

INDICPARAM: Benchmark to evaluate LLMs on low-resource Indic Languages

Ayush Maheshwari, Kaushal Sharma[♣], Vivek Patel[♣], Aditya Maheshwari[♣]

[♣]Indian Institute of Management Indore, India ; [♣]BHARATGEN

ayush.hakmn@gmail.com, {kaushals,vivekp,adityam}@iimidr.ac.in

Abstract

While large language models excel on high-resource multilingual tasks, low- and extremely low-resource Indic languages remain severely under-evaluated. We present INDICPARAM, a human-curated benchmark of over 13,000 multiple-choice questions covering 11 such languages (Nepali, Gujarati, Marathi, Odia as low-resource; Dogri, Maithili, Rajasthani, Sanskrit, Bodo, Santali, Konkani as extremely low-resource) plus SanskritEnglish code-mixing. We evaluated 19 LLMs, both proprietary and open-weights, which reveals that even the top-performing GPT-5 reaches only 45% average accuracy, followed by DeepSeek-3.2 (43.1) and Claude-4.5 (42.7). We additionally label each question as knowledge-oriented or purely linguistic to discriminate factual recall from grammatical proficiency. Further, we assess the ability of LLMs to handle diverse question formats-such as list-based matching, assertion-reason pairs, and sequence ordering-alongside conventional multiple-choice questions. INDICPARAM provides insights into limitations of cross-lingual transfer and establishes a challenging benchmark for Indic languages. The dataset is available at <https://huggingface.co/datasets/bharatgenai/IndicParam>¹.

1 Introduction

Large language models (LLMs) have delivered state-of-the-art results across a range of multilingual tasks, particularly in high- and medium-resource settings such as translation, named entity recognition, and question answering. However, systematic evaluation for low and extremely-low-resource Indic languages remains limited. Existing benchmarks for Indic languages have significant coverage for major languages (*e.g.* Hindi,

Language	#Words (in M)	#Speakers (in M)
Nepali	1642.9	19.4
Marathi	1541.2	99.1
Gujarati	934.1	60.3
Odia	333.8	42.6
Sanskrit	16.9	3.1
Maithili	14.0	14.3
Konkani	2.8	2.6
Santali	0.8	7.7
Bodo	0.8	1.5
Dogri	0.6	2.8
Rajasthani	-	25.8
Total	4480	279.1

Table 1: Number of words (in millions) present in FineWeb2 corpus and number of speakers (in millions) for each language considered in INDICPARAM.

Tamil, Telugu) (Singh et al., 2024; Kakwani et al., 2020), yet they provide only partial or no support for several understudied languages. Additionally, they rarely include coverage for code-mixed usage and evaluation.

India, with a population of approximately 1.4 billion, is home to over 120 languages, at least 30 of which have one million or more speakers (Office of the Registrar General & Census Commissioner, India, 2018). These languages span multiple families, for instance, Marathi and Hindi belong to the Indo-European branch; Kannada and Telugu are Dravidian; Santali is Austro-Asiatic; and Bodo is Sino-Tibetan. The Indian Constitution recognizes 22 languages in the Eighth Schedule, yet most existing benchmarks concentrate on 11 major languages that together account for more than 93% of speakers, largely reflecting the limited availability of web-scale data for others (IAMAI, 2024). With recent advances in multilingual large language models, performance on language understanding has improved for many of these widely used languages (Hu et al., 2020; Qin et al., 2025).

¹Scripts to run benchmark are present at <https://github.com/ayushbits/IndicParam>

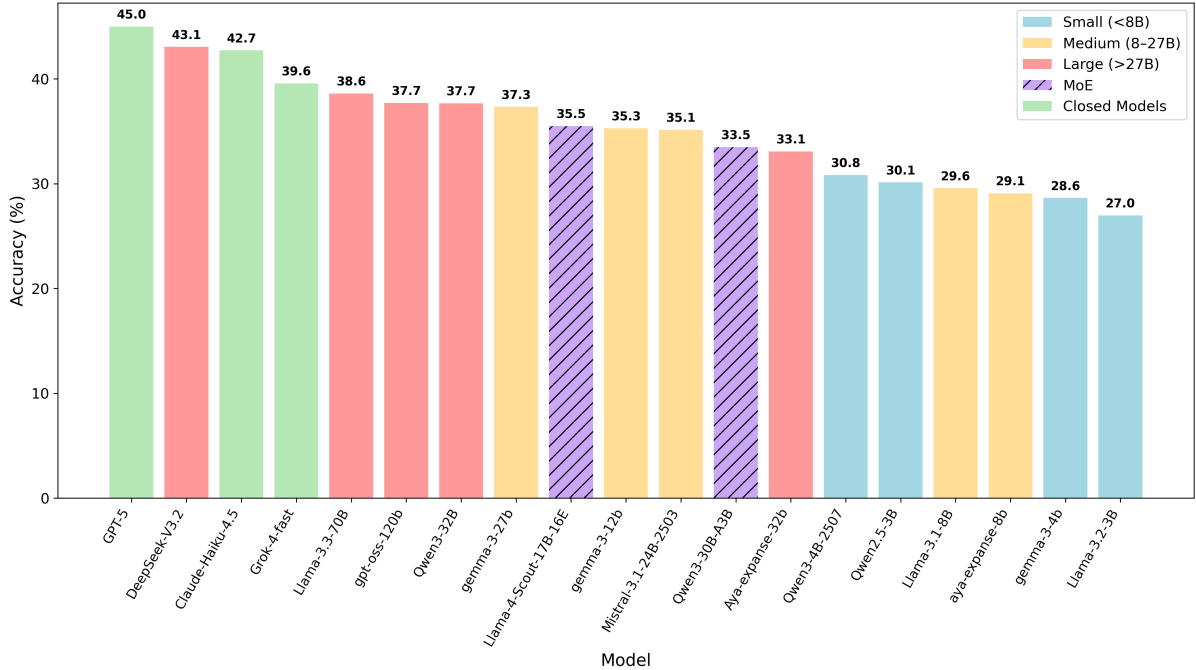


Figure 1: Average performance for all languages on INDICPARAM. Notably, all closed models outperform best performing open-source models.

However, despite a few benchmarks covering the majority set of 11 languages, there exists almost no evaluation resources for other Indic languages (Kumar et al., 2022; Madasu et al., 2023).

In Table 1, we present word counts for the INDICPARAM languages within FineWeb2 (Penedo et al., 2025). FineWeb2 is a cleaned pretraining corpus comprising of approximately 3 trillion words from around 1800 languages (excluding English) gathered from more than 100 Common Crawl snapshots. For the 11 languages analyzed in INDICPARAM, the FineWeb2 corpus contains 4.49 billion words. Four languages-Nepali, Marathi, Gujarati, and Odia-account for 4.45 billion of these words, leaving only 36 million words for the remaining six languages. Collectively, these ten languages represent about 280 million speakers across the Indian subcontinent, yet their aggregate web presence is disproportionately small relative to their speaker base. Based on the availability of web corpus, we classify Marathi, Gujarati, Nepali and Odia as low-resource languages while remaining seven languages as extremely low-resource languages.

To address these limitations, this paper introduces INDICPARAM, a human-supervised evaluation set comprising more than 13,000 questions for 11 low- and extremely low-resource Indic languages. We also include a code-mixed Sanskrit-

English language set evaluation. We hypothesize that despite limited pretraining data in several of these languages, LLMs can exploit cross-lingual transfer among typologically and script-wise related languages to improve performance on these target languages. To separate language understanding from factual recall, each question is labeled as either linguistic (eg, morphology, syntax, semantics, discourse) or knowledge-based (world facts and entities).

Low resource (4): Nepali, Gujarati, Marathi, Odia

Extremely-low resource (7): Dogri, Maithili, Rajasthani, Sanskrit, Bodo, Santali, Konkani

We evaluate 19 LLMs spanning open- and closed-source systems and a range of parameter scales and architectures on INDICPARAM. The evaluated models includes closed models such as GPT-5, Claude-Haiku-4.5, and open weight models such as Llama4-Scout, DeepSeek-V3.2 among others. In Figure 1, we present the average performance of all models. Larger models with substantial exposure to multilingual Indic corpora tend to outperform smaller models. However, no model exceeds a 50% average score on INDICPARAM. Among closed systems, GPT-5 achieves the highest overall accuracy at 45, followed by Claude-

Haiku at 42.7 and Grok-4 at 39.6%, indicating that model robustness on studied languages remains limited even among state-of-the-art proprietary models. Open-weight models narrow the gap with DeepSeek-V3.2 (DeepSeek-AI, 2025a), which is a 685B parameter model, achieving 43.1% trailing GPT-5 by around 2%, while Llama 3.370B attains average accuracy of 39.6%. Other large and medium-size models cluster below these values, reflecting the benefits of larger capacity and Indic exposure. Overall, results indicate that multilingual performance for Indic languages remains challenging across both closed and open models, with substantial room for improvement in cross-lingual generalization. Our contributions can be summarised as follows:

1. We present INDICPARAM, a benchmark consisting of 13K+ questions on 11 low- and extremely-low resource Indic languages, as well as a code-mixed English-Sanskrit variant.
2. We provide granular annotations for each instance in the benchmark, including question-type and question-category labels, to enable systematic analysis of LLM capabilities.
3. We conduct an extensive evaluation encompassing 19 LLMs, including both proprietary closed-source and open-weight models.

