# Evaluating Indian Language Performance in LLMs

Tattle

Aatman Vaidya

# The Growth of Large Language Models

# Multilingual Usage and Claims



ChatGPT now speaks Hindi, Assamese, Bengali and other Indian languages! Here's how to get replies in local languages

OpenAI's ChatGPT, based on the GPT-3.5 language model, can now respond in Hindi and other India... integration...

Meta AI gets multilingual with support for 7 languages including Hindi

Meta is also releasing Llama 3.1, the latest version of its open source large language model that will be available in 8B, 70B and 405B versions.

# Questions

- What are the capabilities of LLMs to understand and generate text in Indian Languages?

# Questions

- What are the capabilities of LLMs to understand and generate text in Indian Languages?

- How are LLMs that support Indian languages developed? What are the key characteristics (license, access, ownership etc)?

# Questions

- What are the capabilities of LLMs to understand and generate text in Indian Languages?

- How are LLMs that support Indian languages developed? What are the key characteristics (license, access, ownership etc)?

- How many resources are available for Indian languages, does this correlate with the number of speakers? How do resource gaps impact LLM performance?

# Methodology

We conduct analysis using existing frameworks proposed by Ecosystem Graphs ([Bommasani et al 2023](#))[1] and by [KJ et al. 2024](#)

### LLMs

### Evaluation

# Methodology

We conduct analysis using existing frameworks proposed by Ecosystem Graphs ([Bommasani et al 2023](#))[1] and by [KJ et al. 2024](#)

## LLMs

| Host |
|:---:|

| Size |
|:---:|

| Training Dataset |
|:---:|

| Access |
|:---:|

| License |
|:---:|

| Type and Architecture |
|:---:|

## Evaluation

# Methodology

We conduct analysis using existing frameworks proposed by Ecosystem Graphs ([Bommasani et al 2023](#))[1] and by [KJ et al. 2024](#)

## LLMs

| Host |
| --- |
| Size |
| Training Dataset |
| Access |
| License |
| Type and Architecture |

