

Energy-Guided Ghost-Cluster Bayesian Semantic Entropy for Hallucination Detection

Anonymous ACL submission

Abstract

Semantic Entropy (SE) is a robust metric for uncertainty quantification but suffers from high computational costs and an inability to detect “consistent hallucinations”, where a model confidently repeats the same incorrect answer. This failure stems from SE’s closed-set assumption, which neglects plausible semantic alternatives outside the sampled outputs. To address these limitations, we propose **GLIB-SE** (Ghost-Logit Integrated Bayesian Semantic Entropy). This Bayesian framework models the semantic distribution with an explicit “Ghost Cluster” for unobserved semantics. Crucially, we utilize the model’s raw logit energy as a dynamic prior to calibrate this cluster, allowing GLIB-SE to expose hidden uncertainty when a model is internally confused despite generating consistent outputs. Furthermore, we derive an adaptive sampling strategy based on posterior entropy variance to optimize the inference budget. Experiments across six benchmarks demonstrate that GLIB-SE significantly outperforms baselines in hallucination detection (AUROC) while reducing sampling costs by over 30% compared to fixed-sample strategies.

1 Introduction

Large Language Models (LLMs) have become foundational to modern AI, excelling in tasks from open-domain question answering to complex reasoning. Despite these capabilities, their reliable deployment in high-stakes domains is severely hindered by hallucinations producing fluent but factually incorrect content (Ji et al., 2023; Zhang et al., 2023). To mitigate this risk, Uncertainty Quantification (UQ) has emerged as a critical defense mechanism, enabling systems to detect errors, trigger abstention, or solicit human intervention (Huang et al., 2023). Ideally, a robust UQ metric should serve as a faithful proxy for the model’s knowledge boundary: exhibiting high uncertainty precisely

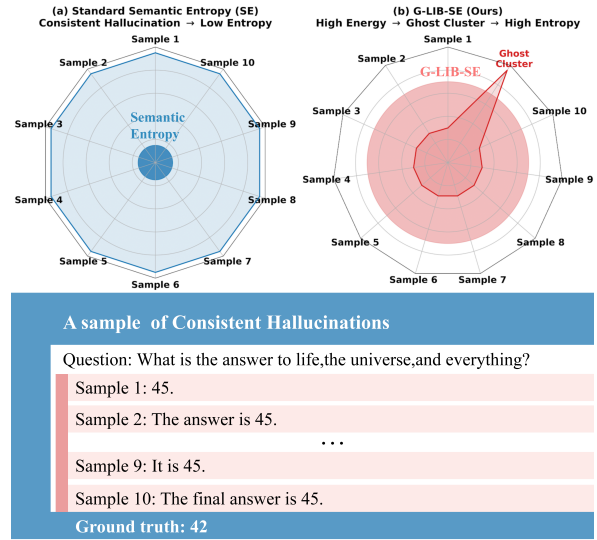


Figure 1: For the same query, the model does not know the correct answer yet produces several semantically consistent wrong responses. Standard Semantic Entropy cannot distinguish this regime, and GLIB-SE uses the **Ghost Cluster** to explicitly flag these hallucinations.

when the model is “confidently wrong” or lacks the necessary information (Kadavath et al., 2022).

Among existing UQ methodologies, Semantic Entropy (SE) has established itself as a prevailing baseline due to its ability to handle semantic equivalence in natural language generation (Farquhar et al., 2024; Kuhn et al., 2023). Unlike traditional token-level entropy, SE computes uncertainty based on the distribution of semantic classes by clustering multiple sampled responses, thereby avoiding inflated entropy scores caused by lexical variations of identical meanings. Yet, as illustrated in Figure 1 this closed-set view collapses under “consistent hallucinations”: the language model can return the same wrong answer with low entropy but high risk. Although SE excels at detecting divergent errors, its reliance on a “closed-set assumption” presents a fundamental flaw when confronting “consistent hallucinations” (Xu et al.,

061	2025). Consistency-based detectors such as Self-	tic Entropy as the root cause of its failure	112
062	CheckGPT (Manakul et al., 2023) share this weak-	on consistent hallucinations and propose a	113
063	ness: when the model repeatedly produces the same	Bayesian solution with an explicit Ghost	114
064	incorrect response, semantic divergence is absent	Cluster to model the unknown.	115
065	and the detector remains silent. This failure stems		
066	from estimating probabilities solely based on lim-	• Energy-Guided Prior Modulation: We pro-	116
067	ited sampled observations, ignoring the true con-	pose an energy-guided dynamic prior that uti-	117
068	confidence information embedded in the model’s in-	lizes raw logits to detect internal model con-	118
069	ternal latent states. Furthermore, achieving stable	fusion, enabling GLIB-SE to distinguish be-	119
070	entropy estimates typically requires a fixed and	tween “justified confidence” (true knowledge)	120
071	substantial number of samples, imposing high com-	and “blind confidence” (hallucination).	121
072	putational overhead during inference and limiting		
073	real-time applicability.	• Pareto-Efficient Inference: We implement a	122
074	To address these challenges, we propose GLIB-	variance-based adaptive sampling algorithm,	123
075	SE (Ghost-Logit Integrated Bayesian Semantic En-	achieving a superior Pareto frontier between	124
076	tropy). Fundamentally, standard SE estimates se-	detection accuracy and computational cost,	125
077	semantic probabilities solely based on observed clus-	reducing sampling overhead by $\sim 30\text{-}50\%$ on	126
078	ter counts; under consistent hallucinations, this	complex reasoning tasks.	127
079	closed-set view falsely assigns zero probability to		
080	unobserved alternatives, leading to severe overcon-	2 Related Work	128
081	fidence. GLIB-SE corrects this bias by adopting a		
082	Bayesian framework that explicitly allocates prob-	2.1 Semantic Uncertainty and Its Limits.	129
083	ability mass to a Ghost Cluster which construct		
084	representing the unobserved semantic space (the	Uncertainty quantification is pivotal for trustwor-	130
085	“unknown unknowns”). Crucially, we utilize the	thy LLMs. While early metrics like perplexity	131
086	model’s raw logit energy as a dynamic prior to cal-	(Min et al., 2020) and recent supervised token-	132
087	ibrate this Ghost mass: when the model exhibits	probability approaches (Quevedo Caballero et al.,	133
088	internal confusion, we inflate the weight of the	2024) are sensitive to lexical variations, Seman-	134
089	Ghost Cluster, effectively uncovering the hidden	tic Entropy (SE) (Farquhar et al., 2024; Kuhn	135
090	uncertainty behind consistent outputs.	et al., 2023) mitigates this by clustering seman-	136
091	A key innovation of GLIB-SE lies in utilizing the	tically equivalent responses. However, SE relies on	137
092	model’s raw logits as a dynamic prior to guide the	a <i>closed-set assumption</i> , equating observed clus-	138
093	weight estimation of the Ghost Cluster. When the	ter frequency with true probability. Consequently,	139
094	model exhibits internal confusion (manifested as a	SE fails under consistent hallucinations , where a	140
095	uniform logit distribution), our mechanism assigns	model confidently repeats the same incorrect an-	141
096	higher probability mass to the Ghost Cluster, even	swer (i.e., a single cluster), causing entropy to	142
097	if the output semantics appear highly consistent.	collapse to zero despite factual error. GLIB-SE	143
098	This effectively calibrates the final entropy to warn	addresses this by explicitly modeling unobserved	144
099	of potential risks. This approach synergizes intrin-	semantics to break the closed-set constraint.	145
100	sic logit information with extrinsic semantic distri-		
101	butions, bridging the gap left by single-perspective	2.2 Consistency Methods and Sampling Cost.	146
102	methods. Moreover, leveraging the Bayesian frame-		
103	work, we derive the variance of the posterior en-	Consistency-based detectors such as SelfCheck-	147
104	tropy to design an adaptive sampling strategy. This	GPT (Manakul et al., 2023) flag divergence across	148
105	strategy enables automatic termination of genera-	generated answers but fail when hallucinations are	149
106	tion upon the convergence of uncertainty estimates,	semantically consistent. High-quality uncertainty	150
107	significantly reducing computational costs without	methods like Deep Ensembles (Lakshminarayanan	151
108	compromising detection performance.	et al., 2017) improve calibration but require train-	152
109	Our main contributions are threefold:	ing multiple models, making them costly. GLIB-	153
		SE instead couples semantic structure with inter-	154
		nal energy and uses adaptive sampling to reach a	155
110	• Ghost-Cluster Bayesian Framework: We	Pareto-efficient frontier without training ensembles	156
111	identify the “closed-set limitation” of Seman-	(Li et al., 2025).	157

2.3 Bayesian Hallucination Detection.

Prior works like Wang et al. (2023); Ciosek et al. (2025) have utilized Bayesian sequential frameworks to optimize query budgets. However, these works lack a mechanism to distinguish between *justified confidence* (knowledge) and *blind confidence* (hallucination). GLIB-SE extends this framework by making the ghost prior *dynamic*, modulated by the model’s internal confusion.