2 Related Work

LLMs have demonstrated substantial improvements in language understanding performance for several Indic languages (Kakwani et al., 2020; Verma et al., 2025; Maji et al., 2025a; Singh et al., 2024). Multiple Indic multilingual benchmarks have been introduced to evaluate such progress, primarily targeting major mid- and high-resource languages (Singh et al., 2025, 2024; Maheshwari et al., 2025). Among these, IndicGenBench (Singh et al., 2024) proposes a benchmark on 29 Indic languages, including those examined in this study. Its evaluation suite includes tasks such as reading comprehension and translation, which are primarily adapted from English benchmark datasets like XQUAD (Artetxe et al., 2020) and FLORES (Goyal et al., 2022). IndicGenBench QA tasks restrict answers to spans supported by the given passage. In contrast, INDICPARAM extends this evaluation paradigm by assessing both the language understanding and the domain-specific general knowledge of LLMs using questions drawn from graduate-level examinations.

IndicGLUE (Kakwani et al., 2020) is a suite of supervised NLP tasks for prominent Indic languages, with a particular emphasis on Hindi, including classification tasks such as news categorization and headline prediction, as well as natural language inference, among others. IndicXTREME (Doddapaneni et al., 2023) broadens task coverage to named entity recognition, question answering, and paraphrase detection, and expands language coverage by translating and curating benchmarks in 11 Indic languages; however, many long-tail languages remain outside its scope. Recent Indic QA benchmarks further diversify evaluation: ParamBench (Maheshwari et al., 2025) focuses on graduate-level questionnaires with varied question types in Hindi, while MILU (Verma et al., 2025) contains 79K questions across 11 Indian languages, including Odia, Gujarati, and Marathi, which are also covered in INDICPARAM. Sanskriti (Maji et al., 2025a) and Pariksha (Watts et al., 2024) are English-language benchmarks designed to assess the socio-cultural alignment of LLMs. Drishtikon (Maji et al., 2025b) further extends this line of work to multimodal cultural understanding across 15 Indic languages.

Despite this progress, most benchmarks concentrate on the majority of 11 widely used Indian languages, with very limited coverage for extremely low-resource languages. In this paper, we introduce graduate-level language understanding questionnaires in 4 low-resource and 7 extremely low-resource languages, enabling targeted assessment of both language competence and general knowledge for underrepresented Indic languages.

3 INDICPARAM

INDICPARAM is a high-quality QA benchmark consisting of 13,207 questions that evaluate language understanding and general knowledge of multilingual LLMs in 11 Indic languages. These languages are spoken by approximately 280 million people across the Indian subcontinent, of whom 10 languages are recognized as scheduled languages in the Indian Constitution. Rajasthani, though not a scheduled language, is spoken by roughly 25 million people. We also include a QA set of code-mixed version of Sanskrit where questions include a mix of English and Sanskrit.

The language-wise distribution of questions is shown in Table 2. The questions are collected from

Language	#Ques	LU(%)	Script	Code
Nepali	1038	18.8	Deva	npi
Marathi	1245	4.7	Deva	mar
Gujarati	1044	0.6	Gujarati	guj
Odia	577	20.3	Orya	ory
Maithili	1286	10.1	Deva	mai
Konkani	1328	2.5	Deva	gom
Santali	873	11.3	Olck	sat
Bodo	1313	-	Deva	brx
Dogri	1027	18.3	Deva	doi
Rajasthani	1190	27.8	Deva	-
Sanskrit	1315	20.7	Deva	san
Sans-Eng	971	11.5	-	-
Total	13207			

Table 2: Distribution of question-answer pairs for different languages in INDICPARAM. ‘Sans-Eng’ denotes a separate set of Sanskrit-English code-mixed question-answer pairs which forms a separate set. LU refers to % of manually classified language-understanding questions; rest are categorized as knowledge-related (*c.f.* Section 3.1).

UGC-NET² which is a nationwide examination administered by a government agency to determine eligibility for PhD admission and for appointment to teaching positions in Indian universities and colleges. The exam is offered in around 85 subjects and is conducted twice annually. Each test consists of two papers composed of multiple-choice questions (MCQs). We selected language specific question papers for each of the languages in INDICPARAM. Further, we developed the dataset by downloading official question papers and answer keys from the official examination website.

Each language comprises multiple PDF question papers, many machine-readable and a subset with non-selectable text. Document layouts vary, with some single-column and most two-column. To ensure uniform accessibility, all PDFs were processed with a proprietary OCR system, and the extracted text was used for downstream curation and annotation. To the best of our knowledge, this corpus has not appeared in prior LLM benchmarking studies and constitutes a newly curated, human-authored dataset explicitly designed for graduate-level evaluation in Indic contexts. We describe annotation setup and team structure in Section 9 in

²<https://ugcnet.nta.ac.in>

Appendix.

3.1 Question Classification

In addition to annotating question-answer pairs in INDICPARAM, each question is assigned to one of two categories: (a) **language understanding** (LU) that includes questions related to linguistics and grammar, and (b) **general knowledge** (GK) targeting fact-based queries. This enables systematic assessment of multilingual LLMs across both language understanding and factual recall capabilities. Given the low- and extremely low-resource nature of the studied languages, this evaluation investigates model performance under data scarcity conditions and examines the extent to which cross-lingual transfer mechanisms can compensate for limited in-language training data.

The distribution of questions across these categories for each language is presented in Table 2. The LU column reports, for each language, the proportion of questions that test language understanding while the remaining questions belonging to the GK category and focusing on factual knowledge. The share of LU questions varies substantially, from under 1% in Gujarati to nearly 28% in Rajasthani, with several other languages such as Odia and Sanskrit also exhibiting over 20% LU coverage.

3.2 Question Type Classification

Following ParamBench (Maheshwari et al., 2025), we annotate each question-answer pair with its corresponding question type. We classify them into six primary categories that capture the diversity of assessment formats (see Table 3): multiple-choice questions (MCQ), assertion and reasoning (A&R), list matching, blank filling, identify incorrect statement (IS), and ordering. The majority (73%) of the questions are MCQs followed by list matching (9%), ordering (6.5%), A&R (6.1%), IS (4.1%) and blank filling (1.1%). The language wise breakdown of all languages remains the same except for Bodo, Gujarati, and Dogri, where the proportion of MCQs remain 36.6%, 33% and 21.7%, respectively, which are significantly lower than other languages in the benchmark. The detailed statistics of language-wise splits of question types is present in Table 6. In Table 11, we present example questions from different languages and their respective question class and question type.

Question Type	#Questions
Multiple-choice	9653
Assertion & Reason	811
List Matching	1185
Fill in the blanks	157
Identify incorrect statement	545
Ordering	856
Total	13207

Table 3: Distribution of questions across all languages in INDICPARAM.

4 Experiments

We evaluated 19 state-of-the-art models spanning both open-weight and proprietary LLMs. The open-weight models range from 3 billion to over 685 billion parameters, including both dense and mixture-of-experts (MoE) architectures. These models collectively represent the state of the art in multilingual reasoning, code understanding, and multimodal generation. All models were used in their instruction-tuned configurations and evaluated under a zero-shot prompting setup. The prompt used during evaluation is present in Table 8. Open-weights models are retrieved from publicly available checkpoints at HuggingFace and use transformers library for inference. We load parameter weights in bf16 precision for all open-weight models except Llama4-Scout which was loaded with 8-bit quantization. We use greedy decoding while generating predictions by setting max-tokens as 50, temperature as 0, do_sample flag to false and batch size is set to 16. We disable thinking mode for Qwen3-MoE and keep reasoning levels at low for GPT-OSS-120B. For Gemma and Mistral, we follow their respective strategies from response generation instead of transformers pipeline function. We use OpenRouter³ to generate predictions for closed models. We next describe the models included in our evaluation.

1. **Qwen Series (Team-Qwen3, 2025):** We evaluated multiple models of the Qwen family, including Qwen2.5-3B, Qwen3-4B-2507, Qwen3-30B-A3B, and Qwen3-32B. The Qwen 3 series integrates substantial architectural and pre-training improvements over the 2.5 generation, trained on roughly 36 trillion multilingual tokens spanning 119 languages. The Qwen3-30B-A3B

variant adopts a mixture-of-experts (MoE) design with 3 billion active parameters per forward pass, providing strong efficiency while maintaining performance parity with larger dense models.

