predictive categorical distribution

$$\mathbf{p}_m = p(y | \mathbf{x}, \mathbf{w}_m) \in \Delta^{C-1}, \quad m \in \{1, \dots, M\}, \quad (8)$$

157 where $\mathbf{p}_m = [p_m(1), \dots, p_m(C)]^\top$ lies in a compact convex region known as the $(C-1)$ -simplex. The per-class
 158 ensemble sample mean and standard deviation are defined as

$$\bar{\mathbf{p}} = \frac{1}{M} \sum_{m=1}^M \mathbf{p}_m, \quad \mathbf{s} = \sqrt{\frac{1}{M-1} \sum_{m=1}^M (\mathbf{p}_m - \bar{\mathbf{p}})^2 + \varepsilon}, \quad M > 1. \quad (9)$$

159 Let $\bar{\mathbf{p}} \geq 0$, $\mathbf{s} \geq 0$, $\mathbf{k} > 0$, and $\varepsilon > 0$ (e.g., 1.0×10^{-8}).

FIX:
VoEx::2

160 **Design rationale.** The variance gate requires a mapping from the signal-to-noise ratio $\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}$ to a gating
 161 weight in $[0, 1)$. We adopt the exponential form $1 - e^{-(\cdot)}$ because it satisfies four desirable properties:
 162 (i) smoothness and differentiability everywhere, enabling end-to-end gradient-based training; (ii) monotonic
 163 increase with mean confidence and monotonic decrease with predictive spread; (iii) bounded output in $[0, 1)$,
 164 ensuring numerical stability during normalization; and (iv) saturation that prevents excessive attenuation of

well-supported predictions. The exponential form is one choice among several possible gating families (e.g., sigmoid, rational, piecewise-linear). Alternative parameterizations that trade-off sensitivity and saturation are discussed as an open direction in Section 6.

The variance gate is defined as

$$\mathbf{\Gamma} = 1 - e^{-\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}}, \quad \mathbf{\Gamma} \in \mathbb{R}^C, \quad 0 \leq \mathbf{\Gamma} < 1 \quad (\text{element-wise}) \quad (10)$$

and the normalized variance-gated member distribution is defined as

$$\mathbf{q}_m = \frac{\mathbf{p}_m \odot \mathbf{\Gamma}}{Z_m}, \quad Z_m = \mathbf{p}_m^\top \mathbf{\Gamma} \quad \mathbf{q}_m \in \mathbb{R}^C, \quad 0 \leq \mathbf{q}_m \leq 1 \quad (\text{element-wise}) \quad (11)$$

where the scalar normalization Z_m ensures $\mathbf{1}^\top \mathbf{q}_m = 1$ and $\mathbf{1} \in \mathbb{R}^C$ denotes the all-ones vector. The per-class $\mathbf{k} > 0$ controls the sensitivity of the gate.

The variance gate modulates ensemble predictions based on the scaled signal-to-noise ratio $\text{SNR} = \bar{\mathbf{p}}/\mathbf{k}\mathbf{s}$ and acts as a smooth reliability correction. This means that classes with high mean confidence and low predictive spread receive larger gate values, while classes that are uncertain or highly variable are suppressed before normalization. In other words, $\mathbf{\Gamma}$ returns a per-class reliability weight close to 1 when the ensemble agrees confidently and close to 0 when the ensemble is unreliable; multiplying each member distribution \mathbf{p}_m by $\mathbf{\Gamma}$ and re-normalizing therefore preserves trusted classes while attenuating untrusted predictions before downstream uncertainty computations. We now examine the sensitivity of the variance gate and gated member distributions to sample mean confidence $\bar{\mathbf{p}}$, predictive spread \mathbf{s} , scaling factor \mathbf{k} , and its effects on gated distributions \mathbf{q}_m .

Proposition 3.1 (Sensitivity to sample mean confidence $\bar{\mathbf{p}}$). *For the exponential gate $\mathbf{\Gamma} = 1 - e^{-\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}}$, the per-class derivative with respect to mean confidence $\partial\mathbf{\Gamma}/\partial\bar{\mathbf{p}} > 0$ and is modulated by predictive spread \mathbf{s} through both an explicit inverse prefactor $1/(k_c s_c)$ and the saturation term $(1 - \Gamma_c)$, which itself depends on s_c through the exponent.*

Proof. Consider a single class $c \in \{1, \dots, C\}$, for which the gate is

$$\Gamma_c = 1 - e^{-\bar{p}_c/k_c s_c}. \quad (12)$$

Differentiating with respect to the mean confidence \bar{p}_c yields

$$\frac{\partial\Gamma_c}{\partial\bar{p}_c} = \frac{1}{k_c s_c} e^{-\bar{p}_c/k_c s_c} = \frac{1 - \Gamma_c}{k_c s_c} > 0. \quad (13)$$

The derivative is strictly positive, showing that increasing the mean confidence for class c always increases its gate value. The dependence on s_c enters through two mechanisms: the explicit prefactor $1/(k_c s_c)$ provides linear suppression, while the saturation term $(1 - \Gamma_c) = e^{-\bar{p}_c/k_c s_c}$ provides nonlinear modulation. As ensemble disagreement s_c increases, both factors contribute to reducing the derivative, meaning that the gate becomes less responsive to changes in mean confidence. Thus, for classes with high predictive variance, increases in \bar{p}_c are suppressed. \square

Remark 3.1.1. *The factor $(1 - \Gamma_c)$ acts as a saturation term. As $\bar{p}_c/(k_c s_c) \rightarrow \infty$, we have $\Gamma_c \rightarrow 1$ and therefore $\partial\Gamma_c/\partial\bar{p}_c \rightarrow 0$. Consequently, once a class is deemed sufficiently reliable, further increases in mean confidence produce diminishing effects. This ensures that the variance-gate is most sensitive in intermediate regions and becomes increasingly insensitive as confidence saturates.*

Proposition 3.2 (Sensitivity to predictive spread \mathbf{s}). *For the exponential gate $\mathbf{\Gamma} = 1 - e^{-\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}}$, the per-class derivative with respect to predictive spread is $\partial\mathbf{\Gamma}/\partial\mathbf{s} < 0$ and is modulated by mean confidence $\bar{\mathbf{p}}$ through both an explicit linear prefactor \bar{p}_c and the saturation term $(1 - \Gamma_c)$, which itself depends on \bar{p}_c through the exponent.*

Proof. Consider a single class $c \in \{1, \dots, C\}$, for which

$$\Gamma_c = 1 - e^{-\bar{p}_c/k_c s_c}. \quad (14)$$

NEW:
pDHD::2
NEW:
VoEx::2
FIX:
pDHD::12

FIX:
pDHD::12

NEW:
VoEx::1

FIX:
pDHD::3

FIX:
pDHD::3

202 Differentiating with respect to the predictive spread s_c yields

$$\frac{\partial \Gamma_c}{\partial s_c} = -\frac{\bar{p}_c}{k_c s_c^2} e^{-\bar{p}_c/k_c s_c} = -\frac{(1 - \Gamma_c) \bar{p}_c}{k_c s_c^2} < 0. \quad (15)$$

203 The derivative is strictly negative, indicating that increasing predictive spread decreases the gate value.
 204 The factor $1/k_c s_c^2$ shows that the magnitude of this effect decays rapidly as s_c grows. Thus, classes with
 205 high ensemble disagreement experience strong suppression, while further increases in large spreads have
 206 diminishing influence. \square

207 **Remark 3.2.1.** *Increasing the predictive spread s_c while holding \bar{p}_c and k_c fixed decreases the gate through*
 208 *two mechanisms: i) a linear dependence on \bar{p}_c and a quadratic decay in s_c ; and ii) the saturation factor*
 209 *$(1 - \Gamma_c)$ that ensures once the gate is already small, additional increases in spread have limited effect. This*
 210 *behavior enforces strong suppression for uncertain classes, while the quadratic decay in s_c ensures diminishing*
 211 *sensitivity as predictive spread grows.*

212 **Proposition 3.3** (Sensitivity to scalar \mathbf{k}). *For the exponential gate $\Gamma = 1 - e^{-\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}}$, where k_c may be*
 213 *user-defined or learned, the derivative with respect to \mathbf{k} is $\partial\Gamma/\partial\mathbf{k} < 0$ and is modulated by mean confidence*
 214 *$\bar{\mathbf{p}}$ through both an explicit linear prefactor \bar{p}_c and the saturation term $(1 - \Gamma_c)$, which itself depends on \bar{p}_c*
 215 *through the exponent.*

216 *Proof.* Consider a single class $c \in \{1, \dots, C\}$, for which

$$\Gamma_c = 1 - e^{-\bar{p}_c/k_c s_c}. \quad (16)$$

217 Differentiating with respect to the scaling parameter k_c yields

$$\frac{\partial \Gamma_c}{\partial k_c} = -\frac{\bar{p}_c}{k_c^2 s_c} e^{-\bar{p}_c/k_c s_c} = -\frac{(1 - \Gamma_c) \bar{p}_c}{k_c^2 s_c} < 0. \quad (17)$$

218 Thus, increasing k_c decreases the gate value. The inverse quadratic dependence on k_c shows that the influence
 219 of the scaling parameter rapidly diminishes as k_c grows, resulting in controlled and saturating sensitivity. \square

220 **Remark 3.3.1.** *Increasing k_c while holding \bar{p}_c and s_c fixed decreases the gate by increasing the effective*
 221 *predictive spread. The quadratic decay in k_c ensures that sensitivity to further increases rapidly diminishes,*
 222 *while the saturation factor $(1 - \Gamma_c)$ prevents excessive attenuation once the gate is already small. Under*
 223 *mild distributional assumptions on ensemble dispersion, the product $k_c s_c$ may be interpreted as a classwise*
 224 *risk-tolerance threshold reflecting typical deviations from the ensemble mean.*

225 3.1.1 Geometric Interpretation

226 We now reinterpret variance-gating geometrically in the probability simplex. Each pair $(\bar{\mathbf{p}}, \{\mathbf{p}_m\})$ defines
 227 a configuration in the simplex Δ^{C-1} . The geometry of this configuration is defined by: (i) Confidence (or
 228 ambiguity): How close (or far) the ensemble sample mean $\bar{\mathbf{p}}$ is positioned near a simplex vertex; and (ii)
 229 Certainty (or uncertainty): How close (or far) the set of members cluster around the ensemble sample mean
 230 $\bar{\mathbf{p}}$ (*i.e.*, ensemble member agreement/disagreement). These geometric effects are modulated by the local
 231 sensitivities of the variance gate, which increases with mean confidence and is progressively attenuated by
 232 predictive spread and the classwise risk-tolerance scale $k_c s_c$. The qualitative combinations of these two
 233 dimensions correspond to four simplex regions that define the space of all ensemble behaviors.

NEW:
cddg::M6

234 **Confident–Certain.** This region is characterized by a near-deterministic ensemble mean and low variance.
 235 The ensemble members form a cluster around a single vertex of the simplex, indicating high confidence and
 236 high agreement:

$$\bar{p}_c \uparrow, s_c \downarrow \implies \Gamma_c = 1 - e^{-\bar{p}_c/k_c s_c} \approx 1. \quad (18)$$

237 **Ambiguous–Certain.** This region occurs when the ensemble mean is near-uniform but the variance is
 238 low. Members concentrate near the simplex barycenter, showing ambiguity but mutual certainty:

$$\bar{p}_c \downarrow, s_c \downarrow \implies \Gamma_c = 1 - e^{-\bar{p}_c/k_c s_c}, \quad \text{depends on } \bar{p}_c/k_c s_c. \quad (19)$$

239 **Confident–Uncertain.** This region is defined by a near-deterministic mean but high variance. Members
 240 radiate outward from a vertex, reflecting confidence but disagreement among ensemble members:

$$\bar{p}_c \uparrow, s_c \uparrow \implies \Gamma_c = 1 - e^{-\bar{p}_c/k_c s_c} \ll 1, \quad \text{with sensitivity suppressed by large } s_c. \quad (20)$$

241 **Ambiguous–Uncertain.** This region is defined by both a near-uniform mean and high variance. Members
 242 are diffused around the barycenter, indicating ambiguity and high disagreement:

$$\bar{p}_c \downarrow, s_c \uparrow \implies \Gamma_c = 1 - e^{-\bar{p}_c/k_c s_c} \approx 0. \quad (21)$$

243 The variance gate functions as a continuous geometric adjustment within the probability simplex. It main-
 244 tains the geometry in confident–certain regions, selectively attenuates confident–uncertain areas, and con-
 245 tracts diffuse ensembles’ predictions in ambiguous–uncertain cases. In contrast, ambiguous–certain regions
 246 remain unaffected, as their geometry already reflects agreement in uncertainty. These four simplex spaces
 247 define the set of ensemble behaviors through which variance-gating is applied.

248 Collectively, the sensitivity properties of the variance gate define a controlled signal-to-noise mechanism.
 249 The gate increases with mean confidence \bar{p}_c , but this effect is progressively attenuated by predictive spread
 250 s_c and the class-wise parameter k_c , with sensitivities that decay quadratically in both quantities. As a
 251 result, confident and consistent ensemble predictions are preserved, while high-confidence but high-variance
 252 predictions are suppressed in a stable and saturating manner (for an alternative risk-based interpretation,
 253 see SI S5).

254 3.2 Variance-Gated Margin Uncertainty

255 Entropy is known to overestimate uncertainty when probability values are spread across many classes and
 256 can underestimate it when a model is highly confident about a few options. Therefore, entropy alone is not
 257 sufficient to provide adequate information for decision-making. We posit the following question: *Can we*
 258 *identify a measure that is sensitive to class separation, while incorporating uncertainty awareness?* We want
 259 to identify a scoring metric that (i) maintains epistemic awareness and (ii) ideally, provides a user-defined
 260 threshold that reflects the degree of acceptable risk.

261 We propose using the margins of the top-2 mean predictions and their corresponding standard deviations.
 262 This approach extends the BvSB score used by Joshi et al. (2009), incorporating a sensitivity scalar k
 263 (applied to all classes). Let \bar{p}_1 and \bar{p}_2 denote the top-1 and top-2 ranked ensemble mean probabilities of $\bar{\mathbf{p}}$,
 264 with s_1 and s_2 , their corresponding standard deviations from \mathbf{s} . Let i denote the class index of the top-1
 265 ranked prediction. We define a prediction rule \hat{y} and derive a margin-based SNR as

$$\hat{y} = \begin{cases} i & \text{if } \bar{p}_1 - k s_1 > \bar{p}_2 + k s_2; \\ \text{abstain} & \text{otherwise} \end{cases}; \quad \text{SNR} = \frac{\bar{p}_1 - \bar{p}_2}{s_1 + s_2 + \varepsilon} > k \quad (22)$$

266 where $\bar{p}_1 - \bar{p}_2$ is the probability margin between the two most likely classes, $s_1 + s_2 + \varepsilon$ is the combined
 267 standard deviations, and “abstain” denotes that the model declines to predict due to insufficient margin. In
 268 principle, the SNR can be interpreted as a binary decision boundary between classes \bar{p}_1 and \bar{p}_2 , restricted
 269 by k . Under mild distributional assumptions on ensemble dispersion, threshold k can be set to reflect the
 270 fraction of samples that require abstentions or human intervention (*i.e.*, by sorting the margin-based SNR
 271 in increasing values). For example, when $k = 1$, only samples with $\text{SNR} > 1$ will be considered; all others
 272 are abstained. However, this criterion fails to capture cases where a model outputs ambiguous and uncertain
 273 predictions. Such outputs artificially inflate the SNR values, leading to misleading classifications. To address
 274 this limitation, we introduce VGMU, using a gating function that rescales the top-ranked model prediction
 275 by incorporating both confidence and variance (*i.e.*, epistemic) information

$$\text{VGMU} = 1 - \gamma \bar{p}_1, \quad \gamma = 1 - e^{-\text{SNR}}. \quad (23)$$

276 The VGMU functions as an uncertainty metric, where small values correspond to confident, well-separated
 277 predictions (low uncertainty), while large values capture ambiguous or uncertain cases (low separation or
 278 high variance). The “variance-gated” designation refers to the role of ensemble standard deviations s_1, s_2
 279 within the SNR-based gate γ . VGMU operates on the ensemble distributions \mathbf{p} or the VGN-transformed
 280 distributions \mathbf{q} , making it a lightweight, post hoc metric applicable at inference time without retraining.

FIX:
pDHd::4

NEW:
pDHd::5
NEW:
cddg::C2

281 **Justification for the top-2 restriction.** The restriction to the top-2 classes is motivated by three
 282 considerations. First, from a decision-theoretic perspective, the practical question in selective prediction
 283 and human-in-the-loop systems is whether the model can distinguish its best prediction from the runner-up;
 284 disagreement about distant classes is irrelevant to the decision at hand. This aligns with the Best-versus-
 285 Second-Best framework of Joshi et al. (2009). Second, restricting to the top-2 classes yields $O(C)$ complexity
 286 compared to $O(M^2C)$ for pairwise divergence measures, enabling the computational speedups demonstrated
 287 in our experiments. Third, the empirical validation in Table 1 shows that on CIFAR-10, VGMU achieves
 288 Spearman $\rho > 0.985$ with full-simplex measures, confirming that the top-2 restriction is consistent in ranking
 289 uncertainty for moderate class counts. On CIFAR-100, the lower correlation with pairwise measures reflects
 290 intentional insensitivity to tail-class disagreement. A property that we view as desirable for decision-focused
 291 uncertainty estimation, where the relevant question is separability of the leading candidates rather than
 292 full-simplex characterization.

293 3.3 Variance-Gated Uncertainty Decomposition

294 Using the variance-gated distributions, we define the ensemble mixture $\bar{\mathbf{q}} = \frac{1}{M} \sum_{m=1}^M \mathbf{q}_m$, where $\bar{\mathbf{q}}$ denotes
 295 the mean of the variance-gated member distributions. The total uncertainty is then measured by

$$\text{TU} := H(\bar{\mathbf{q}}) = -\bar{\mathbf{q}}^\top \log \bar{\mathbf{q}}. \quad (24)$$

296 Following standard ensemble decomposition (Houlsby et al., 2011), we define the gated aleatoric and epistemic
 297 components as

$$\text{AU} := \frac{1}{M} \sum_{m=1}^M H(\mathbf{q}_m) = -\frac{1}{M} \sum_{m=1}^M \mathbf{q}_m^\top \log \mathbf{q}_m, \quad \text{EU} := \text{TU} - \text{AU} = \frac{1}{M} \sum_{m=1}^M D_{\text{KL}}(\mathbf{q}_m \| \bar{\mathbf{q}}). \quad (25)$$

298 To compare our variance-gated decomposition, we computed the corresponding standard decompo-
 299 sition without variance-gating (*i.e.* distributions \mathbf{p}_m) and inter-member disagreements as unbounded
 300 EPKL (Schweighofer et al., 2023), and bounded Expected Pairwise Jensen-Shannon (EPJS) divergences
 301 using each ensemble member pair (i, j) , the midpoint $\mathbf{m}_{ij} = \frac{1}{2}(\mathbf{p}_i + \mathbf{p}_j)$, where pairwise measures are
 302 calculated as

$$\text{EPKL} = \frac{1}{M(M-1)} \sum_{i=1}^M \sum_{\substack{j=1 \\ j \neq i}}^M D_{\text{KL}}(\mathbf{p}_i \| \mathbf{p}_j), \quad D_{\text{KL}}(\mathbf{p}_i \| \mathbf{p}_j) = \mathbf{p}_i^\top (\log \mathbf{p}_i - \log \mathbf{p}_j), \quad (26)$$

$$\text{EPJS} = \frac{1}{M(M-1)} \sum_{i=1}^M \sum_{\substack{j=1 \\ j \neq i}}^M D_{\text{JS}}(\mathbf{p}_i \| \mathbf{p}_j), \quad D_{\text{JS}}(\mathbf{p}_i \| \mathbf{p}_j) = \frac{1}{2} D_{\text{KL}}(\mathbf{p}_i \| \mathbf{m}_{ij}) + \frac{1}{2} D_{\text{KL}}(\mathbf{p}_j \| \mathbf{m}_{ij}). \quad (27)$$

303 4 Analytical Gradients for Variance-Gated Normalization

304 In this section, we introduce reverse-mode differentiation in the ensemble setting, using VGN to capture epis-
 305 temic signals during training. We discuss vector-Jacobian products through the variance-gated normalized
 306 layer, gradients of the gate Γ with respect to $\bar{\mathbf{p}}, \mathbf{s}$ and a learnable per-class parameter \mathbf{k} . By optimizing a
 307 negative log-likelihood objective (*e.g.*, cross-entropy), the model implicitly minimizes predictive uncertainty,
 308 with a particular emphasis on reducing the epistemic component (Equation 5). The propagation of these
 309 gradients back to ensemble members provides a practical approach for variance-aware training within modern
 310 automatic differentiation frameworks.

311 Recall that each ensemble member produces $\mathbf{p}_m \in \Delta^{C-1}$ and all members share the same gating function
 312 $\Gamma = 1 - e^{-\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}}$, with ensemble statistics $\bar{\mathbf{p}}, \mathbf{s}$ and \mathbf{k} defined in Section 3. The gated member distribution
 313 and its normalization constant are

$$\mathbf{q}_m = \frac{\mathbf{p}_m \odot \Gamma}{Z_m}, \quad Z_m = \mathbf{p}_m^\top \Gamma. \quad (28)$$

314 When the training objective is applied to the predicted ensemble mixture distribution (averaged across
315 models), the loss depends only on $\bar{\mathbf{q}}$, where $\mathcal{L} = \mathcal{L}(\bar{\mathbf{q}})$ such that

$$\bar{\mathbf{q}} = \frac{1}{M} \sum_{m=1}^M \mathbf{q}_m \implies \frac{\partial \bar{\mathbf{q}}}{\partial \mathbf{q}_m} = \frac{1}{M} \mathbf{I}. \quad (29)$$

316 Then by the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{q}_m} = \frac{\partial \bar{\mathbf{q}}}{\partial \mathbf{q}_m} \frac{\partial \mathcal{L}}{\partial \bar{\mathbf{q}}} \implies \frac{1}{M} \mathbf{u}, \quad \mathbf{u} = \frac{\partial \mathcal{L}}{\partial \bar{\mathbf{q}}} \quad (\text{upstream gradient}). \quad (30)$$

317 This quantity is obtained by differentiating the loss objective with respect to the variance-gated mixture
318 distribution $\bar{\mathbf{q}}$ and represents the gradient signal backpropagated through the variance-gated normalization
319 layer.

320 In the following analysis we decompose the total gradient of the mixture objective into its independent paths
321 and describe gradients that flow through the ensemble statistics (mean and variance) that parameterize the
322 gate. These local gradients are combined to provide a general reverse-mode differentiation rule for shared
323 ensemble gates. This defines a complete analytical foundation for backpropagating epistemic-aware gradients,
324 in which the supervised cross-entropy loss depends only on $\bar{\mathbf{q}}$. This can be interpreted as optimizing a model
325 by minimizing predictive uncertainty, subject to epistemic constraints (complete analytical derivations and
326 gradient expressions available in SI S2).

327 4.1 Full Gradient Decomposition

328 **Proposition 4.1.** *When the objective depends on the ensemble mixture $\bar{\mathbf{q}} = \frac{1}{M} \sum_{m=1}^M \mathbf{q}_m$, the total gradient
329 received by each ensemble member is*

$$\frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} = \frac{1}{M} \left(\left. \frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \right|_{\Gamma} + \left. \frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \right|_{\bar{\mathbf{p}}} + \left. \frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \right|_{\mathbf{s}} \right). \quad (31)$$

330 *Proof.* Each path represents an independent dependency of \mathcal{L} on \mathbf{p}_m :

- 331 • **Direct normalization** (with Γ and \mathbf{k} fixed): how \mathbf{q}_m changes when \mathbf{p}_m changes. This captures
332 the effect of normalization on the simplex space;
- 333 • **Indirect gating via the mean** (with \mathbf{s} and \mathbf{k} fixed): how \mathbf{q}_m changes when the gate Γ changes
334 through the mean $\bar{\mathbf{p}}$; and
- 335 • **Indirect gating via the variance** (with $\bar{\mathbf{p}}$ and \mathbf{k} fixed): how \mathbf{q}_m changes when the gate Γ changes
336 through the predictive spread \mathbf{s} .