2.4 Logit-Based Energy Signals.

Internal states often reveal uncertainty even when outputs are consistent. Energy-based OOD detection (Liu et al., 2020) establishes the free-energy link for detecting out-of-distribution or erroneous samples via logits. **Semantic Energy** (Ma et al., 2025) leverages the logit landscape as a truthfulness proxy. *Distinction:* While Ma et al. demonstrate the utility of energy, they primarily operate at the sample level or simple aggregation, without integrating energy into the semantic clustering structure. In contrast, GLIB-SE fuses semantic energy directly into the Bayesian framework as a prior. This allows high internal energy to actively expand the probability mass of the Ghost Cluster, effectively detecting consistent hallucinations that output-based clustering alone would miss.

2.5 Kernel-Based Uncertainty.

Moving beyond hard clustering, **Kernel Language Entropy (KLE)** (Nikitin et al., 2024) uses continuous semantic kernels to capture fine-grained similarities. While rigorous, KLE requires computing an $N \times N$ similarity matrix, leading to quadratic complexity $O(N^2)$ that limits scalability. GLIB-SE maintains the linear efficiency $O(N)$ of hard clustering while achieving robust calibration through the energy-guided Ghost Cluster, offering a more Pareto-efficient solution for real-time applications.

3 Methodology

This section formalizes the uncertainty estimation problem, diagnoses the closed-set failure of standard Semantic Entropy (SE), and presents **GLIB-SE**. The framework combines: (1) **Bayesian Semantic Modeling** with an explicit Ghost Cluster for unobserved semantics; (2) a **Logit-Informed Dynamic Prior** driven by Semantic Energy; and (3) a **Variance-Based Adaptive Sampling** strategy. These components are designed jointly for robust-

ness to consistent errors and for **computational efficiency**, decoupling uncertainty estimation from fixed sampling budgets.

3.1 Problem Formulation & Preliminaries

Given an input context \mathbf{x} , an LLM generates a response sequence \mathbf{y} according to $p(\mathbf{y} | \mathbf{x})$. Since multiple sequences can convey the same meaning, we define a semantic clustering function $\mathcal{C}(\cdot)$ such that $\mathcal{C}(\mathbf{y}_i) = c_k$ if \mathbf{y}_i belongs to semantic class c_k .

Semantic Entropy (SE) quantifies uncertainty over meanings. For N sampled sequences $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$ grouped into K clusters with counts $\mathbf{n} = [n_1, \dots, n_K]$, the standard estimator is:

$$\hat{H}_{\text{SE}}(\mathbf{x}) = - \sum_{k=1}^K \frac{n_k}{N} \log \frac{n_k}{N}. \quad (1)$$

Limitation: This frequentist approach suffers from the *closed-set assumption*, assuming the N samples cover the entire semantic space. Under *consistent hallucination*, where the model confidently generates the same incorrect answer N times (i.e., $K = 1$), \hat{H}_{SE} collapses to 0. This fails to capture the *epistemic uncertainty* arising from the finite sample size and ignores the model’s internal confusion.

Consistent Hallucination. We define a consistent hallucination as a scenario where the model generates a set of responses Y that are semantically equivalent ($K = 1, n_1 = N$) and confident (low predictive entropy), yet factually incorrect. In this case, standard SE yields $\hat{H}_{\text{SE}}(\mathbf{x}) = 0$, failing to signal risk. Our goal is to leverage internal states so that low entropy reflects true knowledge (low energy) rather than blind confidence (high energy).

3.2 GLIB-SE Framework

We propose to estimate the semantic distribution π in a Bayesian setting, dynamically guided by the model’s internal state. Figure 2 presents an overview of the proposed GLIB-SE framework.

3.2.1 Bayesian Modeling of Semantic Distribution

We introduce an unobserved **Ghost Cluster** c_g to represent semantics not seen in the sampled responses, extending π to $[\pi_1, \dots, \pi_K, \pi_g]$ following Bayesian SE (Ciosek et al., 2025). With Dirichlet modeling,

$$\pi | Y \sim \text{Dir}(\alpha_{\text{base}} + n_1, \dots, \alpha_{\text{base}} + n_K, \alpha_g), \quad (2)$$

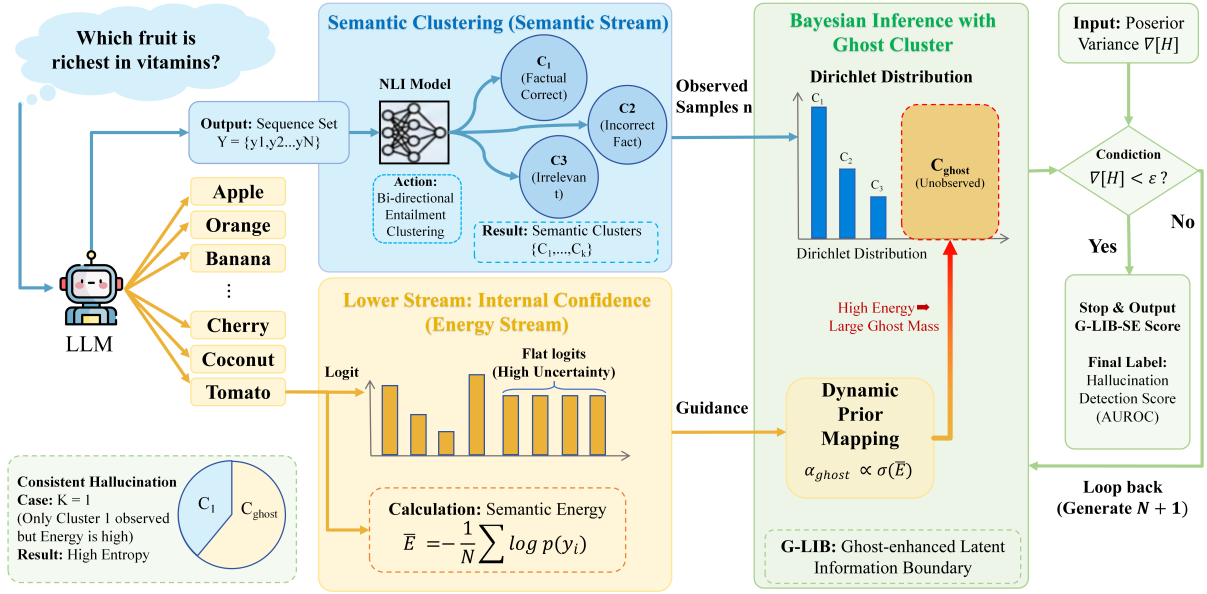


Figure 2: Overview of the GLIB-SE framework. Given an input query, the LLM generates multiple responses. We cluster responses based on semantic equivalence to obtain observed semantic distributions. Simultaneously, we compute Semantic Energy from raw logits to capture the model’s internal confusion. The core GLIB mechanism uses the Semantic Energy to dynamically calibrate the prior weight (α_g) of a Ghost Cluster, representing unobserved semantics. This Bayesian integration allows the framework to detect consistent hallucinations. An adaptive sampling strategy terminates generation when the posterior entropy variance converges, ensuring efficiency.

where α_{base} smooths observed clusters and α_g is a dynamic prior for the ghost mass. The ghost count is always zero ($n_g = 0$), so the posterior parameters are $(\alpha_{\text{base}} + n_1, \dots, \alpha_{\text{base}} + n_K, \alpha_g)$. Posterior expectations smooth the empirical frequencies:

$$\mathbb{E}[\pi_k | Y] = \frac{n_k + \alpha_{\text{base}}}{N + K\alpha_{\text{base}} + \alpha_g}. \quad (3)$$

The ghost count is always zero, so π_g is entirely prior-driven. Uncertainty is measured by the posterior expected entropy:

$$H_{\text{GLIB}}(\mathbf{x}) = \mathbb{E}_{\pi \sim P(\cdot | Y)} \left[- \sum_{c \in \{c_1, \dots, c_K, c_g\}} \pi_c \log \pi_c \right], \quad (4)$$

approximated via Monte Carlo samples from the Dirichlet posterior.

