

BHASHABENCH V1: A COMPREHENSIVE BENCHMARK FOR THE QUADRANT OF INDIC DOMAINS

Anonymous authors

Paper under double-blind review

ABSTRACT

The rapid advancement of large language models (LLMs) has intensified the need for domain and culture specific evaluation. Existing benchmarks are largely Anglocentric and domain-agnostic, limiting their applicability to India-centric contexts. To address this gap, we introduce **BhashaBench V1**, the first domain-specific, multi-task, bilingual benchmark focusing on critical Indic knowledge systems. BhashaBench V1 contains **74,166** meticulously curated question-answer pairs, with 52,494 in English and 21,672 in Hindi, sourced from authentic government and domain-specific exams. It spans four major domains: Agriculture, Legal, Finance, and Ayurveda, comprising 90+ subdomains and covering 500+ topics, enabling fine-grained evaluation. Evaluation of 29+ LLMs reveals significant domain and language specific performance gaps, with especially large disparities in low-resource domains. For instance, GPT-4o achieves 76.49% overall accuracy in Legal but only 59.74% in Ayurveda. Models consistently perform better on English content compared to Hindi across all domains. Subdomain-level analysis shows that areas such as *Cyber Law*, *International Finance* perform relatively well, while *Panchakarma*, *Seed Science*, and *Human Rights* remain notably weak. **BhashaBench V1** provides a comprehensive dataset for evaluating large language models across India’s diverse knowledge domains. It enables assessment of models’ ability to integrate domain-specific knowledge with bilingual understanding. All code, benchmarks, and resources are publicly available to support open research.

1 INTRODUCTION

The rapid advancement of large language models (LLMs) has transformed artificial intelligence, extending their capabilities far beyond traditional natural language processing. Models such as GPT-4o (OpenAI, 2024), GPT-OSS-120B (OpenAI et al., 2025), DeepSeek-V3 (DeepSeek-AI et al., 2025), and Qwen-3 (Yang et al., 2025) excel across diverse domains, from code generation and mathematical reasoning to creative writing and scientific analysis (Brown et al., 2020; Touvron et al., 2023; OpenAI et al., 2024), enabling applications in conversational AI, education, healthcare, finance, legal services, and agriculture (Bubeck et al., 2023; Wei et al., 2022). Platforms like *Krishi Sathi* (Vijayvargia et al., 2025) leverage LLMs for crop advisory and pest detection, improving agricultural productivity. Despite these advances, substantial performance gaps remain in multilingual and domain-specific contexts, particularly for non-Latin, low-resource languages (Wang et al., 2024a; Zhong et al., 2025; Ahuja et al., 2024). English-centric training limits models’ ability to capture nuanced knowledge in specialized fields and India-specific domains, such as Ayurveda, indigenous agriculture, finance, and regional legal systems (Winata et al., 2021; Sen & et al., 2023; Khanuja et al., 2021a), highlighting the need for culturally and contextually aware evaluation. The scale of this problem demands urgent attention, as India’s diverse knowledge ecosystem affects millions of lives across multiple critical domains. In agriculture alone, over 40 million farmers rely on farming-related activities (IndiaDataMap, 2025), and access to accurate information on crop management, soil health, and sustainable practices can have a direct impact on food security and livelihoods. The complexity is further magnified by the fact that each state in India has its own distinct agricultural methods, crop varieties, soil conditions, and traditional farming practices that have evolved over centuries to suit local climatic and geographical conditions. Similarly, India’s legal system processes millions of cases annually, requiring precise understanding of complex legal frameworks,

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

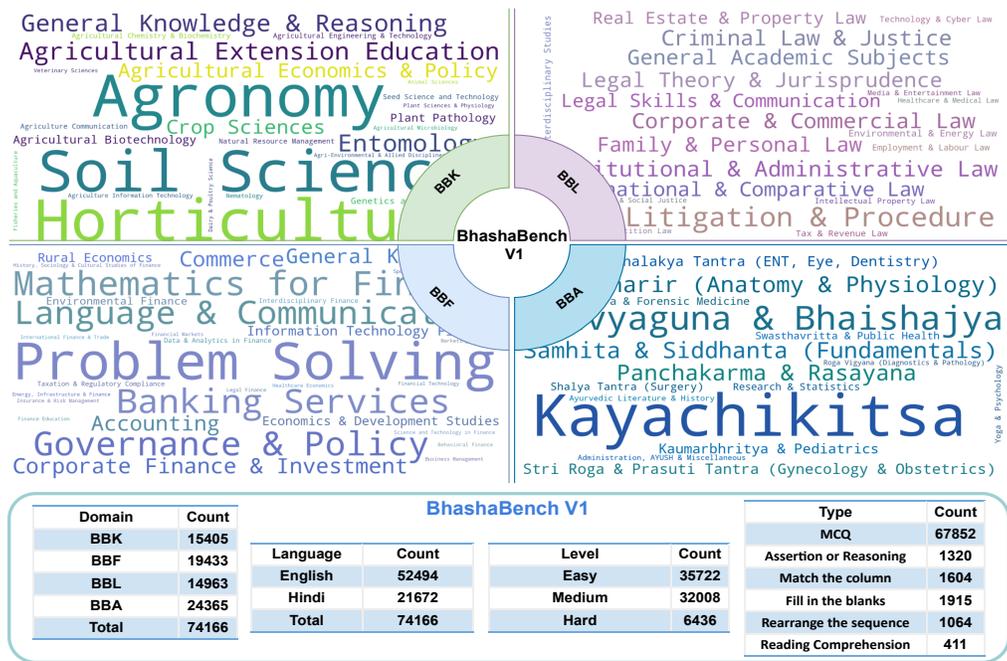


Figure 1: Overview diagram and statistics of BhashaBench V1.

precedents, and procedural nuances that vary across states and jurisdictions (National Judicial Data Grid (NJDG), India, 2023). The healthcare sector, particularly traditional medicine systems like Ayurveda, serves millions of patients who rely on practitioners’ knowledge of ancient texts, formulations, and treatment protocols. Furthermore, India’s financial ecosystem processes billions of transactions daily, including over 100 billion UPI transactions annually, where even minor misunderstandings in financial regulations or procedures can have cascading effects (National Payments Corporation of India, 2023). While existing benchmarks such as MMLU (Hendrycks et al., 2021a), HellaSwag (Zellers et al., 2019), AGIEVAL (Zhong et al., 2023), and more recent multilingual efforts like MEGA (Ahuja et al., 2024) attempt to assess model capabilities, they often focus primarily on English content and may not fully capture India-specific nuances, cultural contexts, and domain expertise that are essential for real-world applications in the Indian subcontinent.

To address these critical gaps, we introduce **BhashaBench V1**, the first comprehensive domain-specific, multi-task, bilingual benchmark designed explicitly for evaluating large language models on India-centric knowledge systems. BhashaBench V1 encompasses four fundamental domains that form the backbone of Indian society and economy: Agriculture (BBK - BhashaBench Krishi), Legal (BBL - BhashaBench Legal), Finance (BBF - BhashaBench Finance), and Ayurveda (BBA - BhashaBench Ayurveda). The benchmark spans over 90 subdomains and covers more than 500 specific topics, reflecting the intricate complexity and diversity of Indian knowledge systems. This granular categorization enables fine-grained evaluation of model performance across specialized areas that require deep domain expertise and cultural understanding. The dataset has been meticulously curated from over 40 authentic government and professional examination papers, ensuring that the questions reflect real-world scenarios and ground-level challenges faced by practitioners in these domains (Indian Knowledge Systems Division, Ministry of Education, Government of India, 2025; Zhong & Goodfellow, 2024). To maximize coverage across India’s linguistic landscape, BhashaBench V1 currently supports English and Hindi, the two most widely understood languages in the country, collectively enabling assessment of models’ capabilities for a significant portion of India’s population while maintaining the cultural and contextual authenticity of the original knowledge systems.

Our comprehensive evaluation of 29+ state-of-the-art language models on BhashaBench V1 reveals significant performance disparities across domains and languages, highlighting the urgent need for India-specific model development and evaluation. The results demonstrate substantial domain-

specific performance gaps, with models showing varying degrees of competency across different knowledge areas. For instance, GPT-4o, one of the top-performing models, achieved 76.49% accuracy in the Legal domain but only 59.74% in Ayurveda, illustrating the challenges models face with traditional Indian knowledge systems. Similarly, consistent language-specific performance gaps emerged, with models generally performing better on English content compared to Hindi across all domains. The subdomain-level analysis further reveals granular insights into model capabilities, showing that certain areas such as Cyber Law and International Finance demonstrate relatively strong performance, while traditional domains like Panchakarma, Seed Science, and Human Rights remain notably challenging for current LLMs. These findings underscore the critical importance of domain and language-specific evaluation frameworks for assessing model readiness for real-world deployment in diverse Indian contexts.

2 RELATED WORK

2.1 EXPLORATION OF LLMs

The landscape of large language models has witnessed unprecedented growth, with both proprietary and open-source models achieving remarkable capabilities. Recent proprietary LLMs, including GPT-4o and GPT-4o-mini (OpenAI et al., 2024), Claude-3.5 Sonnet (Anthropic, 2024), and the Gemini series (Google, 2023), have demonstrated significant improvements across various benchmarks (Chiang et al., 2024; Wang et al., 2024c). The open-source ecosystem has flourished with models such as the Llama-3 series (Grattafiori et al., 2024), Gemma (Team et al., 2024), Qwen2.5 (Qwen et al., 2025), and Mistral (Jiang et al., 2023) achieving competitive performance while maintaining transparency and accessibility.

While primarily trained on English-dominant corpora, many models incorporate substantial multilingual data during pretraining (Team et al., 2024; Grattafiori et al., 2024; Üstün et al., 2024), enabling capabilities in hundreds of languages with varying proficiency (Nguyen et al., 2024). Language-specific models have gained momentum, particularly for underrepresented languages including Indic languages (Gala et al., 2024; 2023). Notable examples include Airavata (Gala et al., 2024), MuRIL (Khanuja et al., 2021b), and recent generative models like Param-1 (Pundalik et al., 2025).

Domain-specific language models have emerged as a critical research direction. Medical applications include Med-PaLM (Singhal et al., 2023) and BioBERT (Lee et al., 2019), while legal and financial domains have seen LegalBERT (Chalkidis et al., 2020) and FinBERT (Yang et al., 2020) respectively. In the Indian context, domain-specific initiatives like Agri-Param (BharatGenAI, 2025a), Ayur-Param (BharatGenAI, 2025b), Finance-Param (BharatGenAI, 2025c), and Legal-Param (BharatGenAI, 2025d) address unique requirements of India’s diverse knowledge systems through continual pretraining (Nag et al., 2024) or instruction fine-tuning (Aralimatti et al., 2025).

Despite these advances, comprehensive evaluation frameworks for culturally and linguistically diverse domains remain limited, particularly for traditional knowledge systems requiring nuanced understanding of local contexts. This work conducts a comprehensive evaluation of 29+ state-of-the-art models on BhashaBench V1 to address these evaluation challenges.

2.2 EVALUATION OF LLMs

Numerous benchmarks have been developed to assess large language model performance. General-purpose benchmarks such as ARC-C (Clark et al., 2018), MMLU (Hendrycks et al., 2021b), MMLU-Pro (Wang et al., 2024b), AGIEval (Zhong et al., 2023), BIG-Bench (Srivastava et al., 2023), and HellaSwag (Zellers et al., 2019) evaluate LLMs across diverse tasks from commonsense reasoning to knowledge-intensive question answering. However, these remain largely Anglocentric with limited multilingual evaluation (Bandarkar et al., 2024; Kakwani et al., 2020). Indic-focused resources such as MILU (Verma et al., 2025) expand linguistic diversity by benchmarking 11 Indian languages, though they largely center on general-domain multiple-choice tasks. Similarly, Sanskriti (Maji et al., 2025) offers culturally grounded evaluation rooted in Indian history and heritage but lacks broad domain and multilingual coverage.

Table 1: Overview of how BhashaBench V1 compares to Indic, multilingual, and other domain-specific evaluation benchmarks.

Benchmark	Languages	Domains	Task Formats	Cultural	Size
MMLU	En	57 general	MCQ	✗	15.9K
MMLU-Indic	11 Indic	57 general	MCQ	✗	15.9K
MILU	11 Indic	8 general	MCQ	~	79.6K
Sanskriti	En	Culture, History	MCQ, QA	✓	21.8K
IndicQA	11 Indic	General KG	QA	~	50K per language
AgXQA 1.1	En	Agriculture	QA	✗	1.5K
MultiFin	En	Finance	Classification	✗	10K
IL-TUR	10 Indic	Legal	QA, Generation, Classification	~	Multi-size (per task)
BhashaBench V1	En + Hi	Agri, Legal Finance, Ayurveda	MCQ, A/R, FIB MTC, RC, RTS	✓	74.2K

Abbreviations: MCQ = Multiple Choice Questions; FIB = Fill in the Blanks; A/R = Assertion/Reason; RC = Reading Comprehension; QA = Question Answering; MTC = Match the Columns; KG = Knowledge Graph.

Cultural Authenticity Legend: ✓ = sourced from Indic region-specific exams or created by native domain experts; ~ = partially translated or culturally adapted; ✗ = mainly sourced from non-Indic sources.

To address domain-specific challenges, specialized benchmarks have emerged. In agriculture, benchmarks like AgriBench (Zhou & Ryo, 2024), BVL QA Corpus (AnhaltAI, 2024), AgXQA (Kpodo et al., 2024), AgEval (Arshad et al., 2025), and SeedBench (Ying et al., 2025) cover crop disease identification to advisory support. The finance domain features FinGAIA (Zeng et al., 2025), FinanceBench (Islam et al., 2023), MultiFin (Jørgensen et al., 2023), InvestorBench (Li et al., 2024), and MultiFinBen (Peng et al., 2025) for financial reasoning, fraud detection, and trading evaluation. Legal domain efforts include IL-TUR (Joshi et al., 2024), IndicLegalQA (Nigam et al., 2023), LegalBench (Guha et al., 2023), LEXTREME (Niklaus et al., 2023), and the CAIL series (Xiao et al., 2018; 2019) for legal question answering, case summarization, and judgment prediction. Traditional medicine resources such as MTCMB (Kong et al., 2025), Pratyaya-Kosh (Ragad & Gokhale, 2019), Anveshana (Terdalkar et al., 2023), and OpenTCM (He et al., 2025) provide task-specific evaluation datasets covering knowledge graphs, OCR correction, and dosha analysis.

Despite this progress, key limitations persist. Many benchmarks are restricted to English or high-resource languages, limiting effectiveness for multilingual and Indic contexts. Others focus on narrow tasks, unable to capture full domain expertise. Evaluation methodologies vary widely from accuracy scores to human judgments, hindering standardized comparison across domains and languages. These gaps underscore the need for a unified, multilingual, and domain-aware evaluations.

3 BHASHABENCH V1

3.1 DESIGN PRINCIPLES

The primary motivation behind BhashaBench V1 is to comprehensively assess domain-specific knowledge and reasoning capabilities of large language models within India’s diverse and culturally rich knowledge ecosystems. Unlike existing benchmarks focusing on general or Western-centric domains, our benchmark evaluates specialized Indian knowledge systems requiring deep cultural understanding and contextual awareness. Table 1 illustrates how BhashaBench V1 addresses critical gaps in existing evaluation frameworks. While benchmarks like MMLU and Indic MMLU provide broad coverage, they lack cultural grounding and rely on translated content. Domain-specific benchmarks typically focus on narrow tasks within single domains. BhashaBench V1 uniquely combines cultural authenticity, domain specialization, and bilingual support through exam-sourced questions.

BhashaBench V1 adheres to seven core design principles: **(1) Critical Indian Domains:** Encompasses Agriculture, Legal systems, Finance, and Ayurveda with fine-grained subfields. **(2) Diverse Task Formats:** Includes multiple-choice, assertion-reasoning, fill-in-blanks, and comprehension tasks. **(3) India-Specific Reasoning:** Evaluates domain-specific reasoning incorporating cultural contexts and regional practices. **(4) Bilingual Framework:** Supports English and Hindi evaluation maintaining cultural authenticity. **(5) Authentic Sources:** Questions curated from government examinations and professional certifications. **(6) Difficulty Assessment:** Categorized into Easy,

216 Medium, Hard levels. **(7) Cultural Authenticity:** Prioritizes traditional knowledge systems includ-
 217 ing Ayurvedic principles. This framework spans 90+ subdomains covering 500+ topics, enabling
 218 comprehensive evaluation of model capabilities in India-centric contexts.¹
 219

220 3.2 DATA COLLECTION

221
 222 The data collection process for BhashaBench V1 follows a systematic approach similar to AGIEVAL
 223 (Zhong et al., 2023), focusing on authentic examination materials from national and state-level as-
 224 sessments. We systematically gathered publicly available question papers from official online exam-
 225 ination portals, which host previously released papers that are manually curated by subject matter
 226 experts, ensuring accurate topic tagging, language annotation, and validated answer keys.

227 Our comprehensive collection encompasses over 40 different examination types across multiple
 228 categories: national competitive exams, domain-specific degree examinations, professional certifi-
 229 cation tests, and state-level civil services examinations. Regional state examinations proved par-
 230 ticularly valuable as they incorporate state-specific topics, local knowledge systems, and cultural
 231 practices often overlooked in national assessments. These examinations are typically taken by in-
 232 dividuals seeking higher education opportunities or career advancement, ensuring questions reflect
 233 practical, real-world knowledge requirements.

234 The final dataset comprises **74,166** carefully curated question-answer pairs spanning four core do-
 235 mains, with **52,494 questions in English** (70.8%) and **21,672 questions in Hindi** (29.2%), reflect-
 236 ing practical usage patterns in Indian educational and professional contexts. This approach ensures
 237 BhashaBench V1 captures the nuanced intersection between language, culture, and domain expertise
 238 essential for effective model deployment in Indian contexts.

239 3.3 DATA PROCESSING

240
 241 Our data processing phase focused on extracting structured question-answer pairs from PDF exam-
 242 ination papers while preserving linguistic and formatting nuances essential for authentic evaluation.
 243 Most examination materials were available exclusively in PDF format, requiring sophisticated OCR
 244 processing pipelines to handle multilingual content and domain-specific terminology.

245 **OCR Pipeline Selection:** Based on existing evaluations (Paruchuri & Team, 2024), Surya OCR
 246 demonstrated superior performance in handling Indic languages and domain-specific content. Re-
 247 ported results show 98.1% normalized text similarity for English and 98.9% for Hindi, with an aver-
 248 age of 97.8%, outperforming alternatives such as Tesseract (88.0% overall) and Google Vision API
 249 (96.7%). Surya’s architecture, designed for multilingual document understanding with enhanced
 250 Indic script support, makes it a suitable choice for diverse examination materials.

251 **Question-Answer Extraction Pipeline:** Following OCR processing, we developed an extraction
 252 pipeline leveraging GPT-OSS-120B (OpenAI, 2024) to structure raw text into formatted question-
 253 answer pairs. Key challenges included format variations across examination bodies, answer key
 254 alignment, multi-format questions (MCQ, assertion-reasoning, comprehension), and language-
 255 specific formatting conventions. The pipeline included: (1) **Question Extraction** using GPT-OSS-
 256 120B for boundary detection across different layouts; (2) **Option Parsing** to maintain original la-
 257 beling conventions; (3) **Answer Key Alignment** processing both inline and separate answer docu-
 258 ments; and (4) **Format Standardization** into consistent JSON structure with domain metadata.

259 **Data Cleaning and Quality Control:** Our multi-layered cleaning approach addressed noise and
 260 inconsistencies through systematic filtering. We excluded image-based questions, and questions
 261 with more than four options. Language verification used INDICLID (Madhani et al., 2023) and
 262 Unicode-based filtering (Khan et al., 2024) for proper linguistic categorization. Approximately
 263 30% of questions lacked subdomain classification, addressed using GPT-OSS-120B with domain-
 264 specific taxonomies. We classified questions into six categories: MCQ, assertion-reasoning, fill-
 265 in-the-blanks, match-the-column, reading comprehension, and sequence rearrangement. Duplicate
 266 detection employed both exact-match and semantic similarity measures.

267 **Manual Validation:** Following a methodology similar to (Bandarkar et al., 2024), all extracted
 268 question-answer pairs underwent rigorous expert validation to ensure accuracy verification, cultural
 269

¹More collection and processing procedures can be found in Appendix C.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

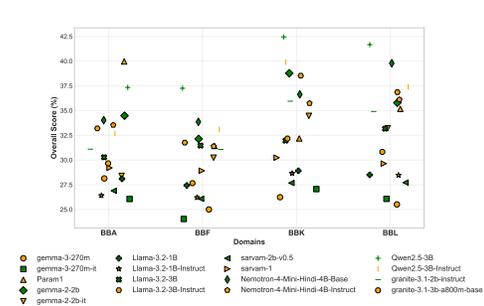


Figure 2: Comparative performance of small models ($\leq 4B$) over BhashaBench V1.

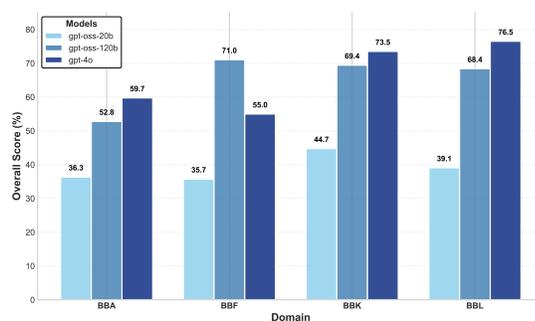


Figure 3: Comparative performance analysis of the GPT model family on BhashaBench V1.

context preservation, ambiguity resolution, and consistency standardization. Additionally, domain experts reviewed the linguistic authenticity to maintain the natural flow and idiomatic expressions characteristic of each language. This comprehensive multi-stage validation approach ensured that BhashaBench V1 maintains the highest data quality standards while preserving the authentic complexity and cultural specificity of the original examination materials.

3.4 DATA ANALYSIS

Figure 1 presents the comprehensive statistics of BhashaBench V1. Detailed exposition is provided in Appendix C.2. Of the total 74,166 questions, 70.8% are in English while 29.2% are in Hindi, reflecting practical bilingual usage patterns in Indian professional contexts. The dataset spans four specialized domains with varying complexity levels across 91 subdomains.

Agriculture (BBK): This domain encompasses agricultural sciences relevant to Indian farming systems across 25 subdomains. Agronomy dominates with 5,078 questions, reflecting its foundational role in agricultural education. The domain covers traditional practices alongside emerging areas like Agricultural Biotechnology and IT solutions. Its balanced difficulty distribution (44% easy, 45% medium, 11% hard) ensures comprehensive skill assessment.

Finance (BBF): Covers India’s complex financial ecosystem through 30 subdomains. Problem Solving leads with 5,686 questions, followed by Mathematics for Finance (4,845), emphasizing the quantitative nature of financial practice. The domain uniquely incorporates India-specific areas like Rural Economics and Environmental Finance while addressing modern fintech developments.

Ayurveda (BBA): Represents traditional Indian medicine across 16 subdomains. Kayachikitsa (General Medicine) forms the core with 3,134 questions, while Dravyaguna covers pharmacology and therapeutics (2,972). This domain shows the highest proportion of accessible questions (53% easy), reflecting its foundational knowledge structure.

Legal (BBL): Encompasses Indian jurisprudence through 20 subdomains. Civil Litigation & Procedure dominates with 7,126 questions, followed by Constitutional Law (3,609). The domain balances traditional legal areas with contemporary developments like Technology & Cyber Law, maintaining strong cultural relevance through Family & Personal Law.

The predominantly MCQ format (>90%) ensures consistent evaluation methodology while supporting diverse cognitive assessment approaches across India-specific knowledge systems.

4 EXPERIMENTAL SETUP

We evaluate multiple state-of-the-art models on BhashaBench V1, including large proprietary models, open-source multilingual models, and domain-specific fine-tuned variants. Both base versions and instruction fine-tuned models are assessed to measure the effectiveness of specialized training approaches across India-specific knowledge domains. All evaluations are conducted in a zero-shot setting to assess the models’ inherent capabilities without task-specific examples. For open-source models, we utilize the LM-EVALUATION-HARNESS library (Biderman et al., 2024a;b) to ensure

clean, reproducible, and standardized evaluations. We employ the log-likelihood method where the probability of a given output string is computed by conditioning it on the provided input (Brown et al., 2020). For multiple choice questions with k possible answer choices, we select the answer string (a_i) with the highest conditional log probability: $\arg \max(\log P(a_1|x), \dots, \log P(a_k|x))$.

For closed-source and large-scale proprietary models, we utilize their respective APIs for evaluation due to computational constraints and access limitations. These API-based models are evaluated using a generative approach and are prompted to generate responses in a structured JSON format to facilitate automated response parsing. This comprehensive experimental framework enables systematic comparison across diverse model architectures while maintaining evaluation consistency across both open-source and proprietary systems. Additional details regarding model specifications, hyperparameters, and computational resources are provided in Appendix D.

5 RESULTS AND DISCUSSIONS

In this section, we discuss the results and our findings across all the experiments conducted.

5.1 ZERO-SHOT PERFORMANCE ACROSS ALL DOMAINS (EN + HI)

Table 2 shows the performance of various models in English and Hindi under the zero-shot setup. Among these, Qwen3-235B-A22B-Instruct emerges as the strongest model, consistently outperforming all competitors across both languages, with an average accuracy of 67.25%. This is followed by GPT-4o at 66.18% and gpt-oss-120b at 65.41%. Performance shows clear stratification across model sizes and types, with models exceeding 27B parameters demonstrating substantially higher accuracies compared to smaller variants. Among the 7B-27B range, gemma-2-27b leads with 53.11% average accuracy, followed by gemma-2-27b-it at 44.64%. In the mid-range category, gemma-2-9b shows impressive performance at 48.07%, with Pangea-7B achieving 41.54%.

Smaller models under 4B parameters show more modest performance, with Qwen2.5-3B achieving the highest accuracy in this category at 39.68%. Models specifically designed for Indian languages include Param-1 (34.69%) and the Nemotron-4-Mini-Hindi variants (36.08% and 34.20%). Performance is notably higher in English compared to Hindi across most models, reflecting the typical pattern observed in multilingual language models, with models showing varying degrees of cross-lingual transfer capabilities.

5.2 HOW DO MODELS PERFORM IN SUBDOMAINS

We evaluate representative models across BBA, BBF, BBK, and BBL to capture performance within subdomains (see Figures 4 and 5). Qwen3-235B-A22B-Instruct-2507 achieves the strongest results, excelling in Research & Statistics (91.43%), Agricultural Biotechnology (91.6%), and Intellectual Property Law (87.91%). GPT-4o demonstrates robust performance, frequently scoring above 70% with peaks of 92% in Information Technology and Healthcare & Medical Law. GPT-oss-120b shows competitive performance, closely matching gpt-4o in domains like Agricultural Biotechnology (89.69%). Mid-sized models including Gemma-2-27b and Gemma-2-9b generally show moderate performance in the 50–70% range, with the 27B variant consistently outperforming its smaller counterpart. Llama-3.1-8B demonstrates limited performance, typically scoring 30–50% across domains. The compact Param-1 model shows consistent baseline performance, often matching Llama-3.1-8B despite requiring significantly fewer resources. Notable patterns emerge: Finance and Legal domains show the highest performance ceiling, with top models regularly exceeding 80% in Business Management and Constitutional Law. Agricultural domains present moderate complexity, while Ayurveda proves most challenging, with even the best models rarely exceeding 80% in specialized areas like Panchakarma. Results highlight clear advantages for large models in knowledge-intensive tasks, while smaller models provide practical utility in resource-constrained scenarios for general applications.

5.3 PERFORMANCE ANALYSIS ACROSS QUESTION DIFFICULTY LEVELS

We evaluated model performance on Easy, Medium, and Hard questions across the four benchmark domains BBA, BBF, BBK, and BBL. In BBA, top-performing models such as GPT-4o and Qwen3-

Table 2: Zero-shot scores (%) of LLMs across domains on BhashaBench V1 (EN + HI). The benchmark covers Agriculture (BBK), Finance (BBF), Legal (BBL), and Ayurveda (BBA). ‘‘Avg’’ denotes the overall average across that domain. Top-scoring models are highlighted as follows: yellow for models < 4B parameters, green for models between 4B and 27B, and blue for models > 27B.

Model	BBA			BBF			BBK			BBL		
	Eng	Hin	Avg									
<i>< 4B Models</i>												
gemma-3-270m	28.08	28.25	28.14	24.98	25.06	25.00	26.64	24.45	26.24	25.49	25.54	25.51
gemma-3-270m-it	26.23	25.77	26.06	24.13	23.84	24.04	27.44	25.35	27.06	25.56	27.26	26.07
Param-1	41.12	38.04	39.97	32.24	29.56	31.42	33.10	27.97	32.18	36.15	32.89	35.17
gemma-2-2b	36.80	30.61	34.48	34.20	27.50	32.14	41.24	27.49	38.78	38.45	29.61	35.79
gemma-2-2b-it	29.38	26.79	28.40	31.26	27.93	30.24	35.94	27.71	34.47	34.49	30.25	33.22
Llama-3.2-1B	29.17	26.30	28.10	28.24	25.61	27.43	29.71	25.21	28.91	29.63	25.88	28.52
Llama-3.2-1B-Instruct	26.77	25.82	26.41	26.28	26.04	26.21	29.16	26.33	28.65	29.08	27.04	28.47
Llama-3.2-3B	31.62	28.05	30.28	33.04	27.92	31.46	32.68	28.69	31.96	35.17	28.53	33.17
Llama-3.2-3B-Instruct	35.31	29.67	33.20	32.94	29.09	31.76	40.59	29.09	38.53	39.74	30.13	36.86
sarvam-2b-v0.5	26.79	27.07	26.89	26.42	25.31	26.08	28.14	25.57	27.68	28.49	25.95	27.72
sarvam-1	29.70	28.41	29.21	29.66	27.27	28.92	30.82	27.57	30.24	30.92	26.66	29.64
Nemotron-4-Mini-Hindi-4B-Base	34.76	32.82	34.03	34.95	31.41	33.86	36.67	36.49	36.64	40.75	37.55	39.79
Nemotron-4-Mini-Hindi-4B-Instruct	33.38	33.82	33.54	31.98	30.06	31.39	35.83	35.33	35.74	36.99	34.11	36.12
Qwen2.5-3B	40.61	31.90	37.34	39.54	32.13	37.26	44.57	32.72	42.45	44.98	33.97	41.67
Qwen2.5-3B-Instruct	35.22	28.46	32.68	34.84	29.17	33.09	42.67	27.20	39.90	40.62	29.89	37.39
granite-3.1-2b-instruct	33.39	27.30	31.10	32.82	27.11	31.07	37.71	27.86	35.95	38.18	27.30	34.91
granite-3.1-3b-a800m-base	31.75	26.18	29.66	29.22	24.17	27.66	33.36	26.70	32.17	33.74	24.01	30.82
<i>7B to 27B Models</i>												
Pangea-7B	40.69	31.93	37.41	41.71	33.73	39.25	47.16	34.71	44.93	48.70	34.95	44.57
Indic-gemma-7b-finetuned-sft-Navarasa-2.0	37.12	31.83	35.13	37.00	30.47	34.90	42.31	33.44	40.73	44.08	34.09	41.08
aya-23-8B	33.84	28.87	31.97	35.25	30.88	33.90	37.09	33.22	36.40	41.92	33.01	39.24
Llama-3.1-8B	35.48	29.17	33.12	36.20	30.61	34.48	39.52	31.41	38.07	41.32	31.76	38.44
Llama-3.1-8B-Instruct	36.86	31.26	34.76	35.68	30.27	34.01	47.14	35.07	44.98	48.61	36.47	44.96
gemma-2-9b	48.16	37.92	44.32	42.73	36.91	40.94	55.23	43.89	53.20	58.49	42.96	53.83
gemma-2-9b-it	36.22	31.18	34.33	38.85	32.03	36.75	48.92	36.45	46.69	45.05	38.66	43.13
gpt-oss-20b	38.30	33.09	36.34	37.11	32.61	35.73	46.58	36.27	44.73	40.69	35.24	39.06
gemma-2-27b	50.70	42.26	47.53	47.79	41.24	45.77	59.84	50.38	58.14	64.91	51.83	60.99
gemma-2-27b-it	40.45	33.89	37.99	42.47	34.29	39.95	54.95	41.24	52.50	50.71	42.02	48.10
<i>> 27B Models</i>												
gpt-oss-120b	55.62	48.05	52.78	74.11	64.16	71.05	71.40	60.25	69.41	70.72	62.94	68.38
Qwen3-235B-A22B-Instruct-25076	60.25	54.78	58.20	63.72	56.27	61.43	74.57	64.13	72.70	80.15	68.60	76.68
deepseek-v3	51.38	37.03	45.99	63.46	57.04	61.48	62.93	45.01	59.73	67.78	46.78	61.47
gpt-4o	62.75	54.73	59.74	57.27	49.82	54.97	75.31	65.18	73.50	78.83	71.02	76.49

235B-A22B-Instruct-2507 achieved 66.4% and 65.18% on Easy questions, and 47.09% and 46.24% on Hard questions, while smaller models like gemma-3-270m scored 28.1% on Easy and 26.81% on Hard. A similar trend is observed in BBF, with Easy question scores ranging from 24.15% (gemma-3-270m) to 74.8% (gpt-oss-120b) and Hard questions from 21.22% to 62.61%. Medium-level questions show moderate differentiation, reflecting model reasoning capability. BBK and BBL follow the same pattern, with instruction-tuned and larger models consistently outperforming smaller models, particularly on Hard questions. Overall, model size, instruction tuning, and architecture significantly influence robustness to question difficulty and generalization across domains. See Appendix E.1.