2. **Llama Series (Team, 2024):** We evaluated Metas Llama-3.2-3B, Llama-3.1-8B, and Llama-3.3-70B. These models are trained on a mixture of high-quality multilingual and code data totaling approximately 15 trillion tokens. We also evaluated on Llama-4-Scout-17B-16E which is an MoE model having a 17 billion active parameters with 16 experts with a total of 109 billion model parameters. This model is trained on 40T tokens specifically designed to improve multilingual text and images.

3. **Gemma Series (Team, 2025):** We evaluated Gemma-3 models at 4B, 12B, and 27B parameter scales. Developed by Google DeepMind, the Gemma 3 family employs hybrid instruction-fine-tuning strategies and trained on diverse multilingual datasets. These models have demonstrated strong alignment and competitive performance across multilingual and reasoning benchmarks.

4. **Aya Expanse (Aryabumi et al., 2024):** The Aya-Expanse-8B and Aya-Expanse-32B models are instruction-tuned multilingual LLMs. These models prioritize cultural and linguistic diversity and are trained on extensive parallel datasets that include over 100 different languages.

5. **Mistral 3.1 (MistralAI):** We evaluated the Mistral-Small-3.1-24B-2503 which is a 24 billion parameter multi-lingual model achieving competitive performance while maintaining fast inference throughput.

6. **DeepSeek V3.2 (DeepSeek-AI, 2025b):** We evaluated DeepSeek-V3.2-Experimental model consisting of 685B parameters. It employs specialized routing and sparsity mechanisms to reduce inference costs while preserving high-quality reasoning, code synthesis, and mathematical reasoning performance.

7. **GPT-OSS-120B (OpenAI et al., 2025):** We examined the GPT-OSS-120B model, an open-weight MoE model by OpenAI, trained on trillions of tokens with a focus on STEM, coding, and general knowledge. It has a total of 116.8B parameters and 5.1B active parameters per token per forward pass.

8. **Frontier Proprietary Models:** For completeness, we also benchmarked several leading closed-source systems accessible through public APIs: GPT-5, Claude-Haiku-4.5, and Grok-4-Fast.

³<https://openrouter.ai/api/v1/chat/completions>

Size	Model	Dogri	Maithili	Rajasthani	Sanskrit	Sans-Eng	Bodo	Santali	Konkani	Average
Small	Qwen2.5-3B	32.1	28.8	28.6	31.9	<u>34.8</u>	<u>28.3</u>	33.1	<u>31.4</u>	30.9
	Llama-3.2-3B	27.9	21.9	26.3	28	<u>27.9</u>	<u>27.8</u>	27.4	28.8	26.9
	Gemma-3-4b	30.5	24.5	27.6	29.1	32.2	25.8	31.7	29.1	28.5
	Qwen3-4B	<u>34.6</u>	<u>31.6</u>	<u>30</u>	<u>33.8</u>	33.6	25.1	<u>33.8</u>	29.6	<u>31.2</u>
Medium	Aya-8b	31.5	29.9	30.8	29.4	29.9	28.1	30.8	26.9	29.5
	Llama-3.1-8B	26.4	29	31.8	32.2	34	26.8	30.6	29.5	30.0
	Gemma-3-12b	37.2	31.6	35	37.1	41.8	<u>30.6</u>	33.7	35.7	35.1
	Mistral3.1-24b	<u>38</u>	<u>37.4</u>	34.8	<u>40</u>	44.1	<u>30.1</u>	35.7	32	36.2
	Gemma-3-27b	37.6	34.4	<u>39.8</u>	38.6	<u>49.1</u>	28.9	<u>36.7</u>	<u>36.5</u>	<u>37.3</u>
Large	Aya-32b	35.2	33.2	34.6	35.2	39	27.2	35.6	33.1	33.7
	Qwen3-32B	<u>43.4</u>	34	33.8	39.2	48.6	<u>31.8</u>	39.7	34.5	33.9
	Llama-3.3-70B	37.3	<u>37.5</u>	<u>38.2</u>	41.7	51.6	<u>28.6</u>	40.1	37.3	37.6
	DeepSeek-3.2	43.3	<u>41.4</u>	<u>40.1</u>	<u>51.3</u>	<u>61.4</u>	36.4	41.5	<u>37.6</u>	<u>43.7</u>
MoE	Qwen3-A3B	35.7	30.9	31.6	36.7	44.2	<u>29.9</u>	35.9	28.3	38.6
	Llama-4-Scout	30.5	31.3	34.8	<u>44</u>	<u>48.7</u>	<u>25.9</u>	<u>36.3</u>	34.6	35.4
	gpt-oss-120b	<u>39.7</u>	<u>35.6</u>	<u>39.4</u>	40.2	48.1	28.3	35.6	<u>35.4</u>	37.4
Closed	GPT-5	45.7	43.4	44.6	54.8	64.6	<u>31.6</u>	40.4	<u>39</u>	45.1
	Grok-4-fast	40.2	38.1	39.1	42.9	54.2	30.6	38.8	37.4	39.7
	Claude-4.5	<u>44.8</u>	<u>42.4</u>	<u>44.1</u>	47.8	58.1	28.1	<u>40.4</u>	42.1	43.0

Table 4: Performance on extremely-low resource languages on INDICPARAM. We **bold** and *italicize* the best overall performance, underline and *italicize* the best performance in each model-size category and underline the second best overall performance.

- *GPT-5* (OpenAI) represents the current flagship frontier model, incorporating multimodal reasoning and tool-use capabilities.
- *Claude-Haiku-4.5* (Anthropic) is a lightweight yet high-accuracy variant of the Claude 4.5 family designed for efficiency-critical environments.
- *Grok-4-Fast* (xAI) is optimized for rapid inference and contextual awareness within streaming conversational environments.

5 Results and discussion

5.1 Overall Model Performance

The evaluation results presented in Figure 1 demonstrate that average performance on INDICPARAM for 19 evaluated models remains moderate, though several interesting patterns emerge when grouped by model scale. Among the small-capacity models (<8B parameters), Qwen3-4B-2507 achieves the highest average accuracy of 30.8%, followed closely by Qwen2.5-3B (30.1). Other models in this category, including Gemma-3-4B, and Llama-3.2-3B, remain around the 27-30% range. This pattern suggests that within low-parameter models, differences in linguistic coverage and pre-training diversity contribute more to Indic language performance than the raw number of parameters.

In the medium-size category (8B-27B),

performance improves considerably, with Gemma-3-27B leading at 37.3%, followed closely by Gemma-3-12B (35.3) and Mistral-Small-3.1-24B (35.1). Other models like Aya-Expanse-8B and Llama-3.1-8B score less than 30%. These findings underscore that model scale beyond 8B yields measurable gains, but pre-training coverage and multilingual representation continue to play a decisive role in performance.

Among the larger and frontier models ($\geq 27B$), overall accuracy rises substantially. The best performance is achieved by GPT-5, which attains an average of 45.0%, followed by DeepSeek-V3.2 (43.1), Claude-Haiku-4.5 (42.7), and Grok-4-Fast (39.6). Open-weight models such as Llama-3.3-70B and GPT-OSS-120B achieve 38.6 and 37.7, respectively, while the Llama-4-Scout-17B-16E models remain slightly lower at 35.5. These results indicate that while data and parameter size enhance overall accuracy, the highest-performing models also benefit from closed fine-tuning cycles and refined alignment procedures on multilingual or Indic-rich corpora.

Overall, the evaluation reveals a clear upward trajectory in accuracy with increasing model scale, yet even the most capable frontier systems such

Size	Model	Nepali	Gujarati	Marathi	Odia	Average
Small	Qwen2.5-3B	27.1	<u>29.5</u>	27.9	28.8	28.2
	Llama3.2-3B	26.2	26.1	28.6	27	27.1
	Gemma-3-4b	<u>29</u>	26.2	29.9	31.5	28.9
	Qwen3-4B	28.7	27.4	<u>30.5</u>	<u>34.7</u>	<u>29.8</u>
Medium	Aya-8b	28.2	27.7	28.3	28.1	28.1
	Llama3.1-8B	25.5	27.2	32.2	29.3	28.7
	Gemma-3-12b	34.8	<u>33.8</u>	38.4	34.5	35.6
	Mistral3.1-24b	34.9	32.6	32.9	27.6	32.6
	Gemma3-27b	<u>39.3</u>	31.1	<u>39.9</u>	<u>39.9</u>	<u>37.4</u>
Large	Aya-32b	31.6	30.7	35.5	34.5	33.0
	Qwen3-32B	33	28.4	33.2	28.9	31.2
	Llama3.3-70B	35	35.9	40.8	40.7	37.9
	DeepSeek-3.2	<u>39.7</u>	<u>39.5</u>	<u>44.6</u>	<u>42.1</u>	<u>41.6</u>
MoE	Qwen3-A3B	<u>37.2</u>	33.5	43.5	<u>40.4</u>	<u>38.7</u>
	Llama-4-Scout	34.6	27.3	<u>43.5</u>	35	35.5
	gpt-oss-120b	37.1	<u>37</u>	39.7	39.9	38.3
Closed	GPT-5	<u>43.4</u>	40	48.4	48.7	44.9
	Grok4-fast	40.9	35.9	39.8	40.7	39.2
	Claude-4.5	43.8	37	43.9	<u>44</u>	<u>42.0</u>

Table 5: Performance on low resource languages on INDICPARAM. We **bold** and *italicize* the best overall performance, underline and *italicize* the best performance in each model-size category and underline the second best overall performance.

as GPT-5 and DeepSeek-V3.2 demonstrate that substantial headroom remains for further enhancement in Indic language understanding. The results affirm that cross-lingual generalization, fine-tuning quality, and the underlying token distribution remain central determinants of downstream accuracy beyond mere parameter count.

