



City Research Online

City St George's, University of London

Citation: Aghazadeh-Chakherloua, R., Guo, Q., Khastgira, S., Popov, P., Zhange, X. & Zhao, X. (2026). A Hierarchical Imprecise Probability Approach to Reliability Assessment of Large Language Models. *Reliability Engineering & System Safety*, 272(Part 2), 112615. doi: 10.1016/j.ress.2026.112615

This is the accepted version of the paper.

This version of the publication may differ from the final published version. To cite this item please consult the publisher's version.

Permanent repository link: <https://openaccess.city.ac.uk/id/eprint/37093/>

Link to published version: <https://doi.org/10.1016/j.ress.2026.112615>

Copyright and Reuse: Copyright and Moral Rights remain with the author(s) and/or copyright holders. Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge, unless otherwise indicated, provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way. For full details of reuse please refer to [City Research Online policy](#).

A Hierarchical Imprecise Probability Approach to Reliability Assessment of Large Language Models

Robab Aghazadeh-Chakherlou^a, Qing Guo^{b,c}, Siddhartha Khastgir^a, Peter Popov^d,
Xiaoge Zhang^e, Xingyu Zhao^{a,*}

^aWMG, University of Warwick, Coventry, United Kingdom

^bCenter for Frontier AI Research, A*STAR, Singapore, Singapore

^cSchool of Computing, National University of Singapore, Singapore, Singapore

^dCenter for Software Reliability, City St George's, University of London, London, United Kingdom

^eDepartment of Industrial and Systems Engineering, The Hong Kong Polytechnic University, Kowloon, Hong Kong

Abstract

Large Language Models (LLMs) are increasingly deployed across diverse domains, raising the need for rigorous reliability assessment methods. Existing benchmark-based evaluations primarily offer descriptive statistics of model accuracy over datasets, providing limited insight into the probabilistic behavior of LLMs under real operational conditions. This paper introduces HIP-LLM, a **H**ierarchical **I**mprecise **P**robability framework for modeling and inferring **LLM** reliability. Building upon the foundations of software reliability engineering, HIP-LLM defines LLM reliability as the probability of failure-free operation over a specified number of future tasks under a given Operational Profile (OP). HIP-LLM represents dependencies across (sub-)domains hierarchically, enabling multi-level inference from subdomain to system-level reliability. HIP-LLM embeds imprecise priors to capture epistemic uncertainty and incorporates OPs to reflect usage contexts. It derives posterior reliability envelopes that quantify uncertainty across priors and data. Experiments on multiple benchmark datasets demonstrate that HIP-LLM offers a more nuanced and standardized reliability characterization than existing benchmark and state-of-the-art approaches. A publicly accessible repository of HIP-LLM is provided.

Keywords:

Large Language Model; Software Reliability; Hierarchical Bayesian Inference; Operational Profile; Epistemic Uncertainty; Imprecise Probability

1. Introduction

Large language models (LLMs) are increasingly applied across a wide range of tasks, from general-purpose reasoning to highly specialized applications. For instance, recent studies [1, 2, 3, 4, 5, 6] employ LLMs as tools to support reliability analysis and safety assurance for safety-critical systems. This expanding scope of application underscores the urgent need for comprehensive and rigorous evaluation methods to assess the *reliability* of LLMs themselves,

*Corresponding author: Xingyu Zhao, xingyu.zhao@warwick.ac.uk

7 given the high-stakes nature of those applications. While numerous studies have proposed
8 domain-specific performance evaluations in areas such as coding [7, 8, 9], medicine [10], and
9 system safety [11, 12, 13], others have focused on broader properties of LLMs, including
10 safety (avoiding harmful content) [14, 15, 16], robustness (resisting adversarial attacks such
11 as jailbreaks) [17, 18], fairness [19], and privacy [20, 21]. However, there remains a notable
12 lack of research dedicated to reliability, which motivates this work.

13 While software *reliability* is highly interrelated with the aforementioned properties (a
14 formal definition and its distinction to other properties will be discussed in later sections), it
15 is a unique property that emphasizes the failure probabilities, specified operational time and
16 conditions, and statistically valid uncertainty quantification. As a “user-centric” property
17 [22, 23, 24], the delivered reliability depends on the behavior of end-users and how the software
18 is expected to be used in practice. This notion is made measurable through the *Operational*
19 *Profile* (OP) [25, 26], which specifies a probability distribution over the types of demands
20 that users actually place on the software. The recent report from OpenAI [27] provide the
21 most comprehensive data on how user-groups use ChatGPT, which aligns with the idea of
22 OPs.

23 Despite software reliability assessment has been studied for decades [28, 29], state-of-the-
24 art LLM evaluation largely relies on benchmarks [30]. This creates several research gaps:

25 Gap-1 (static benchmarks vs dynamic operational use): While we acknowledge the usefulness
26 of benchmarks, they are primarily designed for purposes such as comparing and ranking
27 LLMs on given tasks under fixed datasets, neglecting the variability and dynamic nature
28 of users’ behaviors in real-world applications. Thus, benchmarks cannot serve the
29 purpose of assessing reliability, especially when the OP of the users is different from
30 the data distribution (implicitly) represented by the benchmark dataset.

31 Gap-2 (independent benchmarks vs hierarchical dependencies): While LLMs are generally
32 designed for broad use (covering various task types such as legal reasoning and cod-
33 ing), these tasks naturally exist at different levels of abstraction. It is often useful to
34 organize them hierarchically, dividing the input space into domains and subdomains
35 [31]. High-level domains¹ (e.g., law vs. coding) can be treated as largely independent,
36 whereas subdomains within a domain (e.g., coding in Python vs. coding in C++) are
37 more likely to exhibit dependencies. In such cases, performance in one subdomain may
38 influence users’ confidence in another, whereas cross-domain effects are minimal. Reli-
39 ability assessment should therefore reflect these dependency relationships and support
40 evaluation at multiple abstraction levels, depending on the assessor’s role. For in-
41 stance, a LLM vendor may be interested in system-wide reliability across all domains,
42 a programmer may focus on reliability in coding tasks, and a Python developer may
43 be specifically concerned with the reliability of Python coding. Existing benchmarks,
44 however, typically evaluate LLM capabilities either in isolation or by simply averaging
45 performance across tasks, that obscure these dependency relationships [31].

¹While the precise definition of what constitutes a domain or subdomain is an application-oriented question that may vary case by case reflecting the assessor’s domain knowledge, in this reliability modeling work we do not prescribe specific taxonomies. Instead, we introduce a general hierarchical modeling structure that can accommodate different dependency assumptions (cf. Section 5 for more discussion).

46 Gap-3 (descriptive statistics vs statistical inference): While descriptive statistics summarize
47 the observed dataset by reporting point-estimate benchmark scores, they do not formally
48 model *how the data were generated*. Even when variance is reported as a way of
49 “quantifying uncertainty”, it quantifies variability *within the sample itself*. In contrast,
50 statistical inference treats the observed results as samples from *an underlying stochastic
51 process or population* and uses probabilistic modeling to estimate the parameters of
52 that process [32]. This inferential perspective enables generalization beyond the fixed
53 benchmark *dataset* and supports principled uncertainty quantification about reliability
54 which concerns a underlying *population* of future inputs [31].

55 Indeed, acknowledging these research gaps, recent studies [31, 32] attempt to bridge them
56 by applying both “frequentist” and Bayesian statistical inference methods. However, despite
57 representing important first steps, they do not fully resolve the aforementioned gaps (e.g.,
58 not explicitly incorporate the OP) and introduce new gaps:

59 Gap-4 (failure probability vs probability of future failure-free runs) While both [31, 32] exam-
60 ined the metric of failure probability², they did not extend it to the more practical and
61 widely adopted definition of reliability that predicts the *probability of failure-free runs*
62 over a specified future operational time [33, 34, 35].

63 Gap-5 (non-informative vs informative priors): The work [32] is a frequentist approach that
64 cannot embed prior knowledge explicitly like Bayesian methods. The work [31] employs
65 non-informative priors in its Bayesian models, without facilitating the embedding of any
66 informative prior knowledge that end-users may have.

67 To bridge the five research gaps collectively, we propose HIP-LLM, a hierarchical imprecise
68 probability approach to assess LLM reliability. HIP-LLM devises a more rigorous and
69 versatile assessment with the following key features. First, similar to [31], it structures the
70 model into independent high-level domains, with each domain further divided into depen-
71 dent subdomains. In this way, reliability can be evaluated at multiple levels of detail while
72 preserving statistical dependencies, so that information from one subdomain contributes to
73 the inference about other subdomains within the same domain. We then incorporate OPs as
74 operational weights at each level of the hierarchy (subdomain \rightarrow domain \rightarrow general-purpose
75 LLM) to reflect the dynamic operational use of subdomains and domains in practice. We
76 assume tasks constitute independent and identically distributed (i.i.d.) trials drawn from the
77 specified OP. This assumption is appropriate for reset or single-task scenarios, in which each
78 task is executed in a fresh session without shared context or memory—settings commonly
79 encountered in offline evaluation like benchmarking³ Accordingly, HIP-LLM is intended for
80 LLM usage scenarios in which each task is executed in a new chat session without retaining
81 memory from previous sessions, thereby aligning with the assumption of a “reset” LLM.

82 In contrast to [31], which relies on non-informative priors, HIP-LLM enables the incorpo-
83 ration of informative prior knowledge. Since selecting a single prior in Bayesian inference is

²We later show that this is a special case of our reliability model.

³In contrast, *long-context* or *agentic workflows*, where task outcomes are sequentially dependent through memory or tool use, violate the i.i.d. assumption and require fundamentally different modeling approaches to our HIP-LLM. We provide more discussions later in Remark 1 and Section 5.

84 often controversial due to its encoding of epistemic uncertainty, we adopt *Imprecise Probabil-*
85 *ity* [36, 37] that represents uncertainty about the prior itself without committing to a single
86 distribution. Consequently, reliability metrics (defined in terms of both failure probability
87 and the probability of failure-free runs in future operations) are expressed through *posterior*
88 *distribution envelopes* at different levels of the hierarchy, reflecting both uncertainties from
89 the data and the prior knowledge.

90 In summary, the contributions of this work are as follows:

- 91 • We formally define the reliability assessment problem for LLMs in accordance with
92 established software reliability standards, while delineating its distinction from related
93 properties.
- 94 • We propose HIP-LLM, a hierarchical imprecise probability model that explicitly ad-
95 dresses key research gaps in the current state-of-the-art, by modeling OPs, hierarchical
96 dependencies, statistically principled uncertainty quantification on reliability (both fail-
97 ure probability and the probability future failure-free runs over a specified operational
98 time), and imprecise prior knowledge.
- 99 • We release a public repository containing all experimental data, source code, and models
100 at <https://github.com/aghazadehchakherlou-web/llm-imprecise-bayes>.

101 The remainder of the paper is organized as follows. Section 2 introduces preliminaries and
102 reviews related works on LLM assessment. Section 3 provides the modeling details of HIP-
103 LLM. Section 4 evaluates and illustrates HIP-LLM via experiments on datasets. Section 5
104 discusses the limitations and assumptions, and finally Section 6 concludes the paper.

105 2. Preliminaries and Related Works

106 2.1. Software Reliability

107 According to American National Standards Institute (ANSI), software reliability is defined
108 as [38]:

109 **Definition 1** (Software reliability). *The probability of failure-free software operation for a*
110 *specified period of time in a specified environment.*

111 This probabilistic definition is also largely adopted by software reliability engineering
112 literature [28], which makes reliability amenable to statistical modeling and permits risk-
113 aware aggregation across heterogeneous operational uses. The specific probabilistic metric
114 used to model reliability is domain-dependent and determined by the operational nature of
115 the software. For example, continuous-time systems are continuously operated in the active
116 control of a process, whereas on-demand systems are only invoked upon receipt of discrete
117 demands. In the latter case, such as nuclear power protection systems, the probability of
118 failure on demand (*pdf*) has been adopted in standards [23, 39, 40] and extensively studied
119 [41, 33, 42].

120 Although related, reliability is distinct from safety, security, and accuracy. Reliability fo-
121 cuses on the *probability* of failures (relative to a specification), a stochastic and usage-weighted
122 concept by OPs [42]. Safety concerns *critical failures* with *catastrophic consequences*, i.e.,

123 whether a failure can lead to unacceptable *harm*. For instance, a system can be reliable but
124 unsafe (if frequent but non-harmful failures are tolerated while rare failures produce catas-
125 trophic harm), or safe but unreliable (if it fails frequently but never produce catastrophic
126 harm by fail-safe design)⁴. Security is concerned with resisting *malicious threats*: it normally
127 assumes an explicit threat model and seeks to prevent or mitigate intentional compromises).
128 While security fundamentally requires the analysis of a malicious threat actor, reliability and
129 safety can be studied within benign operational environments where failure originates from
130 sources such as design flaws and implementation bugs. Within the AI/ML community, accu-
131 racy is typically a narrower performance metric, defined as a descriptive statistic (proportion
132 of correct predictions) over a fixed and given dataset. However, such a dataset does not nec-
133 essarily represent the unknown ground-truth population of inputs that the model will face
134 upon deployment, and accuracy does not quantify the uncertainties of failure-free operations
135 over future time like reliability.

136 A central insight from software reliability engineering is that reliability is *user-context*
137 *dependent*: the same software can exhibit very different delivered reliability under different
138 patterns of use. The OP formalizes this dependency, which is defined as [23, 25, 28]:

139 **Definition 2** (Operational Profile (OP)). *An OP is a probability distribution of inputs to*
140 *a software system, representing the relative frequencies with which different input sequences*
141 *expected occur during actual operation.*

142 There are some caveats regarding the term “inputs” in the definition of OP, particularly
143 when considering factors such as software memory and dependencies⁵. We refer readers to
144 [23] for a detailed discussion. For LLMs, OPs serve an analogous role [27]: they specify the
145 distribution of task types (e.g., factual question answering, coding, summarization) and their
146 relative importance in deployment. Incorporating OPs into LLM reliability analysis ensures
147 that evaluation metrics capture not just aggregate benchmark scores, but also how well a
148 model performs across the realistic mix of tasks it is expected to face.

149 2.2. Benchmark-Based Evaluation of LLMs via Descriptive Statistics

150 Most evaluation methods for LLMs report only point estimates such as accuracy scores.
151 Question-answering benchmarks such as MMLU [43] and RACE [44] measure how often a
152 model picks the right answer in multiple-choice questions. Code benchmarks like DS-1000
153 [8], InterCode [7], and HumanEval Pro/MBPP Pro [9] test whether generated programs
154 run correctly. Assistant and autonomy benchmarks such as GAIA [45] and H-CAST [46]
155 look at overall task completion—whether the system can finish a complex job successfully.
156 Model cards and reports for Claude 3 [47], GPT-4o [48], and code models [49] publish these
157 benchmark results as single numbers per model. Broader frameworks like HELM [50] add

⁴We note that for systems in which all failures are catastrophic (e.g., safety-critical systems), safety and reliability assessments do not require different statistical reasoning, such as when measuring the probability that a system will operate safely (i.e., without critical failures) over a given mission.

⁵In practice, because software systems may possess memory, their success/failure often depends on entire sequences of inputs rather than isolated ones. Accordingly, “inputs” may denote a complete stimulus, which may encompass the full sequence of interactions within a task, or even the cumulative sequence of inputs since the system was last reinitialized.

158 other measures such as calibration, robustness, fairness, and efficiency, but still mainly give
159 point values. In real-world use, Clio [21] collects statistics on what kinds of tasks people
160 actually do with AI with a focus on privacy. Method studies like multilingual chain-of-
161 thought [51] and ReAct [52] suggest ways to improve reasoning and interaction, judged by
162 benchmark accuracy.

163 Recent studies widely acknowledge that benchmark-based evaluations for LLMs are valu-
164 able but inherently limited due to the new characteristics and challenges introduced by LLMs.
165 Benchmarks facilitate progress tracking and model comparison, yet they capture only a par-
166 tial view of LLM capabilities. Chang et al. [30] note that traditional benchmarks (reusing
167 fixed task sets for ranking models) have driven early advances but are now approaching their
168 limits due to data repetition, task saturation⁶, and limited assessment of reasoning or interac-
169 tivity. Liu et al. [53] emphasize that standard benchmarks often overlook moral, safety, and
170 social dimensions of model behavior. In domain-specific contexts, Croxford et al. [10] show
171 that medical benchmarks miss crucial aspects such as factual accuracy and clinical reasoning,
172 highlighting the need for domain prior knowledge. Saleh et al. [54] observe that efficiency and
173 performance benchmarks remain fragmented and inconsistent. From a safety perspective, Liu
174 et al. [55] warn that toxicity and robustness tests are often too narrow and easily gamed,
175 recommending flexible, real-world evaluations. Finally, Ye et al. [56] introduce psychometric
176 principles, arguing that current benchmarks frequently lack checks for reliability and fairness.

177 While we acknowledge the usefulness of LLM benchmarks for progress tracking of LLM
178 versions, comparison and ranking of LLMs, they cannot be used to make claims on LLM
179 reliability. Our HIP-LLM approach goes beyond benchmarks, by framing the reliability
180 assessment of LLMs as statistical inference problems.

181 2.3. Statistical Inference for Evaluating LLMs

182 Benchmark-based descriptive statistics simply summarize the collected data, whereas sta-
183 tistical inference draws probabilistic conclusions about the *underlying population or process*
184 that generated those data, thereby enabling generalization, prediction and uncertainty quan-
185 tification. In the AI/ML community, many studies have proposed statistically principled ap-
186 proaches to quantify uncertainty for *individual inputs*, such as text prompts for LLMs, using
187 methods like conformal prediction [57, 58]. However, these approaches focus on instance-level
188 uncertainty and do not capture uncertainty at the operational reliability level. To the best
189 of our knowledge, only two existing works extend uncertainty quantification to this broader
190 reliability perspective for LLMs.

191 In the Anthropic report [32], Miller frames LLM evaluations as statistical experiments and
192 argues for reporting uncertainty alongside benchmark scores. Using classical inference tools,
193 the study constructs confidence intervals via the Central Limit Theorem, applies clustered
194 standard errors for correlated items, recommends paired inference for model comparisons,
195 and develops power analysis and variance reduction techniques. This work emphasizes that
196 benchmark outcomes should be treated as data drawn from an underlying population, moving
197 beyond descriptive statistics.

198 HiBayES [31] applies hierarchical Bayesian generalized linear models that: (a) use a
199 Bayesian statistical inference (multilevel Binomial/Poisson Generalized Linear Models) and

⁶Over time, models are trained on similar data, inflating scores and diminishing benchmark usefulness.

200 captures the nested structure of LLM evaluations; (b) employ partial pooling to account for
 201 dependencies across levels; (c) provide a fully probabilistic framework, yielding full posterior
 202 distributions to quantify uncertainties, rather than relying on point-estimates.

203 Despite these advances, important limitations remain. Neither Miller’s frequentist frame-
 204 work nor HiBayES explicitly defines/models the OP (even though both implicitly assume
 205 that the data represent it). Moreover, neither approach can effectively incorporate (poten-
 206 tially imprecise) prior knowledge: Miller’s work is restricted to classical frequentist inference,
 207 while HiBayES relies on non-informative priors. Both also focus on the narrower notion of
 208 failure probability as reliability, overlooking the more general and standardized reliability
 209 definition in terms of the probability of failure-free runs over future operations. The latter
 210 requires predictive inference that accounts for the propagation of uncertainties into the fu-
 211 ture: a more demanding but practical reliability claim [33, 34]. Miller accounts for clustered
 212 dependence through variance corrections, but does not introduce hierarchical latent variables
 213 or multi-level generative models as in HiBayES or HIP-LLM. Table 1 summarizes the key
 214 features and provides a comparative overview of the three methods (HIP-LLM, HiBayES [31],
 215 and Miller [32]).

216 *2.4. Robust Bayesian Analysis*

217 Robust Bayesian analysis represents a general framework for investigating the sensitivity
 218 of posterior measures to uncertainties in the inputs of Bayesian inference [59, 60]. Regarding
 219 uncertainties in priors, while several dedicated methods have been proposed [61, 33, 62, 63],
 220 Imprecise Probability has emerged as one of the most widely adopted approaches [36, 37,
 221 64, 65, 66]. It addresses the problem of prior uncertainty by avoiding reliance on *a single*
 222 prior distribution, rather representing *a credal set* of plausible priors and deriving posterior
 223 bounds that reflect this epistemic uncertainty.

224 To illustrate the main idea of the Imprecise Probability framework, consider a simple
 225 coin-flipping problem where we wish to estimate the probability of heads θ with data D
 226 ($n = 10$, $k = 3$ heads). As shown in Table 2, instead of a single point estimate or a single
 227 Beta posterior distribution, the *posterior envelope* (i.e., a set of posterior distributions) of
 228 Imprecise Probability⁷ reflects both the observed data and the epistemic uncertainty arising
 229 from imprecise prior knowledge (represented by the set of Beta priors).

Table 2: Classical Bayesian vs Imprecise Probability for a coin-flipping example.

Feature	Classical Bayesian	Imprecise Probability
Prior	$\text{Beta}(\alpha = 2, \beta = 2)$	$\text{Beta}(\alpha, \beta), \alpha \in [1, 3], \beta \in [1, 3]$
Posterior	$\theta D \sim \text{Beta}(5, 9)$	$\theta D \sim \text{Beta}(3 + \alpha, 7 + \beta), \alpha \in [1, 3], \beta \in [1, 3]$
Poster. mean	$\mathbb{E}[\theta D] = 0.36$	$\mathbb{E}[\theta D] \in [0.31, 0.38]$

230 Our hierarchical solution, HIP-LLM, is a robust Bayesian approach that addresses the

⁷For simplicity and illustrative purpose, the example used here is based on the model of *imprecise Linearly Updated Conjugate prior Knowledge (iLUCK)* [65]) which leverages the conjugacy for analytical posterior results. Our HIP-LLM is not using this iLUCK model due to the non-linearity of the hierarchical probabilistic model proposed in Section 3.

Table 1: Comparison of LLM evaluation frameworks highlighting differences in reliability definition, operational profile modeling, and uncertainty treatment

Aspect	Miller [32]	HiBayES [31]	HIP-LLM
Statistical framework	Frequentist	Bayesian	Bayesian
Primary goal	Uncertainty-aware benchmark reporting	Uncertainty-aware, hierarchical accuracy estimation	Uncertainty-aware reliability assessment under OP with embedded priors
Metric under estimation	Failure probability (per task)	Failure probability (per task)	Reliability (probability of failure-free future tasks)
OP	Implicit (dataset assumed representative)	Implicit (dataset assumed representative)	Explicit, multi-level OPs
Hierarchical structure	None	LLM, domains, subdomains	LLM, domains, subdomains
Dependence handling	Clustered standard errors	Partial pooling	Partial pooling
Prior specification	None	Uniform precise priors	Informative imprecise priors (credal sets)
Output	Confidence intervals	Single posterior distribution	Posterior envelopes

231 aforementioned gaps. Similar to HiBayES, it adopts a hierarchical structure, but instead of
 232 yielding a single posterior, it reports posterior envelopes over imprecise priors and uncertain
 233 OPs (as variables).

234 3. The Model: HIP-LLM

235 The proposed method HIP-LLM stands for Hierarchical Imprecise Probability for Large
 236 Language Models reliability assessment. Before introducing the framework, we first formally
 237 define the reliability of LLMs.

238 LLMs are software; therefore, to ensure compatibility with the more general and stan-
 239 dardized definition of software reliability (Def. 1), we define LLM reliability as follows:

240 **Definition 3** (LLM reliability). *The probability that an LLM produces failure-free responses*
 241 *over a specified number of future tasks (sequences of closely related queries), under specified*
 242 *(sub-)domains and operational environment.*

243 The definition of LLM reliability closely parallels the standardized definition of soft-
 244 ware reliability, as both emphasize *probabilistic* reasoning and the requirement of *failure-free*
 245 *performance* in future operations. The key distinctions arise from the nature of LLMs as
 246 “on-demand” software. Whereas classical software reliability is typically framed in terms of
 247 “operational time” that covers continuously operated systems like controllers over clock-time,
 248 LLM reliability is defined more explicitly in terms of a specified number of future discrete
 249 tasks (each consisting of sequences of related queries). Furthermore, while both definitions
 250 account for operational conditions, the LLM context requires explicit reference to domains
 251 and subdomains, reflecting the general-purpose and multi-domain design of LLMs. Finally, in
 252 traditional software, failure is typically defined as a deviation from the specification, whereas
 253 there is no specification for LLMs and thus “failures of LLMs”⁸ are often informally/implic-
 254 itly characterized by factual errors (hallucinations) [67, 68] or divergence from human-expert
 255 answers [69, 70].

256 3.1. Problem Formulation

257 While we provide a table of notations in Appendix C, according to Def. 3, LLM reliability
 258 can be formalised as follows:

259 **Definition 4** (Formalised LLM reliability). *Let \mathcal{X} denote the input-space of all possible tasks*
 260 *for a given LLM, and let π be the OP, i.e., a probability distribution over \mathcal{X} that reflects the*
 261 *likelihood of encountering each task $x \in \mathcal{X}$ in practice. Consider a sequence of $n \geq 1$ tasks*
 262 *i.i.d. according to the OP π , let $\mathbb{I}(x_\tau) \in \{0, 1\}$ indicate success (1) or failure (0) on the τ -th*
 263 *task x_τ , then the LLM reliability is:*

$$R(n, \pi) = \Pr_{x_\tau \sim \pi} \left(\bigcap_{\tau=1}^n \{\mathbb{I}(x_\tau) = 1\} \right) \quad (1)$$

264 Intuitively, $R(n, \pi)$ represents the probability that the LLM will operate failure-free across
 265 the next n i.i.d. tasks according to the OP π . Importantly, we make the following remarks:

266 **Remark 1** (The i.i.d. assumption vs. contextual memory). *A key assumption of the above*
 267 *definition is the task failures/successes are i.i.d. Bernoulli trials. Such modeling is not*
 268 *uncommon in reliability modeling, especially for critical on-demand systems. A typical justi-*
 269 *fication is when demands are rare, and the states/memory of the software and its operational*
 270 *environment are effectively “reset” in-between [23, 71]. In the context of LLMs, we care-*
 271 *fully define our reliability metric in terms of i.i.d “tasks”, rather than individual prompts*
 272 *which are often contextually dependent. A task may consist of a sequence of related prompts*
 273 *aimed at achieving a single task goal. We acknowledge that modern LLMs (e.g., ChatGPT)*
 274 *typically retain chat history as contextual memory, which can violate the i.i.d. assumption*
 275 *between tasks. However, most LLMs provide the option to start a new chat session (for a*
 276 *new single task) without any memory and history from previous chat sessions, aligning with*
 277 *the assumption of “resetting” the LLM. We note such “reset” settings are commonly encoun-*
 278 *tered in LLM offline evaluation like benchmarking [30]. In contrast, long-context or agentic*

⁸A complete and formal characterization of what constitutes a “failure” for LLMs remains an open research question, and out of the scope of this paper. Cf. Section 5 for discussions.

workflows, in which tasks are sequentially dependent through memory retention or tool use, violate this assumption. We note, however, that such dependence is not unique to LLMs and has long been recognized in software reliability modeling, particularly in industrial control software, where sensor readings form continuous and dependent trajectories through the input space. In such settings, two common approaches are typically adopted: either modeling the system at a higher level of abstraction by treating an entire trajectory as a single trial⁹, or employing alternative stateful reliability models (which we leave as a future work).

Remark 2 (Failure probability vs. future reliability). While the reliability metric $R(n, \pi)$ represents the probability of correctly processing the next n future tasks under the OP π , the special case $1 - R(1, \pi)$ (where $n = 1$) represents the failure probability [33, 34, 35] that studied by, e.g., [31, 32]. Accordingly, throughout this paper, we use the term failure probability to denote $1 - R(1, \pi)$ and the probability of failure-free runs (of n future tasks) to denote $R(n, \pi)$, in order to distinguish these two notions. The more general term reliability is used to refer to either quantity when the intended meaning is clear from the surrounding context.

