

# Research Highlights for awardees of

## GKII-Ashoka Breakthrough Grant

Ashoka University and Gupta-Klinsky India Institute at John Hopkins (GKII) launched the 1st round of Breakthrough Grant, a co-funded opportunity to support pilot projects led by faculty members from Ashoka University and Johns Hopkins University that involve interdisciplinary collaborations. The current research grant is focused on health data research.

# Leveraging Language Models and a Common Data Model to Unlock Real-World Evidence from Unstructured Electronic Health Record Data in India

Gautam Ahuja, BS<sup>1,\*</sup>, Himani Balutia, MPharm<sup>2,\*</sup>, Siddhant Poudyal, MS<sup>2,\*</sup>, Bableen Kaur, PhD<sup>2</sup>, Ragul N, BS<sup>1</sup>, Rudra Chinhara, BTech<sup>7</sup>, Sivsanjai GA, BS<sup>1</sup>, Hara Prasad Mishra, MBBS<sup>4</sup>, Tamoghna Ghosh, MBBS<sup>5</sup>, Sanjana Ahuja, BBA<sup>2</sup>, Douluri Pushkala Devi, MSc<sup>2</sup>, Shivangi Singh, MBBS<sup>9</sup>, Upamanyu Das, MBBS<sup>9</sup>, Tanisha Kohli, BS<sup>8</sup>, Angad Singh, BS<sup>1</sup>, Shrey Chhabra, BS<sup>1</sup>, Aarush Kumbhakern, BS<sup>1</sup>, Arundhati Mishra, MSc<sup>2</sup>, Poorva Nandedkar, BDS<sup>11</sup>, Claire Vania, MS<sup>6</sup>, Suditi Arora, BDS<sup>10</sup>, Ashiya Mahran, PGDHM<sup>10</sup>, Divya Laroyia<sup>10</sup>, Shyatto Raha<sup>10</sup>, **Matthew Robinson**, MD<sup>6</sup>, Rintu Kutum, PhD<sup>1,2,3,11,#</sup>

<sup>1</sup>Department of Computer Science, Ashoka University, Sonipat, Haryana, India

<sup>2</sup>Trivedi School of BioSciences, Ashoka University, Sonipat, Haryana, India

<sup>3</sup>Ashoka Mphasis Lab, Ashoka University, Sonipat, Haryana, India

<sup>4</sup>University College of Medical Sciences, Delhi, India <sup>5</sup>All India Institute of Medical Sciences Delhi, Delhi, India

<sup>6</sup>Division of Infectious Diseases, Johns Hopkins School of Medicine, Baltimore, United States

<sup>7</sup>Department of Computer Science, Central University of Haryana, India

<sup>8</sup>Department of Biology, Ashoka University, Sonipat, Haryana, India

<sup>9</sup>Muzaffarnagar Medical College, Uttar Pradesh, India <sup>10</sup>MyHealthcare Technologies Pvt. Ltd, Gurgaon, Haryana, India

<sup>11</sup>Lead contact; \*Contributed equally;

#Correspondence: rintu.kutum@ashoka.edu.in; rintu.kutum@augmented-health-systems.org

#### Objectives

- . To develop natural language processing (NLP) algorithms and leverage open source large language models (osLLMs) tailored to the Indian context to extract and transform clinical features from unstructured clinical notes into the Observational Medical Outcomes Partnership (OMOP) Common Data Model (CDM).
- 2. To determine the prevalence and choice of antibiotic prescription among outpatients with febrile respiratory illness in India and the US.
- 3. To assess the impact of antibiotic choice on treatment outcomes for inpatients with carbapenem-resistant organism pneumonia in India and the US.

#### Methods

- Bridging medication data of OMOP-CDM with the Common Drug Codes for India (CDCI), NRCeS with NLP and open source LLMs.
- Design and development of schema to connect CDCI-NRCeS with OMOP-CDM.
- Creation of high-quality annotation data (medication) to benchmark open-source LLMs with zero-shot and parameter-efficient fine-tuning.
- Validation of the extracted medication-related information with rule-based clinical NLP.
- 2. Mapping medications related to the use of antibiotics with WHO AWaRe (Access, Watch, Reserve) classification.