## Evaluation

| Understanding |
| --- |

| Generation |
| --- |

1 - https://crfm.stanford.edu/ecosystem-graphs/

# Analysis

- We analyze <u>28 models</u> that support Indian languages using the methodology outlined before.
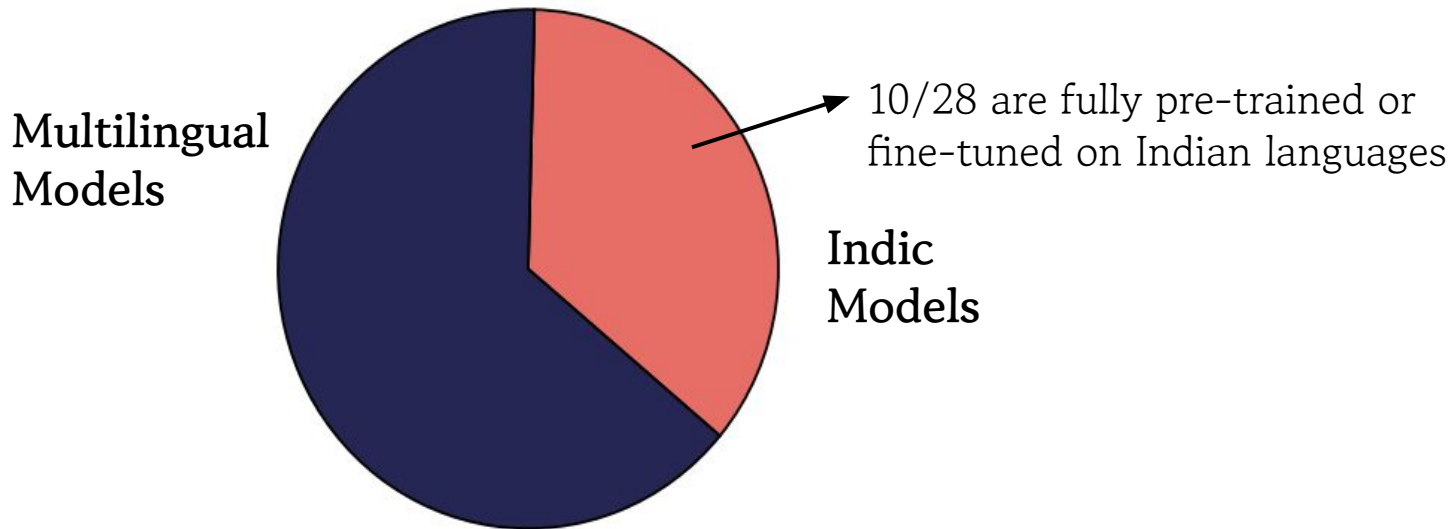- Models trained on multilingual corpora with Indian language data in them were also included.
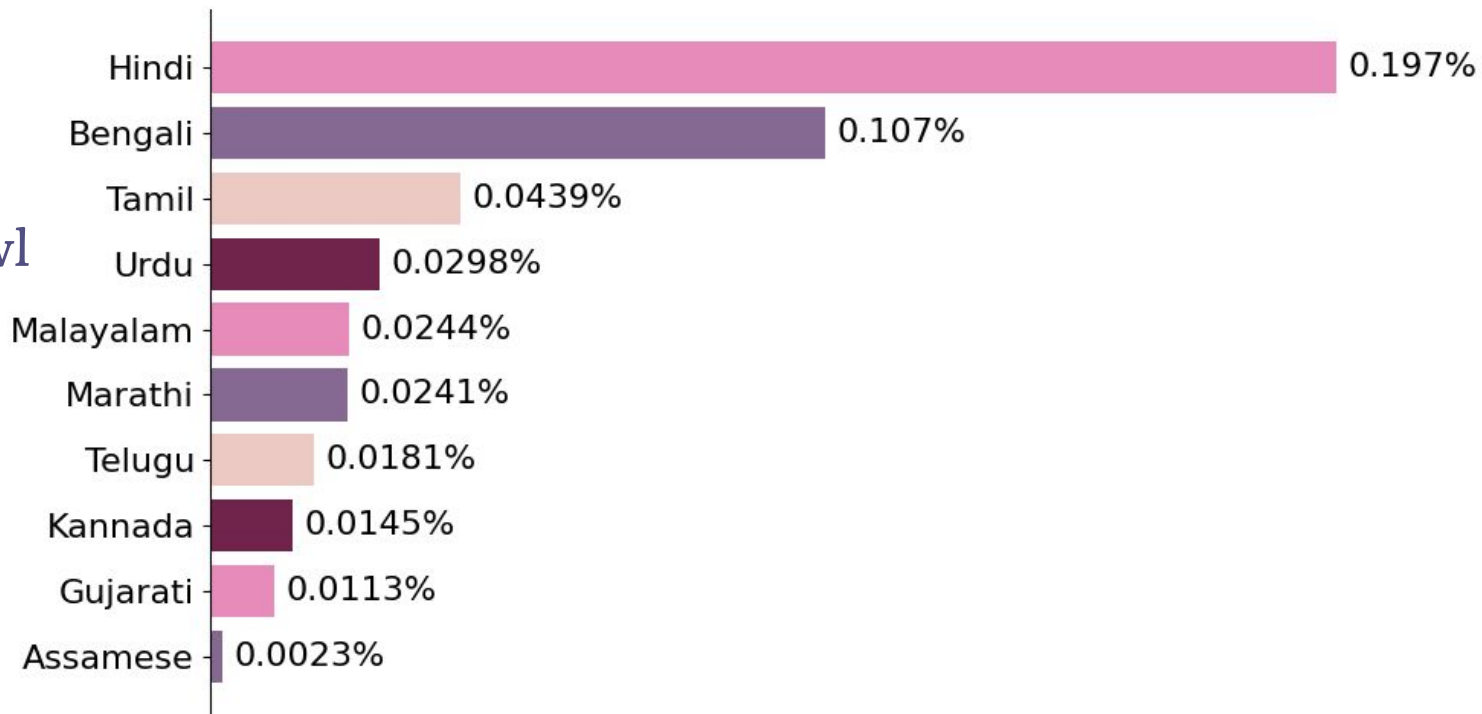
# Analysis

- We analyze <u>28 models</u> that support Indian languages using the methodology outlined before.
- Models trained on multilingual corpora with Indian language data in them were also included.