Remark 3 (General purpose reliability vs. domain-specific reliability). Since the input-space \mathcal{X} represents all possible LLM tasks, so $R(n, \pi)$ is the general-purpose reliability of the LLM under study. For (sub-)domain specific reliability, we need to partition \mathcal{X} and derive “local” OPs; then the lower level (sub-)domain specific reliability can be similarly derived like Eq. (1).

Remark 4 (Binary failure vs. non-binary scoring). In traditional software reliability engineering, failures are naturally defined as binary events, success or failure, according to an explicit system specification. Similarly, most LLM benchmarks introduce task-specific specifications, implemented via automated evaluators or human annotations, which yield binary outcomes for scoring each prompted task. Consistent with this established practice, we model LLM reliability based on binary failure in Def. 4. However, unlike traditional software, what constitutes a “failure” for an LLM can be domain dependent, and either objective with ground truth answers (for, e.g., coding, math, and factual Q&A) or inherently subjective (e.g., for creative writing). The external scoring mechanisms used to determine LLM failures may themselves be noisy or inconsistent, a challenge that is not unique to LLMs but well known in software reliability modeling as the problem of imperfect test oracles. Accordingly, in HIP-LLM Def. 4, failures are conceptually assumed to be well defined, and the framework is designed to operate correctly under a perfect test-oracle. The implications of imperfect or noisy scoring mechanisms (i.e., the imperfect test-oracle problem) are orthogonal to this formulation and are acknowledged as an important direction for future work.

To do statistical inference for the reliability metric defined in Eq. (1) with assumptions in aforementioned remarks, a simplified “textbook” Bayesian model would be the Beta-Binomial one (which also used as one of the baselines in our experiments). This

⁹A widely adopted abstraction in industrial software reliability (e.g., protection systems in nuclear power plants) is to treat a continuous trajectory of dependent sensor readings as a single unit of processing [23, 28], which may result in either success or failure. Under this abstraction, Bernoulli trials remain appropriate despite strong dependence among inputs along the trajectory. The same rationale applies to LLMs, where an interaction session consisting of multiple dependent prompts/tasks can be treated as a single fallible trial.

316 Beta-Binomial estimator applies to a single domain with precise prior knowledge. Let
 317 $\theta := \Pr_{x \sim \pi}(I(x) = 1)$ denote the (unknown) probability of success on a random task drawn
 318 from the OP π . Given N i.i.d. evaluated tasks with C successes and $N - C$ failures, assume
 319 a Binomial likelihood

$$C, N \mid \theta \sim \text{Binomial}(N, \theta) = \theta^C (1 - \theta)^{N-C},$$

320 and a prior $\Pr(\theta)$. By Bayes' rule, the posterior distribution of θ is

$$\Pr(\theta \mid C, N) = \frac{\theta^C (1 - \theta)^{N-C} \Pr(\theta)}{\int_0^1 \theta^C (1 - \theta)^{N-C} \Pr(\theta) d\theta}. \quad (2)$$

321 Similarly for the future reliability of passing n^F tasks:

$$\Pr(R(n^F, \pi) \mid C, N) = \frac{\theta^{n^F} \theta^C (1 - \theta)^{N-C} \Pr(\theta)}{\int_0^1 \theta^C (1 - \theta)^{N-C} \Pr(\theta)} \quad (3)$$

322 If $\Pr(\theta) = \text{Beta}(\alpha, \beta)$, conjugacy yields the closed-form posterior

$$\theta \mid C, N \sim \text{Beta}(\alpha + C, \beta + N - C).$$

323 For future reliability, thanks to the conjugacy again, the posterior mean reliability for n^F
 324 future tasks is therefore:

$$\mathbb{E}[R(n^F, \pi) \mid \mathcal{D}] = \mathbb{E}[\theta^{n^F} \mid \mathcal{D}] = \frac{B(\alpha + C + n^F, \beta + N - C)}{B(\alpha + C, \beta + N - C)},$$

325 where $B(\cdot, \cdot)$ denotes the Beta function. Similarly, the posterior PDF and CDF can also be
 326 derived and we omit them for brevity.

327 To more rigorously assess the formally defined LLM reliability, coping with the aforemen-
 328 tioned Remarks and Gaps, the next subsection introduces our proposed solution HIP-LLM.
 329 It models the LLM as a hierarchical structure consisting of independent domains, each con-
 330 taining statistically dependent subdomains (cf. Fig. 1).

331 **Remark 5** (The need of hierarchical modeling on (sub-)domain (in-)dependencies). *The*
 332 *modeled dependencies and independencies represent the epistemic structure of our hierarchical*
 333 *Bayesian model. That is, observing failures in one (sub-)domain may or may not update our*
 334 *beliefs about the reliability of other (sub-)domains. One possible example of justification¹⁰ is:*
 335 *We model coding and law as independent because they rely on distinct competencies of the*
 336 *LLM—coding on formal, symbolic reasoning and syntax manipulation, and law on narrative*
 337 *understanding and normative interpretation. Since these skills draw from largely separate*
 338 *representations and training data, failures in one domain provide little information about*
 339 *failures in the other. In Bayesian terms, their failure probabilities can be treated as a priori*
 340 *independent parameters, reflecting separate latent skill dimensions of the model. On the*
 341 *other hand, sub-domains failure probabilities are modeled as dependent parameters given their*
 342 *shared LLM competencies.*

¹⁰Cf. Section 5 for discussions on the validity of this hierarchal dependency assumption.

343 Our goal is to infer the posterior distributions of future reliability (and its special case,
 344 the failure probability), at the subdomain, domain, and general-purpose LLM levels, based
 345 on observed correct responses from tasks within each subdomain.

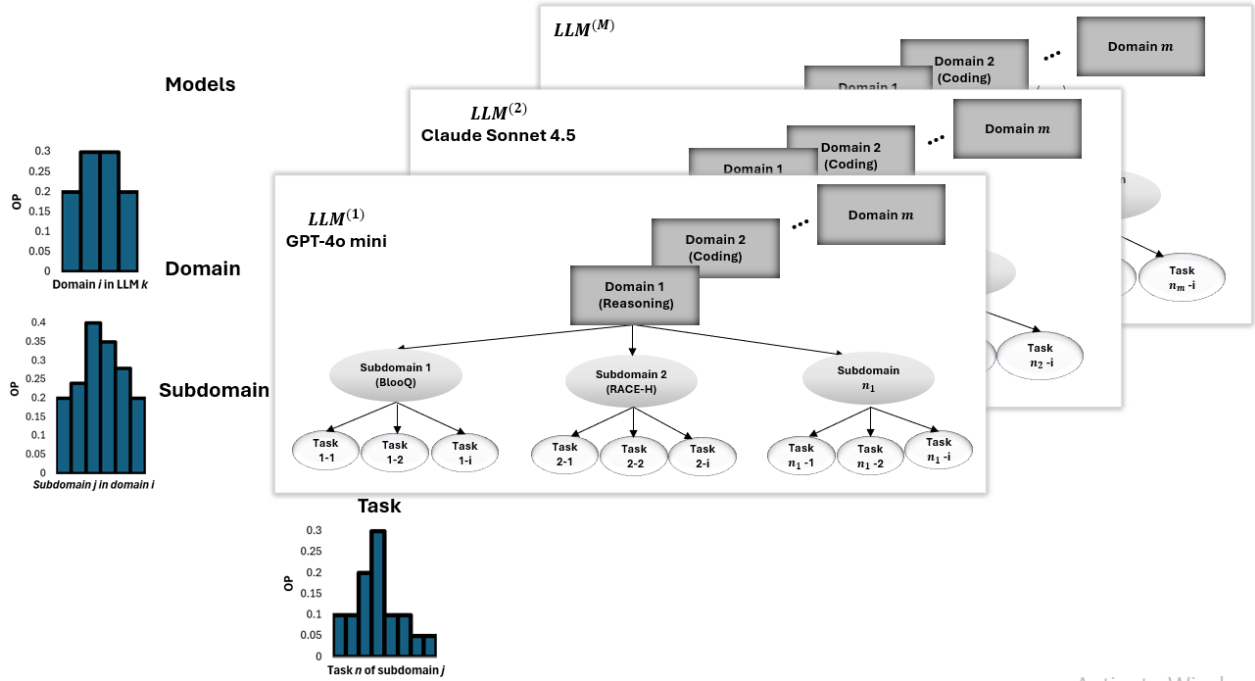


Figure 1: Schematic representation of the hierarchical LLM, domain, and subdomain structure for reliability estimation for M LLM models. Rectangles indicate independent components, while ovals indicate dependent components within the hierarchy.

346 3.2. Proposed Solution

347 Consider a hierarchical structure of an LLM comprising independent domains $D_1, D_2, \dots,$
 348 D_m , where each domain D_i contains statistically dependent subdomains S_{i1}, \dots, S_{in_i} (Fig. 2).
 349 We wish to infer the posterior distributions over subdomain, domain and LLM level reliabil-
 350 ities by observing C_{ij} correct responses out of N_{ij} trials (tasks) in each subdomain.

351 Fig. 2 presents a detailed view of the hierarchical structure (subdomains \rightarrow domains \rightarrow
 352 LLM), with assumed priors and parameters. We assume that subdomain reliabilities within
 353 a domain are *dependent through a shared prior*, and that domain reliabilities are *aggregated*
 354 from their subdomains according to task-specific OPs. Our goal is to construct a principled
 355 hierarchical Bayesian model that supports information sharing across dependent subdomains
 356 through partial pooling (Sec. 3.2.1) and uncertainty quantification via Imprecise Probability
 357 (Sec. 3.2.2).

358 Figures 1 and 2 illustrate a general hierarchical structure comprising multiple LLM instan-
 359 ces ($LLM^{(1)}, LLM^{(2)}, \dots, LLM^{(M)}$). For clarity, however, Theorems 1–6 focus on the
 360 reliability assessment of a *single* LLM system. Accordingly, we omit the superscript (k)
 361 and use unindexed symbols (D_i, S_{ij}, p_L). Extending the framework to multiple LLMs is
 362 straightforward—apply it to each system independently and compare their posterior distri-
 363 butions.

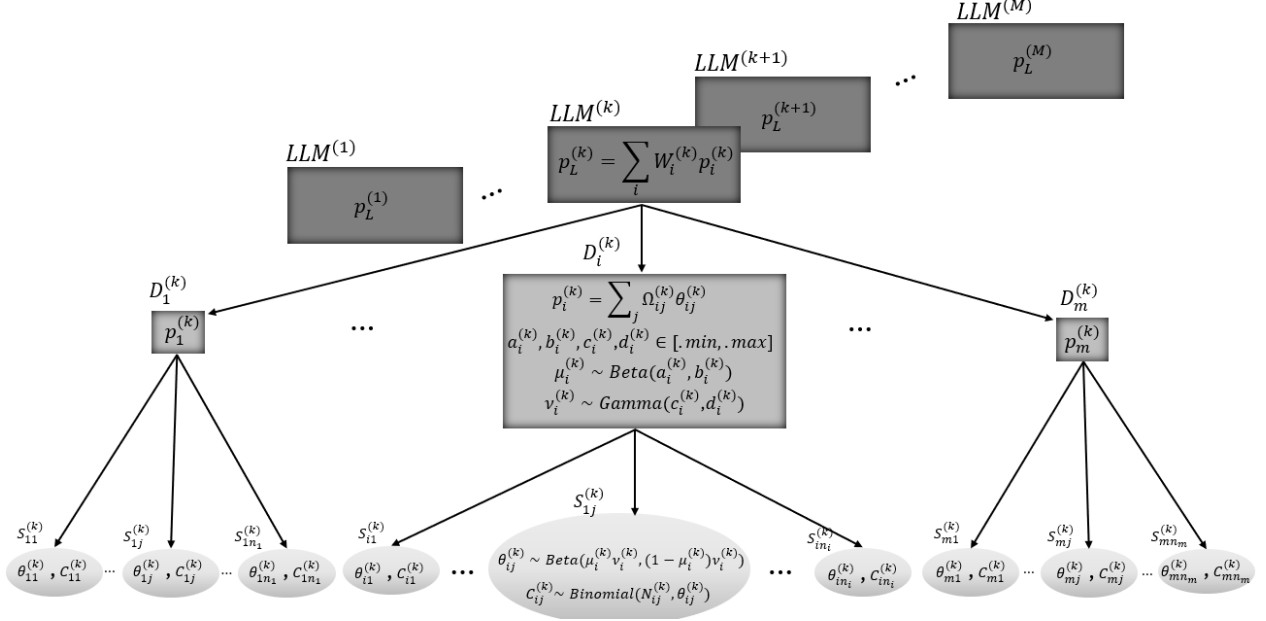


Figure 2: Hierarchical structure with independent domains and dependent subdomains. For readability, parameters and priors are shown only for one subdomain S_{ij} and its parent domain D_i under a representative LLM_k ; the remaining subdomains, domains ($i = 1, \dots, m$), and LLMs ($k = 1, \dots, M$) are identical, differing only in their indices. Rectangles indicate independent components, while ovals indicate dependent components within the hierarchy.

3.2.1. Hierarchical Bayesian Framework for LLM Reliability Modeling

As seen in Fig. 2, inference proceeds hierarchically from the bottom up. At the subdomain level, the observed data consist of the number of correct responses C_{ij} out of total trials N_{ij} . These update the subdomain reliabilities θ_{ij} , which are modeled with a Binomial likelihood¹¹ and a Beta prior¹². The Beta prior is parameterized by domain-level hyperparameters (μ_i, ν_i)¹³, representing the expected reliability and the prior strength within domain D_i . These hyperparameters are in turn governed by domain-specific hyperpriors:

¹¹The assumption of independent Bernoulli trials with constant success probability θ_{ij} may not capture all real dependencies, but it serves as an effective approximation for modeling subdomain outcomes [32, 31, 62]. Aggregating across N_{ij} trials then yields the Binomial likelihood, whose support $\{0, \dots, N_{ij}\}$ matches the possible counts of correct responses observed in subdomain S_{ij} .

$$C_{ij} \mid \theta_{ij}, N_{ij} \sim \text{Binomial}(N_{ij}, \theta_{ij}).$$

¹²We place a Beta prior on subdomain reliability:

$$\theta_{ij} \sim \text{Beta}(\alpha_i, \beta_i).$$

The Beta distribution is the most common choice for probabilities bounded in $[0, 1]$, and it is conjugate to the Binomial likelihood, ensuring closed-form updates and computational stability.

¹³To make priors more intuitive and interpretable, we use a reparameterization. Instead of specifying the Beta prior directly in terms of (α_i, β_i) , we express it as

$$\theta_{ij} \mid \mu_i, \nu_i \sim \text{Beta}(\mu_i \nu_i, (1 - \mu_i) \nu_i),$$

371 $\mu_i \sim \text{Beta}(a_i, b_i)^{14}$, and $\nu_i \sim \text{Gamma}^{15}(c_i, \text{rate}^{16} = d_i)$.

372 Because the hyperparameters (μ_i, ν_i) are shared across all subdomains in domain D_i , the
 373 information is pooled: subdomains with few observations are stabilized by drawing on evi-
 374 dence from other subdomains in the same domain, while subdomains with many observations
 375 are influenced mainly by their own data. This hierarchical setup achieves the desired depen-
 376 dence within domains, while domains remain independent (cf. Section 5 for discussions on
 377 this setup).

378 Once posterior subdomain reliabilities θ_{ij} are inferred, they are aggregated into domain-
 379 level reliabilities $p_i = \sum_j \Omega_{ij} \theta_{ij}$ using OP weights Ω_{ij} , which reflect the practical importance
 380 of subdomains. These domain-level posteriors are then further combined into the overall
 381 LLM reliability, $p_L = \sum_i W_i p_i$, using domain operational weights W_i . Note, OPs at different
 382 levels are represented by the variables Ω_{ij} and W_i , which, for a given LLM use case, can be
 383 instantiated either as fixed constants (when usage is certain) or as probability distributions
 384 that encode uncertainty about how the LLM will be used in practice [72, 73, 74, 75].

385 3.2.2. Uncertainty Handling via Imprecise Probability

386 To address the epistemic uncertainty in prior specification, we adopt an Imprecise Proba-
 387 bility approach: we specify the hyper-hyper-parameters a_i, b_i, c_i, d_i as intervals (Eq. 4) rather
 388 than point values. This produces posterior envelopes (lower and upper bounds) at the sub-
 389 domain, domain, and overall LLM levels.

$$\begin{aligned} a_i &\in [a_i^{\min}, a_i^{\max}], & b_i &\in [b_i^{\min}, b_i^{\max}], \\ c_i &\in [c_i^{\min}, c_i^{\max}], & d_i &\in [d_i^{\min}, d_i^{\max}] \end{aligned} \quad (4)$$

390 At the subdomain level, we compute posterior bounds for each task type by considering
 391 all admissible hyperparameter configurations:

$$\underline{Pr}(\theta_{ij} \mid C_i) \leq Pr(\theta_{ij} \mid C_i) \leq \overline{Pr}(\theta_{ij} \mid C_i) \quad (5)$$

392 Note, C_i denotes the data in domain i . Since we consider dependencies among subdomains,
 393 the posterior θ_{ij} is a function of all data C_i in domain i .

where μ_i denotes the expected reliability (prior mean, $\mu_i = \mathbb{E}[\theta_{ij} \mid \mu_i, \nu_i]$) and ν_i denotes the prior strength, reflecting the confidence in μ_i (equivalent sample size or pseudo-counts, $\nu_i = \alpha_i + \beta_i$). This reparameterization makes prior beliefs easier to specify and justify.

¹⁴This choice reflects the idea that μ_i itself is a probability lying in $(0, 1)$, and the Beta distribution provides a flexible family of shapes that can express different prior beliefs about domain reliability, ranging from diffuse to highly concentrated around particular values (e.g., favoring higher values, lower values, or balanced around 0.5).

¹⁵This treats ν_i as a positive random variable ($\nu_i > 0$) reflecting how tightly subdomains within a domain are assumed to cluster around μ_i . When ν_i is small, the prior is diffuse and allows substantial variation across subdomains (weak pooling). When ν_i is large, the prior concentrates mass near μ_i and subdomains are tightly clustered (strong pooling).

¹⁶The Gamma distribution is commonly written in two equivalent forms: shape–rate and shape–scale. We use the *shape–rate* form, $\nu_i \sim \text{Gamma}(c_i, \text{rate} = d_i)$, with $\mathbb{E}[\nu_i] = c_i/d_i$ and $\text{Var}[\nu_i] = c_i/d_i^2$. If a library expects the *shape–scale* form, set $\theta_i = 1/d_i$ and write $\nu_i \sim \text{Gamma}(c_i, \text{scale} = \theta_i)$; the two parameterizations are mathematically identical under $\theta_i = 1/d_i$. The only practical concern is clarity and reproducibility—accidentally treating a rate as a scale (or vice versa) would change the prior’s mean and variance.

394 Later (Theorem. 2 and Theorem. 3) we will discuss that the closed-form densities (like
 395 subdomain level) do not exist for domain and LLM level reliability. We represent the cumu-
 396 lative distribution function (CDF) envelopes for the domain and overall LLM levels.

397 At the domain level, uncertainty propagates upward through OP weights $p_i = \sum_j \Omega_{ij} \theta_{ij}$,
 398 producing domain reliability bounds:

$$\underline{F}_{p_i}(t | C_i) = \inf_{h_i \in \mathcal{A}_i} F_{p_i}(t | C_i, h_i), \quad \overline{F}_{p_i}(t | C_i) = \sup_{h_i \in \mathcal{A}_i} F_{p_i}(t | C_i, h_i), \quad t \in [0, 1]$$

399 Here and in the following, $t \in [0, 1]$ denotes a generic probability threshold at which the CDF
 400 of the corresponding (non-failure probability or reliability) random variable is evaluated.

401 The domain level posterior p_i only considers C_i , given our assumption on cross-domain
 402 independence.

403 At the LLM level, uncertainty aggregates across all k domains through domain weights
 404 $p_L = \sum_{i=1}^k W_i p_i$, while respecting cross-domain independence, resulting in system-level reli-
 405 ability bounds:

$$\begin{aligned} \underline{F}_{p_L}(t | \text{data}) &= \inf_{\mathcal{H} \in \mathcal{A}_{\text{LLM}}} F_{p_L}(t | \text{data}, \mathcal{H}), \\ \overline{F}_{p_L}(t | \text{data}) &= \sup_{\mathcal{H} \in \mathcal{A}_{\text{LLM}}} F_{p_L}(t | \text{data}, \mathcal{H}) \end{aligned}$$

406 where, $t \in [0, 1]$.

407 While in the next subsection, we develop Theorems 1–3 of deriving posterior sets for
 408 those non-failure probability variables θ_{ij} s, p_i s and p_L , the posterior distribution sets for
 409 future reliability of passing next n^F tasks at each level, e.g., for a domain i :

$$\underline{Pr}\left(\left(\sum_j \Omega_{ij} \theta_{ij}\right)^{n^F} | C_i\right) \leq Pr\left(\left(\sum_j \Omega_{ij} \theta_{ij}\right)^{n^F} | C_i\right) \leq \overline{Pr}\left(\left(\sum_j \Omega_{ij} \theta_{ij}\right)^{n^F} | C_i\right) \quad (6)$$

410 can also be derived, as shown in our Theorems 4–6.

411 3.2.3. Theorems

412 The following theorems are the main mathematical results of HIP-LLM. Intuitively, given
 413 the imprecise prior knowledge encoded by the hyperparameters, the theorems derive the
 414 posterior distributions of non-failure probabilities at different hierarchical levels, as well as the
 415 future reliability, based on the probabilistic reasoning model illustrated in Fig. 2, conditioned
 416 on the observed task failure data.

417 **Theorem 1** (Sub-domain level non-failure probability). *For subdomain S_{ij} in domain D_i ,*
 418 *let $C_i = \{(C_{ik}, N_{ik})\}_{k=1}^{n_i}$ be the observed data. Let the admissible set of hyperparameters be*

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}],$$

419 *and write $h_i = (a_i, b_i, c_i, d_i)$. Then, for any $h_i \in \mathcal{A}_i$, the marginal posterior density of θ_{ij} is*

$$Pr(\theta_{ij} | C_i, h_i) = \frac{f_{\text{marg}}(\theta_{ij}, C_i; h_i)}{Z_{\text{marg}}(h_i)},$$

420 where f_{marg} (unnormalized posterior) and Z_{marg} (normalizing constant) are

$$f_{\text{marg}}(\theta_{ij}, C_i; h_i) = \int_0^1 \int_0^\infty \left[\prod_{k \neq j} \int_0^1 d\theta_{ik} \right] L(\boldsymbol{\theta}_i) Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i) Pr(\mu_i, \nu_i | h_i) d\mu_i d\nu_i$$

421

$$Z_{\text{marg}}(h_i) = \int_0^1 \int_0^\infty Pr(C_i | \mu_i, \nu_i) Pr(\mu_i, \nu_i | h_i) d\mu_i d\nu_i,$$

422 with $L(\boldsymbol{\theta}_i) = Pr(C_i | \boldsymbol{\theta}_i)$.

423 The imprecise marginal posterior is characterized by the lower/upper envelopes

$$\underline{Pr}(\theta_{ij} | C_i) = \inf_{h_i \in \mathcal{A}_i} Pr(\theta_{ij} | C_i, h_i), \quad \overline{Pr}(\theta_{ij} | C_i) = \sup_{h_i \in \mathcal{A}_i} Pr(\theta_{ij} | C_i, h_i).$$

424 The proof of Theorem 1 is presented at Appendix Appendix A.2. In this theorem, the
 425 subdomain posterior $Pr(\theta_{ij} | C_i, h_i)$ has a closed-form density because of conjugacy: the
 426 Beta prior combined with the Binomial likelihood yields a mixture of Beta distributions after
 427 marginalizing over the hyperparameters (μ_i, ν_i) , which can be expressed and evaluated as a
 428 proper probability density function.

429 However, for later Theorems closed-form densities do not exist because $p_i = \sum_j \Omega_{ij} \theta_{ij}$
 430 and $p_L = \sum_i W_i p_i$ are weighted sums of dependent random variables—the distribution of a
 431 sum of Beta random variables has no analytical form except in trivial cases¹⁷. Computing
 432 such densities would require intractable multi-dimensional integrals. The CDF formulation
 433 sidesteps this problem: $F_{p_i}(t | C_i, h_i) = \int_0^1 \int_0^\infty F_{p_i}(t | \mu_i, \nu_i, C_i) Pr(\mu_i, \nu_i | C_i, h_i) d\mu_i d\nu_i$
 434 only requires a two-dimensional integral over (μ_i, ν_i) , where the conditional CDF can be
 435 computed via Monte Carlo sampling of independent Betas. Since CDFs provide all necessary
 436 information for practical reliability assessment (probabilities, quantiles, expectations), they
 437 are the natural representation when densities are unavailable.

438 **Theorem 2** (Domain level posterior non-failure probability). *For domain D_i with local OP*
 439 *weights Ω_{ij} (where $\sum_{j=1}^{n_i} \Omega_{ij} = 1$), let $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$ be the domain-level non-failure prob-*
 440 *ability. Define the admissible set of hyper-hyper-parameters*

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}],$$

441 and write $h_i = (a_i, b_i, c_i, d_i)$.

442 Then, for any $h_i \in \mathcal{A}_i$, the posterior distribution of p_i is characterized by its CDF:

$$\begin{aligned} F_{p_i}(t | C_i, h_i) &= \Pr(p_i \leq t | C_i, h_i) \\ &= \int_0^1 \int_0^\infty F_{p_i}(t | \mu_i, \nu_i, C_i) Pr(\mu_i, \nu_i | C_i, h_i) d\mu_i d\nu_i \end{aligned}$$

443 where

¹⁷Assuming the weights are constants. When the weights are modeled as random variables with their own probability distributions, the same problem persists (if not harder).

- 444 • $F_{p_i}(t \mid \mu_i, \nu_i, C_i)$ is the conditional CDF of $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$ given that $\theta_{ij} \mid \mu_i, \nu_i, C_i \stackrel{ind.}{\sim}$
 445 $Beta(C_{ij} + \mu_i \nu_i, N_{ij} - C_{ij} + (1 - \mu_i) \nu_i)$ for $j = 1, \dots, n_i$,
- 446 • $Pr(\mu_i, \nu_i \mid C_i, h_i)$ is the hyper-posterior obtained via Bayes' rule:

$$Pr(\mu_i, \nu_i \mid C_i, h_i) = \frac{Pr(C_i \mid \mu_i, \nu_i) Beta(\mu_i \mid a_i, b_i) Gamma(\nu_i \mid c_i, rate = d_i)}{\int_0^1 \int_0^\infty Pr(C_i \mid \mu, \nu) Beta(\mu \mid a_i, b_i) Gamma(\nu \mid c_i, rate = d_i) d\mu d\nu}$$

447 The imprecise domain posterior is characterized by CDF envelopes:

$$\underline{F}_{p_i}(t \mid C_i) = \inf_{h_i \in \mathcal{A}_i} F_{p_i}(t \mid C_i, h_i), \quad \overline{F}_{p_i}(t \mid C_i) = \sup_{h_i \in \mathcal{A}_i} F_{p_i}(t \mid C_i, h_i).$$

448 For the proof details, we refer reads to Appendix A.3.

449 **Theorem 3** (LLM-level posterior non-failure probability). *For the LLM system with domain*
 450 *weights W_i (where $\sum_{i=1}^m W_i = 1$), let $p_L = \sum_{i=1}^m W_i p_i$ be the LLM-level failure probability and*
 451 *data = $\{C_1, \dots, C_m\}$ the observed data across all domains. Assume cross-domain indepen-*
 452 *dence.*

453 Define the domain-level admissible sets

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}]$$

454 and write $h_i = (a_i, b_i, c_i, d_i)$ for $i = 1, \dots, m$. Define the LLM-level admissible set as the
 455 Cartesian product

$$\mathcal{A}_{LLM} = \mathcal{A}_1 \times \dots \times \mathcal{A}_m,$$

456 and collect the domain hyperparameters as $\mathcal{H} = (h_1, \dots, h_m) \in \mathcal{A}_{LLM}$.