337 By the multivariate chain rule, these contributions combine additively to provide the total gradient. \square

338 **Remark 4.1.1.** *The term “fixed” in the parameters definition denotes that the corresponding variable is
339 held constant during partial differentiation.*

340 We instantiate these formulations for our variance-gate introduced in Section 3, providing analytic gradients
341 with respect to gating statistics, $(\bar{\mathbf{p}}, \mathbf{s})$ and per-class learnable parameter \mathbf{k} . For the exponential variance-
342 gate $\Gamma = 1 - e^{-\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}}$, gradients with respect to the ensemble mean $\bar{\mathbf{p}}$, predictive spread \mathbf{s} , and sensitivity
343 parameter \mathbf{k} follow directly by the chain rule. Collectively, these derivations establish a complete analytic
344 framework for backpropagating epistemic-aware gradients through ensemble networks, enabling end-to-end
345 optimization under epistemic-sensitive normalization (Figure 2). The variance-gated normalization layer
346 propagates gradients across ensemble members \mathbf{p}_m through a shared gating mechanism parameterized by
347 ensemble statistics $(\bar{\mathbf{p}}, \mathbf{s})$ and a learnable sensitivity parameter \mathbf{k} . Forward activations from each member
348 contribute to the ensemble mean and variance, which adjust the shared gate Γ during backpropagation. By

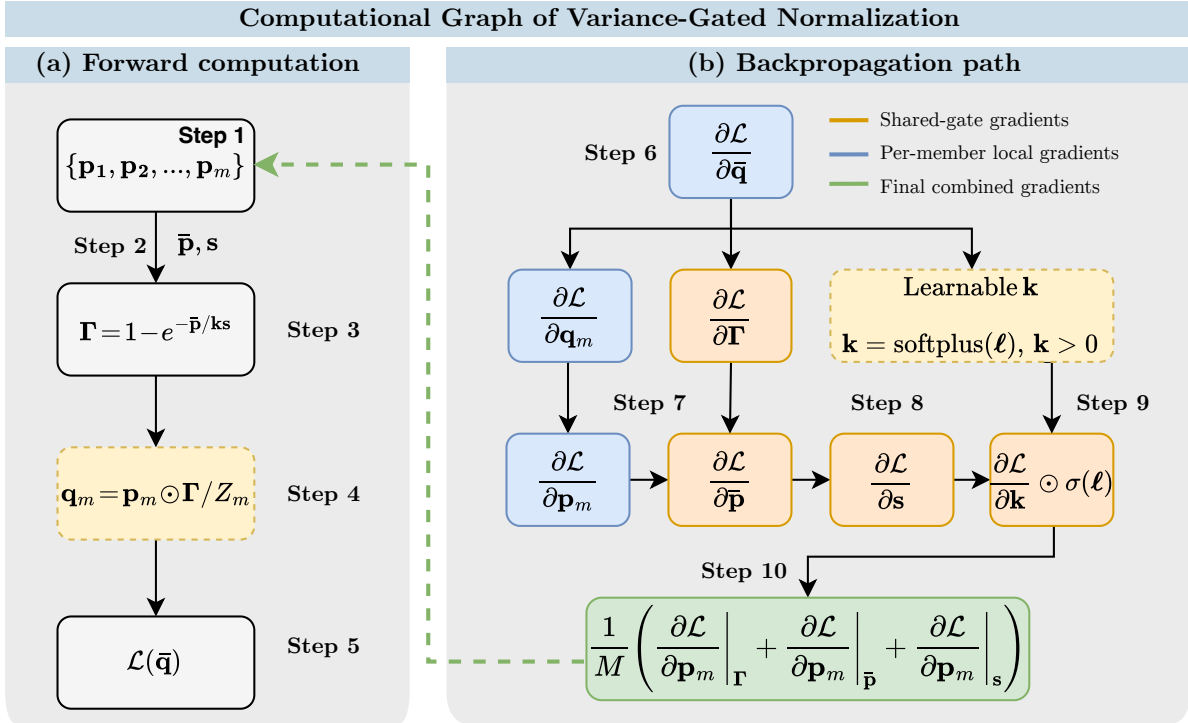


Figure 2: Computational graph of variance-gated normalization. Panel (a): Forward computation in which ensemble predictions are modulated by a shared variance gate and combined into a mixture distribution. Panel (b): Backpropagation path showing how gradients propagate through the normalization layer and shared gate *via* ensemble mean and predictive spread. See below for further step-by-step discussion.

349 including a learnable parameter $\mathbf{k} = \text{softplus}(\ell)$, models can adaptively modulate the strength of epistemic
 350 signals for optimizing predictive uncertainty. See Table S2 for a summary of analytical gradients used for
 351 the variance-gated normalization framework.

352 We summarize the forward and backward passes of VGN in Figure 2 and describe them step-by-step below:

353 **(a) Forward computation.** Each ensemble member $m \in \{1, \dots, M\}$ produces a categorical predictive
 354 distribution $\mathbf{p}_m \in \Delta^{C-1}$ (Step 1). The ensemble sample mean $\bar{\mathbf{p}}$ and predictive spread \mathbf{s} are then computed
 355 across members, summarizing ensemble consensus and disagreement (Step 2). Using these statistics, a
 356 shared variance gate is constructed as $\Gamma = 1 - e^{-\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}}$, where $\mathbf{k} > 0$ controls the gate sensitivity (Step 3).
 357 Each member distribution is modulated by the shared gate and re-normalized as $\mathbf{q}_m = \mathbf{p}_m \odot \Gamma / Z_m$, with
 358 $Z_m = \mathbf{p}_m^\top \Gamma$ (Step 4). Finally, the ensemble mixture $\bar{\mathbf{q}} = \frac{1}{M} \sum_{m=1}^M \mathbf{q}_m$ is formed, and the loss \mathcal{L} is applied to
 359 this mixture distribution (Step 5).

360 **(b) Backpropagation path.** Gradients with respect to the mixture are distributed equally to ensemble
 361 members as $\frac{\partial \mathcal{L}}{\partial \mathbf{q}_m} = \frac{1}{M} \frac{\partial \mathcal{L}}{\partial \bar{\mathbf{q}}}$ (Step 6). The total gradient with respect to each member prediction \mathbf{p}_m then
 362 decomposes into three additive contributions (Proposition 4.1): (i) a direct path through the normalization
 363 with the gate held fixed, $\partial \mathcal{L} / \partial \mathbf{p}_m |_{\Gamma}$; (ii) an indirect path through the ensemble mean $\bar{\mathbf{p}}$; and (iii) an indirect
 364 path through the predictive spread \mathbf{s} (Step 7). Both indirect paths propagate through the shared variance
 365 gate *via* $\partial \Gamma / \partial \bar{\mathbf{p}}$ and $\partial \Gamma / \partial \mathbf{s}$, coupling gradients across ensemble members (Step 8). The sensitivity parameter
 366 \mathbf{k} , reparameterized as $\mathbf{k} = \text{softplus}(\ell)$, receives gradients through its influence on the gate Γ (Step 9). All
 367 three gradient contributions are summed and scaled by $1/M$, yielding the final per-member gradient update
 368 (Step 10).

369 For completeness, all Jacobians, vector–Jacobian products, and gradients with respect to ensemble statistics
 370 and learnable gating parameters are derived in full in SI S2.

NEW:
cddg::M6

NEW:
cddg::M5

5 Experiments

We evaluate the proposed variance-gated ensemble framework on MNIST, SVHN, CIFAR-10, and CIFAR-100 using convolutional backbones and ensemble configurations commonly employed in uncertainty estimation. For MNIST we use a LeNet-5 style network, while for SVHN and CIFAR-10/100 we use WideResNet-28-10. Experiments consider DE, MCD, LLE, and the hybrid variant MCD-LLE, with VGN applied where it is well-defined. VGN is a training-time normalization layer that requires distinct, trainable ensemble members. It applies to DE-VGN and LLE-VGN, but is undefined for MCD (a single network with stochastic forward passes) and is not evaluated for MCD-LLE in this work. Since DE/DE-VGN require training M independent networks, we evaluate them at $M = 5$, consistent with (Lakshminarayanan et al., 2017); LLE/LLE-VGN and MCD/MCD-LLE are evaluated across $M \in \{5, 10, 100\}$. Models are trained using the Adam optimizer with early stopping, and all results are averaged over three trials with different random seeds. We compare VGMU against entropy-based EU (mutual information) and recent information-theoretic baselines, including EPKL and EPJS pairwise divergence measures. For variance-gated variants, the classwise gating parameter \mathbf{k} is learned end-to-end. Uncertainty scores were assessed *via* rank-based agreement with baseline methods (Spearman’s ρ and Kendall’s τ), Cumulative Area Under the Curve (AUC_c) for uncertainty mass concentration, and margin–variance geometry visualizations. Predictive performance and calibration are reported using accuracy, F1-score, Expected Calibration Error (ECE), and ensemble diversity. To quantify practical differences between methods, we report an effect size $(\bar{x}_1 - \bar{x}_2) / \max(\sigma_1, \sigma_2) \geq 1$, where bold values indicate a single best-performing value over the next nearest competitor. Complete network specifications, training protocols, evaluation definitions, and implementation details are provided in SI S3, and additional results for MNIST and SVHN appear in SI S4. MNIST and SVHN are included as complementary benchmarks to provide uncertainty behavior in low-noise image settings. The influence of the principal hyperparameters M (ensemble size), \mathbf{k} (learnable per-class sensitivity), and ensemble type is analyzed in the sections that follow: \mathbf{k} in Section 5.4 (with its theoretical sensitivity established in Proposition 3.3), M in Section 5.5 and Section 5.6, and the ensemble types throughout.

NEW:
cddgNEW:
cddg::M5NEW:
VoEx::4NEW:
cddg::M4

5.1 Rank Consistency with Existing Measures

To validate that VGMU captures uncertainty structure consistent with established measures, we compare its sample-level rankings against information-theoretic baselines. High correlation demonstrates consistency of uncertainty rankings, confirming that the computationally efficient margin-based VGMU score preserves the ordering produced by more expensive methods. Table 1 reports Spearman correlations for CIFAR-10 and CIFAR-100 (additional results in SI S4.1) computed for each testing dataset and ensemble configuration.

On CIFAR-10, VGMU align with all baselines, indicating that the margin-based score captures similar uncertainty structure as pairwise divergence measures and mutual information. For the CIFAR-100 dataset, correlations remain strong for MCD variants but diverged for LLE models (*e.g.*, EPKL, $\rho = 0.566$). This reflects a deliberate design choice rather than a limitation. Pairwise measures capture distributional disagreement across all 100 classes, while VGMU focuses exclusively on the decision-relevant margin between the top-2 predictions. When ensemble members agree on the most likely classes but disagree about the tail distribution, pairwise measures increase substantially while VGMU remains low. For practical decision-making (*e.g.*, whether to trust a prediction or defer to a human for review), disagreement about distant classes is irrelevant; the insensitivity of VGMU to such disagreement is a feature, moving away from information-theoretic approaches (Wimmer et al., 2023).

The addition of VGN improves the correlation (LLE, EPKL) with $\rho = 0.566$ to $\rho = 0.697$, indicating that normalization can partially recover sensitivity to distributional disagreement. Figure 3 illustrates this behavior through rank-rank scatter plots. Points cluster along the diagonal for CIFAR-10, while CIFAR-100 displays off-diagonal deviations for LLE models. These deviations correspond to samples where ensemble members agree on the top predictions but disagree about lower-ranked classes, the cases where decision-focused design of VGMU diverges from entropy-based scores.

Table 1: Spearman rank correlation (ρ) between VGMU and epistemic uncertainty baselines.^{1,2}

Dataset	M	Method	VGMU vs. EPJS	VGMU vs. EPKL	VGMU vs. EU
CIFAR-10	5	DE	0.975 ± 0.011	0.928 ± 0.032	0.962 ± 0.015
		DE-VGN	0.991 ± 0.002	0.977 ± 0.004	0.986 ± 0.003
		LLE	0.992 ± 0.002	0.987 ± 0.003	0.990 ± 0.002
		LLE-VGN	0.988 ± 0.002	0.981 ± 0.004	0.985 ± 0.003
		MCD	0.985 ± 0.003	0.981 ± 0.003	0.982 ± 0.003
		MCD-LLE	0.980 ± 0.004	0.979 ± 0.004	0.979 ± 0.004
	10	LLE	0.992 ± 0.002	0.985 ± 0.004	0.988 ± 0.004
		LLE-VGN	0.993 ± 0.002	0.987 ± 0.003	0.990 ± 0.002
		MCD	0.989 ± 0.002	0.985 ± 0.002	0.986 ± 0.002
		MCD-LLE	0.987 ± 0.002	0.986 ± 0.002	0.986 ± 0.002
	100	LLE	0.997 ± 0.000	0.974 ± 0.004	0.993 ± 0.001
		LLE-VGN	0.997 ± 0.001	0.986 ± 0.006	0.993 ± 0.002
MCD		0.992 ± 0.002	0.990 ± 0.002	0.990 ± 0.002	
MCD-LLE		0.990 ± 0.002	0.989 ± 0.002	0.989 ± 0.002	
CIFAR-100	5	DE	0.783 ± 0.008	0.530 ± 0.029	0.774 ± 0.008
		DE-VGN	0.791 ± 0.021	0.557 ± 0.027	0.781 ± 0.018
		LLE	0.870 ± 0.003	0.763 ± 0.020	0.856 ± 0.004
		LLE-VGN	0.881 ± 0.009	0.809 ± 0.023	0.867 ± 0.010
		MCD	0.868 ± 0.012	0.850 ± 0.011	0.860 ± 0.012
		MCD-LLE	0.863 ± 0.015	0.849 ± 0.015	0.856 ± 0.015
	10	LLE	0.854 ± 0.001	0.566 ± 0.035	0.814 ± 0.004
		LLE-VGN	0.875 ± 0.022	0.697 ± 0.017	0.843 ± 0.023
		MCD	0.874 ± 0.012	0.856 ± 0.011	0.862 ± 0.011
		MCD-LLE	0.874 ± 0.015	0.860 ± 0.016	0.864 ± 0.016
	100	LLE	0.724 ± 0.012	0.276 ± 0.050	0.587 ± 0.023
		LLE-VGN	0.790 ± 0.012	0.279 ± 0.033	0.588 ± 0.014
MCD		0.885 ± 0.015	0.869 ± 0.015	0.872 ± 0.015	
MCD-LLE		0.880 ± 0.016	0.867 ± 0.016	0.869 ± 0.016	

¹ Higher values indicate better alignment with information-theoretic measures. ² Bold values indicate a single best-performing method with an effect size $(\bar{x}_1 - \bar{x}_2) / \max(\sigma_1, \sigma_2) \geq 1$ over the next nearest competitor for each $M \in \{5, 10, 100\}$ and method per VGMU comparison.

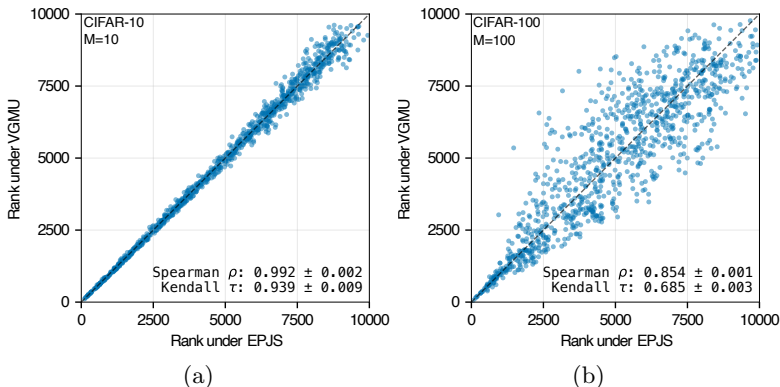


Figure 3: Rank consistency between VGMU and EPJS on CIFAR-10 (a) and CIFAR-100 (b) for the LLE models. The diagonal denotes perfect agreement between uncertainty rankings.

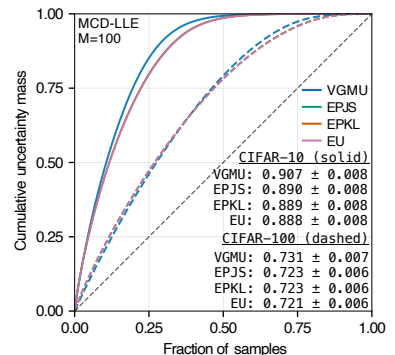


Figure 4: AUCc curves for CIFAR-10/100. The diagonal ($AUC_c = 0.5$), corresponds to no concentration on difficult samples.

Table 2: Cumulative area of the curve (AUC_c) scores for uncertainty mass concentration. Higher values indicate sharper concentration on difficult samples.¹

Dataset	M	Method	VGMU	EPJS	EPKL	EU
CIFAR-10	5	DE	0.759 ± 0.032	0.773 ± 0.030	0.781 ± 0.029	0.775 ± 0.030
		DE-VGN	0.829 ± 0.008	0.838 ± 0.007	0.841 ± 0.007	0.837 ± 0.007
		LLE	0.902 ± 0.010	0.898 ± 0.011	0.902 ± 0.001	0.896 ± 0.011
		LLE-VGN	0.888 ± 0.009	0.887 ± 0.007	0.892 ± 0.007	0.885 ± 0.007
		MCD	0.897 ± 0.007	0.895 ± 0.007	0.895 ± 0.007	0.893 ± 0.007
		MCD-LLE	0.907 ± 0.008	0.901 ± 0.007	0.901 ± 0.007	0.901 ± 0.007
	10	LLE	0.885 ± 0.015	0.873 ± 0.016	0.876 ± 0.015	0.868 ± 0.016
		LLE-VGN	0.885 ± 0.009	0.874 ± 0.010	0.874 ± 0.011	0.870 ± 0.011
		MCD	0.895 ± 0.007	0.883 ± 0.007	0.883 ± 0.007	0.881 ± 0.007
		MCD-LLE	0.907 ± 0.008	0.890 ± 0.008	0.889 ± 0.008	0.888 ± 0.008
	100	LLE	0.881 ± 0.012	0.856 ± 0.009	0.829 ± 0.001	0.834 ± 0.006
		LLE-VGN	0.872 ± 0.016	0.853 ± 0.018	0.838 ± 0.024	0.836 ± 0.020
MCD		0.893 ± 0.007	0.873 ± 0.008	0.871 ± 0.008	0.868 ± 0.008	
MCD-LLE		0.907 ± 0.008	0.890 ± 0.008	0.889 ± 0.008	0.888 ± 0.008	
CIFAR-100	5	DE	0.549 ± 0.002	0.575 ± 0.003	0.600 ± 0.003	0.591 ± 0.003
		DE-VGN	0.554 ± 0.003	0.581 ± 0.004	0.606 ± 0.003	0.596 ± 0.004
		LLE	0.666 ± 0.002	0.691 ± 0.003	0.721 ± 0.002	0.696 ± 0.003
		LLE-VGN	0.677 ± 0.008	0.702 ± 0.009	0.729 ± 0.008	0.706 ± 0.009
		MCD	0.722 ± 0.006	0.739 ± 0.006	0.743 ± 0.006	0.740 ± 0.006
		MCD-LLE	0.733 ± 0.007	0.745 ± 0.006	0.747 ± 0.006	0.745 ± 0.006
	10	LLE	0.628 ± 0.006	0.643 ± 0.008	0.668 ± 0.013	0.651 ± 0.007
		LLE-VGN	0.647 ± 0.021	0.666 ± 0.020	0.684 ± 0.019	0.673 ± 0.019
		MCD	0.719 ± 0.006	0.726 ± 0.005	0.728 ± 0.005	0.725 ± 0.005
		MCD-LLE	0.731 ± 0.007	0.723 ± 0.006	0.723 ± 0.006	0.721 ± 0.006
	100	LLE	0.552 ± 0.006	0.556 ± 0.007	0.592 ± 0.008	0.565 ± 0.006
		LLE-VGN	0.555 ± 0.002	0.558 ± 0.002	0.587 ± 0.003	0.566 ± 0.002
MCD		0.718 ± 0.006	0.713 ± 0.006	0.713 ± 0.006	0.711 ± 0.006	
MCD-LLE		0.731 ± 0.007	0.723 ± 0.006	0.723 ± 0.006	0.721 ± 0.006	

¹ Bold values indicate a single best-performing measure with an effect size $(\bar{x}_1 - \bar{x}_2) / \max(\sigma_1, \sigma_2) \geq 1$ over the next nearest competitor for each $M \in \{5, 10, 100\}$ and method.

5.2 Uncertainty Mass Concentration

For practical deployment in selective prediction or human-in-the-loop systems, uncertainty estimates should concentrate on difficult samples; otherwise they become non-informative for decision-making under uncertainty.

We quantify this through AUC_c , calculated by sorting samples by descending uncertainty and measuring the area under the cumulative distribution. Higher AUC_c values indicate more efficient concentration of uncertainty on difficult samples, while $AUC_c = 0.5$ no concentration. Table 2 summarizes AUC_c across datasets (additional information provided in SI S4.2).

On CIFAR-10, VGMU is comparable across all evaluated ensemble configurations. In several cases, the separation in mean values relative to information-theoretic baselines is larger than the observed run-to-run variability, indicating a tendency toward stronger concentration of uncertainty on difficult samples. From a practical perspective, this behavior implies that VGMU can prioritize challenging samples more efficiently, requiring fewer predictions to capture a fixed proportion of the total uncertainty mass. This trend is less pronounced for CIFAR-100. Information-theoretic measures show better performing AUC_c values than VGMU for all settings, reflecting the insensitivity of VGMU to disagreement across the full class simplex. This observation is consistent with the rank-correlation analysis, where information-theoretic measures capture forms of uncertainty that VGMU intentionally de-emphasizes. However, AUC_c can be restored using the hybrid MCD-LLE configuration (although difference within observed variability). This

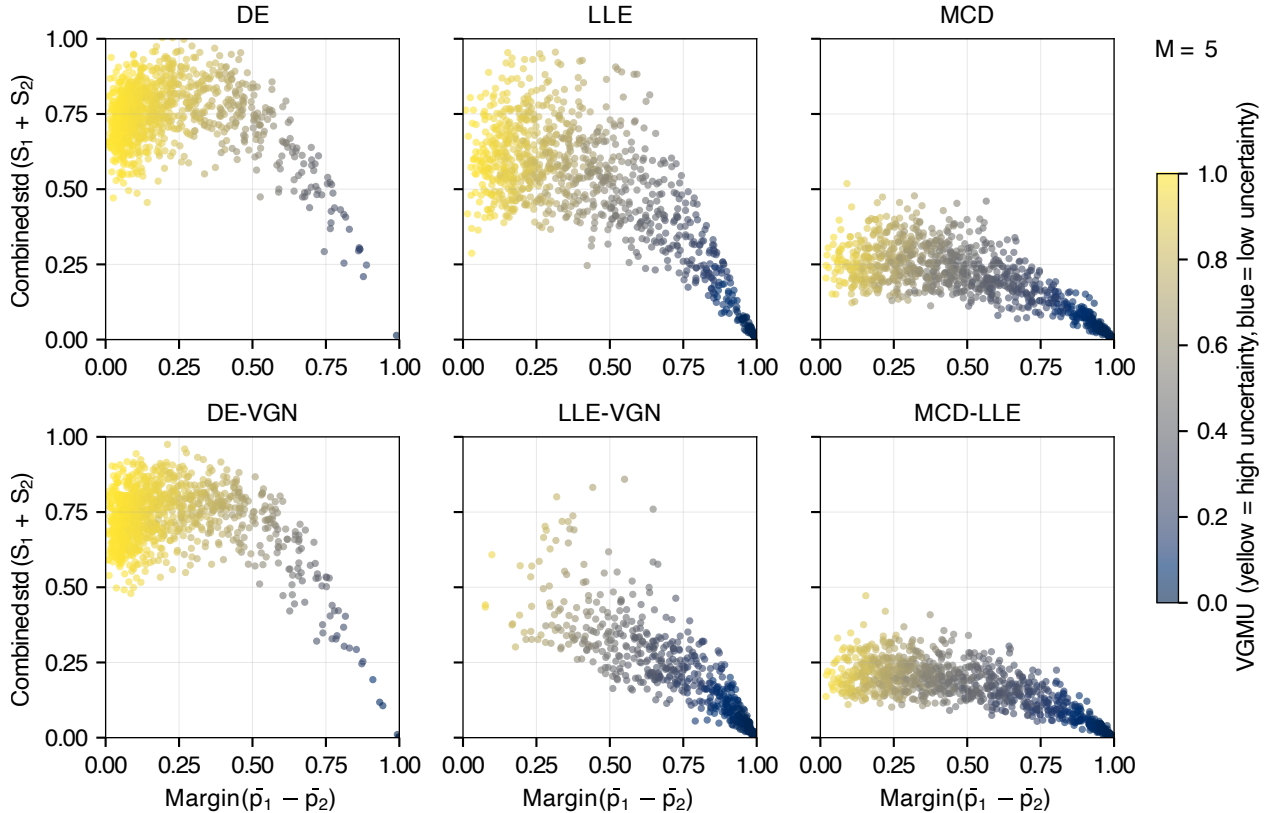


Figure 5: Margin-variance geometry for CIFAR-100 with $M = 5$. Each point represents a test sample; color indicates VG MU value (yellow = high uncertainty, blue = low). DE and DE-VGN show high variance even at large margins, LLE and LLE-VGN show moderate variance, while MCD and MCD-LLE concentrates samples in the high-margin (confident), low-variance (certain) region.

436 pattern suggests that increasing ensemble diversity through the combination of stochastic sampling and
 437 multiple classifier heads may help recover decision-relevant uncertainty structure, allowing margin-based
 438 estimation to remain competitive on more challenging tasks. Figure 4 shows cumulative uncertainty mass
 439 concentration curves for CIFAR-10 and CIFAR-100 using MCD-LLE models ($M = 100$). For CIFAR-10,
 440 uncertainty mass is concentrated among a relatively small fraction of sample (*ca.* 25%), VG MU shows a
 441 small but consistent upward shift relative to the information-theoretic baselines. In this case of CIFAR-100,
 442 the signals are noticeably closer to the diagonal ($AUC_c \approx 0.5$), indicating substantially weaker concentration.
 443 This collapse in concentration is consistent with the increased difficulty of CIFAR-100, where uncertainty
 444 is spread across a broader set of samples and full-simplex disagreement becomes more prominent. SI S4.2
 445 results indicate a general trend where VG MU adopts a conservative uncertainty assessment. In these cases,
 446 uncertainty is distributed more broadly across samples (*i.e.*, $AUC_c \approx 0.5$) rather than being concentrated,
 447 consistent with the margin-based design of VG MU.

NEW:
cddg::M5

448 5.3 Margin-Variance Geometry

449 The VG MU score explicitly couples the predictive margin ($\bar{p}_1 - \bar{p}_2$) with ensemble variance ($s_1 + s_2$), providing
 450 a two-dimensional uncertainty landscape. This geometric perspective provides insight into how different en-
 451 semble methods populate the margin-variance space and how VG MU responds to these configurations. Fig-
 452 ure 5 illustrates the margin-variance landscape for CIFAR-100 (additional visualizations provided in SI S4.3).
 453 Several consistent patterns emerge.

NEW:
cddg::M6

454 **Ensemble method signatures.** Different ensemble strategies show characteristic distributions in margin-
 455 variance space. DE and DE-VGN scatter broadly, showing substantial variance even at high margins. LLE
 456 and LLE-VGN occupy regions with moderate variance and moderate-to-high margins, while MCD and
 457 MCD-LLE produce more compact clusters with lower overall variance.