#### Results

#### A. Schema to connect CDCI-NRCeS with OMOP-CDM Legend ProductMaster SNOMED Indian Extension) (SNOMED Indian Extension) OMOP Non Standard Product Name Brand Name OMOP Standard Product Identifier SupplierMaster Brand Name Supplier Identifier Nomenclature: Generic Identifier 1. drug-brand-name 2. (active-ingredient License Number License Status 4. drug-route 5. drug-form GenericMaster SNOMED CT SubstanceMaster SNOMED Indian Extension) SNOMED International) → SCT-ID (SNOMED International) Generic Name (SNOMED Indian Extension Substance Name → SCT-Name (SNOMED International) CAS Number **OMOP Non Standard** Route of Administration Substance Description SCT-ID Dose Form Molecular Weight SCT-Name Therapeutic Role **OMOP Standard** (RxNoexn) Contra Indication RxN-Name Interaction with Drugs **IUPAC Name**

Molecular Formula

### B. Creation of high quality annotation data (medication)

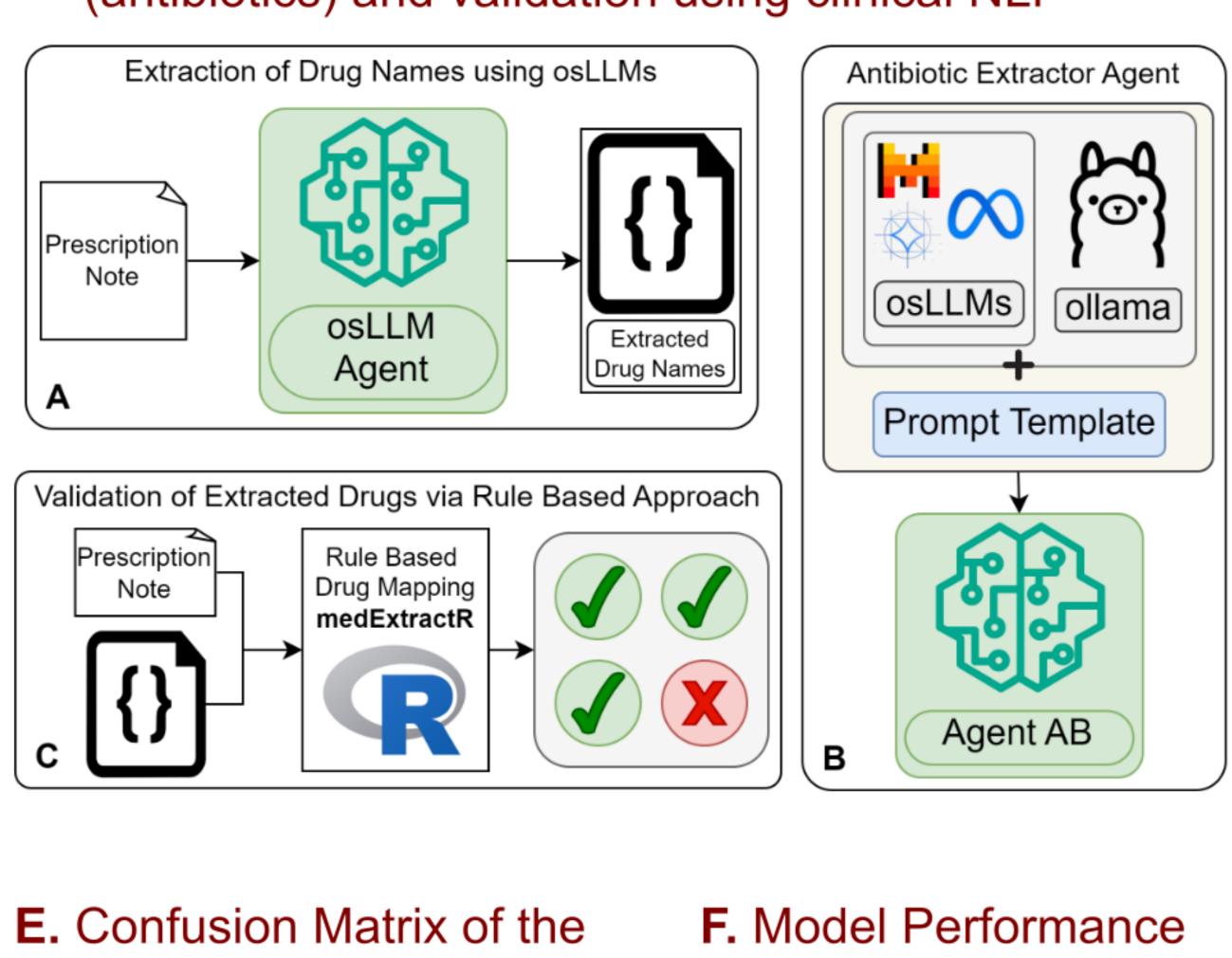
(i) Prescription notes Drez gargle 5 ml in 10ml water thrice a day for 3 days Tablet Mondeslor after dinner for 5 days Tablet Calpol 650 mg 3 to 4 times a day for 3 days Tablet Meftal 250 mg SOS for high grade fever or severe body pain