457 Then, for any $\mathcal{H} \in \mathcal{A}_{LLM}$, the posterior distribution of p_L is characterized by its CDF:

$$F_{p_L}(t \mid data, \mathcal{H}) = Pr(p_L \leq t \mid data, \mathcal{H}) = \int \dots \int G(t \mid \{\mu_i, \nu_i\}_{i=1}^m, data) \prod_{i=1}^m Pr(\mu_i, \nu_i \mid C_i, h_i) \prod_{i=1}^m d\mu_i d\nu_i$$

458 where

- 459 • $G(t \mid \{\mu_i, \nu_i\}_{i=1}^m, data)$ is the conditional CDF of $p_L = \sum_{i=1}^m W_i p_i$ given all hyperparam-
 460 eters, defined as

$$G(t \mid \{\mu_i, \nu_i\}_{i=1}^m, data) = \int_{\mathcal{R}_L(t)} \prod_{i=1}^m f_{p_i}(p_i \mid \mu_i, \nu_i, C_i) dp_1 \dots dp_m,$$

461 where $\mathcal{R}_L(t) := \{(p_1, \dots, p_m) \in (0, 1)^m : \sum_{i=1}^m W_i p_i \leq t\}$, and $f_{p_i}(\cdot \mid \mu_i, \nu_i, C_i)$ is the
 462 conditional density of $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$ under $\theta_{ij} \mid \mu_i, \nu_i, C_i \stackrel{ind.}{\sim} Beta(C_{ij} + \mu_i \nu_i, N_{ij} -$
 463 $C_{ij} + (1 - \mu_i) \nu_i)$,

464 • $Pr(\mu_i, \nu_i | C_i, h_i)$ is the domain-level hyper-posterior for domain i :

$$Pr(\mu_i, \nu_i | C_i, h_i) = \frac{Pr(C_i | \mu_i, \nu_i) \text{Beta}(\mu_i | a_i, b_i) \text{Gamma}(\nu_i | c_i, \text{rate} = d_i)}{\int_0^1 \int_0^\infty Pr(C_i | \mu, \nu) \text{Beta}(\mu | a_i, b_i) \text{Gamma}(\nu | c_i, \text{rate} = d_i) d\mu d\nu}$$

465 • Cross-domain independence ensures $Pr(\{\mu_i, \nu_i\}_{i=1}^m | \text{data}, \mathcal{H}) = \prod_{i=1}^m Pr(\mu_i, \nu_i | C_i, h_i)$.

466 The imprecise LLM posterior is characterized by CDF envelopes:

$$\begin{aligned} \underline{F}_{p_L}(t | \text{data}) &= \inf_{\mathcal{H} \in \mathcal{A}_{LLM}} F_{p_L}(t | \text{data}, \mathcal{H}) \\ \overline{F}_{p_L}(t | \text{data}) &= \sup_{\mathcal{H} \in \mathcal{A}_{LLM}} F_{p_L}(t | \text{data}, \mathcal{H}) \end{aligned}$$

467 For more details, see Section Appendix A.4.

468 Theorems 1, 2, and 3 characterize the posterior distribution of non-failure probability θ_{ij}
 469 or aggregated p_i, p_L . However, in practice, we often care about the reliability over a *specified*
 470 *number of consecutive future operations*, e.g., “what is the probability that an LLM succeeds
 471 on the next 10 tasks in a row?” or “what is the probability that the LLM succeeds on the
 472 next 20 coding tasks?”.

473 The following set of theorems extends the aforementioned theorems to characterize the full
 474 posterior distribution of reliability for n^F consecutive future operations, i.e., the probability
 475 that the LLM operates n^F consecutive failure-free tasks. Conditioning is on the observed
 476 evaluation data across subdomains S_{ij} : C_{ij} correct generations out of N_{ij} prompts.

477 **Theorem 4** (Subdomain posterior reliability for n^F future operations). *For subdomain S_{ij}*
 478 *in domain D_i , let $C_i = \{(C_{ik}, N_{ik})\}_{k=1}^{n_i}$ denote all observed data in the domain. Let the*
 479 *admissible set of hyperparameters be*

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}],$$

480 and write $h_i = (a_i, b_i, c_i, d_i)$. Define the reliability random variable¹⁸ $R_{ij}(n^F) = \theta_{ij}^{n^F}$, repre-
 481 senting the probability of n^F consecutive failure-free operations in subdomain S_{ij} .

482 For any $h_i \in \mathcal{A}_i$, the posterior CDF of $R_{ij}(n^F)$ is:

$$F_{R_{ij}(n^F)}(t | C_i, h_i) = \Pr(\theta_{ij}^{n^F} \leq t | C_i, h_i) = \int_0^{t^{1/n^F}} Pr(\theta_{ij} | C_i, h_i) d\theta_{ij},$$

483 where $Pr(\theta_{ij} | C_i, h_i)$ is the marginal posterior density from Theorem 1.

484 The imprecise posterior distribution is characterized by the CDF envelopes:

$$\begin{aligned} \underline{F}_{R_{ij}(n^F)}(t | C_i) &= \inf_{h_i \in \mathcal{A}_i} F_{R_{ij}(n^F)}(t | C_i, h_i), \\ \overline{F}_{R_{ij}(n^F)}(t | C_i) &= \sup_{h_i \in \mathcal{A}_i} F_{R_{ij}(n^F)}(t | C_i, h_i) \end{aligned}$$

¹⁸Comparing to the reliability definition in Def. 4, the OPs are omitted as we assume they are fixed constants in these theorems.

486 The transformation $\theta_{ij} \mapsto \theta_{ij}^{n^F}$ generally does not yield a closed-form density. We therefore
 487 characterize $R_{ij}(n^F)$ through its CDF. See Appendix A.5 computational methods.

488 **Theorem 5** (Domain posterior reliability for n^F future operations). *For domain D_i with*
 489 *local OP weights Ω_{ij} (where $\sum_{j=1}^{n_i} \Omega_{ij} = 1$), let $C_i = \{(C_{ik}, N_{ik})\}_{k=1}^{n_i}$ denote all observed data*
 490 *in the domain. Let the admissible set of hyperparameters be*

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}],$$

491 and write $h_i = (a_i, b_i, c_i, d_i)$. Define the domain-level reliability

$$p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}, \quad R_i(n^F) = p_i^{n^F},$$

492 representing the probability of n^F consecutive failure-free operations at the domain level.

493 For any $h_i \in \mathcal{A}_i$, the posterior distribution of $R_i(n^F)$ is characterized by its CDF:

$$F_{R_i(n^F)}(t \mid C_i, h_i) = \Pr \left(\left[\sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij} \right]^{n^F} \leq t \mid C_i, h_i \right),$$

494 computed by integrating over the joint posterior $Pr(\boldsymbol{\theta}_i \mid C_i, h_i)$.

495 The imprecise posterior distribution is characterized by the CDF envelopes:

$$\underline{F}_{R_i(n^F)}(t \mid C_i) = \inf_{h_i \in \mathcal{A}_i} F_{R_i(n^F)}(t \mid C_i, h_i),$$

496

$$\overline{F}_{R_i(n^F)}(t \mid C_i) = \sup_{h_i \in \mathcal{A}_i} F_{R_i(n^F)}(t \mid C_i, h_i)$$

497 The CDF can be computed via numerical integration over the joint posterior $Pr(\boldsymbol{\theta}_i \mid$
 498 $C_i, h_i)$ and Monte Carlo sampling, cf. Appendix A.5.

499 **Theorem 6** (LLM posterior reliability for n^F future operations). *For the LLM system with*
 500 *an OP of domain weights W_i (where $\sum_{i=1}^m W_i = 1$), let data = $\{C_1, \dots, C_m\}$ denote all*
 501 *observed data across domains. Assume cross-domain independence.*

502 Define the domain-level admissible sets

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}]$$

503 and write $h_i = (a_i, b_i, c_i, d_i)$ for $i = 1, \dots, m$. Define the LLM-level admissible set as the
 504 Cartesian product

$$\mathcal{A}_{LLM} = \mathcal{A}_1 \times \dots \times \mathcal{A}_m,$$

505 and collect the domain hyperparameters as $\mathcal{H} = (h_1, \dots, h_m) \in \mathcal{A}_{LLM}$.

506 Define the LLM-level reliability

$$p_L = \sum_{i=1}^m W_i p_i, \quad R_L(n^F) = p_L^{n^F},$$

507 representing the probability of n^F consecutive failure-free operations at the LLM level.

508 For any $\mathcal{H} \in \mathcal{A}_{LLM}$, the posterior distribution of $R_L(n^F)$ is characterized by its CDF:

$$F_{R_L(n^F)}(t \mid \text{data}, \mathcal{H}) = \Pr \left(\left[\sum_{i=1}^m W_i p_i \right]^{n^F} \leq t \mid \text{data}, \mathcal{H} \right),$$

509 computed by integrating over $\prod_{i=1}^m \Pr(p_i \mid C_i, h_i)$.

510 The imprecise posterior distribution is characterized by the CDF envelopes:

$$\underline{F}_{R_L(n^F)}(t \mid \text{data}) = \inf_{\mathcal{H} \in \mathcal{A}_{LLM}} F_{R_L(n^F)}(t \mid \text{data}, \mathcal{H}),$$

$$\overline{F}_{R_L(n^F)}(t \mid \text{data}) = \sup_{\mathcal{H} \in \mathcal{A}_{LLM}} F_{R_L(n^F)}(t \mid \text{data}, \mathcal{H})$$

512 Again, while the closed-form density is not available, the CDF can be computed via
513 numerical integration and Monte Carlo sampling. See Appendix A.5 for details.

514 **Remark 6** (Translate vague prior knowledge to hyperparameters). *In HIP-LLM, hyper-*
515 *parameters are not tuning constants but variables used to represent uncertainty about prior*
516 *knowledge. At the lowest level, θ_{ij} indicates “how reliable” (i.e., the probability that the model*
517 *succeeds on a randomly drawn task from the subdomain) of S_{ij} . Subdomains within the same*
518 *domain D_i are assumed to be related, and this dependence is captured by assuming that their*
519 *non-failure probabilities are drawn from a shared Beta distribution governed by two domain-*
520 *level parameters: μ_i and ν_i . Here, μ_i represents the expected reliability of the LLM in domain*
521 *D_i , while ν_i represents how confident we are in this expectation, or equivalently, how strongly*
522 *subdomain reliabilities are expected to cluster around μ_i . Since assessors typically do not*
523 *know the exact values of μ_i and ν_i , HIP-LLM places distributions on them: $\mu_i \sim \text{Beta}(a_i, b_i)$*
524 *and $\nu_i \sim \text{Gamma}(c_i, d_i)$. The hyperparameters (a_i, b_i, c_i, d_i) therefore encode high-level and*
525 *(likely) vague/imperfect expert beliefs rather than performance data.*

526 To ensure that these hyperparameters correspond consistently to the intended prior beliefs,
527 it is convenient to express (a_i, b_i) and (c_i, d_i) through an auxiliary mean–strength parameter-
528 ization. Let: **a**) $r_i \in (0, 1)$ denote the assessor’s belief about the expected reliability in domain
529 D_i , and $s_i > 0$ denote the associated confidence level (interpretable as an equivalent number
530 of prior observations); **b**) $e_i > 0$ denote the assessor’s expected clustering strength (i.e., how
531 strongly subdomain reliabilities are expected to concentrate around the domain mean), and
532 $f_i > 0$ represent the associated confidence in this expectation.

533 The Beta parameters can then be written as follows:

$$a_i = r_i s_i, \quad b_i = (1 - r_i) s_i \tag{7}$$

534 Under this parameterization the prior mean is

$$\mathbb{E}[\mu_i] = \frac{a_i}{a_i + b_i} = r_i,$$

535 while the parameter $s_i = a_i + b_i$ controls the strength (concentration) of the prior. Similarly,
536 the Gamma parameters can be written as

$$c_i = f_i, \quad d_i = \frac{f_i}{e_i}. \tag{8}$$

537 Under this reparameterization, the prior mean of the clustering strength becomes

$$\mathbb{E}[\nu_i] = \frac{c_i}{d_i} = e_i,$$

538 while the parameter f_i controls the concentration of the Gamma prior around e_i . Larger
539 values of f_i correspond to stronger confidence in the expected clustering strength.

540 For example, suppose an assessor believes that the reliability of the LLM in the reasoning
541 domain is likely between 0.8 and 0.9, but with only moderate confidence. This belief can be
542 expressed by specifying $r_i \in [0.8, 0.9]$ and $s_i \in [30, 100]$. The corresponding Beta parameters
543 therefore lie in the ranges $a_i \in [24, 90]$ and $b_i \in [3, 20]$. At the same time, the assessor
544 may believe that the reliabilities of different reasoning subdomains (e.g., logical reasoning,
545 multi-step reasoning etc.) should not differ drastically, but may still vary to some extent.
546 This belief can be expressed by specifying $e_i \in [30, 100]$ and $f_i \in [3, 10]$ which determine
547 the corresponding Gamma hyperparameters through $c_i = f_i$ and $d_i = f_i/e_i$ (i.e., $c_i \in [3, 10]$
548 and $d_i \in [0.03, 0.33]$). Together, the hyperparameters (a_i, b_i, c_i, d_i) therefore encode the asses-
549 sor’s belief that the overall reliability of the domain is high (around 0.8–0.9), while allowing
550 moderate variation between individual subdomains.

551 It is important to note that the parameters s_i and f_i represent different levels of epistemic
552 confidence in the hierarchical model. The parameter s_i reflects the assessor’s confidence in the
553 expected domain reliability μ_i , whereas ν_i controls how similar the reliabilities of individual
554 subdomains are expected to be. Consequently, the model separates two distinct aspects of
555 prior knowledge: beliefs about the overall reliability of the domain and beliefs about the degree
556 of similarity among its subdomains. In this representation, r_i and e_i denote the assessor’s
557 expected values for the domain reliability and the clustering strength, respectively, while s_i
558 and f_i encode the associated confidence levels. This hierarchical structure therefore allows
559 assessors to express both their expectations and their uncertainty about these quantities in a
560 transparent manner.

561 Allowing these quantities to vary within intervals rather than fixing them to single val-
562 ues enables HIP-LLM to represent epistemic uncertainty about prior beliefs following the
563 imprecise probability framework [36, 37]. This uncertainty is then propagated through the
564 hierarchical model, resulting in posterior reliability bounds rather than overconfident point
565 estimates.

566 4. Evaluation

567 To demonstrate and evaluate our HIP-LLM, we empirically investigate five research ques-
568 tions (RQs) in this section.

569 4.1. Research Questions

570 **RQ1 (Effectiveness):** How effectively can HIP-LLM assess and compare posterior reli-
571 ability distributions across different levels of the hierarchy, considering uncertainties prop-
572 agated from subdomains to domains, and finally to general-purpose LLMs? In this RQ, we
573 aim to demonstrate the use case of our HIP-LLM as a reliability assessment tool.

574 **RQ2 (Sensitivity to hyperparameters):** How sensitive are the posterior estimates
575 of HIP-LLM to the hyperparameters (a_i, b_i, c_i, d_i) ? These hyperparameters represent the

576 assessors’ (imprecise) prior knowledge, thus understanding their sensitivity to the posteriors
577 may provide insights on how prior knowledge can be elicited.

578 **RQ3 (Sensitivity to OPs):** How sensitive are the posterior reliability estimates of
579 HIP-LLM to variations in the OPs that characterize operational usage of LLMs at various
580 levels? The delivered and perceived reliability of an LLM depends on how it will be used,
581 i.e. the OP. We hypothesize that general-purpose LLMs may exhibit lower sensitivity to
582 variations in OPs, whereas LLMs trained for specific (sub-)domains are likely to be more
583 sensitive. We investigate and demonstrate how HIP-LLM can characterize and quantify such
584 OP-dependent reliability variations.

585 **RQ4 (Predictability):** How can HIP-LLM predicts future reliability of passing next
586 n^F tasks? While RQ1–RQ3 focus on failure probabilities, which is a special case of reli-
587 ability, we additionally aim to demonstrate HIP-LLM’s capability to predict the probability
588 of successfully completing n^F future tasks and to quantify how this reliability varies as the
589 reliability requirement n^F changes.

590 **RQ5 (Comparison to baselines):** How does HIP-LLM compare to established and
591 state-of-the-art Bayesian reliability estimators? To compare the accuracy of different models
592 against a “ground truth” reliability, we conduct synthetic simulation experiments where the
593 “ground truth” OP and failure probabilities of sub-domains are assumed known.

594 **RQ6 (Failure definitions):** How can HIP-LLM cope with alternative definitions of task
595 success, such as different pass@k criteria? As per Remark 4, the definition of failures of a
596 LLM may vary and subject to uncertainties. While we do not formally study LLM failure
597 definitions, we demonstrate how different failure definitions can be incorporated in HIP-LLM.

598 **RQ7 (Robustness to memory effects):** How robust are HIP-LLM’s posterior esti-
599 mates to violations of the i.i.d. assumption caused by memory-induced dependence during
600 LLM evaluation? While HIP-LLM is designed for reset, single-task scenarios, real-world
601 evaluation data may not be strictly generated under this assumption. It is therefore criti-
602 cal to assess whether its posterior estimates remain meaningful when independence is only
603 approximately satisfied or partially violated by memory-induced dependencies.

604 **RQ8 (Scalability):** How does the computational time and hardware RAM cost of HIP-
605 LLM scale with the number of domains, subdomains, hyperparameter configurations, and
606 Monte Carlo samples? Understanding the scalability of HIP-LLM with these key parameters
607 is essential to determine if it can be applied to large-scale evaluations in practice.

608 4.2. Experimental Setup

609 *OPs.* To emphasize the role of OPs, we “simulate” the OPs by assigning probability distri-
610 butions over these datasets for sampling and by specifying operational weights across (sub-
611)domains. Without loss of generality, and consistent with our three-level hierarchy, we define
612 the task-level OP as a uniform distribution and assign operational weights at the subdo-
613 main and domain levels proportionally to their dataset sizes. That said, assessors may use
614 alternative distributions as the 3-level OPs when additional information is available, such
615 as those approximated from historical usage data [24, 76] or user behavior reports [27]. As
616 demonstrated in RQ3, where we vary the operational weights.

617 **Remark 7** (Simulated OP from benchmarks vs. real-world OP). *Similar to many software*
618 *reliability modeling studies [28], HIP-LLM assumes that the OP is specified and reflects real*

619 *usage conditions. In the absence of real operational data, and following the same experimental*
620 *practice as [31, 32], we leverage existing benchmark datasets to “simulate” OPs for fair*
621 *comparison with baselines. Specifically, benchmark datasets are treated as sampling frames,*
622 *and proportional dataset sizes are adopted as a proxy OP. This choice neither implies that OP*
623 *acquisition is solved nor suggests that such proxies are suitable for real deployments; rather,*
624 *it serves to demonstrate how HIP-LLM integrates OPs when they are provided and to enable*
625 *comparison with existing benchmark-based methods.*

626 *Estimating realistic and evolving OPs from usage data is a well-studied problem in software*
627 *reliability and is usually treated separately from reliability modeling [77, 25]. In practice, an*
628 *OP for an LLM system can be constructed using a simple workflow based on real usage data:*

- 629 1. **Log collection:** *Collect interaction logs or telemetry from deployed systems containing*
630 *user prompts or tasks.*
- 631 2. **Task mapping:** *Assign each task to a domain or subdomain (e.g., coding, reasoning,*
632 *summarization) using predefined rules or automated classifiers.*
- 633 3. **Frequency estimation:** *Count how often each task category occurs to estimate em-*
634 *pirical frequencies.*
- 635 4. **OP construction and update:** *Convert these frequencies into probabilities that de-*
636 *fine the OP, and periodically update the OP as new logs become available.*

637 *While we believe existing OP estimation techniques may be applicable to LLMs, deriving OPs*
638 *from real-world LLM usage data introduces new challenges and warrants dedicated future*
639 *investigation.*

640 *Data.* Specifically, we evaluate our hierarchical framework by simulating LLM operational
641 data from four widely used benchmarks, structured into two domains with two subdomains
642 each (same as HiBayEs [31]), instantiated by the following datasets:

643 • **Domain 1 (Coding)**

- 644 – **Subdomain₁₁:** basic Python programming tasks with unit-test based evaluation
645 (MBPP [9]).
- 646 – **Subdomain₁₂:** data-science oriented Python problems involving libraries such as
647 `pandas` and `numpy` (DS-1000 [8]).

648 • **Domain 2 (Reasoning):**

- 649 – **Subdomain₂₁:** reading comprehension problems where the LLM must answer
650 yes/no questions given short passages (BoolQ [78]).
- 651 – **Subdomain₂₂:** high school level reading comprehension problems with multiple-
652 choice answers (RACE-H [44]).

653 All experiments were conducted using publicly available APIs of GPT-4o, GPT-4o-mini (from
654 OpenAI), Claude Sonnet 4.5, and Claude Haiku 3.5 (from Anthropic)¹⁹. Table 3 reports

¹⁹For academic research purposes only which is permitted under both vendors’ terms of service for research publications.

655 sample accuracies per subdomain, i.e., the proportions of correct responses under Pass@1²⁰.
 656 To operationalize the i.i.d. Bernoulli trial assumption on tasks, all LLM evaluations were
 657 conducted with cleared context between tasks. Each task in MBPP, DS-1000, BoolQ, and
 658 RACE-H was processed as an independent API call with no conversational history, ensuring
 659 that outcomes are not influenced by previous interactions.

Table 3: Evaluation results using Pass@1. Entries are per-subdomain accuracies (C_{ij}/N_{ij}). Rightmost column shows the row mean across models.

Dom	Subdom (Data.)	4o-mini	4o	sonnet-4.5	haiku-3.5	Mean
Dom ₁	Subdom ₁₁ (MBPP)	0.440	0.471	0.450	0.447	0.452
	Subdom ₁₂ (DS-1000)	0.490	0.420	0.493	0.483	0.472
	Mean	0.465	0.446	0.472	0.465	–
Dom ₂	Subdom ₂₁ (BoolQ)	0.890	0.909	0.900	0.883	0.896
	Subdom ₂₂ (RACE-H)	0.820	0.552	0.840	0.859	0.768
	Mean	0.855	0.731	0.87	0.871	–
LLM Mean		0.661	0.585	0.671	0.668	–

660 From the benchmark accuracy scores, all LLMs appear to perform better in the Reason-
 661 ing domain than in the Coding domain. Their performances across most subdomains are
 662 similar, except for RACE-H, where model OpenAI-4o shows a clear weakness. The small
 663 variance between accuracy values for corresponding subdomains across models suggests that
 664 the performance differences are more domain-driven than model-driven.

665 *Hyperparameters.* All the hyperparameters we are using for generating the set of figures in
 666 the next subsection is shown in Appendix B.

667 *Interpretation of Posterior CDF Envelopes.* CDF envelopes are used to jointly represent
 668 data uncertainty and epistemic uncertainty arising from imprecise prior knowledge. Since
 669 closed-form densities are generally unavailable for weighted sums of dependent reliability
 670 parameters, posterior uncertainty is characterized numerically through CDF envelopes, which
 671 support probabilistic queries and risk-aware interpretation across all levels of the hierarchy.
 672 Unless stated otherwise, all figures in this section follow this convention.

673 Throughout the experimental section, where appropriate, we visualize posterior uncer-
 674 tainty using CDF envelopes. These CDFs are constructed from observed benchmark eval-
 675 uation data, consisting of C_{ij} correct responses out of N_{ij} tasks for each subdomain S_{ij} .
 676 Depending on the level of aggregation shown, a CDF represents the posterior distribution
 677 of a non-failure probability at the subdomain level (θ_{ij}), the domain level (p_i), or the over-
 678 all LLM level (p_L). Domain- and LLM-level CDFs are obtained by aggregating lower-level
 679 reliability parameters using the specified OP weights.

680 Monte Carlo sampling is used because the hierarchical aggregation of dependent subdo-
 681 main reliabilities and imprecise priors leads to posterior distributions that are analytically
 682 intractable but can be evaluated efficiently and accurately via numerical sampling.

²⁰Produce 1 solution per task. Score as correct if that solution passes verification.

683 *Baseline Experimental Configuration.* Unless otherwise stated, all empirical experiments are
684 conducted using a fixed baseline configuration with $m = 2$ domains, $n = 2$ subdomains per
685 domain, $K = 160$ hyperparameter configurations per domain, $S = 3000$ Monte Carlo samples
686 per configuration, a (μ, ν) integration grid of size $G = 2000$, a fixed CDF evaluation grid
687 of size $T = 201$, and a capped number of LLM-level configuration pairings $K_{\text{total}} \leq 512$.
688 Domain-level parameters μ_i and ν_i are treated as latent variables and numerically integrated
689 over a grid, while the hyper-hyperparameters (a_i, b_i, c_i, d_i) are sampled from fixed intervals,
690 with $a_i, b_i \in [1, 12]$ and $c_i, d_i \in [1, 25]$, to represent imprecise prior knowledge.

691 All experiments were implemented in Python 3.12.12 using NumPy 2.0.2, pandas 2.2.2,
692 and matplotlib 3.10.0. The experiments were executed in the Google Colab environment on
693 a CPU runtime with an Intel(R) Xeon(R) CPU @ 2.20GHz and approximately 13 GB RAM.
694 Experiments were performed in a single-process setting without explicit parallelization or
695 GPU acceleration.

696 4.3. Results and Analysis

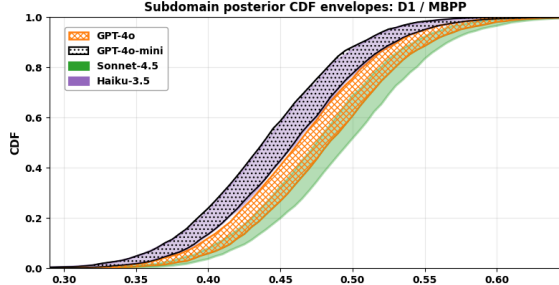
697 This section reports and analyses the empirical results to the RQs.

698 4.3.1. **RQ1**(Effectiveness)

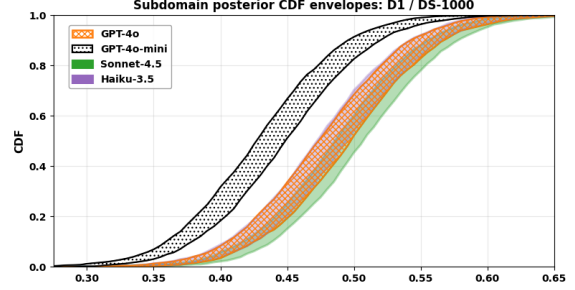
699 Fig. 3 shows posterior CDF envelopes for the four subdomains (MBPP, DS-1000, BoolQ,
700 RACE-H) across the four LLMs where a right-shifted CDF indicates higher reliability and
701 a tighter band indicates greater certainty.

702 *Coding tasks.* On MBPP (top left), the envelopes for GPT-4o-mini and Haiku 3.5 are nearly
703 indistinguishable and lie to the left of GPT-4o, while Sonnet 4.5 is right-most but still
704 overlapping with-4o. On DS-1000 (top right), GPT-4o-mini lies clearly to the left, while
705 Haiku 3.5 overlaps almost entirely with GPT-4o—their envelopes coincide so closely that
706 Haiku 3.5 is barely distinguishable. Sonnet 4.5 remains right-most, with only partial overlap
707 with GPT-4o and Haiku 3.5.

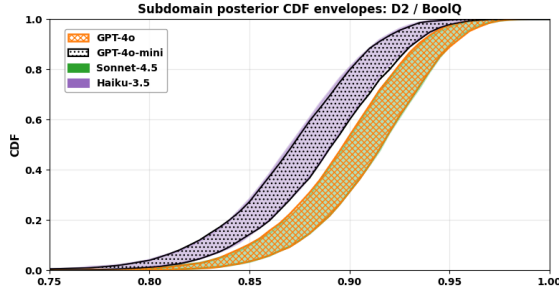
708 *Reasoning tasks.* On BoolQ (bottom left), GPT-4o-mini is almost identical to Haiku 3.5 on
709 the left, while Sonnet 4.5 and GPT-4o nearly coincide on the right. Hence the envelopes
710 separate into two close pairs: mini \approx Haiku 3.5 < Sonnet 4.5 \approx 4o. On RACE-H (bottom
711 right), Haiku 3.5 lies clearly left-most, GPT-4o-mini overlaps partly with Haiku but extends
712 rightward, followed by GPT-4o, and finally Sonnet 4.5 on the far right.