458 **VGMU response.** The behavior of VGMU across this space confirms its sensitivity to both margin and
 459 variance. Low VGMU values (blue) occur predominantly when the predictive margin is large and variance is
 460 low, whereas high values (yellow) arise from either small margins or elevated variance. The variance-gated
 461 formulation induces smooth transitions rather than hard decision boundaries, enabling graded confidence
 462 assessments.

463 **Effect of VGN.** Comparing LLE and LLE-VGN reveals a shift toward lower VGMU values, with sam-
 464 ples concentrating in the high-margin, low-variance region while high-uncertainty regions remain populated.
 465 This reduction in variance is achieved through VGN, which suppresses high-variance class predictions during
 466 training. Similar trends are observed across other datasets and configurations in SI S4.3, although in some
 467 cases the effects are subtle. The margin-variance geometry provides an intuitive interpretation of how, pre-
 468 dictive margin and epistemic disagreement interact. By coupling margin and variance, VGMU assigns low
 469 uncertainty only when predictions are well-separated and consistent across ensemble members, while remain-
 470 ing conservative in regions characterized by ambiguity or disagreement. The effect of VGN is to reshape this
 471 geometry during training by suppressing high-variance class predictions, leading to a greater concentration
 472 of samples in the high-margin, low-variance region (Figure 5: LLE *vs.* LLE-VGN). This geometric behavior
 473 is consistent with the rank-based and concentration analyses, and shows how VGN complements VGMU by
 474 promoting decision-relevant uncertainty with conservative assessments in challenging settings.

475 5.4 Calibration and Performance

476 Table 3 reports calibration metrics for CIFAR-10 and CIFAR-100 (additional results in SI S4.5). For
 477 DE on CIFAR-10, VGN reduces mean ECE by a margin that exceeds the observed run-to-run variability,
 478 while also improving classification accuracy, indicating improved alignment with predictive confidence. In
 479 contrast, for CIFAR-100 and for other ensemble settings, calibration effects are smaller and less consistent,
 480 with differences falling within the variability across runs.

481 Overall, the results indicate that incorporating VGN does not adversely affect calibration and, in several
 482 configurations, provides small but consistent improvements. Similar trends are observed for the MCD-LLE
 483 setting, where VGN tends to improve or preserve calibration rather than degrade it.

484 **Learned k values.** Figure 6 displays the per-class k parameters learned by VGN models (see Figure S5
 485 for CIFAR-100). DE-VGN learns higher values ($\bar{k} \approx 2.5\text{--}4.0$) than LLE-VGN ($\bar{k} \approx 0.75\text{--}1.1$), reflecting the
 486 diversity in deep ensembles. Higher k reduces gate sensitivity (see Proposition 3.3 for additional details),
 487 allowing DE-VGN to tolerate the disagreement among independently trained members. In these examples,
 488 the relative consistency of k across classes suggests that it adapts to ensemble-level diversity rather than
 489 class-specific difficulty. However, results on SVHN (see Figure S4) indicate that k captures both effects,
 490 varying across classes while also reflecting overall ensemble diversity.

491 5.5 Computational Efficiency

492 A key advantage of the VGE framework is computational efficiency. Table 4 summarizes the asymptotic
 493 complexity of each uncertainty measure. VGMU requires only $O(C)$ operations after ensemble moments have
 494 been computed, compared to $O(M^2C)$ for pairwise divergence measures such as EPKL and EPJS. The VGN
 495 uncertainty decomposition (TU, AU, EU) operates at $O(MC)$, matching the cost of standard entropy-based
 496 decompositions. The quadratic scaling of pairwise methods becomes prohibitive for large ensembles. For
 497 example, with $M = 100$ members and $C = 100$ classes, pairwise measures require 10^6 operations per sample,
 498 compared to 10^4 for VGN decomposition and 10^2 for VGMU.

NEW:
VoEx::4NEW:
cddg::C3NEW:
VoEx::4

Table 3: Performance and calibration metrics.¹

Dataset	M	Method	Accuracy	F1-Score	ECE	Diversity ²	
CIFAR-10	5	DE	0.836 ± 0.012	0.836 ± 0.012	0.049 ± 0.014	$1.9 \times 10^{-2} \pm 0.2$	
		DE-VGN	0.875 ± 0.005	0.875 ± 0.006	0.023 ± 0.003	$9.8 \times 10^{-3} \pm 0.5$	
		LLE	0.855 ± 0.002	0.855 ± 0.002	0.062 ± 0.012	$2.4 \times 10^{-3} \pm 0.3$	
		LLE-VGN	0.846 ± 0.008	0.845 ± 0.008	0.067 ± 0.005	$2.3 \times 10^{-3} \pm 0.5$	
		MCD	0.852 ± 0.007	0.851 ± 0.007	0.057 ± 0.002	$1.4 \times 10^{-3} \pm 0.1$	
		MCD-LLE	0.852 ± 0.009	0.852 ± 0.008	0.067 ± 0.004	$3.4 \times 10^{-4} \pm 0.1$	
	10	LLE	0.848 ± 0.003	0.848 ± 0.002	0.058 ± 0.010	$3.6 \times 10^{-3} \pm 0.6$	
		LLE-VGN	0.849 ± 0.004	0.849 ± 0.003	0.066 ± 0.003	$3.2 \times 10^{-3} \pm 0.2$	
		MCD	0.852 ± 0.007	0.851 ± 0.006	0.055 ± 0.003	$1.5 \times 10^{-3} \pm 0.1$	
		MCD-LLE	0.853 ± 0.009	0.852 ± 0.008	0.067 ± 0.004	$3.4 \times 10^{-4} \pm 0.2$	
		100	LLE	0.851 ± 0.004	0.850 ± 0.005	0.065 ± 0.006	$5.8 \times 10^{-3} \pm 0.8$
			LLE-VGN	0.847 ± 0.005	0.848 ± 0.004	0.071 ± 0.006	$5.4 \times 10^{-3} \pm 0.8$
MCD	0.853 ± 0.008		0.853 ± 0.007	0.053 ± 0.004	$1.7 \times 10^{-3} \pm 0.1$		
MCD-LLE	0.853 ± 0.008		0.852 ± 0.008	0.067 ± 0.004	$3.6 \times 10^{-4} \pm 0.2$		
CIFAR-100	5	DE	0.487 ± 0.002	0.481 ± 0.003	0.095 ± 0.006	$3.5 \times 10^{-3} \pm 0.1$	
		DE-VGN	0.507 ± 0.007	0.504 ± 0.007	0.086 ± 0.006	$3.1 \times 10^{-3} \pm 0.2$	
		LLE	0.545 ± 0.002	0.541 ± 0.003	0.074 ± 0.002	$1.6 \times 10^{-3} \pm 0.0$	
		LLE-VGN	0.548 ± 0.005	0.545 ± 0.007	0.111 ± 0.007	$1.3 \times 10^{-3} \pm 0.1$	
		MCD	0.564 ± 0.012	0.562 ± 0.012	0.092 ± 0.001	$3.0 \times 10^{-4} \pm 0.2$	
		MCD-LLE	0.563 ± 0.013	0.562 ± 0.013	0.106 ± 0.001	$1.9 \times 10^{-4} \pm 0.1$	
	10	LLE	0.543 ± 0.005	0.538 ± 0.005	0.062 ± 0.004	$2.4 \times 10^{-3} \pm 0.1$	
		LLE-VGN	0.548 ± 0.014	0.542 ± 0.014	0.095 ± 0.026	$2.0 \times 10^{-3} \pm 0.2$	
		MCD	0.567 ± 0.012	0.565 ± 0.012	0.086 ± 0.001	$3.4 \times 10^{-4} \pm 0.2$	
		MCD-LLE	0.564 ± 0.014	0.563 ± 0.014	0.103 ± 0.003	$2.2 \times 10^{-4} \pm 0.1$	
		100	LLE	0.537 ± 0.010	0.533 ± 0.012	0.059 ± 0.019	$4.3 \times 10^{-3} \pm 0.2$
			LLE-VGN	0.544 ± 0.012	0.542 ± 0.011	0.059 ± 0.009	$4.2 \times 10^{-3} \pm 0.1$
MCD	0.568 ± 0.012		0.566 ± 0.012	0.082 ± 0.002	$3.7 \times 10^{-4} \pm 0.2$		
MCD-LLE	0.565 ± 0.013		0.563 ± 0.013	0.102 ± 0.001	$2.3 \times 10^{-4} \pm 0.1$		

¹ Bold values indicate a single best-performing measure with an effect size $(\bar{x}_1 - \bar{x}_2) / \max(\sigma_1, \sigma_2) \geq 1$ over the next nearest competitor for each $M \in \{5, 10, 100\}$ and method. ² Defined as the ensemble variance averaged across samples and classes, $\mathbb{E}_{i,c}[\text{Var}_M]$.

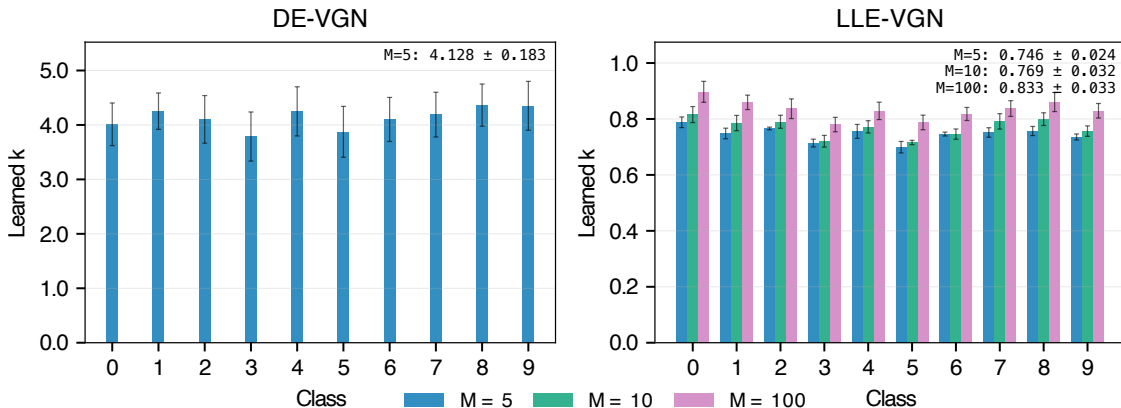


Figure 6: Learned per-class k values for VGN models on CIFAR-10. DE-VGN learns higher values ($\bar{k} \approx 4.1$) than LLE-VGN ($\bar{k} \approx 0.8$), reflecting adaptation to ensemble diversity.

499 **Empirical scaling with M and C .** To validate these asymptotic predictions empirically, we mea-
500 sured wall-clock inference times on synthetic ensemble outputs across a range of ensemble sizes $M \in$
501 $\{2, 5, 10, 20, 50, 100\}$ and class counts $C \in \{10, 100, 1000\}$ on an NVIDIA RTX 4090 GPU (Intel Core i9-

Table 4: Asymptotic complexity of uncertainty measures for M ensemble members and C classes.

Measure	Complexity	Stage
VGMU	$O(C)$	Inference (post-hoc)
VGN decomposition (TU, AU, EU)	$O(MC)$	Training / Inference
Standard entropy decomposition	$O(MC)$	Inference
EPKL / EPJS	$O(M^2C)$	Inference

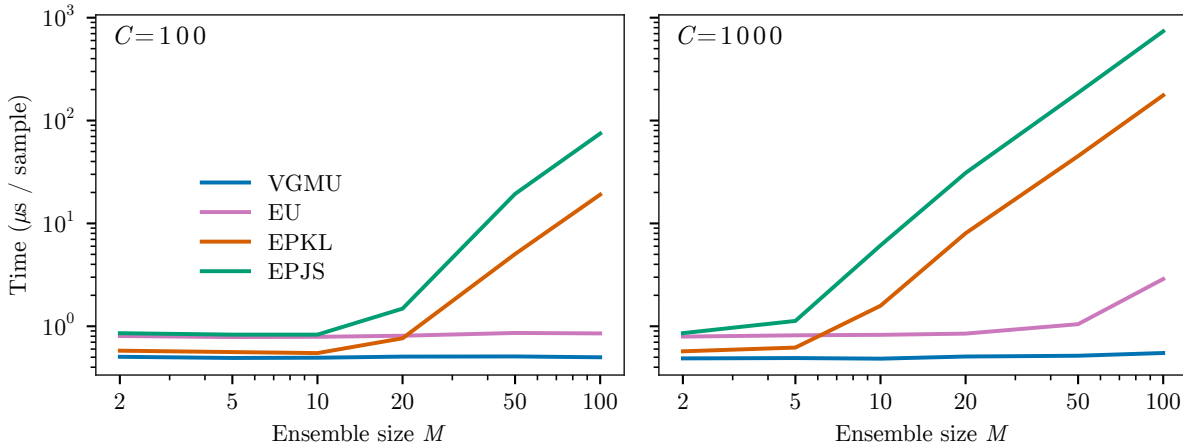


Figure 7: Wall-clock inference time per sample as a function of ensemble size M for $C = 100$ classes and $C = 1000$ classes. VGMU remains constant regardless of M , while pairwise measures (EPKL, EPJS) result in quadratic growth consistent with their $O(M^2C)$ complexity. At $M = 100$, $C = 1000$, VGMU is over $320\times$ faster than EPKL and $1,343\times$ faster than EPJS.

502 13900KF; 32 GB RAM), averaging over 10 trials of 128 samples per configuration. Figure 7 displays the
 503 log-log scaling behavior, and Table 5 reports representative timings for $C = 1000$. VGMU remains near-
 504 constant at approximately $0.5 \mu\text{s}$ per sample regardless of M , confirming its $O(C)$ complexity. In contrast,
 505 EPKL and EPJS exhibit the expected quadratic scaling. At $M = 100$ with $C = 100$, EPKL requires $19 \mu\text{s}$
 506 and EPJS requires $75 \mu\text{s}$, representing $38\times$ and $150\times$ speedup for VGMU, respectively. At $M = 100$ with
 507 $C = 1000$, VGMU achieves a $320\times$ speedup over EPKL and a $1,342\times$ speedup over EPJS. EU scales linearly
 508 in M but remains within an order of magnitude of VGMU for moderate ensemble sizes; however, it diverges
 509 substantially at $M = 100$, $C = 1000$ ($5.2\times$ slower). These results confirm that the computational advan-
 510 tages of VGMU are not merely asymptotic but provide practical speedups of several orders of magnitude
 511 in realistic configurations, enabling real-time uncertainty estimation in deployment scenarios where pairwise
 512 methods would be prohibitive. Full results across all M and C combinations are provided in Table S7.

513 5.6 Out-of-Distribution Detection

514 To evaluate whether VGN and/or the VGMU score provides competitive OOD detection, we train models
 515 on SVHN (in-distribution, ID) and evaluate against CIFAR-10 (OOD). Detection performance is measured
 516 using the area under the ROC curve (AUC) and the false positive rate at 95% true positive rate (FPR@95).
 517 All LLE and LLE-VGN models for SVHN attained *ca.* 95% accuracy with strong calibration (ECE ≈ 0.014).
 518 As in prior sections, we interpret differences conservatively, distinguishing trends whose magnitude exceeds
 519 run-to-run variability.

520 **LLE and LLE-VGN across ensemble sizes.** Figure 8 reports OOD detection results for LLE and LLE-
 521 VGN across $M \in \{5, 10, 100\}$. At $M = 5$, VGMU achieves the highest mean AUC (0.941 ± 0.006 for LLE;
 522 0.942 ± 0.006 for LLE-VGN) and lowest mean FPR@95 (0.173 ± 0.027 ; 0.167 ± 0.023 , respectively). The
 523 AUC differences between VGMU and EPJS or EU fall within overlapping variability, indicating comparable

Table 5: Wall-clock inference time ($\mu\text{s}/\text{sample}$) and VG MU speedup over baselines for representative ensemble sizes M with $C = 1000$ classes.

M	Time ($\mu\text{s}/\text{sample}$)				VG MU Speedup		
	VG MU	EPKL	EPJS	EU	EPKL	EPJS	EU
2	0.5	0.6	0.9	0.8	1.2	1.8	1.6
5	0.5	0.6	1.1	0.8	1.3	2.3	1.7
10	0.5	1.6	6.1	0.8	3.3	13	1.7
20	0.5	8.0	31	0.8	16	61	1.7
50	0.5	45	186	1.0	87	360	2.0
100	0.5	176	738	2.9	320	1343	5.2

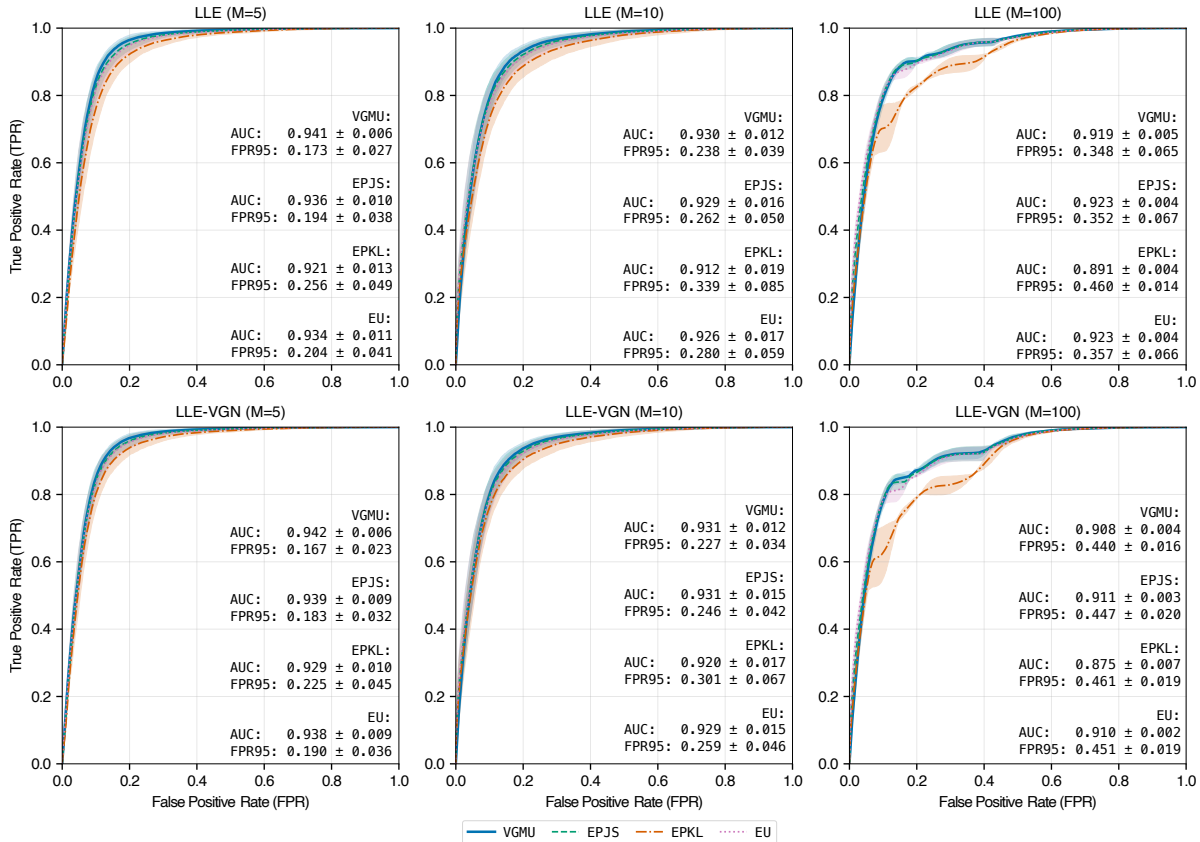


Figure 8: ROC curves for OOD detection (SVHN, ID → CIFAR-10, OOD) across LLE (top) and LLE-VGN (bottom) with $M \in \{5, 10, 100\}$. At small ensemble sizes ($M = 5$), VG MU shows a visible separation from EPKL in both AUC and FPR@95, while differences with EPJS and EU fall within observed variability. As M increases, FPR@95 degrades for all methods. VGN provides consistent but modest FPR@95 improvements at $M = 5$ and $M = 10$, with degradation at $M = 100$.

524 detection performance among these measures. However, the separation from EPKL is more pronounced.
 525 For LLE ($M = 5$), the FPR@95 gap between VG MU (0.173 ± 0.027) and EPKL (0.256 ± 0.049) exceeds the
 526 observed run-to-run variability.

527 At larger M , the ID distribution over EPKL shifts into a range also occupied by a low-uncertainty subset
 528 of OOD samples, increasing ID–OOD overlap in that region (SI Figure S8). With $M = 100$, VG MU
 529 and EPJS perform comparably in AUC (*ca.*, 0.919–0.923), while EPKL degrades more substantially (AUC

530 = 0.891 ± 0.004 for LLE; 0.875 ± 0.007 for LLE-VGN), consistent with the unbounded range of pairwise KL
 531 amplifying tail-disagreement. Incorporating VGN provides a consistent but modest improvement in FPR@95
 532 at small-to-moderate ensemble sizes. At $M = 5$, LLE-VGN reduces mean FPR@95 from 0.173 to 0.167; at
 533 $M = 10$, from 0.238 to 0.227. While these improvements fall within seed-level variability, the direction is
 534 consistent across both ensemble sizes.

535 Across all configurations, VGMMU achieves comparable detection performance while avoiding the $O(M^2C)$
 536 computational cost of pairwise divergence measures, relying instead on ensemble moments with $O(MC)$
 537 complexity.

538 **MCD-LLE: Effect of ensemble configuration.** To examine how the source of ensemble diversity affects
 539 OOD detection, we evaluate MCD-LLE under varying head (H) and sampling (S) configurations (SI S4.7).
 540 VGMMU remains remarkably stable across all configurations, maintaining $\text{AUC} = 0.934 \pm 0.006$ and FPR@95
 541 ≈ 0.216 regardless of the head-sampling configurations. In contrast, information-theoretic baselines improve
 542 as either H or S increases, with mean AUC rising from 0.911 ($M = 5$) to 0.916 ($M = 100$) for EPJS. The
 543 AUC gap between VGMMU and the baselines (≈ 0.018 – 0.023) exceeds the observed run-to-run variability
 544 across all configurations, representing a consistent separation. This stability arises because VGMMU depends
 545 only on the top-2 margin and its associated variance, both of which saturate quickly, even with modest
 546 ensemble sizes. From a practical perspective, this insensitivity to ensemble configuration is advantageous.
 547 VGMMU can provide reliable OOD signals without rigorous tuning of MCD-LLE configurations.

548 **Summary.** Across ensemble types and configurations, VGMMU provides OOD detection performance that is
 549 comparable to or exceeds information-theoretic baselines. In the LLE setting, VGMMU and EPJS/EU produce
 550 similar AUC values with overlapping variability, while VGMMU consistently outperforms EPKL in FPR@95
 551 by a margin exceeding run-to-run variation. In the MCD-LLE setting, the separation between VGMMU and
 552 all baselines is more robust, with AUC differences that consistently exceed the observed standard deviations.
 553 VGN provides modest FPR@95 improvements at small ensemble sizes but does not enhance OOD detection
 554 at large M .

555 6 Discussion

556 6.1 Axiomatic Analysis

557 We introduced VGE as a computationally efficient alternative to entropy- and divergence-based uncertainty
 558 decomposition, with sensitivity to epistemic disagreement. Table 6 positions VGN within the axiomatic
 559 framework of Wimmer et al. (2023). These axioms highlight where standard entropy-based decompositions
 560 fail under finite ensembles and where disagreement-based alternatives improve epistemic sensitivity (for
 561 additional examples and discussion, see SI S6; formal proofs and justifications of axiomatic compliance
 562 for each axiom A0–A5 are provided in SI S7). VGN satisfies non-negativity (A0) and vanishing epistemic
 563 uncertainty for identical ensemble members (A1). This behavior is reflected in stable uncertainty estimates
 564 on low-noise datasets such as MNIST and in confident predictions on CIFAR-10 when ensemble members
 565 agree, consistent with theoretical expectations and information-theoretic baselines. VGN satisfies the EU
 566 component of the mean-preserving variance axiom (A3), consistent with pairwise divergence methods. The
 567 TU component is not satisfied, as suppression of high-variance classes by the gate can reduce total predictive
 568 entropy even as epistemic uncertainty increases (see SI S7). The margin-variance geometry in Section 5.3
 569 further supports this interpretation. Samples with increased ensemble spread are consistently assigned higher
 570 VGMMU values, even when mean confidence remains high. Standard entropy-based decompositions violate
 571 axioms A2 and A4, relating to the behavior under uniform predictions and injected noise. VGN partially
 572 satisfies A2 (with maximal EU attained for vertex-spanning ensembles but not guaranteed for all uniform-
 573 mean configurations due to \mathbf{k} -dependence) and fully satisfies A4 (AU and TU increase under center-shift).
 574 These improvements over entropy-based decompositions help explain the weaker alignment and reduced
 575 uncertainty mass concentration observed for mutual-information-based epistemic uncertainty in Section 5.2.
 576 VGN addresses these limitations by suppressing high-variance class probabilities prior to normalization,

NEW:
pDhd::7

Table 6: Comparison of uncertainty decomposition frameworks. Each row corresponds to an axiom (A0–A5), and each column indicates whether a given framework satisfies an axiom¹.

Axioms Wimmer et al. (2023)	Standard Entropy Decomposition Houlsby et al. (2011)	Expected Pairwise CE/KL Divergence Schweighofer et al. (2023)	Variance-Gated Normalization (proposed)
A0: TU, AU, EU ≥ 0 (non-negativity)	✓	✓	✓
A1: EU = 0 for identical ensemble members.	✓	✓	✓
A2: EU, TU maximal when ensembles are uniform distributions.	✗	●	●
A3: EU [↑] , TU [↑] with mean-preserving variance increases.	✓	✓	●
A4: AU [↑] , TU [↑] with addition of uniform noise.	✗	●	✓
A5: EU invariant to variance-preserving location shifts.	✓	✓	✓

¹ ✓ Satisfied, ✗ violated, ● partially satisfied.