3.2.2 Energy-Guided Prior

A static ghost cluster prior α_g implies a constant belief about the unknown, ignoring the model’s varying internal confidence. We propose to dynamically modulate this prior by interpreting **Semantic Energy** (Ma et al., 2025) as a proxy for *epistemic uncertainty* in the latent space.

Theoretical Justification. Standard Semantic Entropy operates in the discrete output space \mathcal{Y} , which is often sparse. In contrast, the model’s logits \mathbf{z} lie in a continuous latent space. High Semantic Energy, characterized by a flat logit distribution (high entropy in \mathbf{z}), indicates that the model’s internal state is close to the decision boundary or Out-of-Distribution (OOD) regions. From a Bayesian perspective, when the model’s internal state is ambiguous (High Energy), our prior belief in the observed semantic distribution should decrease, and the probability mass assigned to the “unknown” (Ghost Cluster) should increase. Thus, we formulate the ghost prior α_g not as a fixed hyperparameter, but as a function of the internal energy state $E(\mathbf{x})$.

Energy Calculation. For a generated sequence \mathbf{y} of length L , we compute the *sequence-level energy* $E(\mathbf{y})$ as the length-normalized token-level negative log-likelihood under p_θ (approximating the free energy while normalizing for length) (Liu et al., 2020):

$$E(\mathbf{y}) = -\frac{1}{L} \sum_{t=1}^L \log p_\theta(y_t | \mathbf{y}_{<t}, \mathbf{x}). \quad (5)$$

We then average across the N sampled sequences to obtain the *query-level energy* $\bar{E}(\mathbf{x})$, which serves as a proxy for global internal confusion. All experiments and the Appendix A.2 pseudocode use Eq. (5), implemented from token log-probabilities; a raw negative-mean-logit proxy is not used (it differs by the log-sum-exp softmax shift and is not a monotone transform of Eq. (5)).

Dynamic Prior Modulation. To map the continuous energy \bar{E} to the Dirichlet prior α_g , we use a **Soft-Gating Mechanism**:

$$\alpha_g(\bar{E}) = \alpha_{\text{base}} \cdot \left(1 + S_{\text{max}} \cdot \sigma \left(\frac{\bar{E} - \tau}{T} \right) \right), \quad (6)$$

where $\sigma(\cdot)$ is the sigmoid function.

- **Physical Interpretation:** This function models a phase transition. When $\bar{E} \ll \tau$ (High Confidence), the term $\sigma(\cdot) \rightarrow 0$, and α_g collapses to α_{base} , reducing our method to standard Bayesian SE. When $\bar{E} \gg \tau$ (High Uncertainty/Hallucination), $\sigma(\cdot) \rightarrow 1$, and α_g is amplified by S_{max} , effectively forcing the posterior distribution to flatten and shifting mass to the ghost cluster.
- **Role of Parameters:** S_{max} defines the maximum penalty for internal confusion, while T (temperature) controls the sensitivity of the transition.

Unsupervised Self-Calibration via GMM. To set the energy threshold τ without labels, we fit a two-component Gaussian Mixture Model (GMM) on historical energy scores:

$$p(E) = w_1 \mathcal{N}(\mu_1, \sigma_1^2) + w_2 \mathcal{N}(\mu_2, \sigma_2^2). \quad (7)$$

We select τ at the intersection of the two Gaussians (where $P(C_1 | E) = P(C_2 | E)$), enabling model-specific, label-free calibration of what constitutes “High Energy”. A full-step pseudocode description of GLIB-SE is provided in Appendix A.2.

3.3 Variance-Based Adaptive Sampling

Computing SE typically requires a fixed, high number of samples N . We propose an adaptive strategy that stops generation when the uncertainty estimate stabilizes, measured by the variance of the posterior entropy:

$$\mathbb{V}_H = \text{Var}_{\pi \sim P(\pi | Y_n)} [H(\pi)]. \quad (8)$$

In practice we draw S posterior samples $\{\pi^{(s)}\}_{s=1}^S$ from the Dirichlet posterior to estimate

$$\hat{\mathbb{V}}_H(n) = \frac{1}{S-1} \sum_{s=1}^S \left(H(\pi^{(s)}) - \frac{1}{S} \sum_{j=1}^S H(\pi^{(j)}) \right)^2, \quad (9)$$

and iteratively generate samples until $\hat{\mathbb{V}}_H(n) < \epsilon$ or a budget N_{max} is reached. This allocates more compute to ambiguous queries while saving resources on simple ones.

4 Experiments

To comprehensively evaluate the effectiveness of the GLIB-SE framework, we conducted extensive experiments across six benchmark datasets covering diverse task types, including short-form question answering, long-form text generation, and mathematical reasoning. The experiments aim to answer three core research questions: **(1)** Does GLIB-SE outperform existing state-of-the-art methods in detecting LLM hallucinations, particularly in scenarios prone to consistency errors? **(2)** Can the introduced adaptive sampling strategy significantly reduce inference costs while maintaining detection performance? **(3)** Does the logit-based energy guidance mechanism effectively address the failure of semantic entropy in consistent hallucination scenarios?

4.1 Experimental Setup

Datasets and Models. Our evaluation encompasses diverse task scenarios to ensure generalizability. For short-form QA, we used **TriviaQA** (Joshi et al., 2017) and **NQ-Open** (Kwiatkowski et al., 2019); for long-form generation, we selected **BioASQ** (Tsatsaronis et al., 2015) and **CoQA** (Reddy et al., 2019); for mathematical reasoning, we employed **GSM8K** (Cobbe et al., 2021) and **SVAMP** (Patel et al., 2021). We selected mainstream LLMs covering different parameter scales and architectures, including lightweight models (Qwen3-4B-Instruct, Llama-3-8B-Instruct), mid-scale models (Llama-2-13B-Chat), and high-performance models (Llama-3.1-70B-Instruct, Qwen3-32B, Qwen3-30B-A3B, Phi-4) (Touvron et al., 2023; Dubey et al., 2024; Yang et al., 2024; Abdin et al., 2024).

Baselines and Metrics. We compared GLIB-SE against a range of mainstream uncertainty estimation methods: (1) **Token-level:** Naive Entropy and P(True); (2) **Semantic-level:** Semantic Entropy