**Multilingual Models**

10/28 are fully pre-trained or fine-tuned on Indian languages

**Indic Models**

# Training Corpora

## Multilingual Models

Common Crawl

# Training Corpora

## Multilingual Models

Common Crawl



| Language | Percentage |
|----------|-----------|
| English | 42.8% |
| Hindi | 0.197% |
| Bengali | 0.107% |
| Tamil | 0.0439% |
| Urdu | 0.0298% |
| Malayalam | 0.0244% |
| Marathi | 0.0241% |
| Telugu | 0.0181% |
| Kannada | 0.0145% |
| Gujarati | 0.0113% |
| Assamese | 0.0023% |

Log Scale

# Training Corpora

Wikipedia
Articles



| Language | Percentage |
|---|---|
| Urdu | 0.336% |
| Tamil | 0.265% |
| Hindi | 0.256% |
| Bengali | 0.25% |
| Telugu | 0.159% |
| Marathi | 0.154% |
| Malayalam | 0.135% |
| Western Punjabi | 0.115% |
| Punjabi | 0.086% |
| Kannada | 0.052% |
| Gujarati | 0.048% |
| Odia | 0.029% |
| Sindhi | 0.029% |
| Assamese | 0.022% |
| Maithili | 0.022% |
| Sanskrit | 0.019% |
| Santali | 0.018% |
| Bihari (Bhojpuri) | 0.014% |
| Central Tibetan (Lhasa Tibetan) | 0.011% |
| Kashmiri | 0.009% |
| Konkani (Goan Konkani) | 0.006% |
| Tulu | 0.004% |
| Pali | 0.004% |
| Awadhi | 0.004% |
| Dzongkha | 0.00051% |

Log
Scale

# Training Corpora

Common
Crawl

Wikipedia
Articles

Indic Models

# Training Corpora

Common
Crawl

MuRIL - PMINDIA,
Dakshina

Wikipedia
Articles

AryaBhatta-GemmaGenZ -
GenZ_Vikas instruction
dataset

Indic Models

Airavata - Indic Instruct
dataset

# Training Corpora

Common
Crawl

MuRIL - PMINDIA,
Dakshina

AryaBhatta-GemmaGenZ -
GenZ_Vikas instruction
dataset

Wikipedia
Articles

Indic Models

Synthetic Data

Airavata - Indic Instruct
dataset

# Evaluation

Understanding

Generation

(Ahuja et. al 2024) "MEGAVERSE: Benchmarking Large Language Models Across Languages, Modalities, Models and Tasks" - https://arxiv.org/abs/2311.07463

# Evaluation

| Understanding | → | Natural Language Inference, POS Tagging, Text Classification, Sentiment Analysis, Name-Entity Recognition, Paraphrase Detection, Commonsense Reasoning, |
|---|---|---|
| Generation | → | Translation, Summarization, Question-Answering, Image Captioning, Multiple choice Question-Answering, Task Oriented Dialogue |

(Ahuja et. al 2024) "MEGAVERSE: Benchmarking Large Language Models Across Languages, Modalities, Models and Tasks" - https://arxiv.org/abs/2311.07463