577 reshaping ensemble predictions using moment-based statistics rather than relying exclusively on entropy
578 identities.

579 Invariance to variance-preserving location shifts (A5), satisfied by both pairwise divergence methods and
580 VGN, explains why VGMU remains stable under changes in tail-class disagreement. As shown in [Sec-](#)
581 [tion 5.1](#), this invariance accounts for the divergence between VGMU and full-simplex disagreement measures
582 on CIFAR-100. When ensemble members agree on the most likely classes but differ in low-probability regions,
583 pairwise divergences increase substantially, whereas VGMU remains low. VGMU is decision-focused, quan-
584 tifying uncertainty through the margin between the top-ranked classes modulated by predictive variance.
585 Disagreement that does not affect the decision boundary is intentionally de-emphasized, aligning uncertainty
586 estimation with selective prediction and human-in-the-loop decision-making.

587 The axiomatic analysis and experimental results together explain why variance-gated ensembles retain the
588 desirable properties of disagreement-based uncertainty measures while avoiding their principal computational
589 drawbacks, operating at $O(MC)$ for uncertainty decomposition and $O(C)$ for VGMU evaluation. The learned
590 sensitivity parameter \mathbf{k} adapts to ensemble diversity and task difficulty, supporting the interpretation of
591 variance-gating as a data-driven mechanism for epistemic control.

FIX:
cddg::C2

592 6.2 Limitations and Future Work

593 On CIFAR-100, VGMU shows weaker AUC_c concentration and lower rank correlation with full-simplex
594 measures compared to CIFAR-10. This reduction is a direct consequence of the top-2 margin design. On a
595 100-class problem, pairwise measures capture disagreement across all classes, while VGMU intentionally de-
596 emphasizes disagreement about low-ranked classes. This represents an expected trade-off between decision-
597 focused efficiency and full-simplex sensitivity, rather than a failure of the method. Applications requiring
598 sensitivity to distributional disagreement beyond the decision boundary may benefit from hybrid approaches
599 combining margin-based and distributional signals.

600 The proposed framework suggests several additional directions for future work. First, the current formulation
601 relies on first- and second-order ensemble moments. While this enables linear-time computation, higher-order
602 statistics or alternative measures may capture additional structure in multimodal predictive distributions.

603 Second, the gating function itself is defined using a specific exponential form. This choice was motivated
604 by smoothness, monotonicity with respect to confidence and variance, and well-behaved gradients for end-
605 to-end optimization. Alternative gating functions or parameterizations may provide different trade-offs
606 between sensitivity and saturation, and exploring such designs remains an open direction. Finally, extending
607 variance-gating to prediction tasks with severe class imbalance remains an open problem. In these settings,
608 decision margins may be ill-defined and ensemble variance may conflate epistemic uncertainty with class-
609 frequency effects. Exploring interactions between VGN and other epistemic-aware training objectives, such
610 as diversity-promoting loss functions, is another promising direction, as these objectives may interact non-
611 trivially through shared gradient pathways.

612 7 Conclusion

613 We introduced VGE, an epistemic-aware uncertainty framework that leverages ensemble disagreement
614 through a signal-to-noise gating principle. VGE unifies two complementary components operating at dif-
615 ferent stages of the learning pipeline: VGMU, a lightweight decision-based uncertainty score applicable at
616 inference without retraining, and VGN, a differentiable normalization layer that modulates high-variance
617 predictions during training. We derived closed-form vector–Jacobian products that enable end-to-end op-
618 timization of VGN through ensemble sample means and variances, allowing uncertainty sensitivity to be
619 learned directly from data. Empirically, across MNIST, SVHN, CIFAR-10, and CIFAR-100, VGMU exhib-
620 ited strong alignment with information-theoretic uncertainty measures, while revealing a distinct behavior
621 on CIFAR-100. In particular, where pairwise divergence measures emphasize full-simplex disagreement,
622 VGMU prioritizes the decision-relevant margin between top-ranked classes. Despite this difference, VGMU
623 achieved comparable uncertainty ranking and out-of-distribution detection performance at a fraction of the
624 computational cost, enabling real-time uncertainty estimation. Incorporating VGN during training further
625 suppressed high-variance predictions and reshaped ensemble member distributions, with learned sensitivity
626 parameters adapting to ensemble diversity across models. Overall, VGE provides a computationally efficient
627 approach to epistemic uncertainty estimation, supporting post hoc evaluation, training-time distribution
628 shaping, and end-to-end integration. These results suggest that decision-focused margin structure offers a
629 practical alternative to pairwise divergence-based uncertainty, particularly in large-class settings.

630 Broader Impact Statement

631 Reliable uncertainty estimation is important for deploying machine learning systems in risk-sensitive settings.
632 This work introduces a computationally efficient framework for epistemic-aware uncertainty estimation in
633 ensemble models, enabling real-time uncertainty assessment in large ensembles and many-class problems. As
634 with all uncertainty estimates, the outputs should be used as supporting signals rather than as standalone
635 decision criteria. Overall, this work aims to advance the safe, efficient, and responsible deployment of machine
636 learning models by making epistemic uncertainty estimation more accessible and scalable.

637 References

- 638 Arsenii Ashukha, Alexander Lyzhov, Dmitry Molchanov, and Dmitry Vetrov. Pitfalls of in-domain uncer-
639 tainty estimation and ensembling in deep learning. In *International Conference on Learning Representa-*
640 *tions (ICLR)*, 2021. doi: 10.48550/arXiv.2002.06470.
- 641 Stefan Depeweg, José Miguel Hernández-Lobato, Finale Doshi-Velez, and Steffen Udluft. Decomposition of
642 uncertainty in Bayesian deep learning for efficient and risk-sensitive learning. In *International Conference*
643 *on Machine Learning (ICML)*, 2018. doi: 10.48550/arXiv.1710.07283.
- 644 Stanislav Fort, Huiyi Hu, and Balaji Lakshminarayanan. Deep ensembles: A loss landscape perspective.
645 *arXiv*, 2020. doi: 10.48550/arXiv.1912.02757.
- 646 Yarín Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty
647 in deep learning. In *International Conference on Machine Learning (ICML)*, 2016. doi: 10.48550/arXiv.
648 1506.02142.

- 649 Yarin Gal, Jiri Hron, and Alex Kendall. Concrete dropout. *arXiv*, 2017. doi: 10.48550/arXiv.1705.07832.
- 650 Jakob Gawlikowski, Cedrique Rovile Njiteutcheu Tassi, Mohsin Ali, Jongseok Lee, Matthias Humt, Jianxiang
651 Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, Muhammad Shahzad, Wen Yang,
652 Richard Bamler, and Xiao Xiang Zhu. A survey of uncertainty in deep neural networks. *Artificial Intelli-*
653 *gence Review*, 2023. doi: 10.1007/s10462-023-10562-9.
- 654 Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. *arXiv*, 2017. doi:
655 10.48550/arXiv.1705.08500.
- 656 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In
657 *International Conference on Machine Learning (ICML)*, 2017. doi: 10.48550/arXiv.1706.04599.
- 658 James Harrison, John Willes, and Jasper Snoek. Variational Bayesian last layers. In *International Conference*
659 *on Learning Representations (ICLR)*, 2024. doi: 10.48550/arXiv.2404.11599.
- 660 Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classifi-
661 cation and preference learning. *arXiv*, 2011. doi: 10.48550/arXiv.1112.5745.
- 662 Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E. Hopcroft, and Kilian Q. Weinberger. Snapshot
663 ensembles: Train 1, get M for free. In *International Conference on Learning Representations (ICLR)*,
664 2017. doi: 10.48550/arXiv.1704.00109.
- 665 Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: An
666 introduction to concepts and methods. *Machine Learning*, 2021. doi: 10.1007/s10994-021-05946-3.
- 667 Ajay J. Joshi, Fatih Porikli, and Nikolaos Papanikolopoulos. Multi-class active learning for image classifi-
668 cation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. doi: 10.1109/CVPR.
669 2009.5206627.
- 670 Ananya Kumar, Tengyu Ma, Percy Liang, and Aditi Raghunathan. Calibrated ensembles can mitigate
671 accuracy tradeoffs under distribution shift. In *Conference on Uncertainty in Artificial Intelligence (UAI)*,
672 2022. doi: 10.48550/arXiv.2207.08977.
- 673 Kaisar Kushibar, Víctor Manuel Campello, Lidia Garrucho Moras, Akis Linardos, Petia Radeva, and Karim
674 Lekadir. Layer ensembles: A single-pass uncertainty estimation in deep learning for segmentation. *arXiv*,
675 2022. doi: 10.48550/arXiv.2203.08878.
- 676 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncer-
677 tainty estimation using deep ensembles. In *Conference on Neural Information Processing Systems (NIPS)*,
678 2017. doi: 10.48550/arXiv.1612.01474.
- 679 Stefan Lee, Senthil Purushwalkam, Michael Cogswell, David Crandall, and Dhruv Batra. Why M heads are
680 better than one: Training a diverse ensemble of deep networks. *arXiv*, 2015. doi: 10.48550/arXiv.1511.
681 06314.
- 682 Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin, Joshua V. Dillon, Balaji
683 Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive
684 uncertainty under dataset shift. In *Conference on Neural Information Processing Systems (NeurIPS)*,
685 2019. doi: 10.48550/arXiv.1906.02530.
- 686 Neal M. Radford. Bayesian learning for neural networks. PhD Thesis, University of Toronto, 1995.
- 687 Kajetan Schweighofer, Lukas Aichberger, Mykyta Ielanskyi, and Sepp Hochreiter. Introducing an im-
688 proved information-theoretic measure of predictive uncertainty. In *Neural Information Processing Systems*
689 *(NeurIPS)*, Mathematics of Modern Machine Learning Workshop, 2023. doi: 10.48550/arXiv.2311.08309.
- 690 Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification
691 uncertainty. In *Conference on Neural Information Processing Systems (NIPS)*, 2018. doi: 10.48550/arXiv.
692 1806.01768.

- 693 Lewis Smith and Yarin Gal. Understanding measures of uncertainty for adversarial example detection. *arXiv*
694 *preprint*, 2018. doi: 10.48550/arXiv.1803.08533.
- 695 Sophie Steger, Christian Knoll, Bernhard Klein, Holger Fröning, and Franz Pernkopf. Function space diver-
696 sity for uncertainty prediction via repulsive last-layer ensembles. In *International Conference on Machine*
697 *Learning (ICML)*, 2024. doi: 10.48550/arXiv.2412.15758.
- 698 Yeming Wen, Dustin Tran, and Jimmy Ba. BatchEnsemble: An alternative approach to efficient ensemble
699 and lifelong learning. In *International Conference on Learning Representations (ICLR)*, 2020. doi: 10.
700 48550/arXiv.2002.06715.
- 701 Lisa Wimmer, Yusuf Sale, Paul Hofman, Bern Bischl, and Eyke Hüllermeier. Quantifying aleatoric and
702 epistemic uncertainty in machine learning: Are conditional entropy and mutual information appropriate
703 measures? In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2023. doi: 10.48550/arXiv.2209.
704 03302.
- 705 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv*, 2017. doi: 10.48550/arXiv.1605.
706 07146.

Supporting Information

Variance-Gated Ensembles: An Epistemic-Aware Framework for Uncertainty Estimation

Table of Contents

711	S1 Symbols and Abbreviations	S3
712	S2 Analytical Derivations and Gradient Expressions	S4
713	S2.1 Gradients through Ensemble Mean	S5
714	S2.2 Gradients through Ensemble Variance	S5
715	S2.3 Jacobians and Vector–Jacobian Products for Shared Gates	S5
716	S2.4 General Reverse-Mode Differentiation for Shared Ensemble Gates	S7
717	S2.5 Gradients with Respect to Learnable Sensitivity Parameters	S7
718	S3 Implementation Details	S9
719	S3.1 Model Architectures	S9
720	S3.2 Training Protocols	S9
721	S3.3 Evaluation and Uncertainty Metrics	S9
722	S3.4 Variance-Gated Normalization Implementation	S10
723	S3.5 Numerical Stability	S10
724	S4 Additional Experimental Results	S11
725	S4.1 Rank Consistency with Existing Measures	S11
726	S4.2 Uncertainty Mass Concentration	S12
727	S4.3 Margin–Variance Geometry	S13
728	S4.4 Effects of Learned \mathbf{k}	S14
729	S4.5 Performance and Calibration	S15
730	S4.6 Computational Efficiency: Full Scaling Results	S17
731	S4.7 Out-of-Distribution Detection	S18
732	S4.8 Sensitivity of VGN to Hyperparameter \mathbf{k}	S20
733	S5 Risk-Based Interpretation	S22
734	S6 Variance-Gated Behavior Across Axioms	S22
735	S7 Axiomatic Justification of Variance-Gated Normalization	S25
736	S7.1 A0: Non-Negativity	S25
737	S7.2 A1: Vanishing Epistemic Uncertainty for Identical Members	S25

738	S7.3 A2: Maximal EU and TU Under Maximally Disagreeing Uniform-Mean Ensembles	S26
739	S7.4 A3: Monotonicity Under Mean-Preserving Spread	S26
740	S7.5 A4: Monotonicity Under Center-Shift (uniform noise addition)	S27
741	S7.6 A5: Invariance of EU Under Spread-Preserving Location Shifts	S28
742	S7.7 Summary	S28

743 **S1 Symbols and Abbreviations**

Table S1: Symbols and abbreviations used in the variance-gated normalization framework.

Symbol	Domain or Type	Description
C	\mathbb{N}	Number of classes.
M	\mathbb{N}	Number of ensemble members.
Δ^{C-1}	simplex	Probability simplex $\{\mathbf{z} \in \mathbb{R}^C \geq 0 : \mathbf{1}^\top \mathbf{z} = 1\}$.
\mathbf{x}	input	Input features.
\mathbf{y}	categorical	Class label over C classes.
\mathbf{w}_m	parameters	Parameters of ensemble member $m \in \{1, \dots, M\}$.
$\mathbf{1}$	\mathbb{R}^C	All-ones vector.
\mathbf{I}	$\mathbb{R}^{C \times C}$	Identity matrix.
\odot	operator	Hadamard (elementwise) product.
$\text{Diag}(\cdot)$	operator	Diagonal matrix with the given vector on the diagonal.
$\log(\cdot)$	function	Base-2 logarithm (applied elementwise to vectors).
\mathbf{p}_m	Δ^{C-1}	Member- m predictive distribution $p(\mathbf{y} \mathbf{x}, \mathbf{w}_m)$.
$p_m(c)$	$[0, 1]$	Probability for class c from member m .
$\bar{\mathbf{z}}$	Δ^{C-1}	Ensemble sample mean (per-class).
\mathbf{S}	$\mathbb{R}_{\geq 0}^C$	Per-class ensemble standard deviation (predictive spread).
ε	$\mathbb{R}_{> 0}$	Small constant for numerical stability (<i>e.g.</i> , 1.0×10^{-8}).
\mathbf{s}	$\mathbb{R}_{> 0}^C$	Numerically stabilized predictive spread.
k_c	$\mathbb{R}_{> 0}$	Classwise sensitivity scalar for gating.
\mathbf{k}	$\mathbb{R}_{> 0}^C$	Per-class sensitivity vector (learnable) for gating.
$\mathbf{\Gamma}$	$[0, 1)^C$	Variance gate with vector \mathbf{k} .
Z_m	$\mathbb{R}_{> 0}$	Normalization constant for member m .
\mathbf{q}_m	Δ^{C-1}	Normalized variance-gated distribution for member m .
$\bar{\mathbf{q}}$	Δ^{C-1}	Variance-gated ensemble mixture (sample mean over members).
H	$[0, 1]$	Normalized Shannon entropy of $\mathbf{z} \in \Delta^{C-1}$.
TU	$[0, 1]$	Total (predictive) uncertainty.
AU	$[0, 1]$	Aleatoric uncertainty (expected per-member entropy).
EU	$[0, 1]$	Epistemic uncertainty (disagreement).
\mathbf{m}_{ij}	Δ^{C-1}	Midpoint distribution for members (i, j) .
$D_{\text{KL}}(\mathbf{z}_i \mathbf{z}_j)$	$\mathbb{R}_{\geq 0}$	Kullback-Leibler divergence.
$D_{\text{JS}}(\mathbf{z}_i \mathbf{z}_j)$	$[0, 1]$	Jensen-Shannon divergence (normalized).
EPCE	$\mathbb{R}_{\geq 0}$	Expected pairwise cross-entropy.
EPKL	$\mathbb{R}_{\geq 0}$	Expected pairwise KL divergence (disagreement).
EPJS	$[0, 1]$	Expected pairwise Jensen-Shannon divergence (disagreement).
\hat{y}	class index or ‘‘abstain’’	Predicted label under the margin rule.
SNR	\mathbb{R}	Signal-to-noise margin between top-2 classes.
γ	$[0, 1]$	Top-2 variance gate with scalar k .
VGMU	$[0, 1]$	Variance-gated margin uncertainty (lower is more confident).
\mathcal{L}	scalar	Training objective (<i>e.g.</i> , NLL on mixture distribution $\bar{\mathbf{q}}$).
\mathbf{u}	\mathbb{R}^C	Upstream gradient at the mixture.
ℓ	\mathbb{R}^C	Unconstrained parameter with $\mathbf{k} = \text{softplus}(\ell)$.
$\text{softplus}(\ell)$	function	Positive reparameterization for \mathbf{k} .
$\sigma(\ell)$	function	Logistic function; appears in $\partial \mathcal{L} / \partial \ell = (\partial \mathcal{L} / \partial \mathbf{k}) \odot \sigma(\ell)$.
b	index	Batch/sample index in sums (<i>e.g.</i> , \sum_b).

744 **S2 Analytical Derivations and Gradient Expressions**

745 This section provides complete analytical derivations of all gradients required for end-to-end training using
 746 variance-gated normalization. The results in this section support the summary expressions presented in
 747 Section 3 of the main paper. In [Table S2](#), we provide a concise reference for the derived gradients.

Table S2: Summary of analytical gradients for the variance-gated normalization framework. Each gradient describes the partial derivative of the loss \mathcal{L} , the gating function Γ , or intermediate statistics $(\bar{\mathbf{p}}, \mathbf{s})$ and learnable parameter \mathbf{k} (*via* ℓ) with respect to the quantities that affect variance-aware ensemble training.

Gradient	Description
$\frac{\partial \bar{\mathbf{q}}}{\partial \mathbf{q}_m} = \frac{1}{M} \mathbf{I}$	Mixture sensitivity to member distribution.
$\frac{\partial \mathcal{L}}{\partial \mathbf{q}_m} = \frac{1}{M} \mathbf{u}, \quad \mathbf{u} = \frac{\partial \mathcal{L}}{\partial \bar{\mathbf{q}}}$	Upstream gradient distributed equally to members.
$\frac{\partial \mathbf{q}_m}{\partial \Gamma} = \frac{1}{Z_m} [\text{Diag}(\mathbf{p}_m) - \mathbf{q}_m \mathbf{p}_m^\top]$	Jacobian \mathbf{J}_m of normalized gating with respect to the gate Γ .
$\frac{\partial \mathcal{L}_m}{\partial \Gamma} = \frac{1}{MZ_m} [\mathbf{p}_m \odot \mathbf{u} - \mathbf{p}_m (\mathbf{q}_m^\top \mathbf{u})]$	Reverse-mode VJP to the gate (per member).
$\frac{\partial \mathcal{L}}{\partial \Gamma} = \frac{1}{M} \sum_{m=1}^M \frac{1}{Z_m} [\mathbf{p}_m \odot \mathbf{u} - \mathbf{p}_m (\mathbf{q}_m^\top \mathbf{u})]$	Total gradient to shared gate across members.
$\frac{\partial \Gamma}{\partial \bar{\mathbf{p}}} = \frac{1 - \Gamma}{\mathbf{k}\mathbf{s}}$	Gate sensitivity to mean confidence.
$\frac{\partial \Gamma}{\partial \mathbf{s}} = -\frac{(1 - \Gamma) \bar{\mathbf{p}}}{\mathbf{k}\mathbf{s}^2}$	Gate sensitivity to predictive spread.
$\frac{\partial \Gamma}{\partial \mathbf{k}} = -\frac{(1 - \Gamma) \bar{\mathbf{p}}}{\mathbf{k}^2 \mathbf{s}}$	Gate sensitivity scale \mathbf{k} .
$\frac{\partial \mathcal{L}}{\partial \bar{\mathbf{p}}} = \frac{\partial \mathcal{L}}{\partial \Gamma} \odot \frac{1 - \Gamma}{\mathbf{k}\mathbf{s}}$	Backpropagation through gate <i>via</i> mean.
$\frac{\partial \mathcal{L}}{\partial \mathbf{s}} = -\frac{\partial \mathcal{L}}{\partial \Gamma} \odot \frac{(1 - \Gamma) \bar{\mathbf{p}}}{\mathbf{k}\mathbf{s}^2}$	Backpropagation through gate <i>via</i> spread.
$\frac{\partial \mathcal{L}}{\partial \mathbf{k}} = -\sum_b \left(\frac{\partial \mathcal{L}}{\partial \Gamma} \odot \frac{(1 - \Gamma) \bar{\mathbf{p}}}{\mathbf{k}^2 \mathbf{s}} \right)_b$	Backpropagation to \mathbf{k} (sum over classes b).
$\frac{\partial \mathcal{L}}{\partial \ell} = \frac{\partial \mathcal{L}}{\partial \mathbf{k}} \odot \sigma(\ell), \quad \mathbf{k} = \text{softplus}(\ell)$	Through softplus reparameterization of \mathbf{k} .
<i>Total per-member gradient contributions:</i>	
$\frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} = \frac{1}{M} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \Big _{\Gamma} + \frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \Big _{\bar{\mathbf{p}}} + \frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \Big _{\mathbf{s}} \right)$	Sums direct and two indirect paths.
$\frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \Big _{\Gamma} = \frac{1}{MZ_m} [\Gamma \odot \mathbf{u} - \Gamma (\mathbf{q}_m^\top \mathbf{u})]$	Direct (local) path through normalization.
$\frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \Big _{\bar{\mathbf{p}}} = \frac{\partial \mathcal{L}}{\partial \bar{\mathbf{p}}} \odot \frac{1}{M}$	Indirect path <i>via</i> ensemble mean.
$\frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \Big _{\mathbf{s}} = \frac{\mathbf{p}_m - \bar{\mathbf{p}}}{M \mathbf{s}}$	Indirect path <i>via</i> spread.

748 S2.1 Gradients through Ensemble Mean

749 This subsection derives gradients that propagate through the ensemble sample mean, which couples the loss
750 to all ensemble members through the shared mixture distribution.

751 **Proposition S2.1** (Through the mean). *Given $\bar{\mathbf{p}} = \frac{1}{M} \sum_{m=1}^M \mathbf{p}_m$, each member receives*

$$\frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \Big|_{\bar{\mathbf{p}}} = \frac{\partial \mathcal{L}}{\partial \bar{\mathbf{p}}} \frac{\partial \bar{\mathbf{p}}}{\partial \mathbf{p}_m} = \frac{\partial \mathcal{L}}{\partial \bar{\mathbf{p}}} \frac{1}{M} \quad (\text{S1})$$

752 *Proof.* Differentiating $\bar{\mathbf{p}}$ gives $d\bar{\mathbf{p}} = \frac{1}{M} d\mathbf{p}_m$; thus the Jacobian $\partial \bar{\mathbf{p}} / \partial \mathbf{p}_m = \frac{1}{M} \mathbf{I}$. \square

753 **Remark S2.1.1.** *Every ensemble member contributes equally to the gradient path through the mean.*

754 S2.2 Gradients through Ensemble Variance

755 This subsection derives gradients that propagate through the ensemble predictive spread, enabling variance-
756 aware updates that emphasize ensemble members deviating from the mean.

757 **Proposition S2.2** (Through the variance). *Let $\mathbf{S} = \sqrt{\frac{1}{M-1} \sum_{m=1}^M (\mathbf{p}_m - \bar{\mathbf{p}})^2}$ and $\mathbf{s} = \mathbf{S} + \varepsilon$ where $\varepsilon =$
758 1.0×10^{-8} (for numerical stability), then*

$$\frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \Big|_{\mathbf{s}} = \frac{\mathbf{p}_m - \bar{\mathbf{p}}}{M \mathbf{S}} \implies \frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \Big|_{\mathbf{s}} \approx \frac{\mathbf{p}_m - \bar{\mathbf{p}}}{M \mathbf{s}}. \quad (\text{S2})$$

759 *Proof.* Define variance as $\mathbf{v} = \frac{1}{M-1} \sum_{m=1}^M (\mathbf{p}_m - \bar{\mathbf{p}})^2$ such that $\mathbf{S} = \sqrt{\mathbf{v}}$. Hence,

$$\frac{\partial \mathbf{S}}{\partial \mathbf{p}_m} = \frac{1}{2\sqrt{\mathbf{v}}} \frac{\partial \mathbf{v}}{\partial \mathbf{p}_m} = \frac{1}{2\mathbf{S}} \frac{\partial \mathbf{v}}{\partial \mathbf{p}_m}. \quad (\text{S3})$$

760 Using the identity $\mathbf{v} = \mathbb{E}[\mathbf{p}_m^2] - \bar{\mathbf{p}}^2$, we have

$$\frac{\partial \mathbf{v}}{\partial \mathbf{p}_m} = \frac{1}{M} 2\mathbf{p}_m - 2\bar{\mathbf{p}} \frac{\partial \bar{\mathbf{p}}}{\partial \mathbf{p}_m}. \quad (\text{S4})$$

761 Recall $\partial \bar{\mathbf{p}} / \partial \mathbf{p}_m = \frac{1}{M} \mathbf{I}$ therefore,

$$\frac{\partial \mathbf{v}}{\partial \mathbf{p}_m} = \frac{2}{M} (\mathbf{p}_m - \bar{\mathbf{p}}) \implies \frac{\partial \mathbf{S}}{\partial \mathbf{p}_m} = \frac{1}{2\mathbf{S}} \frac{2}{M} (\mathbf{p}_m - \bar{\mathbf{p}}) = \frac{\mathbf{p}_m - \bar{\mathbf{p}}}{M \mathbf{S}} \approx \frac{\mathbf{p}_m - \bar{\mathbf{p}}}{M \mathbf{s}} \quad (\text{S5})$$

762 \square

763 **Remark S2.2.1.** *Members farther from the ensemble mean receive proportionally larger gradients, enabling
764 variance awareness in the backward flow of gradients.*

765 S2.3 Jacobians and Vector–Jacobian Products for Shared Gates

766 This subsection derives the local Jacobians and corresponding vector-Jacobian products for the variance-
767 gated normalization layer when the gate is shared across ensemble members.