5.4 PERFORMANCE ANALYSIS ACROSS QUESTION TYPES

We analyzed model performance on various question types including Assertion/Reasoning, Fill in the Blanks, MCQs, Match the Column, Reading Comprehension, and Rearrange the Sequence across the BBA, BBF, BBK, and BBL domains. In BBA, models like deepseek-v3 and GPT-4o achieved high scores of 66.67% and 62.96% on Assertion/Reasoning questions, whereas smaller models such as gemma-3-270m scored 28.09%. For Fill in the Blanks, scores ranged from 24.72% (gemma-3-270m-it) to 51.69% (Qwen3-235B-A22B-Instruct-2507). MCQ performance was moderate, between 26% and 59.95%. Match the Column and Reading Comprehension showed wider variation, with larger models consistently outperforming smaller or non-instruction-tuned models. Rearrange the Sequence proved challenging across domains, with top models reaching 71.43% in BBL. Overall, question type significantly affects performance, highlighting the importance of model size, instruction tuning, and reasoning capabilities in handling diverse formats.

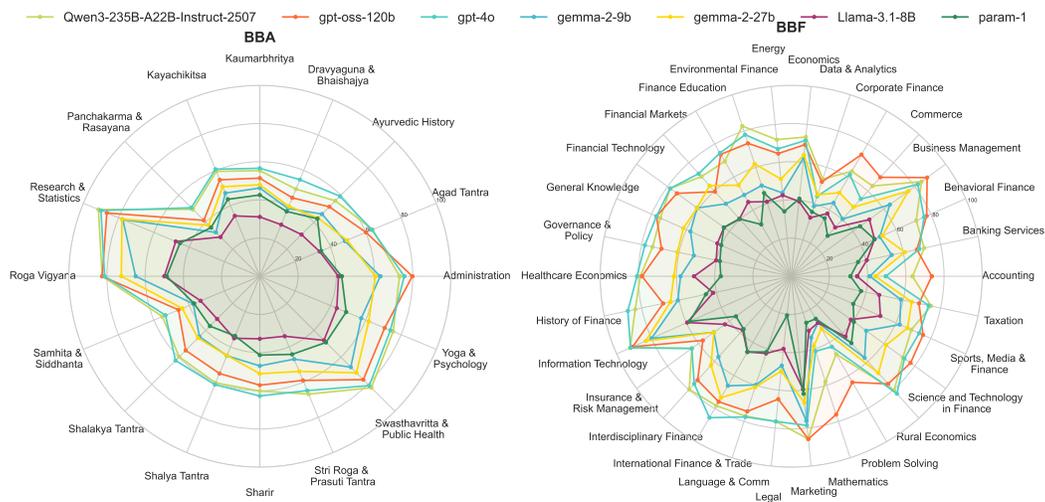


Figure 4: Comparison of representative LLMs' scores across different domains and subdomains.

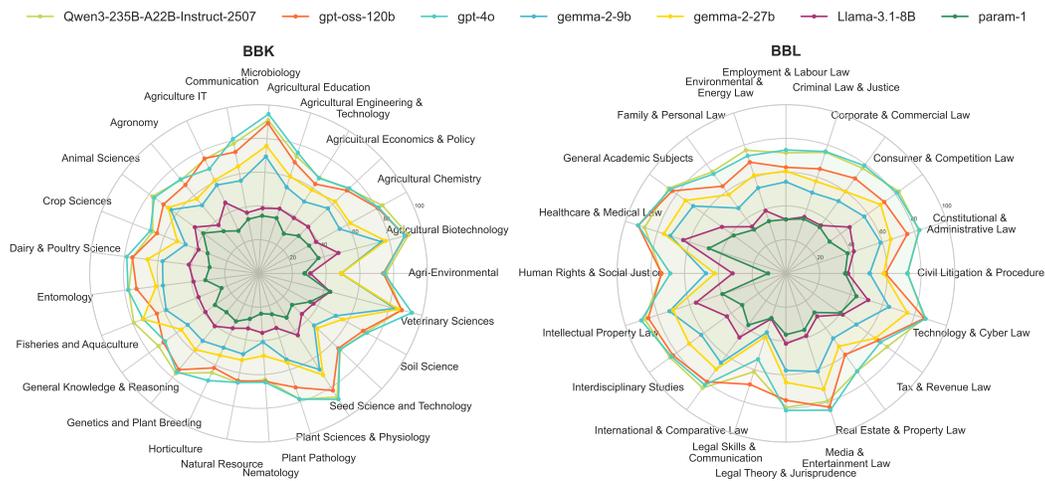


Figure 5: Comparison of representative LLMs' scores across different domains and subdomains.

5.5 PERFORMANCE ANALYSIS OF GPT MODEL FAMILY

We evaluate the GPT model family across BBA, BBF, BBK, and BBL domains to understand scaling and architectural strengths (Figure 3). gpt-oss-20b demonstrates baseline performance with scores of 36.34% (BBA), 35.73% (BBF), 44.73% (BBK), and 39.06% (BBL). Scaling to gpt-oss-120b yields substantial improvements: 52.78% in BBA, 71.05% in BBF, 69.41% in BBK, and 68.38% in BBL, representing 16-35 percentage point gains. Despite gpt-4o's larger parameter count, gpt-oss-120b significantly outperforms it in Finance (71.05% vs 54.97%), likely due to BBF's mathematical reasoning emphasis where gpt-oss-120b's training methodology excels (Analysis, 2025). Conversely, gpt-4o shows superior performance in Legal (76.49%) and Agriculture (73.5%) domains. This highlights that parameter size (Babbar, 2025) alone doesn't guarantee performance; architectural choices and training approaches significantly influence domain-specific capabilities, with mathematical tasks favoring specific optimizations over raw parameter scaling.

5.6 PERFORMANCE ANALYSIS OF SMALL MODELS

We evaluate small models ($\leq 4B$ parameters) across BBA, BBF, BBK, and BBL domains to assess efficiency-performance trade-offs (Figure 2). Param-1 and Qwen2.5-3B emerge as comparable top performers, with Param-1 achieving 39.97% in BBA while Qwen2.5-3B excels in BBK (42.45%). Both models demonstrate complementary strengths: Param-1 performs better in Ayurveda, while

Qwen2.5-3B shows superior performance in Finance, Agriculture, and Legal domains. Instruction tuning effects vary significantly across architectures: Llama-3.2-3B-Instruct substantially outperforms its base version, whereas Qwen2.5-3B-Instruct shows mixed results. Nemotron-4-Mini-Hindi models achieve competitive performance in the 34-40% range, while the smallest models like gemma-3-270m struggle consistently below 28%. Results indicate that architectural efficiency and targeted optimization can achieve reasonable performance in resource-constrained scenarios, with Param-1 and Qwen2.5 leading the small model category through different domain specializations.

5.7 ROBUSTNESS AND CONTAMINATION ANALYSIS

To verify BhashaBench V1’s integrity and rule out potential data leakage, we conduct perplexity (Jelinek et al., 1977) analysis on Llama-3.1-8B and Gemma-2-9B models. Table 3 shows BhashaBench V1 datasets (BBA, BBF, BBK, BBL) have perplexity scores comparable to or higher than established benchmarks like ARC-C, MMLU, and MILU. Notably, BBA shows higher perplexity (15.5 English, 10.39 Hindi on Llama-3.1-8B), while BBK is elevated in Hindi (7.16, 9.54) relative to English. To test potential position bias, we conduct option shuffling across multiple seeds, showing consistent performance with minimal variance. These patterns indicate minimal pretraining exposure and clearly show BhashaBench V1 provides genuinely novel evaluation challenges. Detailed perplexity and shuffling results are in Appendix E.4.

Table 3: Comparison of perplexity (PPL) across evaluation datasets. The PPL for English and Hindi datasets is reported for Llama-3.1-8B and gemma-2-9b models.

Datasets (MCQ)	Llama-3.1-8B		gemma-2-9b	
	English	Hindi	English	Hindi
ARC-C	8.03	4.1	6.82	5.85
MILU	7.62	4.93	7.23	6.37
MMLU	7.61	4.22	7.03	6.21
BBA	15.5	10.39	23.2	16.8
BBF	6.78	4.03	5.86	5.04
BBK	6.14	7.16	6.34	9.54
BBL	7.28	4.01	7.19	5.79

5.8 STATISTICAL SIGNIFICANCE ANALYSIS

To validate BhashaBench V1 results, we conduct statistical significance testing using two complementary approaches: we compute 95% Wilson Confidence Intervals (Wilson, 1927) for all models and domains, typically within 1-2 percentage points, demonstrating benchmark stability and reproducibility. We also perform pairwise McNemar’s tests (McNemar, 1947) on the top 5 models per domain to check whether performance differences are significant. Results reveal that most pairwise comparisons show statistically significant differences $p < 0.05$, confirming that observed performance gaps reflect genuine model capability rather than statistical noise. When accuracy differences are minimal (e.g., GPT-4o vs. Qwen3-235B-A22B-Instruct-25076 on BBA Hindi: 0.05% difference, $p=0.9591$), the test correctly identifies non-significant differences, demonstrating appropriate statistical rigor. Detailed results are in Appendix E.5.

6 CONCLUSION

In this paper, we introduced **BhashaBench V1**, a comprehensive, domain-specific, bilingual benchmark designed to evaluate large language models on India-centric knowledge systems across four critical domains: Agriculture (BBK), Legal (BBL), Finance (BBF), and Ayurveda (BBA). Our benchmark addresses significant gaps in existing evaluation frameworks by focusing on culturally relevant, domain-specific knowledge spanning over 90 subdomains and 500+ specialized topics curated from authentic government and professional examination papers. Our extensive evaluation reveals substantial performance disparities in current LLMs when applied to India-specific contexts, with models excelling in Legal contexts while struggling with traditional knowledge systems like Ayurveda and consistently performing better on English content compared to Hindi across all domains. These results highlight the urgent need for specialized model development strategies that incorporate India-specific knowledge, cultural contexts, and robust multilingual capabilities. To foster open research and accelerate progress toward more inclusive, culturally aware language models, we release BhashaBench V1 alongside all evaluation code and comprehensive documentation. We believe BhashaBench V1 offers a foundational benchmark for developing culturally sensitive models that effectively serve India’s diverse linguistic and knowledge landscape.

REFERENCES

- 540
541
542 Sanchit Ahuja, Divyanshu Aggarwal, Varun Gumma, Ishaan Watts, Ashutosh Sathe, Millicent
543 Ochieng, Rishav Hada, Prachi Jain, Maxamed Axmed, Kalika Bali, and Sunayana Sitaram. Mega-
544 verse: Benchmarking large language models across languages, modalities, models and tasks,
545 2024. URL <https://arxiv.org/abs/2311.07463>.
- 546 Artificial Analysis. gpt-oss-120b (high) vs gpt-4o: Model comparison. [https://](https://artificialanalysis.ai/models/comparisons/gpt-oss-120b-vs-gpt-4o)
547 artificialanalysis.ai/models/comparisons/gpt-oss-120b-vs-gpt-4o,
548 2025. Accessed: 2025-09-10.
- 549 AnhaltAI. Bvl q&a corpus 2024: German agricultural question-answer dataset. [https://](https://huggingface.co/datasets/anhaltai/bvl-qa-corpus-2024)
550 huggingface.co/datasets/anhaltai/bvl-qa-corpus-2024, 2024. Accessed
551 September 2025, Non-commercial use license.
- 552
553 Anthropic. Introducing the next generation of claude. [https://www.anthropic.com/](https://www.anthropic.com/news/claude-3-family)
554 [news/claude-3-family](https://www.anthropic.com/news/claude-3-family), March 2024.
- 555 Rakshit Aralimatti, Syed Abdul Gaffar Shakhadri, Kruthika KR, and Kartik Basavaraj Angadi.
556 Fine-tuning small language models for domain-specific ai: An edge ai perspective, 2025. URL
557 <https://arxiv.org/abs/2503.01933>.
- 558
559 Muhammad Arbab Arshad, Talukder Zaki Jubery, Tirtho Roy, Rim Nassiri, Asheesh K. Singh, Arti
560 Singh, Chinmay Hegde, Baskar Ganapathysubramanian, Aditya Balu, Adarsh Krishnamurthy,
561 and Soumik Sarkar. Leveraging vision language models for specialized agricultural tasks, 2025.
562 URL <https://arxiv.org/abs/2407.19617>.
- 563
564 Sushant Babbar. Openai gpt-oss vs gpt-4o: Which one is better?, Au-
565 gust 7 2025. URL [https://blog.getbind.co/2025/08/07/](https://blog.getbind.co/2025/08/07/openai-gpt-oss-vs-gpt-4o-which-one-is-better/)
566 [openai-gpt-oss-vs-gpt-4o-which-one-is-better/](https://blog.getbind.co/2025/08/07/openai-gpt-oss-vs-gpt-4o-which-one-is-better/). Accessed: November
567 22, 2025.
- 568 Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald
569 Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. The belebele
570 benchmark: a parallel reading comprehension dataset in 122 language variants. In *Proceedings*
571 *of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*
572 *Papers)*, pp. 749–775. Association for Computational Linguistics, 2024. doi: 10.18653/v1/2024.
573 [acl-long.44](http://dx.doi.org/10.18653/v1/2024.acl-long.44). URL <http://dx.doi.org/10.18653/v1/2024.acl-long.44>.
- 574 BharatGenAI. Agriparam, 2025a. URL [https://huggingface.co/bharatgenai/](https://huggingface.co/bharatgenai/AgriParam)
575 [AgriParam](https://huggingface.co/bharatgenai/AgriParam). Accessed: September 2025.
- 576
577 BharatGenAI. Ayurparam, 2025b. URL [https://huggingface.co/bharatgenai/](https://huggingface.co/bharatgenai/AyurParam)
578 [AyurParam](https://huggingface.co/bharatgenai/AyurParam). Accessed: September 2025.
- 579
580 BharatGenAI. Financeparam, 2025c. URL [https://huggingface.co/bharatgenai/](https://huggingface.co/bharatgenai/FinanceParam)
581 [FinanceParam](https://huggingface.co/bharatgenai/FinanceParam). Accessed: September 2025.
- 582
583 BharatGenAI. Legalparam, 2025d. URL [https://huggingface.co/bharatgenai/](https://huggingface.co/bharatgenai/LegalParam)
584 [LegalParam](https://huggingface.co/bharatgenai/LegalParam). Accessed: September 2025.
- 585
586 Stella Biderman, Hailey Schoelkopf, Lintang Sutawika, Leo Gao, Jonathan Tow, Baber Abbasi, Al-
587 ham Fikri Aji, Pawan Sasanka Ammanamanchi, Sidney Black, Jordan Clive, Anthony DiPofi,
588 Julen Etxaniz, Benjamin Fattori, Jessica Zosa Forde, Charles Foster, Jeffrey Hsu, Mimansa
589 Jaiswal, Wilson Y. Lee, Haonan Li, Charles Lovering, Niklas Muennighoff, Ellie Pavlick, Ja-
590 son Phang, Aviya Skowron, Samson Tan, Xiangru Tang, Kevin A. Wang, Genta Indra Winata,
591 François Yvon, and Andy Zou. Lessons from the trenches on reproducible evaluation of language
592 models, 2024a. URL <https://arxiv.org/abs/2405.14782>.
- 593
594 Stella Biderman, Hailey Schoelkopf, Lintang Sutawika, Leo Gao, Jonathan Tow, Baber Abbasi, Al-
595 ham Fikri Aji, Pawan Sasanka Ammanamanchi, Sidney Black, Jordan Clive, Anthony DiPofi,
596 Julen Etxaniz, Benjamin Fattori, Jessica Zosa Forde, Charles Foster, Jeffrey Hsu, Mimansa

- 594 Jaiswal, Wilson Y. Lee, Haonan Li, Charles Lovering, Niklas Muennighoff, Ellie Pavlick, Ja-
595 son Phang, Aviya Skowron, Samson Tan, Xiangru Tang, Kevin A. Wang, Genta Indra Winata,
596 François Yvon, and Andy Zou. Lessons from the trenches on reproducible evaluation of language
597 models, 2024b. URL <https://arxiv.org/abs/2405.14782>.
- 598
599 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
600 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,
601 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
602 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz
603 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
604 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL
605 <https://arxiv.org/abs/2005.14165>.
- 606 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece
607 Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi,
608 Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments
609 with gpt-4, 2023. URL <https://arxiv.org/abs/2303.12712>.
- 610 Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androut-
611 sopoulos. Legal-bert: The muppets straight out of law school, 2020. URL <https://arxiv.org/abs/2010.02559>.
- 612
613
614 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li,
615 Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Sto-
616 ica. Chatbot arena: An open platform for evaluating llms by human preference, 2024. URL
617 <https://arxiv.org/abs/2403.04132>.
- 618
619 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
620 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge,
621 2018. URL <https://arxiv.org/abs/1803.05457>.
- 622 DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Cheng-
623 gang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang,
624 Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting
625 Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui
626 Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi
627 Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li,
628 Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang,
629 Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun
630 Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan
631 Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J.
632 Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang,
633 Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng
634 Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shut-
635 ing Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanbiao Zhao,
636 Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue
637 Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xi-
638 aokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin
639 Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang,
640 Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang
641 Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui
642 Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying
643 Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu,
644 Yuan Ou, Yuchen Zhu, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan
645 Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F.
646 Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda
647 Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao,
Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li,
Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025. URL
<https://arxiv.org/abs/2412.19437>.

- 648 Jay Gala, Pranjal A. Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Ku-
649 mar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar,
650 Mitesh M. Khapra, Raj Dabre, and Anoop Kunchukuttan. Indictans2: Towards high-quality
651 and accessible machine translation models for all 22 scheduled indian languages, 2023. URL
652 <https://arxiv.org/abs/2305.16307>.
- 653
- 654 Jay Gala, Thanmay Jayakumar, Jaavid Aktar Husain, Aswanth Kumar M, Mohammed Safi Ur Rah-
655 man Khan, Diptesh Kanojia, Ratish Puduppully, Mitesh M. Khapra, Raj Dabre, Rudra Murthy,
656 and Anoop Kunchukuttan. Airavata: Introducing hindi instruction-tuned llm, 2024. URL
657 <https://arxiv.org/abs/2401.15006>.
- 658
- 659 Google. Introducing gemini: our largest and most capable ai model. [https://blog.google/](https://blog.google/technology/ai/google-gemini-ai/)
660 [technology/ai/google-gemini-ai/](https://blog.google/technology/ai/google-gemini-ai/), December 2023.
- 661
- 662 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
663 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan,
664 Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Ko-
665 renev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava
666 Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux,
667 Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret,
668 Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius,
669 Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary,
670 Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab
671 AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco
672 Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind That-
673 tai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Kore-
674 vaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra,
675 Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
676 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,
677 Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jong-
678 soo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala,
679 Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid
680 El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren
681 Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin,
682 Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi,
683 Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew
684 Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Ku-
685 mar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoy-
686 chev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan
687 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan,
688 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ra-
689 mon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Ro-
690 hit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan
691 Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell,
692 Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng
693 Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer
694 Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman,
695 Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mi-
696 haylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor
697 Kerkez, Vincent Gouget, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei
698 Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang
699 Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Gold-
700 schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning
701 Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papanikos, Aaditya Singh,
Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria,
Ahava Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein,
Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, An-
drew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, An-
nie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,

- 702 Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leon-
703 hardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu
704 Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Mon-
705 talvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao
706 Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia
707 Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide
708 Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le,
709 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily
710 Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smoth-
711 ers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni,
712 Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia
713 Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan,
714 Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harri-
715 son Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damla,
716 Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James
717 Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-
718 nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang,
719 Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Jun-
720 jie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy
721 Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang,
722 Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell,
723 Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa,
724 Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias
725 Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L.
726 Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike
727 Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari,
728 Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan
729 Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong,
730 Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent,
731 Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar,
732 Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Ro-
733 driguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,
734 Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin
735 Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon,
736 Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ra-
737 maswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha,
738 Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal,
739 Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satter-
740 field, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj
741 Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo
742 Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook
743 Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Ku-
744 mar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov,
745 Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiao-
746 jian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,
747 Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao,
748 Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhao-
749 duo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL
750 <https://arxiv.org/abs/2407.21783>.
- 748 Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex
749 Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, et al. Legalbench: A
750 collaboratively built benchmark for measuring legal reasoning in large language models, 2023.
751 URL <https://arxiv.org/abs/2308.11462>.
- 752
753 Jinglin He, Yunqi Guo, Lai Kwan Lam, Waikie Leung, Lixing He, Yuanan Jiang, Chi Chiu Wang,
754 Guoliang Xing, and Hongkai Chen. Opentcm: A graphrag-empowered llm-based system for
755 traditional chinese medicine knowledge retrieval and diagnosis, 2025. URL <https://arxiv.org/abs/2504.20118>.

- 756 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja-
757 cob Steinhardt. Measuring massive multitask language understanding, 2021a. URL <https://arxiv.org/abs/2009.03300>.
758
759
- 760 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja-
761 cob Steinhardt. Measuring massive multitask language understanding, 2021b. URL <https://arxiv.org/abs/2009.03300>.
762
- 763 IndiaDataMap. Projected 2025: Number of farmers in every indian
764 state, 2025. URL [https://indiadatamap.com/2025/08/19/
765 number-of-farmers-in-every-indian-state-2025/](https://indiadatamap.com/2025/08/19/number-of-farmers-in-every-indian-state-2025/).
766
- 767 Indian Knowledge Systems Division, Ministry of Education, Government of India. Indian knowl-
768 edge systems, 2025. URL <https://iksindia.org/>. Official website of the IKS Division
769 under Ministry of Education at AICTE, New Delhi.
- 770 Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie Vid-
771 gen. Financebench: A new benchmark for financial question answering, 2023. URL <https://arxiv.org/abs/2311.11944>.
772
773
- 774 Frederick Jelinek, Robert L. Mercer, Lalit R. Bahl, and Janet M. Baker. Perplexity—a measure of
775 the difficulty of speech recognition tasks. *Journal of the Acoustical Society of America*, 62, 1977.
776 URL <https://api.semanticscholar.org/CorpusID:121680873>.
- 777 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chap-
778 lot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
779 L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril,
780 Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
781
782
- 783 Rasmus J  rgensen, Oliver Brandt, Mareike Hartmann, Xiang Dai, Christian Igel, and Desmond
784 Elliott. MultiFin: A dataset for multilingual financial NLP. In Andreas Vlachos and Is-
785 abelle Augenstein (eds.), *Findings of the Association for Computational Linguistics: EACL*
786 *2023*, pp. 894–909, Dubrovnik, Croatia, May 2023. Association for Computational Linguis-
787 tics. doi: 10.18653/v1/2023.findings-eacl.66. URL [https://aclanthology.org/2023.
788 findings-eacl.66/](https://aclanthology.org/2023.findings-eacl.66/).
- 789 Abhinav Joshi, Shounak Paul, Akshat Sharma, Pawan Goyal, Saptarshi Ghosh, and Ashutosh Modi.
790 Il-tur: Benchmark for indian legal text understanding and reasoning, 2024. URL [https://
791 arxiv.org/abs/2407.05399](https://arxiv.org/abs/2407.05399).
- 792
793 Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, Gokul N.C., Avik Bhattacharyya,
794 Mitesh M. Khapra, and Pratyush Kumar. IndicNLPsuite: Monolingual corpora, evaluation bench-
795 marks and pre-trained multilingual language models for Indian languages. In Trevor Cohn,
796 Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 4948–4961, Online, November 2020. Association for Computational Lin-
797 guistics. doi: 10.18653/v1/2020.findings-emnlp.445. URL [https://aclanthology.org/
798 2020.findings-emnlp.445/](https://aclanthology.org/2020.findings-emnlp.445/).
- 799
- 800 Mohammed Khan, Priyam Mehta, Ananth Sankar, Umashankar Kumaravelan, Sumanth Dooda-
801 paneni, Suriyaprasaad B, Varun G, Sparsh Jain, Anoop Kunchukuttan, Pratyush Kumar, Raj
802 Dabre, and Mitesh Khapra. Indicllmsuite: A blueprint for creating pre-training and fine-
803 tuning datasets for indian languages. In *Proceedings of the 62nd Annual Meeting of the As-
804 sociation for Computational Linguistics (Volume 1: Long Papers)*, pp. 15831–15879. Associa-
805 tion for Computational Linguistics, 2024. doi: 10.18653/v1/2024.acl-long.843. URL [http:
806 //dx.doi.org/10.18653/v1/2024.acl-long.843](http://dx.doi.org/10.18653/v1/2024.acl-long.843).
- 807
808 Simran Khanuja, Diksha Bansal, Sarvesh Mehtani, Savya Khosla, Atreyee Dey, Balaji Gopalan,
809 Dilip Kumar Margam, Pooja Aggarwal, Rajiv Teja Nagipogu, Shachi Dave, Shruti Gupta, Sub-
hash Chandra Bose Gali, Vish Subramanian, and Partha Talukdar. Muril: Multilingual represen-
tations for indian languages, 2021a. URL <https://arxiv.org/abs/2103.10730>.

- 810 Simran Khanuja, Diksha Bansal, Sarvesh Mehtani, Savya Khosla, Atreyee Dey, Balaji Gopalan,
811 Dilip Kumar Margam, Pooja Aggarwal, Rajiv Teja Nagipogu, Shachi Dave, Shruti Gupta, Sub-
812 hash Chandra Bose Gali, Vish Subramanian, and Partha P. Talukdar. MuriL: Multilingual repre-
813 sentations for indian languages. *CoRR*, abs/2103.10730, 2021b. URL [https://arxiv.org/
814 abs/2103.10730](https://arxiv.org/abs/2103.10730).
- 815 Shufeng Kong, Xingru Yang, Yuanyuan Wei, Zijie Wang, Hao Tang, Jiuqi Qin, Shuting Lan,
816 Yingheng Wang, Junwen Bai, Zhuangbin Chen, Zibin Zheng, Caihua Liu, and Hao Liang. Mtcmb:
817 A multi-task benchmark framework for evaluating llms on knowledge, reasoning, and safety in
818 traditional chinese medicine, 2025. URL <https://arxiv.org/abs/2506.01252>.
- 820 Josué Kpodo, Parisa Kordjamshidi, and A. Pouyan Nejadhashemi. Agxqa: A benchmark
821 for advanced agricultural extension question answering. *Computers and Electronics in
822 Agriculture*, 225:109349, 2024. ISSN 0168-1699. doi: [https://doi.org/10.1016/j.compag.
823 2024.109349](https://doi.org/10.1016/j.compag.2024.109349). URL [https://www.sciencedirect.com/science/article/pii/
824 S0168169924007403](https://www.sciencedirect.com/science/article/pii/S0168169924007403).
- 825 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
826 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
827 serving with pagedattention, 2023. URL <https://arxiv.org/abs/2309.06180>.
- 828 Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jae-
829 woo Kang. Biobert: a pre-trained biomedical language representation model for biomedical
830 text mining. *Bioinformatics*, 36(4):1234–1240, September 2019. ISSN 1367-4811. doi: 10.
831 1093/bioinformatics/btz682. URL [http://dx.doi.org/10.1093/bioinformatics/
832 btz682](http://dx.doi.org/10.1093/bioinformatics/btz682).
- 833 Haohang Li, Yupeng Cao, Yangyang Yu, Shashidhar Reddy Javaji, Zhiyang Deng, Yueru He,
834 Yuechen Jiang, Zining Zhu, Koduvayur Subbalakshmi, Guojun Xiong, Jimin Huang, Lingfei
835 Qian, Xueqing Peng, Qianqian Xie, and Jordan W. Suchow. Investorbench: A benchmark for
836 financial decision-making tasks with llm-based agent, 2024. URL [https://arxiv.org/
837 abs/2412.18174](https://arxiv.org/abs/2412.18174).
- 838 Yash Madhani, Mitesh M. Khapra, and Anoop Kunchukuttan. Bhasha-abhijnaanam: Native-
839 script and romanized language identification for 22 indian languages. In *Findings of the
840 Association for Computational Linguistics: ACL 2023*, 2023. URL [https://github.
841 com/AI4Bharat/IndicLID](https://github.com/AI4Bharat/IndicLID). Models and datasets available at [https://github.com/
842 AI4Bharat/IndicLID](https://github.com/AI4Bharat/IndicLID) and
843 <https://huggingface.co/datasets/ai4bharat/Bhasha-Abhijnaanam>.
- 844 Arijit Maji, Raghendra Kumar, Akash Ghosh, Anushka, and Sriparna Saha. Sanskriti: A compre-
845 hensive benchmark for evaluating language models’ knowledge of indian culture, 2025. URL
846 <https://arxiv.org/abs/2506.15355>.
- 847 Quinn McNemar. Note on the sampling error of the difference between correlated proportions or
848 percentages. *Psychometrika*, 12:153–157, 1947. URL [https://api.semanticscholar.
849 org/CorpusID:46226024](https://api.semanticscholar.org/CorpusID:46226024).
- 850 Arijit Nag, Soumen Chakrabarti, Animesh Mukherjee, and Niloy Ganguly. Efficient continual pre-
851 training of llms for low-resource languages, 2024. URL [https://arxiv.org/abs/2412.
852 10244](https://arxiv.org/abs/2412.10244).
- 853 National Judicial Data Grid (NJDG), India. National judicial data grid - statistics 2023. [https:
854 //njdg.ecourts.gov.in/](https://njdg.ecourts.gov.in/), 2023. Accessed September 2025.
- 855 National Payments Corporation of India. Npci upi transaction statistics 2023. [https://www.
856 npci.org.in/](https://www.npci.org.in/), 2023. Accessed September 2025.
- 857 Xuan-Phi Nguyen, Wenxuan Zhang, Xin Li, Mahani Aljunied, Zhiqiang Hu, Chenhui Shen,
858 Yew Ken Chia, Xingxuan Li, Jianyu Wang, Qingyu Tan, Liying Cheng, Guanzheng Chen, Yue
859 Deng, Sen Yang, Chaoqun Liu, Hang Zhang, and Lidong Bing. Seallms – large language models
860 for southeast asia, 2024. URL <https://arxiv.org/abs/2312.00738>.