# Evaluation

| Tasks | Sub-Tasks | Languages Covered | Papers Referenced |
|---|---|---|---|
| Natural Language Understanding | Natural Language Inference | Hindi, Bengali, Punjabi, Kannada, Gujarati, Malayalam, Marathi, Telugu, Tamil, Oriya, Assamese, code-mixed English-Hindi, Nepali | Ahuja et al. [9], Doddapaneni et al. [26], Aggarwal et al. [7] |
| | Text Classification | Hindi, Kannada, Malayalam, Marathi, Tamil, Telugu, Urdu, Bengali, Gujarati, Assamese, Odia, Punjabi | Doddapaneni et al. [26], Kakwani et al. [41] |
| | Name-Entity Recognition | Urdu, Tamil, Telugu, Hindi, Gujarati, Malayalam, Marathi, Punjabi, Bengali, Kannada | Ahuja et al. [9], Doddapaneni et al. [26], Kakwani et al. [41] |
| Natural Language Generation | Translation | Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Tamil, Telugu, Urdu, Assamese, Bhojpuri, Nepali, Odia, Punjabi, Pashto, Sanskrit, Awadhi, Haryanvi, Tibetan, Bodo, Garhwali, Konkani, Chhattisgarhi, Rajasthani, Maithili, Manipuri, Malvi, Marwari, Santali | Singh et al. [64] |
| | Summarization | Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Tamil, Telugu, Urdu, Assamese, Bhojpuri, Nepali, Odia, Punjabi, Pashto, Sanskrit, Awadhi, Haryanvi, Tibetan, Bodo, Garhwali, Konkani, Chhattisgarhi, Rajasthani, Maithili, Manipuri, Malvi, Marwari, Santali | Singh et al. [64], Ahuja et al. [9], Hada et al. [32], Kumar et al. [47] |
| | Question-Answering | Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Tamil, Telugu, Urdu, Assamese, Odia, Punjabi | Singh et al. [64], Ahuja et al. [9], Doddapaneni et al. [26], Kakwani et al. [41] |
| | Image Captioning | Bengali, Hindi, Telugu | Ahuja et al. [10] |

Table 3. Tasks and Languages Covered in Existing Evaluation Datasets

# Evaluation

| Tasks | Sub-Tasks | Languages Covered | Papers Referenced |
|---|---|---|---|
| Natural Language Understanding | Natural Language Inference | Hindi, Bengali, Punjabi, Kannada, Gujarati, Malayalam, Marathi, Telugu, Tamil, Oriya, Assamese, code-mixed English-Hindi, Nepali | Ahuja et al. [9], Doddapaneni et al. [26], Aggarwal et al. [7] |
| | Text Classification | Hindi, Kannada, Malayalam, Marathi, Tamil, Telugu, Urdu, Bengali, Gujarati, Assamese, Odia, Punjabi | Doddapaneni et al. [26], Kakwani et al. [41] |
| | Name-Entity Recognition | Urdu, Tamil, Telugu, Hindi, Gujarati, Malayalam, Marathi, Punjabi, Bengali, Kannada | Ahuja et al. [9], Doddapaneni et al. [26], Kakwani et al. [41] |
| Natural Language Generation | Translation | Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Tamil, Telugu, Urdu, Assamese, Bhojpuri, Nepali, Odia, Punjabi, Pashto, Sanskrit, Awadhi, Haryanvi, Tibetan, Bodo, Garhwali, Konkani, Chhattisgarhi, Rajasthani, Maithili, Manipuri, Malvi, Marwari, Santali | Singh et al. [64] |
| | Summarization | Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Tamil, Telugu, Urdu, Assamese, Bhojpuri, Nepali, Odia, Punjabi, Pashto, Sanskrit, Awadhi, Haryanvi, Tibetan, Bodo, Garhwali, Konkani, Chhattisgarhi, Rajasthani, Maithili, Manipuri, Malvi, Marwari, Santali | Singh et al. [64], Ahuja et al. [9], Hada et al. [32], Kumar et al. [47] |
| | Question-Answering | Bengali, Gujarati, Hindi, Kannada, Malayalam, Marathi, Tamil, Telugu, Urdu, Assamese, Odia, Punjabi | Singh et al. [64], Ahuja et al. [9], Doddapaneni et al. [26], Kakwani et al. [41] |
| | Image Captioning | Bengali, Hindi, Telugu | Ahuja et al. [10] |