Models	P(True)		Semantic Entropy		Semantic Energy		KLE		GLIB-SE (Ours)		
	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	Avg. Samples
<i>TriviaQA</i>											
Llama-3-8B	0.791	0.702	0.864	0.810	0.832	0.755	0.887	0.830	0.865	0.815	6.07
Llama-2-13B	0.594	0.452	<u>0.825</u>	<u>0.731</u>	0.809	0.715	<u>0.825</u>	0.728	0.832	0.740	6.78
Qwen3-8B	0.872	0.798	0.837	0.765	0.800	0.712	0.848	0.774	<u>0.856</u>	<u>0.782</u>	7.87
Qwen3-4B	0.883	0.812	0.851	0.795	0.851	0.788	<u>0.868</u>	0.801	0.863	<u>0.805</u>	6.92
<i>NQ-Open</i>											
Llama-3-8B	0.761	0.642	0.713	0.605	0.759	0.655	0.777	0.670	<u>0.762</u>	<u>0.658</u>	7.69
Llama-2-13B	0.627	0.510	<u>0.739</u>	<u>0.615</u>	0.701	0.588	0.758	0.625	0.724	0.608	6.11
Qwen3-8B	0.781	0.665	<u>0.779</u>	<u>0.660</u>	0.670	0.535	0.500	0.415	0.761	0.652	6.57
Qwen3-4B	<u>0.769</u>	<u>0.654</u>	0.780	0.672	0.657	0.512	0.500	0.415	0.754	0.648	6.79
<i>BioASQ</i>											
Llama-3.1-70B	0.758	0.685	0.816	0.752	<u>0.827</u>	<u>0.760</u>	0.500	0.450	0.844	0.775	5.62
Seed-Coder-8B	0.672	0.590	<u>0.824</u>	<u>0.765</u>	0.802	0.730	0.500	0.450	0.847	0.788	8.48
Qwen3-30B-A3B	0.640	0.550	<u>0.743</u>	0.680	<u>0.757</u>	<u>0.695</u>	0.500	0.450	0.774	0.712	4.38
Qwen3-32B	0.500	0.450	0.773	0.710	<u>0.794</u>	<u>0.725</u>	0.500	0.450	0.810	0.745	4.57
<i>CoQA</i>											
Llama-3.1-70B	0.798	0.912	0.806	0.909	<u>0.825</u>	0.923	0.500	0.752	0.836	<u>0.922</u>	8.32
Qwen3-30B-A3B	0.739	0.889	0.835	<u>0.932</u>	0.790	0.907	0.500	0.762	<u>0.828</u>	0.938	6.56
Qwen3-32B	0.500	0.420	<u>0.815</u>	<u>0.920</u>	0.809	0.903	0.500	0.756	0.831	0.936	7.33
<i>GSM8K</i>											
Llama-3.1-70B	0.614	0.530	0.831	0.765	0.761	0.690	0.500	0.420	<u>0.819</u>	<u>0.758</u>	7.38
Phi-4	0.496	0.420	<u>0.689</u>	<u>0.610</u>	0.554	0.480	0.500	0.420	0.692	0.615	9.41
Qwen3-32B	0.500	0.420	0.781	0.720	0.656	0.580	0.500	0.420	<u>0.724</u>	<u>0.655</u>	5.25
<i>SVAMP</i>											
Llama-3.1-70B	0.583	0.505	0.840	0.775	<u>0.878</u>	<u>0.815</u>	0.856	0.795	0.901	0.845	8.25
Llama-2-70B	0.480	0.410	<u>0.849</u>	<u>0.785</u>	0.844	0.780	0.850	0.790	0.801	0.740	5.48
Phi-4	0.563	0.485	<u>0.870</u>	<u>0.810</u>	<u>0.891</u>	<u>0.835</u>	<u>0.891</u>	<u>0.835</u>	0.909	0.855	6.73

Table 1: **Hallucination detection: AUROC/AUPR vs. baselines.** (*textbfBold* means per-row best; *Underline* means per-row second-best. All the Non-GLIB methods use a fixed 10 samples; GLIB-SE extraly reports its adaptive average.)

(SE) (Farquhar et al., 2024) and Distributed Semantic Entropy (DSE); (3) **Logit/Energy-based**: Semantic Energy (SEn) (Ma et al., 2025); (4) **Kernel-based**: Kernel Language Entropy (KLE) (Nikitin et al., 2024). To verify the unique contribution of the energy guidance mechanism, we also included **BayesSE** (Bayesian Semantic Entropy without energy guidance) as an ablation baseline. All experiments use two core metrics for hallucination detection: **Area Under the Receiver Operating Characteristic (AUROC)**, which measures threshold-agnostic ranking quality, and **Area Under the Precision-Recall Curve (AUPR)**, which captures detector precision under class imbalance typical of hallucination vs. truthful responses.

Implementation Details. For semantic equivalence judgment, we used the DeBERTa-large-mnli model (He et al., 2021) for bi-directional entailment checking in short-form tasks. In long-form and reasoning tasks, we employed GPT-4o as a judge model to ensure accurate semantic cluster-

ing. In GLIB-SE, the base prior α_{base} was set to 0.1. The energy threshold τ was automatically determined using the GMM-based self-calibration in Section 3.2.2, fitted on 50 unlabelled validation samples; other hyperparameters followed the same global setting across tasks.

4.2 Main Results: Hallucination Detection

Table 1 summarizes AUROC/AUPR across baselines and GLIB-SE with their sampling costs.

Solving Consistent Hallucinations (SVAMP & GSM8K). The most significant improvements are observed in mathematical reasoning tasks, which are notoriously prone to consistent hallucinations (where models confidently repeat the same wrong calculation). On **SVAMP**, GLIB-SE outperforms the standard Semantic Entropy (SE) by a large margin on several models. For instance, on **Llama-3.1-70B-Instruct**, GLIB-SE achieves an AUROC of **0.901**, compared to 0.840 for SE (+6.1%). Similarly, on **Phi-4**, it improves from 0.870 (SE) to **0.909**. This validates our core hypothesis: when

Dataset	Model	BayesSE	GLIB-SE
TriviaQA	Qwen3-4B-Instruct	0.844 / 6.92	0.863 / 6.92
	Qwen3-8B	0.846 / 7.87	0.856 / 7.87
	Llama-2-13B-chat-hf	0.817 / 6.78	0.832 / 6.78
	Llama-3-8B-Instruct	0.853 / 6.07	0.865 / 6.07
NQ-Open	Qwen3-4B-Instruct	0.749 / 6.79	0.754 / 6.79
	Qwen3-8B	0.766 / 6.57	0.761 / 6.57
	Llama-2-13B-chat-hf	0.709 / 6.11	0.724 / 6.11
	Llama-3-8B-Instruct	0.747 / 7.69	0.762 / 7.69
BioASQ	Llama-3.1-70B-Instruct	0.807 / 5.62	0.844 / 5.62
	Seed-Coder-8B-Instruct	0.845 / 8.48	0.847 / 8.48
	Qwen3-32B	0.745 / 4.57	0.810 / 4.57
	Qwen3-30B-A3B	0.700 / 4.38	0.774 / 4.38
CoQA	Qwen3-30B-A3B	0.794 / 6.56	0.828 / 6.56
	Qwen3-32B	0.809 / 7.33	0.831 / 7.33
	Llama-3.1-70B-Instruct	0.834 / 8.32	0.836 / 8.32
GSM8K	Phi-4	0.682 / 9.41	0.692 / 9.41
	Qwen3-32B	0.715 / 5.25	0.724 / 5.25
	Llama-3-70B	0.801 / 7.38	0.819 / 7.38
SVAMP	Phi-4	0.882 / 6.73	0.909 / 6.73
	Llama-2-70B-chat-hf	0.788 / 5.48	0.801 / 5.48
	Llama-3.1-70B-Instruct	0.864 / 8.25	0.901 / 8.25

Table 2: **Ablation: BayesSE vs. GLIB-SE** (AUROC / Avg. Samples).

models are “confidently wrong” (producing consistent outputs but with flat logits), the high Semantic Energy successfully triggers the Ghost Cluster to inflate uncertainty, whereas standard SE collapses to zero.

Robustness on Long-Form Generation (BioASQ). On long-form QA, GLIB-SE consistently surpasses SE. For **Qwen3-32B** on BioASQ, GLIB-SE improves AUROC from 0.773 (SE) to **0.810** (+3.7%).

Short-form factoid tasks (TriviaQA). On TriviaQA the gains are modest because errors are typically either clearly correct or clearly wrong with divergent semantics, where standard SE already issues high entropy. GLIB-SE matches strong baselines here, indicating the Ghost Cluster primarily activates in regimes with consistent-but-wrong generations rather than single-shot factoid mistakes.

Failure Mode Analysis of Baselines. We observe distinct failure modes in existing methods, highlighting the robustness of our approach:

- **KLE Collapse:** On complex tasks like BioASQ and GSM8K, KLE performance degenerates to near-random guessing (AUROC ≈ 0.50). This is attributed to the use of discrete judge models (e.g., GPT-4o) which output hard entailment labels. This collapses KLE’s continuous semantic kernel into a binary matrix, stripping the method of its ability to measure fine-grained semantic distance (Nikitin et al., 2024). GLIB-SE avoids this by operating on the semantic distribution level rather than pairwise kernels.