(ii) Medication annotation JSON format {'drug\_no\_0': {'drug\_name': 'Drez gargle', 'dose': '5 ml', 'dose\_amount': '-', 'dose\_change': '-', 'frequency': '-', 'intake time': '-', 'route': '-', 'duration': '3 days'}, 'drug no 1': {'drug name': 'Tablet Mondeslor', 'dose': '-', 'dose\_amount': '-', 'dose\_change': '-', 'frequency': '-', 'intake\_time': 'after dinner', 'route': '-', 'duration': '5 days'}, 'drug\_no\_2': {'drug\_name': 'Tablet Calpol', 'dose': '650 mg', 'dose\_amount': '-', 'dose\_change': '-', 'frequency': '-', 'intake\_time': '-', 'route': '-', 'duration': '3 days'}, 'drug\_no\_3': {'drug\_name': 'Tablet Meftal', 'dose': '250 mg', 'dose\_amount': '-', 'dose\_change': '-', 'frequency': '-', 'intake\_time': '-', 'route': '-', 'duration': '-'}}

## (iii) Benchmark of osLLMs for medication extraction

| SI No. | osLLM           | Precision | Recall | F1   | Hallucinated | Error rate |
|--------|-----------------|-----------|--------|------|--------------|------------|
| 1      | Gemma2-9b       | 0.91      | 0.9    | 0.9  | 0.21         | 0.09       |
| 2      | Llama2-7b       | 0.77      | 0.77   | 0.77 | 0.32         | 0.22       |
| 3      | Medalpaca-7b    | 0.8       | 0.79   | 0.79 | 0.71         | 0.2        |
| 4      | Meditron-7b     | 0.93      | 0.93   | 0.93 | 0.21         | 0.06       |
| 5      | Llama-3-8b      | 0.93      | 0.91   | 0.92 | 0.18         | 0.08       |
| 6      | Llama-3.1-8b    | 0.95      | 0.94   | 0.94 | 0.12         | 0.05       |
| 7      | Mistral-v0.1-7b | 0.93      | 0.92   | 0.92 | 0.26         | 0.07       |
| 8      | Mistral-v0.3-7b | 0.95      | 0.94   | 0.94 | 0.15         | 0.05       |
| 9      | Qwen2-7b        | 0.95      | 0.93   | 0.93 | 0.18         | 0.06       |