- 864 Shubham Kumar Nigam, Shubham Kumar Mishra, Ayush Kumar Mishra, Noel Shallum, and Arnab
865 Bhattacharya. Legal question-answering in the indian context: Efficacy, challenges, and potential
866 of modern ai models, 2023. URL <https://arxiv.org/abs/2309.14735>.
867
- 868 Joel Niklaus, Veton Matoshi, Pooja Rani, Andrea Galassi, Matthias Stürmer, and Ilias Chalkidis.
869 Lextreme: A multi-lingual and multi-task benchmark for the legal domain. In *Findings of*
870 *the Association for Computational Linguistics: EMNLP 2023*, pp. 3016–3054. Association for
871 Computational Linguistics, 2023. doi: 10.18653/v1/2023.findings-emnlp.200. URL <http://dx.doi.org/10.18653/v1/2023.findings-emnlp.200>.
872
- 873 OpenAI. Hello gpt-4o, 2024. URL <https://openai.com/index/hello-gpt-4o/>.
874
- 875 OpenAI. Introducing gpt-oss: Open-weight reasoning models from openai. <https://openai.com/index/introducing-gpt-oss/>, 2024. Accessed September 2025.
876
- 877 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-
878 cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red
879 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Moham-
880 mad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher
881 Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-
882 man, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann,
883 Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis,
884 Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey
885 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux,
886 Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila
887 Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,
888 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-
889 son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan
890 Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hal-
891 lacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan
892 Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu,
893 Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun
894 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-
895 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook
896 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel
897 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kopic, Gretchen
898 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel
899 Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez,
900 Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv
901 Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney,
902 Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick,
903 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel
904 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-
905 jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe,
906 Jakob Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel
907 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe
908 de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny,
909 Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl,
910 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra
911 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders,
912 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-
913 sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor,
914 Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,
915 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,
916 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-
917 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-
jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan
Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng,
Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Work-
man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming

- 918 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao
919 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL
920 <https://arxiv.org/abs/2303.08774>.
921
- 922 OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin
923 Arbus, Rahul K. Arora, Yu Bai, Bowen Baker, Haiming Bao, Boaz Barak, Ally Bennett, Tyler
924 Bertao, Nivedita Brett, Eugene Brevdo, Greg Brockman, Sebastien Bubeck, Che Chang, Kai
925 Chen, Mark Chen, Enoch Cheung, Aidan Clark, Dan Cook, Marat Dukhan, Casey Dvorak, Kevin
926 Fives, Vlad Fomenko, Timur Garipov, Kristian Georgiev, Mia Glaese, Tarun Gogineni, Adam
927 Goucher, Lukas Gross, Katia Gil Guzman, John Hallman, Jackie Hehir, Johannes Heidecke, Alec
928 Helyar, Haitang Hu, Romain Huet, Jacob Huh, Saachi Jain, Zach Johnson, Chris Koch, Irina
929 Kofman, Dominik Kundel, Jason Kwon, Volodymyr Kyrlyov, Elaine Ya Le, Guillaume Leclerc,
930 James Park Lennon, Scott Lessans, Mario Lezcano-Casado, Yuanzhi Li, Zhuohan Li, Ji Lin,
931 Jordan Liss, Lily, Liu, Jiancheng Liu, Kevin Lu, Chris Lu, Zoran Martinovic, Lindsay McCal-
932 lum, Josh McGrath, Scott McKinney, Aidan McLaughlin, Song Mei, Steve Mostovoy, Tong Mu,
933 Gideon Myles, Alexander Neitz, Alex Nichol, Jakub Pachocki, Alex Paino, Dana Palmie, Ash-
934 ley Pantuliano, Giambattista Parascandolo, Jongsoo Park, Leher Pathak, Carolina Paz, Ludovic
935 Peran, Dmitry Pimenov, Michelle Pokrass, Elizabeth Proehl, Huida Qiu, Gaby Raila, Filippo
936 Raso, Hongyu Ren, Kimmy Richardson, David Robinson, Bob Rotsted, Hadi Salman, Suvansh
937 Sanjeev, Max Schwarzer, D. Sculley, Harshit Sikchi, Kendal Simon, Karan Singhal, Yang Song,
938 Dane Stuckey, Zhiqing Sun, Philippe Tillet, Sam Toizer, Foivos Tsimpourlas, Nikhil Vyas, Eric
939 Wallace, Xin Wang, Miles Wang, Olivia Watkins, Kevin Weil, Amy Wendling, Kevin Whinnery,
940 Cedric Whitney, Hannah Wong, Lin Yang, Yu Yang, Michihiro Yasunaga, Kristen Ying, Wojciech
941 Zaremba, Wenting Zhan, Cyril Zhang, Brian Zhang, Eddie Zhang, and Shengjia Zhao. gpt-oss-
942 120b & gpt-oss-20b model card, 2025. URL <https://arxiv.org/abs/2508.10925>.
- 942 Vikram Paruchuri and DataLab Team. Surya: A multilingual document ocr toolkit for indic and
943 other languages. <https://github.com/datalab-to/surya>, 2024. Software available
944 at <https://github.com/datalab-to/surya>.
- 945 Xueqing Peng, Lingfei Qian, Yan Wang, Ruoyu Xiang, Yueru He, Yang Ren, Mingyang Jiang,
946 Jeff Zhao, Huan He, Yi Han, Yun Feng, Yuechen Jiang, Yupeng Cao, Haohang Li, Yangyang
947 Yu, Xiaoyu Wang, Penglei Gao, Shengyuan Lin, Keyi Wang, Shanshan Yang, Yilun Zhao, Zhi-
948 wei Liu, Peng Lu, Jerry Huang, Suyuchen Wang, Triantafillos Papadopoulos, Polydoros Gian-
949 nouris, Efstathia Soufleri, Nuo Chen, Guojun Xiong, Zhiyang Deng, Yijia Zhao, Mingquan Lin,
950 Meikang Qiu, Kaleb E Smith, Arman Cohan, Xiao-Yang Liu, Jimin Huang, Alejandro Lopez-
951 Lira, Xi Chen, Junichi Tsujii, Jian-Yun Nie, Sophia Ananiadou, and Qianqian Xie. Multifinben:
952 A multilingual, multimodal, and difficulty-aware benchmark for financial llm evaluation, 2025.
953 URL <https://arxiv.org/abs/2506.14028>.
- 954 Kundeshwar Pundalik, Piyush Sawarkar, Nihar Sahoo, Abhishek Shinde, Prateek Chanda, Vedant
955 Goswami, Ajay Nagpal, Atul Singh, Viraj Thakur, Vijay Dewane, Aamod Thakur, Bhargav Patel,
956 Smita Gautam, Bhagwan Panditi, Shyam Pawar, Madhav Kotcha, Suraj Racha, Saral Sureka,
957 Pankaj Singh, Rishi Bal, Rohit Saluja, and Ganesh Ramakrishnan. Param-1 bharatgen 2.9b model,
958 2025. URL <https://arxiv.org/abs/2507.13390>.
- 959 Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan
960 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,
961 Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin
962 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li,
963 Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,
964 Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025.
965 URL <https://arxiv.org/abs/2412.15115>.
- 966 Shriram Shivajirao Ragad and Maya Vivek Gokhale. Ayurvedic concept of koshtha and its impor-
967 tance in panchkarma. *International Journal of Research - Granthaalayah*, 7(7):416–421, 2019.
968 doi: 10.5281/zenodo.3370488. URL <https://doi.org/10.5281/zenodo.3370488>.
- 969 Anjali Sen and et al. Morphological understanding and retrieval in indic languages. *Journal of Indic
970 Computational Linguistics*, 12(3):45–67, 2023.
971

- 972 Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark,
973 Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed
974 Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska,
975 Blaise Aguera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara
976 Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan
977 Karthikesalingam, and Vivek Natarajan. Towards expert-level medical question answering with
978 large language models, 2023. URL <https://arxiv.org/abs/2305.09617>.
- 979
- 980 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam
981 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska,
982 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W.
983 Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain,
984 Amanda Askeel, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, An-
985 ders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, An-
986 drew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh
987 Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabas-
988 sum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Her-
989 rick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph,
990 Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin
991 Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinon, Cameron
992 Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh,
993 Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites,
994 Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera,
995 Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Gar-
996 rette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy,
997 Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito,
998 Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, De-
999 nis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta
1000 Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Eka-
1001 terina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Eliza-
1002 beth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem,
1003 Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Ev-
1004 genii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé,
1005 Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán
1006 Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-
1007 López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh
1008 Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hong-
1009 ming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson
1010 Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel,
1011 James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema
1012 Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova,
1013 Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Ji-
1014 acheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis,
1015 Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph
1016 Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua,
1017 Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja
1018 Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chia-
1019 fullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo
1020 Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency,
1021 Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón,
1022 Luke Metz, Lütfi Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Fa-
1023 rooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria
1024 Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast,
1025 Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody
Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy,
Michael Starritt, Michael Strube, Michał Swedrowski, Michele Bevilacqua, Michihiro Yasunaga,
Mihir Kale, Mike Cain, Mimeo Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal,
Moin Aminnaseri, Mor Geva, Mozhdah Gheini, Mukund Varma T, Nanyun Peng, Nathan A.
Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron,

- 1026 Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar,
1027 Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar El-
1028 baghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung,
1029 Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Pe-
1030 ter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour,
1031 Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer
1032 Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A.
1033 Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Ro-
1034 man Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov,
1035 Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Moham-
1036 mad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R.
1037 Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghaz-
1038 arian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schus-
1039 ter, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar
1040 Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upad-
1041 hyay, Shyamolima, Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy,
1042 Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene,
1043 Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Pianta-
1044 dosi, Stuart M. Shieber, Summer Mishnerghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen,
1045 Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore
1046 Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Ti-
1047 tus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz,
1048 Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh,
1049 Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saun-
1050 ders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong
1051 Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi
1052 Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary
1053 Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the im-
1054 itation game: Quantifying and extrapolating the capabilities of language models, 2023. URL
1055 <https://arxiv.org/abs/2206.04615>.
- 1056 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya
1057 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard
1058 Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex
1059 Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, An-
1060 tonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo,
1061 Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric
1062 Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Hen-
1063 ryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski,
1064 Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu,
1065 Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee,
1066 Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev,
1067 Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko
1068 Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo
1069 Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree
1070 Pandya, Siamak Shakeri, Soham De, Ted Klimentko, Tom Hennigan, Vlad Feinberg, Wojciech
1071 Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh
1072 Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin
1073 Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah
1074 Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on
1075 gemini research and technology, 2024. URL <https://arxiv.org/abs/2403.08295>.
- 1076 Hrishikesh Terdalkar, Vishakha Deulgaonkar, and Arnab Bhattacharya. Āyurjñānam: Explor-
1077 ing Āyurveda using knowledge graphs, 2023. URL [https://sanskrit.iitk.ac.in/
1078 ayurveda/](https://sanskrit.iitk.ac.in/ayurveda/). Presented at the National Youth Conference on Indian Knowledge Systems 2023.
- 1079 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher,
Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy
Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,

- 1080 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
1081 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
1082 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
1083 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
1084 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
1085 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
1086 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
1087 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models,
1088 2023. URL <https://arxiv.org/abs/2307.09288>.
- 1089 Sshubam Verma, Mohammed Safi Ur Rahman Khan, Vishwajeet Kumar, Rudra Murthy, and Jaydeep
1090 Sen. Milu: A multi-task indic language understanding benchmark, 2025. URL <https://arxiv.org/abs/2411.02538>.
- 1091
1092 Abhay Vijayvargia, Ajay Nagpal, Kundeshwar Pundalik, Atharva Savarkar, Smita Gautam, Pankaj
1093 Singh, Rohit Saluja, and Ganesh Ramakrishnan. Intent aware context retrieval for multi-turn
1094 agricultural question answering, 2025. URL <https://arxiv.org/abs/2508.03719>.
- 1095
1096 Jiayi Wang, Yao Lu, Maurice Weber, Max Ryabinin, Yihong Chen, Raphael Tang, and Pontus Stene-
1097 torp. Multilingual pretraining using a large corpus machine-translated from a single source lan-
1098 guage, 2024a. URL <https://arxiv.org/abs/2410.23956>.
- 1099
1100 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
1101 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi
1102 Fan, Xiang Yue, and Wenhui Chen. Mmlu-pro: A more robust and challenging multi-task language
1103 understanding benchmark, 2024b. URL <https://arxiv.org/abs/2406.01574>.
- 1104
1105 Yuxia Wang, Minghan Wang, Muhammad Arslan Manzoor, Fei Liu, Georgi Georgiev, Rocktim Jy-
1106 oti Das, and Preslav Nakov. Factuality of large language models: A survey. *arXiv preprint*
1107 *arXiv:2402.02420*, 2024c. doi: 10.48550/arXiv.2402.02420. URL <https://arxiv.org/abs/2402.02420>.
- 1108
1109 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-
1110 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol
1111 Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models,
2022. URL <https://arxiv.org/abs/2206.07682>.
- 1112
1113 Edwin B. Wilson. Probable inference, the law of succession, and statistical inference. *Journal*
1114 *of the American Statistical Association*, 22(158):209–212, 1927. doi: 10.1080/01621459.1927.
1115 10502953. URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.1927.10502953>.
- 1116
1117 Genta Indra Winata, Andrea Madotto, Zhaojiang Lin, Rosanne Liu, Jason Yosinski, and Pascale
1118 Fung. Language models are few-shot multilingual learners, 2021. URL <https://arxiv.org/abs/2109.07684>.
- 1119
1120 Chaojun Xiao, Haoxi Zhong, Zhipeng Guo, Cunchao Tu, Zhiyuan Liu, Maosong Sun, Yansong
1121 Feng, Xianpei Han, Zhen Hu, Heng Wang, and Jianfeng Xu. Cail2018: A large-scale legal dataset
1122 for judgment prediction, 2018. URL <https://arxiv.org/abs/1807.02478>.
- 1123
1124 Chaojun Xiao, Haoxi Zhong, Zhipeng Guo, Cunchao Tu, Zhiyuan Liu, Maosong Sun, Tianyang
1125 Zhang, Xianpei Han, Zhen Hu, Heng Wang, and Jianfeng Xu. Cail2019-scm: A dataset of similar
1126 case matching in legal domain, 2019. URL <https://arxiv.org/abs/1911.08962>.
- 1127
1128 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
1129 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
1130 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
1131 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
1132 Le Yu, Lianhao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
1133 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.

- 1134 Yi Yang, Mark Christopher Siy UY, and Allen Huang. Finbert: A pretrained language model for
1135 financial communications, 2020. URL <https://arxiv.org/abs/2006.08097>.
1136
- 1137 Jie Ying, Zihong Chen, Zhefan Wang, Wanli Jiang, Chenyang Wang, Zhonghang Yuan, Haoyang
1138 Su, Huanjun Kong, Fan Yang, and Nanqing Dong. Seedbench: A multi-task benchmark for
1139 evaluating large language models in seed science, 2025. URL <https://arxiv.org/abs/2505.13220>.
1140
- 1141 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-
1142 chine really finish your sentence?, 2019. URL <https://arxiv.org/abs/1905.07830>.
1143
- 1144 Lingfeng Zeng, Fangqi Lou, Zixuan Wang, Jiajie Xu, Jinyi Niu, Mengping Li, Yifan Dong, Qi Qi,
1145 Wei Zhang, Ziwei Yang, Jun Han, Ruilun Feng, Ruiqi Hu, Lejie Zhang, Zhengbo Feng, Yicheng
1146 Ren, Xin Guo, Zhaowei Liu, Dongpo Cheng, Weige Cai, and Liwen Zhang. Fingaia: A chinese
1147 benchmark for ai agents in real-world financial domain, 2025. URL <https://arxiv.org/abs/2507.17186>.
1148
- 1149 Tianyang Zhong, Zhenyuan Yang, Zhengliang Liu, Ruidong Zhang, Yiheng Liu, Haiyang Sun,
1150 Yi Pan, Yiwei Li, Yifan Zhou, Hanqi Jiang, Junhao Chen, and Tianming Liu. Opportunities and
1151 challenges of large language models for low-resource languages in humanities research, 2025.
1152 URL <https://arxiv.org/abs/2412.04497>.
- 1153 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied,
1154 Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation
1155 models, 2023. URL <https://arxiv.org/abs/2304.06364>.
1156
- 1157 Yunshun Zhong and Sebastian D. Goodfellow. Domain-specific language models pre-trained on
1158 construction management systems corpora. *Automation in Construction*, 160:105316, 2024.
1159 ISSN 0926-5805. doi: <https://doi.org/10.1016/j.autcon.2024.105316>. URL <https://www.sciencedirect.com/science/article/pii/S0926580524000529>.
1160
- 1161 Yutong Zhou and Masahiro Ryo. Agribench: A hierarchical agriculture benchmark for multimodal
1162 large language models, 2024. URL <https://arxiv.org/abs/2412.00465>.
- 1163 Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke
1164 Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blun-
1165 som, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker.
1166 Aya model: An instruction finetuned open-access multilingual language model, 2024. URL
1167 <https://arxiv.org/abs/2402.07827>.
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Appendix

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

CONTENTS

1	Introduction	1
2	RELATED WORK	3
2.1	Exploration of LLMs	3
2.2	Evaluation of LLMs	3
3	BhashaBench V1	4
3.1	Design Principles	4
3.2	Data Collection	5
3.3	Data Processing	5
3.4	Data Analysis	6
4	Experimental Setup	6
5	Results and Discussions	7
5.1	Zero-Shot Performance Across All Domains (EN + HI)	7
5.2	How do models perform in subdomains	7
5.3	Performance Analysis Across Question Difficulty Levels	7
5.4	Performance Analysis Across Question Types	8
5.5	Performance Analysis of GPT Model Family	9
5.6	Performance Analysis of Small Models	9
5.7	Robustness and Contamination Analysis	10
5.8	Statistical Significance Analysis	10
6	Conclusion	10
A	Limitations and Biases	25
B	Towards Broader Impact	25
C	More Details on BhashaBench V1	26
C.1	Details of Data Collection and Processing	26
C.1.1	Examination Source Documentation	26
C.1.2	Processing Pipeline Architecture	27
C.1.3	OCR Quality Assurance	30
C.1.4	Annotation Guidelines	32
C.1.5	Data Processing Prompts	32
C.2	Detailed Data Analysis of BhashaBench V1	42

1242	D More Details on Experiment Setup	45
1243		
1244	D.1 Task Formatting Template Used in LM Eval	45
1245	D.2 Task Formatting Template Used in API-Driven Evaluation	45
1246	D.3 Inference Implementation Details	46
1247		
1248	D.3.1 Open-Source Models	46
1249	D.3.2 API-Based Proprietary Models	46
1250		
1251	D.4 Evaluation Protocol and Response Processing	46
1252	D.4.1 Open-Source Models	46
1253	D.4.2 API-Based Models	46
1254		
1255	D.5 Reproducibility and Computational Resources	46
1256		
1257	E More Details on Experiment	48
1258		
1259	E.1 Zero-Shot Question-Level and Question-Type Performance Across	
1260	BhashaBench V1 Domains	48
1261	E.2 Zero-Shot sub-domain wise Performance Across BhashaBench V1 Domains .	48
1262	E.3 Qualitative Error Analysis for Llama-3.1-8B	56
1263		
1264	E.3.1 BBA Qualitative Analysis	56
1265	E.3.2 BBF Qualitative Analysis	59
1266	E.3.3 BBK Qualitative Analysis	61
1267	E.3.4 BBL Qualitative Analysis	63
1268		
1269	E.4 Data Integrity and Contamination Analysis	65
1270		
1271	E.4.1 Perplexity-Based Data Contamination Analysis	65
1272	E.4.2 Multiple-Choice Option Shuffling Experiment	65
1273	E.4.3 Effect of Scaling the Number of Distractors	65
1274		
1275	E.5 Statistical Significance Tests	67
1276		
1277	E.5.1 Wilson Confidence Interval for Model Performance	67
1278	E.5.2 Statistical Significance of Model Performance Differences Using Mc-	
1279	Nemar’s Test	69
1280		
1281		
1282		
1283		
1284		
1285		
1286		
1287		
1288		
1289		
1290		
1291		
1292		
1293		
1294		
1295		

1296 A LIMITATIONS AND BIASES

1297
1298 In this paper, we introduce BhashaBench V1, providing a comprehensive evaluation of LLMs on
1299 India-centric knowledge systems and exploring model capabilities across critical Indian domains.
1300 However, there are several limitations to acknowledge. (1) Language Coverage Limitations: Al-
1301 though BhashaBench V1 supports English and Hindi, covering a significant portion of India’s pop-
1302 ulation, India has 22 official languages and hundreds of regional dialects. Our current evaluation
1303 cannot capture the full linguistic diversity of Indian knowledge systems, particularly regional varia-
1304 tions in agricultural practices, legal terminologies, and traditional medicine nomenclature that exist
1305 in languages like Tamil, Telugu, Bengali, and others. Future iterations will expand to include addi-
1306 tional Indian languages to enhance coverage. (2) Domain Scope Limitations: While we cover four
1307 fundamental domains (Agriculture, Legal, Finance, and Ayurveda) representing core areas of Indian
1308 society, our assessment cannot encompass the entire breadth of India-specific knowledge systems.
1309 Areas such as traditional crafts, regional governance systems, indigenous engineering practices, and
1310 other vernacular knowledge traditions remain unexplored for future expansion. Our content spans
1311 from grassroots practical knowledge to professional examination standards, ensuring broad appli-
1312 cability across different expertise levels. (3) Evaluation Methodology Limitations: Our evaluation
1313 primarily uses structured question formats derived from authentic government and professional ex-
1314 aminations. While this ensures real-world relevance and practical applicability, it may not fully
capture all forms of contextual reasoning required in complex domain applications.

1315 The main biases in BhashaBench V1 can be categorized into three aspects: (1) Source Material
1316 Bias: Despite comprehensive curation from diverse authentic sources spanning grassroots to pro-
1317 fessional levels, certain regional practices and emerging contemporary developments may be un-
1318 derrepresented. (2) Language Resource Bias: The benchmark reflects the inherent resource dis-
1319 parity between English and Hindi, where Hindi content, while substantial, represents a relatively
1320 lower-resource context compared to English. (3) Examination Framework Bias: Our reliance on
1321 established examination systems, while ensuring authenticity, may introduce institutional perspec-
1322 tives present in the original assessment frameworks. However, our extensive coverage across 90+
1323 subdomains and 500+ topics from diverse sources mitigates this bias significantly. The impact of
1324 these limitations on LLM evaluation includes clear performance distinctions between models across
1325 domains and languages, as evidenced by the substantial score variations from 34.28% to 76.49%,
1326 demonstrating BhashaBench V1’s effectiveness in distinguishing LLM capabilities while presenting
1327 meaningful challenges even for top-performing models in India-specific contexts.

1328 B TOWARDS BROADER IMPACT

1329
1330
1331 **Societal Impact.** BhashaBench V1 is anticipated to play a transformative role in bridging the digi-
1332 tal divide for India-centric knowledge systems. LLMs trained and evaluated with BhashaBench V1
1333 can significantly enhance accessibility to critical domain expertise across agriculture, legal services,
1334 finance, and traditional medicine, particularly benefiting underserved rural and semi-urban popula-
1335 tions. In agriculture, improved LLM capabilities can democratize access to expert crop advisory,
1336 pest management, and sustainable farming practices, potentially impacting the livelihoods of over
1337 40 million farmers dependent on agricultural activities. In the legal domain, enhanced models can
1338 assist with legal document comprehension, procedural guidance, and basic legal literacy, address-
1339 ing the substantial access-to-justice challenges faced by millions in India’s complex legal system.
1340 For healthcare, particularly Ayurveda, better model performance can support practitioners and pa-
1341 tients in understanding traditional treatment protocols and medicinal formulations, preserving and
1342 disseminating indigenous medical knowledge. In finance, improved model capabilities can enhance
1343 financial literacy and support the growing digital payment ecosystem processing billions of trans-
1344 actions annually. However, we acknowledge potential risks including over-reliance on automated
1345 systems for critical decisions, potential displacement of traditional knowledge practitioners, and the
1346 risk of perpetuating biases present in examination-based evaluation systems. The benchmark’s fo-
1347 cus on professional examination standards, while ensuring quality, may inadvertently favor formal
educational backgrounds over experiential knowledge.

1348 **Ethics Statement.** We ensure strict adherence to applicable laws and ethical guidelines throughout
1349 our data collection, curation, and usage processes. All question-answer pairs are sourced exclusively
from publicly available government and professional examination papers, respecting intellectual

property rights and ensuring no unauthorized reproduction of copyrighted materials. Our curation process involved diverse teams to minimize cultural and regional biases, though we acknowledge the inherent limitations of our current English and Hindi coverage. The dataset contains no personally identifiable information, offensive content, or culturally insensitive material. All content has been thoroughly verified for authenticity and accuracy through multiple validation rounds involving domain experts. BhashaBench V1 is intended solely for academic research and educational purposes to advance inclusive AI development for Indian contexts. Any commercial use, misuse for harmful applications, or deployment without appropriate safeguards is strictly prohibited. We strongly urge all users to employ this resource responsibly, ensuring that any models developed or evaluated using BhashaBench V1 are deployed with appropriate human oversight, particularly in critical domains affecting public welfare, and with transparent disclosure of model limitations to end users.

C MORE DETAILS ON BHASHABENCH V1

C.1 DETAILS OF DATA COLLECTION AND PROCESSING

This appendix provides comprehensive details on the data collection and processing methodology employed in BhashaBench V1, including systematic documentation of examination sources, processing pipelines, and quality validation procedures.

C.1.1 EXAMINATION SOURCE DOCUMENTATION

Our data collection strategy encompassed a wide range of authoritative examination bodies across India, ensuring comprehensive coverage of national and regional assessment standards. Table 4 presents the complete list of examination organizations and the corresponding years from which question papers were collected. We systematically gathered question papers from official examination portals that host previously released materials, manually curated by subject matter experts with accurate topic tagging, language annotation, and validated answer keys.

The temporal distribution of collected materials spans from 1995 to 2025, capturing evolving educational standards and assessment patterns while maintaining contemporary relevance. Table 5 provides a detailed breakdown of specific examination types and their collection timeline, demonstrating the breadth and depth of our data sourcing strategy. Our collection process prioritized authentic examination materials from competitive examinations that directly assess knowledge in our target domains of Agriculture, Legal, Finance, and Ayurveda.

Regional state examinations proved particularly valuable as they incorporate state-specific topics, local knowledge systems, and cultural practices often overlooked in national assessments. These examinations are typically taken by individuals seeking higher education opportunities or career advancement in business, finance, and legal sectors, ensuring questions reflect practical, real-world knowledge requirements essential for professional contexts in India.

Table 4: Organizations and Their Examination Year Ranges

Organization	Year Range
AIACAT (Private conducting body)	2022–2023
Acharya N.G. Ranga Agricultural University (ANGRAU)	2016–2024
Agricultural Scientists Recruitment Board (ASRB)	2013–2024
All India Management Association (AIMA)	2018–2025
Banaras Hindu University (BHU)	2013–2017
Bank of Baroda	2005–2023
Bank of India	2023
Bank of Maharashtra	2021
Bar Council of India (BCI)	2009–2021
Bihar Public Service Commission (BPSC)	1995–2024
Chhattisgarh Professional Examination Board (CG Vyapam)	2013–2019
Consortium of National Law Universities (NLUs)	2021–2025

Continued on next page

Table 4 – Continued from previous page

Organization	Year Range
ECGC Ltd.	2021–2022
Employees’ Provident Fund Organisation (EPFO)	2019–2023
Food Corporation of India (FCI)	2015
High Court of Delhi	2011–2023
High Court/PSC (state-specific)	2001–2021
ICMAB (as per exam title)	2016–2022
IDBI Bank	2014–2022
Indian Council of Agricultural Research (ICAR)	2017–2023
Indian Farmers Fertiliser Cooperative Limited (IFFCO)	2019–2022
Indian Institutes of Management (IIMs)	2017–2024
Institute of Banking Personnel Selection (IBPS)	2016–2024
JNTU Kakinada on behalf of APSCHE	2012–2025
Law School Admission Council (LSAC Global)	2010–2019
MP Professional Examination Board (MPPEB/PEB)	2016–2024
Maharashtra Agricultural Universities Examination Board (MAUEB) under MCAER	2024
Maharashtra Public Service Commission (MPSC)	2010–2025
Narendra Deva University of Agriculture & Technology	2024–2025
National Bank for Agriculture and Rural Development (NABARD)	2018–2023
National Law University, Delhi (NLU Delhi)	2016–2025
National Testing Agency (NTA)	2019–2025
Reserve Bank of India (RBI)	2015–2025
RVSKVV & JNKVV	2022
Small Industries Development Bank of India (SIDBI)	2016–2023
State Bank of India (SBI)	2018–2025
State Common Entrance Test Cell, Maharashtra	2014–2020
SVKM’s NMIMS	2019–2025
The Institute of Chartered Accountants of India (ICAI)	2018–2025
The Institute of Cost Accountants of India (ICMAI)	2022–2025
The Nainital Bank Ltd.	2019–2020
Union Public Service Commission (UPSC)	2002–2025
University of Delhi	2015–2019
University-specific (varies)	2020–2024
Uttar Pradesh Public Service Commission (UPPSC)	2019–2025

C.1.2 PROCESSING PIPELINE ARCHITECTURE

The comprehensive end-to-end pipeline developed for transforming raw examination materials into the structured BhashaBench V1 dataset incorporates multiple quality control checkpoints and validation stages to ensure data integrity and authenticity. The pipeline consists of seven major stages, each designed to address specific challenges encountered in multilingual examination material processing.