- **P(True) Miscalibration:** Simple probability-based metrics fail on reasoning tasks with smaller models (e.g., Phi-4 on GSM8K, AUROC 0.49). This reflects severe *miscalibration*, where the model places high confidence on incorrect reasoning paths. GLIB-SE mitigates this by introducing the Ghost Cluster, which effectively acts as a “calibration buffer” triggered by the internal energy state.

4.3 Efficiency Analysis: Pareto Efficiency

Inference cost is a critical bottleneck for semantic uncertainty methods, which typically require sampling $N = 5$ to 10 sequences. GLIB-SE achieves a superior Pareto frontier by dynamically allocating the sample budget.

Cost-Performance Ratio. As shown in the last column of Table 1 and Figure 3, GLIB-SE maintains SOTA detection performance with significantly fewer sampled sequences. For instance, on **BioASQ** (Qwen3-32B), it requires only **4.57 sampled sequences** on average to surpass the performance of fixed-sample SE ($N = 10$). This represents a **54.3% reduction** in computational overhead. Similarly, on **TriviaQA**, GLIB-SE converges in ~ 6 sampled sequences.

Adaptive Mechanism. This efficiency stems from our variance-based stopping criterion. The model exits early (in 2–3 samples) for unambiguous queries where the posterior $\mathbb{V}[H]$ vanishes quickly, while reserving the budget for high-entropy, ambiguous queries. Figure 3 illustrates this dominance, where GLIB-SE’s curve lies strictly to the upper-left (better utility, lower cost) of the fixed- N baselines.

4.4 Mechanism Analysis

Ablation Study. Table 2 details the BayesSE vs. GLIB-SE ablation. On **SVAMP** (Llama-3.1-70B), adding energy guidance improves AUROC from 0.864 (BayesSE) to 0.901 (GLIB-SE). This confirms that the performance gain is not solely due to the Bayesian formulation but specifically stems from the *Logit-Informed Dynamic Prior*, which correctly up-weights the Ghost Cluster when internal confidence is low.

Visualizing the Correction. As visualized in our Density-Entropy analysis (Figure 4), standard SE fails for samples in the “low entropy, high energy” region (red dots, representing consistent hallucinations). GLIB-SE utilizes the vertical axis (energy) to distinguish these from true correct an-

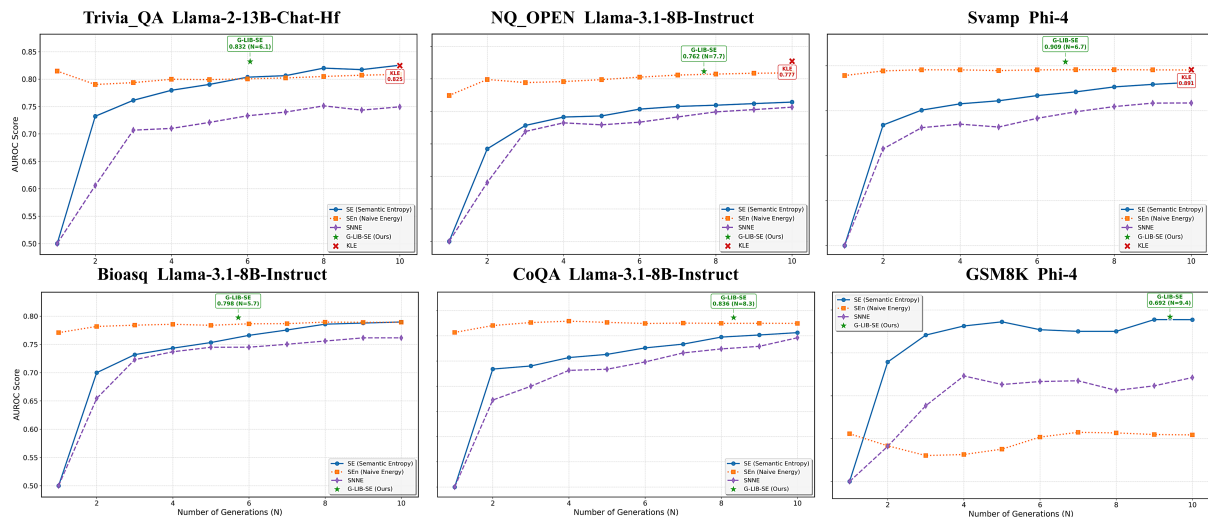


Figure 3: AUROC versus sampled sequences for SE, SEn, SNNE, KLE, and GLIB-SE across all datasets. GLIB-SE curves consistently sit toward the upper-left region of the Pareto frontier, indicating that our method dynamically balances cost and effectiveness to achieve stronger performance with fewer samples.

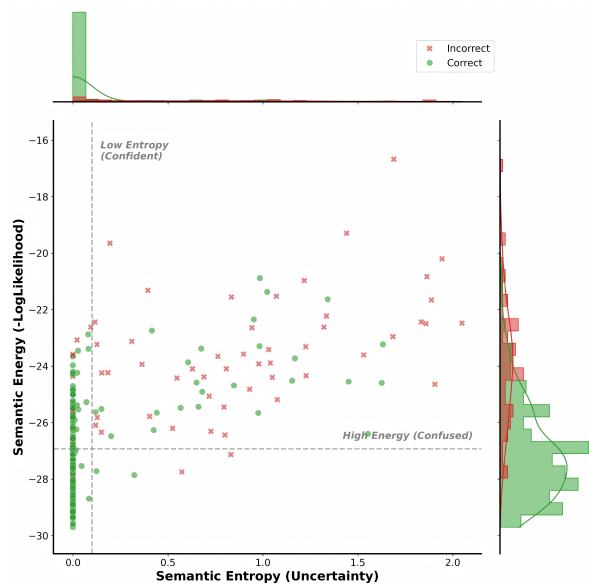


Figure 4: Semantic entropy (x-axis) versus semantic energy (y-axis) for correct (green) and incorrect (red) answers. Points concentrated in the lower-left region (low entropy, low energy) correspond to correct, non-hallucinatory responses, while incorrect answers exhibit higher entropy or energy.

swers (green dots), effectively “pulling” them into the high-uncertainty regime.

5 Conclusion

In this paper, we presented **GLIB-SE**, a probabilistic framework that synergizes internal model confidence (Logit Energy) with external semantic consistency (Semantic Entropy) for robust hal-

lucination detection. The core innovation lies in modeling the unobserved semantic space via a dynamically weighted “Ghost Cluster”, which effectively resolves the critical failure mode of standard semantic entropy under “consistent hallucinations”—scenarios where models are confident yet factually incorrect. Furthermore, our variance-based adaptive sampling strategy successfully decouples uncertainty estimation from fixed computational budgets, establishing a superior Pareto frontier between detection accuracy and inference efficiency on reasoning and generation tasks.

6 Future Work

Despite its state-of-the-art performance, GLIB-SE currently relies on white-box access to output logits, which limits its deployment in certain environments. Future research will prioritize two avenues to broaden its applicability:

- **Black-box Adaptation:** We aim to explore proxy metrics for semantic energy, such as verbalized confidence or perplexity approximations, to extend the framework to closed-source API models.
- **Multimodal Expansion:** We plan to investigate the correlation between “internal confusion” and “semantic consistency” in Vision-Language Models (VLMs), thereby adapting the Ghost Cluster mechanism to address hallucinations in image captioning and visual reasoning domains.

544 Limitations

545 While GLIB-SE demonstrates robust performance,
546 it has three primary limitations:

547 **1. White-box dependency.** The approach relies
548 on access to output logits to compute Semantic
549 Energy. Closed-source API-only models that hide
550 log-probabilities cannot be directly covered unless
551 they expose token-level scores.

552 **2. NLI bottleneck.** Semantic clustering de-
553 pends on an NLI model (e.g., DeBERTa or GPT-
554 4o). For extremely high-throughput deployments,
555 the latency from pairwise NLI checks can dominate
556 total runtime, motivating lighter-weight semantic
557 judges.

558 Ethics Statement

559 This work focuses on improving the reliability and
560 safety of Large Language Models by detecting hal-
561 lucinations.