768 **Proposition S2.3.** *Let the normalized gated member probabilities be $\mathbf{q}_m = (\mathbf{p}_m \odot \Gamma) / Z_m$, $Z_m = \mathbf{p}_m^\top \Gamma$.
769 Then the Jacobian $\mathbf{J}_m = \partial \mathbf{q}_m / \partial \Gamma = [\text{Diag}(\mathbf{p}_m) - \mathbf{q}_m \mathbf{p}_m^\top] / Z_m$.*

770 *Proof.* With \mathbf{p}_m fixed $\partial(\mathbf{p}_m \odot \Gamma) = \mathbf{p}_m \odot \partial \Gamma$ and $\partial(\mathbf{p}_m^\top \Gamma) = \mathbf{p}_m^\top \partial \Gamma$, we differentiate with respect to Γ ,

$$d\mathbf{q}_m = \frac{Z_m \partial(\mathbf{p}_m \odot \Gamma) - (\mathbf{p}_m \odot \Gamma) \partial Z}{Z_m^2} = \frac{Z_m (\mathbf{p}_m \odot \partial \Gamma) - (\mathbf{p}_m \odot \Gamma) \mathbf{p}_m^\top \partial \Gamma}{Z_m^2} = \frac{\mathbf{p}_m \odot \partial \Gamma}{Z_m} - \frac{\mathbf{p}_m \odot \Gamma}{Z_m^2} \odot \mathbf{p}_m^\top \partial \Gamma \quad (\text{S6})$$

771 Substitute $\mathbf{q}_m = (\mathbf{p}_m \odot \Gamma) / \mathbf{p}_m^\top \Gamma$ and $\text{Diag}(\mathbf{p}_m) \odot \partial \Gamma = \mathbf{p}_m \partial \odot \Gamma$,

$$\partial \mathbf{q}_m = \frac{1}{Z_m} \left[\mathbf{p}_m \odot \partial \Gamma - \mathbf{q}_m (\mathbf{p}_m^\top \partial \Gamma) \right] \implies \mathbf{J}_m = \frac{\partial \mathbf{q}_m}{\partial \Gamma} = \frac{1}{Z_m} \left[\text{Diag}(\mathbf{p}_m) - \mathbf{q}_m \mathbf{p}_m^\top \right] \in \mathbb{R}^{C \times C} \quad (\text{S7})$$

772

□

773 **Corollary S2.3.1** (Jacobian-vector product, forward-mode). *For any direction vector $d\Gamma \in \mathbb{R}^C$,*

$$d\mathbf{q}_m = \left(\frac{\partial \mathbf{q}_m}{\partial \Gamma} \right) d\Gamma = \frac{1}{Z_m} \left[\mathbf{p}_m \odot d\Gamma - \mathbf{q}_m (\mathbf{p}_m^\top d\Gamma) \right]. \quad (\text{S8})$$

774 *This calculates the directional derivative $d\mathbf{q}_m$ for $d\Gamma$ in $O(C)$ time.*

775 **Corollary S2.3.2** (Vector-Jacobian product, reverse mode). *Since the loss is evaluated on the ensemble*
 776 *mixture $\bar{\mathbf{q}}$, each ensemble member receives a scaled upstream gradient $\frac{1}{M} \mathbf{u}$, where $\mathbf{u} = \partial \mathcal{L} / \partial \bar{\mathbf{q}}$. The reverse-*
 777 *mode gradient of \mathcal{L}_m with respect to the gate is therefore*

$$\frac{\partial \mathcal{L}_m}{\partial \Gamma} = \frac{1}{M} \left(\frac{\partial \mathbf{q}_m}{\partial \Gamma} \right)^\top \mathbf{u} = \frac{1}{MZ_m} \left[\mathbf{p}_m \odot \mathbf{u} - \mathbf{p}_m (\mathbf{q}_m^\top \mathbf{u}) \right]. \quad (\text{S9})$$

778 *This calculates the per-member reverse-mode gradient $\partial \mathcal{L}_m / \partial \Gamma$ in $O(C)$ time through the variance-gated*
 779 *normalization layer.*

780 **Remark S2.3.1.** *The upstream gradient \mathbf{u} is identical for all ensemble members. Each ensemble member*
 781 *receives a scaled contribution $\frac{1}{M} \mathbf{u}$ when gradients are backpropagated.*

782 **Proposition S2.4.** *Let the normalized gated member probabilities be $\mathbf{q}_m = (\mathbf{p}_m \odot \Gamma) / Z_m$, $Z_m = \mathbf{p}_m^\top \Gamma$.*
 783 *Then the loss with respect to member m probabilities is*

$$\left. \frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \right|_{\Gamma} = \frac{1}{MZ_m} \left[\Gamma \odot \mathbf{u} - \Gamma (\mathbf{q}_m^\top \mathbf{u}) \right]. \quad (\text{S10})$$

784 *Proof.* We define the following, apply the quotient rule on \mathbf{q}_m , and substitute values

$$\mathbf{p}_m \odot \Gamma = Z_m \mathbf{q}_m, \quad d(\mathbf{p}_m \odot \Gamma) = \text{Diag}(\Gamma) d\mathbf{p}_m, \quad dZ_m = \Gamma^\top d\mathbf{p}_m, \quad \Gamma \odot \mathbf{u} = \text{Diag}(\Gamma) \mathbf{u} \quad (\text{S11})$$

$$d\mathbf{q}_m = \frac{Z_m d(\mathbf{p}_m \odot \Gamma) - (\mathbf{p}_m \odot \Gamma) dZ_m}{Z_m^2} \implies d\mathbf{q}_m = \frac{1}{Z_m} \left[\text{Diag}(\Gamma) d\mathbf{p}_m - \mathbf{q}_m (\Gamma^\top d\mathbf{p}_m) \right]. \quad (\text{S12})$$

785 Recall $\partial \mathcal{L} / \partial \mathbf{q}_m = \frac{1}{M} \mathbf{u}$, therefore

$$d\mathcal{L} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{q}_m} \right)^\top d\mathbf{q}_m = \left(\frac{1}{M} \mathbf{u} \right)^\top d\mathbf{q}_m \implies d\mathcal{L} = \frac{1}{MZ_m} \mathbf{u}^\top \left[\text{Diag}(\Gamma) - \mathbf{q}_m \Gamma^\top \right] d\mathbf{p}_m, \quad (\text{S13})$$

786

$$\left. \frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \right|_{\Gamma} = \frac{1}{MZ_m} \left[\text{Diag}(\Gamma) - \Gamma \mathbf{q}_m^\top \right] \mathbf{u} \implies \left. \frac{\partial \mathcal{L}}{\partial \mathbf{p}_m} \right|_{\Gamma} = \frac{1}{MZ_m} \left[\Gamma \odot \mathbf{u} - \Gamma (\mathbf{q}_m^\top \mathbf{u}) \right]. \quad (\text{S14})$$

787

□

788 **Remark S2.4.1.** *There is no sum over members in $\partial \mathcal{L} / \partial \mathbf{p}_m|_{\Gamma}$. Each member contributes independently*
 789 *through its own normalization Z_m .*

790 S2.4 General Reverse-Mode Differentiation for Shared Ensemble Gates

791 This subsection combines the individual gradient paths into a reverse-mode differentiation rule for ensemble
792 layers with shared gating mechanisms.

793 **Proposition S2.5.** *When the ensemble shares a common gate Γ , the total gradient accumulates over all*
794 *members. Let $\mathbf{J}_m = \partial \mathbf{q}_m / \partial \Gamma$ and $\mathbf{u} = \partial \mathcal{L} / \partial \bar{\mathbf{q}}$, where $\bar{\mathbf{q}} = \frac{1}{M} \sum_{m=1}^M \mathbf{q}_m$ is the ensemble mixture distribution.*
795 *Then the total gradient of the loss with respect to Γ is*

$$\frac{\partial \mathcal{L}}{\partial \Gamma} = \frac{1}{M} \sum_{m=1}^M \mathbf{J}_m^\top \mathbf{u} = \frac{1}{M} \sum_{m=1}^M \frac{1}{Z_m} \left[\mathbf{p}_m \odot \mathbf{u} - \mathbf{p}_m (\mathbf{q}_m^\top \mathbf{u}) \right], \quad Z_m = \mathbf{p}_m^\top \Gamma. \quad (\text{S15})$$

796 *Proof.* For a single member $\mathbf{q}_m = (\mathbf{p}_m \odot \Gamma) / Z_m$, the differential is

$$d\mathbf{q}_m = \frac{1}{Z_m} \left[\mathbf{p}_m \odot d\Gamma - \mathbf{q}_m (\mathbf{p}_m^\top d\Gamma) \right]. \quad (\text{S16})$$

797 With the upstream gradient \mathbf{u} , the change in the loss satisfies

$$d\mathcal{L}_m = \frac{1}{M} \mathbf{u}^\top d\mathbf{q}_m = \frac{1}{MZ_m} \left[\mathbf{u}^\top (\mathbf{p}_m \odot d\Gamma) - \mathbf{u}^\top [\mathbf{q}_m (\mathbf{p}_m^\top d\Gamma)] \right] \quad (\text{S17})$$

798 Since $\mathbf{p}_m \odot d\Gamma = \text{Diag}(\mathbf{p}_m) d\Gamma$, $\mathbf{u}^\top (\mathbf{p}_m \odot d\Gamma) = \mathbf{u}^\top \text{Diag}(\mathbf{p}_m) d\Gamma$. Using the product transpose of matrices
799 property $a^\top Bx = (B^\top a)^\top x$, we have $\mathbf{u}^\top \text{Diag}(\mathbf{p}_m) d\Gamma = (\text{Diag}(\mathbf{p}_m)^\top \mathbf{u})^\top d\Gamma$. Given a diagonal matrix is
800 symmetric, $\text{Diag}(\mathbf{p}_m)^\top = \text{Diag}(\mathbf{p}_m) \implies \text{Diag}(\mathbf{p}_m) \mathbf{u} = \mathbf{p}_m \odot \mathbf{u}$. Therefore,

$$d\mathcal{L}_m = \frac{1}{M} \mathbf{u}^\top d\mathbf{q}_m = \frac{1}{MZ_m} \left[(\mathbf{p}_m \odot \mathbf{u})^\top d\Gamma - (\mathbf{q}_m^\top \mathbf{u}) \mathbf{p}_m^\top d\Gamma \right]. \quad (\text{S18})$$

801 By identification with the total differential $d\mathcal{L} = (\partial \mathcal{L} / \partial \Gamma)^\top d\Gamma$, we obtain

$$\frac{\partial \mathcal{L}}{\partial \Gamma} = \frac{1}{M} \sum_{m=1}^M \frac{1}{Z_m} \left[\mathbf{p}_m \odot \mathbf{u} - \mathbf{p}_m (\mathbf{q}_m^\top \mathbf{u}) \right], \quad (\text{S19})$$

802

□

803 **Remark S2.5.1.** *Although the proof begins from the forward differential $d\mathbf{q}_m = (\partial \mathbf{q}_m / \partial \Gamma) d\Gamma$, the final*
804 *expression corresponds to the reverse-mode vector–Jacobian product used in gradient backpropagation, where*
805 *$\mathbf{u} = \partial \mathcal{L} / \partial \bar{\mathbf{q}}$.*

806 **Remark S2.5.2.** *Since the gate Γ is shared across all ensemble members, the total gradient $\partial \mathcal{L} / \partial \Gamma$ accu-*
807 *mulates contributions from each member distribution \mathbf{q}_m . The factor $\frac{1}{M}$ arises from differentiating the loss*
808 *with respect to the mixture distribution $\bar{\mathbf{q}} = \frac{1}{M} \sum_m \mathbf{q}_m$.*

809 S2.5 Gradients with Respect to Learnable Sensitivity Parameters

810 This subsection derives gradients with respect to the learnable sensitivity parameters modulating the strength
811 of variance-gating, including the softplus reparameterization used to enforce positivity.

812 **Proposition S2.6.** *For the variance-gating function $\Gamma = 1 - e^{-\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}}$, the gating scalar $\mathbf{k} > 0$ is learned*
813 *via a softplus reparameterization $\mathbf{k} = \text{softplus}(\boldsymbol{\ell})$, where $\boldsymbol{\ell} \in \mathbb{R}^C$ are unconstrained learnable parameters.*
814 *The gating function $\Gamma \in \mathbb{R}^{B \times C}$ is defined per sample and class, whereas $\mathbf{k}, \mathbf{s} \in \mathbb{R}^C$ are classwise parameters*
815 *shared across the batch. Then the gradients of the loss are*

$$\frac{\partial \mathcal{L}}{\partial \bar{\mathbf{p}}} = \frac{\partial \mathcal{L}}{\partial \Gamma} \frac{1 - \Gamma}{\mathbf{k}\mathbf{s}}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{s}} = -\frac{\partial \mathcal{L}}{\partial \Gamma} \frac{(1 - \Gamma) \bar{\mathbf{p}}}{\mathbf{k}\mathbf{s}^2}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{k}} = -\sum_b \left(\frac{\partial \mathcal{L}}{\partial \Gamma} \frac{(1 - \Gamma) \bar{\mathbf{p}}}{\mathbf{k}^2 \mathbf{s}} \right)_b. \quad (\text{S20})$$

816 *Proof.* Recall the partial derivatives

$$\frac{\partial \Gamma}{\partial \bar{\mathbf{p}}} = \frac{1 - \Gamma}{\mathbf{k}\mathbf{s}}, \quad \frac{\partial \Gamma}{\partial \mathbf{s}} = -\frac{(1 - \Gamma)\bar{\mathbf{p}}}{\mathbf{k}\mathbf{s}^2}, \quad \frac{\partial \Gamma}{\partial \mathbf{k}} = -\frac{(1 - \Gamma)\bar{\mathbf{p}}}{\mathbf{k}^2\mathbf{s}}. \quad (\text{S21})$$

817 Hence, after applying the chain rule, we obtain

$$\frac{\partial \mathcal{L}}{\partial \bar{\mathbf{p}}} = \frac{\partial \mathcal{L}}{\partial \Gamma} \frac{\partial \Gamma}{\partial \bar{\mathbf{p}}}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{s}} = \frac{\partial \mathcal{L}}{\partial \Gamma} \frac{\partial \Gamma}{\partial \mathbf{s}}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{k}} = \sum_b \left(\frac{\partial \mathcal{L}}{\partial \Gamma} \frac{\partial \Gamma}{\partial \mathbf{k}} \right)_b, \quad \frac{\partial \mathcal{L}}{\partial \ell} = \frac{\partial \mathcal{L}}{\partial \mathbf{k}} \frac{\partial \mathbf{k}}{\partial \ell} = \frac{\partial \mathcal{L}}{\partial \mathbf{k}} \sigma(\ell), \quad (\text{S22})$$

818 where $\partial \mathcal{L} / \partial \ell$ is the gradient propagated through softplus reparameterization, which enforces positivity of \mathbf{k}
 819 and modulates the gradient magnitude by the logistic factor $\sigma(\ell) = 1 / (1 + e^{-\ell})$. \square

820 **Remark S2.6.1.** *The partial derivatives reveal a complementary effect of the gating parameters. Increasing*
 821 *$\bar{\mathbf{p}}$ results with a corresponding increase in Γ , reinforcing confident predictions, whereas larger \mathbf{s} or \mathbf{k} reduce*
 822 *Γ , suppressing highly variable predictions. Through this mechanism, the variance-gating function adaptively*
 823 *balances aleatoric and epistemic contributions.*

824 S3 Implementation Details

825 **Reporting and Interpretation** All results are presented as mean \pm standard deviation over three inde-
 826 pendent training runs with different random seeds. These statistics are intended to reflect training variability
 827 rather than to establish formal statistical significance. We interpret differences conservatively, focusing on
 828 (i) consistency of trends across datasets and ensembles, and (ii) differences whose magnitude exceeds the
 829 typical run-to-run variability observed across random seeds.

830 S3.1 Model Architectures

831 **MNIST** For the MNIST dataset, we used a LeNet-5 style convolutional neural network (CNN) as the
 832 base architecture for DE, MCD, LLE, and MCD-LLE models. The network consists of a shared feature
 833 extractor with two convolutional blocks using 5×5 kernels, mapping from 1 to 6 channels in the first block
 834 and from 6 to 16 channels in the second. Each block applies a Tanh activation, followed by 2×2 average
 835 pooling and spatial dropout ($p = 0.1$). The convolutional blocks ends with an adaptive average pooling
 836 layer, producing a fixed 16-dimensional representation that is subsequently flattened. This representation
 837 was passed to a multilayer perceptron (MLP) consisting of a linear layer (16 to 120 units), a Tanh activation,
 838 and a dropout ($p = 0.1$). Classification is performed by a task-specific head implemented as a two-layer MLP
 839 with a $120 \rightarrow 84 \rightarrow 10$ structure, using Tanh activations and dropout between layers. For LLE variants,
 840 the classifier head H is replaced by a list of 100 independent classifier heads, while all preceding layers are
 841 shared.

842 **SVHN, CIFAR-10, and CIFAR-100** For SVHN, CIFAR-10, and CIFAR-100, we employ a ResNet-18
 843 and WideResNet-28-10 (Zagoruyko & Komodakis, 2017) as the base model for DE, MCD, LLE, and MCD-
 844 LLE configurations. The network follows the standard BasicBlock for both the WideResNet networks, with
 845 spatial dropout applied within convolutional blocks. The networks applies batch normalization, ReLU, and
 846 global average pooling. The pooled representation is passed through an additional fully connected block
 847 consisting of a linear layer with batch normalization, ReLU activation, and dropout, prior to classification.
 848 Dropout is applied with probability $p = 0.1$ for SVHN and $p = 0.3$ for CIFAR-10/100. For last-layer ensemble
 849 variants, the classifier module is instantiated as a list of independent classifier heads H , while all preceding
 850 layers are shared, analogous to the MNIST setup.

851 S3.2 Training Protocols

852 Models were trained using the Adam optimizer with learning rate 1.0×10^{-3} for SVHN, CIFAR-10/100 and
 853 1.0×10^{-4} for MNIST with a batch size 128 and the cross-entropy loss objective. For multihead architectures,
 854 losses were averaged across heads prior to backpropagation. Training was done with early stopping based
 855 on the validation loss using a minimum delta of 1.0×10^{-4} and a patience of 5 epochs. Each dataset was
 856 normalized using its respective mean and standard deviation prior to training. Training was performed in
 857 triplicate using different random seeds under fully deterministic settings, including cuDNN and cuBLAS
 858 routines. Experiments were conducted with ensemble sizes of 5, 10, or 100.

859 S3.3 Evaluation and Uncertainty Metrics

860 We compare the proposed variance-gated ensemble framework against standard entropy-based uncertainty
 861 decompositions and recent information-theoretic approaches (Schweighofer et al., 2023). We report both
 862 variance-gated and non-gated variants of entropy- and divergence-based uncertainty measures, together with
 863 the proposed variance-gated margin uncertainty. We compare VGMU against: (i) standard entropy-based
 864 epistemic uncertainty computed as mutual information Houlisby et al. (2011); (ii) Expected Pairwise KL
 865 divergence Schweighofer et al. (2023); and (iii) Expected Pairwise Jensen-Shannon divergence (this work).
 866 For variance-gated variants, we report results using learned per-class \mathbf{k} .

867 We include a diversity measure $\mathbb{E}_{i,c}[\text{Var}_M]$, defined as the ensemble variance averaged across samples and
 868 classes. In addition to uncertainty and diversity metrics, we evaluate predictive performance and calibration
 869 using accuracy, F1-score, and expected calibration error.

870 We assess rank consistency *via* Spearman’s rank correlation ρ and Kendall’s τ . Spearman’s ρ measures
 871 whether samples ranked higher under one measure tend to also be ranked higher under another, capturing
 872 overall agreement between rankings. Kendall’s τ quantifies pairwise ordering agreement by measuring the
 873 fraction of sample pairs whose relative ordering is preserved between two rankings, making it sensitive to local
 874 rank inversions. Together, these metrics provide complementary views of global and fine-grained ranking
 875 stability across uncertainty scores.

876 In addition to rank correlation, we report the cumulative area under the curve (AUCc), which summarizes
 877 how rapidly scores accumulate when samples are ordered from highest to lowest. AUCc is high when a small
 878 number of top-ranked samples account for a large fraction of the total score mass, and low when scores are
 879 more evenly distributed. This provides a complementary perspective by characterizing the concentration
 880 and sharpness of uncertainty estimates beyond ranking agreement alone.

881 Monte Carlo Dropout (i.e, MCD and MCD-LLE) inference sampling (S) was done using a dropout rate of
 882 $p = 0.1$ with $H = 1$ and $S \in \{5, 10, 100\}$ or $H \in \{1, 10\}$ and $S = 10$. For $H > 1$, feature representations
 883 were re-used to improve compute efficiency.

884 S3.4 Variance-Gated Normalization Implementation

885 The VGN module is applied on top of ensemble predictions during training. Given an ensemble of M
 886 per-member probability distributions $\mathbf{P} \in \mathbb{R}^{B \times M \times C}$, the VGN computes outputs \mathbf{Q} *via* the variance-based
 887 gating mechanism. The gated output is computed as $\mathbf{Q} = (\mathbf{P} \odot \mathbf{\Gamma}) / \mathbf{Z}$, where \mathbf{Z} normalizes each member’s
 888 distribution. The final prediction is obtained by averaging the gated distributions: $\bar{\mathbf{Q}} = \frac{1}{M} \sum_{m=1}^M \mathbf{Q}_m$.
 889 Gradients flow through the VGN *via* a custom backward pass that accounts for contributions through the
 890 normalization constant \mathbf{Z} , the mean $\bar{\mathbf{P}}$, and variance \mathbf{s} pathways. The per-class gate parameters \mathbf{k} are re-
 891 parameterized as $\mathbf{k} = \text{softplus}(\ell) + \epsilon$ where ℓ is initialized to 0, resulting with $\mathbf{k} \approx 0.693$ at initialization.
 892 This ensures strictly positive gate parameters and non-saturating initialization. For sensitivity of VGN to
 893 fixed *vs.* learned hyperparameter k , see SI S4.8

894 S3.5 Numerical Stability

895 All operations susceptible to numerical instability include a small constant ϵ (*e.g.*, 1.0×10^{-8}) to prevent
 896 underflow, overflow, and division-by-zero errors. Ensemble standard deviations were stabilized as $\mathbf{s} = \mathbf{S} + \epsilon$.
 897 The variance-gating function $\mathbf{\Gamma} = 1 - e^{-\bar{\mathbf{P}}/\mathbf{ks}}$ is clamped below by ϵ to ensure finite gradients. The per-class
 898 scaling parameter \mathbf{k} is learned using a softplus reparameterization and clamped to a minimum value of
 899 1.0×10^{-3} to prevent gate saturation (*i.e.*, excessively large signal-to-noise ratios) and to maintain stability
 900 as $\mathbf{k} \rightarrow 0$. Normalization constants $Z_m = \mathbf{p}_m^\top \mathbf{\Gamma}$ are likewise lower-bounded by ϵ ; when this bound is active,
 901 gradients are passed through unchanged to avoid disrupting backpropagation.