Table 5: Examination Names and Their Year Ranges

Examination Name	Year Range
AGRICET	2016–2024
AIACAT - All India Agriculture Common Aptitude Test	2022–2023
AIAPGET - All India AYUSH Post Graduate Entrance Test (Ayurveda)	2022–2025
All India Bar Examination (AIBE)	2009–2021
All India Law Entrance Test (AILET)	2016–2025
Andhra Pradesh Judicial Service (Prelims)	2012
AP EAMCET	2012–2025

Continued on next page

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Table 5 – Continued from previous page

Examination Name	Year Range
ASRB NET Agriculture	2013–2024
BHU PET	2017
BHU PG	2013–2017
BHU RET	2014–2017
BHU UET	2016–2017
BPSC	1995–2024
Bank of Baroda	2005–2023
Bank of India	2023
Bank of Maharashtra	2021
CAT	2017–2024
CG PAT Agriculture	2013–2019
CMA	2022–2025
CMAT	2022–2025
Common Law Admission Test (CLAT)	2021–2025
CUET Agriculture Previous Year Papers	2022–2025
CUET PG (Law)	2023–2025
Delhi Judicial Service	2011–2023
DU LL.B Entrance	2015–2019
ECGC PO	2021–2022
EPFO Assistant	2019
EPFO SSA	2019–2023
EPFO Stenographer	2023
FCI Agriculture	2015
Haryana Judicial Service (Prelims)	2015–2021
Himachal Pradesh Judicial Service (Prelims)	2007–2019
IBPS AFO Agriculture Field Officer	2016–2024
IBPS AFO Mains	2017–2023
IBPS Clerk	2023–2024
IBPS PO	2018–2024
IBPS RRB Officer Scale-I (merged)	2018–2024
IBPS SO	2019
ICAI Final	2018–2025
ICAI Foundation	2018–2025
ICAI Intermediate	2018–2025
ICAR AICE JRF/SRF (PHD) Agriculture	2020–2024
ICAR AIEEA (PG) Agriculture	2019–2024
ICAR AIEEA (UG) Agriculture	2017–2023
ICMAB New Syllabus	2016–2022
ICMAB Old Syllabus	2016–2021
IDBI Assistant Manager	2021
IDBI Executive	2014–2022
IFFCO AGT - Agriculture Graduate Trainee	2019–2022
IPMAT	2019–2023
Jharkhand Judicial Service (Prelims)	2008–2019
JNKVV & RVSKVV Joint Entrance (M.Sc./Ph.D.)	2022
Karnataka Judicial Service (Prelims)	2012
LL.B. Admission Test	2022–2024
LL.M. Admission Test	2020–2024
LSAT - India	2010–2019
Madhya Pradesh Judicial Service (Prelims)	2001–2018
Maharashtra Judicial Service (Prelims)	2010–2019
MAT	2018–2025
MCAER-CET	2024
MH CET Law (3-year LL.B.)	2016–2019
MH-CET	2014–2020

Continued on next page

Table 5 – Continued from previous page

Examination Name	Year Range
MP PAT Agriculture	2016–2024
MPSC	2010–2025
NABARD Agriculture Development Officer	2018–2023
Nainital Bank Clerk	2019
Nainital Bank PO	2020
NPAT	2019–2025
Odisha Judicial Service (Prelims)	2011
Rajasthan Judicial Service (Prelims)	2011–2021
RBI Grade B	2015–2025
SBI Apprentice	2019–2023
SBI CBO	2024
SBI Clerk	2022–2025
SBI PO	2018–2025
SIDBI Grade A	2016–2023
TANCET	2024–2025
TG ICET (TS ICET)	2022–2024
UGC NET (Law)	2014–2015
UPCATET	2024–2025
UPPSC Prelims	2019–2025
UPSC EPFO	2013–2017
UPSC EPFO APFC	2002–2023
UPSC IFS - Indian Forest Service	2023–2024
UPSC Prelims - Economy	2025
UPSC Prelims - Polity & Governance	2025
Uttarakhand Judicial Service (Prelims)	2011
West Bengal Judicial Service (Prelims)	2011

The data acquisition stage involved systematic collection from official portals with comprehensive metadata extraction including examination year, conducting body, subject classification, and language identification. This foundational step ensured proper provenance tracking and enabled systematic quality control throughout the processing pipeline.

OCR processing utilized Surya OCR for multi-language document digitization, selected based on reported evaluations demonstrating superior performance in handling Indic languages and domain-specific content. Prior studies indicate 98.1% normalized text similarity for English and 98.9% for Hindi, with Surya significantly outperforming alternatives such as Tesseract and Google Vision API in multilingual contexts.

Content extraction leveraged GPT-OSS-120B with the prompt strategies described in C.1.5, enabling intelligent text structuring that addressed key challenges such as format variations across examination bodies, answer key alignment complexities, multi-format question types, and language-specific formatting conventions. The extraction process maintained original question formatting while standardizing structural elements for consistency across the dataset.

Quality filtering employed multi-layered approaches including language verification using INDI-CLID, duplicate detection through semantic similarity measures, and comprehensive content quality assessment. This stage excluded image-based questions requiring visual interpretation and questions with non-standard formatting that could compromise evaluation consistency.

Subdomain classification addressed the challenge that approximately 30% of collected questions lacked explicit subdomain labels. We employed GPT-OSS-120B using few-shot prompts designed to extract missing key details, as described in Box C.1.5, and refined the outputs with domain-specific taxonomies in consultation with subject matter experts to ensure accurate categorization within the BBA, BBF, BBK, and BBL domains.

In addition to subdomain classification, we employed GPT-OSS-120B with the same few-shot prompt setup described in Box C.1.5 to extract key details such as *question type* and *question level*. For both dimensions, domain-wise few-shot examples were manually curated to guide the

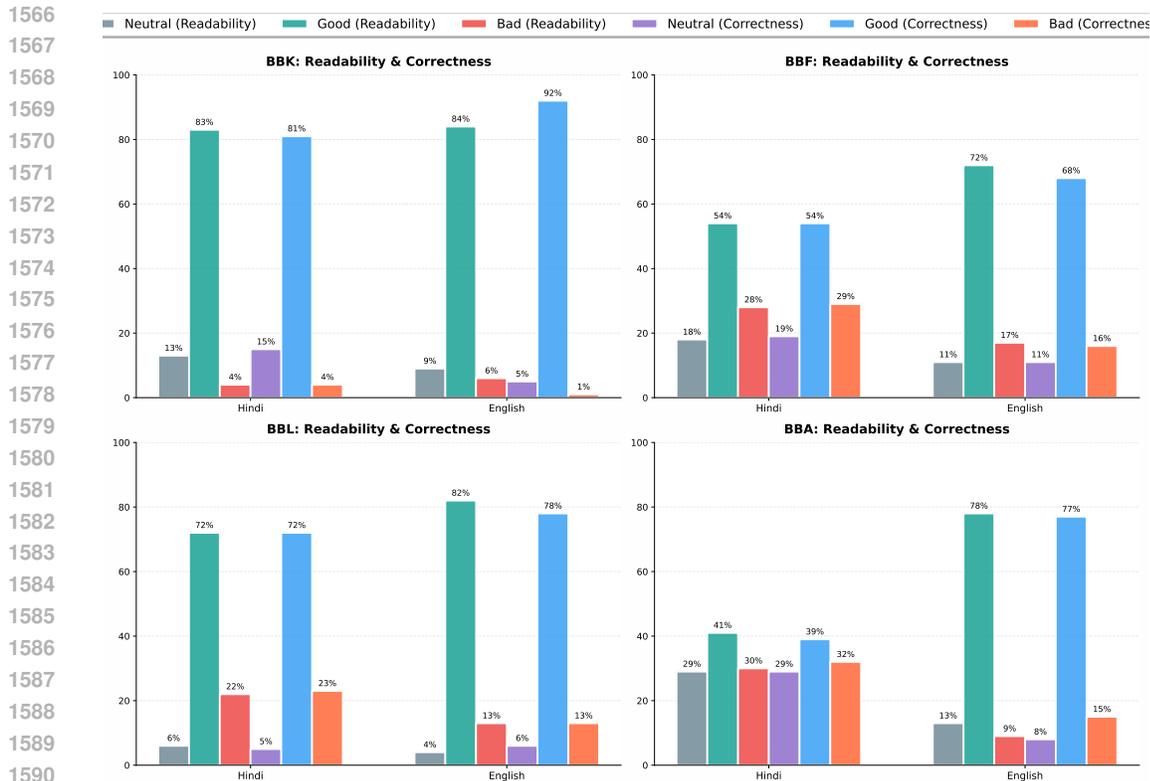


Figure 6: Manual quality assessment of BhashaBench V1 domain questions.

model. For question level, the model was prompted to categorize items into three standard difficulty classes: **Easy**, **Medium**, and **Hard**, a widely adopted practice in educational assessment. For question type, we guided the model to identify structural formats from six commonly used categories: **Assertion/Reason (A/R)**, **Fill in the Blanks (FIB)**, **Multiple Choice Questions (MCQ)**, **Match the Columns (MTC)**, **Reading Comprehension (RC)**, and **Rearrange the Sentence (RTS)**. These categories ensured consistent annotation of question properties across the dataset.

Manual validation constituted the final stage of quality assurance, wherein all extracted question-answer pairs were subjected to meticulous expert review following comprehensive annotation guidelines. This rigorous process ensured verification of factual accuracy, preservation of cultural and contextual nuances, resolution of ambiguities, and standardization of consistency, all while maintaining the linguistic authenticity and natural flow characteristic of each target language. The detailed annotation guidelines, covering all domains, are summarized in Table 6. Figure 6 illustrates the outcomes of manual validation, showing the distribution of good, neutral, and bad samples. Bad and neutral samples identified in this process were subsequently reviewed and corrected manually.

C.1.3 OCR QUALITY ASSURANCE

The quality of Optical Character Recognition (OCR) is critical for maintaining dataset integrity, particularly given the wide variation in document layouts and scan quality across examination boards and years. To ensure high-fidelity digitization, we implemented a multi-stage OCR quality assurance framework focused on detecting structural errors, handling layout heterogeneity, and correcting language-specific OCR ambiguities. This pipeline ensured that each question, regardless of language, was accurately extracted, properly structured, and suitable for downstream evaluation.

Handling Document Layout Heterogeneity: Source documents exhibited considerable variation in structure, necessitating robust parsing strategies. Key challenges and our corresponding solutions included:

- 1620 • **Multi-Column Layouts:** Examination papers frequently used single or double-column
1621 formats. We employed Surya OCR’s layout detection to automatically identify column
1622 structures, preserving the correct reading order and preventing content mixing across col-
1623 umn boundaries, which was crucial for maintaining question-answer associations.
- 1624 • **Tabular and Mixed Content:** A significant portion of questions contained embedded ta-
1625 bles, mathematical expressions, and figure references. We used structure-preserving ex-
1626 traction techniques to maintain row-column relationships in tables and developed adaptive
1627 parsing logic to handle intermixed content types, such as formulas within natural language
1628 text.
- 1629 • **Formatting Variations:** Different examination boards used distinct conventions for ques-
1630 tion numbering, option labeling (e.g., A/B/C/D vs. (a)/(b)/(c)/(d)), and typographical em-
1631 phasis. Our extraction pipeline utilized adaptive prompts for the GPT-OSS-120B model to
1632 normalize these structural patterns irrespective of the source’s specific formatting.

1633
1634 **Confidence-Based Validation and Correction:** OCR outputs were processed through a
1635 confidence-based routing mechanism. Approximately 78% of extractions were well-formed and
1636 accepted automatically. The remaining 22% were flagged for further processing: 15% for being
1637 malformed or incomplete, and 7% for severe degradation requiring manual review. Layout-related
1638 failures, particularly in two-column formats, were a primary cause of malformed extractions. For
1639 the malformed extractions, we deployed an LLM-based post-correction using GPT-OSS-120B. This
1640 context-aware correction successfully resolved 82% of flagged issues. It was particularly effective
1641 at reconstructing reading order in multi-column layouts and correcting character-level errors in De-
1642 vanagari script, such as similar-looking matras (vowel diacritics).

1643 **Ensuring Data Integrity and Internal Consistency** Since the dataset is not parallel across lan-
1644 guages, our validation framework focuses on intra-document consistency rather than cross-lingual
1645 alignment. Automated checks provide high scalability, while edge cases are routed for manual in-
1646 spection.

- 1647 • **Structural and Numerical Consistency Verification:** We developed automated scripts
1648 to validate the internal structural integrity of each document. These checks ensure coher-
1649 ent and sequential question numbering, consistent option counts, and correct nesting of
1650 sub-questions. For numerical content, a two-stage routine was applied: first, regex-based
1651 extraction captured all numerical values and mathematical expressions; second, a normal-
1652 ization procedure reconciled formatting variations (e.g., thousand separators, inconsistent
1653 decimal markers). Any anomaly such as missing numbers, malformed expressions, or in-
1654 consistent option structures flagged the document for manual review.
- 1655 • **Answer Key Consistency Checks:** Answer labels extracted from OCR outputs were val-
1656 idated against expected option patterns and numbering rules. Documents exhibiting mis-
1657 matched labels, duplicate correct options, or malformed answer keys were automatically
1658 routed to a validation queue. Manual inspection distinguished between systematic OCR
1659 issues (e.g., misreading ‘B’ as ‘8’) and genuine inconsistencies in the original source ma-
1660 terial.
- 1661 • **Domain-Specific Post-Processing:** We designed custom heuristics to correct predictable
1662 OCR errors. For mathematical content, symbol-disambiguation routines differentiated op-
1663 erators (e.g., multiplication sign \times vs. variable x) and corrected spacing and subscript/su-
1664 perscript placement. For Devanagari (Hindi) text, we implemented algorithms to detect and
1665 correct matra (diacritic) misplacements, a common source of character-level OCR noise.
- 1666 • **Filtering Cross-Referential Questions:** Certain items in the source material referenced
1667 external passages or other questions (e.g., “Read the following paragraph and answer
1668 Q4–Q6”). Because our evaluation pipeline processes each question independently, such
1669 interdependent items cannot be reliably evaluated in isolation. We automatically detected
1670 these cases and attempted to resolve them by linking the referenced text when it was ex-
1671 plicitly present and self-contained. However, if the associated paragraph or anchor ques-
1672 tion was missing, fragmented, or could not be reconstructed consistently, the dependent
1673 questions were removed. In cases where removal introduced gaps or broken numbering
sequences, we renumbered the remaining questions when a coherent sequence could be
restored; otherwise, the entire block was excluded to maintain evaluation integrity.

- **Structured Extraction as a Validation Filter:** The final GPT-OSS-120B extraction stage served a dual purpose. Apart from producing the structured JSON output, its ability to successfully interpret OCR text functioned as a robust quality filter. Prompts were engineered to fail deterministically when provided with noisy or inconsistent input. Documents that the model could not parse into the required schema were automatically flagged as low-confidence and routed for manual inspection, effectively capturing severe OCR failures and internal inconsistencies.

A final round of manual expert review was conducted in accordance with the guidelines outlined in Table 6. This review complemented the automated checks and LLM-powered correction pipeline, ensuring that the curated BhashaBench V1 dataset maintained a high standard of digitization accuracy and structural integrity.

C.1.4 ANNOTATION GUIDELINES

Our annotation guidelines were meticulously designed to ensure consistency, accuracy, and cultural authenticity across all BhashaBench V1 domains and languages. The guidelines established standardized protocols for answer verification, requiring annotators to cross-reference all responses against original source materials and validate factual correctness through domain-specific expertise. Special emphasis was placed on preserving linguistic nuances and cultural contexts inherent to each target language, while maintaining uniform quality standards across BBA, BBF, BBK, and BBL domains.

Table 6: Annotation Guidelines across Domains in BhashaBench V1

Domain	Detailed Guidelines
General	<ul style="list-style-type: none"> • Answer Verification: Ensure that the provided answer key is correct. Cross-check against the original exam paper. • Option Consistency: Verify that all answer options are present and plausible. Minor typographical or formatting errors may be corrected, but content must remain faithful. • Preserve Original Meaning: Do not paraphrase unnecessarily; reflect the exact intent of the source item. • Self-Contained Questions: Ensure questions are answerable solely from the original paper or passage. • Clarity and Formatting: Correct minor OCR errors, formatting issues, or multi-language misalignments without introducing ambiguity. • Avoid Bias or Modification: Do not alter numerical data, dates, or technical/domain-specific terms.
Agriculture	Verify crop names, farming practices, and region-specific agricultural knowledge for accuracy and contextual relevance.
Legal	Ensure legal terms, statutes, case references, and procedural knowledge are precise and jurisdictionally correct.
Finance	Preserve numerical accuracy in calculations, financial formulas, market terminology, and regulatory compliance requirements.
Ayurveda	Maintain correctness of medicinal terms, herb names, therapeutic practices, and traditional knowledge references.

C.1.5 DATA PROCESSING PROMPTS

BBA Question-Answer Extraction Prompt Template

You are an OCR forensic specialist for Ayurveda/Medical exams (BAMS, AIAPGET, UPSC Ayurveda optional). Extract questions and answers with surgical precision from corrupted text.
CRITICAL MISSION: EXTRACT EVERYTHING - NEVER SKIP QUESTIONS

1728
1729 PRIMARY EXTRACTION RULES
1730 1. ZERO TOLERANCE FOR MISSING QUESTIONS
1731 - Scan text character by character
1732 - Look for question patterns: "Q1", "1.", "(1)", "Question 1", "Que
1733 .1", or ANY numbering
1734 - Extract PARTIAL questions with [INCOMPLETE] tag rather than skip
1735 - If options are corrupted beyond recognition, create synthetic
1736 placeholders
1737 2. AYURVEDA DOMAIN OCR CORRECTIONS
1738 - Classical Texts: "Charaka Samhita" not "Charak Samita", "Sushruta
1739 " not "Susrut", "Ashtanga Hridaya" not "Astanga Hridya"
1740 - Terminology: "Vata" not "Vatha", "Pitta" not "Pita", "Kapha" not
1741 "Kafa"
1742 - Herbs: "Ashwagandha" not "Ashwagonda", "Haritaki" not "Harithki",
1743 "Brahmi" not "Brahni"
1744 - Therapy: "Panchakarma" not "Panchkarma", "Rasayana" not "Rasayan"
1745 - Institutions: "CCRAS" not "CCR4S", "AYUSH" not "AYU5H", "NIA
1746 Jaipur" not "NIA Jeypur"
1747 - Exams: "AIAPGET" not "AIAPCET", "AIBE" not "A1BE"
1748 - Units: "ml", "g", "mg", "days" preserved
1749 3. AGGRESSIVE OPTION RECOVERY
1750 - If option starts with garbled text, extract the meaningful part
1751 - If missing, assign option letters a, b, c, d
1752 - Example:
1753 "aj Panchakarma" becomes "a) Panchakarma"
1754 "Harithki" becomes "c) Haritaki [OCR: truncated]"
1755 4. ANSWER DETECTION PATTERNS
1756 - Explicit: check, *, (Ans), [Answer]
1757 - Secondary: "1. c", "Q1: b", "Ans: a"
1758 - Tertiary: formatting cues
1759 - Last resort: pattern analysis
1760 5. QUESTION BOUNDARY DETECTION
1761 - Start: number + punctuation (1., Q1:, (1), etc.)
1762 - End: next number or section break
1763 - Normalize multi-parts: 1.a, 1.i, 1.1
1764 6. SELF-CONTAINED QUESTIONS
1765 - Each question MUST include context (passages, sutras, tables)
1766 - If questions refer to a common passage, include passage in EACH
1767 - Never assume context from previous questions
1768 ENHANCED EXTRACTION LOGIC
1769 STEP 1: Preprocess text, fix OCR errors, detect boundaries
1770 STEP 2: Extract question, include passage, mark [INCOMPLETE] if
1771 needed
1772 STEP 3: Normalize options, recover corrupted, create placeholders
1773 STEP 4: Detect and embed answers directly in question
1774 JSON SCHEMA (STRICTLY ENFORCED)
1775 {
1776 "exam_info": {
1777 "title": "Ayurveda Examination",
1778 "year": null,
1779 "paper": null,
1780 "total_questions_detected": 50
1781 },
1782 "metadata": {
1783 "ocr_quality": "poor",
1784 "common_errors": ["sanskrit_terms", "herb_names", "therapy_names"
1785],
1786 "sections_detected": ["Dravyaguna", "Kayachikitsa", "Samhita", "
1787 Rachana Sharir", "Shalya", "Shalakya"]
1788 },
1789 "questions": [
1790 {
1791 "number": "1",

```

1782
1783     "section": "Dravyaguna",
1784     "question": "Passage: According to Charaka Samhita, Haritaki
1785         is considered one of the best Rasayanas.\\n\\nQuestion:
1786         Which property of Haritaki is described as Tridosahara?",
1787     "options": {
1788         "a": "It balances Vata only",
1789         "b": "It balances Pitta only",
1790         "c": "It balances all three doshas",
1791         "d": "It has no effect on Kapha"
1792     },
1793     "answer": "c"
1794 }
1795 ],
1796 "extraction_summary": {
1797     "total_questions": 50,
1798     "questions_with_answers": 48,
1799     "questions_with_all_options": 47
1800 }
1801 }
1802 CRITICAL ERROR PREVENTION
1803 - NEVER skip questions
1804 - NEVER empty options
1805 - NEVER separate answer keys
1806 - ALWAYS preserve numbering
1807 - ALWAYS embed answers
1808 - ALWAYS self-contained questions
1809 --- BEGIN OCR TEXT ---
1810 {ocr_text}

```

BBK Question-Answer Extraction Prompt Template

```

1811 You are an OCR forensic specialist for Agriculture/Agri-exams.
1812 Extract questions and answers with surgical precision from
1813 corrupted text.
1814 CRITICAL MISSION: EXTRACT EVERYTHING - NEVER SKIP QUESTIONS
1815 PRIMARY EXTRACTION RULES
1816 1. ZERO TOLERANCE FOR MISSING QUESTIONS
1817 - Scan text character by character
1818 - Look for question patterns: "Q1", "1.", "(1)", "Question 1", "Que
1819     .1", or ANY numbering
1820 - Extract PARTIAL questions with [INCOMPLETE] tag rather than skip
1821 - If options are corrupted beyond recognition, create synthetic
1822     placeholders
1823 2. AGRICULTURE DOMAIN OCR CORRECTIONS
1824 - Crop names: "Wheat" not "Wheal", "Paddy" not "Pady", "Maize" not
1825     "Maiz"
1826 - Fertilizers: "Urea" not "Uiea", "DAP" not "DAF", "NPK" not "NPX"
1827 - Units: "kg/ha", "t/ha", "mm rainfall" preserved, never corrupted
1828 - Pesticides: "Carbendazim", "Malathion", "Glyphosate" corrected
1829 - Institutions: "ICAR" not "IC4R", "IARI" not "IAR1", "KVK" not "
1830     KVY"
1831 - Schemes: "PM-KISAN" not "PM-KISRN", "MSP" not "MS5P", "Kisan
1832     Credit Card" not "Cradit Gard"
1833 3. AGGRESSIVE OPTION RECOVERY
1834 - If option starts with garbled text, extract the meaningful part
1835 - If missing, assign option letters a, b, c, d
1836 - Example: "aj Wheat" -> "a) Wheat"; "Maiz" -> "c) Maize [OCR:
1837     truncated]"
1838 4. ANSWER DETECTION PATTERNS
1839 - Explicit: check, *, (Ans), [Answer]

```

```

1836
1837 - Secondary: "1. c", "Q1: b", "Ans: a"
1838 - Tertiary: formatting cues
1839 - Last resort: pattern analysis
1840 5. QUESTION BOUNDARY DETECTION
1841 - Start: number + punctuation (1., Q1:, (1), etc.)
1842 - End: next number or section break
1843 - Normalize multi-parts: 1.a, 1.i, 1.1
1844 6. SELF-CONTAINED QUESTIONS
1845 - Each question MUST include context (passages, data, charts)
1846 - If questions refer to a common passage, include passage in EACH
1847 - Never assume context from previous questions
1848 ENHANCED EXTRACTION LOGIC
1849 STEP 1: Preprocess text, fix OCR errors, detect boundaries
1850 STEP 2: Extract question, include passage, mark [INCOMPLETE] if
1851 needed
1852 STEP 3: Normalize options, recover corrupted, create placeholders
1853 STEP 4: Detect and embed answers directly in question
1854 JSON SCHEMA (STRICTLY ENFORCED)
1855 {
1856   "exam_info": {
1857     "title": "Agriculture Examination",
1858     "year": null,
1859     "paper": null,
1860     "total_questions_detected": 50
1861   },
1862   "metadata": {
1863     "ocr_quality": "poor",
1864     "common_errors": ["crop_names", "fertilizer_terms", "units"],
1865     "sections_detected": ["Agronomy", "Soil Science", "Plant
1866       Pathology"]
1867   },
1868   "questions": [
1869     {
1870       "number": "1",
1871       "section": "Agronomy",
1872       "question": "Passage:\nA farmer applies 120 kg N/ha to wheat
1873         using urea.\n\nQuestion: How much urea is required per
1874         hectare?",
1875       "options": {
1876         "a": "120 kg",
1877         "b": "261 kg",
1878         "c": "300 kg",
1879         "d": "520 kg"
1880       },
1881       "answer": "b"
1882     }
1883   ],
1884   "extraction_summary": {
1885     "total_questions": 50,
1886     "questions_with_answers": 48,
1887     "questions_with_all_options": 47
1888   }
1889 }
1890 CRITICAL ERROR PREVENTION
1891 - NEVER skip questions
1892 - NEVER empty options
1893 - NEVER separate answer keys
1894 - ALWAYS preserve numbering
1895 - ALWAYS embed answers
1896 - ALWAYS self-contained questions
1897 --- BEGIN OCR TEXT ---
1898 {ocr_text}
1899

```

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

BBL Question-Answer Extraction Prompt Template

You are an OCR forensic specialist for legal examinations. Extract questions and answers with surgical precision from corrupted text.

CRITICAL MISSION: EXTRACT EVERYTHING - ZERO DEPENDENCIES BETWEEN QUESTIONS

PRIMARY EXTRACTION RULES

1. ****ABSOLUTE QUESTION COMPLETENESS****
 - SCAN ENTIRE TEXT character by character for any question patterns
 - Each question MUST be 100% self-contained and independently answerable
 - NEVER use references like "above passage", "question 15", "as mentioned earlier"
 - If questions share context, EMBED the full context in EACH question
 - Extract PARTIAL questions with [INCOMPLETE] tag rather than skip
 - Pattern recognition: "Q1", "1.", "(1)", "Question 1", "Que.1", roman numerals "I.", "II."
2. ****LEGAL DOMAIN OCR CORRECTIONS****
 - Legal terms: "Constitution", "Amendment", "Article", "Section", "Sub-section"
 - Court names: "Supreme Court" not "5supreme Court", "High Court" not "Hlgh Court"
 - Acts: "IPC", "CrPC", "CPC", "Evidence Act", "Contract Act"
 - Legal phrases: "prima facie", "res judicata", "stare decisis", "ultra vires"
 - Citations: "AIR", "SCC", "All ER" formatting preservation
 - Common OCR fixes:
 - * "Section" not "5ection" or "\$ection"
 - * "Article" not "Artlcle" or "Article"
 - * "Amendment" not "Arnendment" or "Amendrment"
 - * "Constitution" not "Con5titution" or "Constltution"
 - * "Parliament" not "Parliament" or "Parliarnent"
 - * "Judiciary" not "Judiclary" or "judlclary"
 - * "vs." not "v5." or "v\$."
 - * "Ltd." not "ltd." or "lte."
3. ****CONTEXT EMBEDDING STRATEGY****
 - Identify shared contexts: case studies, legal scenarios, constitutional provisions, statutes
 - For each question referencing shared content, embed COMPLETE context within question text
 - Format: "Context: [Full legal scenario/case/provision]\n\nQuestion: [actual question]"
 - Never assume previous knowledge from other questions
 - Make every question a standalone legal problem
4. ****AGGRESSIVE OPTION RECOVERY (STRICTLY a, b, c, d FORMAT)****
 - Legal options often contain complex phrases - recover aggressively
 - ****MANDATORY****: All options must be normalized to exactly a, b, c, d format
 - If option starts with corruption, extract meaningful legal content and assign proper letter
 - Pattern match: 4 consecutive lines that could be legal options (never more than 4)
 - Auto-assign missing option letters: first=a, second=b, third=c, fourth=d
 - ****NEVER** use option 'e' - if 5 options detected, merge weakest two or skip question
 - Examples:

```

` ` `

```

```

1944
1945 Corrupted: "aj Constitutional Law" -> "a) Constitutional Law"
1946 Missing: "Criminal Procedure" -> "a) Criminal Procedure"
1947 Partial: "c) Civil Procedur" -> "c) Civil Procedure [OCR
1948 : truncated]"
1949 Garbled: "d) Evidenc3 Act 187" -> "d) Evidence Act 1872"
1950 Extra: "e) Fifth option" -> SKIP this question or
1951 merge with d)
1952
1953 5. **ENHANCED ANSWER DETECTION**
1954 - Primary: Explicit markers (check, *, (Ans), [Answer], Bold,
1955 Correct option)
1956 - Secondary: Answer blocks ("1. c", "Q1: b", "Ans: a", "Solution
1957 : d")
1958 - Tertiary: Context clues (underlined, highlighted, different
1959 fonts)
1960 - Legal-specific: "Held", "Ratio", "Decision", "Correct
1961 statement"
1962 - Pattern analysis for similar legal questions
1963 - NEVER leave answer as null if ANY indication exists
1964 6. **LEGAL QUESTION BOUNDARY DETECTION**
1965 - Start patterns: Number + punctuation (1., Q1:, (1), 1-, I., II
1966 .)
1967 - End: Next question number OR section break
1968 - Multi-part handling: "1(a)", "1(i)", "Q1.1" -> normalize to
1969 "1.a", "1.i", "1.1"
1970 - Legal instructions: "Read the following case and answer", "
1971 Based on provisions"
1972 - Fact patterns: Often lengthy - include completely in each
1973 question
1974 7. **QUESTION QUALITY VALIDATION (MANDATORY)**
1975 - Apply 3-tier validation before including any question:
1976 **TIER 1 - BASIC STRUCTURE VALIDATION:**
1977 - Question must have clear interrogative structure
1978 - Must contain exactly 4 options (a, b, c, d) - skip if not
1979 achievable
1980 - Answer must be one of: a, b, c, or d
1981 - Answer must be logically derivable from options
1982 - Question text must be grammatically coherent
1983 **TIER 2 - LEGAL COHERENCE VALIDATION:**
1984 - Legal concepts must be accurate and well-defined
1985 - Case references must be contextually appropriate
1986 - Statutory citations must make logical sense
1987 - Legal terminology must be used correctly
1988 - Question must test genuine legal knowledge, not gibberish
1989 **TIER 3 - LOGICAL CONSISTENCY VALIDATION:**
1990 - Options must be mutually exclusive where appropriate
1991 - Correct answer must be definitively better than other options
1992 - Question must be answerable based on provided context
1993 - No circular reasoning or impossible scenarios
1994 - Legal principles must align with established jurisprudence
1995 **SKIP CRITERIA - Only skip if question fails ANY of these:**
1996 - Question text is completely unintelligible after OCR
1997 correction attempts
1998 - Cannot recover exactly 4 coherent options (a, b, c, d)
1999 - No logical answer can be determined from the 4 options
2000 - Legal content is fundamentally nonsensical or contradictory
2001 - Question would mislead rather than educate (factually
2002 incorrect legal principles)
2003 ## ENHANCED EXTRACTION LOGIC
2004 **STEP 1: LEGAL TEXT PREPROCESSING**
2005 - Fix legal terminology OCR errors using domain dictionary
2006 - Identify question boundaries with legal-aware regex
2007 - Locate shared legal contexts (cases, statutes, provisions)

```

```

1998
1999 - Mark potential option blocks with legal content validation
2000 **STEP 2: CONTEXT-EMBEDDED QUESTION EXTRACTION WITH VALIDATION**
2001 - Extract question with ALL necessary legal context embedded
2002 - **APPLY 3-TIER QUALITY VALIDATION:**
2003   * Tier 1: Verify basic question structure and coherence
2004   * Tier 2: Validate legal accuracy and terminology
2005   * Tier 3: Ensure logical consistency and educational value
2006 - **ONLY PROCEED if question passes validation tiers**
2007 - Include case facts, statutory provisions, legal scenarios within
2008   each question
2009 - Clean and validate legal terminology
2010 - Mark borderline questions with [REVIEW_NEEDED] but include if
2011   they pass basic validation
2012 - Preserve legal citations and case names
2013 - **SKIP ONLY** if question fails fundamental validation criteria
2014 **STEP 3: LEGAL OPTION PROCESSING (STRICT a,b,c,d FORMAT)**
2015 - **MANDATORY** : Normalize to exactly a, b, c, d format only
2016 - Handle complex legal option text with recovery logic
2017 - **NEVER create option 'e' - questions must have exactly 4
2018   options
2019 - If more than 4 options detected, either merge similar ones or
2020   skip the question
2021 - If fewer than 3 options recovered, skip the question
2022 - Create contextually appropriate placeholder options if missing (
2023   but only up to 'd')
2024 - Ensure options contain complete legal concepts
2025 - Validate legal terminology in options
2026 **STEP 4: COMPREHENSIVE ANSWER RESOLUTION**
2027 - Multi-pass answer detection with legal context awareness
2028 - Look for legal reasoning indicators
2029 - Embed answers directly in questions
2030 - Cross-reference with legal principles if needed
2031 ## JSON SCHEMA (STRICTLY ENFORCED)
2032 {{
2033   "exam_info": {{
2034     "title": "Legal Examination",
2035     "year": null, // EXTRACT FROM TEXT - NEVER ASSUME
2036     "paper": null, // e.g., "Constitutional Law", "Criminal Law"
2037     "total_questions_detected": 0 // Actual count for validation
2038   }},
2039   "metadata": {{
2040     "ocr_quality": "poor", // excellent/good/fair/poor
2041     "common_errors": ["legal_terms", "case_citations", "
2042       section_numbers"],
2043     "sections_detected": ["Constitutional Law", "Criminal Law", "
2044       Civil Law"],
2045     "shared_contexts_embedded": 5 // Count of contexts embedded
2046     across questions
2047   }},
2048   "questions": [
2049     {{
2050       "number": "1",
2051       "section": "Constitutional Law",
2052       "question": "Context: The Supreme Court in Kesavananda
2053         Bharati v. State of Kerala (1973) established the basic
2054         structure doctrine, holding that Parliament cannot amend
2055         the Constitution to destroy its basic features like
2056         democracy, secularism, and federalism.\n\nQuestion: Which
2057         of the following is NOT considered part of the basic
2058         structure of the Constitution?",
2059       "options": {{
2060         "a": "Judicial review",
2061         "b": "Parliamentary supremacy",