5.2 Performance on Low- and Extremely Low-Resource Indic Languages

Table 4 and Table 5 compare contemporary LLMs on eight extremely low-resource Indic languages (Dogri, Maithili, Rajasthani, Sanskrit, Sanskrit-English code-mixed, Bodo, Santali, and Konkani) and four higher (yet still low-resource) Indic languages (Nepali, Gujarati, Marathi, and Odia). A clear performance hierarchy and a significant difficulty gap between the two tiers are observed.

On extremely low-resource languages, GPT-5 achieves the highest average of 45.1%, followed by the open-source DeepSeek-3.2 at 43.7. Claude-4.5 records 43.0, while Grok-4-fast attains 39.7. No model exceeds 41.5 on the most underrepresented languages (Bodo and Santali),

with most systems remaining below 36%. GPT-5 shows exceptional strength on Sanskrit (54.8) and Sanskrit-English code-mixed text (64.6).

In the low-resource cohort, scores rise substantially due to greater pre-training exposure. GPT-5 again leads with 44.9%, ahead of Claude-4.5 (42.0) and DeepSeek-3.2 (41.6). Grok-4-fast scores 39.2, trailing the top closed models by 3–5 points. Among open models, Gemma-3-27b (37.4) and the MoE Qwen3-A3B (38.7) perform best in their respective size categories.

The average gap between extremely low-resource and low-resource performance exceeds 2 points for most frontier models, highlighting the enduring challenge of languages with minimal digital presence. Within the extremely low-resource set, Sanskrit-English code-mixing is anomalously easier (top scores 58–64%), whereas Bodo and Santali remain the hardest for all architectures. Although scaling has significantly improved performance on moderately low-resource Indic languages, extremely low-resource languages still impose a strict ceiling below 45% for even the strongest current systems.

5.3 Performance on General Knowledge and Language Understanding Tasks

The evaluation distinguishes between two complementary capabilities: (i) *General Knowledge* questions that require factual reasoning and culturally grounded world knowledge (Table 9), and (ii) *Language-Specific Understanding* questions that probe grammatical, lexical, morphosyntactic, and discourse-level proficiency within each Indic language (Table 10).

5.3.1 General Knowledge

Frontier models exhibit clear separation from the rest. GPT-5 achieves the highest average (45.8), followed by Claude-4.5 (43.5) and DeepSeek-3.2 (43.3). Among open-weight systems, Llama-3.3-70B (39.3) and gpt-oss-120b (38) perform competitively, while Grok-4-fast attains 40.3. Large gains are observed with scale: medium-sized models (8B–27B) top out at 37.9 (Gemma-3-27b), and small models (<8B) remain below 32 (Qwen3-4B: 31.3). As shown in Table 9, GPT-5 dominates Sanskrit (57.8) and Sanskrit-English code-mixed (66.3) contexts, reflecting superior classical-language representation.

5.3.2 Language-Specific Understanding

Performance on pure linguistic tasks is substantially lower and reveals greater variance, underscoring that morphosyntactic mastery lags behind factual recall in most current systems. Even the strongest models rarely exceed 60% on any single language (see Table 10). GPT-5 again leads (52.53 average), followed by Claude-4.5 (49.8) and DeepSeek-3.2 (46.3) while Grok-4-fast records 44.9. Several medium-scale models punch above their parameter count: Llama-4-Scout (41.8, especially strong on Marathi) and Gemma-3-27b (40.9) occasionally outperform larger dense models on morphologically complex languages. Smaller models struggle severely, with many scores in the 20–35% range; as evident in Table 10, Sanskrit, Sanskrit+Eng, Dogri, Nepali and Santali prove particularly challenging across the board (often <50% even for frontier systems).

The gap between general knowledge (Table 9) and language-specific understanding (Table 10) shrinks markedly at the frontier: GPT-5, Claude-4.5, and DeepSeek-3.2 exhibit the average differences of 10–15 percentage points between the two categories, whereas smaller models

drop 5–10 points when moving from factual to linguistic tasks. These results indicate that while scaling and refined post-training have significantly improved Indic factual reasoning, deep grammatical and discursive proficiency in low- and extremely low-resource Indic languages remains a critical frontier for further advancement.

5.4 Cross-Dimensional Insights

A comparative analysis of general knowledge (Table 9) and language-specific understanding (Table 10) reveals distinct scaling behaviours. General knowledge accuracy correlates strongly with model scale (Pearson $r = 0.91$ across the 19 evaluated systems), whereas language-specific understanding shows a weaker but still substantial correlation ($r = 0.73$). The correlation between the two dimensions is moderately high ($r = 0.79$), indicating that top-tier factual reasoning usually co-occurs with stronger linguistic proficiency, yet meaningful divergences remain.

Several models achieve notably balanced performance. Gemma-3-27b (GK: 37.9, LU: 40.9) and Mistral3.1-24b (GK: 35.3, LU: 42.7) outperform many larger dense models on linguistic tasks, demonstrating that continued pre-training on diverse Indic corpora and instruction tuning can compensate for parameter count in morphosyntactically demanding settings. Similarly, the MoE-based Llama-4-Scout (17B/16E) reaches 41.8 on language understanding while scoring only 35.6 on general knowledge, highlighting the value of sparse architectures for grammar-heavy evaluation.

At the frontier, GPT-5, Claude-4.5, DeepSeek-3.2, and Grok-4-fast maintain the average gaps (10–15 percentage points) between the two dimensions, suggesting that extensive post-training and reinforcement learning from multilingual feedback help align linguistic mastery with factual recall. Models with <8B parameters routinely exhibit drops of 510 points when shifting from knowledge to linguistic questions.

These results underscore that future progress on low- and extremely low-resource Indic languages will require not only continued scaling but also (i) higher-quality and more diverse Indic pre-training data, (ii) dedicated cross-lingual continuation training, and (iii) architectural innovations that prioritise fine-grained morphological and syntactic representation.

6 Conclusion

We present INDICPARAM consisting of >13k questions for 11 low and extremely-low resource Indic languages. Our benchmark reveals that even the strongest contemporary LLMs, including frontier proprietary systems, remain far from reliable on graduate-level questions in low- and extremely low-resource Indic languages, with no model surpassing 50% average accuracy. While larger closed models such as GPT-5, Claude-4.5 and open-weight model, DeepSeek-3.2 consistently outperform open-weight alternatives. These findings indicate INDICPARAM as a challenging, human-curated evaluation suite that fills a key gap in existing Indic benchmarks and underline the need for more balanced pretraining and adaptation for under-represented Indic languages.

7 Limitations

While INDICPARAM substantially broadens evaluation for low- and extremely low-resource Indic languages, the benchmark is derived exclusively from UGC-NET language papers in linguistics and literature, which biases the content toward academic exam styles and may not fully represent everyday or domain-general language use. Further, our evaluation is conducted using a log-likelihood-based multiple-choice setup in a zero-shot configuration, which may yield different results from alternative evaluation schemes such as few-shot prompting, or chain-of-thought prompting, and therefore does not capture the full range of behaviors that these models might exhibit under richer prompting strategies. We use a single, uniform prompt template across all languages and models; different prompt wordings, instructions, or language choices (for example, prompting in English versus the target Indic language) could lead to substantially different performance, so our reported results should be interpreted as one consistent evaluation setting rather than an exhaustive exploration of prompting space.

8 Ethics

All data described in this work was scraped from publicly available resources. INDICPARAM is released for non-commercial research use; code and annotation scripts are licensed under MIT. All annotation was conducted by consenting native speakers compensated at local rates, and no personally identifiable exam-candidate information

appears in the benchmark. Users interested in commercial use should contact UGC-NTA for relevant permissions.

