

# Operator perspectives on Datasets for AI/ML in telecom networks

Aditya Jain

# Data Intelligence works at incredible scale..



## Coverage

~37.9%

Revenue market share<sup>1</sup>



## Scale

400M+

Subscribers India<sup>2</sup>



2.8+ T

Records Processed  
Daily

900+

Data Streams  
Integrated



30PB+

Capacity

30TB+

Incremental  
Data Daily



50M+

Transactions per second

1200+

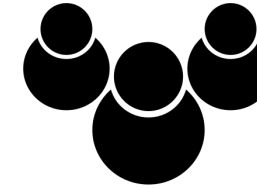
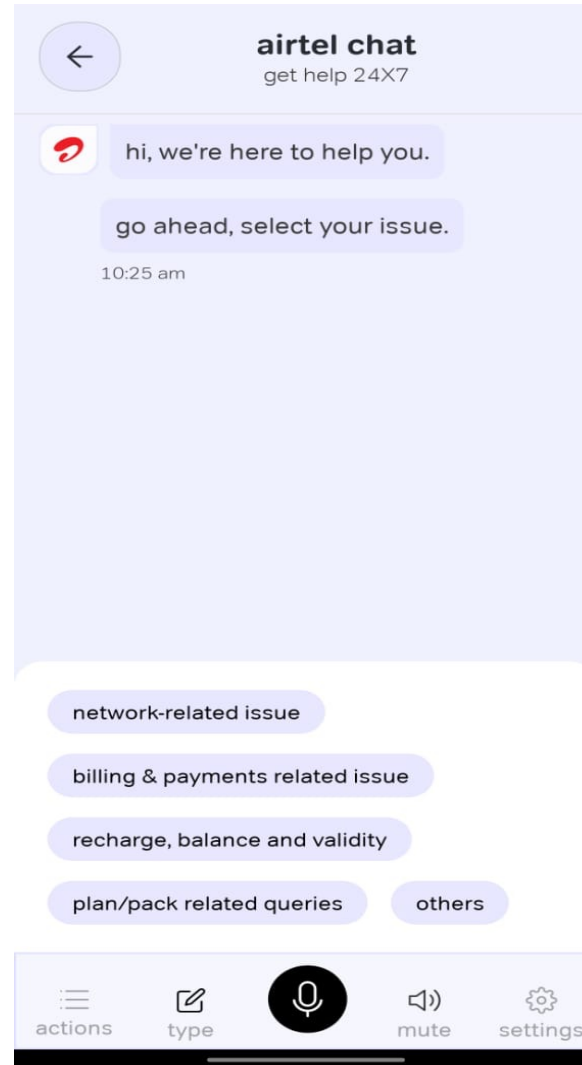
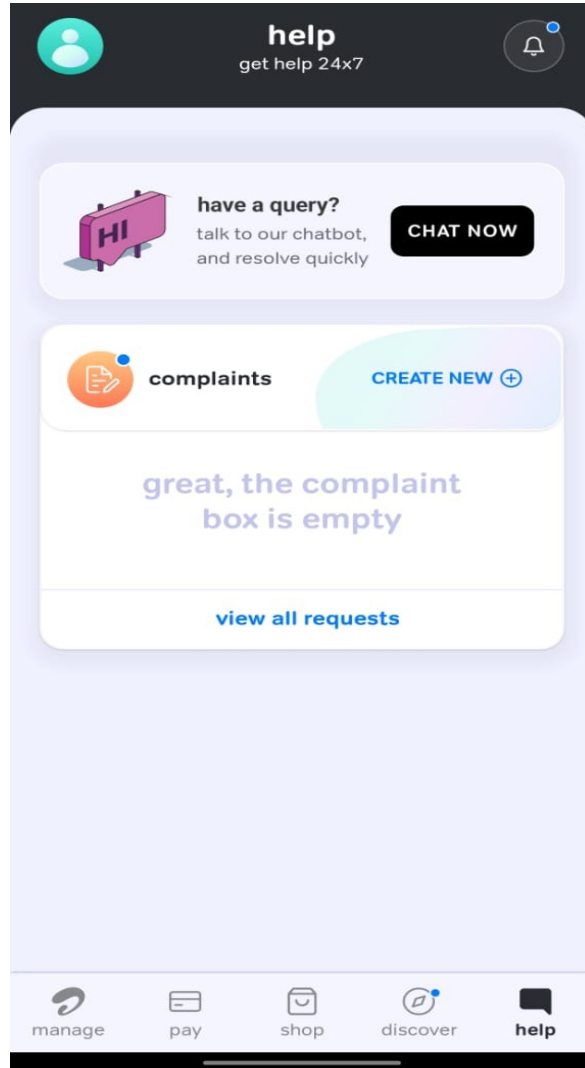
Active business  
users

<sup>1</sup> ET News Release [Link](#)