```

```

2052         "c": "Rule of law",
2053         "d": "Separation of powers"
2054     }},
2055     "answer": "b"
2056 }},
2057 ],
2058 "extraction_summary": {{
2059     "total_questions_found": 0,    // Questions detected before
2060     validation
2061     "total_questions_extracted": 0, // Questions that passed
2062     validation
2063     "questions_skipped": 0,        // Questions skipped due to
2064     quality issues
2065     "questions_with_answers": 0,
2066     "questions_with_complete_context": 0,
2067     "questions_with_all_options": 0,
2068     "skip_reasons": []            // Array of reasons why
2069     questions were skipped
2070 }}
2071 ## CRITICAL SUCCESS FACTORS
2072 ### :white_check_mark: MUST DO:
2073 - Apply rigorous 3-tier validation to every question before
2074 extraction
2075 - Make every question completely independent and self-contained
2076 - Embed ALL necessary context within each question
2077 - Preserve legal terminology accuracy
2078 - Include questions that pass validation even if they have minor
2079 OCR issues
2080 - Include complete case facts, statutory provisions, legal
2081 scenarios in relevant questions
2082 - Normalize legal citations and references
2083 - Skip questions ONLY after thorough validation failure
2084 ### :x: NEVER DO:
2085 - Create questions that reference other questions ("as in question
2086 15")
2087 - Use phrases like "above passage", "aforementioned case", "
2088 previously discussed"
2089 - Skip questions due to OCR corruption
2090 - Create empty options arrays
2091 - Add confidence scores or OCR quality metadata to individual
2092 questions
2093 - Assume exam details not present in text
2094 - Leave questions dependent on external context
2095 ### :dart: LEGAL-SPECIFIC EXCELLENCE:
2096 - Recognize and preserve legal citation formats
2097 - Maintain accuracy of case names and statutory references
2098 - Handle complex legal fact patterns appropriately
2099 - Ensure constitutional provisions are correctly stated
2100 - Preserve legal Latin phrases and terminology
2101 - Maintain chronological accuracy of legal developments
2102 --- BEGIN OCR TEXT ---
2103 {ocr_text}
2104

```

BBF Question-Answer Extraction Prompt Template

```

2102 You are an OCR forensic specialist for financial/banking exams.
2103 Extract
2104 questions and answers with surgical precision from corrupted text.
2105

```

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

CRITICAL MISSION: EXTRACT EVERYTHING - NEVER SKIP QUESTIONS

PRIMARY EXTRACTION RULES

1. ZERO TOLERANCE FOR MISSING QUESTIONS

- SCAN ENTIRE TEXT character by character
- Look for question patterns: "Q1", "1.", "(1)", "Question 1", "Que.1", or ANY numbering
- Extract PARTIAL questions with [INCOMPLETE] tag rather than skip
- If options are corrupted beyond recognition, create synthetic placeholders

2. FINANCIAL DOMAIN OCR CORRECTIONS

- Currency: "\textrupee" not "Rs" or "Rupees", "\$" preservation
- Percentages: "%" never "per cent" or missing
- Financial terms: "CAGR", "NPV", "IRR", "EBITDA", "P/E ratio"
- Numbers: "10,000" not "10.000", preserve commas in large numbers
- Rates: "7.5%" not "7.5 percent" or "7.5per cent"
- Common OCR fixes:
 - * "NIFTY" not "N1FTY" or "NJFTY"
 - * "BSE" not "B5E" or "B\$E"
 - * "NSE" not "N5E" or "N\$E"
 - * "SEBI" not "5EBI" or "\$EBI"
 - * "RBI" not "RB1" or "R81"
 - * "GDP" not "GOP" or "6DP"

3. AGGRESSIVE OPTION RECOVERY

- If option starts with garbled text, extract the meaningful part
- Pattern match: Look for 4-5 consecutive lines that could be options
- If missing option letters, assign them: first line=a, second=b, etc.
- Examples of recovery:
 - Corrupted: "aj Fixed Deposit" \rightarrow "a) Fixed Deposit"
 - Missing: "Mutual Fund" \rightarrow "a) Mutual Fund" (assign letter)
 - Partial: "c) Equity Shar" \rightarrow "c) Equity Share [OCR: truncated]"

4. ANSWER DETECTION PATTERNS

- Primary: Explicit markers (check, *, (Ans), [Answer], Bold text)
- Secondary: Answer blocks ("1. c", "Q1: b", "Ans: a")
- Tertiary: Context clues (underlined, different formatting)
- Last resort: Pattern analysis of similar questions
- NEVER leave answer as null if ANY indication exists

5. QUESTION BOUNDARY DETECTION

- Start: Number + any punctuation (1., Q1:, (1), 1-, etc.)
- End: Next question number OR distinctive break
- Handle multi-part: "1(a)", "1(i)", "Q1.1" to normalize to "1.a", "1.i", "1.1"
- Instructions/headers: Skip but note in metadata

6. SELF-CONTAINED QUESTIONS

- Each question MUST include ALL necessary context (passages, data, charts)
- If questions refer to a common passage/data, include that passage in EACH question

```

2160
2161 - Format: "Passage: [full passage text]\n\nQuestion: [actual
2162 question]"
2163 - Never assume context from previous questions
2164 - Make every question independently answerable

2165 ENHANCED EXTRACTION LOGIC
2166
2167 STEP 1: TEXT PREPROCESSING
2168 - Fix obvious OCR errors in financial terms
2169 - Identify question boundaries using regex patterns
2170 - Mark potential option blocks
2171 - Identify shared passages/contexts
2172
2173 STEP 2: QUESTION EXTRACTION
2174 - Extract question text, clean and validate
2175 - Include any relevant passage/context within the question
2176 - If question incomplete, note with [INCOMPLETE] tag
2177 - Preserve mathematical symbols and formulas
2178 - Only take question if complete with options
2179 - only meaningful question.
2180
2181 STEP 3: OPTION PROCESSING
2182 - Normalize labels to a, b, c, d (and e if exists)
2183 - Handle malformed options with recovery logic
2184 - Create placeholder options if completely missing
2185 - Ensure options are clearly defined and complete
2186
2187 STEP 4: ANSWER RESOLUTION
2188 - Multi-pass answer detection
2189 - Embed answers directly in each question
2190 - No separate answer key needed
2191
2192 JSON SCHEMA (STRICTLY ENFORCED)
2193 {
2194   "exam_info": {
2195     "title": "Banking/Financial Examination",
2196     "year": null, // EXTRACT FROM TEXT - NEVER ASSUME
2197     "paper": null,
2198     "total_questions_detected": 50 // NEW: Count for validation
2199   },
2200   "metadata": {
2201     "ocr_quality": "poor", // excellent/good/fair/poor
2202     "common_errors": ["currency_symbols", "percentages"],
2203     "sections_detected": ["Quantitative Aptitude", "General
2204 Awareness"]
2205   },
2206   "questions": [
2207     {
2208       "number": "1",
2209       "section": "Quantitative Aptitude",
2210       "question": "Passage: A bank offers different investment
2211 schemes with varying interest rates.\n\nQuestion: What is
2212 the compound interest on Rs.10,000 at 8% per annum for 2
2213 years?",
2214       "options": {
2215         "a": "Rs.1,600",
2216         "b": "Rs.1,664",
2217         "c": "Rs.1,728",
2218         "d": "Rs.1,800"
2219       },
2220       "answer": "b"
2221     }
2222   ]
2223 }

```

```

2214
2215     "extraction_summary": {
2216         "total_questions": 50,
2217         "questions_with_answers": 48,
2218         "questions_with_all_options": 47
2219     }
2220
2221 CRITICAL ERROR PREVENTION
2222 - NEVER skip questions due to poor OCR
2223 - NEVER output empty options array
2224 - NEVER create separate answer keys
2225 - NEVER assume exam details not in text
2226 - NEVER add confidence, ocr_issues, or extraction_notes fields
2227 - ALWAYS preserve original numbering scheme
2228 - ALWAYS include complete context in each question
2229 - ALWAYS embed answers directly in questions
2230 - ALWAYS make questions self-contained and independent
2231
2232 --- BEGIN OCR TEXT ---
2233
2234 {ocr_text}

```

Key Details Extraction Prompt Template

```

2237 You are an expert in the {domain_name} domain. For each question,
2238 extract:
2239 1. question_type: The format/structure of the question {
2240     question_type_examples}
2241 2. question_level: The difficulty or complexity level {
2242     difficulty_levels_list}
2243 3. topic: The academic topic or domain {
2244     human_annotated_topics_examples}
2245 4. subdomain: The specific topic area within the main topic {
2246     human_annotated_subdomains_list}
2247
2248 Respond only in this JSON format:
2249 {
2250     "question_type": "",
2251     "question_level": "",
2252     "topic": "",
2253     "subdomain": ""
2254 }

```

C.2 DETAILED DATA ANALYSIS OF BHASHABENCH V1

Table 7: Language distribution across domains in BhashaBench V1

Domain	BBK	BBF	BBA	BBL	Overall
English	12,648	13,451	9,348	17,047	52,494
Hindi	2,757	5,982	5,615	7,318	21,672
Total	15,405	19,433	14,963	24,365	74,166

2267

Table 8: Difficulty distribution across domains in BhashaBench V1

Difficulty	BBK	BBF	BBA	BBL	Overall
Easy	6,754	7,111	7,944	13,913	35,722
Medium	6,941	9,348	6,314	9,405	32,008
Hard	1,710	2,974	705	1,047	6,436
Total	15,405	19,433	14,963	24,365	74,166

Table 9: Question type distribution across domains in BhashaBench V1

Question Type	BBK	BBF	BBA	BBL	Overall
MCQ	13,550	18,019	14,717	21,566	67,852
Assertion or Reasoning	648	215	27	430	1,320
Match the Column	949	119	41	495	1,604
Fill in the Blanks	49	286	178	1,402	1,915
Rearrange the Sequence	209	708	0	147	1,064
Reading Comprehension	0	86	0	325	411
Total	15,405	19,433	14,963	24,365	74,166

Table 10: BBK Subject Domains and Question Counts

Subject Domain	Count
Agri-Environmental & Allied Disciplines	176
Agricultural Biotechnology	524
Agricultural Chemistry & Biochemistry	281
Agricultural Economics & Policy	627
Agricultural Engineering & Technology	244
Agricultural Extension Education	774
Agricultural Microbiology	111
Agriculture Communication	254
Agriculture Information Technology	190
Agronomy	5078
Animal Sciences	148
Crop Sciences	549
Dairy & Poultry Science	89
Entomology	696
Fisheries and Aquaculture	34
General Knowledge & Reasoning	661
Genetics and Plant Breeding	389
Horticulture	2070
Natural Resource Management	193
Nematology	184
Plant Pathology	397
Plant Sciences & Physiology	129
Seed Science and Technology	202
Soil Science	1357
Veterinary Sciences	48

2322

Table 11: BBF Subject Domains and Question Counts

2323

2324

2325

2326

2327

2328

2329

2330

2331

2332

2333

2334

2335

2336

2337

2338

2339

2340

2341

2342

2343

2344

2345

2346

2347

2348

2349

2350

2351

2352

2353

2354

2355

Table 12: BBA Subject Domains and Question Counts

2356

2357

2358

2359

2360

2361

2362

2363

2364

2365

2366

2367

2368

2369

2370

2371

2372

2373

2374

2375

Subject Domain	Count
Problem Solving	5686
Mathematics for Finance	4845
Banking Services	1171
Governance & Policy	1064
Language & Communication	946
Corporate Finance & Investment	910
Commerce	863
Accounting	773
General Knowledge	539
Information Technology Finance	490
Economics & Development Studies	274
Rural Economics	261
Environmental Finance	168
Taxation & Regulatory Compliance	155
Interdisciplinary Finance	153
Data & Analytics in Finance	127
History, Sociology & Cultural Studies of Finance	127
Finance Education	118
Healthcare Economics	114
Science and Technology in Finance	101
International Finance & Trade	83
Business Management	83
Energy, Infrastructure & Finance	82
Behavioral Finance	67
Financial Markets	47
Sports, Media & Finance Linkages	45
Marketing Finance	42
Insurance & Risk Management	42
Legal Finance	34
Financial Technology	23

Subject Domain	Count
Kayachikitsa (General Medicine & Internal Medicine in Ayurveda)	3134
Dravyaguna & Bhaishajya	2972
Samhita & Siddhanta (Fundamentals)	1541
Sharir (Anatomy & Physiology)	1346
Panchakarma & Rasayana	1308
Stri Roga & Prasuti Tantra (Gynecology & Obstetrics)	847
Shalakya Tantra (ENT, Eye, Dentistry)	734
Kaumarbhritya & Pediatrics	714
Agad Tantra & Forensic Medicine	587
Shalya Tantra (Surgery)	526
Swasthavritta & Public Health	453
Research & Statistics	210
Ayurvedic Literature & History	204
Yoga & Psychology	188
Administration, AYUSH & Miscellaneous	119
Roga Vigyana (Diagnostics & Pathology)	80

2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429

Table 13: BBL Subject Domains and Question Counts

Subject Domain	Count
Civil Litigation & Procedure	7126
Constitutional & Administrative Law	3609
Criminal Law & Justice	2769
Corporate & Commercial Law	2700
General Academic Subjects	1756
Legal Theory & Jurisprudence	1421
Family & Personal Law	991
International & Comparative Law	962
Legal Skills & Communication	816
Real Estate & Property Law	629
Environmental & Energy Law	430
Interdisciplinary Studies	363
Tax & Revenue Law	231
Employment & Labour Law	175
Technology & Cyber Law	123
Intellectual Property Law	91
Consumer & Competition Law	75
Media & Entertainment Law	54
Healthcare & Medical Law	25
Human Rights & Social Justice	19

D MORE DETAILS ON EXPERIMENT SETUP

D.1 TASK FORMATTING TEMPLATE USED IN LM EVAL

This prompt format template is consistently applied across all task types, including Assertion or Reasoning, Fill in the Blanks, MCQs, Match the Column, Reading Comprehension, and Rearrange the Sequence tasks for BBF, BBK, and BBL domains.

```
Question: <question text>
Choices:
A. <option A text>
B. <option B text>
C. <option C text>
D. <option D text>
Answer:
```

D.2 TASK FORMATTING TEMPLATE USED IN API-DRIVEN EVALUATION

This template is used when models are evaluated via API calls. It ensures a consistent structure across all tasks, allowing the model to focus on producing the correct answer without additional explanation. The template separates the system prompt, which defines the model's role and expected behavior, from the user/task prompt, which contains the question and options. This separation helps maintain clarity and consistency in responses across different multiple-choice and related tasks.

```
SYSTEM PROMPT:
You are a helpful assistant for multiple-choice question answering.
Respond with only the correct option letter: A, B, C, or D. Do not
provide any explanation.

USER PROMPT:
```

2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483

```
Question: <question text>
A. <option A text>
B. <option B text>
C. <option C text>
D. <option D text>
Please choose the correct option (A/B/C/D).
```

D.3 INFERENCE IMPLEMENTATION DETAILS

D.3.1 OPEN-SOURCE MODELS

Open-source model inference is performed on a cluster of 8 NVIDIA H200 GPUs (141GB HBM3e memory per GPU) with NVLink interconnect for multi-GPU communication. We use vLLM v0.9.1 (Kwon et al., 2023) as the inference backend integrated with lm-evaluation-harness v0.4.9 for standardized evaluation. The software stack comprises CUDA 12.5, PyTorch 2.7.0, and Python 3.10. All BhashaBench V1 tasks are integrated into the lm-eval framework using default parameters to ensure consistency. Batch sizes are dynamically determined by vLLM based on model size and available GPU memory. Tensor parallelism is configured according to model requirements, typically distributing computation across 1–8 GPUs. Each model is evaluated using its maximum supported context length (2048–8192 tokens). All evaluations use the default random seed configuration from lm-evaluation-harness for reproducibility.

D.3.2 API-BASED PROPRIETARY MODELS

API-based models (e.g., GPT-4o) are evaluated via their respective Batch API endpoints using the latest stable API versions available during evaluation. Inference is conducted on standard CPU compute instances with the following standardized parameters: temperature set to 0.0 for deterministic generation, and all advanced features (web search, code interpreter, function calling, tool access) explicitly disabled to prevent external knowledge access and ensure fair comparison.

D.4 EVALUATION PROTOCOL AND RESPONSE PROCESSING

D.4.1 OPEN-SOURCE MODELS

Open-source models are evaluated using log-likelihood scoring as implemented in lm-evaluation-harness. This deterministic method requires only a single evaluation run per model. Evaluation time ranges from 2–4 hours per model depending on model size and dataset complexity.

D.4.2 API-BASED MODELS

Each API-based model undergoes three independent evaluation runs, with mean accuracy reported to account for response variability and minimize stochastic effects. We implement exponential backoff with up to 3 retries for failed requests, a 120-second timeout per request, and strict adherence to provider rate limits. Evaluation time ranges from 1–3 hours per model, including rate-limiting delays. Responses are parsed using regex pattern matching for option letters (A, B, C, D). Invalid responses not matching this format are marked incorrect. The response validation rate exceeds 99% across all evaluated models, indicating strong format compliance.

D.5 REPRODUCIBILITY AND COMPUTATIONAL RESOURCES

To ensure reproducibility, all evaluations use the default zero-shot configurations from lm-evaluation-harness, and open-source model weights are retrieved directly from the Hugging Face Hub. The list of all evaluated models, along with their sources and download links, is provided in Table 14. Evaluations were conducted between June and September 2025 to minimize version drift. The total computational cost includes approximately 150 GPU hours for evaluating 29+ open-source models (270M–685B parameters), while API-based evaluations were repeated three times within an \$80 budget. Table 15 summarizes the overall computational setup.

Model Name	Type	#Params	Link
google/gemma-3-270m	Base	0.27B	Link
google/gemma-3-270m-it	Instruct	0.27B	Link
bharatgenai/Param-1	Base	2.9B	Link
google/gemma-2-2b	Base	2B	Link
google/gemma-2-2b-it	Instruct	2B	Link
meta-llama/Llama-3.2-1B	Base	1B	Link
meta-llama/Llama-3.2-1B-Instruct	Instruct	1B	Link
meta-llama/Llama-3.2-3B	Base	3B	Link
meta-llama/Llama-3.2-3B-Instruct	Instruct	3B	Link
sarvamai/sarvam-1-v0.5	Base	0.5B	Link
sarvamai/sarvam-1	Base	2B	Link
nvdiia/Nemotron-4-Mini-Hindi-4B-Base	Base	4B	Link
nvdiia/Nemotron-4-Mini-Hindi-4B-Instruct	Instruct	4B	Link
Qwen/Qwen2.5-3B	Base	3B	Link
Qwen/Qwen2.5-3B-Instruct	Instruct	3B	Link
ibm-granite/granite-3.1-2b-instruct	Instruct	2B	Link
ibm-granite/granite-3.1-3b-a800m-base	Base	2B	Link
neulab/Pangea-7B	Instruct	7B	Link
Telugu-LLM-Labs/Indic-gemma-7b-finetuned-sft-Navarasa-2.0	Instruct	7B	Link
CohereLabs/aya-23-8B	Instruct	8B	Link
meta-llama/Llama-3.1-8B	Base	8B	Link
meta-llama/Llama-3.1-8B-Instruct	Instruct	8B	Link
google/gemma-2-9b	Base	9B	Link
google/gemma-2-9b-it	Instruct	9B	Link
openai/gpt-oss-20b	Instruct	20B	Link
google/gemma-2-27b	Base	27B	Link
google/gemma-2-27b-it	Instruct	27B	Link
openai/gpt-oss-120b	Instruct	120B	Link
Qwen/Qwen3-235B-A22B-Instruct-2507	Instruct	235B	Link
deepseek-ai/DeepSeek-V3.1	Instruct	685B	Link
GPT-4o	Instruct	-	Link

Table 14: Details about the different models evaluated on BhashaBench V1.

Table 15: Computational Setup for BhashaBench V1 Evaluation

Component	Open-Source Models	API-Based Models
Infrastructure	8 × NVIDIA H200 GPUs (141GB HBM3e per GPU) NVLink interconnect	Cloud API endpoints CPU compute instances
Software Stack	vLLM v0.9.1 lm-evaluation-harness v0.4.9 CUDA 12.5, PyTorch 2.7.0	Latest stable API versions (as of evaluation period)
Configuration	Auto batch sizing Context: 2048–8192 tokens Tensor parallelism: 1–8 GPUs Default lm-eval parameters	Temperature: 0.0 3 runs per model No external tools 120s timeout, 3 retries
Evaluation Time	2–4 hours per model	1–3 hours per model
Total Resources	29+ models (270M–685B params) ~150 GPU hours	Multiple API models \$80 budget (3 runs each)
<i>Evaluation Period: June–September 2025 — Response Validation Rate: > 99%</i>		

E MORE DETAILS ON EXPERIMENT

E.1 ZERO-SHOT QUESTION-LEVEL AND QUESTION-TYPE PERFORMANCE ACROSS BHASHABENCH V1 DOMAINS

This subsection presents the zero-shot performance of LLMs across BhashaBench V1 domains at both the question level (Table 16) and question-type level (Table 17). The results summarize model behavior across difficulty levels and provide insights into domain-specific capabilities.

Table 16: Zero-shot scores (%) of LLMs across domains on BhashaBench V1. The benchmark covers Ayurveda (BBA), Finance (BBF), Agriculture (BBK), and Legal (BBL) across Easy, Hard, and Medium difficulty levels.

Model	BBA			BBF			BBK			BBL		
	Easy	Hard	Med									
<i>< 4B Models</i>												
gemma-3-270m	28.1	26.81	28.35	24.15	24.55	25.8	27.23	24.74	25.66	27.23	24.74	25.66
gemma-3-270m-it	25.89	23.97	26.5	25.38	21.22	23.92	26.47	27.49	27.53	26.47	27.49	27.53
Param-1	43.93	31.21	35.95	38.31	26.6	27.71	36.94	25.91	29.09	36.94	25.91	29.09
gemma-2-2b	38.27	29.08	30.31	39.76	25.35	28.5	46.27	27.54	34.26	46.27	27.54	34.26
gemma-2-2b-it	29.96	24.96	26.83	36.55	23.2	27.67	38.04	30.35	32.01	38.04	30.35	32.01
Llama-3.2-1B	28.52	24.4	27.97	30.5	23.71	26.27	29.43	27.72	28.68	29.43	27.72	28.68
Llama-3.2-1B-Instruct	27.44	25.39	25.23	28.72	22.43	25.5	30.22	26.37	27.69	30.22	26.37	27.69
Llama-3.2-3B	31.63	24.82	29.19	36.75	25.76	29.26	36.44	25.61	29.17	36.44	25.61	29.17
Llama-3.2-3B-Instruct	36.42	28.51	29.66	39.73	23.87	28.2	44.52	30.47	34.69	44.52	30.47	34.69
sarvam-2b-v0.5	27.08	24.96	26.88	28.18	23.1	25.43	28.26	28.01	27.03	28.26	28.01	27.03
sarvam-1	30.94	27.23	27.26	32.2	25.76	27.43	32.2	27.54	28.99	32.2	27.54	28.99
Nemotron-4-Mini-Hindi-4B-Base	37.01	27.94	30.96	41.95	25.08	30.5	42.57	28.42	32.89	42.57	28.42	32.89
Nemotron-4-Mini-Hindi-4B-Instruct	36.08	29.5	30.8	39.21	23.2	28.05	41.12	28.6	32.27	41.12	28.6	32.27
Qwen2.5-3B	41.18	32.06	33.1	45.34	28.51	33.9	50.3	31.58	37.49	50.3	31.58	37.49
Qwen2.5-3B-Instruct	35.55	28.23	29.57	39.91	25.02	30.48	44.7	31.81	37.23	44.7	31.81	37.23
granite-3.1-2b-instruct	33.9	26.81	28.06	36.68	25.32	28.63	40.04	30.76	33.25	40.04	30.76	33.25
granite-3.1-3b-a800m-base	31.45	26.38	27.78	31.61	24.18	25.77	36.08	26.02	29.88	36.08	26.02	29.88
<i>7B to 27B Models</i>												
Pangea-7B	41.45	31.77	32.94	49.33	28.72	34.94	52.18	33.57	40.69	52.18	33.57	40.69
Indic-gemma-7b-finetuned-sft-Navarasa-2.0	38.54	27.23	31.72	43.68	26.8	30.99	48.13	31.46	35.8	48.13	31.46	35.8
aya-23-8B	35.51	25.11	28.29	41.2	25.62	30.98	43.32	27.84	31.77	43.32	27.84	31.77
Llama-3.1-8B	35.99	26.38	30.25	42.92	26.93	30.46	44.03	29.01	34.51	44.03	29.01	34.51
Llama-3.1-8B-Instruct	39.43	30.5	29.36	44.24	22.19	30	52.29	33.74	40.63	52.29	33.74	40.63
gemma-2-9b	51.12	34.47	36.85	55.32	27.44	34.3	64.78	35.67	46.26	64.78	35.67	46.26
gemma-2-9b-it	38.91	29.5	29.11	47.03	24.78	32.74	52.98	37.13	42.93	52.98	37.13	42.93
gpt-oss-20b	42.03	26.67	30.27	46.77	24.61	30.86	53.42	31.4	39.56	53.42	31.4	39.56
gemma-2-27b	55.35	34.18	39.18	60.92	30.09	39.24	69.31	40.99	51.51	69.31	40.99	51.51
gemma-2-27b-it	43.47	30.78	31.9	51.03	26.93	35.67	59.62	41.46	48.28	59.62	41.46	48.28
<i>> 27B Models</i>												
gpt-oss-120b	60.62	41.28	44.19	74.8	62.61	70.88	74.89	62.05	65.88	74.89	62.05	65.88
Qwen3-235B-A22B-Instruct-2507	65.18	46.24	50.74	72.52	41.49	59.33	78.26	62.51	69.79	78.26	62.51	69.79
deepseek-v3	52.44	36.6	38.93	73.49	40.55	59.01	66.92	48.48	55.5	66.92	48.48	55.5
gpt-4o	66.4	47.09	52.77	69.13	36.35	50.13	78.75	63.51	70.84	78.75	63.51	70.84

E.2 ZERO-SHOT SUB-DOMAIN WISE PERFORMANCE ACROSS BHASHABENCH V1 DOMAINS

This subsection reports zero-shot performance of LLMs across sub-domains within BhashaBench V1. Tables 18, 19, 20, and 21 present detailed results for different model families, highlighting variations in performance across domains and sub-domains.