References

- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. On the cross-lingual transferability of monolingual representations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637.
- Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Jon Ander Campos, Yi Chern Tan, Kelly Marchisio, Max Bartolo, Sebastian Ruder, Acyr Locatelli, Julia Kreutzer, Nick Frosst, Aidan Gomez, Phil Blunsom, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. 2024. [Aya 23: Open weight releases to further multilingual progress](#).
- DeepSeek-AI. 2025a. Deepseek-v3.2-exp: Boosting long-context efficiency with deepseek sparse attention.
- DeepSeek-AI. 2025b. Deepseek-v3.2-exp: Boosting long-context efficiency with deepseek sparse attention.
- Sumanth Doddapaneni, Rahul Aralikkatte, Gowtham Ramesh, Shreya Goyal, Mitesh M Khapra, Anoop Kunchukuttan, and Pratyush Kumar. 2023. Towards leaving no indic language behind: Building monolingual corpora, benchmark and models for indic languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12402–12426.
- Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, MarcAurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation. In *International conference on machine learning*, pages 4411–4421. PMLR.
- IAMAI. 2024. [Internet in india 2024](#). Research report, Internet and Mobile Association of India (IAMAI). Prepared in collaboration with Kantar.
- Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, Gokul NC, Avik Bhattacharyya, Mitesh M Khapra, and Pratyush Kumar. 2020. Indicnlp suite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for indian

- languages. In *Findings of the association for computational linguistics: EMNLP 2020*, pages 4948–4961.
- Aman Kumar, Himani Shrotriya, Prachi Sahu, Amogh Mishra, Raj Dabre, Ratish Puduppully, Anoop Kunchukuttan, Mitesh M Khapra, and Pratyush Kumar. 2022. Indicnlg benchmark: Multilingual datasets for diverse nlg tasks in indic languages. In *Proceedings of the 2022 conference on empirical methods in natural language processing*, pages 5363–5394.
- Lokesh Madasu, Gopichand Kanumolu, Nirimal Surange, and Manish Shrivastava. 2023. Mukhyansh: A headline generation dataset for indic languages. In *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation*, pages 620–634.
- Ayush Maheshwari, Kaushal Sharma, Vivek Patel, and Aditya Maheshwari. 2025. Parambench: A graduate-level benchmark for evaluating llm understanding on indic subjects. *arXiv preprint arXiv:2508.16185*.
- Arijit Maji, Raghvendra Kumar, Akash Ghosh, Sriparna Saha, et al. 2025a. Sanskriti: A comprehensive benchmark for evaluating language models’ knowledge of indian culture. *arXiv preprint arXiv:2506.15355*.
- Arijit Maji, Raghvendra Kumar, Akash Ghosh, Nemil Shah, Abhilekh Borah, Vanshika Shah, Nishant Mishra, Sriparna Saha, et al. 2025b. Drishtikon: A multimodal multilingual benchmark for testing language models’ understanding on indian culture. *arXiv preprint arXiv:2509.19274*.
- Team MistralAI. [Mistralai/mistral-small-3.1-24b-instruct-2503 ð hugging face](#).
- Office of the Registrar General & Census Commissioner, India. 2018. [Census of india 2011: Paper 1 of 2018 – language, india, states and union territories \(table c-16\)](#). Census of India 2011, Paper 1 of 2018, Language Data (Table C-16).
- OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, and Rahul K. 2025. [gpt-oss-120b & gpt-oss-20b model card](#).
- Guilherme Penedo, Hynek Kydlíek, Vinko Sabolec, Bettina Messmer, Negar Foroutan, Amir Hossein Kargaran, Colin Raffel, Martin Jaggi, Leandro Von Werra, and Thomas Wolf. 2025. [Fineweb2: One pipeline to scale them all – adapting pre-training data processing to every language](#).
- Libo Qin, Qiguang Chen, Yuhang Zhou, Zhi Chen, Yinghui Li, Lizi Liao, Min Li, Wanxiang Che, and Philip S. Yu. 2025. [A survey of multilingual large language models](#). *Patterns*, 6(1).
- Abhishek Kumar Singh, Vishwajeet Kumar, Rudra Murthy, Jaydeep Sen, Ashish Mittal, and Ganesh Ramakrishnan. 2025. Indic qa benchmark: A multilingual benchmark to evaluate question answering capability of llms for indic languages. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 2607–2626.
- Harman Singh, Nitish Gupta, Shikhar Bharadwaj, Dinesh Tewari, and Partha Talukdar. 2024. Indicgenbench: A multilingual benchmark to evaluate generation capabilities of llms on indic languages. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11047–11073.
- Gemma 3 Team. 2025. [Gemma 3 technical report](#).
- Llama 3 Team. 2024. [The llama 3 herd of models](#).
- Team-Qwen3. 2025. [Qwen3 technical report](#).
- Sshubam Verma, Mohammed Safi Ur Rahman, Vishwajeet Kumar, Rudra Murthy Venkataramana, and Jaydeep Sen. 2025. Milu: A multi-task indic language understanding benchmark. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics*.
- Ishaan Watts, Varun Gumma, Aditya Yadavalli, Vivek Seshadri, Manohar Swaminathan, and Sunayana Sitaram. 2024. Pariksha: A large-scale investigation of human-llm evaluator agreement on multilingual and multi-cultural data. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7900–7932.

Appendix

9 Annotation Setup

For each language, we first collected official UGC-NET language papers and corresponding answer keys from the official website⁴ and converted all files into a uniform machine-readable format. Question papers appeared in double-column layouts, and many PDFs contained non-selectable text. Hence, all documents were processed through a proprietary OCR pipeline to extract text consistently across heterogeneous layouts and scan qualities. During this stage, each question was tagged with its language, subject, exam session, and year, and answer keys were aligned to their corresponding items.

Following OCR, human annotators manually corrected recognition errors, restored missing characters and diacritics, and normalised punctuation and spacing so that questions and options follow a pre-specified parsing format. Beyond surface correction, annotators standardised structural formattings such as option markers, numbering schemes, and line breaks. Each final entry in INDICPARAM is stored with explicit fields for question text, question type, four answer options, correct option, language, subject, year, and unique question identifier.

Annotation was carried out by native speakers having expertise in the target language and having prior experience reading graduate-level exam material. Since many papers mix language-focused and knowledge-focused content, annotators were required to demonstrate strong proficiency in the relevant Indic script. All annotators received detailed written guidelines describing how to handle incomplete questions, ambiguous answer keys, and exam-specific conventions, as well as how to assign labels for language understanding (LU) versus knowledge-related (KR) items and for the six question types used in the benchmark. Prior to full-scale annotation, annotators were given a sample dataset along with worked examples. They were then assigned trial files, which were reviewed and corrected by the reviewers. Feedback was provided iteratively, and only after demonstrating consistent accuracy were annotators assigned larger batches of data. For question-answer pairs, annotation quality followed a focused two-step review pipeline. In the first step, annotators corrected OCR errors in the question and options, aligned each item with the official answer key, and ensured that exactly one option was marked as correct. In the second step, a senior reviewer re-checked these corrected pairs against the original exam papers, verifying that the question text and the designated correct option was faithful to the source. Annotators were compensated at a rate of \$1 per 10 questions.

10 Question Types-wise results

Table 6 presents the distribution of question types across subjects, and Table 7 provides the corresponding performance results for these question types.

10.1 Question Types

Table 6 provides a detailed breakdown of the types of questions present in each language. MCQs dominate the distribution, but a substantial number of list matching, AssertionReasoning(A&R), identification of incorrect statements (IS), ordering and blank filling questions contribute to different question types.

10.2 Question type-wise results

Table 7 shows the performance of each model on INDICPARAM by question type. In the smaller-sized models, Qwen3-4B demonstrates better performance in a majority formats of questions, which includes MCQ, IS, fill in the blanks and Ordering. This model outperforms alternative 34B models, including Llama3.2-3B and Gemma3-4B. In the medium-size model category, Gemma3-27B shows the best overall performance, with considerable improvements in MCQ, A&R, and fill in the blanks. At the same time, Gemma3-12B and Mistral3.1-24B also show competitive results in List Matching and Ordering, respectively. Large models show a considerable improvement in their ability to handle different question types. Llama3.3-70B model performs best on most question types, especially MCQ, list matching, A&R and IS. DeepSeek-3.2 performs best on MCQ, list matching, A&R shows good results compared to other

⁴<https://ugcnet.nta.nic.in/>

Language	MCQ	A&R	List Matching	Blank Filling	IS	Ordering	Total
Nepali	957	1	38	5	25	12	1038
Marathi	933	83	139	1	54	35	1245
Gujarati	345	172	208	0	128	191	1044
Odia	379	40	70	1	22	65	577
Maithili	1265	0	0	21	0	0	1286
Konkani	1069	30	134	19	45	31	1328
Santali	772	36	46	11	2	6	873
Bodo	481	288	270	17	0	257	1313
Dogri	223	153	217	1	231	202	1027
Rajasthani	1114	4	11	14	13	34	1190
Sanskrit	1210	1	39	34	11	20	1315
Sans-Eng	905	3	13	33	14	3	971
Total	9653	811	1185	157	545	856	13207

Table 6: Different question types across all languages in INDICPARAM. A&R refers to Assertion and Reason type question, IS refers to identification of incorrect statement. ‘Sans-Eng’ denotes a separate set of Sanskrit-English code-mixed question-answer pairs which forms a separate set.

large-sized models in the category. Within MoE models, Qwen3-A3B performs best on MCQ, IS and fill in the blanks type questions. On the other hand, gpt-oss-120b shows the best results on list matching, A&R, and ordering related questions. Closed-source models consistently outperform their open-source counterparts. GPT-5, achieves best scores across almost all question categories, including MCQ, fill in the blanks, list matching, A&R and ordering. In contrast, Claude-4.5 demonstrates best performance for IS questions. Grok4-fast performance is consistently lower than that of GPT-5 and Claude-4.5.