#### C. Agent based workflow to extract medication (antibiotics) and validation using clinical NLP



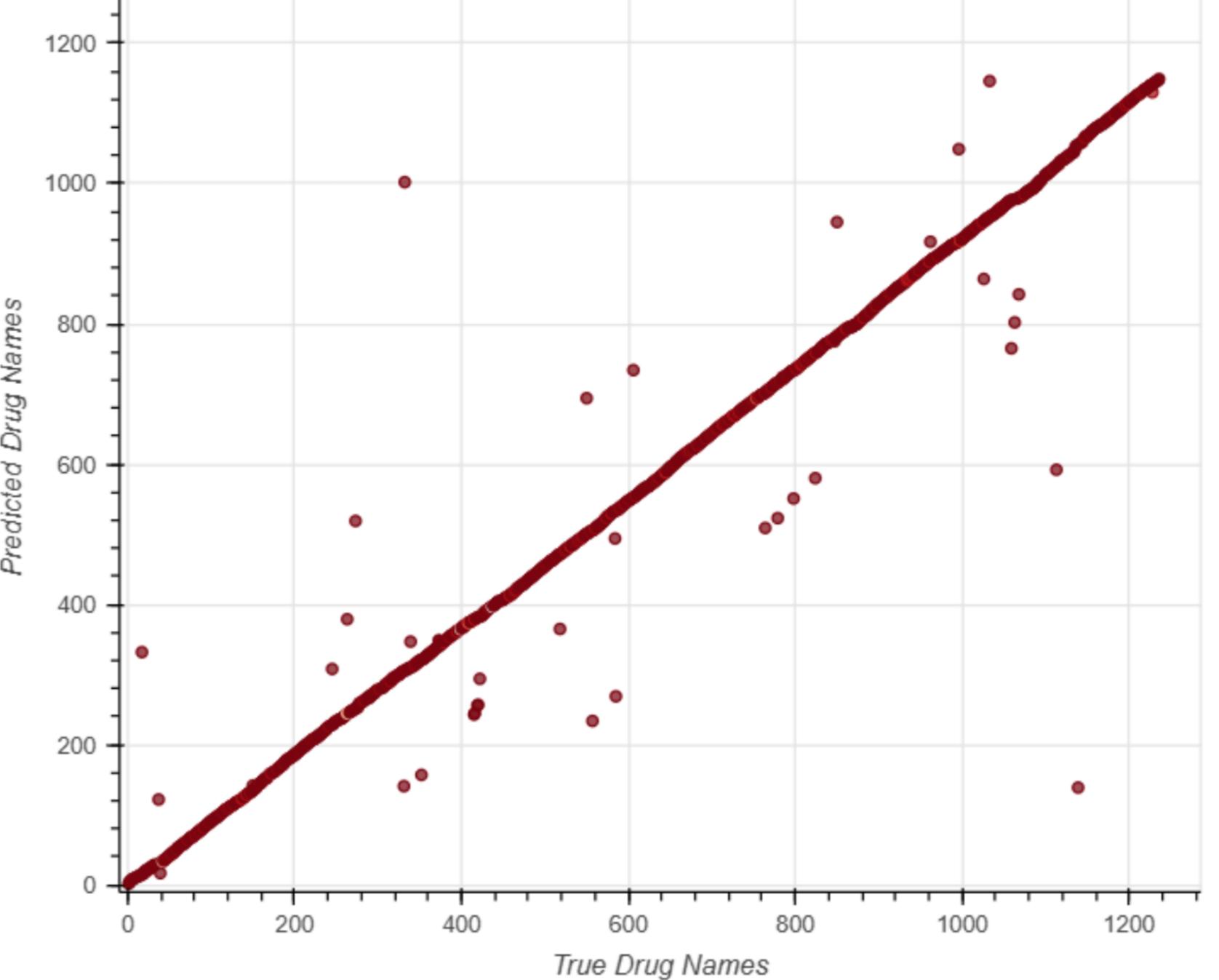


| best r   | nodel             |                       |                     |
|----------|-------------------|-----------------------|---------------------|
|          | Predicted<br>Drug | Predicted<br>Non-Drug | Drug 1/0            |
| Drug     | 1790              | 88                    | Drug vs<br>Non-Drug |
| Non-Drug | 166               | 0                     |                     |