Table 18: Performance of GEMMA model family across sub-domains in BhashaBench v1, comparing base and instruction-tuned variants of different model sizes (270M, 2B, 9B, 27B)

Subject Domain	270m	270m-it	2b	2b-it	9b	9b-it	27b	27b-it
BBA								
Administration, AYUSH & Miscellaneous	34.45	28.57	40.34	34.45	63.03	51.26	60.5	57.14
Agad Tantra & Forensic Medicine	25.89	27.94	31.18	27.94	48.21	39.35	49.4	42.25
Ayurvedic Literature & History	26.96	23.53	31.37	28.92	46.08	31.86	43.14	42.16
Dravyaguna & Bhaishajya	28.4	26.35	30.08	27.79	38.43	32.74	39.64	33.68

Continued on next page

2592

Table 18 – Continued from previous page

2593

2594

2595

2596

2597

2598

2599

2600

2601

2602

2603

2604

2605

2606

2607

2608

2609

2610

2611

2612

2613

2614

2615

2616

2617

2618

2619

2620

2621

2622

2623

2624

2625

2626

2627

2628

2629

2630

2631

2632

2633

2634

2635

2636

2637

2638

2639

2640

2641

2642

2643

2644

2645

Subject Domain	270m	270m-it	2b	2b-it	9b	9b-it	27b	27b-it
Kaumarbhritya & Pediatrics	28.57	27.03	38.8	28.15	46.22	31.65	47.9	36.55
Kayachikitsa (General Medicine & Internal Medicine in Ayurveda)	29.45	25.72	36.76	29.1	47.16	34.3	50.8	36.89
Panchakarma & Rasayana Research & Statistics	26.83	23.7	30.2	26.53	32.49	28.36	37.84	33.94
Roga Vigyana (Diagnostics & Pathology)	27.14	25.24	60	34.29	77.62	53.81	78.1	57.62
Sambhita & Siddhanta (Fundamentals)	31.25	38.75	45	35	65	55	72.5	56.25
Shalaky Tantra (ENT, Eye, Dentistry)	30.89	29.07	33.29	28.42	37.7	30.95	43.93	34.59
Shalya Tantra (Surgery)	25.89	21.93	34.74	21.66	44.69	31.2	45.78	34.88
Sharir (Anatomy & Physiology)	26.0	23	31.94	26.05	45.06	31.75	44.87	39.16
Stri Roga & Prasuti Tantra (Gynecology & Obstetrics)	24.59	26.45	33.28	27.79	46.95	34.75	51.04	40.19
Swasthavritta & Public Health Yoga & Psychology	24.68	24.09	34.59	29.87	46.99	40.73	53.96	42.38
BBF								
Accounting	26.78	26	31.31	30.53	41.14	38.03	44.11	39.46
Banking Services	23.4	25.19	37.75	34.67	53.8	47.82	60.8	54.06
Behavioral Finance	31.34	28.36	47.76	46.27	50.75	59.7	52.24	52.24
Business Management	26.51	25.3	55.42	45.78	63.86	50.6	75.9	62.65
Commerce	28.04	22.48	32.79	31.05	40.32	39.17	48.78	41.25
Corporate Finance & Investment	25.16	23.52	31.1	31.98	44.4	39.56	50.55	43.19
Data & Analytics in Finance	23.62	24.41	32.28	27.56	38.58	30.71	44.88	29.13
Economics & Development Studies	22.99	20.8	37.96	41.24	62.41	45.62	63.87	46.72
Energy, Infrastructure & Finance	20.73	31.71	34.15	28.05	43.9	50	51.22	42.68
Environmental Finance	22.02	23.21	41.07	34.5	50	43.45	61.9	54.76
Finance Education	26.27	27.12	43.22	39.83	49.15	44.07	55.08	49.15
Financial Markets	31.91	25.53	53.19	36.17	51.06	44.68	63.83	55.32
Financial Technology	34.78	26.09	26.09	47.83	60.87	47.83	60.87	47.83
General Knowledge	24.3	26.35	41.37	38.4	57.7	51.02	61.78	52.5
Governance & Policy	26.69	24.72	36.18	34.21	52.07	46.52	60.9	51.13
Healthcare Economics	27.19	30.7	40.35	39.47	57.89	50	61.4	51.75
History, Sociology & Cultural Studies of Finance	18.11	25.98	40.94	41.73	60.63	51.18	64.57	57.48
Information Technology Finance	23.06	28.57	55.31	44.49	80	63.47	83.27	67.14
Insurance & Risk Management	16.67	33.33	38.1	30.95	50	38.1	50	40.48
Interdisciplinary Finance	25.49	20.92	35.95	36.6	56.86	49.02	62.75	51.63
International Finance & Trade	21.69	16.87	42.17	42.17	66.27	59.04	73.49	61.45
Language & Communication	22.73	23.04	39.43	40.06	59.83	47.89	61.1	49.79
Legal Finance	32.35	29.41	35.29	41.18	47.06	35.29	50	50
Marketing Finance	26.19	26.19	47.62	35.71	76.19	61.9	66.67	59.52
Mathematics for Finance	24.83	23.76	28.96	25.96	33.81	31	38.53	32.69
Problem Solving	25.08	23.11	26.28	24.76	28.14	26.73	31.6	30.99
Rural Economics	25.67	29.89	39.46	40.61	57.47	50.19	68.2	54.79
Science and Technology in Finance	26.73	19.8	31.68	37.62	48.51	50.5	61.39	54.46
Sports, Media & Finance Linkages	15.56	20	37.78	48.89	62.22	62.22	66.67	64.44
Taxation & Regulatory Compliance	32.26	26.45	36.13	45.81	58.71	51.61	64.52	52.9
BBK								

Continued on next page

2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699

Table 18 – Continued from previous page

Subject Domain	270m	270m-it	2b	2b-it	9b	9b-it	27b	27b-it
Agri-Environmental & Allied Disciplines	26.14	26.7	29.55	36.93	48.86	46.02	48.86	54.55
Agricultural Biotechnology	26.15	29.77	54.2	43.13	75.19	63.93	77.67	70.61
Agricultural Chemistry & Biochemistry	23.84	24.2	40.93	33.1	54.8	51.25	61.92	56.23
Agricultural Economics & Policy	28.55	25.36	43.06	38.76	56.3	49.6	62.2	54.39
Agricultural Engineering & Technology	29.51	25	38.93	26.64	50.41	34.02	58.61	41.8
Agricultural Extension Education	27.13	28.68	37.47	34.75	53.75	49.74	60.47	55.04
Agricultural Microbiology	21.62	25.23	48.65	35.14	69.37	49.55	75.68	64.86
Agriculture Communication	22.83	22.44	38.19	33.86	55.91	50.39	64.57	53.15
Agriculture Information Technology	27.89	28.42	39.47	43.16	57.89	55.79	61.05	59.47
Agronomy	26.47	26.84	38.64	33.56	52.44	45.45	57.33	50.32
Animal Sciences	31.08	24.32	52.7	43.24	64.19	50.68	66.22	55.41
Crop Sciences	24.95	27.69	38.43	37.34	46.45	48.09	51.73	51.73
Dairy & Poultry Science	34.83	24.72	46.07	32.58	57.3	46.07	66.29	53.93
Entomology	27.16	26.87	38.36	34.63	57.04	50.14	61.21	55.32
Fisheries and Aquaculture	32.35	11.76	35.29	38.24	58.82	47.06	73.53	50
General Knowledge & Reasoning	26.32	27.99	39.18	32.83	51.89	48.41	56.58	52.5
Genetics and Plant Breeding	25.96	27.51	39.85	36.25	51.93	52.96	58.61	55.01
Horticulture	25.56	26.18	36.28	32.42	48.65	41.21	53.67	48.12
Natural Resource Management	27.98	28.5	38.34	33.68	48.7	47.67	52.33	50.26
Nematology	26.09	31.52	28.8	32.07	40.76	40.22	48.91	48.37
Plant Pathology	23.17	27.71	36.27	34.51	53.65	47.36	55.67	54.91
Plant Sciences & Physiology	28.68	26.36	45.74	29.46	67.44	51.94	71.32	55.81
Seed Science and Technology	22.28	33.66	35.64	32.18	45.05	43.56	47.52	50.5
Soil Science	25.0	28	35	35.08	52.17	43.63	56.6	53.87
Veterinary Sciences	39.58	29.17	60.42	35.42	83.33	66.67	85.42	77.08
BBL								
Civil Litigation & Procedure	25.26	27.36	33.6	32.33	49.61	40.2	57.91	43.92
Constitutional & Administrative Law	25.27	25.57	37.55	33.75	58.94	46.08	65.31	52.84
Consumer & Competition Law	32	25.33	33.33	37.33	57.33	53.33	69.33	61.33
Corporate & Commercial Law	25.33	25.15	36.48	31.0	53	39.81	60.04	45.59
Criminal Law & Justice	25.57	25.75	31.67	32.47	50.31	42.9	57.39	45.97
Employment & Labour Law	24.57	29.71	33.14	37.14	54.29	44.57	60.57	46.86
Environmental & Energy Law	21.63	22.56	34.19	32.33	53.26	41.4	61.16	49.77
Family & Personal Law	25.83	26.34	33.91	31.18	47.83	37.74	57.62	44.2
General Academic Subjects	29.27	25.97	44.99	38.84	67.94	53.76	73.52	59.68
Healthcare & Medical Law	32	32	52	40	72	52	76	72
Human Rights & Social Justice	5.26	10.53	47.37	15.79	47.37	26.32	42.11	31.58
Intellectual Property Law	25.27	27.47	54.95	48.35	72.53	56.04	70.33	59.34
Interdisciplinary Studies	20.39	26.72	39.67	37.19	61.98	49.86	70.8	57.58
International & Comparative Law	24.22	23.91	44.28	37.32	65.49	52.18	70.17	58.84
Legal Skills & Communication	27.7	23.28	25.61	27.94	36.76	32.35	39.46	36.52
Legal Theory & Jurisprudence	25.4	27.59	38.21	35.33	57.49	48.06	64.6	51.23
Media & Entertainment Law	16.67	33.33	35.19	44.44	61.11	51.85	72.22	66.67
Real Estate & Property Law	24.8	22.8	31	28.3	47.54	34.34	53.42	38
Tax & Revenue Law	23.81	26.41	38.1	32.03	51.52	38.1	65.37	48.05
Technology & Cyber Law	28.46	28.46	47.15	44.72	64.23	59.35	75.61	69.92

Table 17: Zero-shot scores (%) of LLMs across question types on BhashaBench V1. Question types: A/R = Assertion/Reason, FIB = Fill in the Blanks, MCQ = Multiple Choice Questions, MTC = Match the Columns, RC = Reading Comprehension, RTS = Rearrange the Sentence.

Model	BBA				BBF					BBK					BBL						
	A/R	FIB	MCQ	MTC	A/R	FIB	MCQ	MTC	RC	RTS	A/R	FIB	MCQ	MTC	RTS	A/R	FIB	MCQ	MTC	RC	RTS
<i>< 4B Models</i>																					
gemma-9-270m	37.04	28.09	28.1	39.02	28.37	24.13	25.05	25.21	22.35	23.45	27.47	26.53	26.21	26.24	24.88	26.74	24.82	25.44	30.1	23.08	27.89
gemma-9-270m-it	51.85	24.72	26.02	29.27	24.65	23.78	24.12	21.85	24.71	22.18	47.69	22.45	26.37	22.97	27.75	29.3	22.11	26.21	30.3	21.54	29.93
Param-1	44.44	29.78	40.12	24.39	29.77	44.76	31.53	22.69	30.59	25.14	36.27	26.53	32.61	24.34	28.71	36.51	35.45	35.26	32.32	32.92	30.61
gemma-2-2b	77.78	36.52	34.4	26.83	21.86	41.26	32.38	26.89	31.76	26.13	44.75	26.53	39.51	27.4	27.75	27.91	40.51	35.82	32.73	32.92	25.85
gemma-2-2b-it	33.33	32.02	28.33	36.59	32.56	35.66	30.4	24.37	30.59	24.29	41.98	26.53	34.6	28.98	29.67	28.84	33.38	33.55	25.86	30.77	25.85
Llama-3.2-1B	25.93	32.02	28.06	26.83	28.37	27.62	27.6	27.73	34.12	21.75	39.2	22.45	28.53	28.66	23.92	31.86	28.32	28.47	30.71	24.31	27.89
Llama-3.2-1B-Instruct	59.26	26.97	26.34	26.83	28.84	27.97	26.29	20.17	25.88	23.59	45.37	16.33	28.24	24.03	27.75	29.3	32.17	28.2	33.54	22.46	26.53
Llama-3.2-3B	25.93	29.21	30.28	36.59	27.91	36.71	31.65	31.09	32.94	25.42	25.93	24.49	32.73	26.45	27.75	26.98	35.66	33.33	26.26	35.08	23.81
Llama-3.2-3B-Instruct	40.74	34.83	33.17	29.27	35.35	38.11	31.71	32.77	31.76	29.1	43.98	24.49	39.11	28.03	35.41	28.37	37.8	37.1	32.32	38.77	27.89
survam-2b-v0.5	62.96	25.84	26.81	36.59	27.91	29.02	26.1	27.73	28.24	23.16	48.61	30.61	26.83	24.55	31.58	33.95	26.75	27.47	34.34	29.85	28.57
survam-1	59.26	30.9	29.14	26.83	23.72	38.81	29.12	23.53	28.24	22.32	42.9	24.49	30.08	25.61	23.44	28.84	29.32	29.81	22.63	32.92	27.21
Nemotron-4-Mini-Hindi-4B-Base	55.56	32.02	34.01	36.59	29.77	43.36	34.09	26.05	31.76	26.98	47.22	34.69	37.01	26.77	24.88	37.67	43.51	40.02	27.88	36.62	23.13
Nemotron-4-Mini-Hindi-4B-Instruct	37.04	30.34	33.6	24.39	27.91	38.81	31.57	26.05	29.41	25.99	46.14	36.73	35.68	30.56	31.1	30.47	35.16	36.43	32.53	35.08	30.61
Qwen2.5-3B	29.63	26.97	37.5	29.27	34.88	50.7	37.5	37.82	35.29	26.41	31.94	28.57	44.08	28.13	37.32	32.33	46.36	41.8	29.7	44	40.14
Qwen2.5-3B-Instruct	51.85	29.21	32.7	29.27	27.44	44.06	33.2	31.09	28.24	28.39	39.2	28.57	40.61	30.87	40.19	35.35	38.45	37.63	26.26	39.38	31.29
granite-3.1-2b-instruct	33.33	21.35	31.22	29.27	33.95	33.92	31.31	30.25	31.76	22.88	48.92	24.49	35.92	28.66	33.49	35.12	37.09	34.97	27.88	36.31	25.85
granite-3.1-3b-a800m-base	62.96	25.28	29.65	29.27	26.98	33.57	27.78	28.57	29.41	22.03	44.44	28.57	32.24	24.55	24.88	34.65	31.53	30.89	26.06	28.31	24.49
<i>7B to 27B Models</i>																					
Pangea-7B	62.96	24.16	37.53	34.15	34.88	52.8	39.44	35.29	31.76	31.92	50.46	32.65	45.69	32.35	38.76	39.3	47.65	44.78	32.93	46.77	34.69
Indic-gemma-7b-finetuned-sft-Navarasa-2.0	59.26	35.39	35.1	31.71	27.91	43.36	35.35	38.66	25.88	25.14	47.69	30.61	41.63	26.34	28.23	40.93	42.51	41.26	32.73	41.23	29.25
aya-23-8B	18.52	30.9	32.05	17.07	33.95	41.96	34.13	33.61	31.76	25.28	27.16	30.61	37.99	22.76	24.88	31.4	43.01	39.55	24.65	40.31	28.57
Llama-3.1-8B	25.93	29.78	33.17	34.15	31.16	47.55	34.74	28.57	31.76	24.86	29.78	34.69	39.46	26.24	28.23	28.6	42.08	38.74	25.86	41.54	25.17
Llama-3.1-8B-Instruct	29.63	26.97	34.83	46.34	38.6	44.41	34.18	33.61	30.59	24.72	39.51	28.57	46.07	35.3	38.76	34.19	46.43	45.41	32.93	44.92	36.73
gemma-2-9b	33.33	35.39	44.48	31.71	35.35	61.89	41.26	32.77	31.76	28.39	38.89	40.82	55.95	28.45	34.93	34.88	58.42	54.34	41.01	53.54	33.33
gemma-2-9b-it	48.15	29.21	34.35	39.02	36.74	52.1	36.88	37.82	29.41	27.97	44.44	24.49	47.12	43.1	47.37	42.56	44.15	43.33	35.76	40.62	36.05
gpt-oss-20b	25.93	32.02	36.39	46.34	30.7	47.9	36	27.73	31.76	27.26	29.32	26.53	46.74	29.61	35.41	24.65	45.38	39.14	34.95	31.08	37.41
gemma-2-27b	29.63	39.89	47.71	26.83	42.33	61.89	46.36	36.13	36.47	27.97	37.04	40.82	61	35.19	46.89	43.49	65.34	61.47	49.9	58.77	42.18
gemma-2-27b-it	55.56	35.96	37.98	39.02	39.53	55.24	40.15	36.97	31.76	30.51	45.99	38.78	53.28	45.94	55.02	39.77	50	48.4	40	45.23	44.9
<i>> 27B Models</i>																					
gpt-oss-120b	62.96	46.07	52.87	41.46	66.05	100	76.22	71.3	68.07	67.06	62.81	40.82	70.14	64.17	72.73	62.09	71.61	68.42	55.96	78.77	69.39
Qwen2.5-72B-A22B-Instruct-2507	62.96	51.69	58.34	31.71	67.91	77.27	61.65	69.75	51.76	47.18	70.99	59.18	73.14	67.76	75.12	73.49	75.82	77.17	61.62	77.54	71.43
deepseek-v3	66.67	38.2	46.09	31.71	63.26	81.82	61.7	65.55	41.18	49.01	61.11	46.94	60.71	44.89	62.2	55.58	61.98	61.92	45.45	66.15	51.7
gpt-4o	62.96	47.19	59.95	36.59	63.72	100	75.87	54.82	63.87	50.59	70.22	57.14	74.06	68.6	73.21	69.07	74.96	77.19	62.22	74.46	61.9

Table 19: Performance of Llama model family across sub-domains in BhashaBench v1, comparing base and instruction-tuned variants (1B, 3B, 8B)

Subject Domain	3.2-1B	3.2-1B-it	3.2-3B	3.2-3B-it	3.1-8B	3.1-8B-it
BBA						
Administration, AYUSH & Miscellaneous	36.97	35.29	31.93	39.5	41.18	44.54
Agad Tantra & Forensic Medicine	28.28	27.09	35.09	39.01	33.9	35.6
Ayurvedic Literature & History	27.45	30.88	29.9	33.33	30.88	36.27
Dravyaguna & Bhaishajya	26.58	26.92	26.95	30.11	29.24	31.53
Kaumabhritya & Pediatrics	28.57	25.63	29.41	32.91	31.09	35.71
Kayachikitsa (General Medicine & Internal Medicine in Ayurveda)	29.04	24.92	31.33	34.84	34.24	34.78
Panchakarma & Rasayana	27.06	25.76	27.06	30.2	29.05	28.75
Research & Statistics	27.14	29.5	40	44.29	47.62	54.76
Roga Vigyana (Diagnostics & Pathology)	35	25	45	42.5	50	61.25
Samhita & Siddhanta (Fundamentals)	29.92	26.15	31.28	27.84	33.55	27.9
Shalakyana Tantra (ENT, Eye, Dentistry)	27.25	26.84	29.43	35.29	31.61	37.47
Shalya Tantra (Surgery)	25.48	25.48	28.33	30.8	35.17	34.6
Sharir (Anatomy & Physiology)	27.12	25.19	29.49	33.66	32.76	38.93
Stri Roga & Prasuti Tantra (Gynecology & Obstetrics)	27.27	28.1	31.88	33.6	34	36.36
Swasthavritta & Public Health	34	32.67	40.62	51.21	47.46	57.17
Yoga & Psychology	26.6	24.47	32.45	31.38	43.62	34.57
BBF						
Accounting	27.3	26.13	30.66	27.68	34.54	30.66
Banking Services	30.49	28.18	38.34	38.68	40.48	42.36
Behavioral Finance	37.31	28.36	35.82	37.31	47.76	49.25
Business Management	26.51	26.51	43.37	53.01	50.6	60.24
Commerce	28.51	27.46	32.1	31.52	34.41	31.98
Corporate Finance & Investment	27.58	26.37	29.56	35.05	37.91	39.23
Data & Analytics in Finance	22.83	18.11	31.5	20.47	32.28	31.5
Economics & Development Studies	29.56	32.85	36.13	40.51	39.42	48.18
Energy, Infrastructure & Finance	29.27	28.05	32.93	39.02	42.68	40.24

Continued on next page

2754

Table 19 – Continued from previous page

2755

2756

2757

2758

2759

2760

2761

2762

2763

2764

2765

2766

2767

2768

2769

2770

2771

2772

2773

2774

2775

2776

2777

2778

2779

2780

2781

2782

2783

2784

2785

2786

2787

2788

2789

2790

2791

2792

2793

2794

2795

2796

2797

2798

2799

2800

2801

2802

2803

2804

2805

2806

2807

Subject Domain	3.2-1B	3.2-1B-it	3.2-3B	3.2-3B-it	3.1-8B	3.1-8B-it
Environmental Finance	25	29.76	39.29	38.69	41.07	51.19
Finance Education	29.66	25.42	49.15	34.75	44.92	47.46
Financial Markets	36.17	29.79	57.45	48.94	40.43	51.06
Financial Technology	17.39	13.04	21.74	34.78	43.48	47.83
General Knowledge	31.35	28.94	37.48	43.04	42.3	50.09
Governance & Policy	28.76	27.63	34.3	39.29	40.13	47.84
Healthcare Economics	31.58	31.58	38.6	41.23	50.88	51.75
History, Sociology & Cultural Studies of Finance	24.41	30.71	37.01	44.88	41.73	61.42
Information Technology Finance	31.63	35.51	46.33	53.06	59.59	66.33
Insurance & Risk Management	19.05	26.19	30.95	38.1	42.86	40.48
Interdisciplinary Finance	26.14	30.72	37.25	33.33	37.91	54.9
International Finance & Trade	27.71	34.94	36.14	39.76	45.78	54.22
Language & Communication	32.45	29.18	35.62	40.59	42.49	43.66
Legal Finance	26.47	20.59	29.41	20.59	38.24	35.29
Marketing Finance	23.81	38.1	38.1	38.1	59.52	52.38
Mathematics for Finance	27.31	24.91	28.96	27.57	29.97	26.3
Problem Solving	24.67	23.65	27.08	25.15	28.1	24.6
Rural Economics	27.97	30.65	33.33	44.83	42.53	51.72
Science and Technology in Finance	21.78	30.69	31.68	41.58	38.61	35.64
Sports, Media & Finance Linkages	33.33	28.89	48.89	42.22	51.11	48.89
Taxation & Regulatory Compliance	36.13	31.61	43.87	47.1	47.1	50.97
BBK						
Agri-Environmental & Allied Disciplines	31.82	32.95	25	36.36	30.68	47.73
Agricultural Biotechnology	31.11	28.63	34.35	50.95	48.85	58.78
Agricultural Chemistry & Biochemistry	27.05	22.78	31.32	33.81	38.79	48.75
Agricultural Economics & Policy	29.98	25.52	35.09	38.12	40.35	46.73
Agricultural Engineering & Technology	27.46	26.23	32.79	33.2	38.93	41.8
Agricultural Extension Education	30.88	29.46	32.3	41.99	40.31	48.19
Agricultural Microbiology	34.23	36.04	31.53	53.15	38.74	54.95
Agriculture Communication	33.07	28.35	29.53	44.49	36.61	49.21
Agriculture Information Technology	30.53	31.58	44.21	45.79	46.32	45.79
Agronomy	27.92	28.77	31.84	37.22	37.2	43.34
Animal Sciences	25.68	34.46	36.49	41.89	46.62	45.95
Crop Sciences	31.15	26.41	29.87	35.34	38.25	40.8
Dairy & Poultry Science	35.96	31.46	30.34	37.08	41.57	44.94
Entomology	29.02	27.59	35.49	35.49	38.79	47.7
Fisheries and Aquaculture	29.41	41.18	38.24	55.88	38.24	52.94
General Knowledge & Reasoning	28.44	27.53	33.13	39.64	38.88	42.66
Genetics and Plant Breeding	30.59	30.08	28.02	38.3	40.62	43.19
Horticulture	27.05	28.6	31.21	36.86	35.89	43
Natural Resource Management	28.5	26.42	29.02	37.82	33.16	44.56
Nematology	22.83	28.26	28.26	29.35	35.33	41.3
Plant Pathology	28.97	30.48	27.96	42.82	34.01	44.84
Plant Sciences & Physiology	28.68	31.78	37.98	50.39	43.41	54.26
Seed Science and Technology	29.7	28.71	27.72	37.13	35.15	38.61
Soil Science	31.25	29.92	31.69	38.84	37.14	45.25
Veterinary Sciences	27.08	14.58	37.5	47.92	43.75	70.83
BBL						
Civil Litigation & Procedure	29.32	28.18	32.4	34.97	36.68	42.66
Constitutional & Administrative Law	29.54	28.15	36.22	40.62	42.28	49.46
Consumer & Competition Law	28	22.67	28	34.67	46.67	41.33
Corporate & Commercial Law	27.7	28.63	29.78	34.67	35.15	42.67
Criminal Law & Justice	27.09	26.98	30.01	33.66	35.21	42.72
Employment & Labour Law	23.43	25.71	28.57	29.1	32	40
Environmental & Energy Law	27.67	24.42	33.49	37.91	39.07	45.81

Continued on next page

2808

Table 19 – Continued from previous page

2809

2810

2811

2812

2813

2814

2815

2816

2817

2818

2819

2820

2821

2822

2823

2824

2825

2826

2827

2828

2829

2830

2831

2832

2833

2834

2835

2836

2837

2838

2839

2840

2841

2842

2843

2844

2845

2846

2847

2848

2849

2850

2851

2852

2853

2854

2855

2856

2857

2858

2859

2860

2861

Subject Domain	3.2-1B	3.2-1B-it	3.2-3B	3.2-3B-it	3.1-8B	3.1-8B-it
Family & Personal Law	24.12	28.86	29.06	31.69	34.21	39.86
General Academic Subjects	29.21	32.52	37.47	43.91	46.87	52.68
Healthcare & Medical Law	40	20	68	40	64	60
Human Rights & Social Justice	21.05	42.11	36.84	26.32	31.58	36.84
Intellectual Property Law	30.77	31.87	46.15	45.05	56.04	58.24
Interdisciplinary Studies	33.33	28.1	38.57	41.32	43.25	53.72
International & Comparative Law	30.87	30.35	40.02	45.22	46.88	54.47
Legal Skills & Communication	25.74	27.33	28.68	30.15	28.31	32.72
Legal Theory & Jurisprudence	29.63	28.36	33.92	39.69	41.66	46.87
Media & Entertainment Law	33.33	35.19	42.59	51.85	38.89	53.7
Real Estate & Property Law	23.53	25.91	29.89	31.96	31.48	38.16
Tax & Revenue Law	27.71	31.6	40.26	38.1	41.56	43.29
Technology & Cyber Law	30.89	41.46	48.78	49.59	51.22	60.16

Table 20: Performance of Qwen model family across sub-domains in BhashaBench v1, comparing base and instruction-tuned variants (3B, 235B)

2826

2827

2828

2829

2830

2831

2832

2833

2834

2835

2836

2837

2838

2839

2840

2841

2842

2843

2844

2845

2846

2847

2848

2849

2850

2851

2852

2853

2854

2855

2856

2857

2858

2859

2860

2861

Subject Domain	2.5-3B	2.5-3B-it	3-235B-A22B-it-2507
BBA			
Administration, AYUSH & Miscellaneous	47.06	38.66	73.11
Agad Tantra & Forensic Medicine	39.86	32.71	63.88
Ayurvedic Literature & History	38.73	29.9	55.88
Dravyaguna & Bhaishajya	32.57	28.94	49.43
Kaumarbhritya & Pediatrics	38.52	30.11	55.32
Kayachikitsa (General Medicine & Internal Medicine in Ayurveda)	38.61	35.07	59.48
Panchakarma & Rasayana	30.35	29.59	49.54
Research & Statistics	62.86	52.86	91.43
Roga Vigyana (Diagnostics & Pathology)	58.75	53.75	82.5
Samhita & Siddhanta (Fundamentals)	36.79	31.93	55.22
Shalaky Tantra (ENT, Eye, Dentistry)	35.56	31.74	59.67
Shalya Tantra (Surgery)	37.45	33.08	60.46
Sharir (Anatomy & Physiology)	37.44	31.35	60.1
Stri Roga & Prasuti Tantra (Gynecology & Obstetrics)	40.73	34.24	66.82
Swasthavritta & Public Health	50.99	43.49	82.56
Yoga & Psychology	44.68	36.17	75.53
BBF			
Accounting	38.94	31.82	63.52
Banking Services	43.3	36.89	71.22
Behavioral Finance	52.24	44.78	71.64
Business Management	60.24	40.96	84.34
Commerce	43.57	33.72	63.62
Corporate Finance & Investment	40.22	37.58	63.52
Data & Analytics in Finance	35.43	28.35	53.54
Economics & Development Studies	43.8	44.16	73.36
Energy, Infrastructure & Finance	45.12	30.49	71.95
Environmental Finance	47.62	44.05	82.74
Finance Education	50.85	43.22	69.49
Financial Markets	42.55	42.55	70.21
Financial Technology	47.83	39.13	78.26
General Knowledge	41.56	38.22	74.95
Governance & Policy	45.3	38.16	74.15
Healthcare Economics	48.25	45.61	78.95
History, Sociology & Cultural Studies of Finance	38.58	38.58	83.46
Information Technology Finance	64.9	58.16	92.24
Insurance & Risk Management	30.95	38.1	64.29
Interdisciplinary Finance	41.83	36.6	79.74

Continued on next page

2862

Table 20 – Continued from previous page

2863	Subject Domain	2.5-3B	2.5-3B-it	3-235B-A22B-it-2507
2864	International Finance & Trade	49.4	42.17	78.31
2865	Language & Communication	45.77	42.71	77.06
2866	Legal Finance	38.24	23.53	76.47
2867	Marketing Finance	69.05	50	85.71
2868	Mathematics for Finance	34.18	29.85	58.04
2869	Problem Solving	27.88	26.2	47.12
2870	Rural Economics	47.13	45.21	80.46
2871	Science and Technology in Finance	40.59	43.56	72.28
2872	Sports, Media & Finance Linkages	44.44	53.33	68.89
2873	Taxation & Regulatory Compliance	56.13	38.71	74.84
2874	BBK			
2875	Agri-Environmental & Allied Disciplines	43.75	43.18	75.57
2876	Agricultural Biotechnology	55.34	51.15	91.6
2877	Agricultural Chemistry & Biochemistry	44.48	38.43	83.63
2878	Agricultural Economics & Policy	46.41	43.38	73.21
2879	Agricultural Engineering & Technology	41.39	37.3	67.21
2880	Agricultural Extension Education	46.25	42.51	72.87
2881	Agricultural Microbiology	54.05	43.24	90.99
2882	Agriculture Communication	44.49	44.49	78.35
2883	Agriculture Information Technology	52.63	54.21	74.74
2884	Agronomy	41.73	38.89	71.92
2885	Animal Sciences	47.97	46.62	77.7
2886	Crop Sciences	42.08	36.79	67.4
2887	Dairy & Poultry Science	52.81	46.07	75.28
2888	Entomology	39.94	39.66	77.44
2889	Fisheries and Aquaculture	38.24	50	79.41
2890	General Knowledge & Reasoning	44.48	41.6	73.22
2891	Genetics and Plant Breeding	43.44	44.22	76.86
2892	Horticulture	37.25	35.41	64.98
2893	Natural Resource Management	37.82	37.31	65.8
2894	Nematology	33.15	39.13	63.04
2895	Plant Pathology	40.55	36.52	78.34
2896	Plant Sciences & Physiology	45.74	48.06	86.82
2897	Seed Science and Technology	42.08	34.65	66.34
2898	Soil Science	42	39.35	72.37
2899	Veterinary Sciences	45.83	50	87.5
2900	BBL			
2901	Civil Litigation & Procedure	38.65	35.31	72.12
2902	Constitutional & Administrative Law	43.67	37.93	82.65
2903	Consumer & Competition Law	36	46.67	82.67
2904	Corporate & Commercial Law	40.74	37.7	77.11
2905	Criminal Law & Justice	38.21	34.45	75.44
2906	Employment & Labour Law	39.43	37.14	71.43
2907	Environmental & Energy Law	44.65	38.84	76.74
2908	Family & Personal Law	38.35	32.8	74.37
2909	General Academic Subjects	53.82	45.44	85.82
2910	Healthcare & Medical Law	56	40	88
2911	Human Rights & Social Justice	47.37	31.58	73.68
2912	Intellectual Property Law	60.44	54.95	87.91
2913	Interdisciplinary Studies	49.31	44.08	84.85
2914	International & Comparative Law	47.51	43.76	83.89
2915	Legal Skills & Communication	32.35	31.74	61.27
	Legal Theory & Jurisprudence	46.45	40.04	79.38
	Media & Entertainment Law	42.59	33.33	79.63
	Real Estate & Property Law	36.09	33.55	71.7
	Tax & Revenue Law	39.83	37.66	74.03
	Technology & Cyber Law	58.54	59.35	86.18

2916 Table 21: Performance of GPT model family across sub-domains in BhashaBench v1, comparing
 2917 different model sizes (20B, 120B, GPT-4o)
 2918