11 Performance on KR and LU questions

The trends observed in Table 9 and Table 10 show that accuracy improves consistently with increasing model scale, with frontier models achieving the strongest results. As discussed in Section 5.3.1, general knowledge tasks display steady but moderate gains, whereas Section 5.3.2 highlights far greater variability in language-specific tasks, reflecting the difficulty of deep multilingual grounding. Architectural improvements, including MoE routing, enable mid-sized models narrow the performance gap with larger systems.

12 Zero shot prompting template

In Table 8, we provide the zero-shot prompt used while evaluating all the models.

Size	Model	MCQ	List Matching	A&R	Blanks	IS	Ordering
Small-size	Qwen2.5-3B	30.3	32.6	25.6	26.1	27.7	30.5
	Llama3.2-3B	26.9	25.2	29.3	28.7	22.4	31.2
	Gemma3-4b	28.7	27.6	<u>30.8</u>	28	27.5	28.7
	Qwen3-4B	<u>31.5</u>	<u>27.8</u>	27.1	<u>31.8</u>	<u>28.4</u>	<u>32.6</u>
Medium-size	Aya-8b	29.1	28.9	32.3	34.4	25.5	27.9
	Llama3.1-8B	30.5	23.2	32.7	35	24.4	27.6
	Gemma3-12b	35.1	<u>35.9</u>	36.7	33.1	<u>37.1</u>	34.8
	Mistral3.1-24b	35.7	34.3	27.1	33.8	35.8	<u>37</u>
	Gemma3-27b	<u>38.1</u>	32.8	<u>37.2</u>	<u>38.9</u>	36.3	35.7
Large-size	Aya-32b	33.8	36.4	31.4	<u>36.9</u>	27.2	31.3
	Qwen3-32B	34.6	24.6	32.8	29.9	31.9	29.8
	Llama3.3-70B	<u>37</u>	<u>39.7</u>	<u>41.6</u>	34.4	<u>40.7</u>	<u>38.1</u>
	DeepSeek-3.2	43.2	47.2	44.5	<u>45.9</u>	38	<u>37.1</u>
MoE	Qwen3-A3B	<u>39.9</u>	33.1	37.9	<u>40.1</u>	<u>35.6</u>	34
	Llama-4-Scout	37.5	28.3	34.6	35	30.8	26.4
	gpt-oss-120b	37.9	<u>38</u>	<u>38.7</u>	38.2	32.5	<u>37.4</u>
Closed	GPT-5	46.1	<u>41.3</u>	<u>41.4</u>	47.8	<u>43.3</u>	41.6
	Grok4-fast	40.4	39.2	33.7	36.9	38.5	36.9
	Claude-4.5	<u>43.9</u>	39	40	37.6	45.7	36.6

Table 7: Question-type wise performance of all models on INDICPARAM.

Zero-Shot Prompting
<p>Prompt = f"""Question: {'question_text'}</p> <p>Options:</p> <p>A) {'option_a'}</p> <p>B) {'option_b'}</p> <p>C) {'option_c'}</p> <p>D) {'option_d'}</p> <p>The above question is written in {language} language. Please analyze the question and options carefully, and select the correct answer. Respond ONLY with one letter (A, B, C, or D) corresponding to the correct option. Do not provide any explanation or additional text.””</p>

Table 8: Zero-Shot prompt applied across all models for evaluation

Model	Gujarati	Konkani	Maithili	Marathi	Oriya	Rajasthani	Sanskrit	Sans-Eng	Dogri	Nepali	Santali
Qwen2.5-3B	<u>29.3</u>	<u>31.5</u>	29.2	28.2	29.1	27.1	31.7	<u>35.5</u>	32.2	27.3	33.5
Llama3.2-3B	26.2	28.2	21.1	28.3	27.4	26.3	27.8	27.9	28	25.1	27.3
Gemma-3-4b	26.3	29.1	23.8	29.2	30.4	25.3	29.5	32.3	30.4	28.2	31.5
Qwen3-4B	27.3	29.3	<u>31.5</u>	<u>29.7</u>	<u>32.6</u>	<u>27.4</u>	<u>34.5</u>	33.6	<u>35.6</u>	<u>28.9</u>	<u>33.9</u>
Aya-8b	27.6	26.6	29	28.3	27.8	31.5	29.7	30.7	31.2	27.2	30.1
Llama-3.1-8B	27.3	29.1	27.9	31.8	28.9	29.3	32.4	34.2	25.9	24.6	30.4
Gemma-3-12b	<u>33.9</u>	35.5	29.7	37.5	33	32.6	39	43.1	36.7	34.8	33.5
Mistral3.1-24b	32.3	31.7	<u>36.3</u>	32.1	28	32.3	<u>41.3</u>	45.2	38	34.8	35.8
Gemma-3-27b	31.3	<u>36</u>	33.2	<u>38.9</u>	<u>37.8</u>	<u>36.9</u>	39.9	<u>49.4</u>	<u>37.8</u>	<u>39.4</u>	<u>36.7</u>
Aya-32b	28.4	32.5	31.7	32.5	28.3	33	36.3	40.6	35.3	32.5	35.5
Qwen3-32B	35.9	34.2	31.5	40	38.5	29.5	40.1	49.5	43.1	33.6	39.7
Llama-3.3-70B	33.6	37.2	35.2	42.7	<u>39.3</u>	36.6	42.8	52.1	37.5	35.8	39
DeepSeek-3.2	<u>39.6</u>	<u>37.3</u>	<u>39.7</u>	<u>43.9</u>	38.7	<u>37.4</u>	<u>53.5</u>	<u>62.9</u>	<u>42.8</u>	<u>39.5</u>	40.7
Qwen3-A3B	30.5	28.3	29.6	34.9	33.9	29.1	38.2	44.7	35.8	30.8	35.3
Llama-4-Scout	27.4	34.4	30	<u>42.1</u>	32.2	32.2	<u>44.3</u>	49.4	30	34	35.7
gpt-oss-120b	<u>37</u>	<u>34.9</u>	<u>33.9</u>	38.7	<u>36.3</u>	<u>35.5</u>	41.5	<u>50.2</u>	<u>38.7</u>	<u>36.1</u>	<u>35.7</u>
GPT-5	40.1	<u>38.5</u>	41.4	47.7	44.1	40.8	57.8	66.3	45.1	<u>41.5</u>	40.1
Grok-4-fast	35.9	37.1	36.5	38.8	38.3	37.1	45.2	55.9	39.9	40.6	38.1
Claude-4.5	37	41.9	<u>40.3</u>	42.9	39.8	39.7	48.8	59.5	<u>44.9</u>	43.4	<u>40.3</u>

Table 9: Performance of models on knowledge related question category.