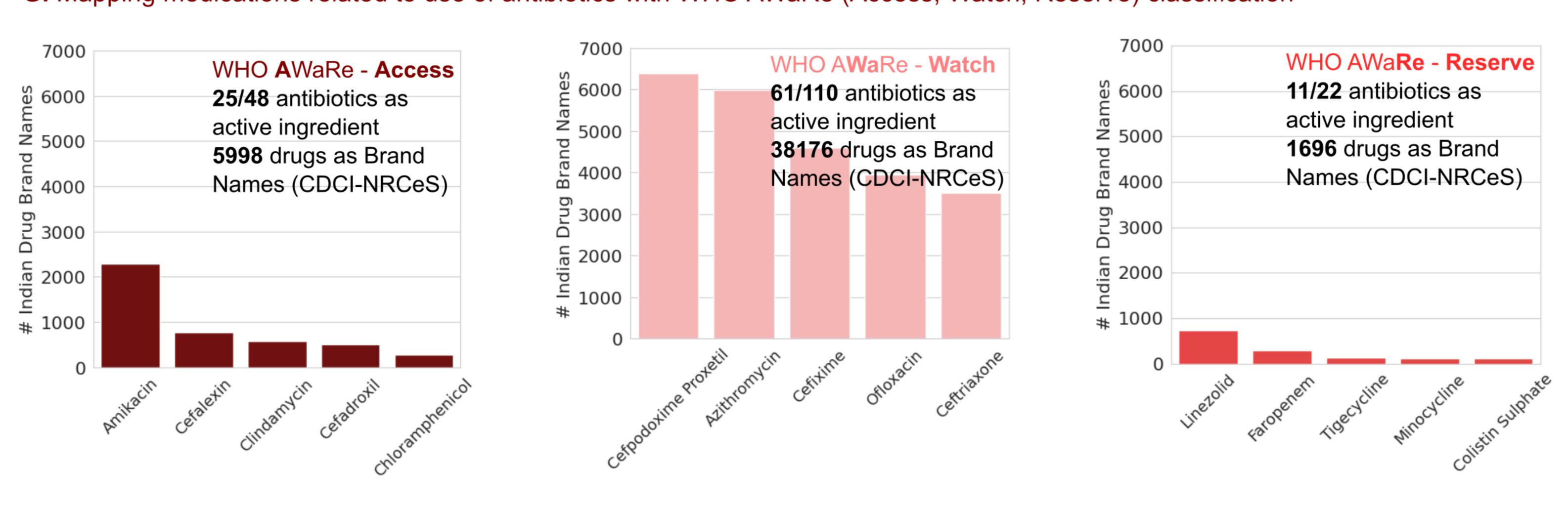
Classification of Drugs

|                     | Precision | Recall | F1   |
|---------------------|-----------|--------|------|
| Drug vs<br>Non-Drug | 0.95      | 0.91   | 0.93 |

#### D. Visual representation of the performance of agent based medication extraction with osLLM



## G. Mapping medications related to use of antibiotics with WHO AWaRe (Access, Watch, Reserve) classification



## **Future Direction**

Mapping medications related to antibiotics with symptoms and disease conditions to determine the prevalence and use of antibiotics.

## Acknowledgments

RK would like to thank the GKII-Ashoka Breakthrough grant and Koita Centre for Digital Health-Ashoka (KCDH-A) for the funding. Also, RK acknowledges the compute infrastructure support from the Mphasis F1 Foundation.

### Machine Learning Methods for Enhanced Forecasting of Antiretroviral Therapy Demand in India

Bhavesh Neekhra<sup>1</sup>, Swapnanil Mukherjee<sup>1</sup>, Kshitij Kapoor<sup>1</sup>, Debayan Gupta<sup>1</sup>, Manish Bamrotiya<sup>2</sup>, Sunil S. Solomon<sup>2</sup>, Steven J. Clipman<sup>2</sup> <sup>1</sup>Ashoka University, Sonipat, Haryana, India