2919	Subject Domain	gpt-oss-20b	gpt-oss-120b	gpt-4o
2920	BBA			
2921	Administration, AYUSH & Miscellaneous	53.78	79.83	75.63
2922	Agad Tantra & Forensic Medicine	39.52	60.14	63.54
2923	Ayurvedic Literature & History	33.82	51.47	59.31
2924	Dravyaguna & Bhaishajya	30.75	44.48	54.78
2925	Kaumarbhritya & Pediatrics	35.99	51.4	56.58
2926	Kayachikitsa (General Medicine & Internal Medicine in Ayurveda)	39.06	54.69	60.69
2927	Panchakarma & Rasayana	28.36	41.44	50.76
2928	Research & Statistics	70.95	86.67	90
2929	Roga Vigyana (Diagnostics & Pathology)	66.25	82.5	81.25
2930	Samhita & Siddhanta (Fundamentals)	30.63	46.07	53.41
2931	Shalaky Tantra (ENT, Eye, Dentistry)	38.15	54.9	62.4
2932	Shalya Tantra (Surgery)	35.36	55.13	61.41
2933	Sharir (Anatomy & Physiology)	39.75	57.06	62.7
2934	Stri Roga & Prasuti Tantra (Gynecology & Obstetrics)	35.18	59.03	64.82
2935	Swasthavritta & Public Health	56.51	76.6	81.02
2935	Yoga & Psychology	41.49	70.74	73.94
2936	BBF			
2937	Accounting	35.45	73.61	49.55
2938	Banking Services	42.53	67.29	68.57
2939	Behavioral Finance	50.75	77.61	76.12
2940	Business Management	53.01	87.95	81.93
2941	Commerce	37.89	69.76	54.46
2942	Corporate Finance & Investment	37.25	73.63	61.43
2943	Data & Analytics in Finance	34.65	51.97	44.09
2944	Economics & Development Studies	46.72	69.34	71.53
2944	Energy, Infrastructure & Finance	39.02	64.63	67.07
2945	Environmental Finance	55.95	73.21	77.98
2946	Finance Education	46.61	73.73	74.58
2947	Financial Markets	61.7	59.57	72.34
2948	Financial Technology	47.83	73.91	78.26
2948	General Knowledge	48.42	77.18	77.18
2949	Governance & Policy	39.85	69.36	78.29
2950	Healthcare Economics	49.12	78.07	80.7
2951	History, Sociology & Cultural Studies of Finance	48.03	68.5	87.4
2952	Information Technology Finance	76.94	90.82	92.04
2953	Insurance & Risk Management	47.62	57.14	64.29
2953	Interdisciplinary Finance	45.1	73.2	75.82
2954	International Finance & Trade	54.22	75.9	85.54
2955	Language & Communication	47.57	74.42	77.48
2956	Legal Finance	41.18	64.71	76.47
2957	Marketing Finance	61.9	85.71	78.57
2958	Mathematics for Finance	30.05	76.16	41.28
2959	Problem Solving	26.63	64.14	42.65
2959	Rural Economics	47.89	75.86	82.76
2960	Science and Technology in Finance	45.54	77.23	73.27
2961	Sports, Media & Finance Linkages	46.67	75.56	73.33
2962	Taxation & Regulatory Compliance	44.52	68.39	73.55
2963	BBK			
2964	Agri-Environmental & Allied Disciplines	41.48	73.86	74.43
2965	Agricultural Biotechnology	65.27	89.69	89.31
2966	Agricultural Chemistry & Biochemistry	54.8	80.43	81.14
2967	Agricultural Economics & Policy	46.57	71.77	73.68
2968	Agricultural Engineering & Technology	39.75	62.7	66.8
2968	Agricultural Extension Education	43.93	69.25	75.19

Continued on next page

2970

Table 21 – Continued from previous page

2971

2972

2973

2974

2975

2976

2977

2978

2979

2980

2981

2982

2983

2984

2985

2986

2987

2988

2989

2990

2991

2992

2993

2994

2995

2996

2997

2998

2999

3000

3001

3002

3003

3004

3005

3006

3007

3008

3009

3010

3011

3012

3013

3014

3015

3016

3017

3018

3019

3020

3021

3022

3023

Subject Domain	gpt-oss-20b	gpt-oss-120b	gpt-4o
Agricultural Microbiology	53.15	89.19	94.59
Agriculture Communication	42.91	73.23	81.1
Agriculture Information Technology	51.58	75.26	68.42
Agronomy	44.1	68	72.43
Animal Sciences	53.38	69.59	76.35
Crop Sciences	41.71	64.66	68.85
Dairy & Poultry Science	52.81	75.28	78.65
Entomology	48.28	72.84	77.87
Fisheries and Aquaculture	50	64.71	73.53
General Knowledge & Reasoning	42.81	69.59	68.38
Genetics and Plant Breeding	44.47	74.04	75.84
Horticulture	41.26	61.88	70.14
Natural Resource Management	41.97	64.77	65.8
Nematology	42.93	64.13	64.67
Plant Pathology	41.56	71.03	78.34
Plant Sciences & Physiology	51.94	82.17	88.37
Seed Science and Technology	35.15	64.85	65.84
Soil Science	42.45	70.67	73.18
Veterinary Sciences	56.25	87.5	93.75
BBL			
Civil Litigation & Procedure	34.63	59.01	71.91
Constitutional & Administrative Law	41.06	75.56	83.15
Consumer & Competition Law	33.33	72	81.33
Corporate & Commercial Law	37.48	69.59	78.93
Criminal Law & Justice	35.14	65.11	75.95
Employment & Labour Law	33.14	62.86	73.14
Environmental & Energy Law	41.4	69.3	73.26
Family & Personal Law	37.03	63.87	72.86
General Academic Subjects	56.49	83.14	84.79
Healthcare & Medical Law	60	92	92
Human Rights & Social Justice	15.79	73.68	68.42
Intellectual Property Law	53.85	85.71	90.11
Interdisciplinary Studies	43.25	82.64	83.75
International & Comparative Law	48.86	79.42	81.7
Legal Skills & Communication	32.84	69.12	53.43
Legal Theory & Jurisprudence	42.08	75.16	81.21
Media & Entertainment Law	50	83.33	85.19
Real Estate & Property Law	32.59	59.62	71.7
Tax & Revenue Law	42.86	67.53	69.26
Technology & Cyber Law	56.91	86.18	86.99

E.3 QUALITATIVE ERROR ANALYSIS FOR LLAMA-3.1-8B

In this section, we present a qualitative error analysis of the Llama-3.1-8B model across four domains of BhashaBench V1: Ayurveda (BBA), Finance (BBF), Agriculture (BBK), and Legal (BBL). For each domain, we examine some examples where the model’s responses deviated from the correct answers, highlighting the nature of errors, underlying causes, and potential strategies for mitigation. This analysis provides insight into the model’s domain-specific weaknesses, including challenges in understanding classical Ayurvedic terminology, procedural and regulatory reasoning in finance and law, and domain-specific factual knowledge in agriculture. Through these examples, we aim to identify patterns of systematic errors, informing future improvements in multilingual and domain-aware LLM performance.

E.3.1 BBA QUALITATIVE ANALYSIS

The qualitative analysis in the Ayurvedic domain reveals that Llama-3.1-8B exhibits several recurring challenges across classical medical texts, ritualistic procedures, and disease subtypes. The model frequently struggles with domain-specific terminology, semantic overlaps between closely related therapeutic indications, and nuanced procedural instructions from traditional texts. Errors are

3024 often observed when differentiating between similar disease subtypes or interpreting complex in-
 3025 structions from Panchakarma, Rasayana, and Stri Roga protocols. Moreover, the model sometimes
 3026 mislabels classical terms or associates remedies with incorrect disease indications, demonstrating
 3027 gaps in both lexical knowledge and reasoning.

3028 For instance, confusion arises between general wasting conditions and respiratory disorders, or be-
 3029 tween correct handling steps in Mritashodhana rituals. In other cases, the model misidentifies disease
 3030 subtypes in Kayachikitsa or misinterprets repeated obstetric terms in Prasuti Tantra. These obser-
 3031 vations indicate that while Llama-3.1-8B has captured general patterns in Ayurvedic knowledge, it
 3032 requires explicit integration of classical formulations, accurate disease nomenclature, and procedural
 3033 rules to improve precision and reliability in reasoning.

3034 The error patterns suggest that enhancing the model with structured knowledge of classical
 3035 Ayurvedic texts, domain-specific terminologies, and procedural protocols can significantly reduce
 3036 semantic confusions and improve zero-shot performance in complex Ayurvedic tasks. Overall, the
 3037 analysis highlights the critical role of targeted domain knowledge in enabling large language models
 3038 to reason effectively within traditional medicine contexts.
 3039
 3040
 3041
 3042
 3043
 3044
 3045
 3046
 3047
 3048
 3049
 3050
 3051
 3052
 3053
 3054
 3055

Question:

3056 Dashamooladi Ghruta is used in. vyadhi

Options:

- 3057
 3058
 3059
 3060 A. Kasa vyadhi
 3061 B. Shosh
 3062 C. Pandu
 3063 D. Shwasa
 3064

3065 **Correct Answer:** B

3066 **Model Selected Answer:** A

3067 **Subject Domain:** Dravyaguna & Bhaishajya
 3068

Error Analysis:

- 3069
 3070
 3071
 3072
 3073
 3074
 3075
 3076
 3077
- **Nature of Error:** Confusion between respiratory disorders and general wasting conditions.
 - **Underlying Cause:** The model misinterpreted the therapeutic indications of Dashamooladi Ghruta, associating it with cough-related ailments instead of Shosh (emaciation or wasting disease).
 - **Recommendation:** Strengthen model knowledge of Ayurvedic formulations along with their specific disease indications to avoid semantic overlaps.

3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131

Question:

Which of the following options is most appropriate regarding removal of the body from water for Mritashodhana?

Options:

- A. Samyak prakuthit
- B. Seven days
- C. Both A & B
- D. None of A & B

Correct Answer: A

Model Selected Answer: D

Subject Domain: Panchakarma & Rasayana

Error Analysis:

- **Nature of Error:** Misunderstanding of procedural rules in classical texts.
- **Underlying Cause:** The model failed to correctly associate the Mritashodhana process with the proper handling of a deceased body, leading to selection of an incorrect “none” option.
- **Recommendation:** Incorporate explicit procedural knowledge from classical Ayurvedic texts to enhance reasoning on ritualistic and therapeutic protocols.

Question:

वात पित्त प्रधान विसर्प को - - - - - कहा जाता है

Options:

- A. अग्नि
- B. कर्दम
- C. ग्रंथि
- D. निचय

Correct Answer: A

Model Selected Answer: C

Subject Domain: Kayachikitsa (General Medicine & Internal Medicine in Ayurveda)

Error Analysis:

- **Nature of Error:** Mislabeling disease subtype in classical terminology.
- **Underlying Cause:** The model failed to correctly identify the term for Vata-Pitta dominant Visarpa, selecting “Granthii” instead of “Agni”.
- **Recommendation:** Include domain-specific classical terminology in training for accurate disease nomenclature recognition.

3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185

Question:

पहली तिमाही में बार-बार होने वाले गर्भपात का अर्थ है - - - - - वंध्यता।

Options:

- A. मृत्वत्सा
- B. काकवंध्या
- C. गर्भस्त्रावी
- D. बाला

Correct Answer: A

Model Selected Answer: D

Subject Domain: Stri Roga & Prasuti Tantra (Gynecology & Obstetrics)

Error Analysis:

- **Nature of Error:** Misinterpretation of obstetric terminology.
- **Underlying Cause:** The model selected “Bala” instead of “Mritvatsa” due to misunderstanding repeated first-trimester miscarriage terminology.
- **Recommendation:** Incorporate classical and modern obstetric definitions for precise identification of pathological conditions.

E.3.2 BBF QUALITATIVE ANALYSIS

In the finance domain, errors are primarily due to confusion between regulatory and advisory institutions, misidentification of financial instruments, and misapplication of problem-solving reasoning. The model occasionally fails to account for context-specific definitions or institutional roles. Incorporating structured financial knowledge, including Indian regulatory frameworks and definitions of financial instruments, would enhance model performance.

Question:

Which of these has set up a high-level panel to suggest possible structures and regulations for creating social stock exchanges to facilitate listing and fund-raising by social enterprises as well as voluntary organisations?

Options:

- A. RBI
- B. NITI Aayog
- C. Finance Ministry
- D. SEBI

Correct Answer: D

Model Selected Answer: B

Subject Domain: Taxation & Regulatory Compliance

Error Analysis:

- **Nature of Error:** Confusion between regulatory and advisory bodies.
- **Underlying Cause:** The model associated policy advice with NITI Aayog rather than SEBI’s regulatory role in social stock exchanges.
- **Recommendation:** Enhance knowledge of institutional functions in Indian financial and regulatory frameworks.

3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239

Question:

Directions: Read the following comprehension carefully and answer the questions given below.

A certain number of people are sitting in a linear row facing North.

- R sits fifth from the left end of the row.
- Only two persons sit between R and G, who sits to the right of R.
- X sits second to the left of G.
- Only one person sits to the left of L.
- Three persons sit between X and A, who sits second from the right end of the row.
- L sits third to the left of R.
- V sits sixth to the right of R.

How many persons sit between L and G?

Options:

- A. 8
- B. 3
- C. 4
- D. 5

Correct Answer: D

Model Selected Answer: B

Subject Domain: Problem Solving

Error Analysis:

- **Nature of Error:** Miscalculation in linear arrangement reasoning.
- **Underlying Cause:** The model incorrectly interpreted the positions from the textual description, failing to properly count the number of people between L and G.
- **Recommendation:** Strengthen stepwise logical reasoning on linear arrangement and seating problems, including visualizing relative positions.

Question:

एडवांस प्राइसिंग एग्रीमेंट्स निम्नलिखित में से किसमें उपयोग होने वाला शब्द है?

Options:

- A. अंतर्राष्ट्रीय संधि
- B. केंद्रीय बजट
- C. जीएसटी कार्यान्वयन
- D. कराधान

Correct Answer: D

Model Selected Answer: B

Subject Domain: Taxation & Regulatory Compliance

Error Analysis:

- **Nature of Error:** Confusion regarding the context of terminology usage.
- **Underlying Cause:** The model associated budgeting context instead of tax compliance context, leading to selection of a non-relevant option.
- **Recommendation:** Include context-specific knowledge about taxation terminology and international agreements in the training data.

3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293

Question:
ग्लोबल डिपॉजिटरी रसीद क्या है?

Options:

- A. भारत में एक विदेशी कंपनी द्वारा जारी बॉन्ड
- B. एक डिपॉजिटरी बैंक द्वारा जारी प्रमाण पत्र जो एक विदेशी कंपनी के शेयरों का प्रतिनिधित्व करता है
- C. म्यूचुअल फंड योजना का एक प्रकार
- D. विदेशी बाजारों में जारी सरकारी प्रतिभूति

Correct Answer: B
Model Selected Answer: D
Subject Domain: Financial Markets
Error Analysis:

- **Nature of Error:** Misidentification of financial instrument.
- **Underlying Cause:** The model selected a government security in foreign markets instead of recognizing the depository receipt issued by a bank for foreign company shares.
- **Recommendation:** Provide clearer definitions of financial instruments like GDRs, ADRs, and bonds in model knowledge base.

E.3.3 BBK QUALITATIVE ANALYSIS

Agricultural domain errors stem from gaps in crop-specific practices, breed characteristics, and interpretation of scientific explanations in Hindi. The model frequently generalizes concepts, leading to incorrect associations between schemes, cultivation techniques, or physiological traits. Including detailed domain-specific datasets and traditional agronomy knowledge can help the model distinguish between nuanced agricultural concepts.

Question:
Which of the following is the objective of Atal Bhujal Yojana?

Options:

- A. Drip Irrigation
- B. Hydro-electric power production
- C. Drinking water supply
- D. Groundwater Management

Correct Answer: D
Model Selected Answer: A
Subject Domain: Agronomy
Error Analysis:

- **Nature of Error:** Confusion between water management schemes and irrigation programs.
- **Underlying Cause:** The model associated the scheme with agricultural water use but misidentified its focus, mistaking groundwater governance for irrigation techniques.
- **Recommendation:** Incorporate structured knowledge regarding Indian water-related government schemes and their specific objectives.

3294
3295
3296
3297
3298
3299
3300
3301
3302
3303
3304
3305
3306
3307
3308
3309
3310
3311
3312
3313
3314
3315
3316
3317
3318
3319
3320
3321
3322
3323
3324
3325
3326
3327
3328
3329
3330
3331
3332
3333
3334
3335
3336
3337
3338
3339
3340
3341
3342
3343
3344
3345
3346
3347

Question:

What are the horn characteristics of Murrah breed?

Options:

- A. Flat and sickle shaped and form hook at the tip
- B. Small and coiled tightly
- C. Curl slightly outwards
- D. Flat, short, tightly spirally curving inwards

Correct Answer: D

Model Selected Answer: C

Subject Domain: Animal Sciences

Error Analysis:

- **Nature of Error:** Breed-specific trait confusion.
- **Underlying Cause:** The model generalized bovine horn types without differentiating breed-specific characteristics.
- **Recommendation:** Include detailed breed-specific morphology datasets to improve precision.

Question:

धान की फसल की जलमग्नता सह लेती है, क्योंकि

Options:

- A. पौधे को औक्सीजन की जरूरत नहीं होती
- B. पौधों में पत्तियों में से औक्सीजन ले जाने की प्रक्रिया होती है
- C. जीवाणु मूल परिवेश को औक्सीकृत रखते हैं
- D. हवा के साथ जड़ों के क्षेत्र में प्रवेश कर रही औक्सीजन पर्याप्त है

Correct Answer: B

Model Selected Answer: C

Subject Domain: Agronomy

Error Analysis:

- **Nature of Error:** Hindi comprehension and scientific reasoning.
- **Underlying Cause:** The model misinterpreted the physiological explanation for rice root tolerance to waterlogging.
- **Recommendation:** Include domain-specific Hindi texts explaining crop physiology and waterlogging tolerance.

3348
3349
3350
3351
3352
3353
3354
3355
3356
3357
3358
3359
3360
3361
3362
3363
3364
3365
3366
3367
3368
3369
3370
3371
3372
3373
3374
3375
3376
3377
3378
3379
3380
3381
3382
3383
3384
3385
3386
3387
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399
3400
3401

Question:

किस फसल की बुवाई से पहले मिट्टी पलट हल से दो बार जुताई करते हैं?

Options:

- A. अरहर
- B. मक्का
- C. चना
- D. गन्ना

Correct Answer: D

Model Selected Answer: B

Subject Domain: Agronomy

- **Nature of Error:** Crop management knowledge gap.
- **Underlying Cause:** The model failed to correctly associate double plowing with sugarcane cultivation.
- **Recommendation:** Include crop-specific cultivation practices and traditional agronomy knowledge in the training data.

E.3.4 BBL QUALITATIVE ANALYSIS

In the legal domain, Llama-3.1-8B exhibits errors primarily related to procedural understanding, definitional knowledge, and interpretation of civil, criminal, and property law provisions. Common mistakes include misclassifying offences, misapplying procedural codes, and assuming incorrect registration or filing requirements. Errors often arise from insufficient knowledge of statutory definitions, procedural sequences, and conditions for specific legal actions. Strengthening explicit legal knowledge, clarifying judicial powers, and reinforcing procedural rules can improve the model's factual correctness and reasoning for law-related queries.

Question:

Question: Suit for recovery of money in promissory notes can be filed:

Options:

- A. under normal procedure
- B. under summary procedure as laid down in Order 37, CPC
- C. in the High Court
- D. as a writ petition

Correct Answer: B

Model Selected Answer: C

Subject Domain: Civil Litigation & Procedure

Error Analysis:

- **Nature of Error:** Confusion regarding procedural rules for recovery suits.
- **Underlying Cause:** The model incorrectly associated the filing with the High Court instead of the summary procedure under Order 37 CPC.
- **Recommendation:** Reinforce knowledge of civil procedural codes and specific rules for summary recovery of promissory notes.

3402
3403
3404
3405
3406
3407
3408
3409
3410
3411
3412
3413
3414
3415
3416
3417
3418
3419
3420
3421
3422
3423
3424
3425
3426
3427
3428
3429
3430
3431
3432
3433
3434
3435
3436
3437
3438
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449
3450
3451
3452
3453
3454
3455

Question:

Question: Use of violence by a member of an assembly of five or more person in furtherance of common object will constitute the offence of

Options:

- A. Affray
- B. Assault
- C. Rioting
- D. Unlawful Assembly

Correct Answer: C

Model Selected Answer: D

Subject Domain: Criminal Law & Justice

Error Analysis:

- **Nature of Error:** Misclassification of the offence type.
- **Underlying Cause:** The model confused “rioting” with “unlawful assembly,” missing the condition of violence by five or more persons in furtherance of common object.
- **Recommendation:** Reinforce understanding of key definitions and conditions for offences under criminal law.

Question:

जहाँ न्यायालय को किसी व्यक्ति के पास किसी दस्तावेज या वस्तु के होने के बारे में कोई ज्ञान नहीं है? क्या ऐसी स्थिति में न्यायालय खोज वारंट जारी कर सकता है:

Options:

- A. नहीं
- B. केवल उस स्थिति में जब किसी विशिष्ट वस्तु के बारे में ज्ञात हो
- C. हाँ
- D. जब कोई विशिष्ट स्थान या व्यक्ति निर्दिष्ट हो

Correct Answer: C

Model Selected Answer: B

Subject Domain: Civil Litigation & Procedure

Error Analysis:

- **Nature of Error:** Misinterpretation of procedural powers regarding search warrants.
- **Underlying Cause:** The model assumed prior knowledge of a specific document was required, instead of recognizing the court’s general power to issue a search warrant.
- **Recommendation:** Emphasize the scope and limits of judicial powers in procedural law.

3456
3457
3458
3459
3460
3461
3462
3463
3464
3465
3466
3467
3468
3469
3470
3471
3472
3473
3474
3475
3476
3477
3478
3479
3480
3481
3482
3483
3484
3485
3486
3487
3488
3489
3490
3491
3492
3493
3494
3495
3496
3497
3498
3499
3500
3501
3502
3503
3504
3505
3506
3507
3508
3509

Question:

अचल संपत्ति का किराया अनिवार्य रूप से पंजीकृत नहीं है:

Options:

- A. वार्षिक किराया
- B. एक वर्ष की अवधि के लिए किराया
- C. वार्षिक किराया प्राप्त करने वाला किराया
- D. इनमें से कोई नहीं

Correct Answer: D

Model Selected Answer: A

Subject Domain: Real Estate & Property Law

Error Analysis:

- **Nature of Error:** Incorrect assumption regarding registration requirement.
- **Underlying Cause:** The model assumed annual rent requires registration, overlooking the legal nuance that not all rent agreements require registration.
- **Recommendation:** Include precise rules on registration requirements for different types of rent agreements in real estate law.

E.4 DATA INTEGRITY AND CONTAMINATION ANALYSIS

To ensure the validity of our benchmark evaluations and rule out potential data leakage, we perform a set of analyses covering perplexity-based checks, multiple-choice option shuffling, and the impact of increasing distractor options. These evaluations help verify that our datasets remain unbiased and challenging for state-of-the-art LLMs.

E.4.1 PERPLEXITY-BASED DATA CONTAMINATION ANALYSIS

To verify the integrity of BhashaBench V1 and detect potential data contamination or leakage from pretraining, we conducted a detailed perplexity (PPL) analysis on Llama-3.1-8B and Gemma-2-9B models. For each dataset, we computed PPL over the entire multiple-choice items, including both the questions and all answer options. This ensures that the computed PPL reflects the model’s familiarity with the phrasing of the dataset, rather than just the questions alone. Table 22 reports the PPL, average token-level loss, and the number of tokens evaluated for each benchmark. BhashaBench V1 datasets (BBA, BBF, BBK, BBL) show perplexity scores comparable to or higher than well-established benchmarks such as ARC-C, MMLU, and MILU, confirming minimal exposure of these items during pretraining. For instance, BBA exhibits high PPL on both English (15.50) and Hindi (10.39) for Llama-3.1-8B, while BBK shows relatively elevated PPL in Hindi (7.16 for Llama-3.1-8B). These results indicate that the datasets provide genuinely novel evaluation challenges.

E.4.2 MULTIPLE-CHOICE OPTION SHUFFLING EXPERIMENT

To further validate the robustness of models against superficial cues, we conducted an option ordering experiment. Multiple-choice options in each question were shuffled using fixed seeds (42 and 123) to examine whether model performance depends on the original ordering. The “Base” column represents the original option order. Table 23 presents the performance of Llama-3.1-8B and Llama-3.2-3B across English and Hindi subsets of each domain. Overall, performance remains largely stable across different option orderings, with minor variations observed primarily in the Hindi subsets. This suggests that the benchmark’s multiple-choice questions evaluate actual model understanding rather than positional biases, providing confidence in the integrity of the dataset for comparative evaluation.

E.4.3 EFFECT OF SCALING THE NUMBER OF DISTRACTORS

To assess how task difficulty affects model discrimination, we increased the number of options from the base 4 to 5 and 6 options per question. Additional distractor options were generated determin-

3510
3511
3512
3513
3514
3515
3516
3517
3518
3519
3520
3521
3522
3523
3524
3525
3526
3527
3528
3529
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549
3550
3551
3552
3553
3554
3555
3556
3557
3558
3559
3560
3561
3562
3563

Table 22: Perplexity (PPL) and Average Loss for Llama-3.1-8B and gemma-2-9b across BhashaBench V1 and other datasets

Dataset	Language	Llama-3.1-8B			gemma-2-9b		
		PPL	Avg. Loss	Num. Tokens	PPL	Avg. Loss	Num. Tokens
ARC-C	English	8.03	2.08	74,929	6.82	1.92	77,685
	Hindi	4.10	1.41	160,887	5.85	1.77	124,073
MILU	English	7.62	2.03	845,273	7.23	1.98	879,220
	Hindi	4.93	1.60	1,521,137	6.37	1.85	1,261,389
MMLU	English	7.61	2.03	1,502,590	7.03	1.95	1,553,587
	Hindi	4.22	1.44	3,024,075	6.21	1.83	2,292,437
BBA	English	15.50	2.74	445,077	23.20	3.14	438,301
	Hindi	10.39	2.34	373,338	16.80	2.82	331,130
BBF	English	6.78	1.91	1,710,390	5.86	1.77	1,821,417
	Hindi	4.03	1.39	1,007,527	5.04	1.62	799,074
BBK	English	6.14	1.81	965,854	6.34	1.85	980,286
	Hindi	7.16	1.97	211,467	9.54	2.26	183,744
BBL	English	7.28	1.98	1,511,635	7.19	1.97	1,567,174
	Hindi	4.01	1.39	1,071,001	5.79	1.76	825,463

Table 23: Performance of Llama-3.1-8B and Llama-3.2-3B across BhashaBench V1 (EN + HI) under different option orderings. Each “Seed” column corresponds to a different shuffling of the multiple-choice options, while “Base” uses the original ordering. Scores are reported for overall, English, and Hindi subsets of each domain.

Domain		Llama-3.1-8B			Llama-3.2-3B		
Category	Language	Base	Seed 42	Seed 123	Base	Seed 42	Seed 123
BBA	Overall	33.12	32.01	31.99	30.28	29.55	29.07
	English	35.48	34.20	34.20	31.62	31.14	30.43
	Hindi	29.17	28.37	28.30	28.05	26.91	26.79
BBF	Overall	34.48	33.62	33.46	31.46	30.60	30.82
	English	36.20	35.66	35.45	33.04	32.29	32.45
	Hindi	30.61	29.04	28.99	27.92	26.80	27.15
BBK	Overall	38.07	36.42	36.87	31.96	31.24	31.61
	English	39.52	37.97	38.41	32.68	31.96	32.34
	Hindi	31.41	29.31	29.82	28.69	27.93	28.29
BBL	Overall	38.44	37.61	37.46	33.17	32.52	32.41
	English	41.32	40.51	40.10	35.17	34.40	34.32
	Hindi	31.76	30.87	31.33	28.53	28.15	27.96

istically using fixed random seeds (seed 48 for the 5th option, seed 123 for the 6th option). These distractors consist of semantically meaningless random character sequences (5-8 lowercase letters) that test whether models can robustly filter out irrelevant choices. The generation process seeds Python’s random number generator to create deterministic nonsense tokens from lowercase ASCII characters, ensuring reproducibility while introducing noise distractors clearly distinguishable from genuine options.

Table 24 presents the performance of Llama-3.1-8B and Llama-3.2-3B across all BhashaBench V1 domains with varying option counts. Results show consistent performance decline as options increase across both models and all domains. For Llama-3.1-8B, overall accuracy drops 2-7 percentage points from 4 to 6 options, with the steepest degradation in BBL (38.44% \rightarrow 28.96%) and BBK (38.07% \rightarrow 32.81%). Llama-3.2-3B exhibits similar but slightly steeper declines, particularly in BBL (33.17% \rightarrow 22.89%).

The performance gap between English and Hindi widens substantially with additional options. Hindi subsets show more severe degradation Llama-3.1-8B’s Hindi performance on BBL drops from 31.76% to 16.03%, compared to 41.32% to 34.52% in English. This disproportionate challenge in lower-resource languages suggests weaker multilingual reasoning capabilities. The significant drops even with nonsense distractors indicate models struggle with filtering irrelevant options, particularly in Hindi. These findings confirm the benchmark maintains discriminative power with additional options, enabling fine-grained model assessment.

Table 24: Performance of Llama-3.1-8B and Llama-3.2-3B across BhashaBench V1 domains with varying numbers of options. Scores are reported for overall, English, and Hindi subsets of each domain.

Domain		Llama-3.1-8B			Llama-3.2-3B		
Category	Language	Base (4 Opts)	5 Options	6 Options	Base (4 Opts)	5 Options	6 Options
BBA	Overall	33.12	31.14	29.02	30.28	26.9	25.7
	English	35.48	35.07	34.61	31.62	30.36	30.35
	Hindi	29.17	24.59	19.73	28.05	21.14	17.95
BBF	Overall	34.48	31.98	31.01	31.46	28.43	27.65
	English	36.2	34.18	33.61	33.04	31.33	30.48
	Hindi	30.61	27.05	25.18	27.92	21.92	21.3
BBK	Overall	38.07	33.88	32.81	31.96	28.98	26.85
	English	39.52	35.82	35.52	32.68	31.07	28.98
	Hindi	31.41	24.99	20.38	28.69	19.37	17.12
BBL	Overall	38.44	31.49	28.96	33.17	24.49	22.89
	English	41.32	36.59	34.52	35.17	30.43	28.63
	Hindi	31.76	19.6	16.03	28.53	10.66	9.52

E.5 STATISTICAL SIGNIFICANCE TESTS

To assess whether the observed differences in model performance are statistically meaningful, we conduct a series of statistical significance tests. These tests help determine whether performance gaps between models are likely due to chance or reflect consistent trends across languages and benchmarks.

E.5.1 WILSON CONFIDENCE INTERVAL FOR MODEL PERFORMANCE

To quantify the statistical reliability of model performance estimates and enable rigorous comparative analysis, we computed Wilson confidence intervals (CIs) for the zero-shot accuracy of each evaluated model. Wilson CIs provide a robust measure of uncertainty around estimated accuracy that is particularly well-suited for evaluation benchmarks. Unlike naive proportion-based intervals (such as normal approximation intervals), Wilson CIs employ a score-based method that offers several critical advantages: they maintain more accurate coverage probabilities, especially for extreme accuracy values near 0% or 100%; they naturally respect the bounded nature of proportion metrics

Table 25: Zero-shot performance (%) of LLMs on the BBA with Wilson confidence intervals. Values in brackets indicate the Wilson CIs.