562 **Data Usage and Privacy.** All datasets used in
563 this study (TriviaQA, NQ-Open, BioASQ, CoQA,
564 GSM8K, SVAMP) are publicly available bench-
565 marks widely used in the NLP community. Our
566 experiments do not involve the collection of new
567 personal data or crowdsourcing that would require
568 IRB approval. We have adhered to the intended
569 use licenses of these datasets and the open-weights
570 models (Llama, Qwen, Phi) employed.

571 **Potential Risks.** While GLIB-SE significantly
572 improves hallucination detection, it is not infallible.
573 We caution against deploying this method as the
574 sole arbiter of truth in high-stakes domains (e.g.,
575 medical diagnosis or legal advice) without human-
576 in-the-loop verification. False negatives remain a
577 possibility, aligning with theoretical findings on
578 the impossibility of perfect automated detection
579 without external ground truth.

580 References

581 Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan,
582 Jyoti Aneja, Ahmed Awadallah, Hany Awad, Nguyen
583 Awad, Nguyen Bach, Amit Bahree, Arash Bakhtiari,
584 and 1 others. 2024. Phi-3 technical report: A highly
585 capable language model locally on your phone. *arXiv*
586 *preprint arXiv:2404.14219*.

587 Kamil Ciosek, Nicolò Felicioni, and Sina Ghiassian.
588 2025. Hallucination detection on a budget: Efficient
589 bayesian estimation of semantic entropy. *Transac-*
590 *tions on Machine Learning Research*.

591 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
592 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias

Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, and 1 others. 2021. Training verifiers
to solve math word problems. *arXiv preprint*
arXiv:2110.14168. 593
594
595
596

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, 597
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letak, 598
Akhil Mathur, Alan Schelten, Amy Yang, Angela 599
Fan, and 1 others. 2024. The Llama 3 herd of models. 600
arXiv preprint arXiv:2407.21783. 601

Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and 602
Yarin Gal. 2024. Detecting hallucinations in large 603
language models using semantic entropy. *Nature*, 604
630:625–630. 605

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and 606
Weizhu Chen. 2021. Deberta: Decoding-enhanced 607
BERT with disentangled attention. In *International* 608
Conference on Learning Representations. 609

Yuxin Huang, Rui Li, Liu Zhang, S Yan, and 1 others. 610
2023. A survey on uncertainty quantification in large 611
language models. *arXiv preprint arXiv:2311.12324*. 612

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan 613
Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea 614
Madotto, and Pascale Fung. 2023. Survey of halluci- 615
nation in natural language generation. *ACM Comput-* 616
ing Surveys, 55(12):1–38. 617

Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke 618
Zettlemoyer. 2017. TriviaQA: A large scale distantly 619
supervised challenge dataset for reading comprehen- 620
sion. In *Proceedings of the 55th Annual Meeting of* 621
the Association for Computational Linguistics, pages 622
1601–1611. 623

Saurav Kadavath, Tom Conerly, Amanda Askell, Tom 624
Henighan, Dawn Drain, Ethan Perez, Nicholas 625
Schiefer, Zac Hatfield-Dodds, and 1 others. 2022. 626
Language models (mostly) know what they know. 627
arXiv preprint arXiv:2207.05221. 628

Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. 629
Semantic uncertainty: Linguistic invariances for un- 630
certainty estimation in natural language generation. 631
In *International Conference on Learning Representa-* 632
tions (ICLR). 633

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red- 634
field, Michael Collins, Ankur Parikh, Chris Alberti, 635
Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken- 636
ton Lee, and 1 others. 2019. Natural questions: a 637
benchmark for question answering research. *Trans-* 638
actions of the Association for Computational Linguis- 639
tics, 7:453–466. 640

Balaji Lakshminarayanan, Alexander Pritzel, and 641
Charles Blundell. 2017. Simple and scalable pre- 642
dictive uncertainty estimation using deep ensembles. 643
In *Advances in Neural Information Processing Sys-* 644
tems, volume 30. 645

646	Jiuck Li, A. Magesh, and V. V. Veeravalli. 2025.	large-scale biomedical semantic indexing and ques-	701
647	Principled detection of hallucinations in large lan-	tion answering competition. <i>BMC Bioinformatics</i> ,	702
648	guage models via multiple testing. <i>arXiv preprint</i>	16:1–28.	703
649	<i>arXiv:2508.18473</i> .		
650	Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan	Xinyu Wang, Yifan Yan, Lifu Huang, Xipeng Zheng,	704
651	Li. 2020. Energy-based out-of-distribution detection.	and Xuanjing Huang. 2023. Hallucination detection	705
652	In <i>Advances in Neural Information Processing Sys-</i>	for generative large language models by Bayesian	706
653	<i>tems (NeurIPS)</i> .	sequential estimation . In <i>Proceedings of the 2023</i>	707
654	Huan Ma, Jiadong Pan, Jing Liu, Yan Chen, Joey Tianyi	<i>Conference on Empirical Methods in Natural Lan-</i>	708
655	Zhou, Guangyu Wang, Qinghua Hu, Hua Wu,	<i>guage Processing</i> , pages 15361–15371.	709
656	Changqing Zhang, and Haifeng Wang. 2025. Se-	Fangyuan Xu, Xiting Hu, Ziyi Yu, L. Lin, X. Zhang,	710
657	semantic energy: Detecting LLM hallucination beyond	Y. Zhang, W. Zhou, J. Gu, and Xiaojun Wan. 2025.	711
658	entropy. <i>arXiv preprint arXiv:2508.14496</i> .	HAD: HAllucination detection language models	712
659	Potsawee Manakul, Adian Liusie, and Mark J. F. Gales.	based on a comprehensive hallucination taxonomy.	713
660	2023. SelfCheckGPT: Zero-resource black-box hal-	<i>arXiv preprint arXiv:2510.19318</i> .	714
661	lucination detection for generative large language	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,	715
662	models. In <i>Proceedings of the 2023 Conference on</i>	Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan	716
663	<i>Empirical Methods in Natural Language Processing</i> ,	Li, Dayiheng Liu, Fei Huang, and 1 others.	717
664	pages 9004–9017.	2024. Qwen2.5 technical report. <i>arXiv preprint</i>	718
665	Sewon Min, Julian Michael, Hannaneh Hajishirzi, and	<i>arXiv:2412.15115</i> .	719
666	Luke Zettlemoyer. 2020. Ambitious and rubbish!	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,	720
667	Revisiting the evaluation of open domain question	Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,	721
668	answering. In <i>Proceedings of the 58th Annual Meet-</i>	Yulong Chen, and 1 others. 2023. Siren’s song in	722
669	<i>ing of the Association for Computational Linguistics</i> ,	the AI ocean: A survey on hallucination in large	723
670	pages 6015–6024.	language models. <i>arXiv preprint arXiv:2309.01219</i> .	724
671	Alexander Nikitin, Jannik Kossen, Yarin Gal, and		
672	Pekka Marttinen. 2024. Kernel language en-		
673	tropy: Fine-grained uncertainty quantification for		
674	LLMs from semantic similarities. <i>arXiv preprint</i>		
675	<i>arXiv:2405.20003</i> .		
676	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.		
677	2021. Are NLP models really able to solve sim-		
678	ple math word problems? In <i>Proceedings of the</i>		
679	<i>2021 Conference of the North American Chapter of</i>		
680	<i>the Association for Computational Linguistics</i> , pages		
681	2080–2094.		
682	E. Quevedo Caballero, A. S. Abdelfattah, A. Rodriguez,		
683	and T. Černý. 2024. Detecting hallucinations in large		
684	language model generation: A token probability ap-		
685	proach. <i>arXiv preprint arXiv:2405.19648</i> .		
686	Siva Reddy, Danqi Chen, and Christopher D Manning.		
687	2019. CoQA: A conversational question answering		
688	challenge. <i>Transactions of the Association for Com-</i>		
689	<i>putational Linguistics</i> , 7:249–266.		
690	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-		
691	bert, Amjad Almahairi, Yasmine Babaei, Nikolay		
692	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti		
693	Bhosale, and 1 others. 2023. Llama 2: Open foun-		
694	dation and fine-tuned chat models. <i>arXiv preprint</i>		
695	<i>arXiv:2307.09288</i> .		
696	George Tsatsaronis, Georgios Balikas, Prodromos		
697	Malakasiotis, Ioannis Partalas, Matthias Zschunke,		
698	Michael Alvers, Dirk Weissenborn, Anastasia		
699	Krithara, Sergios Petridis, Dimitris Polychronopou-		
700	los, and 1 others. 2015. An overview of the BIOASQ		

A Appendix

A.1 Implementation Details

Hyperparameters. We set the base Dirichlet prior $\alpha_{\text{base}} = 0.1$, scaling factor $S_{\text{max}} = 100$, and temperature $T = 0.1$. The energy threshold τ is determined via the GMM self-calibration in Section 3.2.2 on 50 unlabeled queries per model.