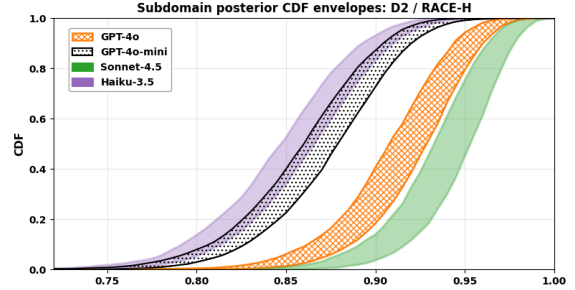
(a) Posterior CDF envelopes of the non-failure probability for Subdom₁₁ (MBPP dataset in Dom₁) across four models.



(b) Posterior CDF envelopes of the non-failure probability for Subdom₁₂ (DS-1000 dataset in Dom₁) across four models.



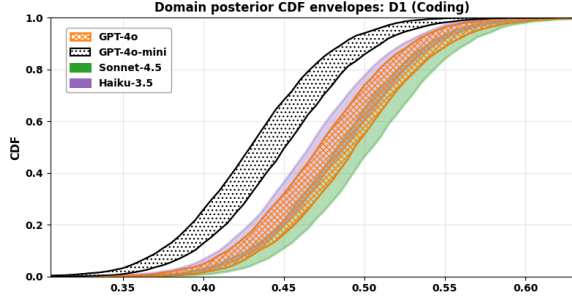
(c) Posterior CDF envelopes of the non-failure probability for Subdom₂₁ (BoolQ dataset in Dom₂) across four models.



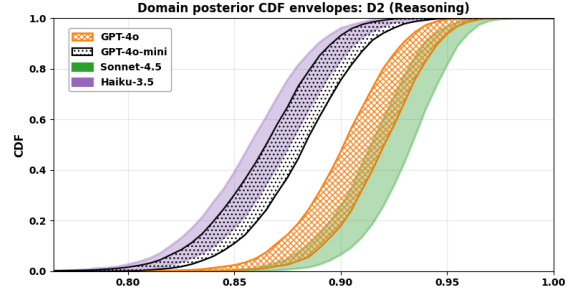
(d) Posterior CDF envelopes of the non-failure probability for Subdom₂₂ (RACE-H dataset in Dom₂) across four models.

Figure 3: Posterior CDF envelopes of non-failure probability at the subdomain level. Rows correspond to domains (D₁ (coding): MBPP, DS-1000; D₂ (reasoning): BoolQ, RACE-H). Subdomains within the same domain are statistically dependent through shared domain-level hyperparameters (μ_i, ν_i), where μ_i represents the average domain reliability and ν_i controls the strength of coupling (partial pooling) among subdomains. As a result, observations from one subdomain inform the inferred reliability of the others. Domains are assumed statistically independent. A question-oriented interpretation of figure: After observing the evaluation data, how do the four models differ in terms of their subdomain-level non-failure probability?

713 Fig. 4 aggregates subdomains within each domain via the operational weights (Ω_{ij}) and
 714 reports the posterior CDFs of the domain reliabilities (p_i). Note, for simplicity, we just
 715 assign the operational weights (Ω_{ij}) proportionally according to the dataset sizes of the
 716 two sub-domains in each domain i . Similarly to sub-domain results, we may observe and
 717 compare domain-level non-failure probabilities p_i . Again, comparing to accuracy scores (point
 718 estimates) and HiBayEs (single posterior distributions), our HIP-LLM yields distribution
 719 envelopes that considers more types of uncertainties.



(a) Posterior CDF envelopes of non-failure probability $p_1 = \sum_j \Omega_{1j} \theta_{1j}$ for Dom_1 (Coding) across four LLMs, with subdomain-level operational weights $\Omega_1 = (0.204, 0.796)$.



(b) Posterior CDF envelopes of non-failure probability $p_2 = \sum_j \Omega_{2j} \theta_{2j}$ for Dom_2 (reasoning) across four models, with subdomain-level operational weights $(0.483, 0.517)$.

Figure 4: Posterior CDF envelopes of non-failure probability of domain-level ($p_i = \sum_j \Omega_{ij} \theta_{ij}$). A question-oriented interpretation of figure: After observing the evaluation data, how do the four models differ in terms of their domain level non-failure probability?

720 Fig. 5 aggregates both domains into the overall LLM reliability p_L using the cross-domain
 721 operational weights $W = [0.149, 0.851]$ for [Coding, Reasoning]. As before, the weights are
 722 simply assigned proportionally according to the dataset sizes of domains. The figure shows,
 723 4o-mini and Haiku bands overlap almost completely, indicating near-equivalent reliability,
 724 while GPT-4o and Sonnet 4.5 partially overlap, reflecting moderate but consistent uncer-
 725 tainty between them. Overall, Sonnet 4.5 remains most reliable (which is consistent with
 726 Table 3).

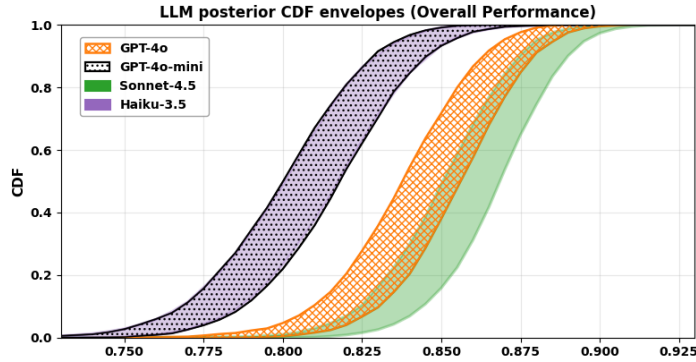
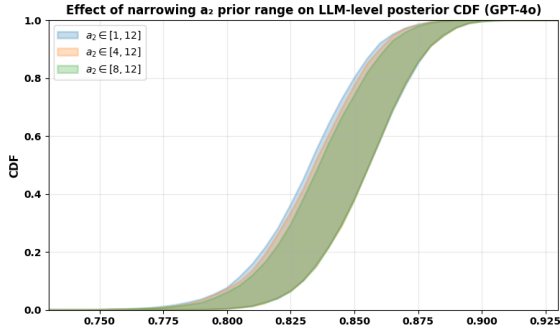
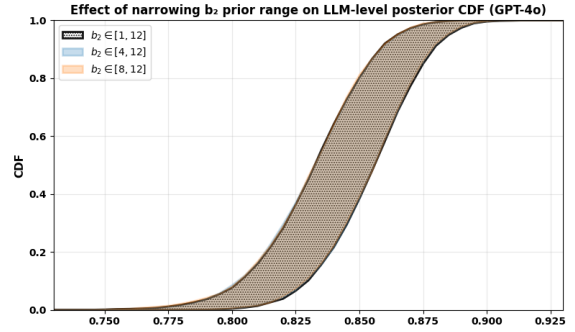


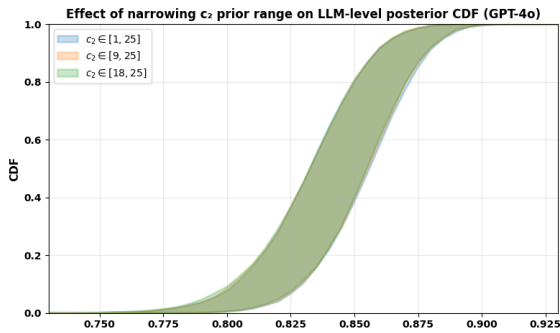
Figure 5: Posterior CDF envelope of non-failure probability of overall LLM level ($p_L = \sum_i W_i p_i$) with domain weights $W = (0.149, 0.851)$ across four models. A question-oriented interpretation of figure: *Engineering question:* After observing the evaluation data, how do the four LLMs differ in terms of their LLM-level non-failure probability?



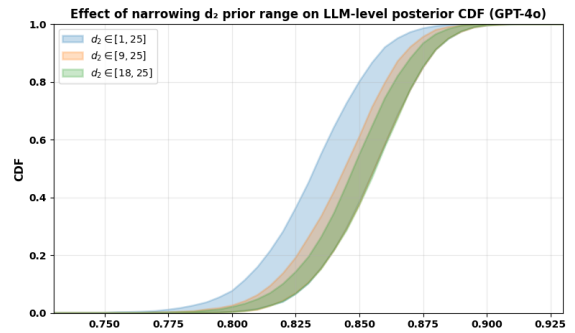
(a) Sensitivity of posterior CDF envelopes of non-failure probability ($p_L = \sum_{i=1}^m W_i p_i$) to a .



(b) Sensitivity of posterior CDF envelopes of non-failure probability of LLM ($p_L = \sum_{i=1}^m W_i p_i$) to b .



(c) Sensitivity of posterior CDF envelopes of non-failure probability of LLM ($p_L = \sum_{i=1}^m W_i p_i$) to c .



(d) Sensitivity of posterior CDF envelopes of non-failure probability of LLM ($p_L = \sum_{i=1}^m W_i p_i$) to d .

Figure 6: Effect of variations in the hyperparameter values (a , b , c , and d) on the posterior CDF envelopes of non-failure probability of LLM ($p_L = \sum_{i=1}^m W_i p_i$). A question-oriented interpretation of the figures: “What are the LLM-level posterior CDF envelopes of the non-failure probability as domain-level hyperparameters vary within their admissible ranges, after observing the testing data?”

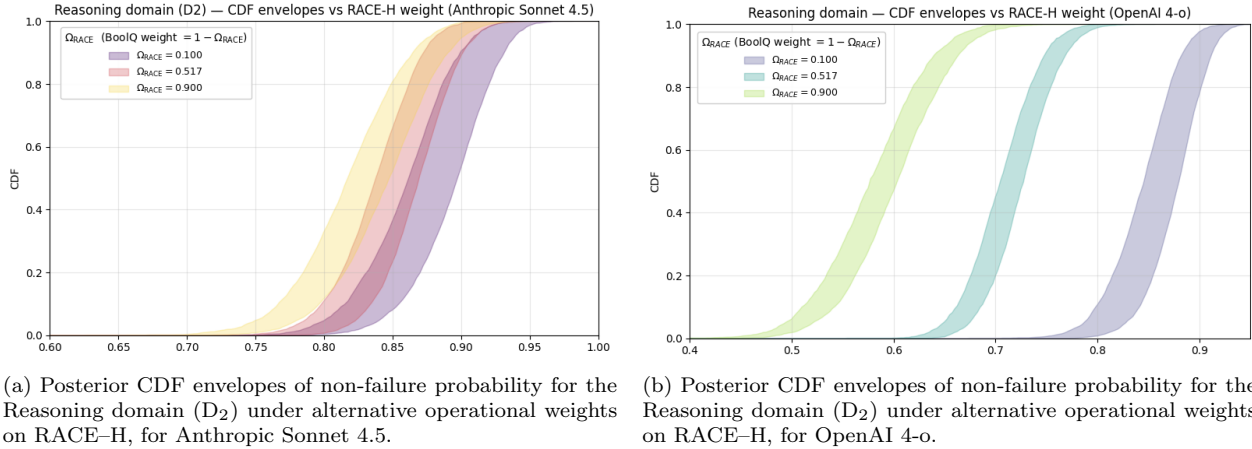
728 Corresponding to Remark 6, the sensitivity analysis in Fig. 6 highlights that different
 729 hyperparameters encode qualitatively different aspects of expert elicitation and therefore
 730 require different levels of care when specified. Variations in a_i and b_i reflect differences in
 731 expert expectations about the average domain-level reliability and mainly affect the location
 732 of the posterior reliability envelopes. As such, specifying these parameters requires careful
 733 consideration of how optimistic or pessimistic prior beliefs are justified by available evidence.
 734 In contrast, the hyperparameters c_i and d_i encode expert beliefs about confidence and pooling
 735 strength across subdomains, and directly control the width of the posterior envelopes.

736 Fig. 6a illustrates the sensitivity of the LLM-level posterior reliability to the hyperpa-
 737 rameter a_i in the prior $\mu_i \sim \text{Beta}(a_i, b_i)$. For example, keeping $b_i = 4$, an expert belief cor-
 738 responding to $a_i = 8$ yields a prior mean $\mu_i = 0.67$, while a more optimistic belief expressed
 739 by $a_i = 16$ increases the mean to $\mu_i = 0.80$. When propagated through the hierarchical
 740 model and updated with the same evaluation data, these beliefs lead to a rightward shift of
 741 the posterior CDF envelope, indicating a higher central estimate of p_L , while the envelope
 742 width remains similar due to the constraining effect of the data. Fig. 6b shows the comple-
 743 mentary effect of varying b_i . With $a_2 = 12$, allowing $b_2 \in [1, 12]$ corresponds to prior
 744 means $\mu_2 \in [0.50, 0.92]$, whereas narrowing the range to $b_2 \in [4, 12]$ and $[8, 12]$ restricts μ_2 to

745 [0.50, 0.75] and [0.50, 0.60], respectively.

746 Fig. 6c and Fig. 6d illustrate the sensitivity of the LLM-level posterior reliability dis-
 747 tributions to the hyperparameters c_i and d_i . Allowing wide ranges such as $c_2, d_2 \in [1, 25]$
 748 corresponds to admitting both weak and strong confidence scenarios, while progressively
 749 narrowing these ranges (e.g., $c_2 \in [9, 25]$ or $[18, 25]$, and similarly for d_2) excludes low- or
 750 high-confidence assumptions. When propagated through the hierarchical model and updated
 751 with the same evaluation data, these changes leave the central location of the posterior
 752 CDF largely unchanged but systematically tighten or widen the envelope, indicating that
 753 c_i and d_i primarily control epistemic uncertainty rather than the expected reliability level.
 754 The stronger effect of d_i occurs because it directly reduces the pooling strength ν_i , while c_i
 755 mainly affects how concentrated this belief is. As a result, changing d_i more strongly weakens
 756 pooling across subdomains and leads to larger changes in the width of the posterior reliability
 757 envelopes.

758 4.3.3. RQ3 (Sensitivity to OPs)



(a) Posterior CDF envelopes of non-failure probability for the Reasoning domain (D_2) under alternative operational weights on RACE-H, for Anthropic Sonnet 4.5.

(b) Posterior CDF envelopes of non-failure probability for the Reasoning domain (D_2) under alternative operational weights on RACE-H, for OpenAI 4-o.

Figure 7: Posterior CDF envelopes of non-failure probability for the Reasoning domain (D_2) under alternative operational weights on RACE-H, $\Omega_{\text{RACE}} \in \{0.10, 0.517, 0.90\}$, with $\Omega_{\text{BoolQ}} = 1 - \Omega_{\text{RACE}}$. (a) Sonnet 4.5; (b) OpenAI 4-o. A question-oriented interpretation of figure: How does the inferred domain-level non-failure probability change as the operational profile shifts emphasis between RACE-H and BoolQ?

759 Fig. 7 shows posterior CDF envelopes for the Reasoning-domain reliability p_2 under three
 760 OPs, $\Omega_{\text{RACE}} \in \{0.10, 0.517, 0.90\}$ with $\Omega_{\text{BoolQ}} = 1 - \Omega_{\text{RACE}}$. Fig. 7a reports Anthropic
 761 Sonnet 4.5; Fig. 7b reports OpenAI 4-o.

762 (a) *Sonnet 4.5*. As Ω_{RACE} increases the CDF shifts slightly left (lower delivered p_2),
 763 with strong overlap among bands. This weak sensitivity to the OP is consistent with the
 764 subdomain accuracies being both high and similar ($\theta_{\text{BoolQ}} \approx 0.90$ vs. $\theta_{\text{RACE-H}} \approx 0.84$): the
 765 mixture $p_2 = \Omega_{\text{BoolQ}}\theta_{\text{BoolQ}} + \Omega_{\text{RACE}}\theta_{\text{RACE-H}}$ changes only modestly as weight moves from
 766 BoolQ to RACE-H. The mild widening of the envelope at higher Ω_{RACE} reflects slightly
 767 larger posterior uncertainty on RACE-H relative to BoolQ.

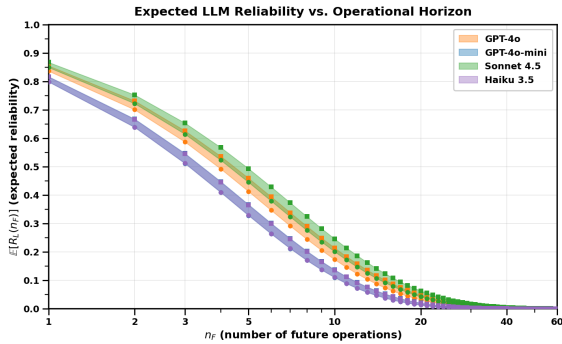
768 (b) *OpenAI 4-o*. In contrast, the envelopes are well separated and move substantially
 769 left as Ω_{RACE} increases. This strong dependence on the operational weighting arises because
 770 4-o’s subdomain accuracies differ widely ($\theta_{\text{BoolQ}} \approx 0.91$ vs. $\theta_{\text{RACE-H}} \approx 0.55$), so the mixture
 771 p_2 is dominated by whichever subdomain receives more weight.

772 Overall, the posterior CDF envelopes for OpenAI 4-o exhibit a strong dependence on
 773 the BoolQ/RACE-H weighting, whereas Sonnet 4.5 is comparatively invariant to Ω_{RACE} .
 774 This arises because GPT-4-o’s subdomain accuracies differ widely (0.91 vs 0.55), produc-
 775 ing a hierarchical posterior with pronounced weight sensitivity. These results suggest that
 776 GPT-4-o’s reasoning performance is more uneven across benchmark types, while Sonnet-4.5
 777 demonstrates domain-level robustness.

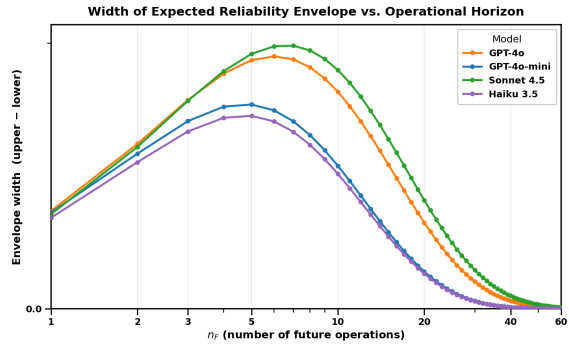
778 From a practical point of view, these results show that changing the operational weights
 779 can meaningfully change the reliability experienced by users, especially for OpenAI 4-o. The
 780 operational weights represent how often different types of tasks occur in real use, so changing
 781 them corresponds to a change in user behavior or application context. For 4-o, different task
 782 mixes lead to clearly different reliability levels, meaning that the same model can appear
 783 reliable in one application but much less so in another. In contrast, Sonnet 4.5 shows relatively
 784 stable reliability across different task mixes, indicating more uniform performance. This
 785 highlights that understanding the expected task distribution is important when deploying
 786 LLMs, particularly for models whose performance varies strongly across subdomains.

787 **4.3.4. RQ4 (Predictability)**

788 Fig. 8a presents the expected LLM-level reliability $\mathbb{E}[R_L(n^F)]$ as a function of the opera-
 789 tional horizon n^F , where n^F denotes the number of consecutive operations in the future. The
 790 horizontal axis displays n^F on a logarithmic scale, while the vertical axis shows the expected
 791 reliability—the posterior mean probability that the LLM successfully completes all next n^F
 792 consecutive tasks without failure.



(a) Imprecise posterior envelopes of the expected LLM reliability ($\mathbb{E}[R_L(n^F)]$) versus the required number of failure-free future tasks n^F for reliability claims. Imprecise posterior envelopes show decay in reliability over consecutive tasks across four models.



(b) Envelope widths (upper–lower) of expected reliability $\mathbb{E}[R_L(n^F)]$. Each curve represents the spread (upper–lower bound) of expected reliability as the number of future operations increases. The bell-shaped profile indicates that uncertainty grows with operational horizon before decaying as overall reliability converges toward zero.

Figure 8: Expected reliability envelopes and their width as functions of number of failure-free future tasks required by assessors (i.e., $\mathbb{E}[R_L(n^F)]$ and [upper bound of envelop – lower bound of envelop]). Question-oriented interpretation of the figure: “what is the posterior expected reliability of passing next n^F future tasks, after seeing testing data?”

793 At $n^F = 1$ all models exhibit high expected reliability, $\mathbb{E}[R_L(1)] \approx 0.80\text{--}0.87$, i.e., a high
 794 probability of producing a correct answer on the very next task (the non-failure probability).
 795 As the horizon length n^F increases, expected reliability drops sharply for every model. Across
 796 all horizons the ordering is stable: Sonnet-4.5 achieves the highest reliability, GPT-4o follows

797 closely, while Haiku-3.5 and GPT-4o-mini form the lower pair (with GPT-4o-mini typically
 798 the lowest). However, as shown in As Fig. 8b, the uncertainty (measured by the interval
 799 width) of Sonnet-4.5’s reliability prediction is also the highest.

800 As Fig. 8b shows, at $n^F = 1$ the envelopes are quite tight, indicating low uncertainty. The
 801 bands widen modestly for small-to-medium horizons (i.e., when n^F increases from 2 to ~ 6),
 802 as a result of compounded prediction uncertainties. As n^F increases, all curves converge
 803 toward zero, indicating reduced uncertainty in the reliability prediction. Intuitively, this
 804 suggests greater confidence that the LLMs are unlikely to pass all n^F future tasks.

805 From a user perspective, Fig. 8 illustrates how reliability claims vary with the required
 806 number of future failure-free tasks n^F : as n^F increases, the claim becomes more stringent,
 807 leading to a corresponding decrease in confidence. This provides practical guidance for se-
 808 lecting models based on intended usage. For example, in healthcare decision support, an
 809 assessor may require a very high reliability of failure-free performance for a small number of
 810 critical diagnostic queries, favoring models with stronger short-horizon reliability guarantees
 811 (i.e., a higher but more steeply declining curve). In legal document analysis, a compliance
 812 team may comprise reliability a bit over a longer sequence of tasks (i.e., a lower but flatter
 813 curve), such as cross-checking clauses or summarizing precedents. In this way, Fig. 8 illus-
 814 trates how model reliability scales with the required number of future failure-free tasks and
 815 supports informed deployment decisions across different application scenarios.

816 **4.3.5. RQ5 (Comparison to baselines)**

817 We evaluate the proposed HIP-LLM method against three Bayesian baselines using
 818 controlled synthetic experiments in which the *ground-truth system reliability* (GT) is known.
 819 The system consists of two domains, each comprising two subdomains. Each subdomain is
 820 characterized by a reliability θ_{ij} and an OP weight OP_{ij} . The overall system reliability is
 821 defined as the OP-weighted average of subdomain reliabilities²¹,

$$p_L = \sum_{i,j} \text{OP}_{ij} \theta_{ij}. \tag{9}$$

822 In Bayesian inference, the posterior distribution reflects a combination of prior beliefs and
 823 observed data, with their relative influence governed by the amount of available evidence. As
 824 the sample size increases, the posterior becomes increasingly dominated by the data rather
 825 than the prior²². Thus, to illustrate HIP-LLM’s ability to incorporate informative prior and
 826 to examine the influence of priors on posterior estimates, we consider two distinct sample-size
 827 regimes:

²¹In the paper, the OP is defined hierarchically, with W_i representing the importance of each domain and Ω_{ij} representing the relative importance of subdomain j within domain i , so that the system reliability is given by $p_L = \sum_i W_i \sum_j \Omega_{ij} \theta_{ij}$. In the experiments, this structure is written in an equivalent but simpler form by specifying directly the combined subdomain weights $\text{OP}_{ij} = W_i \Omega_{ij}$. The system reliability is then computed as $p_L = \sum_{i,j} \text{OP}_{ij} \theta_{ij}$, which is algebraically identical to the hierarchical expression. Under this representation, domain weights are not fixed or ignored: they are implicitly given by $W_i = \sum_j \text{OP}_{ij}$, and the within-domain weights are recovered as $\Omega_{ij} = \text{OP}_{ij}/W_i$. This choice simplifies the experimental setup by defining the OP at the subdomain level, while remaining fully consistent with the hierarchical formulation used in the paper.

²²In the limiting cases, when no data are observed the posterior coincides with the prior, whereas with an infinite amount of data the posterior converges to “frequentist” estimates driven solely by the data.

- 828 • **Small- N** : [100, 500, 1000, 300] samples per subdomain, where priors may have more
829 influence on the posterior.
- 830 • **Large- N** : [1000, 5000, 10000, 3000] samples per subdomain, where data are expected
831 to dominate more than the priors.

832 For each regime, synthetic binomial observations are generated from the same ground-truth
833 subdomain reliabilities. Inference is performed under three OP scenarios:

- 834 • **Dataset-based OP (OP^{data})**: OP weights are set proportional to dataset sizes, in-
835 troducing a structural mismatch with respect to the ground truth OP. This setting
836 represents a naive scenario in which the OP is ignored and assessors rely solely on
837 (benchmark) datasets for reliability assessment.
- 838 • **Approximated OP (OP^{approx})**: The ground truth OP is perturbed by $\pm 20\%$ noise
839 and renormalized, representing a scenario in which assessors recognize the importance
840 of OPs and apply existing OP estimation methods to approximate them, rather than
841 relying solely on evaluation datasets (i.e., OP^{data}), while still incurring estimation er-
842 rors.
- 843 • **Ground-truth OP (OP^{GT})**: The ground truth OP is used directly, representing an
844 idealized oracle scenario in which the assessors know the OP for certain.

845 We compare the following inference methods:

- 846 • **BB-UnInf**: Independent Beta–Binomial model with non-informative priors (Beta(1,1)).
- 847 • **BB-Inf**: Independent Beta–Binomial model with informative priors centered on the
848 GT reliabilities.
- 849 • **HiBayES ([31])**: A hierarchical Bayesian model with partial pooling across subdo-
850 mains.
- 851 • **HIP-LLM**: The proposed imprecise-probability approach, producing an envelope over
852 posterior distributions induced by a family of priors.

853 The criteria for comparison are as follows:

- 854 • **Median (\hat{p}_L^{med})**: For methods producing a single posterior distribution (BB-UnInf,
855 BB-Inf, HiBayES), the median is defined as the empirical median of the posterior
856 samples of p_L . For HIP-LLM, which induces a *set of posterior distributions* indexed by
857 hyperparameter configurations $h \in \mathcal{H}$, we first compute the median of each posterior,

$$m_h = Q_{0.50}(p_L^{(h)}),$$

858 and then report the envelope of posterior medians,

$$\hat{p}_L^{\text{med}} \in \left[\min_{h \in \mathcal{H}} m_h, \max_{h \in \mathcal{H}} m_h \right].$$

- 859 • **Error:** For single-posterior methods, the error is defined as the absolute deviation
 860 between the posterior median and the ground-truth reliability,

$$\text{Error} = |\hat{p}_L^{\text{med}} - p_L^{\text{GT}}|.$$

861 For HIP-LLM, the error is computed *per posterior envelope* as

$$e_h = |m_h - p_L^{\text{GT}}|,$$

862 and the reported error corresponds to the envelope of admissible errors,

$$\text{Error} \in \left[\min_{h \in \mathcal{H}} e_h, \max_{h \in \mathcal{H}} e_h \right].$$

- 863 • **Credible Interval (CI):** For single-posterior methods, a Bayesian 90% credible inter-
 864 val is defined as

$$\text{CI}_{0.90} = [Q_{0.05}(p_L), Q_{0.95}(p_L)],$$

865 where $Q_q(\cdot)$ denotes the empirical quantile.

866 For HIP-LLM, which produces a family of posteriors $\{p_L^{(h)}\}_{h \in \mathcal{H}}$, we report the widest
 867 credible interval compatible with the assumed prior uncertainty,

$$\left[\min_{h \in \mathcal{H}} Q_{0.05}(p_L^{(h)}), \max_{h \in \mathcal{H}} Q_{0.95}(p_L^{(h)}) \right].$$

868 Tables 4 and 5 illustrate how system reliability estimates are influenced by the OP and
 869 prior assumptions. In the Small- N regime, where data are limited, all methods slightly un-
 870 derestimate the ground-truth reliability ($p_L^{\text{GT}} = 0.5860$), with posterior medians typically
 871 between 0.57 and 0.58. In this setting, both OP mismatch and prior choice matter: using
 872 an approximate or dataset-based OP increases error compared to the ground-truth OP, and
 873 informative priors can reduce point error when they are correctly aligned with the GT sub-
 874 domain reliabilities. This reduction, however, is due to favorable prior specification rather
 875 than inherent robustness and would not persist if the prior were misspecified.