<sup>2</sup>Johns Hopkins University School of Medicine, Baltimore, MD, USA

#### **Problem Statement**

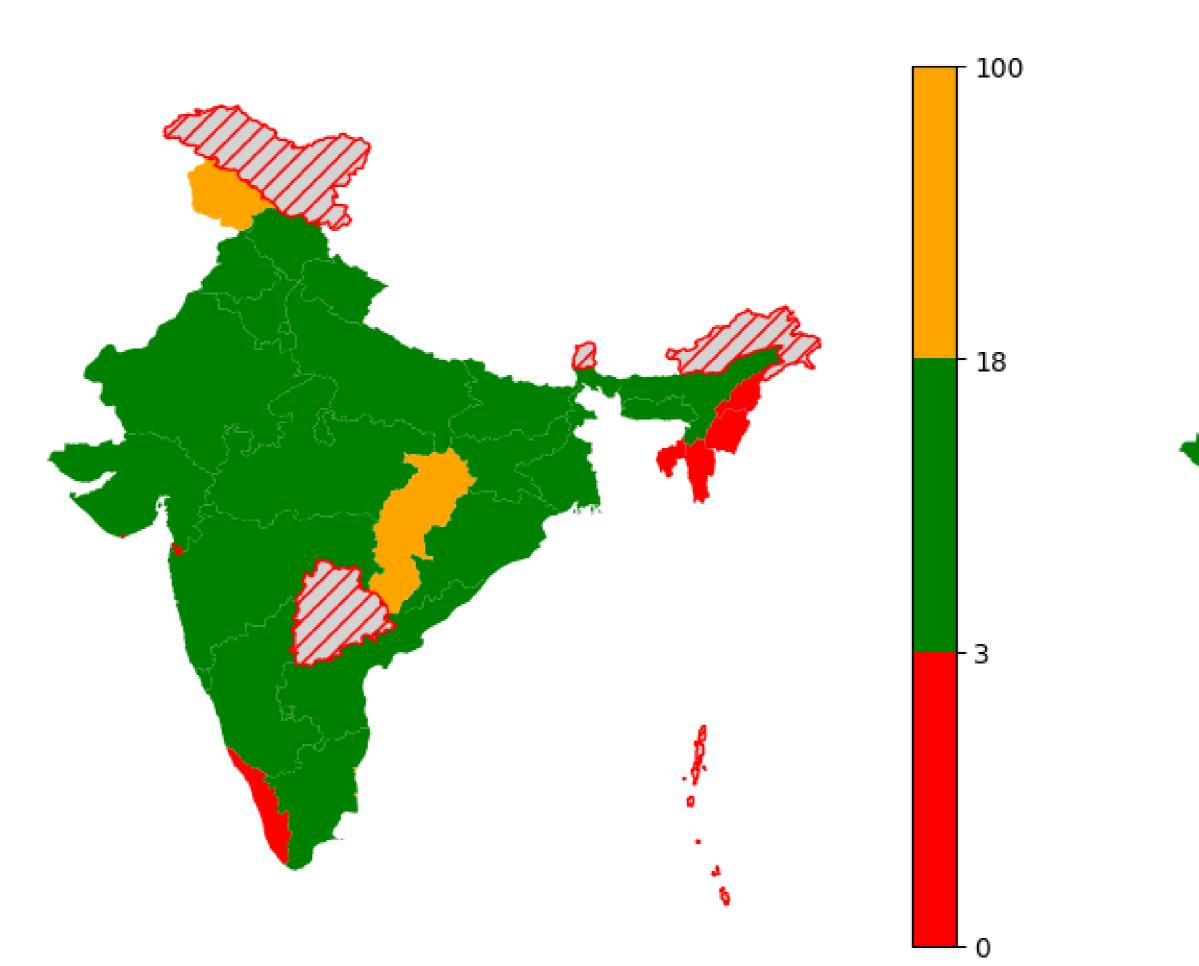
The primary objective is to develop and implement a machine learning-based model to enhance ART demand prediction in India. The specific aims are:

- Aim 1: Develop a machine learning model to predict the demand for each ART drug in India for the coming year.
- Aim 2: Validate the machine learning predictive model using historical data and compare its performance with current prediction methods.
- Aim 3: Develop a user-friendly interface for healthcare administrators to use the predictive model.

## Methods and Exploratory Data Analysis

- Exploratory data analysis (EDA) to understand the data patterns, identify missing values, outliers, relationship between features etc. to select the most suitable model.
- Experimented with various models including Support Vector Machine (SVM), Long Short Term Memory (LSTM), Simple Moving Average (SMA), AutoRegressive Integrated Moving Average (ARIMA), and Timesfm.