Adaptive sampling. The posterior variance threshold is $\epsilon = 1\text{e-}4$ with a maximum sample budget $N_{\text{max}} = 10$.

Compute resources. Experiments use NVIDIA H200 (144GB) GPUs. Semantic equivalence uses microsoft/deberta-large-mnli for short-form tasks and GPT-4o for long-form or reasoning tasks.

A.2 GLIB-SE Pseudocode

Algorithm 1: GLIB-SE: Energy-Guided Ghost-Cluster Semantic Entropy

Input: Query \mathbf{x} , budget N_{max} , variance threshold γ , base prior α_{base} , mapping params Θ

Output: Uncertainty scores H_{glib} , H_{base} ; sample count n

```
1 Initialize  $n \leftarrow 0$ , semantic clusters  $\mathcal{C} \leftarrow \emptyset$ , energy history  $\mathcal{E} \leftarrow \emptyset$ .
2 while  $n < N_{\text{max}}$  do
3    $n \leftarrow n + 1$ .
4   Sample response  $\mathbf{y}_n \sim p(\mathbf{y}|\mathbf{x})$ ; compute energy  $E_n \leftarrow -\frac{1}{L} \sum_{t=1}^L \log p_{\theta}(y_t | \mathbf{y}_{<t}, \mathbf{x})$  (Eq. (5)).
5   Update clusters  $\mathcal{C}$  and counts  $\mathbf{n}$  with  $\mathbf{y}_n$  via bi-directional entailment.
6   Update query energy  $\bar{E} \leftarrow \text{mean}(\mathcal{E} \cup \{E_n\})$ .
7   Compute  $\alpha_g \leftarrow \alpha_{\text{base}} (1 + S_{\text{max}} \cdot \sigma(\frac{\bar{E} - \text{offset}}{\text{temp}}))$  via the chosen mapping  $f(\cdot)$ .
8   Form posterior  $\alpha \leftarrow [\mathbf{n} + \alpha_{\text{base}}, \alpha_g]$ ; estimate  $H_{\text{glib}}$  and variance  $\hat{\mathbb{V}}_H$  with MC samples from  $\text{Dir}(\alpha)$ .
9   if  $\hat{\mathbb{V}}_H < \gamma$  then
10    break
11 Compute baseline  $H_{\text{base}}$  using posterior  $[\mathbf{n} + \alpha_{\text{base}}$  (no ghost).
12 return  $H_{\text{glib}}$ ,  $H_{\text{base}}$ ,  $n$ .
```

A.3 Adaptive Sampling Distribution

Figure 5 visualizes the empirical distribution of sample counts that GLIB-SE uses across datasets and models.

A.4 Sensitivity to Max Scale and Prior Slope

Interpretation. Figures 6 and 7 summarize a grid search over the ghost prior mapping, varying the maximum scale S_{max} and the sigmoid sensitivity σ

(slope). Brighter cells indicate higher AUROC; performance is stable across a broad region, with mild gains when $S_{\text{max}} \in [100, 500]$ and $\sigma \in [0.1, 5.0]$. This suggests the method is robust to moderate mis-specification of these hyperparameters.

748
749
750
751
752

Adaptive Sampling Distribution (GLIB-SE)

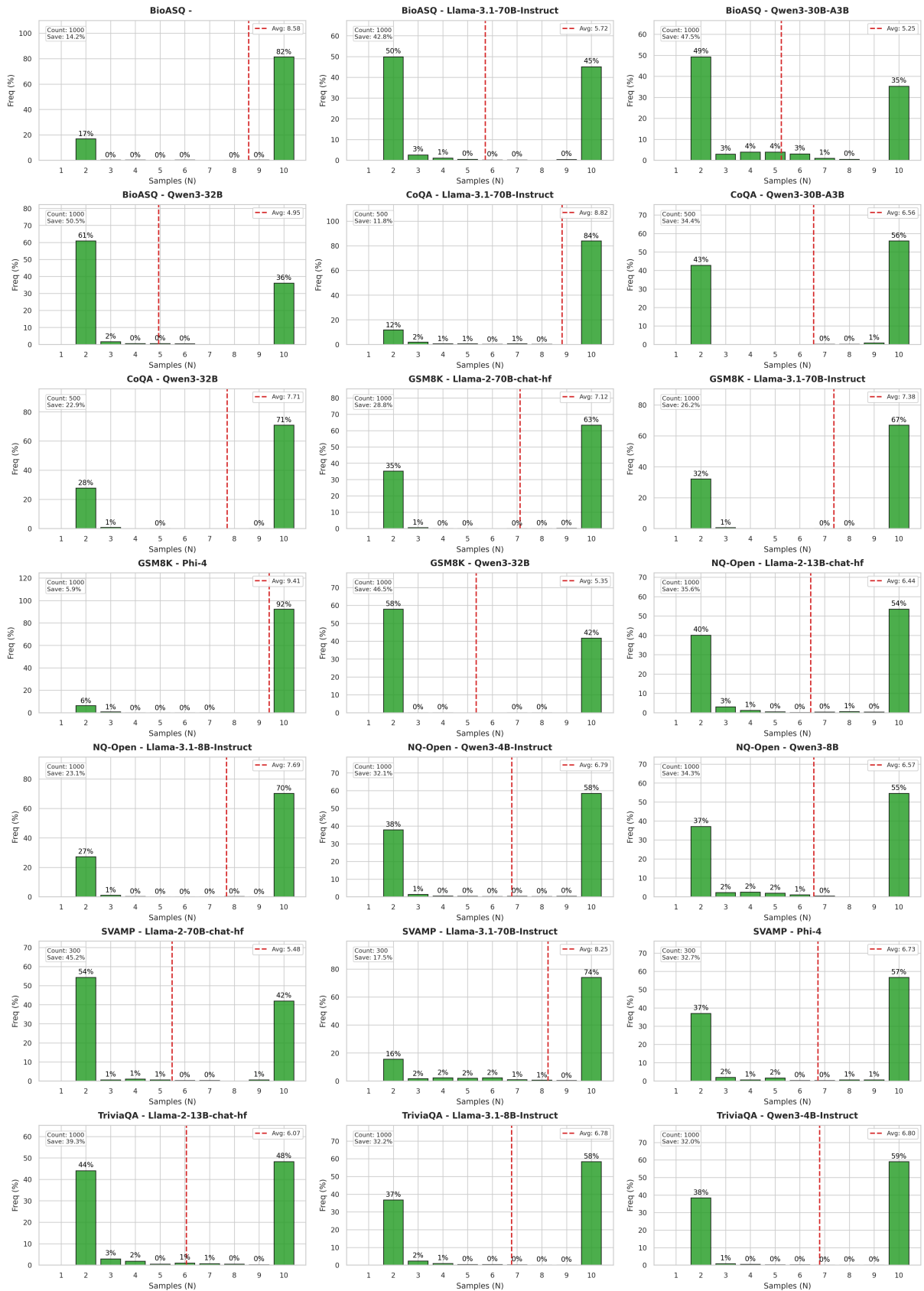


Figure 5: Adaptive sampling histogram for GLIB-SE across datasets and models. Bars indicate the fraction of queries that stop at each sample count N ; the red dashed line denotes the average samples used. The distribution shows aggressive early stopping on easy cases and budget allocation on harder queries.