Table 4: Posterior estimates under the Small- N ([100, 500, 1000, 300]) regime ($p_L^{\text{GT}} = 0.5860$).

OP	Method	Median value/ bound	Error value/bound	90% Interval
OP ^{data}	BB-UnInf	0.5733	0.0127	[0.5549, 0.5917]
	BB-Inf	0.5749	0.0111	[0.5569, 0.5924]
	HiBayES	0.5732	0.0128	[0.5552, 0.5915]
	HIP-LLM	[0.5713, 0.5751]	[0.0109, 0.0147]	[0.5528, 0.5937]
OP ^{approx}	BB-UnInf	0.5740	0.0120	[0.5460, 0.6009]
	BB-Inf	0.5756	0.0104	[0.5509, 0.5995]
	HiBayES	0.5712	0.0148	[0.5438, 0.5966]
	HIP-LLM	[0.5698, 0.5774]	[0.0086, 0.0162]	[0.5421, 0.6048]
OP ^{GT}	BB-UnInf	0.5782	0.0078	[0.5486, 0.6062]
	BB-Inf	0.5796	0.0064	[0.5538, 0.6042]
	HiBayES	0.5750	0.0110	[0.5459, 0.6011]
	HIP-LLM	[0.5738, 0.5816]	[0.0044, 0.0122]	[0.5446, 0.6100]

876 HIP-LLM behaves differently in the Small- N regime by reporting intervals rather than
877 single reliability values. For example, it produces ranges of plausible medians such as
878 [0.5713, 0.5751] under the dataset-based OP and [0.5738, 0.5816] under the ground-truth
879 OP. These wider bounds reflect epistemic uncertainty about prior assumptions and make ex-
880 plicit that multiple reliability estimates are consistent with the data. Hierarchical modeling
881 through HiBayES also smooths estimates across subdomains, but can introduce additional
882 bias when data are scarce.

883 In the Large- N regime, the effect of priors and modeling choices largely disappears. All
884 methods produce nearly identical posterior medians across OP scenarios, and uncertainty
885 intervals become very narrow. For HIP-LLM, the median interval collapses to a tight range
886 (e.g., [0.5860, 0.5866] under the dataset-based OP), showing that prior uncertainty no longer
887 plays a significant role.

888 When comparing methods, HIP-LLM intervals should be interpreted as robustness bounds
889 rather than standard credible intervals; for conservative decision making, the lower bound of
890 the HIP-LLM median interval provides a worst-case reliability estimate consistent with the
891 available data.

892 4.3.6. **RQ6** (*Sensitivity to the definition of success/failure*)

893 A key modeling assumption in HIP-LLM is the use of binary success/failure labels at
894 the subdomain level, summarized by the observed success counts C_{ij} and total trials N_{ij} .
895 In practice, however, the notion of “success” is not unique and may depend on the chosen
896 evaluation criterion.

897 To assess the robustness of HIP-LLM to alternative success definitions, we conduct a con-
898 trolled experiment on a single dataset (MBPP) and a single model (Claude Sonnet 4.5). We
899 compare two common criteria: *Pass@1*, where only the first generated solution is evaluated,
900 and *Pass@3*, where a task is considered successful if at least one out of three independent
901 generations passes all tests. All other experimental conditions are held fixed.

Table 5: Posterior estimates under the Large-N ([1000, 5000, 10000, 3000]) regime ($p_L^{\text{GT}} = 0.5860$).

OP	Method	Median value/ bound	Error value/bound	90% Interval
OP ^{data}	BB-UnInf	0.5907	0.0047	[0.5849, 0.5965]
	BB-Inf	0.5908	0.0048	[0.5849, 0.5965]
	HiBayES	0.5909	0.0049	[0.5846, 0.5969]
	HIP-LLM	[0.5860, 0.5866]	[0.0000, 0.0006]	[0.5846, 0.5968]
OP ^{approx}	BB-UnInf	0.5887	0.0027	[0.5801, 0.5973]
	BB-Inf	0.5885	0.0025	[0.5799, 0.5970]
	HiBayES	0.5887	0.0027	[0.5800, 0.5980]
	HIP-LLM	[0.5752, 0.5764]	[0.0096, 0.0108]	[0.5794, 0.5978]
OP ^{GT}	BB-UnInf	0.5928	0.0068	[0.5836, 0.6017]
	BB-Inf	0.5926	0.0066	[0.5837, 0.6015]
	HiBayES	0.5929	0.0069	[0.5838, 0.6025]
	HIP-LLM	[0.5784, 0.5797]	[0.0063, 0.0076]	[0.5832, 0.6024]

902 Figure 9 shows the subdomain-level posterior CDF envelopes for the true success proba-
 903 bility $\frac{C}{N}$ under Pass@1 and Pass@3. As expected, the more permissive Pass@3 criterion yields
 904 a higher observed success count and a right-shifted posterior envelope. However, the two en-
 905 velopes exhibit substantial overlap, indicating that the inferred reliability is only moderately
 906 sensitive to the choice of success definition.

907 4.3.7. **RQ7** (*Robustness to Memory Effects*)

908 Our formal reliability definition models task outcomes as i.i.d. Bernoulli trials under a
 909 fixed OP. In practice, however, LLM usage may violate independence because the model can
 910 condition on previously observed conversational context (i.e., *memory*). To assess the practi-
 911 cal impact of such dependence on HIP-LLM, we conducted two complementary experiments:
 912 (i) an empirical evaluation to quantify memory growth when memory is not reset after each
 913 task and its relation to dependence strength; and (ii) a targeted sensitivity analysis that in-
 914 jects memory-induced dependence into one of subdomains (Subdom₂₁-BoolQ) and examines
 915 how the resulting dependence shifts posterior envelopes across hierarchy levels.

916 We evaluated one subdomain (Subdom₂₁-BoolQ)²³ in a single continuous session in which
 917 the entire conversation history was retained in every API call. For each query, we recorded
 918 the API-reported memory. In this setting, memory grows monotonically with the number
 919 of questions and provides an operational proxy for the strength of cross-item dependence
 920 introduced by persistent context (Fig. 10a). After 299 questions, the retained context reached
 921 $\approx 2.06 \times 10^5$ bytes, i.e., ≈ 201 KB, illustrating that long-context evaluation can induce non-
 922 negligible memory accumulation even within a single benchmark run.

923 Using the above proxy, we observed that the dependence strength induced by memory

²³We conducted this experiment on 300 tasks from the BoolQ dataset. To preserve the OP proportions across subdomains, we kept the same relative weights for the other datasets as in the main evaluation.

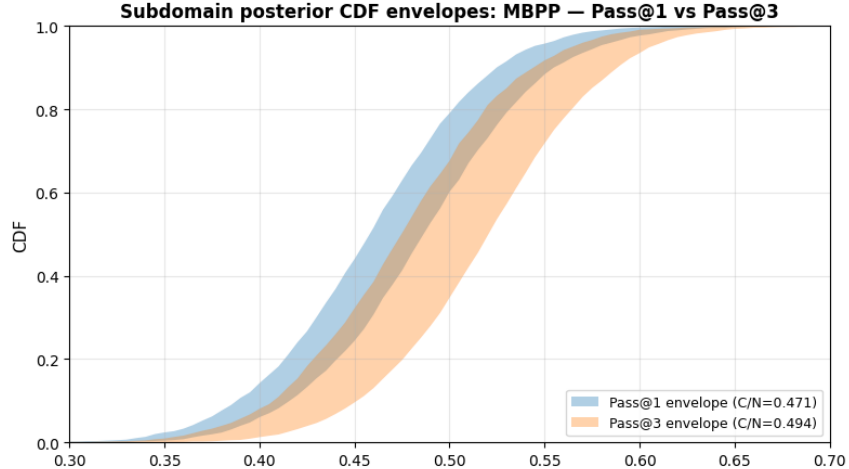


Figure 9: Subdomain-level posterior CDF envelopes of non-failure probability for MBPP under alternative success definitions. Results are shown for Claude Sonnet 4.5 on MBPP ($N = 257$), using identical tasks, prompts, and tests. Pass@1 yields $\frac{C}{N} = 0.471$, while Pass@3 yields $\frac{C}{N} = 0.494$. Envelopes reflect imprecise hierarchical inference with identical hyper-hyperparameter intervals; only the success definition differs. A question-oriented interpretation of figure: How does the inferred subdomain-level non-failure probability change when alternative success (failure) definitions, such as Pass@1 versus Pass@3, are used on the same evaluation data?

924 increases approximately linearly with the retained context size. This motivates a parsimo-
 925 nious sensitivity model in which the dependence parameter is a linear function of memory.
 926 The linear mapping is used only as a stress-test mechanism to quantify how departures from
 927 independence could propagate into the inferred reliability distributions; it does not change
 928 the definition of reliability itself, but probes the robustness of posterior conclusions when the
 929 i.i.d. approximation is imperfect.

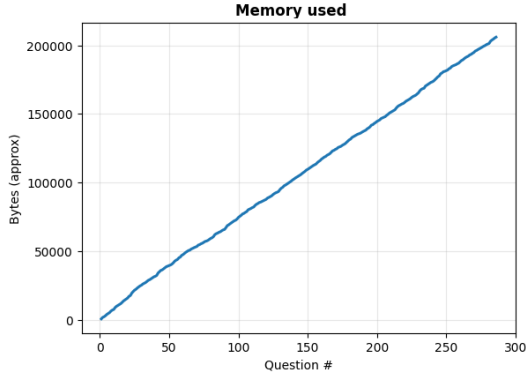
930 We then injected memory-induced dependence into BoolQ according to the linear depen-
 931 dence model and compared posterior CDF envelopes under three BoolQ θ settings,

$$\theta_{\text{BoolQ}} \in \{0.915, 0.940, 0.945\},$$

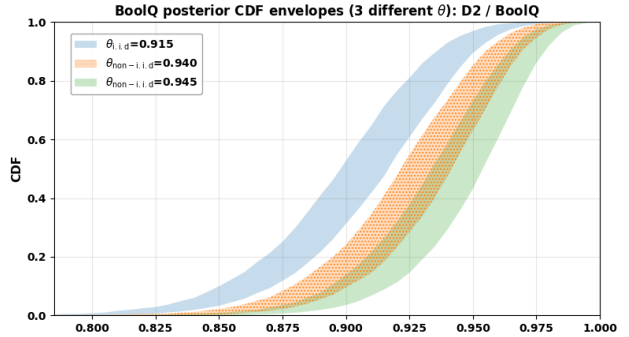
932 where $\theta_{\text{BoolQ}} = 0.915$ corresponds to the *i.i.d. assumption* (no memory-induced dependence),
 933 while $\theta_{\text{BoolQ}} = 0.940$ and 0.945 represent non-i.i.d. regimes induced by increasing retained
 934 context. The remaining subdomain accuracies (MBPP, DS-1000, RACE-H) were kept fixed,
 935 and the same hierarchical structure and operational-profile aggregation were preserved (sub-
 936 domain \rightarrow domain via Ω_{ij} , and domain \rightarrow LLM via W_i). For each θ_{BoolQ} , we propagated
 937 uncertainty through the hierarchy under the same imprecise hyperparameter intervals and
 938 computed CDF envelopes at three levels: (i) BoolQ subdomain, (ii) Reasoning domain (D_2),
 939 and (iii) overall LLM level.

940 Across all three levels, the envelopes shift smoothly with θ_{BoolQ} : increasing θ_{BoolQ} yields
 941 a consistent rightward shift (stochastically larger reliability), while decreasing θ_{BoolQ} shifts
 942 the envelopes leftward.

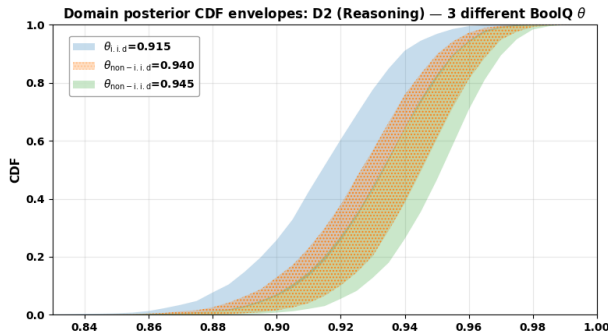
943 These experiments support two practical conclusions: (1) persistent context can induce
 944 measurable dependence during benchmark-style evaluation, and the effect scales with mem-
 945 ory; and (2) within a realistic range around the observed reliability levels, our posterior



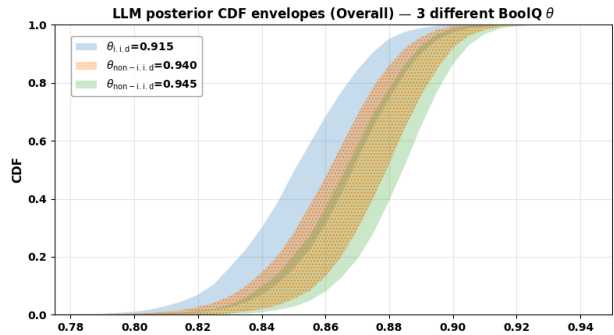
(a) Growth of retained conversational context as a function of the number of BoolQ tasks evaluated in a single continuous session. Memory size (in bytes) increases monotonically and serves as an operational proxy for the strength of memory-induced dependence.



(b) Posterior CDF envelopes of non-failure probability for the BoolQ subdomain under three accuracy settings: $\theta_{\text{BoolQ}} = 0.915$ (i.i.d.), and $\theta_{\text{BoolQ}} = 0.940, 0.945$ (non-i.i.d.). Increasing memory-induced dependence produces smooth rightward shifts of the envelope.



(c) Posterior CDF envelopes of non-failure probability at the Reasoning domain level (Dom2) resulting from propagation of BoolQ dependence through subdomain aggregation with weights Ω_{ij} .



(d) Posterior CDF envelopes of non-failure probability at the LLM level after domain aggregation with operational-profile weights W_i , illustrating the attenuated but consistent impact of subdomain-level dependence.

Figure 10: Sensitivity of posterior CDF envelopes of non-failure probability to memory-induced non-i.i.d. dependence. The top-left panel (10a) quantifies memory accumulation during same-session evaluation, motivating a linear dependence proxy. The remaining panels show how increasing dependence in the BoolQ subdomain propagates through the hierarchical reliability model, yielding smooth and controlled shifts of envelopes at subdomain, domain, and LLM levels.

946 envelope conclusions remain stable, with dependence primarily producing controlled shifts
947 rather than qualitative changes in domain- and LLM-level reliability characterization.

948 4.3.8. **RQ8** (*Scalability*)

949 This section characterizes the computational cost of the proposed HIP-LLM inference
950 pipeline. The analysis is performed at the level of the algorithmic procedure (rather than for
951 any single experimental figure) and is structured to (i) derive theoretical time and memory
952 scaling laws from the inference steps, and (ii) empirically validate these scaling laws by
953 controlled parameter sweeps.

954 Let m denote the number of domains and \bar{n} the average number of subdomains per
955 domain²⁴. For each domain i , HIP-LLM represents hyperparameters on a discretized grid over
956 (μ_i, ν_i) of size $G = n_\mu n_\nu$. Imprecision at the hyperparameter level is handled by evaluating
957 a finite set of K hyperparameter configurations (or candidates) over the grid. To propagate
958 uncertainty from subdomain to domain and from domain to system level, the method uses
959 Monte Carlo sampling with S samples per configuration. Posterior envelopes are reported
960 on a fixed evaluation grid (e.g., a CDF grid) of size T .

961 The inference is decomposed into three hierarchical stages:

- 962 1. **Subdomain level:** For each domain and each hyperparameter configuration, HIP-
963 LLM evaluates subdomain posteriors conditional on (μ_i, ν_i) on the grid and constructs
964 the corresponding imprecise posterior set.
- 965 2. **Domain level:** For each configuration, HIP-LLM propagates subdomain uncertainty
966 to the domain level via Monte Carlo sampling by drawing subdomain reliability vari-
967 ables and aggregating them using OP weights Ω_{ij} .
- 968 3. **LLM level:** HIP-LLM aggregates domain-level reliability to the system level using
969 domain weights W_i , while preserving imprecision by computing lower/upper envelopes
970 across the explored hyperparameter configurations.

971 Theoretical time complexity is obtained by counting the dominant operations required
972 by the above stages as functions of (m, \bar{n}, G, K, S, T) . Subdomain-level cost is driven by
973 grid-based hyperposterior evaluation and any intra-domain coupling required to compute
974 subdomain posteriors under shared hyperparameters. Domain-level cost is driven by Monte
975 Carlo propagation, which scales with the number of subdomains drawn per sample and the
976 number of samples. LLM-level cost is driven by repeating aggregation across hyperparameter
977 configurations and computing envelope statistics on the evaluation grid. Memory complexity
978 is characterized by the storage required for (i) grid-based posterior arrays per domain/con-
979 figuration, (ii) Monte Carlo samples (if stored) or streaming statistics (if not stored), and
980 (iii) envelope representations on the evaluation grid.

981 To validate the theoretical scaling, we use controlled experiments in which all numerical
982 parameters are fixed to baseline values and then varied one at a time. Specifically, we run
983 repeated executions while sweeping m , \bar{n} , K , S , and G independently, and record wall-clock

²⁴Although different domains may have different numbers of subdomains, the computational cost is mainly determined by the typical workload of a domain. For this reason, we describe the scaling in terms of the average number of subdomains \bar{n} . This keeps the formulas simple, reflects the observed runtime in practice, and avoids overly pessimistic worst-case estimates.

984 time and peak memory usage for each run. This isolates the marginal effect of each parameter
985 and enables direct comparison with the predicted scaling laws. To quantify scaling exponents,
986 we fit a power-law model of the form $\text{time} = c x^\alpha$ (and analogously for memory), where x is
987 the swept parameter²⁵ and α ²⁶ is the inferred scaling exponent. Here c is a constant which
988 sets the scale.

989 Many steps of HIP-LLM can be computed independently across domains and hyperparam-
990 eter configurations, making the method naturally parallelizable. Parallel execution reduces
991 wall-clock runtime in practice, while the underlying asymptotic computational complexity
992 remains unchanged.

993 The following experiment empirically validates the theoretical computational complexity
994 analysis of the HIP-LLM inference pipeline. We report controlled scalability experiments that
995 measure wall-clock runtime and peak memory usage while varying key algorithmic parameters
996 one at a time. The results are used to confirm the predicted asymptotic scaling laws²⁷ and
997 to identify the dominant computational stages in practice.

998 *Scalability plots.* (Fig. 11a- Fig. 11e) showing runtime as a function of each swept parameter
999 (m , \bar{n} , K , S , and G), together with fitted scaling laws. Across all sweeps, the empirical
1000 results closely match the theoretical scaling predictions and confirm that subdomain-level
1001 inference dominates the computational cost of HIP-LLM in practice.

1002 *Baseline timing breakdown.* (Fig. 11f) reporting the relative contribution of each inference
1003 stage under the reference configuration. The results show that subdomain-level posterior
1004 computation dominates runtime, accounting for over 99% of total execution time, while
1005 domain- and LLM-level aggregation contributes a negligible fraction. This explains the weak
1006 empirical dependence on S and highlights the subdomain inference stage as the primary
1007 computational bottleneck.

1008 *Memory Complexity.* According to Tab. 6 empirical results confirm the expected memory
1009 behavior of HIP-LLM. Peak memory usage scales approximately linearly with the number
1010 of domains m , the average number of subdomains \bar{n} , the number of hyperparameter con-
1011 figurations K , and the number of Monte Carlo samples S , due to the storage of posterior
1012 quantities and sampled reliabilities. In contrast, memory usage is effectively independent
1013 of the integration grid size G , indicating that the (μ, ν) grid is reused across configurations
1014 rather than stored per sample.

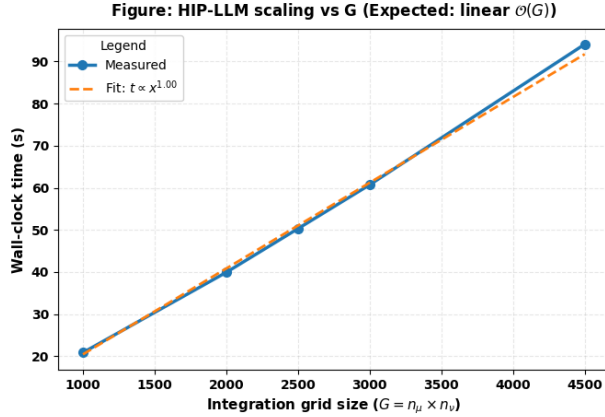
1015 5. Discussion

1016 In this section, we discuss the assumptions and limitations of the proposed HIP-LLM.

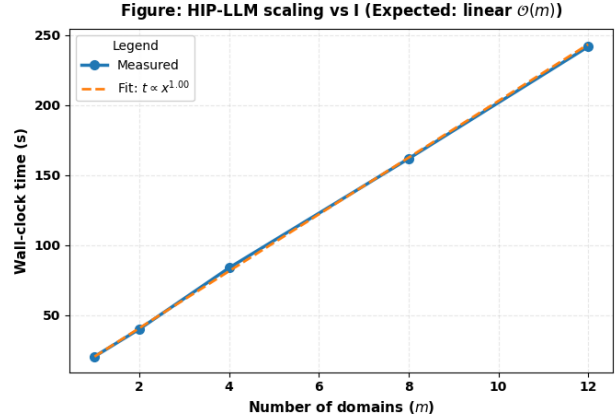
²⁵The parameter we change (e.g. number of subdomains, samples, domains, etc.)

²⁶The parameter that shows how fast time grows when x increases. for instance, $\alpha = 1$ corresponds to linear growth, while $\alpha = 2$ corresponds to quadratic growth, and so on.

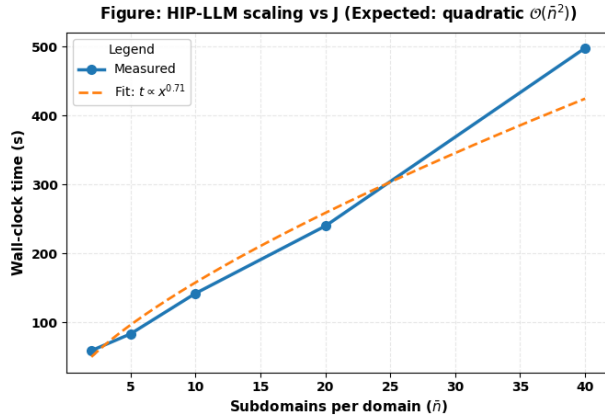
²⁷Do the measured runtimes grow the way the theory says they should?



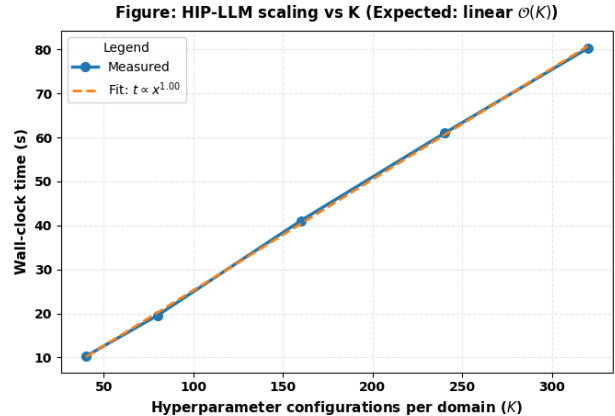
(a) Scaling with the (μ, ν) integration grid size $G = n_\mu n_\nu$, showing linear growth in wall-clock time, consistent with the theoretical $\mathcal{O}(G)$ complexity.



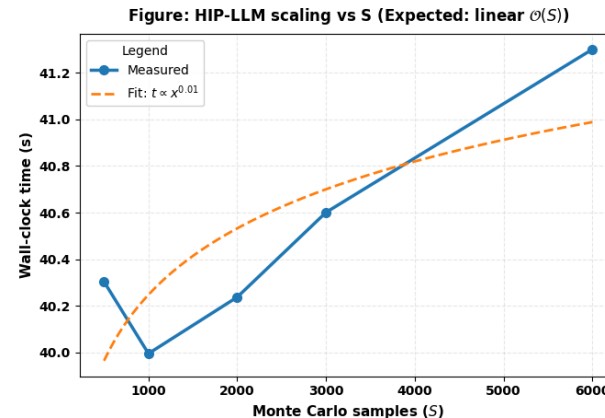
(b) Scaling with the number of domains m , demonstrating linear runtime growth due to independent per-domain inference, in agreement with the expected $\mathcal{O}(m)$ behavior.



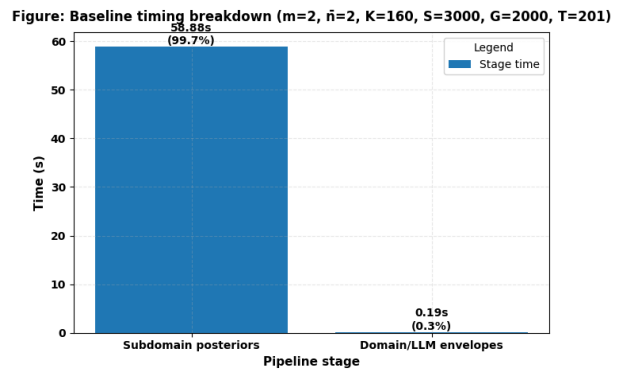
(c) Scaling with the average number of subdomains per domain \bar{n} , exhibiting sub-quadratic growth over the tested range, reflecting practical efficiencies relative to the worst-case $\mathcal{O}(\bar{n}^2)$ bound.



(d) Scaling with the number of hyperparameter configurations per domain K , confirming linear complexity due to independent evaluation across configurations.



(e) Scaling with the number of Monte Carlo samples S , showing weak dependence on S and indicating that Monte Carlo propagation is not the dominant computational cost in the tested regime.



(f) Baseline timing breakdown for the full HIP-LLM pipeline, showing that subdomain posterior computation dominates the total runtime, while domain- and LLM-level envelope aggregation contributes a negligible fraction.

Figure 11: Empirical scalability of the HIP-LLM inference pipeline under controlled parameter sweeps. Each subfigure varies a single parameter while keeping all others fixed at the baseline configuration ($m = 2$, $\bar{n} = 2$, $K = 160$, $S = 3000$, $G = 2000$, $T = 201$, $K_{\text{total}} \leq 512$). Measured runtimes are compared against power-law fits to validate theoretical complexity predictions and to identify dominant computational stages.

Table 6: Peak memory usage under controlled scalability experiments. Each row reports the measured peak memory while sweeping one parameter and keeping all others fixed at the baseline configuration.

Swept parameter	Range	Peak memory (MB)
Number of domains m	1 \rightarrow 12	24.6 \rightarrow 146.7
Subdomains per domain \bar{n}	2 \rightarrow 40	35.8 \rightarrow 314.2
Monte Carlo samples S	500 \rightarrow 6000	7.6 \rightarrow 69.4
Hyperparameter configs K	40 \rightarrow 320	19.1 \rightarrow 57.8
Grid size G	1000 \rightarrow 4500	\approx 35.7 (constant)

1017 *Assumptions on hierarchical dependencies of failure probabilities.* Inspired by HiBayES [31],
 1018 we introduce a hierarchical dependency structure where dependencies are defined at the level
 1019 of failure probabilities across related task groups. Specifically, high-level domains (e.g., cod-
 1020 ing vs. law) are modeled as independent because their underlying competencies and training
 1021 data sources are substantially distinct. In contrast, subdomains within a domain (e.g., coding
 1022 in C++ vs. coding in Python) are modeled as having dependent failure probabilities, reflect-
 1023 ing shared latent skills, representations, or reasoning patterns that influence performance
 1024 across related tasks. This abstraction allows tractable modeling of correlated reliability
 1025 while acknowledging that the true dependencies among knowledge areas are more complex
 1026 and fluid. While we do not prescribe specific taxonomies of domain/sub-domains for LLMs
 1027 in this paper, our formulation is intended to be without loss of generality: users of HIP-LLM
 1028 can redefine their own “domains” to reflect their own assumptions about dependence. For
 1029 example, if users believe that two domains (say, A and B) are not fully independent, they
 1030 may treat them as two subdomains under a newly introduced higher-level synthetic domain.
 1031 Such flexibility allows HIP-LLM to adapt to different operational interpretations or empirical
 1032 evidence of dependency, reflecting the assessor’s domain knowledge.