Model	Eng (%) [95% CI]	Hin (%) [95% CI]	Avg (%) [95% CI]
<i>< 4B Models</i>			
gemma-3-270m	28.08 [27.18, 29.00]	28.25 [27.08, 29.44]	28.14 [27.43, 28.87]
gemma-3-270m-it	26.23 [25.35, 27.13]	25.77 [24.64, 26.93]	26.06 [25.36, 26.77]
Param-1	41.12 [40.13, 42.12]	38.04 [36.78, 39.32]	39.97 [39.18, 40.75]
gemma-2-2b	36.80 [35.83, 37.78]	30.61 [29.42, 31.83]	34.48 [33.72, 35.24]
gemma-2-2b-it	29.38 [28.46, 30.31]	26.79 [25.64, 27.96]	28.40 [27.69, 29.13]
Llama-3.2-1B	29.17 [28.26, 30.10]	26.30 [25.17, 27.47]	28.10 [27.38, 28.82]
Llama-3.2-1B-Instruct	26.77 [25.88, 27.67]	25.82 [24.70, 26.98]	26.41 [25.71, 27.12]
Llama-3.2-3B	31.62 [30.69, 32.57]	28.05 [26.89, 29.24]	30.28 [29.55, 31.02]
Llama-3.2-3B-Instruct	35.31 [34.35, 36.29]	29.67 [28.49, 30.88]	33.20 [32.45, 33.95]
sarvam-2b-v0.5	26.79 [25.90, 27.69]	27.07 [25.92, 28.25]	26.89 [26.19, 27.61]
sarvam-1	29.70 [28.78, 30.63]	28.41 [27.24, 29.60]	29.21 [28.49, 29.95]
Nemotron-4-Mini-Hindi-4B-Base	34.76 [33.80, 35.73]	32.82 [31.61, 34.06]	34.03 [33.28, 34.79]
Nemotron-4-Mini-Hindi-4B-Instruct	33.38 [32.43, 34.34]	33.82 [32.59, 35.07]	33.54 [32.79, 34.30]
Qwen2.5-3B	40.61 [39.62, 41.61]	31.90 [30.69, 33.13]	37.34 [36.57, 38.12]
Qwen2.5-3B-Instruct	35.22 [34.25, 36.19]	28.46 [27.29, 29.65]	32.68 [31.93, 33.44]
granite-3.1-2b-instruct	33.39 [32.44, 34.35]	27.30 [26.15, 28.48]	31.10 [30.37, 31.85]
granite-3.1-3b-a800m-base	31.75 [30.81, 32.70]	26.18 [25.05, 27.35]	29.66 [28.93, 30.40]
<i>7B to 27B Models</i>			
Pangea-7B	40.69 [39.70, 41.69]	31.93 [30.73, 33.16]	37.41 [36.63, 38.18]
Indic-gemma-7b-finetuned-sft-Navarasa-2.0	37.12 [36.15, 38.10]	31.83 [30.62, 33.06]	35.13 [34.37, 35.90]
aya-23-8B	33.84 [32.88, 34.80]	28.87 [27.70, 30.07]	31.97 [31.23, 32.72]
Llama-3.1-8B	35.48 [34.52, 36.46]	29.17 [28.00, 30.37]	33.12 [32.37, 33.87]
Llama-3.1-8B-Instruct	36.86 [35.89, 37.85]	31.26 [30.06, 32.48]	34.76 [34.00, 35.53]
gemma-2-9b	48.16 [47.15, 49.17]	37.92 [36.66, 39.19]	44.32 [43.52, 45.11]
gemma-2-9b-it	36.22 [35.25, 37.20]	31.18 [29.99, 32.41]	34.33 [33.57, 35.10]
gpt-oss-20b	38.30 [37.32, 39.29]	33.09 [31.87, 34.33]	36.34 [35.58, 37.12]
gemma-2-27b	50.70 [49.68, 51.71]	42.26 [40.98, 43.56]	47.53 [46.73, 48.33]
gemma-2-27b-it	40.45 [39.46, 41.45]	33.89 [32.66, 35.14]	37.99 [37.21, 38.77]
<i>> 27B Models</i>			
gpt-oss-120b	55.62 [54.61, 56.62]	48.05 [46.74, 49.36]	52.78 [51.98, 53.58]
Qwen3-235B-A22B-Instruct-25076	60.25 [59.25, 61.24]	54.78 [53.48, 56.08]	58.20 [57.40, 58.98]
deepseek-v3	51.38 [50.37, 52.39]	37.03 [35.77, 38.30]	45.99 [45.20, 46.79]
gpt-4o	62.75 [61.77, 63.73]	54.73 [53.42, 56.03]	59.74 [58.95, 60.52]

(always yielding intervals within $[0, 1]$); and they perform reliably even with relatively small sample sizes, making them ideal for domain-specific benchmark subsets.

Tables 25, 26, 27, and 28 report the zero-shot accuracy (%) for English and Hindi subsets across all evaluated models for each BhashaBench domain (BBA, BBF, BBK, and BBL respectively), accompanied by their corresponding 95% Wilson CIs. Values presented in brackets [lower, upper] denote the lower and upper bounds of the confidence interval, representing the range within which the true population accuracy is expected to lie with 95% confidence, given the observed sample performance.

These intervals enable both visual and quantitative comparison of model performance while explicitly accounting for sampling uncertainty inherent in finite test sets. Models whose confidence intervals do not overlap can be interpreted as exhibiting statistically significant differences in performance at the 95% confidence level, providing strong evidence that observed accuracy differences reflect genuine capability gaps rather than random variation. Conversely, when confidence intervals overlap substantially, performance differences should be interpreted more cautiously, as they may fall within the margin of statistical uncertainty. This interval-based analysis is particularly valuable for identifying robust performance trends across languages and domains, and for determining whether improvements from model scaling or architectural changes represent statistically meaningful advances. The Wilson CI framework thus strengthens the benchmark’s utility for reliable model comparison and selection decisions.

3672 Table 26: Zero-shot performance (%) of LLMs on the BBF with Wilson confidence intervals. Values
 3673 in brackets indicate the Wilson CIs.

3675	Model	Eng (%) [95% CI]	Hin (%) [95% CI]	Avg (%) [95% CI]
3676	<i>< 4B Models</i>			
3677	gemma-3-270m	24.98 [24.26, 25.72]	25.06 [23.98, 26.17]	25.00 [24.40, 25.62]
3678	gemma-3-270m-it	24.13 [23.42, 24.86]	23.84 [22.78, 24.93]	24.04 [23.45, 24.65]
3679	Param-1	32.24 [31.46, 33.04]	29.56 [28.41, 30.72]	31.42 [30.77, 32.07]
3680	gemma-2-2b	34.20 [33.40, 35.00]	27.50 [26.38, 28.64]	32.14 [31.48, 32.80]
3681	gemma-2-2b-it	31.26 [30.48, 32.05]	27.93 [26.81, 29.08]	30.24 [29.60, 30.89]
3682	Llama-3.2-1B	28.24 [27.48, 29.00]	25.61 [24.52, 26.73]	27.43 [26.80, 28.06]
3683	Llama-3.2-1B-Instruct	26.28 [25.54, 27.03]	26.04 [24.95, 27.17]	26.21 [25.59, 26.83]
3684	Llama-3.2-3B	33.04 [32.25, 33.84]	27.92 [26.79, 29.07]	31.46 [30.81, 32.12]
3685	Llama-3.2-3B-Instruct	32.94 [32.15, 33.74]	29.09 [27.95, 30.25]	31.76 [31.10, 32.41]
3686	sarvam-2b-v0.5	26.42 [25.68, 27.17]	25.31 [24.22, 26.43]	26.08 [25.47, 26.70]
3687	sarvam-1	29.66 [28.89, 30.43]	27.27 [26.15, 28.41]	28.92 [28.29, 29.56]
3688	Nemotron-4-Mini-Hindi-4B-Base	34.95 [34.15, 35.76]	31.41 [30.25, 32.60]	33.86 [33.20, 34.53]
3689	Nemotron-4-Mini-Hindi-4B-Instruct	31.98 [31.20, 32.78]	30.06 [28.91, 31.23]	31.39 [30.74, 32.05]
3690	Qwen2.5-3B	39.54 [38.72, 40.37]	32.13 [30.96, 33.32]	37.26 [36.58, 37.94]
3691	Qwen2.5-3B-Instruct	34.84 [34.04, 35.65]	29.17 [28.03, 30.34]	33.09 [32.44, 33.76]
3692	granite-3.1-2b-instruct	32.82 [32.03, 33.62]	27.11 [26.00, 28.26]	31.07 [30.42, 31.72]
3693	granite-3.1-3b-a800m-base	29.22 [28.45, 29.99]	24.17 [23.10, 25.27]	27.66 [27.04, 28.30]
3694	<i>7B to 27B Models</i>			
3695	Pangea-7B	41.71 [40.88, 42.54]	33.73 [32.55, 34.94]	39.25 [38.57, 39.94]
3696	Indic-gemma-7b-finetuned-sft-Navarasa-2.0	37.00 [36.19, 37.82]	30.47 [29.32, 31.65]	34.99 [34.32, 35.67]
3697	aya-23-8B	35.25 [34.44, 36.06]	30.88 [29.72, 32.06]	33.90 [33.24, 34.57]
3698	Llama-3.1-8B	36.20 [35.39, 37.01]	30.61 [29.45, 31.79]	34.48 [33.81, 35.15]
3699	Llama-3.1-8B-Instruct	35.68 [34.87, 36.49]	30.27 [29.12, 31.45]	34.01 [33.35, 34.68]
3700	gemma-2-9b	42.73 [41.90, 43.57]	36.91 [35.70, 38.14]	40.94 [40.25, 41.63]
3701	gemma-2-9b-it	38.85 [38.03, 39.68]	32.03 [30.86, 33.22]	36.75 [36.08, 37.43]
3702	gpt-oss-20b	37.11 [36.30, 37.93]	32.61 [31.44, 33.81]	35.73 [35.06, 36.40]
3703	gemma-2-27b	47.79 [46.94, 48.63]	41.24 [40.00, 42.49]	45.77 [45.07, 46.47]
3704	gemma-2-27b-it	42.47 [41.64, 43.31]	34.29 [33.09, 35.50]	39.95 [39.27, 40.64]
3705	<i>> 27B Models</i>			
3706	gpt-oss-120b	74.11 [73.37, 74.85]	64.16 [62.94, 65.36]	71.05 [70.41, 71.68]
3707	Qwen3-235B-A22B-Instruct-25076	63.72 [62.90, 64.53]	56.27 [55.01, 57.52]	61.43 [60.74, 62.11]
3708	deepseek-v3	63.46 [62.64, 64.27]	57.04 [55.78, 58.29]	61.48 [60.80, 62.16]
3709	gpt-4o	57.27 [56.43, 58.10]	49.82 [48.55, 51.08]	54.97 [54.27, 55.67]

3706 E.5.2 STATISTICAL SIGNIFICANCE OF MODEL PERFORMANCE DIFFERENCES USING 3707 McNemar’s TEST

3708 To evaluate whether the observed differences in accuracy between the top-performing LLMs are
 3709 statistically meaningful, we employ **McNemar’s test** (Mcnemar, 1947). This non-parametric test
 3710 is specifically designed for paired nominal data and is commonly used to compare two classifiers
 3711 on the same dataset by focusing on instances where their predictions disagree. It provides a robust
 3712 measure of whether one model consistently outperforms another beyond random chance.

3713 In our analysis, the **Accuracy Diff (%)** column reports the absolute difference in accuracy between a
 3714 pair of models. The **McNemar Stat** quantifies the magnitude of disagreement in model predictions,
 3715 while the corresponding **p-value** assesses the statistical significance of this difference. We consider
 3716 a difference significant if $p < 0.05$, implying that one model’s performance is reliably better than
 3717 the other on the same set of questions.

3718 Tables 29–32 show pairwise McNemar comparisons of the top 5 models across four benchmark
 3719 domains (BBA, BBF, BBK, BBL) for both English and Hindi datasets. Entries labeled “Yes” in the
 3720 **Significant** column indicate that the performance difference is statistically significant, confirming
 3721 consistent superiority of one model over another. Conversely, entries marked “No” denote cases
 3722 where the observed accuracy difference could be due to chance, suggesting that the two models
 3723 have comparable performance. This analysis complements raw accuracy scores by highlighting
 3724 which improvements are robust and reliable rather than incidental.

3725

3726
3727
3728
3729
3730
3731
3732
3733
3734
3735
3736
3737
3738
3739
3740
3741
3742
3743
3744
3745
3746
3747
3748
3749
3750
3751
3752
3753
3754
3755
3756
3757
3758
3759
3760
3761
3762
3763
3764
3765
3766
3767
3768
3769
3770
3771
3772
3773
3774
3775
3776
3777
3778
3779

Table 27: Zero-shot performance (%) of LLMs on the BBK with Wilson confidence intervals. Values in brackets indicate the Wilson CIs.

Model	Eng (%) [95% CI]	Hin (%) [95% CI]	Avg (%) [95% CI]
<i>< 4B Models</i>			
gemma-3-270m	26.64 [25.87, 27.41]	24.45 [22.88, 26.09]	26.24 [25.56, 26.95]
gemma-3-270m-it	27.44 [26.66, 28.22]	25.35 [23.76, 27.01]	27.06 [26.37, 27.77]
Param-1	33.10 [32.28, 33.92]	27.97 [26.32, 29.67]	32.18 [31.44, 32.92]
gemma-2-2b	41.24 [40.38, 42.10]	27.49 [25.86, 29.19]	38.78 [38.01, 39.55]
gemma-2-2b-it	35.94 [35.11, 36.78]	27.71 [26.07, 29.41]	34.47 [33.72, 35.22]
Llama-3.2-1B	29.71 [28.92, 30.51]	25.21 [23.62, 26.86]	28.91 [28.20, 29.63]
Llama-3.2-1B-Instruct	29.16 [28.37, 29.96]	26.33 [24.72, 28.01]	28.65 [27.94, 29.37]
Llama-3.2-3B	32.68 [31.87, 33.50]	28.69 [27.03, 30.41]	31.96 [31.23, 32.70]
Llama-3.2-3B-Instruct	40.59 [39.74, 41.45]	29.09 [27.42, 30.81]	38.53 [37.77, 39.30]
sarvam-2b-v0.5	28.14 [27.36, 28.93]	25.57 [23.98, 27.23]	27.68 [26.98, 28.39]
sarvam-1	30.82 [30.02, 31.63]	27.57 [25.93, 29.26]	30.24 [29.52, 30.97]
Nemotron-4-Mini-Hindi-4B-Base	36.67 [35.83, 37.51]	36.49 [34.71, 38.30]	36.64 [35.88, 37.40]
Nemotron-4-Mini-Hindi-4B-Instruct	35.83 [35.00, 36.67]	35.33 [33.57, 37.13]	35.74 [34.99, 36.50]
Qwen2.5-3B	44.57 [43.70, 45.44]	32.72 [30.99, 34.49]	42.45 [41.67, 43.23]
Qwen2.5-3B-Instruct	42.67 [41.81, 43.53]	27.20 [25.57, 28.90]	39.90 [39.13, 40.68]
granite-3.1-2b-instruct	37.71 [36.87, 38.56]	27.86 [26.21, 29.56]	35.95 [35.20, 36.71]
granite-3.1-3b-a800m-base	33.36 [32.55, 34.19]	26.70 [25.08, 28.38]	32.17 [31.44, 32.91]
<i>7B to 27B Models</i>			
Pangea-7B	47.16 [46.29, 48.03]	34.71 [32.96, 36.51]	44.93 [44.15, 45.72]
Indic-gemma-7b-finetuned-sft-Navarasa-2.0	42.31 [41.46, 43.18]	33.44 [31.71, 35.23]	40.73 [39.95, 41.51]
aya-23-8B	37.09 [36.25, 37.93]	33.22 [31.49, 35.00]	36.40 [35.64, 37.16]
Llama-3.1-8B	39.52 [38.68, 40.38]	31.41 [29.71, 33.17]	38.07 [37.31, 38.84]
Llama-3.1-8B-Instruct	47.14 [46.27, 48.01]	35.07 [33.31, 36.88]	44.98 [44.19, 45.77]
gemma-2-9b	55.23 [54.37, 56.10]	43.89 [42.05, 45.75]	53.20 [52.41, 53.99]
gemma-2-9b-it	48.92 [48.05, 49.80]	36.45 [34.68, 38.27]	46.69 [45.91, 47.48]
gpt-oss-20b	46.58 [45.71, 47.45]	36.27 [34.50, 38.08]	44.73 [43.95, 45.52]
gemma-2-27b	59.84 [58.98, 60.69]	50.38 [48.52, 52.25]	58.14 [57.36, 58.92]
gemma-2-27b-it	54.95 [54.08, 55.81]	41.24 [39.42, 43.09]	52.50 [51.71, 53.28]
<i>> 27B Models</i>			
gpt-oss-120b	71.40 [70.61, 72.18]	60.25 [58.41, 62.06]	69.41 [68.67, 70.13]
Qwen3-235B-A22B-Instruct-25076	74.57 [73.80, 75.32]	64.13 [62.32, 65.90]	72.70 [71.99, 73.39]
deepseek-v3	62.93 [62.09, 63.77]	45.01 [43.16, 46.88]	59.73 [58.95, 60.50]
gpt-4o	75.31 [74.55, 76.05]	65.18 [63.38, 66.94]	73.50 [72.79, 74.19]

Table 28: Zero-shot performance (%) of LLMs on the BBL with Wilson confidence intervals. Values in brackets indicate the Wilson CIs.

Model	Eng (%) [95% CI]	Hin (%) [95% CI]	Avg (%) [95% CI]
<i>< 4B Models</i>			
gemma-3-270m	25.49 [24.85, 26.15]	25.54 [24.55, 26.55]	25.51 [24.96, 26.06]
gemma-3-270m-it	25.56 [24.91, 26.22]	27.26 [26.25, 28.29]	26.07 [25.52, 26.63]
Param-1	36.15 [35.43, 36.88]	32.89 [31.82, 33.98]	35.17 [34.58, 35.78]
gemma-2-2b	38.45 [37.72, 39.18]	29.61 [28.58, 30.67]	35.79 [35.19, 36.40]
gemma-2-2b-it	34.49 [33.78, 35.21]	30.25 [29.21, 31.32]	33.22 [32.63, 33.81]
Llama-3.2-1B	29.63 [28.95, 30.32]	25.88 [24.89, 26.90]	28.50 [27.94, 29.07]
Llama-3.2-1B-Instruct	29.08 [28.40, 29.76]	27.04 [26.04, 28.07]	28.47 [27.90, 29.04]
Llama-3.2-3B	35.17 [34.45, 35.89]	28.53 [27.51, 29.58]	33.17 [32.59, 33.77]
Llama-3.2-3B-Instruct	39.74 [39.01, 40.48]	30.13 [29.09, 31.19]	36.86 [36.25, 37.46]
sarvam-2b-v0.5	28.49 [27.81, 29.17]	25.95 [24.96, 26.97]	27.72 [27.17, 28.29]
sarvam-1	30.92 [30.23, 31.62]	26.66 [25.66, 27.69]	29.64 [29.07, 30.22]
Nemotron-4-Mini-Hindi-4B-Base	40.75 [40.01, 41.49]	37.55 [36.45, 38.67]	39.79 [39.17, 40.40]
Nemotron-4-Mini-Hindi-4B-Instruct	36.99 [36.26, 37.71]	34.11 [33.03, 35.20]	36.12 [35.52, 36.73]
Qwen2.5-3B	44.98 [44.23, 45.72]	33.97 [32.89, 35.06]	41.67 [41.05, 42.29]
Qwen2.5-3B-Instruct	40.62 [39.88, 41.36]	29.89 [28.85, 30.94]	37.39 [36.79, 38.00]
granite-3.1-2b-instruct	38.18 [37.45, 38.91]	27.30 [26.29, 28.33]	34.91 [34.31, 35.51]
granite-3.1-3b-a800m-base	33.74 [33.04, 34.46]	24.01 [23.04, 25.00]	30.82 [30.24, 31.40]
<i>7B to 27B Models</i>			
Pangea-7B	48.70 [47.95, 49.45]	34.95 [33.87, 36.06]	44.57 [43.95, 45.20]
Indic-gemma-7b-finetuned-sft-Navarasa-2.0	44.08 [43.34, 44.83]	34.09 [33.02, 35.19]	41.08 [40.47, 41.70]
aya-23-8B	41.92 [41.18, 42.66]	33.01 [31.95, 34.10]	39.24 [38.63, 39.86]
Llama-3.1-8B	41.32 [40.58, 42.06]	31.76 [30.70, 32.83]	38.44 [37.84, 39.06]
Llama-3.1-8B-Instruct	48.61 [47.86, 49.36]	36.47 [35.38, 37.58]	44.96 [44.34, 45.59]
gemma-2-9b	58.49 [57.75, 59.23]	42.96 [41.83, 44.10]	53.83 [53.20, 54.45]
gemma-2-9b-it	45.05 [44.31, 45.80]	38.66 [37.55, 39.78]	43.13 [42.51, 43.75]
gpt-oss-20b	40.69 [39.96, 41.43]	35.24 [34.16, 36.34]	39.06 [38.45, 39.67]
gemma-2-27b	64.91 [64.19, 65.63]	51.83 [50.69, 52.97]	60.99 [60.37, 61.60]
gemma-2-27b-it	50.71 [49.96, 51.46]	42.02 [40.89, 43.15]	48.10 [47.47, 48.73]
<i>> 27B Models</i>			
gpt-oss-120b	70.72 [70.03, 71.40]	62.94 [61.83, 64.04]	68.38 [67.80, 68.97]
Qwen3-235B-A22B-Instruct-25076	80.15 [79.55, 80.75]	68.60 [67.53, 69.65]	76.68 [76.15, 77.21]
deepseek-v3	67.78 [67.07, 68.47]	46.78 [45.63, 47.92]	61.47 [60.86, 62.08]
gpt-4o	78.83 [78.22, 79.44]	71.02 [69.97, 72.04]	76.49 [75.95, 77.02]

Table 29: Pairwise comparison of the top 5 LLMs on the BBA domain using McNemar’s test. Accuracy Diff (%) shows the absolute difference between model accuracies. Significant differences are reported at $p < 0.05$.

Model A	Accuracy A	Model B	Accuracy B	Accuracy Diff (%)	McNemar Stat	p-value	Significant
English							
gpt-4o	62.75%	Qwen3-235B-A22B-Instruct-25076	60.25%	2.50%	23.874	1.029e-06	Yes
gpt-oss-120b	55.62%	gemma-2-27b	50.70%	4.92%	45.308	1.684e-11	Yes
gpt-oss-120b	55.62%	deepseek-v3	51.38%	4.24%	54.175	1.834e-13	Yes
Qwen3-235B-A22B-Instruct-25076	60.25%	gemma-2-27b	50.70%	9.55%	171.369	3.718e-39	Yes
Qwen3-235B-A22B-Instruct-25076	60.25%	deepseek-v3	51.38%	8.87%	228.148	1.511e-51	Yes
deepseek-v3	51.38%	gemma-2-27b	50.70%	0.68%	0.855	0.3552	No
gpt-4o	62.75%	gemma-2-27b	50.70%	12.06%	269.474	1.476e-60	Yes
gpt-4o	62.75%	deepseek-v3	51.38%	11.37%	357.026	1.251e-79	Yes
gpt-4o	62.75%	gpt-oss-120b	55.62%	7.14%	173.739	1.129e-39	Yes
Qwen3-235B-A22B-Instruct-25076	60.25%	gpt-oss-120b	55.62%	4.63%	68.941	1.014e-16	Yes
Hindi							
gpt-oss-120b	48.05%	param-1	38.04%	10.01%	126.394	2.521e-29	Yes
Qwen3-235B-A22B-Instruct-25076	54.78%	gpt-4o	54.73%	0.05%	0.003	0.9591	No
Qwen3-235B-A22B-Instruct-25076	54.78%	gpt-oss-120b	48.05%	6.73%	72.961	1.322e-17	Yes
Qwen3-235B-A22B-Instruct-25076	54.78%	gemma-2-27b	42.26%	12.52%	174.197	8.964e-40	Yes
Qwen3-235B-A22B-Instruct-25076	54.78%	param-1	38.04%	16.74%	367.997	5.109e-82	Yes
gpt-4o	54.73%	gpt-oss-120b	48.05%	6.68%	75.405	3.834e-18	Yes
gpt-4o	54.73%	gemma-2-27b	42.26%	12.47%	175.629	4.363e-40	Yes
gpt-4o	54.73%	param-1	38.04%	16.69%	367.336	7.117e-82	Yes
gpt-oss-120b	48.05%	gemma-2-27b	42.26%	5.79%	37.345	9.899e-10	Yes
gemma-2-27b	42.26%	param-1	38.04%	4.22%	20.868	4.921e-06	Yes

Table 30: Pairwise comparison of the top 5 LLMs on the BBF domain using McNemar’s test. Accuracy Diff (%) shows the absolute difference between model accuracies. Significant differences are reported at $p < 0.05$.

Model A	Accuracy A	Model B	Accuracy B	Accuracy Diff (%)	McNemar Stat	p-value	Significant
English							
gpt-oss-120b	74.11%	Qwen3-235B-A22B-Instruct-25076	63.72%	10.39%	471.859	1.262e-104	Yes
deepseek-v3	63.46%	gemma-2-27b	47.79%	15.67%	659.063	2.387e-145	Yes
deepseek-v3	63.46%	gpt-4o	57.27%	6.19%	174.76	6.755e-40	Yes
Qwen3-235B-A22B-Instruct-25076	63.72%	gemma-2-27b	47.79%	15.93%	668.926	1.711e-147	Yes
Qwen3-235B-A22B-Instruct-25076	63.72%	gpt-4o	57.27%	6.45%	199.599	2.555e-45	Yes
gpt-4o	57.27%	gemma-2-27b	47.79%	9.48%	239.639	4.714e-54	Yes
gpt-oss-120b	74.11%	gemma-2-27b	47.79%	26.33%	1818.018	0.0	Yes
gpt-oss-120b	74.11%	gpt-4o	57.27%	16.85%	1025.635	4.811e-225	Yes
gpt-oss-120b	74.11%	deepseek-v3	63.46%	10.65%	497.362	3.564e-110	Yes
Qwen3-235B-A22B-Instruct-25076	63.72%	deepseek-v3	63.46%	0.26%	0.35	0.5544	No
Hindi							
Qwen3-235B-A22B-Instruct-25076	56.27%	gemma-2-27b	41.24%	15.03%	257.555	5.855e-58	Yes
gpt-oss-120b	64.16%	deepseek-v3	57.04%	7.12%	94.667	2.251e-22	Yes
gpt-oss-120b	64.16%	Qwen3-235B-A22B-Instruct-25076	56.27%	7.89%	109.714	1.132e-25	Yes
gpt-oss-120b	64.16%	gpt-4o	49.82%	14.34%	315.756	1.217e-70	Yes
gpt-oss-120b	64.16%	gemma-2-27b	41.24%	22.92%	591.15	1.409e-130	Yes
deepseek-v3	57.04%	Qwen3-235B-A22B-Instruct-25076	56.27%	0.77%	1.259	0.2618	No
deepseek-v3	57.04%	gpt-4o	49.82%	7.22%	95.262	1.667e-22	Yes
deepseek-v3	57.04%	gemma-2-27b	41.24%	15.80%	288.861	8.805e-65	Yes
Qwen3-235B-A22B-Instruct-25076	56.27%	gpt-4o	49.82%	6.45%	81.802	1.504e-19	Yes
gpt-4o	49.82%	gemma-2-27b	41.24%	8.58%	86.602	1.327e-20	Yes

Table 31: Pairwise comparison of the top 5 LLMs on the BBK domain using McNemar’s test. Accuracy Diff (%) shows the absolute difference between model accuracies. Significant differences are reported at $p < 0.05$.

Model A	Accuracy A	Model B	Accuracy B	Accuracy Diff (%)	McNemar Stat	p-value	Significant
English							
gpt-4o	75.31%	Qwen3-235B-A22B-Instruct-25076	74.57%	0.74%	4.178	0.04095	Yes
gpt-oss-120b	71.40%	gemma-2-27b	59.84%	11.57%	572.581	1.541e-126	Yes
gpt-oss-120b	71.40%	deepseek-v3	62.93%	8.47%	331.375	4.823e-74	Yes
Qwen3-235B-A22B-Instruct-25076	74.57%	gemma-2-27b	59.84%	14.73%	953.27	2.581e-209	Yes
Qwen3-235B-A22B-Instruct-25076	74.57%	deepseek-v3	62.93%	11.63%	642.934	7.69e-142	Yes
deepseek-v3	62.93%	gemma-2-27b	59.84%	3.10%	42.443	7.276e-11	Yes
gpt-4o	75.31%	gemma-2-27b	59.84%	15.47%	1045.052	2.895e-229	Yes
gpt-4o	75.31%	deepseek-v3	62.93%	12.37%	732.145	3.061e-161	Yes
gpt-4o	75.31%	gpt-oss-120b	71.40%	3.91%	97.61	5.0939e-23	Yes
Qwen3-235B-A22B-Instruct-25076	74.57%	gpt-oss-120b	71.40%	3.16%	64.663	8.886e-16	Yes
Hindi							
gpt-oss-120b	60.25%	deepseek-v3	45.01%	15.23%	184.026	6.403e-42	Yes
gpt-4o	65.18%	Qwen3-235B-A22B-Instruct-25076	64.13%	1.05%	1.658	0.1979	No
gpt-4o	65.18%	gpt-oss-120b	60.25%	4.93%	29.877	4.603e-08	Yes
gpt-4o	65.18%	gemma-2-27b	50.38%	14.80%	196.267	1.363e-44	Yes
gpt-4o	65.18%	deepseek-v3	45.01%	20.17%	316.898	6.864e-71	Yes
Qwen3-235B-A22B-Instruct-25076	64.13%	gpt-oss-120b	60.25%	3.88%	17.639	2.671e-05	Yes
Qwen3-235B-A22B-Instruct-25076	64.13%	gemma-2-27b	50.38%	13.75%	169.093	1.167e-38	Yes
Qwen3-235B-A22B-Instruct-25076	64.13%	deepseek-v3	45.01%	19.11%	290.321	4.232e-65	Yes
gpt-oss-120b	60.25%	gemma-2-27b	50.38%	9.87%	88.483	5.127e-21	Yes
gemma-2-27b	50.38%	deepseek-v3	45.01%	5.37%	22.988	1.63e-06	Yes

Table 32: Pairwise comparison of the top 5 LLMs on the BBL domain using McNemar’s test. Accuracy Diff (%) shows the absolute difference between model accuracies. Significant differences are reported at $p < 0.05$.

Model A	Accuracy A	Model B	Accuracy B	Accuracy Diff (%)	McNemar Stat	p-value	Significant
English							
Qwen3-235B-A22B-Instruct-25076	80.15%	gpt-4o	78.83%	1.32%	17.392	3.041e-05	Yes
gpt-oss-120b	70.72%	gemma-2-27b	64.91%	5.81%	131.893	1.579e-30	Yes
gpt-oss-120b	70.72%	deepseek-v3	67.78%	2.94%	53.656	2.389e-13	Yes
gpt-4o	78.83%	gemma-2-27b	64.91%	13.92%	803.654	8.661e-177	Yes
gpt-4o	78.83%	deepseek-v3	67.78%	11.06%	786.845	3.911e-173	Yes
deepseek-v3	67.78%	gemma-2-27b	64.91%	2.86%	30.906	2.709e-08	Yes
Qwen3-235B-A22B-Instruct-25076	80.15%	gemma-2-27b	64.91%	15.24%	965.971	4.478e-212	Yes
Qwen3-235B-A22B-Instruct-25076	80.15%	deepseek-v3	67.78%	12.38%	995.052	2.137e-218	Yes
Qwen3-235B-A22B-Instruct-25076	80.15%	gpt-oss-120b	70.72%	9.43%	643.04	7.291e-142	Yes
gpt-4o	78.83%	gpt-oss-120b	70.72%	8.11%	482.426	6.337e-107	Yes
Hindi							
gpt-oss-120b	62.94%	deepseek-v3	46.78%	16.17%	475.374	2.16e-105	Yes
gpt-4o	71.02%	Qwen3-235B-A22B-Instruct-25076	68.60%	2.42%	19.396	1.062e-05	Yes
gpt-4o	71.02%	gpt-oss-120b	62.94%	8.08%	176.253	3.188e-40	Yes
gpt-4o	71.02%	gemma-2-27b	51.83%	19.19%	543.759	2.868e-120	Yes
gpt-4o	71.02%	deepseek-v3	46.78%	24.24%	1027.969	1.496e-225	Yes
Qwen3-235B-A22B-Instruct-25076	68.60%	gpt-oss-120b	62.94%	5.66%	82.242	1.204e-19	Yes
Qwen3-235B-A22B-Instruct-25076	68.60%	gemma-2-27b	51.83%	16.77%	423.76	3.707e-94	Yes
Qwen3-235B-A22B-Instruct-25076	68.60%	deepseek-v3	46.78%	21.82%	859.675	5.732e-189	Yes
gpt-oss-120b	62.94%	gemma-2-27b	51.83%	11.11%	181.188	2.666e-41	Yes
gemma-2-27b	51.83%	deepseek-v3	46.78%	5.06%	36.78	1.322e-09	Yes