Model	Gujarati	Konkani	Maithili	Marathi	Oriya	Rajasthani	Sanskrit	Sans-Eng	Dogri	Nepali	Santali
Qwen2.5-3B	<u>66.7</u>	27.3	25.4	20.7	27.4	32.4	<u>32.6</u>	<u>29.7</u>	<u>31.9</u>	26.2	30.3
Llama3.2-3B	16.7	<u>51.5</u>	28.5	34.5	25.6	26.4	28.8	27.9	27.7	30.8	28.3
Gemma-3-4b	16.7	30.3	30.8	43.1	35.9	33.6	27.7	31.5	30.9	<u>32.3</u>	33.3
Qwen3-4B	50	39.4	<u>33.1</u>	<u>46.6</u>	<u>42.7</u>	<u>36.7</u>	30.7	33.3	29.8	27.7	<u>33.3</u>
Aya-8b	50	39.4	38.5	27.6	29.1	29.1	28.4	23.4	33	32.8	<u>36.4</u>
Llama-3.1-8B	16.7	45.5	39.2	39.7	30.8	38.2	31.4	32.4	28.7	29.7	32.3
Gemma-3-12b	16.7	42.4	<u>49.2</u>	56.9	40.2	41.5	29.5	31.5	<u>39.4</u>	34.9	35.4
Mistral3.1-24b	83.3	45.5	46.9	48.3	25.6	41.2	<u>34.8</u>	35.1	37.8	35.4	35.4
Gemma-3-27b	0	<u>57.6</u>	44.6	<u>60.3</u>	<u>47.9</u>	<u>47.6</u>	33.3	<u>46.8</u>	36.7	<u>39</u>	36.4
Aya-32b	<u>33.3</u>	<u>57.6</u>	46.2	46.6	31.6	38.8	31.1	27	34.6	35.4	36.4
Qwen3-32B	33.3	45.5	56.2	56.9	49.6	44.8	36	41.4	44.7	41	40.4
Llama-3.3-70B	16.7	42.4	<u>57.7</u>	58.6	44.4	42.4	37.5	47.7	36.2	<u>43.1</u>	48.5
DeepSeek-3.2	<u>16.7</u>	<u>48.5</u>	56.9	<u>58.6</u>	<u>55.6</u>	<u>47</u>	<u>42.8</u>	<u>49.5</u>	<u>45.7</u>	40.5	<u>47.5</u>
Qwen3-A3B	<u>50</u>	30.3	43.1	48.3	36.8	38.2	30.7	40.5	35.6	34.9	40.4
Llama-4-Scout	16.7	45.5	42.3	70.7	46.2	41.5	<u>42.8</u>	<u>43.2</u>	32.4	36.9	<u>41.4</u>
gpt-oss-120b	33.3	<u>54.5</u>	<u>50.8</u>	60.3	<u>53.8</u>	<u>49.7</u>	34.8	31.5	<u>44.1</u>	<u>41.5</u>	35.4
GPT-5	<u>33.3</u>	60.6	60.8	<u>63.8</u>	66.7	<u>54.5</u>	<u>43.2</u>	51.4	48.4	51.8	43.4
Grok-4-fast	33.3	51.5	52.3	60.3	50.4	44.2	33.7	40.5	41.5	42.6	<u>44.4</u>
Claude-4.5	33.3	51.5	<u>60.8</u>	63.8	<u>60.7</u>	55.8	43.6	46.8	44.1	<u>45.6</u>	41.4

Table 10: Performance of models on language understanding (LU) question category.

13 Examples questions in *INDICPARAM*

Table 11 presents examples of six distinct types of questions for each language used in *INDICPARAM*.

Question	option(a)	option(b)	option(c)	option(d)	Ans	Type	Class	Lang
1. बालिवधस्य वर्णनमस्ति रामायणस्य	सुन्दरकाण्डे	किष्किन्धा काण्डे	अरण्यकाण्डे	बालकाण्डे	b	MCQ	G	san
2. 'मन्त्रिपरिषदं द्वादशमात्यान्कुर्वीत' इति कस्य मान्यता ?	बार्हस्पत्यानाम्	कौटिल्यस्य	औशनसाम्	मानवानाम्	d	MCQ	G	san
3. अधोऽडिकतानां समीचीनमुत्तरं चिनुत - (a) सरमा-पणि 1. बृहदारण्यकोप-सम्वादः निषत् (b) स्वाध्यायान्मा 2. ऋग्वेदस्य प्रमदः दशममण्डले (c) कल्पः 3. तैत्तिरीयोपनिषत् (d) आत्मनस्तु 4. हस्तः कामाय सर्वं प्रियं भवति (a) (b) (c) (d)	1 3 2 4	4 2 3 1	2 3 4 1	3 2 1 4	c	Order	G	san
4. भाशा दी सभर्ने थमां लौहकी इकाई ऐ :	ध्वनि	ध्वनिग्राम	रूपग्राम	वाक्य	a	MCQ	L	doi
5. पैहली चंदी च दिती गे दिये प्रविष्टियों दा दूई चंदी दिये प्रविष्टियें कन्नै स्हेई मिलान करो : चंदी-1 चंदी-2 (अ) डोगरी काव्य-चर्चा (i) प्रो. रामनाथ शास्त्री (ब) परख-पड़ताल (ii) शिवनाथ (स) डोगरी साहित्य दा इतिहास (iii) प्रो. चम्पा शर्मा (द) बाबा जित्तो (iv) ओम गोस्वामी कोड : (अ) (ब) (स) (द)	(ii) (i) (iv) (ii) (iii) (iv) (i) (iii) (iv) (ii)	(iii) (iv) (i) (iii) (ii) (i) (iv) (ii)	(i) (iii) (iv) (ii)	(i) (iii) (iv) (ii)	c	LM	G	doi
6. 'इक लड़ाई होर' ते 'बंजर' दे रचेता न :	मदन मोहन शर्मा ते मोहन सिंह ।	मोहन सिंह ते नरसिंह देव जम्वाल ।	मोहन सिंह ते मदन मोहन शर्मा ।	जितेन्द्र शर्मा ते दीनू भाई पंत ।	c	Order	G	doi

7. બન્ને યાદીની વિગતો સર-ખાવી સાચાં જોડકાં બનાવો : (a) શેષાદ્રિ (i) વિજયરાય વૈદ્ય (b) રામ વૃંદાવની (ii) ત્રિભુ-વનદાસ લુહાર (c) ત્રિશુળ (iii) રાજેન્દ્ર શાહ (d) મયૂરાનંદ (iv) ખબરદાર (a) (b) (c) (d)	(i) (iv) (iii) (ii)	(iv) (iii) (ii) (i)	(ii) (iv) (iii) (i)	(iii) (ii) (i) (iv)	b	LM	G	guj
8. નીચેનાં વિધાનોને કાર્યકાર-ણસંબંધે તપાસો : (A) 'વદતોવ્યાઘાત' મધ્યમપદલોપી સમાસ છે. (R) મધ્યમપદલોપી સમાસ સર્વપદપ્રધાન છે.	(A) અને (R) બન્ને સાચાં છે.	(A) અને (R) બન્ને ખોટાં છે.	(A) સાચું છે અને (R) ખોટું છે.	(A) ખોટું છે અને (R) સાચું છે.	b	A&R	G	guj
9. નીચેનામાંથી સુસંગત વિગત-જૂથ જણાવો :	વીસમી સદી, વૈશ્વાનર, નટમંડળ	ગુજરાત વિદ્યા- સભા, ગુજરાતી સાહિત્ય પરિષદ, ગુજરાત સાહિત્ય અકાદમી	ઊડણ ચરકલડી, સોયનું નાકું, મનીષા	ખીચડી, અંતઃસ્ત્રોતા, આવતી- કાલનો સૂરજ	b	MCQ	G	guj
10. ફાવો તો પર્યાય વૈંચૂન કાઢૂન વાક્ય પૂર્ણ કરાત. સોळाव्या शेंकड्यातल्या रामायणाच्या हातबरपांत	ठांयीं ठांयीं पुर्तुगेज उतरां मेळटात.	कांयच पुर्तुगेज उतरां मेळनात.	संस्कृत उतरां भरसून पुर्तुगेजीक कोंकणीची सया मारता.	देशी आनी अपभ्रंशी भारतीय भासो मेळटात.	b	Blank fill- ing	G	gom
11. भाशीक नादांचो भाशेचे बांदावळीचे नदरेंतल्यान अभ्यास करपी व्याकरणाक किर्ते म्हणतात ?	अर्थविज्ञान	नादविज्ञान	वाक्य- विचार	स्वनीम- विचार	d	MCQ	G	gom
12. 'काळोकिट' हो समास हांतलो खंयचो ?	उपमावाचक कर्मधारय समास	विशेशण उभयपद समास	द्वंद्व समास	बहुव्रीहि समास	b	MCQ	L	gom
13. हुनकासँ भेंट भेल छल' लिखने छथि	मणिपद्म	अमर	सुमन	मधुप	a	MCQ	G	mai
14. "शिव की विष पचाय यदि लितथि नहि माथे ।"	देखितथि	बुझितथि	करितथि	सकितथि	d	Blank fill- ing	G	mai
15. 'एकावली परिणय' में अंगीरस अछि	शान्त	अद्भुत	हास्य	श्रृंगार	d	MCQ	L	mai