1033 Note, such (in)dependencies at the parameter level (i.e., correlations or independence
 1034 between the failure probabilities of domains or subdomains) are different from task-level
 1035 outcome dependencies, which we discuss next.

1036 *Assumptions of i.i.d. Bernoulli Trails for Task Outcomes.* It is important to distinguish the
 1037 assumption of independent task outcomes from the (in)dependencies among failure prob-
 1038 abilities discussed earlier. The latter concern parameter-level relationships, i.e., how the
 1039 long-run failure rates across domains or subdomains may co-vary due to shared competen-
 1040 cies, whereas the former concerns instance-level independence of individual task results, given
 1041 those parameters. Even if the failure probabilities of two sub-domains are correlated random
 1042 variables, we can still assume that, conditional on those probabilities, each task outcome
 1043 within a sub-domain is an independent Bernoulli trial.

1044 As discussed earlier in Def. 4 and Remark 1, the i.i.d. assumption is pragmatically justified
 1045 in the context of LLM evaluations when tasks are executed in isolation, such as *initiating new*
 1046 *chat sessions without contextual memory*. Nevertheless, caution is warranted: in practice, the
 1047 way tasks are designed, sequenced, or batched can introduce subtle dependencies that violate
 1048 this idealization, e.g., in long-context, agentic workflows. Assessors should therefore ensure
 1049 testing procedures (e.g., resetting each task session) align with the independence assumption,
 1050 a common setup in offline benchmarking.

1051 Indeed, the i.i.d. Bernoulli trial assumption remains a simplification, motivated by
1052 mathematical convenience and interpretability. However, this assumption is not introduced
1053 uniquely by our HIP-LLM, rather that models with the i.i.d. assumption have a long history
1054 of being used in software reliability assessment. Thayer *et al.*'s [79] model was one of the
1055 earliest, used in early works on random testing [80]. Regulatory bodies also have recom-
1056 mended using the i.i.d. assumption in reliability assessments when appropriate [40]. Having
1057 said that, weakening the i.i.d. assumption has been a long-standing research topic in the
1058 software reliability community. Numerous models have been proposed to capture depen-
1059 dence among successive tests, including a binary Markov chain model [81], a Markov renewal
1060 process [82, 83], the Markov-based model with benign-failure states [84], and Bayesian mod-
1061 els examining how sensitive reliability claims are to potential dependencies of test outcomes
1062 [71, 85].

1063 A road-map for extending HIP-LLM beyond i.i.d. task outcomes can build directly on
1064 aforementioned, established software reliability models. More concretely, a first step is treat-
1065 ing a multi-turn interaction as the unit of analysis and allowing within-session dependence
1066 while modeling across-session dynamics via latent states. One can introduce Markov/hidden-
1067 state structures, e.g., hierarchical hidden Markov model or Markov renewal formulations
1068 where latent “interaction states” evolve over time and govern task-level failure probabilities.

1069 *What constitutes a failure for LLMs.* Reliability, by definition, concerns “failures”. In tra-
1070 ditional software systems, failures are relatively straightforward to define, as they can be
1071 explicitly linked to violations of formal specifications or requirements. However, for AI sys-
1072 tems, and particularly for LLMs, no such explicit specification exists. As a result, a complete
1073 and formal characterization of what constitutes a failure for LLMs remains an open research
1074 challenge. The notion of “failure” can vary across applications and stakeholders. For in-
1075 stance, factual errors, inconsistency²⁸, harmful or biased content, or divergence from human
1076 judgment may all be considered failures depending on the context [30].

1077 The choice of scoring criteria for defining LLM failures is inherently application depen-
1078 dent and should align with domain-specific specifications (if any) and evaluation norms. For
1079 example, in code generation, failures are typically defined in terms of functional correctness,
1080 such as compilation success or passing predefined test cases. In mathematical problem solv-
1081 ing, correctness is often determined by agreement with a unique ground-truth answer or a
1082 formally verifiable solution. In contrast, in legal analysis, failure definitions may rely on
1083 expert judgment regarding factual accuracy, consistency with legal principles, or omission of
1084 critical considerations. For more open-ended tasks such as email drafting or creative writ-
1085 ing, failures are often defined relative to stylistic, pragmatic, or preference-based criteria,
1086 frequently assessed through human evaluation. These examples illustrate that failure defini-
1087 tions for LLMs are not universal but must be specified in accordance with the intended use
1088 and domain requirements, reinforcing the need to interpret reliability estimates as conditional
1089 on the chosen scoring criteria.

1090 That said, in this work, our focus is not on prescribing a specific definition of LLM

²⁸Most existing LLM evaluation studies focus on one-shot accuracy metrics (e.g., Pass@1), whereas real-world usage often involves iterative prompting. Inconsistencies across multiple outputs can undermine user trust and thus regarded as a failure.

1091 failure but on *providing a reliability assessment modeling framework that is agnostic to failure*
1092 *definitions and failure detections (i.e., defining the test oracle)*. For example, the reliability
1093 definition in Def. 4 can be readily extended from binary failures to non-binary scores by
1094 replacing the indicator function with a continuous scoring function (e.g., taking values in
1095 $[0, 1]$ through probabilistic or multi-rater failure models), thereby capturing scoring noise
1096 or subjectivity (commonly referred to as the imperfect test oracle problem in traditional
1097 software reliability). HIP-LLM can accommodate any well-defined notion of failure, enabling
1098 users to instantiate the model with the criteria most relevant to their domain or evaluation
1099 objectives.

1100 *How can OPs be specified in practice.* OPs are fundamental to software reliability assessment,
1101 and the same principle applies to LLM reliability. HIP-LLM provides explicit hierarchical
1102 “interfaces” for integrating OPs; however, in the absence of access to real operational data,
1103 this study relies on simulated OPs derived from existing benchmark datasets. Specifically,
1104 benchmark datasets are treated as sampling frames that approximate the unknown opera-
1105 tional data distribution (cf. Remark 7), with task-level distributions assumed within each
1106 subdomain and operational weights assigned across subdomains and domains. This experi-
1107 mental choice is made solely to demonstrate how HIP-LLM incorporates OPs and to enable
1108 fair comparison with existing benchmark-based methods, and *should not be interpreted as a*
1109 *solution to OP acquisition*. Accurately estimating OPs from real usage data is beyond the
1110 scope of this reliability modeling study and is a well-recognized, separate research problem
1111 in software reliability engineering, for which numerous methodologies have been proposed
1112 (e.g., [25, 26, 28]). The concept of OP has also been discussed and measured in AI software
1113 contexts [24, 76, 86, 87].

1114 Importantly, the challenge of modeling LLM OPs does not undermine the proposed frame-
1115 work; rather, it highlights why OPs are essential to reliability assessment. In established
1116 software reliability practice, OPs are progressively refined as new information becomes avail-
1117 able (e.g., execution logs, usage telemetry, or expert-informed adjustments) without requiring
1118 changes to the underlying reliability model [25, 26, 28, 77]. As more representative datasets,
1119 user interaction logs, or task-distribution estimates for LLMs become available [27], OPs can
1120 be updated accordingly, allowing reliability estimates to be recalibrated toward actual field
1121 usage. Such separation and “decoupling design” of HIP-LLM between evidence (testing re-
1122 sults) and usage assumptions (OPs) enables HIP-LLM to evolve naturally with deployment
1123 conditions, improving fidelity without *redesigning* the model.

1124 While we acknowledge that our present treatment of OP uncertainty remains simpli-
1125 fied (for demonstration/comparison), more complex modeling on OP uncertainties, e.g., dy-
1126 namic OPs, will form our important future works. Notably, a range of reliability models
1127 in traditional software engineering explicitly account for OP uncertainty and dynamics, in-
1128 cluding Bayesian updating, adaptive OP adjustment, and black-box or time-varying usage
1129 models [72, 26, 73, 74, 75]. These established techniques provide a foundation that can be
1130 retrofitted into HIP-LLM to support uncertain and dynamical OPs in LLM deployments.

1131 *How to embed prior knowledge.* HIP-LLM models uncertainty at three levels: the subdomain
1132 success rate (θ_{ij}), how much domains share information (μ_i, ν_i), and the top-level parameters
1133 (a_i, b_i, c_i, d_i) that control the importance of those priors.

1134 As a Bayesian framework, HIP-LLM requires priors to initiate inference, and these pri-
1135 ors may substantially influence the resulting posteriors, particularly when data are limited
1136 [88, 89]. The challenge of eliciting meaningful priors is *not unique* to HIP-LLM but is com-
1137 mon to all Bayesian models. Unlike HiBayES [31], which employs noninformative priors,
1138 HIP-LLM follows the principles of robust Bayesian analysis [59, 60, 90] by adopting the
1139 Imprecise Probability framework. This approach enables the representation of uncertainty
1140 about prior beliefs through sets of priors rather than committing to a single distribution.
1141 Prior research has shown that Imprecise Probability can facilitate expert elicitation by sup-
1142 porting structured “thought experiments” [65, 91, 92] allowing experts to express vague,
1143 partial, or interval-valued prior beliefs. Similarly, HIP-LLM’s prior layer can be informed by
1144 such thought experiments and expert judgment. While we acknowledge that eliciting priors
1145 in practice remains challenging, HIP-LLM at least offers a practical “what-if” analysis tool
1146 by allowing assessors to explore how different prior assumptions might influence posterior
1147 reliability claims, for example: if my prior belief is A, then after observing evidence B, the
1148 posterior reliability would be C.

1149 *On the practicality of HIP-LLM.* HIP-LLM is designed to replace current LLM assessment
1150 practices that rely primarily on benchmark scores. As highlighted by the 5 gaps identified
1151 in the introduction, benchmark-based evaluations often ignore OPs, overlook dependencies
1152 between related domains, and cannot incorporate prior knowledge, limiting their relevance
1153 to real-world use. In contrast, HIP-LLM produces decision-oriented reliability assessments
1154 by integrating OPs, hierarchical domain dependencies, and epistemic uncertainty and pri-
1155 ors. Rather than reporting isolated scores (which, strictly speaking, cannot be regarded as
1156 “reliability” under standardized definitions such as ANSI [38]), HIP-LLM enables rigorous
1157 reliability claims such as “*the confidence that a model will complete a specified number of fu-*
1158 *ture tasks without failure under a given usage profile*”. This allows stakeholders to compare
1159 LLMs within specific domains, aggregate reliability across domains, or assess how reliability
1160 changes across different user groups or usage patterns. Such capabilities support informed
1161 offline decision-making, including model selection, risk assessment, and regulatory certifica-
1162 tion prior to deployment. While this work focuses on offline assessment, HIP-LLM could
1163 potentially be extended to online reliability monitoring with appropriate engineering opti-
1164 mizations for efficiency (e.g., caching and sliding data windows), which we leave for future
1165 work.

1166 6. Conclusion

1167 This paper presents HIP-LLM, a hierarchical Bayesian framework with imprecise probabil-
1168 ity for evaluating the reliability of LLMs. While HIP-LLM acknowledges the general-purpose
1169 nature of contemporary LLMs, it also adheres to the standardized definition of software re-
1170 liability which inherently depends on specific applications and operational contexts. We
1171 reconcile these seemingly conflicting aspects by introducing a multi-level hierarchical struc-
1172 ture of evaluation tasks that explicitly models dependencies among subdomains, integrates
1173 information through partial pooling, and incorporates OPs at multiple levels of abstraction.
1174 Moreover, as a Bayesian inference method, HIP-LLM respects the difficulty of selecting pri-
1175 ors. Instead of simply using non-informative priors, HIP-LLM provides the mechanism to

1176 embed imprecise prior knowledge and reports posterior envelopes (rather than single dis-
1177 tributions). Our experiments on common LLMs demonstrate that these features enable a
1178 more nuanced and standardized reliability estimation compared to existing benchmarks and
1179 state-of-the-art methods.

1180 While we believe HIP-LLM represents an important first step toward a principled reli-
1181 ability assessment framework for LLMs (and more broadly, for multimodal generative AI)
1182 several key limitations remain (cf. Section 5). These include the challenge of formally defining
1183 what constitutes a failure for an LLM, accurately estimating OPs, and systematically elic-
1184 iting imprecise prior knowledge. Moreover, the applicability of HIP-LLM is limited to reset or
1185 single-task usage scenarios conforming to the i.i.d. assumption (cf. Remark 1); and should
1186 not be overgeneralized to long-context or agentic LLM workflows (unless such workflows are
1187 modeled at a higher level of abstraction, treating the entire interaction trajectory as a single
1188 trial), which require dependency-aware reliability models—an important future work. More
1189 future extensions should aim to develop methods for handling uncertain and dynamically
1190 evolving OPs, and integrating more refined mechanisms for prior elicitation.

1191 **Acknowledgment**

1192 SK and XZ have received funding from the European Union’s EU Framework Program
1193 for Research and Innovation Europe Horizon (grant agreement No 101202457). XZ’s contri-
1194 bution is also supported by the UK EPSRC New Investigator Award [EP/Z536568/1]. SK’s
1195 contribution is supported by the UKRI Future Leaders Fellowship Grant [MR/S035176/1].

1196 Views and opinions expressed are those of the authors only and do not necessarily reflect
1197 those of the European Union or European Research Executive Agency (REA). Neither the
1198 European Union nor the granting authority can be held responsible for them.

1199 **References**

- 1200 [1] Z. Pang, Y. Luan, J. Chen, T. Li, Parinfoqpt: An LLM-based two-stage framework for
1201 reliability assessment of rotating machine under partial information, *Reliability Engi-
1202 neering & System Safety* 250 (2024) 110312.
- 1203 [2] X. Xiao, P. Chen, B. Qi, H. Zhao, J. Liang, J. Tong, H. Wang, Krail: A knowledge-driven
1204 framework for human reliability analysis integrating ideas-data and large language
1205 models, *Reliability Engineering & System Safety* (2025) 111585.
- 1206 [3] S. Zheng, K. Pan, J. Liu, Y. Chen, Empirical study on fine-tuning pre-trained large lan-
1207 guage models for fault diagnosis of complex systems, *Reliability Engineering & System
1208 Safety* 252 (2024) 110382.
- 1209 [4] U. Gohar, M. C. Hunter, R. R. Lutz, M. B. Cohen, Codefeater: Using llms to find
1210 defeaters in assurance cases, in: *proceedings of the 39th IEEE/ACM International Con-
1211 ference on Automated Software Engineering*, 2024, pp. 2262–2267.
- 1212 [5] A. Murugesan, I. Wong, J. Arias, R. Stroud, S. Varadarajan, E. Salazar, G. Gupta,
1213 R. Bloomfield, J. Rushby, Automating semantic analysis of system assurance cases using
1214 goal-directed asp, *Theory and Practice of Logic Programming* 24 (4) (2024) 805–824.

- 1215 [6] B. Sultan, L. Apvrille, AI-driven consistency of SysML diagrams, in: proceedings of the
1216 ACM/IEEE 27th International Conference on Model Driven Engineering Languages and
1217 Systems, 2024, pp. 149–159.
- 1218 [7] J. Yang, A. Prabhakar, K. Narasimhan, S. Yao, Intercode: Standardizing and bench-
1219 marking interactive coding with execution feedback, *Advances in Neural Information
1220 Processing Systems* 36 (2023) 23826–23854.
- 1221 [8] Y. Lai, C. Li, Y. Wang, T. Zhang, R. Zhong, L. Zettlemoyer, W.-t. Yih, D. Fried,
1222 S. Wang, T. Yu, DS-1000: A natural and reliable benchmark for data science code
1223 generation, in: *International Conference on Machine Learning*, PMLR, 2023, pp. 18319–
1224 18345.
- 1225 [9] Z. Yu, Y. Zhao, A. Cohan, X.-P. Zhang, HumanEval pro and MBPP pro: Evalu-
1226 ating Large Language Models on self-invoking code generation, *arXiv preprint
1227 arXiv:2412.21199* (2024).
- 1228 [10] E. Croxford, Y. Gao, N. Pellegrino, K. Wong, G. Wills, E. First, F. Liao, C. Goswami,
1229 B. Patterson, M. Afshar, Current and future state of evaluation of large language models
1230 for medical summarization tasks, *Npj health systems* 2 (1) (2025) 6.
- 1231 [11] S. Charalampidou, A. Zeleskidis, I. M. Dokas, Hazard analysis in the era of AI: Assessing
1232 the usefulness of ChatGPT4 in stpa hazard analysis, *Safety Science* 178 (2024) 106608.
- 1233 [12] G. K. Kaya, D. Bovell, M. Sujan, G. Braithwaite, Large language models powered system
1234 safety assessment: applying STPA and FRAM, *Safety Science* 191 (2025) 106960.
- 1235 [13] Y. Qi, X. Zhao, S. Khastgir, X. Huang, Safety analysis in the era of large language
1236 models: a case study of STPA using ChatGPT, *Machine Learning with Applications* 19
1237 (2025) 100622.
- 1238 [14] Z. Zhang, L. Lei, L. Wu, R. Sun, Y. Huang, C. Long, X. Liu, X. Lei, J. Tang, M. Huang,
1239 SafetyBench: Evaluating the safety of large language models, in: proceedings of the
1240 62nd Annual Meeting of the Association for Computational Linguistics, *ACL*, 2024, pp.
1241 15537–15553.
- 1242 [15] X. Liu, Y. Zhu, J. Gu, Y. Lan, C. Yang, Y. Qiao, Mm-safetybench: A benchmark for
1243 safety evaluation of multimodal large language models, in: *European Conference on
1244 Computer Vision*, Springer, 2024, pp. 386–403.
- 1245 [16] Y. Mou, S. Zhang, W. Ye, Sg-bench: Evaluating llm safety generalization across diverse
1246 tasks and prompt types, *Advances in Neural Information Processing Systems* 37 (2024)
1247 123032–123054.
- 1248 [17] Center for AI Standards and Innovation (CAISI), Evaluation of deepseek AI models,
1249 Report, National Institute of Standards and Technology (2025).
- 1250 [18] T. Yu, Y. Jing, X. Zhang, W. Jiang, W. Wu, Y. Wang, W. Hu, B. Du, D. Tao, Bench-
1251 marking reasoning robustness in large language models, *arXiv preprint arXiv:2503.04550*
1252 (2025).

- 1253 [19] S. Wang, P. Wang, T. Zhou, Y. Dong, Z. Tan, J. Li, Ceb: Compositional evaluation
1254 benchmark for fairness in Large Language Models, in: International Conference on Rep-
1255 resentation Learning, Vol. 2025, 2025, pp. 22627–22668.
- 1256 [20] Q. Li, J. Hong, C. Xie, J. Tan, R. Xin, J. Hou, X. Yin, Z. Wang, D. Hendrycks, Z. Wang,
1257 B. Li, B. He, D. Song, LLM-PBE: Assessing data privacy in Large Language Models,
1258 Proceedings of the VLDB Endow. 17 (11) (2024) 3201–3214.
- 1259 [21] A. Tamkin, M. McCain, K. Handa, E. Durmus, L. Lovitt, A. Rathi, S. Huang, A. Mount-
1260 field, J. Hong, S. Ritchie, et al., Clio: Privacy-preserving insights into real-world AI use,
1261 arXiv preprint arXiv:2412.13678 (2024).
- 1262 [22] B. Littlewood, L. Strigini, Software reliability and dependability: a roadmap, in: Pro-
1263 ceedings of the Conference on the Future of Software Engineering, 2000, pp. 175–188.
- 1264 [23] L. Strigini, B. Littlewood, Guidelines for statistical testing, Project Report
1265 PASCAN/WO6-CCN2/TN12, City University London (1997).
- 1266 [24] Y. Dong, W. Huang, V. Bharti, V. Cox, A. Banks, S. Wang, X. Zhao, S. Schewe,
1267 X. Huang, Reliability assessment and safety arguments for machine learning components
1268 in system assurance, ACM Transactions on Embedded Computing Systems 22 (3) (2023)
1269 1–48.
- 1270 [25] J. Musa, Operational profiles in software-reliability engineering, IEEE Software 10 (2)
1271 (1993) 14–32.
- 1272 [26] J. D. Musa, Adjusting measured field failure intensity for operational profile variation,
1273 in: Proceedings of IEEE Int. Symposium on Software Reliability Engineering, 1994, pp.
1274 330–333.
- 1275 [27] A. Chatterji, T. Cunningham, D. J. Deming, Z. Hitzig, C. Ong, C. Y. Shan, K. Wadman,
1276 How people use chatgpt, Tech. rep., National Bureau of Economic Research (2025).
- 1277 [28] M. R. Lyu, et al., Handbook of software reliability engineering, Vol. 222, IEEE computer
1278 society press Los Alamitos, 1996.
- 1279 [29] B. Littlewood, L. Strigini, Validation of ultra-high dependability for software-based sys-
1280 tems, Comm. of the ACM 36 (1993) 69–80.
- 1281 [30] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang,
1282 Y. Wang, et al., A survey on evaluation of large language models, ACM transactions on
1283 intelligent systems and technology 15 (3) (2024) 1–45.
- 1284 [31] L. Luettgau, H. Coppock, M. Dubois, C. Summerfield, C. Ududec, Hibayes: A Hi-
1285 erarchical Bayesian Modeling Framework for AI Evaluation Statistics, arXiv preprint
1286 arXiv:2505.05602 (2025).
- 1287 [32] E. Miller, Adding error bars to evals: A statistical approach to language model evalua-
1288 tions, arXiv preprint arXiv:2411.00640 (2024).

- 1289 [33] L. Strigini, A. Povyakalo, Software fault-freeness and reliability predictions, in: Com-
1290 puter Safety, Reliability, and Security, Vol. 8153 of LNCS, Springer Berlin Heidelberg,
1291 2013, pp. 106–117.
- 1292 [34] P. Bishop, A. Povyakalo, L. Strigini, Bootstrapping confidence in future safety from past
1293 safe operation, in: IEEE 33rd Int.Symp. on Software Reliability Engineering, IEEE,
1294 2022, pp. 97–108.
- 1295 [35] X. Zhao, A. Banks, J. Sharp, V. Robu, D. Flynn, M. Fisher, X. Huang, A safety frame-
1296 work for critical systems utilising deep neural networks, in: International Conference on
1297 Computer Safety, Reliability, and Security, Springer, 2020, pp. 244–259.
- 1298 [36] T. Augustin, F. P. Coolen, G. De Cooman, M. C. Troffaes, Introduction to imprecise
1299 probabilities, John Wiley & Sons, 2014.
- 1300 [37] M. C. Troffaes, Decision making under uncertainty using imprecise probabilities, Inter-
1301 national journal of approximate reasoning 45 (1) (2007) 17–29.
- 1302 [38] ANSI/IEEE, Standard glossary of software engineering terminology, Tech. rep., STD-
1303 729-1991 (1991).
- 1304 [39] IEC, IEC61508, Functional Safety of Electrical/ Electronic/Programmable Electronic
1305 Safety Related Systems, 2010.
- 1306 [40] C. Atwood, J. LaChance, H. Martz, D. Anderson, M. Englehardt, D. Whitehead,
1307 T. Wheeler, Handbook of parameter estimation for probabilistic risk assessment, Re-
1308 port NUREG/CR-6823, U.S. Nuclear Regulatory Commission, Washington, DC (2003).
- 1309 [41] B. Littlewood, J. Rushby, Reasoning about the reliability of diverse two-channel systems
1310 in which one channel is ‘possibly perfect’, IEEE Tran. on Software Engineering 38 (5)
1311 (2012) 1178–1194.
- 1312 [42] X. Zhao, B. Littlewood, A. Povyakalo, L. Strigini, D. Wright, Modeling the probability
1313 of failure on demand (pfd) of a 1-out-of-2 system in which one channel is “quasi-perfect”,
1314 Reliability Engineering & System Safety 158 (2017) 230–245.
- 1315 [43] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, J. Steinhardt, Measur-
1316 ing massive multitask language understanding, arXiv preprint arXiv:2009.03300 (2020).
- 1317 [44] G. Lai, Q. Xie, H. Liu, Y. Yang, E. Hovy, Race: Large-scale reading comprehension
1318 dataset from examinations, arXiv preprint arXiv:1704.04683 (2017).
- 1319 [45] G. Mialon, C. Fourrier, T. Wolf, Y. LeCun, T. Scialom, Gaia: a benchmark for general
1320 ai assistants, in: The Twelfth International Conference on Learning Representations,
1321 2023.
- 1322 [46] D. Rein, J. Becker, A. Deng, S. Nix, C. Canal, D. O’Connel, P. Arnott, R. Bloom,
1323 T. Broadley, K. Garcia, et al., HCAST: Human-calibrated autonomy software tasks,
1324 arXiv preprint arXiv:2503.17354 (2025).

- 1325 [47] A. Anthropic, The claude 3 model family: Opus, sonnet, haiku, Claude-3 Model Card
1326 1 (1) (2024) 4.
- 1327 [48] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Os-
1328 trow, A. Welihinda, A. Hayes, A. Radford, et al., Gpt-4o system card, arXiv preprint
1329 arXiv:2410.21276 (2024).
- 1330 [49] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. D. O. Pinto, J. Kaplan, H. Edwards,
1331 Y. Burda, N. Joseph, G. Brockman, et al., Evaluating large language models trained on
1332 code, arXiv preprint arXiv:2107.03374 (2021).
- 1333 [50] P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang,
1334 D. Narayanan, Y. Wu, A. Kumar, et al., Holistic evaluation of language models, arXiv
1335 preprint arXiv:2211.09110 (2022).
- 1336 [51] F. Shi, M. Suzgun, M. Freitag, X. Wang, S. Srivats, S. Vosoughi, H. W. Chung, Y. Tay,
1337 S. Ruder, D. Zhou, et al., Language models are multilingual chain-of-thought reasoners,
1338 arXiv preprint arXiv:2210.03057 (2022).
- 1339 [52] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, Y. Cao, React: Synergizing
1340 reasoning and acting in language models, in: International Conference on Learning
1341 Representations (ICLR), 2023.
- 1342 [53] Y. Liu, Y. Yao, J.-F. Ton, X. Zhang, R. Guo, H. Cheng, Y. Klochkov, M. F. Taufiq,
1343 H. Li, Trustworthy llms: a survey and guideline for evaluating large language models'
1344 alignment, arXiv preprint arXiv:2308.05374 (2023).
- 1345 [54] Y. Saleh, M. Abu Talib, Q. Nasir, F. Dakalbab, Evaluating large language models: a
1346 systematic review of efficiency, applications, and future directions, *Frontiers in Computer
1347 Science* 7 (2025) 1523699.
- 1348 [55] S. Liu, C. Li, J. Qiu, X. Zhang, F. Huang, L. Zhang, Y. Hei, P. S. Yu, The Scales
1349 of Justitia: A Comprehensive Survey on Safety Evaluation of LLMs, arXiv preprint
1350 arXiv:2506.11094 (2025).
- 1351 [56] H. Ye, J. Jin, Y. Xie, X. Zhang, G. Song, Large language model psychometrics: A system-
1352 atic review of evaluation, validation, and enhancement, arXiv preprint arXiv:2505.08245
1353 (2025).
- 1354 [57] V. Quach, A. Fisch, T. Schuster, A. Yala, J. H. Sohn, T. S. Jaakkola, R. Barzilay, Con-
1355 formal Language Modeling, in: The 12th International Conference on Learning Repre-
1356 sentations (ICLR'24).
- 1357 [58] Z. Wang, J. Duan, L. Cheng, Y. Zhang, Q. Wang, X. Shi, K. Xu, H. T. Shen, X. Zhu,
1358 ConU: Conformal Uncertainty in Large Language Models with Correctness Coverage
1359 Guarantees, in: Findings of the Association for Computational Linguistics: EMNLP
1360 2024, Association for Computational Linguistics, Miami, Florida, USA, 2024, pp. 6886–
1361 6898.