<sup>2</sup> PRESS RELEASE Bharti Airtel Q4 FY24 consolidated results

# Looking at customer engagement ...

Working at tremendous scale



Active App  
15M +  
Daily Users



200K + Chat and  
Voice Queries

Postpaid

Prepaid

Broadband

DTH

# Huge diversity in interactions

Voice

22 National Language

100+ Language

10k+ Dialects

Chat

Transliteration

Understanding Domain

Multi-lingual

O t p not coming

ಕರೆ ಡ್ರಾಪ್ ಎದುರಿಸುತ್ತಿದೆ

मैं कॉल ड्रॉप face कर  
रहा हूँ

आइ वांट टू कैंसल plan

ऐसएमएस नहीं जा रहे

Signals breaking

# Historically -> NLP : To detect intent of conversations



Real Data Challenges and Approach

## Challenges

Multi-Lingual Support

Understand Transliteration

Identifying domain specific word

Handling Garbage Queries

Low Latency



## Solution Approach

Active Learning



Deep Learning  
Transformer Based Multi-Lingual  
Model



Serviceable FastAPI

# Now -> LLMs works better in high resource languages..



Real Data Challenges and Approach

## Challenges

Multi-Lingual Support

Understand Transliteration

Identifying domain specific word

Handling Garbage Queries

Low Latency



## Solution Approach

Fine-tuned LLMs

# Challenges of LLMs with low resource languages



- ✓ Performance of LLMs for Low Resource Languages (LRLs) is hindered by **unavailability of high-quality open-source large-scale data** required for pre-training and fine tuning<sup>1</sup>
- ✓ **Multi-lingual models take advantage of cross-lingual transfer up until a point**, after which the overall performance on monolingual and cross-lingual benchmarks degrades.<sup>1</sup>
- ✓ LLMs fine-tuned specifically for low resource languages outperform base GPT models<sup>2</sup> but require heavy computation cost and time

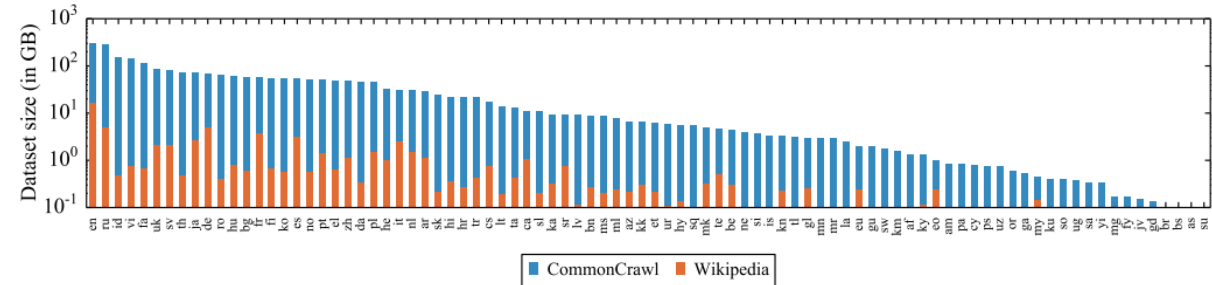


Figure 1: Amount of data in GiB (log-scale) for the 88 languages that appear in both the Wiki-100 corpus used for mBERT and XLM-100, and the CC-100 used for XLM-R. CC-100 increases the amount of data by several orders of magnitude, in particular for low-resource languages.

Model	Eng		Zho		Vie		Ind		Tha	
	GSM8K	MATH	GSM8K	MATH	GSM8K	MATH	GSM8K	MATH	GSM8K	MATH
ChatGPT-3.5	<b>80.8</b>	34.1	48.2	21.5	55.0	26.5	64.3	26.4	35.8	18.1
Qwen1.5-7B-chat	56.8	15.3	40.0	2.7	37.7	9.0	36.9	7.7	21.9	4.7
SeaLLM-7B-v2	78.2	27.5	<b>53.7</b>	17.6	69.9	23.8	<b>71.5</b>	24.4	59.6	22.4
SeaLLM-7B-v2.5	78.5	<b>34.9</b>	51.3	<b>22.1</b>	<b>72.3</b>	<b>30.2</b>	<b>71.5</b>	<b>30.1</b>	<b>62.0</b>	<b>28.4</b>

Table 4: GSM8K and MATH scores (Cobbe et al., 2021; Hendrycks et al., 2021b) and their translated-versions in Chinese, Vietnamese, Indonesian and Thai, under zero-shot chain-of-thought prompting for different models.

1. <https://arxiv.org/pdf/1911.02116> (Unsupervised Cross-lingual Representation Learning at Scale)  
2. <https://arxiv.org/pdf/2312.00738> (SeaLLMs - Large Language Models for Southeast Asia)

- ❖ Creation of benchmark datasets :
  - ❖ Multi modal in nature
  - ❖ Assess Network performance and perform Root Cause Analysis
  - ❖ Across low/high resource languages
  
- ❖ Creation of benchmark performance metrics
  - ❖ To measure Accuracy holistically across RCA, quality of response, etc.

***Benchmark Dataset Instances in other Industry around Reasoning***

- *MATH (Mathematics for Machine Learning)*
- *GSM8K (Grade School Math 8K)*
- -
- -

***Performance Metrics***

- *Accuracy, F1-Score*
- *BLEU Score*
- *ROGUE Score*
- -
- -



# Thank You !!

Aditya Jain

*aditya1.jain@airtel.com*