G-LIB-SE Parameter Sensitivity Analysis (AUROC)

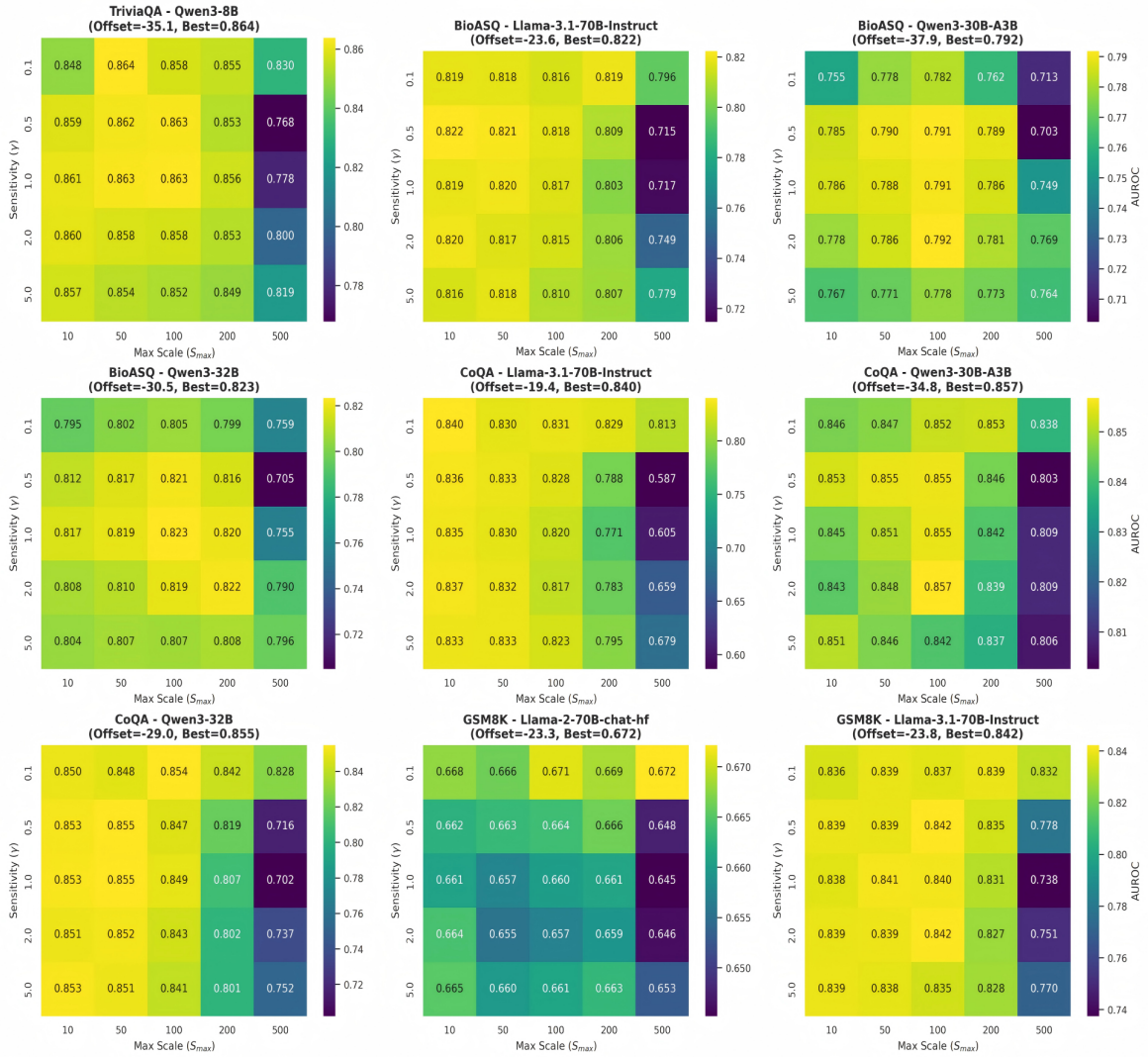


Figure 6: Heatmap analysis (part 1) of AUROC under different ghost prior settings.

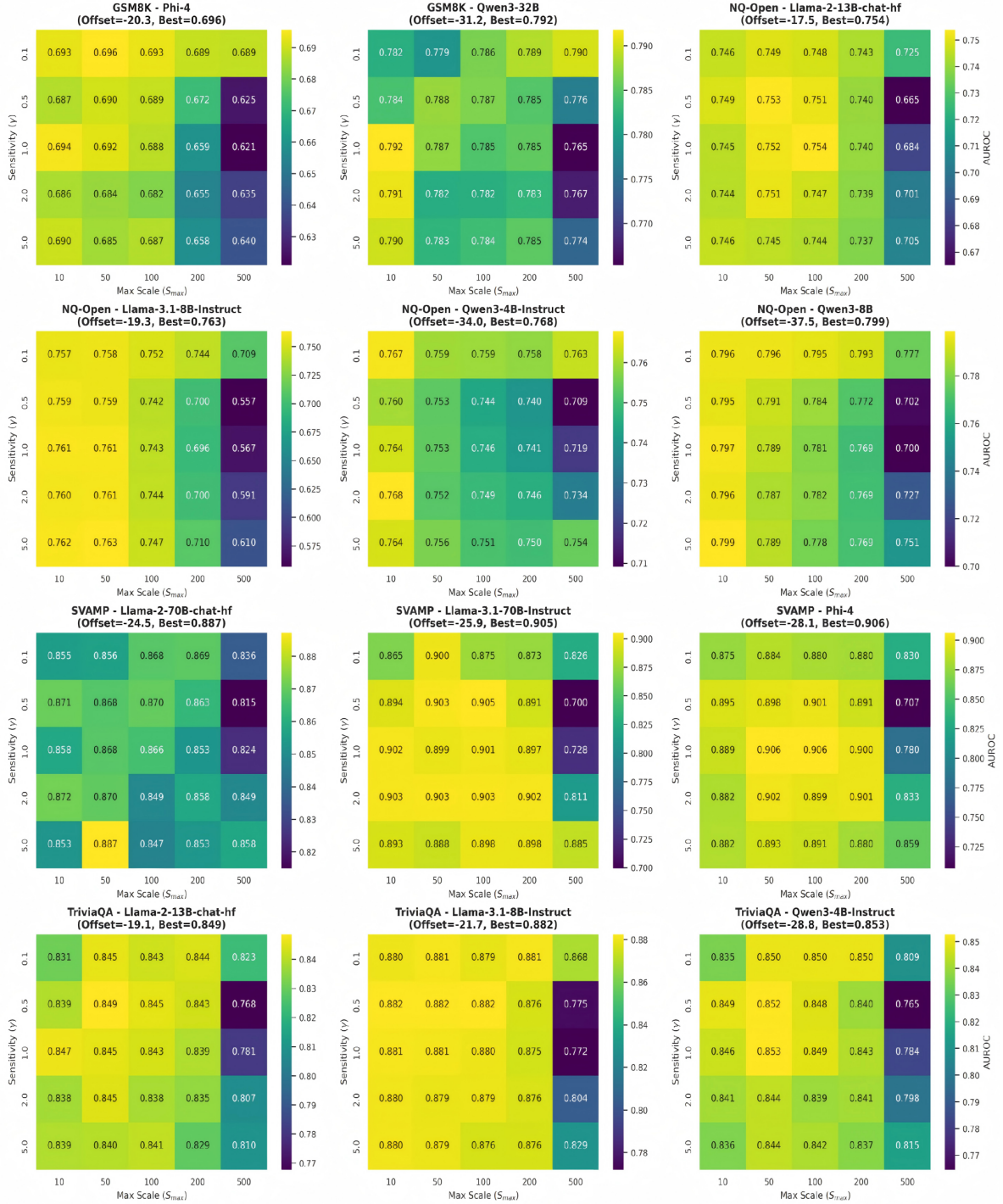


Figure 7: Heatmap analysis (part 2) of AUROC under different ghost prior settings. Rows sweep the sigmoid sensitivity σ (how sharply the prior reacts to logit energy), columns sweep the maximum scale S_{max} (upper bound on ghost mass inflation). The large yellow plateau shows GLIB-SE maintains strong performance over a wide hyperparameter range.