Table 3. Tasks and Languages Covered in Existing Evaluation Datasets

# Evaluation

| Language | Performance of LLM on Tasks | | | Number of Speakers in India |
|---|---|---|---|---|
| | Understanding | Generation | Evaluation Studies Referenced | |
| Hindi | HIGH | HIGH | | 528,347,193 |
| Bengali | HIGH | HIGH | | 97,237,669 |
| Marathi | HIGH | HIGH | Singh et al. [64], Ahuja et al. [9], Aggarwal et al. [7], Doddapaneni et al. [26], Kakwani et al. [41] | 83,026,680 |
| Telugu | HIGH | HIGH | | 81,127,740 |
| Tamil | HIGH | HIGH | | 69,026,881 |
| Gujarati | MEDIUM | MEDIUM | | 55,492,554 |
| Urdu | HIGH | HIGH | Singh et al. [64], Ahuja et al. [9], Doddapaneni et al. [26] | 50,772,631 |
| Kannada | MEDIUM | MEDIUM | | 43,706,512 |
| Oriya | LOW | LOW | | 37,521,324 |
| Malayalam | MEDIUM | MEDIUM | Singh et al. [64], Ahuja et al. [9], Aggarwal et al. [7], Doddapaneni et al. [26], Kakwani et al. [41] | 34,838,819 |
| Punjabi | LOW | LOW | | 33,124,726 |
| Assamese | MEDIUM | MEDIUM | | 15,311,351 |

Table 4. Overall Performance of LLMs on Evaluation Tasks for Indian Languages

No: of Speakers Source - https://en.wikipedia.org/wiki/Languages_of_India

# Evaluation

⚠️ This ranking is <u>relative to each other</u> and <u>NOT universal</u>

| Langauge | Performance of LLM on Tasks | | | Number of Speakers in India |
|---|---|---|---|---|
| | Understanding | Generation | Evaluation Studies Referenced | |
| Hindi | HIGH | HIGH | | 528,347,193 |
| Bengali | HIGH | HIGH | | 97,237,669 |
| Marathi | HIGH | HIGH | Singh et al. [64], Ahuja et al. [9], Aggarwal et al. [7], Doddapaneni et al. [26], Kakwani et al. [41] | 83,026,680 |
| Telugu | HIGH | HIGH | | 81,127,740 |
| Tamil | HIGH | HIGH | | 69,026,881 |
| Gujarati | MEDIUM | MEDIUM | | 55,492,554 |
| Urdu | HIGH | HIGH | Singh et al. [64], Ahuja et al. [9], Doddapaneni et al. [26] | 50,772,631 |
| Kannada | MEDIUM | MEDIUM | | 43,706,512 |
| Oriya | LOW | LOW | | 37,521,324 |
| Malayalam | MEDIUM | MEDIUM | Singh et al. [64], Ahuja et al. [9], Aggarwal et al. [7], Doddapaneni et al. [26], Kakwani et al. [41] | 34,838,819 |
| Punjabi | LOW | LOW | | 33,124,726 |
| Assamese | MEDIUM | MEDIUM | | 15,311,351 |

Table 4. Overall Performance of LLMs on Evaluation Tasks for Indian Languages

No: of Speakers Source - https://en.wikipedia.org/wiki/Languages_of_India

# Takeaways

- **Evaluation for low-resource languages is task, language and context specific.**
  - Difficult to generalise about coverage across tasks and datasets.

# Takeaways

- Evaluation for low-resource languages is task, language and context specific.
  - Difficult to generalise about coverage across tasks and datasets.

- The performance of models in a language is not directly correlated with the number of people who speak the language.

# Takeaways

- **Evaluation for low-resource languages is task, language and context specific.**
  - Difficult to generalise about coverage across tasks and datasets.

- **The performance of models in a language is not directly correlated with the number of people who speak the language.**