We also performed model validation by comparing our model's prediction with actual demand values for a period of 15 months.



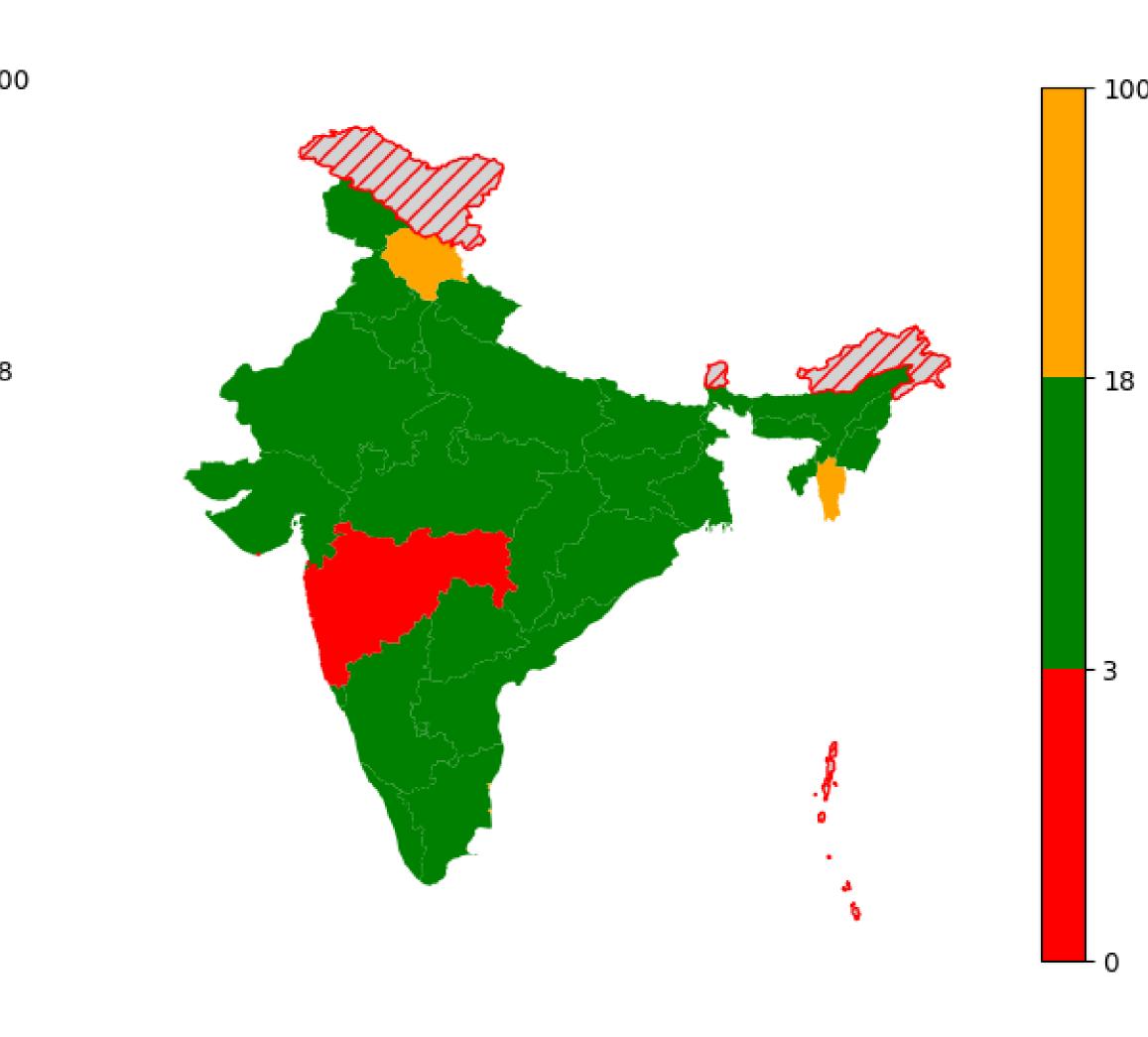


Figure 3. ABC/3TC Pediatric Stock, March

. ABC/3TC Pediatric Stock, January