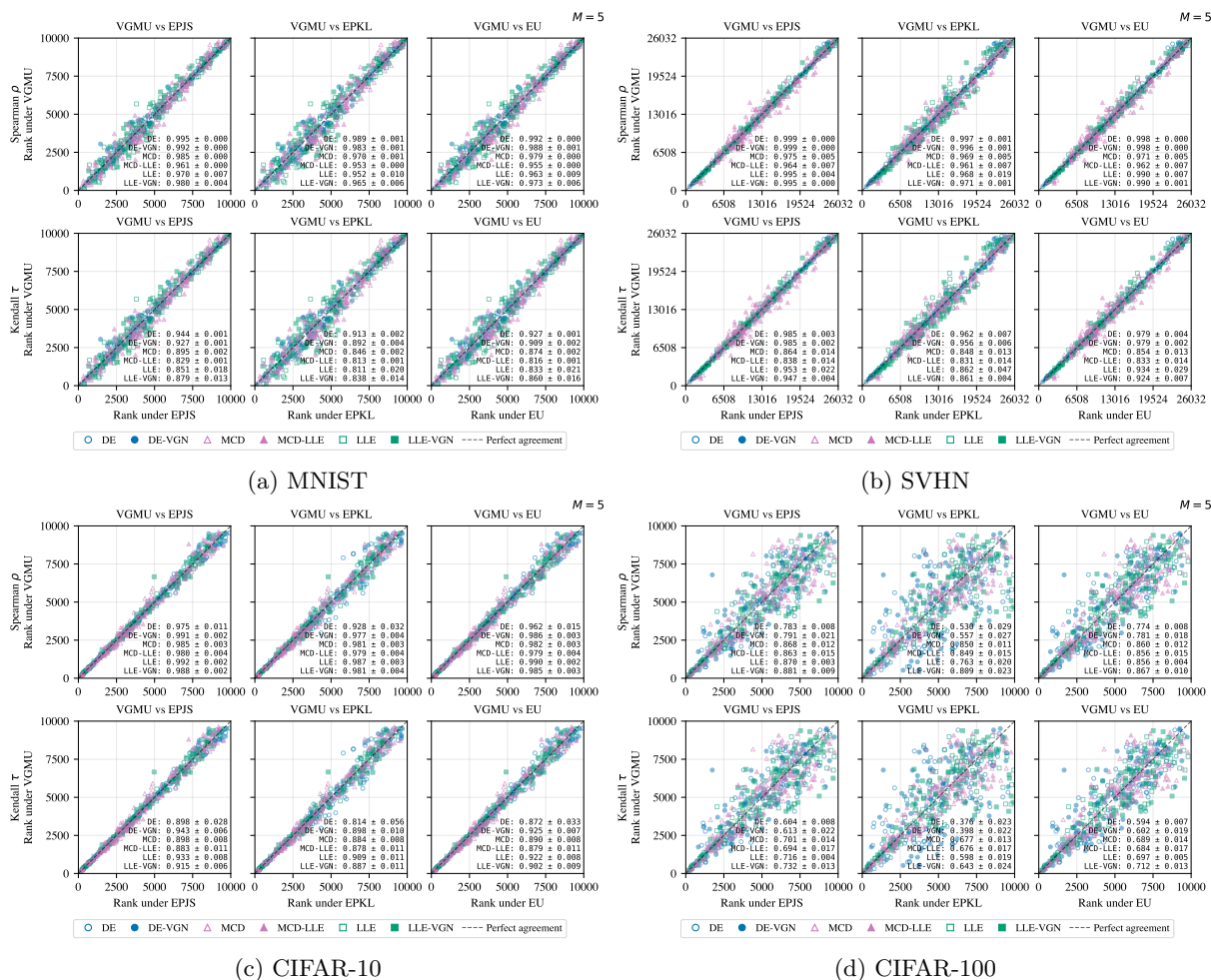
902 **S4 Additional Experimental Results**903 **S4.1 Rank Consistency with Existing Measures**

Figure S1: Spearman's ρ and Kendall's τ rank correlations between VGUM and entropy- and divergence-based epistemic uncertainty measures (EPJS, EPKL, and mutual information). Results are reported across datasets and ensemble configurations, illustrating the degree of alignment between margin-based and full-simplex uncertainty rankings.

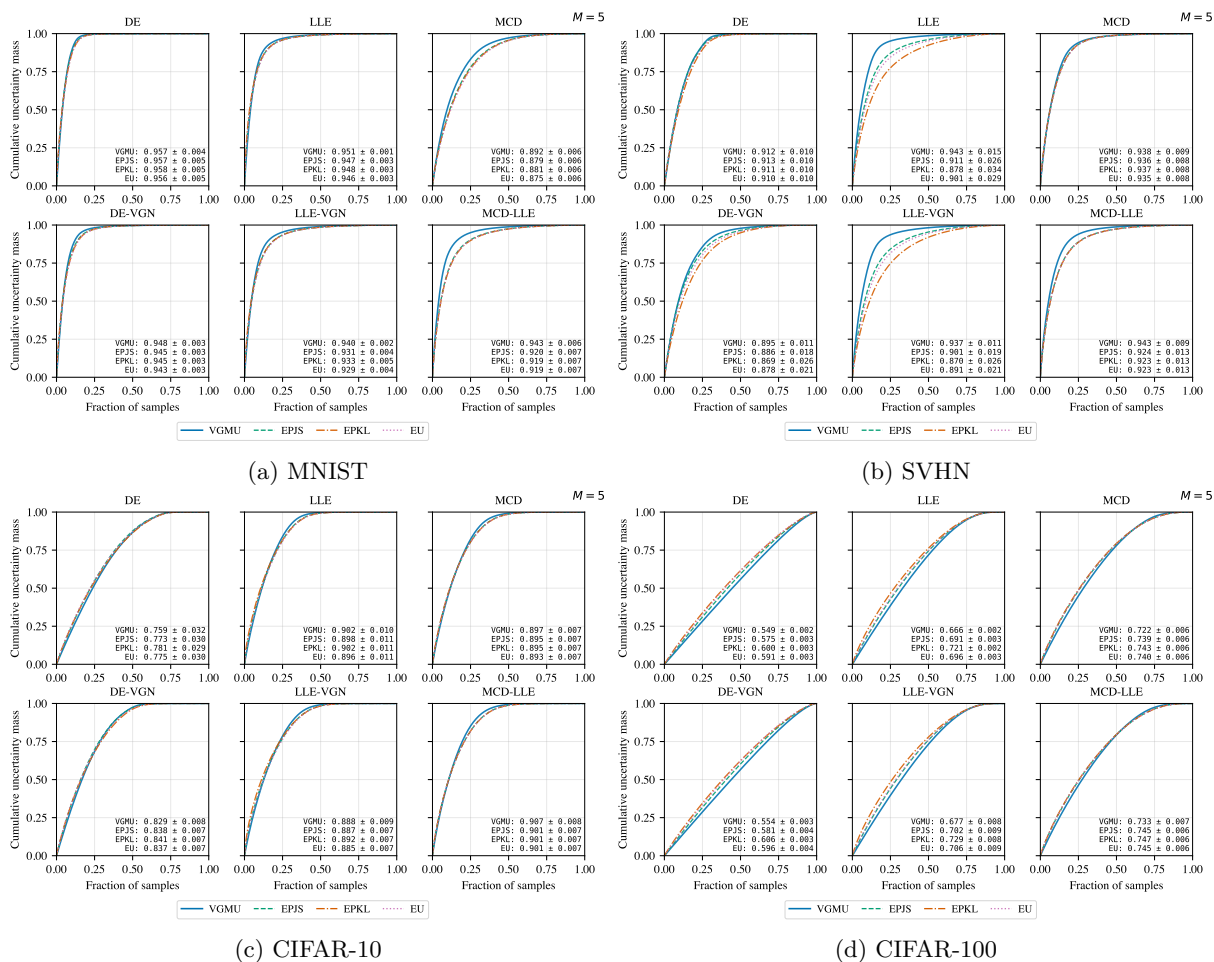
904 **S4.2 Uncertainty Mass Concentration**

Figure S2: Cumulative uncertainty mass concentration curves and corresponding AUC_c values for VGMU and information-theoretic baselines. Samples are sorted by descending uncertainty, and higher AUC_c indicates stronger concentration of uncertainty on difficult samples. Results are shown across datasets and ensemble configurations.

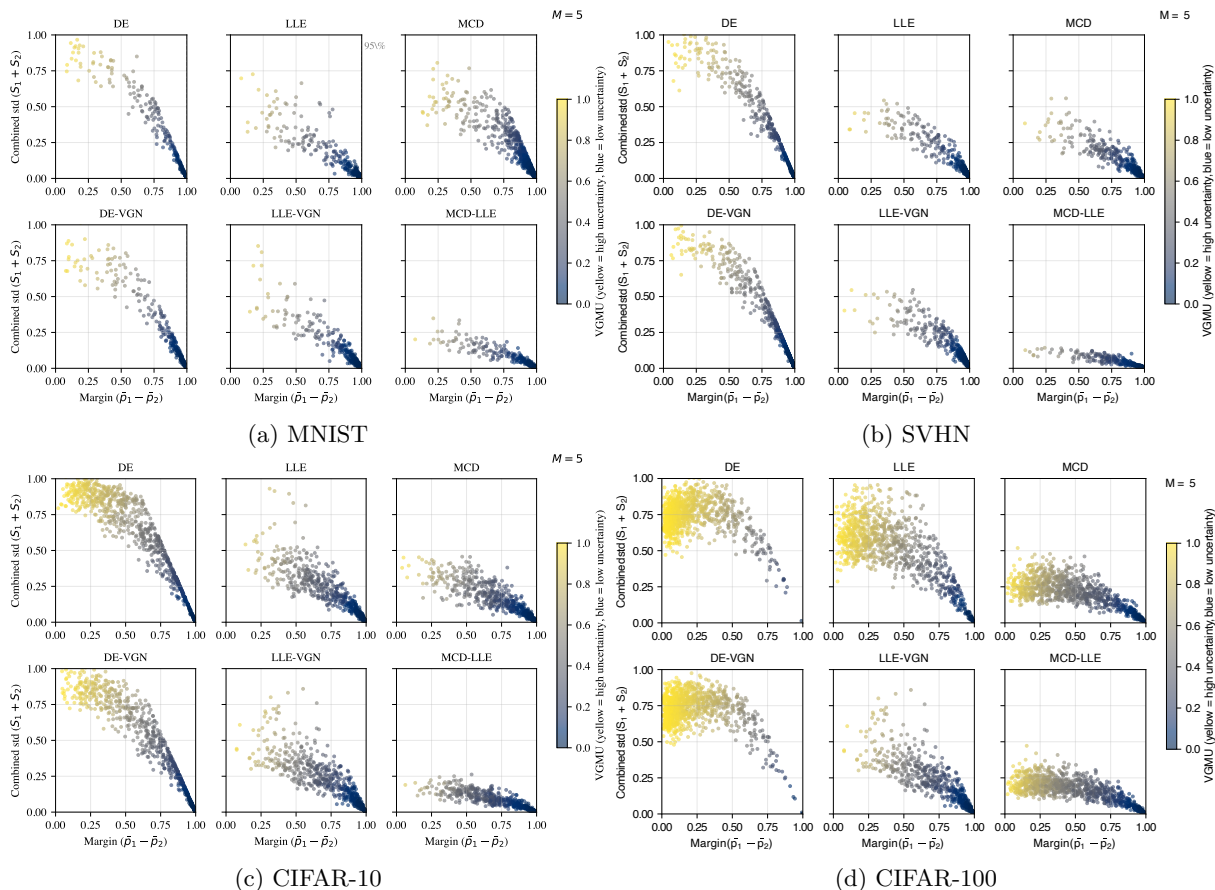
905 **S4.3 Margin-Variance Geometry**

Figure S3: Visualization of predictive margin ($\bar{p}_1 - \bar{p}_2$) and combined predictive spread ($S_1 + S_2$) for different ensemble methods. Each point corresponds to a test sample and is colored by its VG MU value, illustrating how margin-based uncertainty couples class separability with epistemic disagreement.

906 **S4.4 Effects of Learned k**

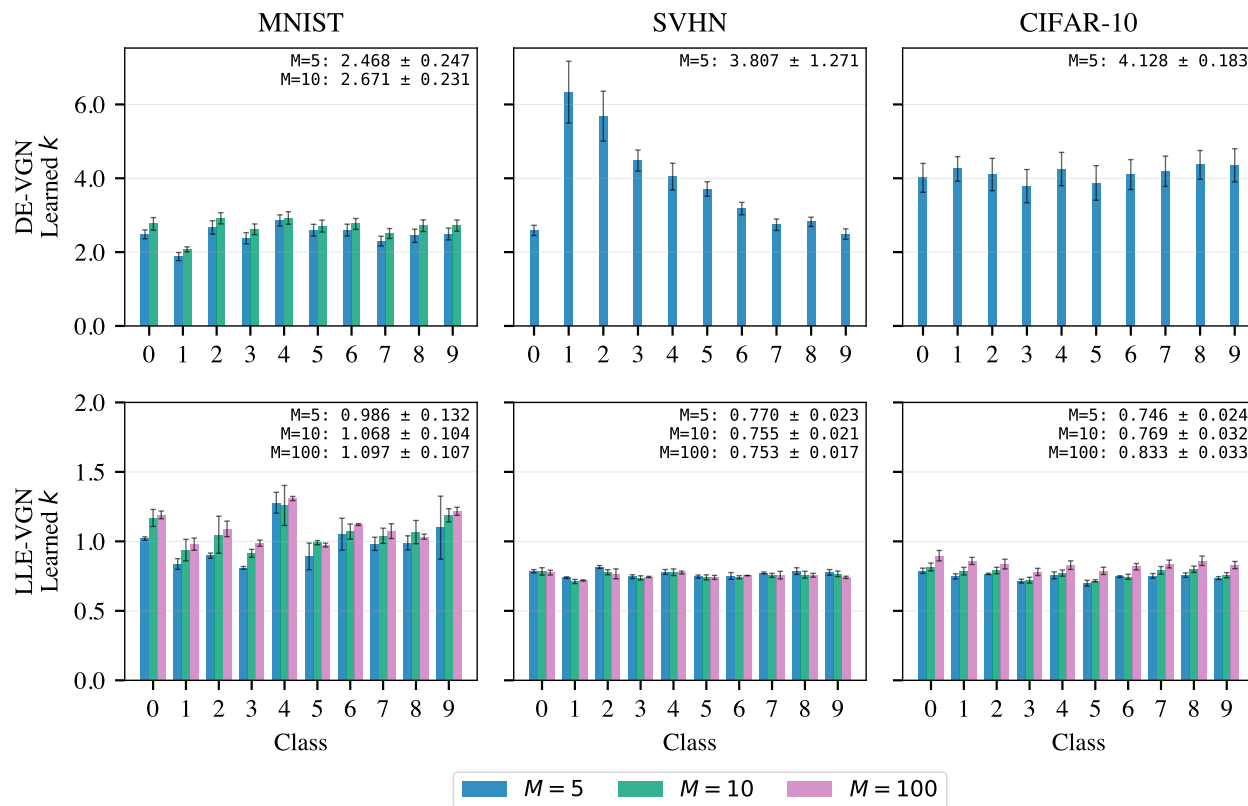


Figure S4: Learned k -values.

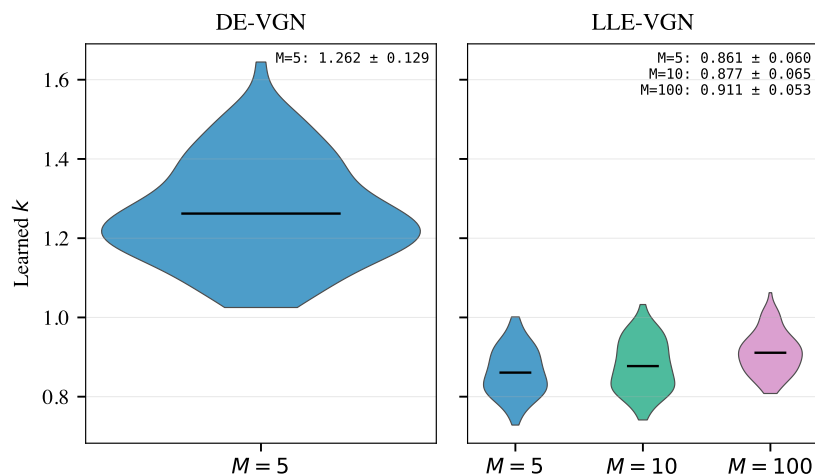


Figure S5: Learned k -values for CIFAR-100.

907 **S4.5 Performance and Calibration**

Table S3: Performance and calibration metrics for ensemble models and the MNIST dataset.

Network	Accuracy	F1-Score	ECE	Diversity
$M = 1$	0.973 ± 0.001	0.973 ± 0.001	0.004 ± 0.001	—
MCD; $M = 5$	0.971 ± 0.001	0.970 ± 0.001	0.040 ± 0.003	$2.7 \times 10^{-3} \pm 0.1$
MCD; $M = 10$	0.971 ± 0.001	0.971 ± 0.001	0.041 ± 0.005	$2.9 \times 10^{-3} \pm 0.1$
MCD; $M = 100$	0.974 ± 0.001	0.973 ± 0.001	0.041 ± 0.004	$3.0 \times 10^{-3} \pm 0.2$
MCD-LLE; $M = 5$	0.973 ± 0.001	0.973 ± 0.001	0.006 ± 0.002	$1.7 \times 10^{-4} \pm 0.0$
MCD-LLE; $M = 10$	0.973 ± 0.001	0.973 ± 0.001	0.006 ± 0.002	$1.9 \times 10^{-4} \pm 0.1$
MCD-LLE; $M = 100$	0.973 ± 0.001	0.973 ± 0.001	0.006 ± 0.002	$2.0 \times 10^{-4} \pm 0.1$
MCD-LLE; $M = 1000$	0.973 ± 0.001	0.973 ± 0.001	0.006 ± 0.002	$2.0 \times 10^{-4} \pm 0.1$
DE; $M = 5$	0.981 ± 0.001	0.981 ± 0.001	0.012 ± 0.002	$2.3 \times 10^{-3} \pm 0.2$
DE-VGN; $M = 5$	0.982 ± 0.001	0.982 ± 0.001	0.007 ± 0.002	$1.4 \times 10^{-3} \pm 0.1$
LLE; $M = 5$	0.971 ± 0.003	0.971 ± 0.003	0.005 ± 0.001	$6.0 \times 10^{-4} \pm 0.2$
LLE; $M = 10$	0.972 ± 0.001	0.971 ± 0.001	0.005 ± 0.001	$7.6 \times 10^{-4} \pm 0.9$
LLE; $M = 100$	0.970 ± 0.001	0.970 ± 0.001	0.005 ± 0.002	$1.0 \times 10^{-3} \pm 0.2$
LLE-VGN; $M = 5$	0.970 ± 0.002	0.970 ± 0.002	0.005 ± 0.001	$6.8 \times 10^{-4} \pm 0.4$
LLE-VGN; $M = 10$	0.971 ± 0.002	0.971 ± 0.002	0.005 ± 0.001	$8.9 \times 10^{-4} \pm 0.1$
LLE-VGN; $M = 100$	0.970 ± 0.003	0.969 ± 0.003	0.005 ± 0.001	$1.1 \times 10^{-3} \pm 0.1$

Table S4: Performance and calibration metrics for ensemble models and the SVHN dataset.

Network	Accuracy	F1-Score	ECE	Diversity
$M = 1$	0.948 ± 0.002	0.943 ± 0.003	0.017 ± 0.004	—
MCD; $M = 5$	0.946 ± 0.003	0.942 ± 0.003	0.010 ± 0.004	$7.3 \times 10^{-4} \pm 0.8$
MCD; $M = 10$	0.947 ± 0.003	0.942 ± 0.003	0.009 ± 0.004	$8.1 \times 10^{-4} \pm 1.0$
MCD; $M = 100$	0.947 ± 0.003	0.943 ± 0.003	0.008 ± 0.003	$9.0 \times 10^{-4} \pm 1.1$
MCD-LLE; $M = 5$	0.947 ± 0.003	0.943 ± 0.003	0.016 ± 0.004	$7.3 \times 10^{-5} \pm 1.6$
MCD-LLE; $M = 10$	0.947 ± 0.003	0.943 ± 0.003	0.016 ± 0.004	$8.2 \times 10^{-5} \pm 1.9$
MCD-LLE; $M = 100$	0.948 ± 0.003	0.943 ± 0.003	0.016 ± 0.004	$8.5 \times 10^{-5} \pm 1.9$
MCD-LLE; $M = 1000$	0.947 ± 0.003	0.943 ± 0.003	0.016 ± 0.004	$8.5 \times 10^{-5} \pm 2.0$
DE; $M = 5$	0.955 ± 0.002	0.952 ± 0.002	0.026 ± 0.006	$5.8 \times 10^{-3} \pm 0.6$
DE-VGN; $M = 5$	0.960 ± 0.000	0.957 ± 0.000	0.016 ± 0.003	$4.0 \times 10^{-3} \pm 0.4$
LLE; $M = 5$	0.950 ± 0.002	0.946 ± 0.002	0.013 ± 0.006	$9.7 \times 10^{-4} \pm 1.0$
LLE; $M = 10$	0.949 ± 0.003	0.945 ± 0.003	0.017 ± 0.003	$1.3 \times 10^{-3} \pm 0.1$
LLE; $M = 100$	0.951 ± 0.004	0.947 ± 0.004	0.012 ± 0.002	$3.2 \times 10^{-3} \pm 0.3$
LLE-VGN; $M = 5$	0.953 ± 0.005	0.949 ± 0.006	0.014 ± 0.002	$8.1 \times 10^{-4} \pm 0.2$
LLE-VGN; $M = 10$	0.952 ± 0.002	0.949 ± 0.002	0.012 ± 0.001	$1.1 \times 10^{-3} \pm 0.0$
LLE-VGN; $M = 100$	0.951 ± 0.004	0.947 ± 0.005	0.015 ± 0.004	$2.5 \times 10^{-3} \pm 0.1$

Table S5: Performance and calibration metrics for ensemble models and the CIFAR-10 dataset.

Network	Accuracy	F1-Score	ECE	Diversity
$M = 1$	0.853 ± 0.008	0.852 ± 0.008	0.070 ± 0.005	—
MCD; $M = 5$	0.852 ± 0.007	0.851 ± 0.007	0.057 ± 0.002	$1.4 \times 10^{-3} \pm 0.1$
MCD; $M = 10$	0.852 ± 0.007	0.851 ± 0.006	0.055 ± 0.003	$1.5 \times 10^{-3} \pm 0.1$
MCD; $M = 100$	0.853 ± 0.008	0.853 ± 0.007	0.053 ± 0.004	$1.7 \times 10^{-3} \pm 0.1$
MCD-LLE; $M = 5$	0.852 ± 0.009	0.852 ± 0.008	0.067 ± 0.004	$3.1 \times 10^{-4} \pm 0.1$
MCD-LLE; $M = 10$	0.853 ± 0.009	0.852 ± 0.008	0.067 ± 0.004	$3.5 \times 10^{-4} \pm 0.2$
MCD-LLE; $M = 100$	0.853 ± 0.008	0.852 ± 0.008	0.067 ± 0.004	$3.6 \times 10^{-4} \pm 0.2$
MCD-LLE; $M = 1000$	0.853 ± 0.008	0.852 ± 0.008	0.067 ± 0.004	$3.7 \times 10^{-4} \pm 0.2$
DE; $M = 5$	0.836 ± 0.012	0.836 ± 0.012	0.049 ± 0.014	$1.9 \times 10^{-2} \pm 0.2$
DE-VGN; $M = 5$	0.875 ± 0.005	0.875 ± 0.006	0.023 ± 0.003	$9.8 \times 10^{-3} \pm 0.5$
LLE; $M = 5$	0.855 ± 0.002	0.855 ± 0.002	0.062 ± 0.012	$2.4 \times 10^{-3} \pm 0.3$
LLE; $M = 10$	0.848 ± 0.003	0.848 ± 0.002	0.058 ± 0.010	$3.6 \times 10^{-3} \pm 0.5$
LLE; $M = 100$	0.851 ± 0.004	0.850 ± 0.005	0.065 ± 0.006	$5.9 \times 10^{-3} \pm 0.8$
LLE-VGN; $M = 5$	0.846 ± 0.008	0.845 ± 0.008	0.067 ± 0.005	$2.3 \times 10^{-3} \pm 0.5$
LLE-VGN; $M = 10$	0.849 ± 0.004	0.849 ± 0.003	0.066 ± 0.003	$3.2 \times 10^{-3} \pm 0.2$
LLE-VGN; $M = 100$	0.847 ± 0.005	0.848 ± 0.004	0.071 ± 0.006	$5.4 \times 10^{-3} \pm 0.8$

Table S6: Performance and calibration metrics for ensemble models and the CIFAR-100 dataset.

Network	Accuracy	F1-Score	ECE	Diversity
$M = 1$	0.565 ± 0.012	0.563 ± 0.012	0.122 ± 0.002	—
MCD; $M = 5$	0.564 ± 0.012	0.562 ± 0.012	0.092 ± 0.001	$3.0 \times 10^{-4} \pm 0.2$
MCD; $M = 10$	0.567 ± 0.012	0.565 ± 0.012	0.086 ± 0.001	$3.4 \times 10^{-4} \pm 0.2$
MCD; $M = 100$	0.568 ± 0.012	0.566 ± 0.012	0.082 ± 0.002	$3.7 \times 10^{-4} \pm 0.2$
MCD-LLE; $M = 5$	0.563 ± 0.013	0.562 ± 0.013	0.106 ± 0.001	$1.9 \times 10^{-4} \pm 0.1$
MCD-LLE; $M = 10$	0.564 ± 0.014	0.563 ± 0.014	0.103 ± 0.003	$2.2 \times 10^{-4} \pm 0.1$
MCD-LLE; $M = 100$	0.565 ± 0.013	0.563 ± 0.013	0.102 ± 0.001	$2.3 \times 10^{-4} \pm 0.1$
MCD-LLE; $M = 1000$	0.565 ± 0.013	0.564 ± 0.013	0.102 ± 0.001	$2.3 \times 10^{-4} \pm 0.1$
DE; $M = 5$	0.487 ± 0.002	0.481 ± 0.003	0.095 ± 0.006	$3.5 \times 10^{-3} \pm 0.0$
DE-VGN; $M = 5$	0.507 ± 0.007	0.504 ± 0.007	0.086 ± 0.006	$3.1 \times 10^{-3} \pm 0.2$
LLE; $M = 5$	0.545 ± 0.002	0.541 ± 0.003	0.074 ± 0.002	$1.6 \times 10^{-3} \pm 0.0$
LLE; $M = 10$	0.543 ± 0.005	0.538 ± 0.005	0.062 ± 0.004	$2.4 \times 10^{-3} \pm 0.1$
LLE; $M = 100$	0.537 ± 0.010	0.533 ± 0.012	0.059 ± 0.019	$4.3 \times 10^{-3} \pm 0.2$
LLE-VGN; $M = 5$	0.548 ± 0.005	0.545 ± 0.007	0.111 ± 0.007	$1.3 \times 10^{-3} \pm 0.1$
LLE-VGN; $M = 10$	0.545 ± 0.014	0.542 ± 0.014	0.095 ± 0.026	$2.0 \times 10^{-3} \pm 0.2$
LLE-VGN; $M = 100$	0.544 ± 0.012	0.542 ± 0.011	0.059 ± 0.009	$4.2 \times 10^{-3} \pm 0.1$

908 **S4.6 Computational Efficiency: Full Scaling Results**Table S7: Wall-clock inference time (μs / sample) and VGMM speedup across all ensemble sizes M and classes C .