- **Data Contamination (Ahuja et al. 2023)**

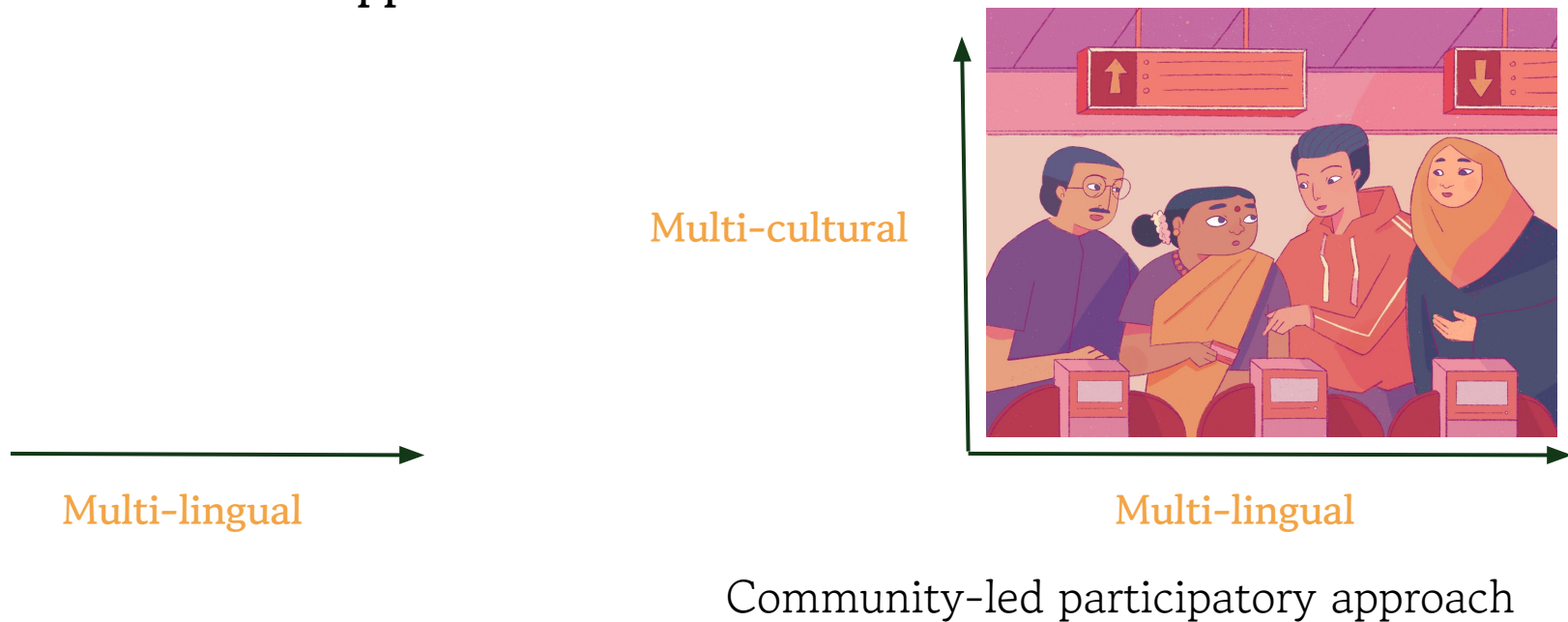| Dataset | Card Fill | Data Acc. w/o Down. | Release Date |
|---|---|---|---|
| XNLI | Full | Yes | September 2019 |
| Indic-XNLI | Full | Yes | April 2022 |
| PAWS-X | Full | Yes | August 2019 |
| XCOPA | Partial | Yes | April 2020 |
| XStoryCloze | Partial | No | May 2023 |
| XQuAD | Full | Yes | October 2019 |
| MLQA | Full | Yes | October 2019 |
| TyDiQA-GoldP | Full | Yes | February 2020 |
| IndicQA | Partial | Yes | September 2022 |
| PAN-X | Full | Yes | July 2017 |
| UDPOS | Full | Yes | March 2020 |
| XLSum | Partial | Yes | June 2021 |
| Jigsaw | None | No | February 2020 |
| GLUECos NLI | None | No | June 2020 |
| EN-ES-CS | None | No | May 2016 |

# Takeaways

- Evaluation has been 1-D so far

Multi-lingual →

# Takeaways

- Evaluating Indian language performance effectively requires a multi-cultural and contextual approach



Multi-cultural

Multi-lingual

Multi-lingual

Community-led participatory approach

# Acknowledgments

Tattle

## Supported and Funded By

ML Commons

Analysis of Indic Language Capabilities in LLMs

arXiv