- 1362 [59] J. Berger, E. Moreno, L. Pericchi, M. Bayarri, J. Bernardo, J. Cano, J. Horra, J. Martín,
1363 D. Rios, B. Betrò, A. Dasgupta, P. Gustafson, L. Wasserman, J. Kadane, C. Srinivasan,
1364 M. Lavine, A. O’Hagan, W. Polasek, C. Robert, S. Sivaganesan, An overview of robust
1365 Bayesian analysis, *Test* 3 (1994) 5–124.
- 1366 [60] D. Insua, F. Ruggeri, *Robust Bayesian Analysis*, Lecture Notes in Statistics, Springer
1367 New York, 2012.
- 1368 [61] P. Bishop, R. Bloomfield, B. Littlewood, A. Povyakalo, D. Wright, Toward a formal-
1369 ism for conservative claims about the dependability of software-based systems, *IEEE*
1370 *Transactions on Software Engineering* 37 (5) (2010) 708–717.
- 1371 [62] X. Zhao, K. Salako, L. Strigini, V. Robu, D. Flynn, Assessing safety-critical systems
1372 from operational testing: A study on autonomous vehicles, *Information and Software*
1373 *Techn.* 128 (2020) 106393.
- 1374 [63] R. A. Chakherlou, L. Strigini, Impact of Prior Beliefs on Dependability Prediction for a
1375 Changed System Using Pre-change Operational Evidence, in: *24th International Con-*
1376 *ference on Software Quality, Reliability and Security*, IEEE, 2024, pp. 572–583.
- 1377 [64] L. V. Utkin, F. P. A. Coolen, Imprecise probabilistic inference for software run reliability
1378 growth models., *Journal of Uncertain Systems.* 12 (4) (2018) 292–308.
- 1379 [65] G. Walter, T. Augustin, Imprecision and prior-data conflict in generalized bayesian
1380 inference, *Journal of Statistical Theory and Practice* 3 (1) (2009) 255–271.
- 1381 [66] X. Zhao, R. Calinescu, S. Gerasimou, V. Robu, D. Flynn, Interval change-point detec-
1382 tion for runtime probabilistic model checking, in: *Proceedings of the 35th IEEE/ACM*
1383 *International Conference on Automated Software Engineering*, 2020, pp. 163–174.
- 1384 [67] S. Farquhar, J. Kossen, L. Kuhn, Y. Gal, Detecting hallucinations in large language
1385 models using semantic entropy, *Nature* 630 (8017) (2024) 625–630.
- 1386 [68] M. Dahl, V. Magesh, M. Suzgun, D. E. Ho, Hallucinating law: Legal mistakes with large
1387 language models are pervasive, *Law, regulation, and policy* (2024).
- 1388 [69] X. Huang, W. Ruan, et al., A survey of safety and trustworthiness of large language
1389 models through the lens of verification and validation, *Artificial Intelligence Review*
1390 57 (7) (2024) 175.
- 1391 [70] Y. Zhang, Y. Tang, W. Ruan, X. Huang, S. Khastgir, P. Jennings, X. Zhao, Protip:
1392 Probabilistic robustness verification on text-to-image diffusion models against stochastic
1393 perturbation, in: *European Conference on Computer Vision*, Springer, 2024, pp. 455–
1394 472.
- 1395 [71] K. Salako, X. Zhao, The unnecessary of assuming statistically independent tests in
1396 bayesian software reliability assessments, *IEEE Transactions on Software Engineering*
1397 49 (4) (2023) 2829–2838.

- 1398 [72] P. Bishop, A. Povyakalo, Deriving a frequentist conservative confidence bound for prob-
1399 ability of failure per demand for systems with different operational and test profiles,
1400 *Reliability Engineering & System Safety* 158 (2017) 246–253.
- 1401 [73] R. Pietrantuono, P. Popov, S. Russo, Reliability assessment of service-based software un-
1402 der operational profile uncertainty, *Reliability Engineering & System Safety* 204 (2020)
1403 107193.
- 1404 [74] P. Popov, Why black-box bayesian safety assessment of autonomous vehicles is prob-
1405 lematic and what can be done about it?, *IEEE Transactions on Intelligent Vehicles*In
1406 press (2025).
- 1407 [75] P. Popov, Dynamic safety assessment of Autonomous Vehicle based on Multivariate
1408 Bayesian Inference (DyAVSA), *Journal of Reliable Intelligent Environments* 11 (3)
1409 (2025) 1–23.
- 1410 [76] W. Huang, X. Zhao, A. Banks, V. Cox, X. Huang, Hierarchical distribution-aware testing
1411 of deep learning, *ACM Transactions on Software Engineering and Methodology* 33 (2)
1412 (2023) 1–35.
- 1413 [77] C. Smidts, C. Mutha, M. Rodríguez, M. J. Gerber, Software testing with an operational
1414 profile: OP definition, *ACM Computing Surveys (CSUR)* 46 (3) (2014) 1–39.
- 1415 [78] C. Clark, K. Lee, M.-W. Chang, T. Kwiatkowski, M. Collins, K. Toutanova, BoolQ:
1416 Exploring the Surprising Difficulty of Natural Yes/No Questions, in: *Proceedings of*
1417 *the Conference of the North American Chapter of the Association for Computational*
1418 *Linguistics: Human Language Technologies, Vol. 1, 2019*, pp. 2924–2936.
- 1419 [79] R. A. Thayer, M. Lipow, E. C. Nelson, *Software Reliability*, North-Holland, 1978.
- 1420 [80] J. W. Duran, S. C. Ntafos, An Evaluation of Random Testing, *IEEE Transactions on*
1421 *Software Engineering SE-10* (4) (1984) 438–444.
- 1422 [81] S. Chen, S. Mills, A binary Markov process model for random testing, *IEEE Transactions*
1423 *on Software Engineering* 22 (3) (1996) 218–223.
- 1424 [82] K. Goseva-Popstojanova, K. S. Trivedi, Failure correlation in software reliability models,
1425 *IEEE Transactions on Reliability* 49 (1) (2000) 37–48.
- 1426 [83] Y. Dai, M. Xie, K. Poh, Modeling and analysis of correlated software failures of multiple
1427 types, *IEEE Transactions on Reliability* 54 (1) (2005) 100–106.
- 1428 [84] A. Bondavalli, S. Chiaradonna, F. Di Giandomenico, L. Strigini, Dependability models
1429 for iterative software considering correlation between successive inputs, in: *Proceedings*
1430 *of IEEE International. Computer Performance and Dependability Symposium, IEEE,*
1431 *Erlangen, Germany, 1995*, pp. 13–21.
- 1432 [85] K. Salako, X. Zhao, Demonstrating software reliability using possibly correlated tests:
1433 Insights from a conservative bayesian approach, *Quality and Reliability Engineering*
1434 *International* 40 (3) (2024) 1197–1220.

- 1435 [86] X. Zhao, W. Huang, A. Banks, V. Cox, D. Flynn, S. Schewe, X. Huang, Assessing
1436 the reliability of deep learning classifiers through robustness evaluation and operational
1437 profiles, in: AISafety-IJCAI'21, Vol. 2916, CEUR-WS, 2021, p. 16.
- 1438 [87] X. Zhao, R. Aghazadeh-Chakherlou, C.-H. Cheng, P. Popov, L. Strigini, On the need for
1439 a statistical foundation in scenario-based testing of autonomous vehicles, in: Proceedings
1440 of the IEEE International Conference on Intelligent Transportation Systems (ITSC),
1441 Gold Coast, Australia, 2025.
- 1442 [88] P. Bishop, R. Bloomfield, B. Littlewood, A. Povyakalo, D. Wright, Toward a formal-
1443 ism for conservative claims about the dependability of software-based systems, IEEE
1444 Transactions on Software Engineering 37 (5) (2011) 708–717.
- 1445 [89] X. Zhao, V. Robu, D. Flynn, K. Salako, L. Strigini, Assessing the safety and reliability
1446 of autonomous vehicles from road testing, in: International Symposium on Software
1447 Reliability Engineering (ISSRE), 2019, pp. 13–23.
- 1448 [90] G. Walter, L. J. M. Aslett, F. P. A. Coolen, Bayesian nonparametric system reliability
1449 using sets of priors, International Journal of Approximate Reasoning 80 (2017) 67–88.
- 1450 [91] X. Zhao, V. Robu, D. Flynn, F. Dinmohammadi, M. Fisher, M. Webster, Probabilistic
1451 model checking of robots deployed in extreme environments, in: Proceedings of the 33rd
1452 AAI Conference on Artificial Intelligence, Vol. 33, Honolulu, Hawaii, USA, 2019, pp.
1453 8076–8084.
- 1454 [92] X. Zhao, S. Gerasimou, R. Calinescu, C. Imrie, V. Robu, D. Flynn, Bayesian learning for
1455 the robust verification of autonomous robots, Communications Engineering 3 (1) (2024)
1456 18.

1457 **Appendix A. Mathematical Derivations**

1458 This appendix contains the proofs and supplementary information for the theorems pre-
 1459 sented in the main text.

1460 *Appendix A.1. Common Framework and Notation*

1461 In a LLM, for a domain D_i with subdomains S_{i1}, \dots, S_{in_i} . For each subdomain $j =$
 1462 $1, \dots, n_i$:

- 1463 • **Parameter:** $\boldsymbol{\theta}_i = (\theta_{i1}, \dots, \theta_{in_i}) \in (0, 1)^{n_i}$ (subdomain reliability)
- 1464 • **Parameter:** $n^F =$ number of trials in the future
- 1465 • **Data:** $C_i = \{(C_{ik}, N_{ik})\}_{k=1}^{n_i}$ (observed correct responses and trial counts)
- 1466 • **Hierarchical prior.**

$$\begin{aligned} C_{ij} \mid \theta_{ij}, N_{ij} &\sim \text{Binomial}(N_{ij}, \theta_{ij}), \\ \theta_{ij} \mid \mu_i, \nu_i &\stackrel{\text{iid}}{\sim} \text{Beta}(\mu_i \nu_i, (1 - \mu_i) \nu_i), \\ \mu_i &\sim \text{Beta}(a_i, b_i), \quad \nu_i \sim \text{Gamma}(c_i, \text{rate} = d_i), \end{aligned}$$

1467 with $\mu_i \in (0, 1)$, $\nu_i > 0$ (domain-level parameters), and hyper- hyperparameters $h_i =$
 1468 (a_i, b_i, c_i, d_i) .

- 1469 • **Admissible hyperparameters (imprecise prior).**

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}].$$

- 1470 • **Common Likelihood:**

$$L(\boldsymbol{\theta}_i) = Pr(C_i \mid \boldsymbol{\theta}_i) = \prod_{j=1}^{n_i} \binom{N_{ij}}{C_{ij}} \theta_{ij}^{C_{ij}} (1 - \theta_{ij})^{N_{ij} - C_{ij}}$$

1471 *Appendix A.2. Proof of Theorem 1*

1472 **Theorem 1** For subdomain S_{ij} in domain D_i , let $C_i = \{(C_{ik}, N_{ik})\}_{k=1}^{n_i}$ be the observed
 1473 data. Let the admissible set of hyperparameters be

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}],$$

1474 and write $h_i = (a_i, b_i, c_i, d_i)$. Then, for any $h_i \in \mathcal{A}_i$, the marginal posterior density of θ_{ij} is

$$Pr(\theta_{ij} \mid C_i, h_i) = \frac{f_{\text{marg}}(\theta_{ij}, C_i; h_i)}{Z_{\text{marg}}(h_i)},$$

1475 where f_{marg} (unnormalized posterior) and Z_{marg} (normalizing constant) are

$$\begin{aligned} f_{\text{marg}}(\theta_{ij}, C_i; h_i) = \\ \int_0^1 \int_0^\infty \left[\prod_{k \neq j} \int_0^1 d\theta_{ik} \right] L(\boldsymbol{\theta}_i) Pr(\boldsymbol{\theta}_i \mid \mu_i, \nu_i) Pr(\mu_i, \nu_i \mid h_i) d\mu_i d\nu_i \end{aligned}$$

1476

$$Z_{\text{marg}}(h_i) = \int_0^1 \int_0^\infty Pr(C_i | \mu_i, \nu_i) Pr(\mu_i, \nu_i | h_i) d\mu_i d\nu_i,$$

1477 with $L(\boldsymbol{\theta}_i) = Pr(C_i | \boldsymbol{\theta}_i)$.

1478 The imprecise marginal posterior is characterized by the lower/upper envelopes

$$\underline{Pr}(\theta_{ij} | C_i) = \inf_{h_i \in \mathcal{A}_i} Pr(\theta_{ij} | C_i, h_i), \quad \overline{Pr}(\theta_{ij} | C_i) = \sup_{h_i \in \mathcal{A}_i} Pr(\theta_{ij} | C_i, h_i).$$

1479 *Proof of Theorem 1.* We proceed step by step.1480 **Step 1 — Likelihood and hierarchical prior densities.** The joint data likelihood
1481 factorizes over subdomains:

$$L(\boldsymbol{\theta}_i) = Pr(C_i | \boldsymbol{\theta}_i) = \prod_{k=1}^{n_i} \binom{N_{ik}}{C_{ik}} \theta_{ik}^{C_{ik}} (1 - \theta_{ik})^{N_{ik} - C_{ik}}. \quad (\text{A.1})$$

1482 Given (μ_i, ν_i) , the prior on $\boldsymbol{\theta}_i$ factorizes as

$$Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i) = \prod_{k=1}^{n_i} \frac{\Gamma(\nu_i)}{\Gamma(\mu_i \nu_i) \Gamma((1 - \mu_i) \nu_i)} \theta_{ik}^{\mu_i \nu_i - 1} (1 - \theta_{ik})^{(1 - \mu_i) \nu_i - 1}.$$

1483 Given the domain-level parameters (μ_i, ν_i) , the subdomain reliabilities are conditionally
1484 independent and identically distributed:

$$\theta_{ij} | \mu_i, \nu_i \stackrel{\text{iid}}{\sim} \text{Beta}(\mu_i \nu_i, (1 - \mu_i) \nu_i), \quad j = 1, \dots, n_i.$$

1485 Here $\mu_i \in (0, 1)$ is the prior mean and $\nu_i > 0$ is the prior strength (concentration).1486 For each k , the Beta pdf (Single-coordinate density) is:

$$Pr(\theta_{ik} | \mu_i, \nu_i) = \frac{\Gamma(\nu_i)}{\Gamma(\mu_i \nu_i) \Gamma((1 - \mu_i) \nu_i)} \theta_{ik}^{\mu_i \nu_i - 1} (1 - \theta_{ik})^{(1 - \mu_i) \nu_i - 1}, \quad \theta_{ik} \in (0, 1).$$

1487 Since $\{\theta_{ik}\}_{k=1}^{n_i}$ are conditionally independent given (μ_i, ν_i) ,

$$\begin{aligned} Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i) &= \prod_{k=1}^{n_i} Pr(\theta_{ik} | \mu_i, \nu_i) \\ &= \prod_{k=1}^{n_i} \text{Beta}(\theta_{ik} | \mu_i \nu_i, (1 - \mu_i) \nu_i) \\ &= \left(\frac{\Gamma(\nu_i)}{\Gamma(\mu_i \nu_i) \Gamma((1 - \mu_i) \nu_i)} \right)^{n_i} \prod_{k=1}^{n_i} \theta_{ik}^{\mu_i \nu_i - 1} (1 - \theta_{ik})^{(1 - \mu_i) \nu_i - 1} \end{aligned} \quad (\text{A.2})$$

1488 Given hyperparameters $h_i = (a_i, b_i, c_i, d_i)$, we take μ_i and ν_i to be independent:

$$Pr(\mu_i, \nu_i | h_i) = Pr(\mu_i | h_i) Pr(\nu_i | h_i).$$

1489 We assume:

$$\mu_i \sim \text{Beta}(a_i, b_i), \quad \mu_i \in (0, 1), \quad \nu_i \sim \text{Gamma}(c_i, \text{rate} = d_i), \quad \nu_i > 0.$$

1490 Explicitly,

$$Pr(\mu_i | h_i) = \frac{\Gamma(a_i + b_i)}{\Gamma(a_i)\Gamma(b_i)} \mu_i^{a_i-1} (1 - \mu_i)^{b_i-1}, \quad Pr(\nu_i | h_i) = \frac{d_i^{c_i}}{\Gamma(c_i)} \nu_i^{c_i-1} e^{-d_i \nu_i}$$

1491 Thus,

$$\begin{aligned} Pr(\mu_i, \nu_i | h_i) &= \text{Beta}(\mu_i | a_i, b_i) \text{Gamma}(\nu_i | c_i, \text{rate} = d_i) \\ &= \frac{\Gamma(a_i + b_i)}{\Gamma(a_i)\Gamma(b_i)} \frac{d_i^{c_i}}{\Gamma(c_i)} \mu_i^{a_i-1} (1 - \mu_i)^{b_i-1} \nu_i^{c_i-1} e^{-d_i \nu_i} \\ &= \frac{1}{B(a_i, b_i)} \frac{d_i^{c_i}}{\Gamma(c_i)} \mu_i^{a_i-1} (1 - \mu_i)^{b_i-1} \nu_i^{c_i-1} e^{-d_i \nu_i} \end{aligned} \quad (\text{A.3})$$

1492 where $B(a_i, b_i) = \frac{\Gamma(a_i)\Gamma(b_i)}{\Gamma(a_i + b_i)}$.

1493 The hyperprior factorizes:

$$Pr(\mu_i, \nu_i | h_i) = \frac{\Gamma(a_i + b_i)}{\Gamma(a_i)\Gamma(b_i)} \mu_i^{a_i-1} (1 - \mu_i)^{b_i-1} \cdot \frac{d_i^{c_i}}{\Gamma(c_i)} \nu_i^{c_i-1} e^{-d_i \nu_i}.$$

1494 **Step 2 — Normalizing constant.** By definition,

$$Z_{\text{marg}}(C_i, h_i) = \int_0^1 Pr(\theta_{ij} | C_i, h_i) d\theta_{ij} = \int_0^1 \int_0^\infty Pr(C_i | \mu_i, \nu_i) Pr(\mu_i, \nu_i | h_i) d\mu_i d\nu_i,$$

1495 where

$$Pr(C_i | \mu_i, \nu_i) = \int_{(0,1)^{n_i}} L(\theta_i) Pr(\theta_i | \mu_i, \nu_i) d\theta_i.$$

1496 Start from the definition:

$$Pr(\theta_{ij} | C_i, h_i) = \frac{f_{\text{marg}}(\theta_{ij}, C_i; h_i)}{Z_{\text{marg}}(C_i, h_i)}$$

1497 where

$$\begin{aligned} f_{\text{marg}}(\theta_{ij}, C_i; h_i) &= \\ &= \int_0^1 \int_0^\infty \left[\prod_{k \neq j} \int_0^1 d\theta_{ik} \right] L(\theta_i) Pr(\theta_i | \mu_i, \nu_i) Pr(\mu_i, \nu_i | h_i) d\mu_i d\nu_i \end{aligned}$$

1498 The posterior integrates to 1:

$$\int_0^1 Pr(\theta_{ij} | C_i, h_i) d\theta_{ij} = 1 \implies \frac{1}{Z_{\text{marg}}(C_i, h_i)} \int_0^1 f_{\text{marg}}(\theta_{ij}, C_i; h_i) d\theta_{ij} = 1$$

1499 Hence

$$Z_{\text{marg}}(C_i, h_i) = \int_0^1 f_{\text{marg}}(\theta_{ij}, C_i; h_i) d\theta_{ij}$$

1500 By expanding the integral:

$$\begin{aligned}
Z_{\text{marg}}(C_i, h_i) &= \int_0^1 \int_0^1 \int_0^\infty \left[\prod_{k \neq j} \int_0^1 d\theta_{ik} \right] L(\boldsymbol{\theta}_i) Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i) Pr(\mu_i, \nu_i | h_i) d\mu_i d\nu_i d\theta_{ij} \\
&= \int_0^1 \int_0^\infty \left[\int_{(0,1)^{n_i}} L(\boldsymbol{\theta}_i) Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i) d\boldsymbol{\theta}_i \right] Pr(\mu_i, \nu_i | h_i) d\mu_i d\nu_i. \tag{A.4}
\end{aligned}$$

1501 By substituting Eqs. A.1, A.2, A.3 into A.4 we have:

$$\begin{aligned}
Z_{\text{marg}}(C_i, h_i) &= \left[\prod_{j=1}^{n_i} \binom{N_{ij}}{C_{ij}} \right] \frac{d_i^{c_i}}{B(a_i, b_i) \Gamma(c_i)} \int_0^1 \int_0^\infty \mu_i^{a_i-1} (1-\mu_i)^{b_i-1} \nu_i^{c_i-1} e^{-d_i \nu_i} \\
&\quad \times \left(\frac{\Gamma(\nu_i)}{\Gamma(\alpha_i) \Gamma(\beta_i)} \right)^{n_i} \prod_{j=1}^{n_i} \frac{\Gamma(C_{ij} + \alpha_i) \Gamma(N_{ij} - C_{ij} + \beta_i)}{\Gamma(N_{ij} + \nu_i)} d\nu_i d\mu_i \tag{A.5}
\end{aligned}$$

1502 **Step 3 — Imprecise posterior.** For any fixed $h_i = (a_i, b_i, c_i, d_i) \in \mathcal{A}_i$, Steps 1–2 give
1503 the precise (pointwise) posterior density

$$Pr(\theta_{ij} | C_i, h_i) = \frac{f_{\text{marg}}(\theta_{ij}, C_i; h_i)}{Z_{\text{marg}}(C_i, h_i)},$$

1504 The hyperparameter vector h_i is not known exactly but only to lie within the admissible
1505 hyperrectangle

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}].$$

1506 Hence the family $\mathcal{F}_i = \{Pr(\theta_{ij} | C_i, h_i) : h_i \in \mathcal{A}_i\}$ represents all posterior densities consistent
1507 with our prior ignorance about h_i .

1508 Following the theory of imprecise probabilities, we define the *lower* and *upper* posterior
1509 densities (envelopes) as the pointwise infimum and supremum over this family:

$$\underline{Pr}(\theta_{ij} | C_i) = \inf_{h_i \in \mathcal{A}_i} Pr(\theta_{ij} | C_i, h_i), \quad \overline{Pr}(\theta_{ij} | C_i) = \sup_{h_i \in \mathcal{A}_i} Pr(\theta_{ij} | C_i, h_i).$$

1510 Operationally, these envelopes are obtained by evaluating the closed-form posterior $Pr(\theta_{ij} |$
1511 $C_i, h_i)$ at the extremal corners of \mathcal{A}_i or through numerical optimization if the extrema occur
1512 in the interior. The resulting pair $(\underline{Pr}, \overline{Pr})$ bounds all admissible precise posteriors and fully
1513 characterizes the imprecise posterior belief about θ_{ij} given the uncertainty in (a_i, b_i, c_i, d_i) .
1514 □

1515 Appendix A.3. Proof of Theorem 2

1516 **Theorem 2** For domain D_i with local OP weights Ω_{ij} (where $\sum_{j=1}^{n_i} \Omega_{ij} = 1$), let $p_i =$
1517 $\sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$ be the domain-level non-failure probability. Define the admissible set of hyper-
1518 hyper-parameters

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}],$$

1519 and write $h_i = (a_i, b_i, c_i, d_i)$.

1520 Then, for any $h_i \in \mathcal{A}_i$, the posterior distribution of p_i is characterized by its CDF:

$$\begin{aligned} F_{p_i}(t | C_i, h_i) &= \Pr(p_i \leq t | C_i, h_i) \\ &= \int_0^1 \int_0^\infty F_{p_i}(t | \mu_i, \nu_i, C_i) Pr(\mu_i, \nu_i | C_i, h_i) d\mu_i d\nu_i \end{aligned}$$

1521 where

- 1522 • $F_{p_i}(t | \mu_i, \nu_i, C_i)$ is the conditional CDF of $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$ given that $\theta_{ij} | \mu_i, \nu_i, C_i \stackrel{\text{ind.}}{\sim}$
- 1523 Beta($C_{ij} + \mu_i \nu_i, N_{ij} - C_{ij} + (1 - \mu_i) \nu_i$) for $j = 1, \dots, n_i$,
- 1524 • $Pr(\mu_i, \nu_i | C_i, h_i)$ is the hyper-posterior obtained via Bayes' rule:

$$\begin{aligned} Pr(\mu_i, \nu_i | C_i, h_i) &= \\ &= \frac{Pr(C_i | \mu_i, \nu_i) \text{Beta}(\mu_i | a_i, b_i) \text{Gamma}(\nu_i | c_i, \text{rate} = d_i)}{\int_0^1 \int_0^\infty Pr(C_i | \mu, \nu) \text{Beta}(\mu | a_i, b_i) \text{Gamma}(\nu | c_i, \text{rate} = d_i) d\mu d\nu} \end{aligned}$$

1525 The imprecise domain posterior is characterized by CDF envelopes:

$$\underline{F}_{p_i}(t | C_i) = \inf_{h_i \in \mathcal{A}_i} F_{p_i}(t | C_i, h_i), \quad \overline{F}_{p_i}(t | C_i) = \sup_{h_i \in \mathcal{A}_i} F_{p_i}(t | C_i, h_i).$$

1526 The domain reliability $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$ is a weighted sum of subdomain reliabilities,
 1527 whose joint posterior distribution $Pr(\boldsymbol{\theta}_i | C_i, h_i)$ has no tractable closed form. However, the
 1528 hierarchical structure of our model provides a natural decomposition: we can express this
 1529 joint posterior as a mixture over the shared hyperparameters (μ_i, ν_i) , where conditionally
 1530 on these hyperparameters, the subdomains become independent with Beta posteriors. This
 1531 decomposition, formalized in Lemma 7, is the key to deriving the domain-level CDF. The
 1532 proof proceeds by (1) expressing the domain CDF as an integral of the joint posterior over a
 1533 weighted-sum constraint region, (2) applying the hierarchical decomposition from the lemma,
 1534 (3) exchanging the order of integration via Fubini's theorem, and (4) recognizing the inner
 1535 integral as a conditional CDF, yielding a tractable mixture representation.