M	C	Time (μs / sample)				VGMM Speedup		
		VGMM	EPKL	EPJS	EU	EPKL	EPJS	EU
2	10	0.5	0.6	0.9	0.8	1.1	1.6	1.5
2	100	0.5	0.6	0.9	0.8	1.1	1.7	1.6
2	1000	0.5	0.6	0.9	0.8	1.2	1.8	1.6
5	10	0.5	0.6	0.8	0.8	1.1	1.7	1.6
5	100	0.5	0.6	0.8	0.8	1.1	1.7	1.6
5	1000	0.5	0.6	1.1	0.8	1.3	2.3	1.7
10	10	0.5	0.6	0.9	0.8	1.2	1.7	1.6
10	100	0.5	0.5	0.8	0.8	1.1	1.7	1.6
10	1000	0.5	1.6	6.1	0.8	3.3	13	1.7
20	10	0.5	0.6	0.9	0.8	1.1	1.7	1.6
20	100	0.5	0.8	1.5	0.8	1.5	2.9	1.6
20	1000	0.5	8.0	31	0.8	16	61	1.7
50	10	0.5	0.7	1.3	0.8	1.4	2.5	1.6
50	100	0.5	5.0	19	0.9	9.9	38	1.7
50	1000	0.5	45	186	1.0	87	360	2.0
100	10	0.5	1.8	6.8	0.8	3.4	13	1.6
100	100	0.5	19	75	0.9	38	150	1.7
100	1000	0.5	176	738	2.9	320	1342	5.2

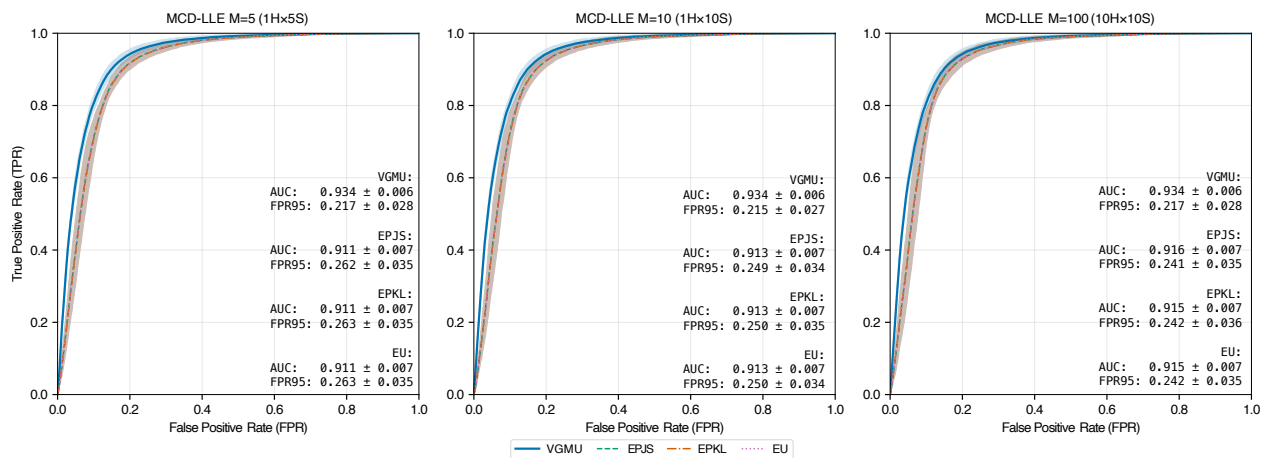
909 **S4.7 Out-of-Distribution Detection**FIX:
pDHd:10

Figure S6: ROC curves for OOD detection (SVHN \rightarrow CIFAR-10) with MCD-LLE under varying head (H) \times sample (S) configurations with increasing stochastic samples: $1H \times 5S$ ($M = 5$), $1H \times 10S$ ($M = 10$), and $10H \times 10S$ ($M = 100$). VGMU curves are nearly identical across all configurations.

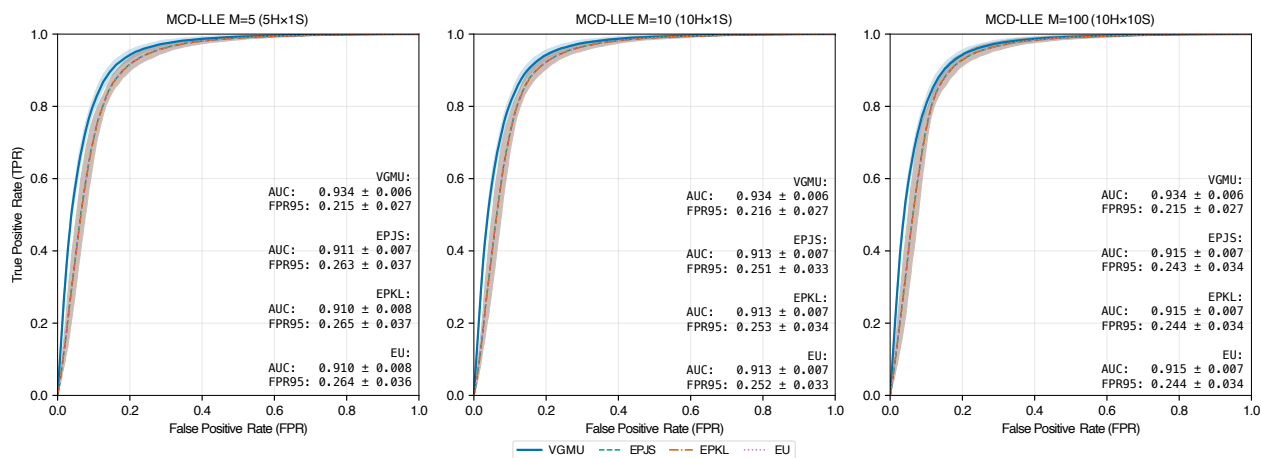


Figure S7: ROC curves for OOD detection (SVHN \rightarrow CIFAR-10) with MCD-LLE under varying head (H) \times sample (S) configurations with increasing classifier heads: $5H \times 1S$ ($M = 5$), $10H \times 1S$ ($M = 10$), and $10H \times 10S$ ($M = 100$). The rightmost panel ($M = 100$) is shared with [Figure S6](#).

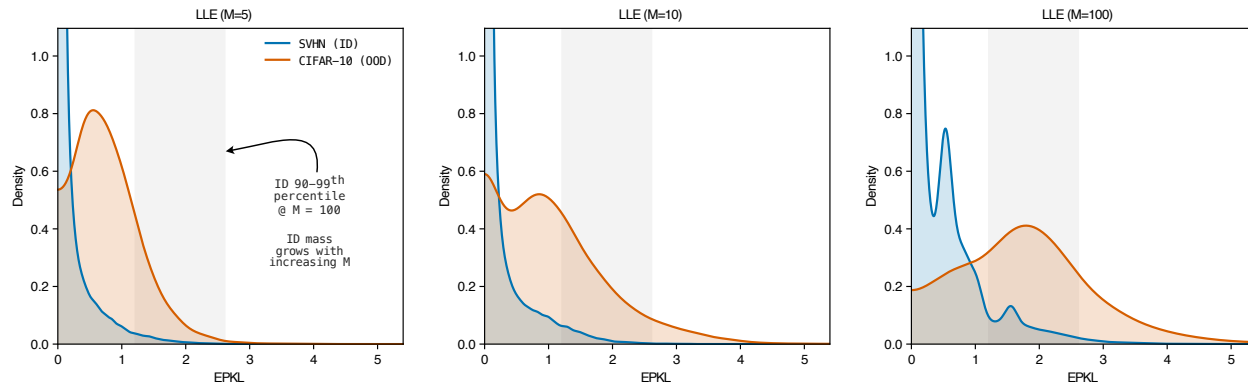


Figure S8: Per-sample EPKL density estimates on SVHN (ID) and CIFAR-10 (OOD) for LLE with $M \in \{5, 10, 100\}$. The shaded band identifies the ID 90–99th percentile range at $M = 100$, fixed across panels. At $M = 5$, the OOD distribution is unimodal at low EPKL with the ID concentrated near zero, resulting in sufficient separation. As M increases, the ID distribution accumulates mass. The ID fraction grows from 1.7% ($M = 5$) to 3.0% ($M = 10$) to 9.0% ($M = 100$), while the OOD fraction grows from 17.9% to 32.5% to 52.2%. At $M = 100$, the ID distribution becomes visibly multimodal. The compression of the OOD-to-ID mass ratio from $10.9\times \rightarrow 5.8\times$, translates into the flat segment of the ROC at $\text{FPR} \in [0.15, 0.30]$ visible in the EPKL curve of Figure 8 (LLE, $M = 100$ panel).

910 S4.8 Sensitivity of VGN to Hyperparameter k

911 In the main text, the per-class parameter k in VGN is learned end-to-end during training. A natural question
 912 is whether this learnable parameterization is necessary, or whether a fixed k suffices. To investigate, we
 913 train LLE-VGN on CIFAR-10 (WideResNet-28-10, $M=10$) with k fixed at three values spanning an order
 914 of magnitude: $k \in \{0.127, 0.693, 4.018\}$, corresponding to $\text{softplus}(\ell) \in \{-2.0, 0.0, 4.0\}$. All other
 915 hyperparameters, including the training protocol and early stopping criteria, are held constant. Results are
 916 averaged over three random seeds.

917 **Classification performance.** Table S8 reports accuracy, F1-score, ECE, and diversity. Classification
 918 accuracy is largely insensitive to the choice of k , with all configurations achieving 0.842–0.849 accuracy. The
 919 learnable- k variant (0.849) and fixed $k = 0.127$ (0.849) perform comparably. ECE varies modestly across
 920 settings; the non-VGN baseline achieves the lowest ECE (0.058), while fixed $k = 4.018$ resulted with the
 highest (0.076).

Table S8: Classification performance for LLE and LLE-VGN on CIFAR-10 (WideResNet-28-10, $M=10$). LLE reports ensemble predictions (P); VGN variants report variance-gated predictions (Q).

Method	k	Accuracy	F1	ECE	Diversity
LLE	–	0.848 ± 0.003	0.848 ± 0.002	0.058 ± 0.010	$3.6 \times 10^{-3} \pm 0.6$
LLE-VGN	0.127	0.849 ± 0.006	0.850 ± 0.005	0.065 ± 0.002	$3.3 \times 10^{-3} \pm 0.4$
LLE-VGN	0.693	0.845 ± 0.005	0.845 ± 0.005	0.061 ± 0.004	$3.6 \times 10^{-3} \pm 0.4$
LLE-VGN	4.018	0.842 ± 0.004	0.843 ± 0.003	0.076 ± 0.009	$3.5 \times 10^{-3} \pm 0.6$
LLE-VGN	learned	0.849 ± 0.004	0.849 ± 0.003	0.066 ± 0.003	$3.2 \times 10^{-3} \pm 0.2$

921

922 **Uncertainty quality.** Table S9 reports AUCc scores, which measure the concentration of uncertainty
 923 mass on difficult samples. Here, the effect of k is more pronounced. Fixed $k = 0.693$ matches the non-VGN
 924 baselines, indicating that this value of k effectively removes any benefits of the variance gate. In contrast,
 925 fixed $k = 0.127$ and $k = 4.018$ both improve over the baseline, with $k = 4.018$ achieving the highest AUCc
 926 across all uncertainty metrics. The learnable- k variants falls between the two best fixed values, indicating
 that it recovers near-optimal performance without requiring a hyperparameter search.

Table S9: Uncertainty mass concentration (AUCc) for LLE and LLE-VGN on CIFAR-10 (WideResNet-28-10, $M=10$). Higher values indicate sharper concentration on difficult samples.

Method	k	VGMU	EU	EPJS	EPKL
LLE	–	0.885 ± 0.015	0.868 ± 0.016	0.873 ± 0.016	0.876 ± 0.015
LLE-VGN	0.127	0.894 ± 0.008	0.879 ± 0.009	0.884 ± 0.008	0.885 ± 0.009
LLE-VGN	0.693	0.885 ± 0.006	0.868 ± 0.007	0.873 ± 0.006	0.876 ± 0.006
LLE-VGN	4.018	0.899 ± 0.012	0.886 ± 0.012	0.891 ± 0.012	0.892 ± 0.010
LLE-VGN	learned	0.897 ± 0.009	0.881 ± 0.010	0.886 ± 0.010	0.886 ± 0.010

927

928 **Rank correlations.** Table S10 and Table S11 report Spearman and Kendall rank correlations between
 929 VGMU and epistemic uncertainty baselines. The pattern mirrors the AUCc results: fixed $k = 0.693$ matches
 930 the baseline, while $k = 4.018$ and learnable k achieve the highest correlations. All VGN variants maintain
 931 strong rank agreement ($\rho > 0.986$, $\tau > 0.907$).

932 **Summary.** The sensitivity analysis reveals that VGN’s benefit depends on the choice of k , with perfor-
 933 mance varying non-monotonically across the tested range. An unfavorable choice (*e.g.*, $k = 0.693$) can

Table S10: Spearman rank correlation (ρ) between VG MU and epistemic uncertainty baselines on CIFAR-10 (WideResNet-28-10, $M = 10$).

Method	k	VG MU vs. EU	VG MU vs. EPJS	VG MU vs. EPKL
LLE	–	0.988 ± 0.004	0.992 ± 0.002	0.985 ± 0.004
LLE-VGN	0.127	0.991 ± 0.002	0.994 ± 0.001	0.988 ± 0.003
LLE-VGN	0.693	0.989 ± 0.001	0.993 ± 0.001	0.986 ± 0.002
LLE-VGN	4.018	0.993 ± 0.002	0.995 ± 0.001	0.990 ± 0.002
LLE-VGN	learned	0.992 ± 0.002	0.995 ± 0.001	0.990 ± 0.002

Table S11: Kendall rank correlation (τ) between VG MU and epistemic uncertainty baselines on CIFAR-10 (WideResNet-28-10, $M = 10$).

Method	k	VG MU vs. EU	VG MU vs. EPJS	VG MU vs. EPKL
LLE	–	0.918 ± 0.011	0.939 ± 0.009	0.905 ± 0.012
LLE-VGN	0.127	0.928 ± 0.007	0.946 ± 0.005	0.917 ± 0.009
LLE-VGN	0.693	0.920 ± 0.003	0.940 ± 0.003	0.907 ± 0.003
LLE-VGN	4.018	0.934 ± 0.007	0.951 ± 0.006	0.924 ± 0.010
LLE-VGN	learned	0.932 ± 0.008	0.948 ± 0.007	0.923 ± 0.010

⁹³⁴ remove the effect of variance gating, while a well-chosen fixed k (*e.g.*, 4.018) can slightly outperform the
⁹³⁵ learnable variant. However, the learnable parameterization consistently achieves near-optimal uncertainty
⁹³⁶ calibration without requiring prior knowledge of the appropriate k value, making it the recommended default
⁹³⁷ for practical use.

938 S5 Risk-Based Interpretation

939 The term \mathbf{ks} represents a classwise tolerance scale describing the dispersion of ensemble member probabili-
 940 ties around the ensemble mean. Under mild distributional assumptions on ensemble dispersion, $\mathbf{k} \in \{1, 2, 3\}$
 941 corresponds to *ca.* 68%, 95%, or 99.7% of members within one, two, or three standard deviations of the
 942 mean, respectively. When ensemble predictions are consistent ($\mathbf{ks} \ll \bar{\mathbf{p}}$), the gate satisfies $\Gamma \approx 1$, leaving
 943 probabilities unchanged ($\mathbf{q}_m \approx \mathbf{p}_m$). When disagreement is significant ($\mathbf{ks} \gg \bar{\mathbf{p}}$), Γ decreases; however, after
 944 normalization, probability mass is redistributed such that high-variance (uncertain) classes are suppressed.
 945 The quadratic decay of sensitivity with respect to \mathbf{s} and \mathbf{k} ensures that further increases in disagreement
 946 produce diminishing changes to the gate. This variance-gating effect is not uniform suppression. The gating
 947 selectively down-weights inconsistent ensemble predictions while preserving those supported by ensemble
 948 agreement. Normalization amplifies the relative impact of smaller gates, ensuring that ensemble disagree-
 949 ment yields conservative, bounded, and well-calibrated predictive distributions. This behavior ensures that
 950 epistemic disagreement translates into conservative distributions, thereby enforcing the consistency and
 951 boundedness properties required by Wimmer’s Axioms A0–A5.

952 S6 Variance-Gated Behavior Across Axioms

953 [Table S12](#) and [Figure S9](#) provide illustrations of variance-gated ensemble behavior under the axiomatic
 954 framework introduced by [Wimmer et al. \(2023\)](#). [Table S12](#) defines ensemble configurations designed to
 955 test each axiom (A2–A5), while [Figure S9](#) visualizes the corresponding simplex geometry and uncertainty
 956 decompositions as the gating parameter k varies.

957 **A2 (Maximal at Uniform Distribution).** When ensemble members span the simplex vertices (P_0),
 958 the variance-gated decomposition correctly assigns maximal epistemic uncertainty ($\text{EU} = 1.0$), as members
 959 are maximally disagreeing despite the uniform ensemble mean. When members collapse to identical uniform
 960 distributions (Q_0), EU vanishes ($\text{EU} = 0.0$) while total and aleatoric uncertainty remain maximal ($\text{TU} = \text{AU}$
 961 $= 1.0$). This configuration is invariant to gating since Q_0 has zero variance. The results confirm that VGN
 962 distinguishes between distributional ambiguity (high AU) and ensemble disagreement (high EU), consistent
 963 with the partial satisfaction noted in [Table 6](#).

964 **A3 (Mean-Preserving Spread).** The transformation from P_0 (identical members with $\boldsymbol{\mu} =$
 965 $[0.70, 0.20, 0.10]$, $\boldsymbol{\sigma} = \mathbf{0}$) to Q_0 (spread members with preserved mean, $\boldsymbol{\sigma} = [0.17, 0.10, 0.10]$) demonstrates
 966 that VGN correctly increases EU when ensemble disagreement grows without changing the predictive mean.
 967 At $k = 0$, EU increases from 0.000 to 0.070. [Figure S9](#) shows that increasing k progressively attenuates this
 968 epistemic signal, with EU decreasing from 0.070 ($k = 0$) to 0.055 ($k = 1$) to 0.045 ($k = 2$), illustrating the
 969 saturation behavior described in [Proposition 3.2](#). Correspondingly, the VGMU increases from 0.300 to 0.412
 970 ($k = 0$), reflecting decision-relevant uncertainty under ensemble disagreement.

971 **A4 (Center-Shift).** Shifting ensemble mass toward the simplex barycenter (from $\boldsymbol{\mu} = [0.70, 0.20, 0.10]$
 972 in P_0 to $\boldsymbol{\mu} = [0.40, 0.35, 0.25]$ in Q_0) while preserving spread increases both TU ($0.730 \rightarrow 0.984$) and AU
 973 ($0.714 \rightarrow 0.971$). EU remains approximately stable across this transformation ($0.016 \rightarrow 0.013$ at $k = 0$),
 974 demonstrating that VGN correctly attributes the increased uncertainty to aleatoric rather than epistemic
 975 sources. The VGMU increases substantially ($0.325 \rightarrow 0.887$), reflecting the reduced class separability near the
 976 barycenter. This behavior persists across k values, confirming that variance-gating preserves the distinction
 977 between location-induced ambiguity and spread-induced disagreement.

978 **A5 (Spread-Preserving Location Shift).** Moving ensembles toward a vertex (from $\boldsymbol{\mu} = [0.70, 0.20, 0.10]$
 979 in P_0 to $\boldsymbol{\mu} = [0.90, 0.05, 0.05]$ in Q_0) while preserving spread decreases TU ($0.730 \rightarrow 0.359$) and AU ($0.714 \rightarrow$
 980 0.314). At $k = 0$, EU shows sensitivity to the location shift ($0.016 \rightarrow 0.045$), which represents a deviation
 981 from strict axiom compliance. However, as k increases, the gap between P and Q epistemic uncertainty
 982 values narrows substantially: from 0.029 at $k = 0$, to 0.015 at $k = 1$, to 0.006 at $k = 2$. This progressive
 983 reduction suggests that variance-gating enforces approximate invariance to location shifts, with stronger
 984 gating (larger k) providing closer adherence to axiom A5. The corresponding VGMU values ($0.325 \rightarrow 0.103$
 985 at $k = 0$) reflect the improved class separability near the vertex.

Table S12: Illustrative ensembles used to evaluate Wimmer’s axiom (A2–A5). Each pair of ensembles P_0 and Q_0 displays a transformation corresponding to a specific axiom. These examples are used to test whether a given uncertainty measure respects the qualitative properties required by each axiom.

Axioms	Ensemble P_0	Ensemble Q_0
A2: Maximal at Uniform Distribution	$P_0 = \begin{bmatrix} 1.00 & 0.00 & 0.00 \\ 0.00 & 1.00 & 0.00 \\ 0.00 & 0.00 & 1.00 \end{bmatrix}$ $\mu = [0.33 \quad 0.33 \quad 0.33]$ $\sigma = [0.58 \quad 0.58 \quad 0.58]$	$Q_0 = \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$ $\mu = [0.33 \quad 0.33 \quad 0.33]$ $\sigma = [0.00 \quad 0.00 \quad 0.00]$
A3: Mean-Preserving Spread	$P_0 = \begin{bmatrix} 0.70 & 0.20 & 0.10 \\ 0.70 & 0.20 & 0.10 \\ 0.70 & 0.20 & 0.10 \end{bmatrix}$ $\mu = [0.70 \quad 0.20 \quad 0.10]$ $\sigma = [0.00 \quad 0.00 \quad 0.00]$	$Q_0 = \begin{bmatrix} 0.90 & 0.10 & 0.00 \\ 0.60 & 0.30 & 0.10 \\ 0.60 & 0.20 & 0.20 \end{bmatrix}$ $\mu = [0.70 \quad 0.20 \quad 0.10]$ $\sigma = [0.17 \quad 0.10 \quad 0.10]$
A4: Center-Shift	$P_0 = \begin{bmatrix} 0.80 & 0.15 & 0.05 \\ 0.70 & 0.20 & 0.10 \\ 0.60 & 0.25 & 0.15 \end{bmatrix}$ $\mu = [0.70 \quad 0.20 \quad 0.10]$ $\sigma = [0.10 \quad 0.05 \quad 0.05]$	$Q_0 = \begin{bmatrix} 0.50 & 0.30 & 0.20 \\ 0.40 & 0.35 & 0.25 \\ 0.30 & 0.40 & 0.30 \end{bmatrix}$ $\mu = [0.40 \quad 0.35 \quad 0.25]$ $\sigma = [0.10 \quad 0.05 \quad 0.05]$
A5: Spread-Preserving Location Shift	$P_0 = \begin{bmatrix} 0.80 & 0.15 & 0.05 \\ 0.70 & 0.20 & 0.10 \\ 0.60 & 0.25 & 0.15 \end{bmatrix}$ $\mu = [0.70 \quad 0.20 \quad 0.10]$ $\sigma = [0.10 \quad 0.05 \quad 0.05]$	$Q_0 = \begin{bmatrix} 1.00 & 0.00 & 0.00 \\ 0.90 & 0.05 & 0.05 \\ 0.80 & 0.10 & 0.10 \end{bmatrix}$ $\mu = [0.90 \quad 0.05 \quad 0.05]$ $\sigma = [0.10 \quad 0.05 \quad 0.05]$

986 These results provide examples of the axiomatic compliance summarized in [Table 6](#) of the main text. The simple
987 visualizations in [Figure S9](#) complement the geometric interpretation in [Section 3.1.1](#), illustrating how the
988 four ensemble simple regions (*i.e.*, confident–certain, ambiguous–certain, confident–uncertain, ambiguous–
989 uncertain) manifest under controlled axiom-testing transformations. The progressive reduction in uncer-
990 tainty measures with increasing k demonstrates the risk-tolerance interpretation discussed in [Section S5](#),
991 where larger k values accommodate greater ensemble disagreement before suppressing predictions. The
992 partial deviations observed, particularly for A5 at low k , reflect the finite-ensemble and exponential gating
993 effects inherent to the VGN formulation, while the convergence with increasing k supports the properties
994 claimed in the main text.