16. पुढीलपैकी कोणते नाटक विजय तेंडुलकर यांनी लिहिले नाही ?	गृहस्थ	सूर्यास्त	बेबी	श्रीमंत	b	MCQ	G	mar
17. पुढीलपैकी कोणती कादंबरी ऐतिहासिक नाही ?	वज्राघात	कालिकामूर्ति	सावळ्या तांडेल	दुर्देवी रंगू	b	IS	G	mar
18. पुढील साहित्यकृतींचा कालानुसार क्रम लावा.	धग, किडलेली माणसे, औदुंबर, स्मृतिचित्रे,	किडलेली माणसे, औदुंबर, स्मृतिचित्रे, धग,	धग, औदुंबर, स्मृतिचित्रे, किडलेली माणसे,	औदुंबर, स्मृतिचित्रे, किडलेली माणसे, धग,	d	Order	G	mar
19. शास्त्रीय मार्क्सवादी दृष्टिमा आइडियोलजी भन्नाले के बुझिन्छ ?	भ्रमपूर्ण चेतना ।	सैद्धान्तिक आस्था	प्रचलित विचारधारा ।	पुस्तकहरू पढेर पाइने ज्ञान ।	a	MCQ	G	nep
20. वाक्यमा शब्दहरू माझ रहने सम्बन्धलाई के भनिन्छ ?	विन्यासक्रमी सम्बन्ध ।	सहचारक्रमी सम्बन्ध ।	व्यतिरेकी सम्बन्ध ।	परस्परव्यापि सम्बन्ध ।	a	MCQ	L	nep
21. नेपाली भाषामा प्रयोग हुने 'चोलो', 'पटुका', 'मुन्धुम', 'बाउसे' शब्दहरू कुन भाषा परिवारबाट आएका हुन् ?	भारोपेली ।	भोट-बर्मेली ।	आग्नेय ।	द्रविड ।	b	MCQ	L	nep
22. "राजस्थानी भाषा रौ विगसाव नागर अपभ्रंश सूं हुयौ है।" ओ कथन है -	डॉ. एल.पी. टेस्सिटोरी	डॉ. ग्रियर्सन	डॉ. सुनीतिकुमार चटर्जी	डॉ. सुकुमार सेन	b	MCQ	G	Raj
23. राजस्थानी अनुवाद री दीठ सूं किसौ अनुवाद रौ जोड़ौ गलत है ?	गांधीजी री आत्मकथा - आईदानसिंह भाटी	भरथरी शतक - मनोहर शर्मा	मेघदूत - नारायणसिंह भाटी	रचाव - चेतन स्वामी	d	IS	G	Raj
24. बुगचौ' सबद रौ सही अर्थ है :	लुगाइयां रै कपड़ा राखण रौ कपड़े रौ थेलौ	लुगाइयां री लकड़ी री पेटी	लुगाइयां री कपड़ा राखण री लोह री पेटी	मड़दां रै कपड़ा राखण री लकड़ी री पेटी	a	MCQ	L	Raj
25. स हि कदाचिद् वाच्ये विधिरूपे प्रतिषेधरूपः' अत्र स इत्यनेन कोऽभिप्रेतः ?	वाच्यार्थः	लक्ष्यार्थः	व्यङ्ग्यार्थः	तात्पर्यार्थः	c	MCQ	L	sans-mix
26. कौषीतकि उपनिषद् कस्य वेदस्य अस्ति ? कौषीतकि उपनिषद् किस वेद की है ? Of which Veda is कौषीतकि उपनिषद् ?	अथर्ववेदस्य	सामवेदस्य	ऋग्वेदस्य	कृष्णयजुर्वेदस्य	c	MCQ	G	sans-mix

27. ଅଧୋଲିଖିତେଷୁ କେନ ସହ କସ୍ୟ ସମ୍ବନ୍ଧ: ? ତାଲିକା ଚିନୁତ - (a) ସତ୍କାର୍ଯ୍ୟବାଦ: (i) ନ୍ୟାୟବୈଶେଷିକାଣାମ୍ (b) ପରମାଣୁବାଦ: (ii) ବୌଦ୍ଧାନାମ୍ (c) ବିବର୍ତ୍ତବାଦ: (iii) ସଂଖ୍ୟାନାମ୍ (d) ବିଜ୍ଞାନବାଦ: (iv) ଅଦ୍ୱୈତବେଦାନ୍ତିନାମ୍ (a) (b) (c) (d)	(iii) (ii) (i) (iv)	(iii) (iv) (i) (ii)	(iii) (i) (iv) (ii)	(iii) (i) (ii) (iv)	c	LM	G	sans-mix
28. ସେଦାୟ ହାପଡ଼ାମକୋ କାଥା ଲେକାତେ ହୋଝ ହୋପୋନକୋବାକ୍ ପାରିସ ଦୋ ଓକୋ ଜାୟଗା ରେ କୋ ହାଟିଜ ଲେଦା ?	ଚାମ୍ପା	ଚାୟ	ସାସାଡବେଡା	ଜାରପୀ	c	MCQ	G	sat
29. ଲାତାର ରେ ଓଲ ଆକାନ ଥାପିତ କାଥା (A) ଆର ଓଜେ (R) ରେନାକ୍ ଓକା ଗାଦେଲ ବାଝିତୀ ସୁହିୟା ? ଥାପିତ (A) : ଭାରତୀୟ ନୋଜୋର ତେ ଗାୟାନ ସାଁବହେତୁ ଦୋ ସୁକ ରେ ମୁଘାତ ଲାକତୀୟାନା । ଓଜେ (R) : ଚେଦାକ୍ ଜେ ନୋବା କାଥା ହୋଁ ପଞ୍ଚିମ ଦିସୋମ ରେନ ସାଁପହେତୁହିୟା ନୋଜୋର ତେ ହୋଁ ନାପାୟା ।	(A) ଆର (R) ବାଝସୁହିୟା	(A) ସୁହିୟା, (R) ବାଝସୁହିୟା	(A) ବାଝସୁହିୟା, (R) ସୁହିୟା	(A) ଆର (R) ବାନାର ସୁହିୟା	d	A&R	G	sat
30. ଓକା ପାରସୀ ଘାରୋଚ୍ଛୁ ରେନାକ୍ ଓକା ପାରସୀ କାନା ? ଓନା ସୁହି ମିଲାବ ମେଁ । ସୂଚୀ - I ସୂଚୀ - II a.ଆର୍ଯ୍ୟ ଭାଷା ପରିବାର 1. ତିବ୍ବତୀ b. ଅଗ୍ନେୟ ଭାଷା ପରିବାର 2. ମାଲ୍ତୋ c. ଦ୍ରାବିଡ଼ ଭାଷା ପରିବାର 3. ମୁଣ୍ଡା d. ଚୀନୀ ଭାଷା ପରିବାର 4. ବଂଗଳା କୂଟ : a b c d	1 2 3 4	4 3 2 1	3 1 4 2	2 4 3 1	b	LM	G	sat
31. ମୁକୁ ଅକ୍ଷର ଶେଷରେ କଣ ଥାଏ ?	ଅକ୍ଷର	ସୂରଧୂନି	ବ୍ୟଞ୍ଜନଧୂନି	ସୂର୍ଜଧୂନି	b	MCQ	L	ory
32. ପ୍ରକାଶକାଳ ଅନୁଯାୟୀ ଉପଯୁକ୍ତ କ୍ରମ ନିର୍ଦ୍ଧାରଣ କର ।	ଉଷା, ଚିଲିକା, ପାର୍ବତୀ, ମହାଯାତ୍ରା	ଚିଲିକା, ଉଷା, ପାର୍ବତୀ, ମହାଯାତ୍ରା	ପାର୍ବତୀ, ଚିଲିକା, ଉଷା, ମହାଯାତ୍ରା	ଉଷା, ପାର୍ବତୀ, ଚିଲିକା, ମହାଯାତ୍ରା	d	Order	G	ory
33. କେଉଁଟି ବ୍ୟୋମକେଶ ତ୍ରିପାଠୀଙ୍କ ନାଟକ ନୁହେଁ ?	ସିଂହଦ୍ୱାର	ଜାଗରଣ	ମୁଠାଏମାଟି	ସୁନାଫରୁଆ	c	IS	G	ory

34. बाहायनाने गाहायाव होनाय मेथाइफोरनि बाहायजानाय सम फारि साजायनायखौ सायख - I. बैसागो मेथाइ II. दोहोरोम सिबिनाय मेथाइ III. बेलाड IV. मैगं खानायाव खननाय मेथाइ	II III IV I	III III IV	IV III I II	I III IV III	b	Order	G	brx
35. अलंबार' लाइसिखौ सोर सुजुर्दोमोन	मदाराम ब्रह्म	सुकुमार बसुमतारी	सतिस चन्द्र बसुमतारी	प्रम'द चन्द्र ब्रह्म	d	MCQ	G	brx
36. गाहायनि बुंफोरनाय (A) आरो जाहोन (R) खौ गेबें ना गोरान्थि बाहायनानै सायख : बुंफोरनाय (A) : आथिखालाव मोननाय बासिराम जोहोलाव बेलाडखौ थारै आबुं बेलाड एबा natural ballad बुंथावा । जाहोन (R) : मानोना बे बेलाडा लिरनाय महराव फोसावजानाय ।	(A)-आरों (R) मोननैबो गेबें	(A)-आ गेबें आरो (R)-आ गोरान्थि	(A) आरो (R) मोननैबो गोरान्थि	(A)-आ गोरान्थि आरो (R) आ गेबें	b	A&R	G	brx

Table 11: Sample questions from all languages in INDICPARAM along with their options, answers, question types, and question classes, with languages indicated using ISO-639 codes. We use *Order* to denote ordering-type questions, and the abbreviations *G* and *L* refers to General Knowledge (GK) and Language Understanding (LU) questions respectively.