1536 **Lemma 7** (Hierarchical decomposition of subdomain joint posterior). *Under the hierarchical*
 1537 *model with hyperparameters $h_i = (a_i, b_i, c_i, d_i)$, the joint posterior of subdomain reliabilities*
 1538 *$\boldsymbol{\theta}_i = (\theta_{i1}, \dots, \theta_{in_i})$ admits the decomposition*

$$Pr(\boldsymbol{\theta}_i | C_i, h_i) = \int_0^1 \int_0^\infty Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i, C_i) Pr(\mu_i, \nu_i | C_i, h_i) d\mu_i d\nu_i, \quad (\text{A.6})$$

1539 where:

- 1540 1. The conditional distribution factorizes into independent Betas:

$$Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i, C_i) = \prod_{j=1}^{n_i} \text{Beta}(\theta_{ij} | C_{ij} + \mu_i \nu_i, N_{ij} - C_{ij} + (1 - \mu_i) \nu_i).$$

1541 2. The hyper-posterior is obtained via Bayes' rule:

$$\begin{aligned} Pr(\mu_i, \nu_i | C_i, h_i) &= \frac{Pr(C_i | \mu_i, \nu_i) Pr(\mu_i, \nu_i | h_i)}{Pr(C_i | h_i)}, \\ Pr(\mu_i, \nu_i | h_i) &= Beta(\mu_i | a_i, b_i) \times Gamma(\nu_i | c_i, rate = d_i), \end{aligned} \quad (\text{A.7})$$

1542 where the marginal likelihood is

$$Pr(C_i | \mu_i, \nu_i) = \prod_{j=1}^{n_i} \binom{N_{ij}}{C_{ij}} \frac{B(C_{ij} + \mu_i \nu_i, N_{ij} - C_{ij} + (1 - \mu_i) \nu_i)}{B(\mu_i \nu_i, (1 - \mu_i) \nu_i)}, \quad (\text{A.8})$$

1543 and the evidence is

$$Pr(C_i | h_i) = \int_0^1 \int_0^\infty Pr(C_i | \mu, \nu) Beta(\mu | a_i, b_i) Gamma(\nu | c_i, rate = d_i) d\mu d\nu. \quad (\text{A.9})$$

1544 *Proof of Lemma 7.* The decomposition follows directly from the law of total probability and
 1545 the hierarchical prior structure. By the tower property of conditional expectation applied to
 1546 densities:

$$\begin{aligned} Pr(\boldsymbol{\theta}_i | C_i, h_i) &= \\ & \int_0^1 \int_0^\infty Pr(\boldsymbol{\theta}_i, \mu_i, \nu_i | C_i, h_i) d\mu_i d\nu_i = \\ & \int_0^1 \int_0^\infty Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i, C_i) Pr(\mu_i, \nu_i | C_i, h_i) d\mu_i d\nu_i \end{aligned}$$

1547 The conditional independence of $\{\theta_{ij}\}$ given (μ_i, ν_i, C_i) follows from the hierarchical prior:

$$\theta_{ij} | \mu_i, \nu_i \stackrel{\text{iid}}{\sim} Beta(\mu_i \nu_i, (1 - \mu_i) \nu_i), \quad j = 1, \dots, n_i,$$

1548 combined with independent Binomial likelihoods $C_{ij} | \theta_{ij}, N_{ij} \sim \text{Binomial}(N_{ij}, \theta_{ij})$. By
 1549 conjugacy, the Beta-Binomial update yields:

$$\theta_{ij} | \mu_i, \nu_i, C_i \sim Beta(C_{ij} + \mu_i \nu_i, N_{ij} - C_{ij} + (1 - \mu_i) \nu_i).$$

1550 The hyper-posterior (A.7) follows from Bayes' rule. The marginal likelihood (A.8) is
 1551 the Beta-Binomial distribution for each subdomain, obtained by integrating the Binomial
 1552 likelihood against the Beta prior. The evidence (A.9) is the normalization constant ensuring
 1553 the hyper-posterior integrates to one. \square

1554 *Proof of Theorem 2.* We derive the posterior CDF of the domain reliability $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$
 1555 using the hierarchical decomposition from Lemma 7.

1556 **Step 1: CDF definition.** For any $t \in [0, 1]$, the posterior CDF of p_i is

$$F_{p_i}(t | C_i, h_i) = \Pr(p_i \leq t | C_i, h_i) = \Pr\left(\sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij} \leq t \mid C_i, h_i\right).$$

1557 Define the constraint region

$$\mathcal{R}_i(t) := \left\{ \boldsymbol{\theta}_i \in (0, 1)^{n_i} : \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij} \leq t \right\}.$$

1558 Then

$$F_{p_i}(t | C_i, h_i) = \int_{\mathcal{R}_i(t)} Pr(\boldsymbol{\theta}_i | C_i, h_i) d\boldsymbol{\theta}_i.$$

1559 **Step 2: Apply hierarchical decomposition.** Substituting the decomposition from
1560 Lemma 7 (equation (A.6)):

$$F_{p_i}(t | C_i, h_i) = \int_{\mathcal{R}_i(t)} \left[\int_0^1 \int_0^\infty Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i, C_i) Pr(\mu_i, \nu_i | C_i, h_i) d\mu_i d\nu_i \right] d\boldsymbol{\theta}_i.$$

1561 **Step 3: Change of integration order.** By Fubini's theorem (applicable since all inte-
1562 grands are non-negative and integrate to finite values):

$$F_{p_i}(t | C_i, h_i) = \int_0^1 \int_0^\infty \left[\int_{\mathcal{R}_i(t)} Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i, C_i) d\boldsymbol{\theta}_i \right] Pr(\mu_i, \nu_i | C_i, h_i) d\mu_i d\nu_i.$$

1563 **Step 4: Conditional CDF.** The inner integral is the conditional CDF of p_i given (μ_i, ν_i)
1564 and C_i :

$$F_{p_i}(t | \mu_i, \nu_i, C_i) = \int_{\mathcal{R}_i(t)} Pr(\boldsymbol{\theta}_i | \mu_i, \nu_i, C_i) d\boldsymbol{\theta}_i = \Pr \left(\sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij} \leq t \mid \mu_i, \nu_i, C_i \right),$$

1565 where $\theta_{ij} | \mu_i, \nu_i, C_i$ are independent Beta random variables as specified in Lemma 7.

1566 **Step 5: Mixture representation.** Combining Steps 3 and 4:

$$F_{p_i}(t | C_i, h_i) = \int_0^1 \int_0^\infty F_{p_i}(t | \mu_i, \nu_i, C_i) Pr(\mu_i, \nu_i | C_i, h_i) d\mu_i d\nu_i.$$

1567 where

- 1568 • $F_{p_i}(t | \mu_i, \nu_i, C_i)$ is the conditional CDF of $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$ given that $\theta_{ij} | \mu_i, \nu_i, C_i \stackrel{\text{ind.}}{\sim}$
1569 $\text{Beta}(C_{ij} + \mu_i \nu_i, N_{ij} - C_{ij} + (1 - \mu_i) \nu_i)$ for $j = 1, \dots, n_i$,
- 1570 • $Pr(\mu_i, \nu_i | C_i, h_i)$ is the hyper-posterior obtained via Bayes' rule:

$$Pr(\mu_i, \nu_i | C_i, h_i) = \frac{Pr(C_i | \mu_i, \nu_i) \text{Beta}(\mu_i | a_i, b_i) \text{Gamma}(\nu_i | c_i, \text{rate} = d_i)}{\int_0^1 \int_0^\infty Pr(C_i | \mu, \nu) \text{Beta}(\mu | a_i, b_i) \text{Gamma}(\nu | c_i, \text{rate} = d_i) d\mu d\nu}$$

1571 This expresses the domain posterior CDF as a weighted average of conditional CDFs,
 1572 where the weights are given by the hyper-posterior $Pr(\mu_i, \nu_i | C_i, h_i)$.

1573 **Step 6: Imprecise probability bounds.** For each $h_i \in \mathcal{A}_i$, the mixture formula defines
 1574 a valid CDF $F_{p_i}(\cdot | C_i, h_i)$. The imprecise posterior is characterized by pointwise envelopes
 1575 over the admissible set:

$$\underline{F}_{p_i}(t | C_i) = \inf_{h_i \in \mathcal{A}_i} F_{p_i}(t | C_i, h_i), \quad \overline{F}_{p_i}(t | C_i) = \sup_{h_i \in \mathcal{A}_i} F_{p_i}(t | C_i, h_i), \quad t \in [0, 1]$$

1576

□

1577 *Appendix A.4. Proof of Theorem 3*

1578 **Theorem 3** For the LLM system with domain weights W_i (where $\sum_{i=1}^m W_i = 1$), let
 1579 $p_L = \sum_{i=1}^m W_i p_i$ be the LLM-level failure probability and $\text{data} = \{C_1, \dots, C_m\}$ the observed
 1580 data across all domains. Assume cross-domain independence.

1581 Define the domain-level admissible sets

$$\mathcal{A}_i = [a_i^{\min}, a_i^{\max}] \times [b_i^{\min}, b_i^{\max}] \times [c_i^{\min}, c_i^{\max}] \times [d_i^{\min}, d_i^{\max}]$$

1582 and write $h_i = (a_i, b_i, c_i, d_i)$ for $i = 1, \dots, m$. Define the LLM-level admissible set as the
 1583 Cartesian product

$$\mathcal{A}_{\text{LLM}} = \mathcal{A}_1 \times \dots \times \mathcal{A}_m,$$

1584 and collect the domain hyperparameters as $\mathcal{H} = (h_1, \dots, h_m) \in \mathcal{A}_{\text{LLM}}$.

1585 Then, for any $\mathcal{H} \in \mathcal{A}_{\text{LLM}}$, the posterior distribution of p_L is characterized by its CDF:

$$F_{p_L}(t | \text{data}, \mathcal{H}) = Pr(p_L \leq t | \text{data}, \mathcal{H}) = \int \dots \int G(t | \{\mu_i, \nu_i\}_{i=1}^m, \text{data}) \prod_{i=1}^m Pr(\mu_i, \nu_i | C_i, h_i) \prod_{i=1}^m d\mu_i d\nu_i$$

1586 where

1587 • $G(t | \{\mu_i, \nu_i\}_{i=1}^m, \text{data})$ is the conditional CDF of $p_L = \sum_{i=1}^m W_i p_i$ given all hyperpa-
 1588 rameters, defined as

$$G(t | \{\mu_i, \nu_i\}_{i=1}^m, \text{data}) = \int_{\mathcal{R}_L(t)} \prod_{i=1}^m f_{p_i}(p_i | \mu_i, \nu_i, C_i) dp_1 \dots dp_m,$$

1589 where $\mathcal{R}_L(t) := \{(p_1, \dots, p_m) \in (0, 1)^m : \sum_{i=1}^m W_i p_i \leq t\}$, and $f_{p_i}(\cdot | \mu_i, \nu_i, C_i)$ is the
 1590 conditional density of $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$ under $\theta_{ij} | \mu_i, \nu_i, C_i \stackrel{\text{ind.}}{\sim} \text{Beta}(C_{ij} + \mu_i \nu_i, N_{ij} -$
 1591 $C_{ij} + (1 - \mu_i) \nu_i)$,

1592 • $Pr(\mu_i, \nu_i | C_i, h_i)$ is the domain-level hyper-posterior for domain i :

$$Pr(\mu_i, \nu_i | C_i, h_i) = \frac{Pr(C_i | \mu_i, \nu_i) \text{Beta}(\mu_i | a_i, b_i) \text{Gamma}(\nu_i | c_i, \text{rate} = d_i)}{\int_0^1 \int_0^\infty Pr(C_i | \mu, \nu) \text{Beta}(\mu | a_i, b_i) \text{Gamma}(\nu | c_i, \text{rate} = d_i) d\mu d\nu}$$

1593 • Cross-domain independence ensures $Pr(\{\mu_i, \nu_i\}_{i=1}^m \mid \text{data}, \mathcal{H}) = \prod_{i=1}^m Pr(\mu_i, \nu_i \mid C_i, h_i)$.

1594 The imprecise LLM posterior is characterized by CDF envelopes:

$$\begin{aligned} \underline{F}_{p_L}(t \mid \text{data}) &= \inf_{\mathcal{H} \in \mathcal{A}_{\text{LLM}}} F_{p_L}(t \mid \text{data}, \mathcal{H}) \\ \overline{F}_{p_L}(t \mid \text{data}) &= \sup_{\mathcal{H} \in \mathcal{A}_{\text{LLM}}} F_{p_L}(t \mid \text{data}, \mathcal{H}) \end{aligned}$$

1595 *Proof of Theorem 3.* We derive the posterior distribution of the LLM reliability $p_L = \sum_{i=1}^m W_i p_i$
1596 by characterizing its CDF, building on the domain-level results from Theorem 2.

1597 **Step 1: CDF definition.** For any $t \in [0, 1]$, the posterior CDF of p_L is

$$F_{p_L}(t \mid \text{data}, \mathcal{H}) = Pr(p_L \leq t \mid \text{data}, \mathcal{H}) = Pr\left(\sum_{i=1}^m W_i p_i \leq t \mid \text{data}, \mathcal{H}\right).$$

1598 Define the region

$$\mathcal{R}_L(t) := \left\{ (p_1, \dots, p_m) \in (0, 1)^m : \sum_{i=1}^m W_i p_i \leq t \right\}.$$

1599 Then

$$F_{p_L}(t \mid \text{data}, \mathcal{H}) = \int_{\mathcal{R}_L(t)} \prod_{i=1}^m Pr(p_i \mid C_i, h_i) dp_1 \cdots dp_m,$$

1600 where cross-domain independence ensures the joint density factorizes as

$$Pr(p_1, \dots, p_m \mid \text{data}, \mathcal{H}) = \prod_{i=1}^m Pr(p_i \mid C_i, h_i).$$

1601 **Step 2: Domain-level mixture representation.** From Theorem 2, each domain posterior
1602 can be written as

$$Pr(p_i \mid C_i, h_i) = \int_0^1 \int_0^\infty f_{p_i}(p_i \mid \mu_i, \nu_i, C_i) Pr(\mu_i, \nu_i \mid C_i, h_i) d\mu_i d\nu_i,$$

1603 where $f_{p_i}(\cdot \mid \mu_i, \nu_i, C_i)$ is the conditional density of $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$ given that

$$\theta_{ij} \mid \mu_i, \nu_i, C_i \stackrel{\text{ind.}}{\sim} \text{Beta}(C_{ij} + \mu_i \nu_i, N_{ij} - C_{ij} + (1 - \mu_i) \nu_i), \quad j = 1, \dots, n_i.$$

1604 **Step 3: Hierarchical decomposition across domains.** Substituting the domain-level
1605 mixture into the LLM-level CDF:

$$\begin{aligned} F_{p_L}(t \mid \text{data}, \mathcal{H}) &= \\ & \int_{\mathcal{R}_L(t)} \prod_{i=1}^m \left[\int_0^1 \int_0^\infty f_{p_i}(p_i \mid \mu_i, \nu_i, C_i) Pr(\mu_i, \nu_i \mid C_i, h_i) d\mu_i d\nu_i \right] dp_1 \cdots dp_m \end{aligned}$$

1606 **Step 4: Change of integration order.** By Fubini’s theorem (applicable because all
 1607 integrands are non-negative and integrate to finite values), we can exchange the order of
 1608 integration:

$$F_{p_L}(t \mid \text{data}, \mathcal{H}) = \int_0^1 \cdots \int_0^1 \int_0^\infty \cdots \int_0^\infty \left[\int_{\mathcal{R}_L(t)} \prod_{i=1}^m f_{p_i}(p_i \mid \mu_i, \nu_i, C_i) dp_1 \cdots dp_m \right] \\ \times \prod_{i=1}^m Pr(\mu_i, \nu_i \mid C_i, h_i) \prod_{i=1}^m d\mu_i d\nu_i.$$

1609 **Step 5: Conditional CDF.** The inner integral is the conditional CDF of p_L given all
 1610 hyperparameters:

$$G(t \mid \{\mu_i, \nu_i\}_{i=1}^m, \text{data}) := \int_{\mathcal{R}_L(t)} \prod_{i=1}^m f_{p_i}(p_i \mid \mu_i, \nu_i, C_i) dp_1 \cdots dp_m \\ = Pr \left(\sum_{i=1}^m W_i p_i \leq t \mid \{\mu_i, \nu_i\}_{i=1}^m, \text{data} \right),$$

1611 where $p_i \mid \mu_i, \nu_i, C_i$ are independent across domains, each distributed according to the
 1612 weighted sum of independent Betas as specified in Step 2.

1613 **Step 6: Mixture representation.** Combining Steps 4 and 5:

$$F_{p_L}(t \mid \text{data}, \mathcal{H}) = \int \cdots \int G(t \mid \{\mu_i, \nu_i\}_{i=1}^m, \text{data}) \prod_{i=1}^m Pr(\mu_i, \nu_i \mid C_i, h_i) \prod_{i=1}^m d\mu_i d\nu_i \quad (\text{A.10})$$

1614 where the integrals range over $\mu_i \in (0, 1)$ and $\nu_i \in (0, \infty)$ for all $i = 1, \dots, m$.

1615 This expresses the LLM posterior CDF as a weighted average of conditional CDFs, where
 1616 the weights are given by the product of independent domain-level hyper-posteriors.

1617 **Step 7: Imprecise probability bounds.** For each $\mathcal{H} \in \mathcal{A}_{\text{LLM}}$, the formula above defines a
 1618 valid CDF $F_{p_L}(\cdot \mid \text{data}, \mathcal{H})$. The imprecise posterior is characterized by pointwise envelopes:

$$\underline{F}_{p_L}(t \mid \text{data}) = \inf_{\mathcal{H} \in \mathcal{A}_{\text{LLM}}} F_{p_L}(t \mid \text{data}, \mathcal{H}), \\ \overline{F}_{p_L}(t \mid \text{data}) = \sup_{\mathcal{H} \in \mathcal{A}_{\text{LLM}}} F_{p_L}(t \mid \text{data}, \mathcal{H})$$

1619 where, $t \in [0, 1]$ □

1620 Appendix A.5. Reliability for Multiple Consecutive Operations

1621 Theorems 1, 2, and 3 establish posterior distributions for the reliability parameters θ_{ij} , p_i ,
 1622 and p_L at the subdomain, domain, and LLM levels, respectively. These parameters represent
 1623 the probability of success on a *single* future operation. However, in practical reliability
 1624 assessment, we often need to evaluate performance over *multiple consecutive operations*.

1625 This appendix clarifies the mathematical framework for extending the case $n^F = 1$ to
 1626 $n^F > 1$ (reliability with the required number of future failure-free tasks), which is formalized
 1627 in Theorems 4, 5, and 6. We explain:

- 1628 1. Why reliability over n^F operations takes the form $\theta_{ij}^{n^F}$ (or $p_i^{n^F}, p_L^{n^F}$)
 1629 2. Why we characterize these distributions through CDFs rather than closed-form densi-
 1630 ties
 1631 3. How to derive the CDF integral formulas
 1632 4. How this framework extends across all hierarchical levels

1633 *Appendix A.5.1. Definition: Reliability with the required n^F failure-free future tasks*

1634 At the subdomain level, we define reliability for n^F consecutive future operations as:

$$R_{ij}(n^F) = \theta_{ij}^{n^F}$$

1635 Assume that:

- 1636 • Each task in subdomain S_{ij} succeeds independently with probability θ_{ij}
- 1637 • Operations are identically distributed (i.i.d. assumption)
- 1638 • We observe n^F consecutive operations

1639 Under these assumptions, the probability that *all* n^F operations succeed is:

$$\Pr(\text{all } n^F \text{ operations succeed}) = \underbrace{\theta_{ij} \times \theta_{ij} \times \cdots \times \theta_{ij}}_{n^F \text{ times}} = \theta_{ij}^{n^F}$$

1640 Because θ_{ij} is uncertain—it is a random variable with posterior distribution $Pr(\theta_{ij} | C_i, h_i)$
 1641 from Theorem 1—the n^F -operation reliability $R_{ij}(n^F) = \theta_{ij}^{n^F}$ is also a random variable. We
 1642 must therefore characterize its full posterior distribution.

1643 *Appendix A.5.2. Why the CDF Approach*

1644 Unlike the posteriors in Theorems 1–3, which characterize reliability for a *single* future
 1645 operation, the distribution of $R_{ij}(n^F) = \theta_{ij}^{n^F}$ generally does not have a closed-form probability
 1646 density function, even when θ_{ij} follows a Beta distribution. This is because:

- 1647 • The transformation $\theta_{ij} \mapsto \theta_{ij}^{n^F}$ is nonlinear for $n^F \neq 1$
- 1648 • At higher hierarchical levels (domain, LLM), we have weighted sums of dependent
 1649 random variables raised to powers, making closed forms even less tractable

1650 Therefore, we characterize the distribution of $R_{ij}(n^F)$ through its *CDF*, which has a
 1651 tractable integral representation and is sufficient for all practical reliability calculations (e.g.,
 1652 computing probabilities, quantiles, expected values).

1653 *Appendix A.5.3. Deriving the CDF Formula*

1654 What does CDF answers: “What is the probability that n^F -operation reliability is at
1655 most r ?”

$$F_{R_{ij}(n^F)}(t | C_i, h_i) = \Pr(R_{ij}(n^F) \leq t | C_i, h_i) = \Pr(\theta_{ij}^{n^F} \leq t | C_i, h_i)$$

1656 To compute this probability, we use the following key observation. If $\theta_{ij}^{n^F} \leq t$, then taking
1657 the n^F -th root of both sides:

$$\theta_{ij} \leq t^{1/n^F}$$

1658 Since $\theta_{ij} \in (0, 1)$ and $t \in [0, 1]$, both sides are positive, and the function $x \mapsto x^{1/n^F}$ is
1659 monotonically increasing on $[0, 1]$ for $n^F > 0$. Therefore, the inequality is preserved under
1660 this transformation.

1661 It follows that:

$$\Pr(\theta_{ij}^{n^F} \leq t | C_i, h_i) = \Pr(\theta_{ij} \leq t^{1/n^F} | C_i, h_i)$$

1662 The right-hand side is simply the CDF of the posterior distribution $Pr(\theta_{ij} | C_i, h_i)$ from
1663 Theorem 1, evaluated at t^{1/n^F} :

$$\Pr(\theta_{ij} \leq t^{1/n^F} | C_i, h_i) = \int_0^{t^{1/n^F}} Pr(\theta_{ij} | C_i, h_i) d\theta_{ij}$$

1664 This gives us the formula in Theorem 4.

1665 *Appendix A.5.4. Extension to Domain and LLM Levels*

1666 The same logic extends to domain and LLM levels:

- 1667 • **Domain level:** $R_i(n^F) = p_i^{n^F}$ where $p_i = \sum_{j=1}^{n_i} \Omega_{ij} \theta_{ij}$
- 1668 • **LLM level:** $R_L(n^F) = p_L^{n^F}$ where $p_L = \sum_{i=1}^m W_i p_i$

1669 At these levels, closed forms are even less available because we must:

- 1670 1. Compute the distribution of weighted sums of dependent Beta random variables (from
1671 Theorems 2 and 3)
- 1672 2. Apply the power transformation to obtain $R_i(n^F)$ or $R_L(n^F)$

1673 Both steps lack analytical solutions, so we compute CDFs via numerical integration or
1674 Monte Carlo sampling.

1675 **Appendix B. Global numerical settings used in all figures**

1676 All figures in this paper are produced under the same imprecise hierarchical Bayes setup
1677 and numerical settings:

- 1678 • **Domains/subdomains.** D1 (Coding) = {MBPP, DS-1000}, D2 (Reasoning) =
1679 {BoolQ, RACE-H}.

- 1680 • **Aggregation weights.** Within-domain weights $\Omega_1 = [0.204, 0.796]$ (MBPP, DS-1000)
 1681 and $\Omega_2 = [0.483, 0.517]$ (BoolQ, RACE-H); LLM-level weights $W = [0.149, 0.851]$ (D1,
 1682 D2).

- 1683 • **Hyperpriors and imprecision ranges (per domain i).**
 1684 $\mu_i \sim \text{Beta}(a_i, b_i)$, $\nu_i \sim \text{Gamma}(c_i, \text{rate} = d_i)$, with $a_i, b_i \in [1, 12]$ and $c_i, d_i \in [1, 25]$.

- 1685 • **Simulation** For each domain and each sampled hyperparameter configuration: (1)
 1686 sample (μ_i, ν_i) from the discretized posterior; (2) draw subdomain accuracies θ_{ij} from
 1687 their Beta posteriors; (3) aggregate to p_i and p_L ; (4) compute either empirical CDFs
 1688 of $Z \in \{\theta_{ij}, p_i, p_L\}$ on a fixed grid, or Monte Carlo expectations such as $\widehat{\mathbb{E}}[R_L(n^F)] =$
 1689 $\frac{1}{S} \sum_{s=1}^S (p_L^{(s)})^{n^F}$. Repeating over K hyperparameter configurations yields a family of
 1690 CDFs/expectations; the plotted bands are the *pointwise min-max* across configurations
 1691 (epistemic uncertainty).

- 1692 • **CDF envelopes.** For any quantity $Z \in \{\theta_{ij}, p_i, p_L\}$, we compute an empirical CDF for
 1693 each configuration, and plot the pointwise min-max envelope $[F_{\min}(t), F_{\max}(t)]$ across
 1694 configurations.

Appendix C. List of Mathematical Notations

Symbol	Description
\mathcal{X}	Input-space of all possible tasks for an LLM
π	Operational Profile (OP), probability distribution over tasks
n	Sequence of tasks
$I(x_\tau)$	Indicator function: 1 for success, 0 for failure on task τ
x_τ	The τ -th task in a sequence of n tasks
$R(n, \pi)$	LLM reliability: probability of failure-free operation over n tasks
M	Number of LLM models being evaluated
m	Number of independent domains
D_i	Domain i ($i = 1, \dots, m$)
n_i	Number of subdomains in domain i
S_{ij}	Subdomain j in domain i ($j = 1, \dots, n_i$)
C_{ij}	Number of correct responses in subdomain S_{ij}
N_{ij}	Number of trials (tasks) in subdomain S_{ij}
C_i	Set of observed data in domain i : $\{(C_{ik}, N_{ik})\}_{k=1}^{n_i}$
data	Observed data across all domains: $\{C_1, \dots, C_m\}$
θ_{ij}	Subdomain reliability (success probability) for S_{ij}
$\boldsymbol{\theta}_i$	Vector of subdomain reliabilities: $(\theta_{i1}, \dots, \theta_{in_i})$
p_i	Domain-level non-failure probability (posterior)
p_L	LLM-level non-failure probability (posterior)
$R_{ij}(n^F)$	Subdomain reliability for n^F consecutive tasks
$R_i(n^F)$	Domain reliability for n^F consecutive tasks
$R_L(n^F)$	LLM reliability for n^F consecutive tasks
Ω_{ij}	Operational weight for subdomain j in domain i
W_i	Operational weight for domain i
μ_i	Expected reliability (prior mean) for domain i
ν_i	Prior strength (concentration parameter) for domain i
a_i, b_i	Hyperparameters of the Beta prior for μ_i
c_i, d_i	Hyperparameters of Gamma prior for ν_i
h_i	Domain-level hyperparameter tuple (a_i, b_i, c_i, d_i)
\mathcal{A}_i	Admissible set of domain-level hyperparameters for domain i
\mathcal{H}	Tuple of all domain hyperparameters: (h_1, \dots, h_m)
\mathcal{A}_{LLM}	LLM-level admissible hyperparameter set: $\mathcal{A}_1 \times \dots \times \mathcal{A}_m$
$F_{p_i}(t C_i, h_i)$	Posterior CDF of domain-level non-failure probability p_i
$[\underline{F}_{p_i}(t C_i), \overline{F}_{p_i}(t C_i)]$	Interval-valued CDF of domain-level non-failure probability p_i
$F_{p_L}(t \text{data}, \mathcal{H})$	Posterior CDF of LLM-level non-failure probability p_L
$[\underline{F}_{p_L}(t \text{data}), \overline{F}_{p_L}(t \text{data})]$	Interval-valued CDF of domain-level non-failure probability p_L
$F_{R_{ij}(n^F)}(t C_i, h_i)$	Posterior CDF of subdomain reliability for n^F consecutive tasks
$F_{R_i(n^F)}(t C_i, h_i)$	Posterior CDF of domain-level reliability for n^F consecutive tasks
$F_{R_L(n^F)}(t \text{data}, \mathcal{H})$	Posterior CDF of LLM-level reliability for n^F consecutive tasks

Symbol	Description
$f_{\text{marg}}(\theta_{ij}, C_i; h_i)$	Unnormalized marginal posterior for θ_{ij}
$Z_{\text{marg}}(h_i)$	Normalizing constant for marginal posterior
$f_{p_i}(\cdot \mu_i, \nu_i, C_i)$	Conditional density of domain reliability
$G(t \{\mu_i, \nu_i\}_{i=1}^m, \text{data})$	conditional CDF of p_L given all domain hyperparameters
$\mathcal{R}_i(t)$	Constraint region $\{(p_1, \dots, p_m) : \sum_j \Omega_{ij} p_{ij} \leq t\}$ for domain i
$\mathcal{R}_L(t)$	Constraint region induced by threshold t for LLM-level CDF evaluation
$f_{\text{marg}}(\theta_{ij}, C_i; h_i)$	Unnormalized marginal posterior for θ_{ij}
$Z_{\text{marg}}(h_i)$	Normalizing constant for marginal posterior
$f_{p_i}(\cdot \mu_i, \nu_i, C_i)$	Conditional density of domain-level non-failure probability p_i
$G(t \{\mu_i, \nu_i\}_{i=1}^m, \text{data})$	Conditional CDF of LLM-level non-failure probability p_L given all domain hyperparameters
$\mathcal{R}_i(t)$	Constraint region for domain-level CDF evaluation of p_i
$\mathcal{R}_L(t)$	Constraint region for LLM-level CDF evaluation of p_L
t	Probability threshold at which the corresponding CDF is evaluated