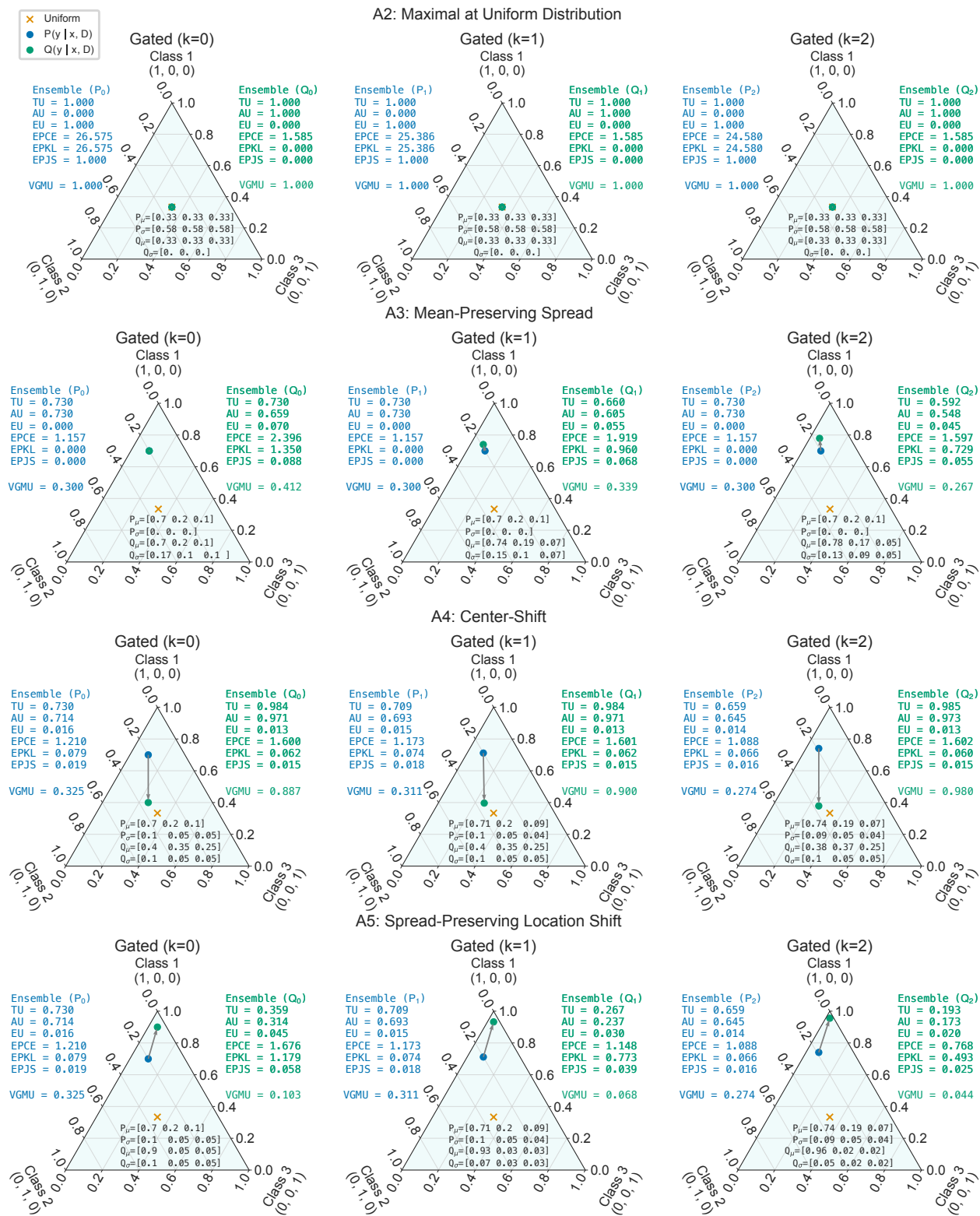


Figure S9: Variance-gated behavior across Wimmer’s axioms (A2–A5). Annotations report TU, AU, and EU uncertainty, as well as divergence-based metrics (EPCE, EPKL, EPJS). Gating progressively reduces epistemic contributions, demonstrating compliance with the axioms under controlled variance attenuation. The orange cross denotes the simplex uniform distribution, while circles indicate distributions before (P , blue) and after (Q , green) gating.

995 S7 Axiomatic Justification of Variance-Gated Normalization

996 This section provides formal proofs and justifications of the axiomatic properties claimed in Table 6 for the
997 VGN framework. We adopt the axiomatic framework introduced by Wimmer et al. (2023) and evaluate each
998 axiom (A0–A5) with respect to the VGN uncertainty decomposition.

999 Recall the relevant definitions from Section 3.3. Given an ensemble of M members producing categorical
1000 distributions $\mathbf{p}_m \in \Delta^{C-1}$, the per-class ensemble sample mean and standard deviation are $\bar{\mathbf{p}}$ and \mathbf{s} , respec-
1001 tively, where \mathbf{s} is stabilized by a small constant $\varepsilon > 0$ to ensure $s_c > 0$ for all c .

1002 The variance gate is $\Gamma = 1 - e^{-\bar{\mathbf{p}}/\mathbf{k}\mathbf{s}}$, the gated member distribution is $\mathbf{q}_m = (\mathbf{p}_m \odot \Gamma)/Z_m$ with $Z_m = \mathbf{p}_m^\top \Gamma$,
1003 and the gated mixture is $\bar{\mathbf{q}} = \frac{1}{M} \sum_{m=1}^M \mathbf{q}_m$. The VGN uncertainty decomposition is

$$\begin{aligned} \text{TU} &:= H(\bar{\mathbf{q}}) = -\bar{\mathbf{q}}^\top \log \bar{\mathbf{q}} \\ \text{AU} &:= \frac{1}{M} \sum_{m=1}^M H(\mathbf{q}_m) \\ \text{EU} &:= \text{TU} - \text{AU} = \frac{1}{M} \sum_{m=1}^M D_{\text{KL}}(\mathbf{q}_m \parallel \bar{\mathbf{q}}). \end{aligned}$$

1004 S7.1 A0: Non-Negativity

1005 **Proposition S7.1** (Non-negativity). *For any ensemble $\{\mathbf{p}_m\}_{m=1}^M$ with $\mathbf{p}_m \in \Delta^{C-1}$ and any $\mathbf{k} > 0$, the
1006 VGN decomposition satisfies $\text{TU} \geq 0$, $\text{AU} \geq 0$, and $\text{EU} \geq 0$.*

1007 *Proof.* We proceed by verifying that each gated distribution \mathbf{q}_m lies in the probability simplex Δ^{C-1} , from
1008 which non-negativity of TU, AU, and EU follows from standard information-theoretic inequalities.

1009 **Step 1: \mathbf{q}_m is a valid probability distribution.** Since $\mathbf{p}_m \in \Delta^{C-1}$, we have $p_m(c) \geq 0$ for all c
1010 and $\sum_c p_m(c) = 1$. The gate satisfies $0 \leq \Gamma_c < 1$ for all c (since the exponential is strictly positive).
1011 Therefore $p_m(c)\Gamma_c \geq 0$ for all c , which gives $q_m(c) = p_m(c)\Gamma_c/Z_m \geq 0$. The normalization constant
1012 $Z_m = \sum_c p_m(c)\Gamma_c > 0$ whenever at least one class c satisfies $p_m(c) > 0$ and $\Gamma_c > 0$, which holds for any
1013 ensemble with $\bar{\mathbf{p}} \neq \mathbf{0}$ (ensured by $\varepsilon > 0$ in \mathbf{s}). By construction, $\sum_c q_m(c) = Z_m/Z_m = 1$. Hence $\mathbf{q}_m \in \Delta^{C-1}$.

1014 **Step 2: $\bar{\mathbf{q}}$ is a valid probability distribution.** As a convex combination of distributions in Δ^{C-1} , the
1015 mixture $\bar{\mathbf{q}} = \frac{1}{M} \sum_m \mathbf{q}_m$ satisfies $\bar{q}(c) \geq 0$ for all c and $\sum_c \bar{q}(c) = 1$.

1016 **Step 3: Non-negativity of TU, AU, EU.** Since $\bar{\mathbf{q}} \in \Delta^{C-1}$, Shannon entropy gives $\text{TU} = H(\bar{\mathbf{q}}) \geq 0$, with
1017 equality if and only if $\bar{\mathbf{q}}$ is a point mass. Similarly, $\text{AU} = \frac{1}{M} \sum_m H(\mathbf{q}_m) \geq 0$ as an average of non-negative
1018 entropies. Finally, since $-\log$ is convex, Jensen’s inequality gives for each m :

$$D_{\text{KL}}(\mathbf{q}_m \parallel \bar{\mathbf{q}}) = -\sum_c q_m(c) \log \frac{\bar{q}(c)}{q_m(c)} \geq -\log \left(\sum_c q_m(c) \frac{\bar{q}(c)}{q_m(c)} \right) = -\log(1) = 0, \quad (\text{S23})$$

1019 with equality if and only if $\mathbf{q}_m = \bar{\mathbf{q}}$. Hence $\text{EU} = \frac{1}{M} \sum_m D_{\text{KL}}(\mathbf{q}_m \parallel \bar{\mathbf{q}}) \geq 0$. \square

1020 S7.2 A1: Vanishing Epistemic Uncertainty for Identical Members

1021 **Proposition S7.2** (Vanishing EU under consensus). *If all ensemble members produce identical predictions,
1022 $\mathbf{p}_m = \mathbf{p}$ for all $m \in \{1, \dots, M\}$, then $\text{EU} = 0$.*

1023 *Proof.* When $\mathbf{p}_m = \mathbf{p}$ for all m , the ensemble statistics are $\bar{\mathbf{p}} = \mathbf{p}$ and $\mathbf{s} = \boldsymbol{\varepsilon}$ (the stabilization vector). The
1024 gate becomes $\Gamma_c = 1 - e^{-p_c/k_c \boldsymbol{\varepsilon}}$ for each class c . Since the gate depends only on shared statistics $(\bar{\mathbf{p}}, \mathbf{s})$, the
1025 vector Γ is identical for all members. Because all members share the same \mathbf{p} and Γ

$$\mathbf{q}_m = \frac{\mathbf{p} \odot \Gamma}{\mathbf{p}^\top \Gamma} = \mathbf{q}, \quad \text{for all } m. \quad (\text{S24})$$

1026 Consequently, $\bar{\mathbf{q}} = \frac{1}{M} \sum_m \mathbf{q} = \mathbf{q}$, and

$$\text{EU} = \frac{1}{M} \sum_{m=1}^M D_{\text{KL}}(\mathbf{q} \| \mathbf{q}) = 0. \quad (\text{S25})$$

1027

□

1028 **Remark S7.2.1.** *This result is independent of \mathbf{k} , ε , the specific form of \mathbf{p} , and the number of ensemble*
 1029 *members M . The gate Γ is a shared function of ensemble statistics; when members agree, it applies an*
 1030 *identical transformation to all members, preserving consensus.*

1031 S7.3 A2: Maximal EU and TU Under Maximally Disagreeing Uniform-Mean Ensembles

1032 **Proposition S7.3** (Partial satisfaction of A2). *Let the ensemble consist of $M = C$ members, each placing*
 1033 *all mass on a distinct class (i.e., $\mathbf{p}_m = \mathbf{e}_m$, the m -th standard basis vector). Then the VGN decomposition*
 1034 *gives $\text{TU} = \log C$ and $\text{EU} = \log C$ (i.e., both are maximal). However, the gate introduces a dependence*
 1035 *on \mathbf{k} for intermediate configurations, which prevents a global maximality guarantee under all uniform-mean*
 1036 *ensembles.*

1037 *Proof.*

1038 **Case 1: Vertex-spanning ensemble ($\mathbf{p}_m = \mathbf{e}_m$).** The ensemble mean is $\bar{p}_c = 1/C$ for all c . The sample
 1039 standard deviation satisfies $s_c > 0$ for all c (since members alternate between 0 and 1 for each class). The
 1040 gate is $\Gamma_c = 1 - e^{-\bar{p}_c/k_c s_c} > 0$ for all c . By symmetry across classes, $\Gamma_c = \Gamma$ for all c .

1041 For member m with $\mathbf{p}_m = \mathbf{e}_m$, only the m -th component is nonzero, so $q_m(c) = p_m(c) \Gamma_c / Z_m$. Since
 1042 $p_m(c) = \delta_{mc}$ and $\Gamma_c = \Gamma$ for all c

$$q_m(c) = \frac{\delta_{mc} \Gamma}{\Gamma} = \delta_{mc}. \quad (\text{S26})$$

1043 Thus $\mathbf{q}_m = \mathbf{e}_m$ (the gated distributions remain one-hot). The mixture is $\bar{q}(c) = 1/C$ for all c . Therefore

$$\text{TU} = H(\bar{\mathbf{q}}) = \log C \quad \text{AU} = \frac{1}{C} \sum_m H(\mathbf{e}_m) = 0 \quad \text{EU} = \log C. \quad (\text{S27})$$

1044 These are the maximum possible values for entropy of a C -class distribution.

1045 **Case 2: Identical uniform members ($\mathbf{p}_m = \mathbf{u} = [1/C, \dots, 1/C]^\top$).** By Proposition S7.2 (A1), $\text{EU} = 0$.
 1046 Meanwhile $\text{TU} = H(\bar{\mathbf{q}})$. Since all members are identical and uniform, $\mathbf{s} = \varepsilon$ and the gate saturates ($\Gamma_c \rightarrow 1$
 1047 as $\varepsilon \rightarrow 0$), resulting in $\mathbf{q}_m \approx \mathbf{u}$ and $\text{TU} \approx \log C$. Thus TU is maximal but EU is zero.

1048 A2 requires that EU and TU are maximal when the ensemble mean is uniform. Case 1 satisfies this when
 1049 members are maximally spread (vertex-spanning). However, for intermediate configurations where the mean
 1050 is uniform but members are not at the vertices, the gate parameter \mathbf{k} modulates the degree of suppression.
 1051 In particular, larger k values reduce the effective gate, attenuating disagreement signals and potentially
 1052 decreasing EU below its non-gated maximum. This \mathbf{k} -dependence prevents a maximality guarantee across
 1053 all uniform-mean ensemble configurations, resulting in partial satisfaction. □

1054 S7.4 A3: Monotonicity Under Mean-Preserving Spread

1055 **Proposition S7.4** (EU increases with mean-preserving spread). *Consider an ensemble $P = \{\mathbf{p}_m\}$ with*
 1056 *non-identical members, mean $\bar{\mathbf{p}}$, and deviations $\mathbf{d}_m = \mathbf{p}_m - \bar{\mathbf{p}}$. Let $P' = \{\bar{\mathbf{p}} + \alpha \mathbf{d}_m\}$ for $\alpha > 1$ such that all*
 1057 *members remain in Δ^{C-1} . Then $\text{EU}(P') > \text{EU}(P)$.*

1058 *Proof.*

1059 **Case 1: $s_c = 0$ for all c .** Let P have identical members, so $\text{EU}(P) = 0$ by A1. Since P' has strictly greater
 1060 variance on at least one class, P' contains non-identical members. It remains to show that distinct \mathbf{p}'_m yield
 1061 distinct gated distributions \mathbf{q}'_m . For every class with $\bar{p}_c > 0$, the gate satisfies $\Gamma'_c = 1 - e^{-\bar{p}_c/k_c s'_c} > 0$;

1062 classes with $\bar{p}_c = 0$ have $p'_m(c) = 0$ for all m (non-negative terms summing to zero) and do not affect
 1063 the gated distributions. Since $q'_m(c) = p'_m(c) \Gamma'_c / Z'_m$ with $\Gamma'_c > 0$ for all classes with $\bar{p}_c > 0$, the map
 1064 $\mathbf{p}'_m \mapsto \mathbf{q}'_m$ is an elementwise rescaling by a shared positive vector followed by normalization. Such a map
 1065 preserves distinctness, where the ratios between entries of distinct probability vectors are unchanged by
 1066 positive rescaling, so $\mathbf{p}'_m \neq \mathbf{p}'_n \implies \mathbf{q}'_m \neq \mathbf{q}'_n$. Therefore

$$\text{EU}(P') = \frac{1}{M} \sum_m D_{\text{KL}}(\mathbf{q}'_m \parallel \bar{\mathbf{q}}') > 0 = \text{EU}(P). \quad (\text{S28})$$

1067 **Case 2: $s_c > 0$ in reference ensemble P .** Since $q_m(c) = p_m(c) \Gamma_c / Z_m$ and $\bar{q}(c) = \frac{1}{M} \sum_m q_m(c)$, the KL
 1068 divergence in EU depends on the log-ratios $\log(q_m(c) / \bar{q}(c))$. A key structural property is that Γ_c cancels
 1069 from every log-ratio, since $\bar{q}(c) = \Gamma_c \cdot \frac{1}{M} \sum_n p_n(c) / Z_n$, so

$$\log \frac{q_m(c)}{\bar{q}(c)} = \log \frac{p_m(c) / Z_m}{\frac{1}{M} \sum_n p_n(c) / Z_n}. \quad (\text{S29})$$

1070 Scaling \mathbf{d}_m by $\alpha > 1$ increases the spread of $\{p'_m(c)\}_m$ around the fixed mean \bar{p}_c , which directly amplifies
 1071 these log-ratios. Meanwhile, the reduced gate $\mathbf{\Gamma}'$ suppresses high-variance classes, partially counteracting
 1072 this effect. However, since Γ_c cancels from the log-ratios, the spread of member values directly determines
 1073 the KL divergences, while the gate suppression affects EU only indirectly through $q_m(c)$ and normalization
 1074 constants Z_m . As a result, $\text{EU}(P') > \text{EU}(P)$. \square

1075 **Remark S7.4.1.** Note that A3 as stated by [Wimmer et al. \(2023\)](#) also requires TU monotonicity under
 1076 mean-preserving spread. For VGN, increased spread reduces gate values, which can decrease AU by more
 1077 than EU increases, resulting in lower TU (e.g., [Figure S9](#) shows TU decreasing from 0.730 to 0.592 at $k = 2$).
 1078 VGN therefore satisfies the EU component of A3 but not the TU component.

1079 S7.5 A4: Monotonicity Under Center-Shift (uniform noise addition)

1080 **Proposition S7.5** (AU and TU increase under center-shift). Consider a transformation that shifts the
 1081 ensemble mean $\bar{\mathbf{p}}$ toward the barycenter $\mathbf{u} = [1/C, \dots, 1/C]^\top$ while preserving the per-class standard deviation
 1082 \mathbf{s} . Then AU and TU both increase.

1083 *Proof.* Let $P = \{\mathbf{p}_m\}$ be the original ensemble and $P' = \{\mathbf{p}'_m\}$ the shifted ensemble, with $\bar{\mathbf{p}}' = (1 - \alpha)\bar{\mathbf{p}} + \alpha\mathbf{u}$
 1084 for some $\alpha \in (0, 1]$ and $\mathbf{s}' = \mathbf{s}$.

1085 **Effect on the gate.** Since $\mathbf{s}' = \mathbf{s}$ and the mean shifts toward the barycenter: for classes c where $\bar{p}_c > 1/C$
 1086 (high-probability classes), $\bar{p}'_c < \bar{p}_c$, so $\Gamma'_c < \Gamma_c$; for classes c where $\bar{p}_c < 1/C$ (low-probability classes), $\bar{p}'_c > \bar{p}_c$,
 1087 so $\Gamma'_c > \Gamma_c$. The gate values become more uniform across classes, suppressing all classes more equally.

1088 **AU increases.** When the gate is class-independent (i.e., $\Gamma_c = \Gamma$ for all c), it cancels in the normalization
 1089 and $\mathbf{q}_m = \mathbf{p}_m$ exactly, since $q_m(c) = p_m(c) \Gamma / (\sum_j p_m(j) \Gamma) = p_m(c)$. In the limit $\bar{\mathbf{p}}' \rightarrow \mathbf{u}$, the gate
 1090 satisfies $\Gamma'_c \rightarrow 1 - e^{-1/Ck_c s_c}$ uniformly across classes and $\mathbf{q}'_m \rightarrow \mathbf{p}'_m$. Since each \mathbf{p}'_m is closer to the
 1091 barycenter than \mathbf{p}_m and has higher entropy, the more uniform gate preserves this higher entropy in \mathbf{q}'_m ,
 1092 giving $\text{AU}(P') = \frac{1}{M} \sum_m H(\mathbf{q}'_m) \geq \text{AU}(P)$.

1093 **TU increases.** When the gate is class-independent, $\bar{\mathbf{q}} = \bar{\mathbf{p}}$ exactly, the uniform gate cancels during normal-
 1094 ization and the mixture mean is unaffected. As center-shift makes the gate more uniform, $\bar{\mathbf{q}}'$ approximates
 1095 $\bar{\mathbf{p}}'$ more closely than $\bar{\mathbf{q}}$ approximates $\bar{\mathbf{p}}$. Since $\bar{\mathbf{p}}' = (1 - \alpha)\bar{\mathbf{p}} + \alpha\mathbf{u}$ is closer to the barycenter than $\bar{\mathbf{p}}$, the
 1096 concavity of Shannon entropy gives

$$H(\bar{\mathbf{p}}') \geq (1 - \alpha) H(\bar{\mathbf{p}}) + \alpha \log C \geq H(\bar{\mathbf{p}}), \quad (\text{S30})$$

1097 with strict inequality when $\bar{\mathbf{p}} \neq \mathbf{u}$. Since the gate becomes more uniform under center-shift, $\bar{\mathbf{q}}'$ closely
 1098 approximates $\bar{\mathbf{p}}'$, giving $\text{TU}(P') = H(\bar{\mathbf{q}}') \geq H(\bar{\mathbf{q}}) = \text{TU}(P)$. \square

1099 **Remark S7.5.1.** *This result distinguishes VGN from standard entropy decomposition, which violates A4.*
 1100 *The key mechanism is that the variance gate depends on $\bar{\mathbf{p}}$, so it responds to center-shifts by suppressing*
 1101 *all classes more equally. This allows the increased distributional ambiguity to manifest in both AU and TU.*
 1102 *In standard entropy decomposition, the mixture $\bar{\mathbf{p}}$ is more uniform than the individual members, so $H(\bar{\mathbf{p}})$*
 1103 *may not increase under center-shift even as the uncertainty of individual members increases, causing TU to*
 1104 *underestimate the increased aleatoric ambiguity.*

1105 S7.6 A5: Invariance of EU Under Spread-Preserving Location Shifts

1106 **Proposition S7.6** (Approximate EU invariance under spread-preserving location shifts). *Consider a trans-*
 1107 *formation that shifts the ensemble mean $\bar{\mathbf{p}}$ to $\bar{\mathbf{p}}'$ while preserving the per-class standard deviation $\mathbf{s}' = \mathbf{s}$. The*
 1108 *VGN epistemic uncertainty satisfies $\text{EU}(P') = \text{EU}(P)$ when the relative disagreement structure is preserved,*
 1109 *and achieves controlled approximate invariance in the general case via the learnable parameter \mathbf{k} .*

1110 *Proof. Exact invariance under proportional shifts.* Consider the special case where the location shift
 1111 preserves the per-class signal-to-noise ratio, *i.e.*, $\bar{p}'_c/s'_c = \bar{p}_c/s_c$ for all c . This occurs, for example, when
 1112 both $\bar{\mathbf{p}}$ and \mathbf{s} are scaled by the same factor, but since $\mathbf{s}' = \mathbf{s}$, this case is trivial and corresponds to no shift.
 1113 In the general case, a location shift with preserved \mathbf{s} changes $\bar{p}_c/(k_c s_c)$ and therefore Γ_c , so exact invariance
 1114 does not hold.

1115 **Approximate invariance mechanism.** Despite the absence of exact invariance, VGN achieves controlled
 1116 approximate invariance through two mechanisms:

1117 (i) *Saturation of the gate.* For classes with $\bar{p}_c/(k_c s_c) \gg 1$, the gate satisfies $\Gamma_c \approx 1$ regardless of the specific
 1118 value of \bar{p}_c . In this region, the gated distributions $\mathbf{q}_m \approx \mathbf{p}_m/(\mathbf{p}_m^\top \mathbf{1}) = \mathbf{p}_m$, and the decomposition reduces
 1119 to the ungated case. Since the ungated entropy decomposition satisfies A5 exactly (Table 6), VGN inherits
 1120 this invariance in the saturated region.

1121 (ii) *Controlled sensitivity via \mathbf{k} .* The derivative of Γ_c with respect to \bar{p}_c is $\partial\Gamma_c/\partial\bar{p}_c = (1 - \Gamma_c)/(k_c s_c)$. As k_c
 1122 increases, this sensitivity decreases, making the gate less responsive to location shifts. In the limit $k_c \rightarrow \infty$,
 1123 $\Gamma_c \rightarrow 0$ uniformly for all classes, and $\mathbf{q}_m \rightarrow \mathbf{p}_m$ after renormalization with a gate that is approximately
 1124 proportional across classes.

1125 For any $\eta > 0$, there exists a k^* such that for all $k_c \geq k^*$

$$|\text{EU}(P') - \text{EU}(P)| < \eta. \quad (\text{S31})$$

1126 This follows from the continuity of the KL divergence and the uniform convergence of Γ as \mathbf{k} increases.

1127 **Empirical confirmation.** The illustrative ensembles in Table S12 confirm this behavior. For the A5 test
 1128 configuration, the EU gap between P_0 and Q_0 decreases from 0.029 at $k = 0$ to 0.006 at $k = 2$ (Section S6),
 1129 demonstrating progressive convergence toward exact invariance with increasing k . \square

1130 **Remark S7.6.1.** Recall A2, where no choice of \mathbf{k} can guarantee maximality across all uniform-mean
 1131 configurations, A5 admits a formal convergence guarantee. For any $\eta > 0$, there exists k^* such that
 1132 $|\text{EU}(P') - \text{EU}(P)| < \eta$ for all $k_c \geq k^*$. In addition, in the saturated gate region, VGN recovers the standard
 1133 entropy decomposition, which satisfies A5 exactly. The full satisfaction mark (\checkmark) for A5 in Table 6 reflects
 1134 this distinction. VGN contains an A5-satisfying decomposition as a limiting case, with a \mathbf{k} -controllable bound
 1135 on the deviation from exact invariance. In practice, the learned values of \mathbf{k} (Section S4.4) place the gate
 1136 in a region where approximate invariance holds, as confirmed by the progressive reduction in EU gap with
 1137 increasing k shown above.

1138 S7.7 Summary

1139 Table S13 summarizes the formal status of each axiom. Axioms A0, A1, A4, and A5 are satisfied (A5 with
 1140 controlled approximation *via* \mathbf{k}). A2 and A3 are partially satisfied. A2 achieves maximality for the vertex-
 1141 spanning ensemble but is \mathbf{k} -dependent for intermediate configurations; A3 satisfies EU monotonicity but TU

1142 can decrease due to gate suppression of high-variance classes. These results are consistent with the claims
 1143 in [Table 6](#) and the empirical illustrations in [Section S6](#).

Table S13: Summary of axiomatic proofs for variance-gated normalization.

Axiom	Status	Proof Basis
A0: Non-negativity	✓	Entropy and KL divergence non-negativity
A1: Vanishing EU (identical members)	✓	Shared gate $\Rightarrow \mathbf{q}_m = \mathbf{q} \Rightarrow D_{\text{KL}} = 0$
A2: Maximal at uniform	●	Exact for vertex-spanning; \mathbf{k} -dependent otherwise
A3: Mean-preserving spread	●	EU proven; TU limited by gate suppression
A4: Center-shift	✓	Gate becomes uniform; concavity of H
A5: Location-shift invariance	✓	Gate saturation and \mathbf{k} -controlled convergence