Figure 2. ABC/3TC Pediatric Stock,

February 2017

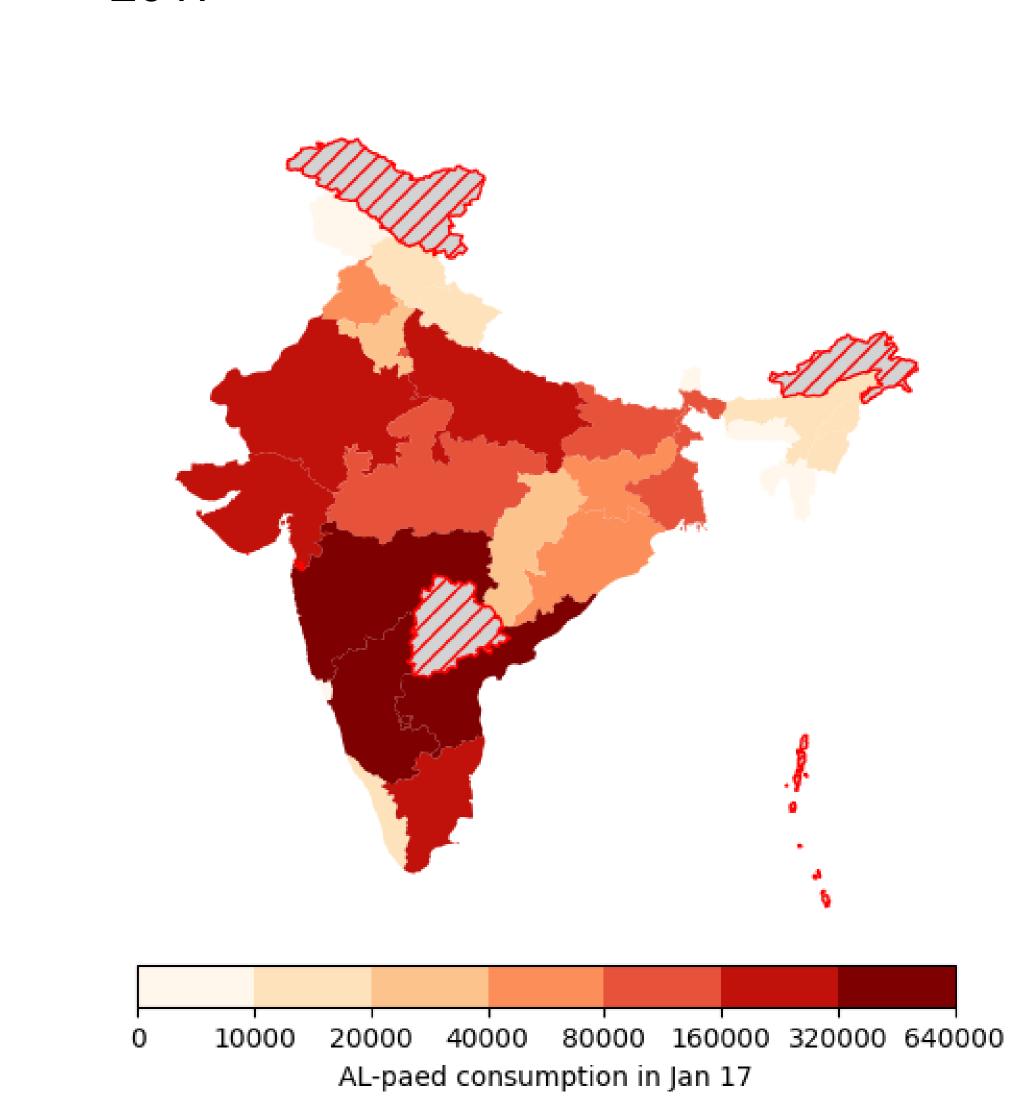


Figure 4. ABC/3TC Pediatric Consumption, January 2015

AL-paed consumption in Jan 15

10000 20000 40000 80000 160000 320000 640000

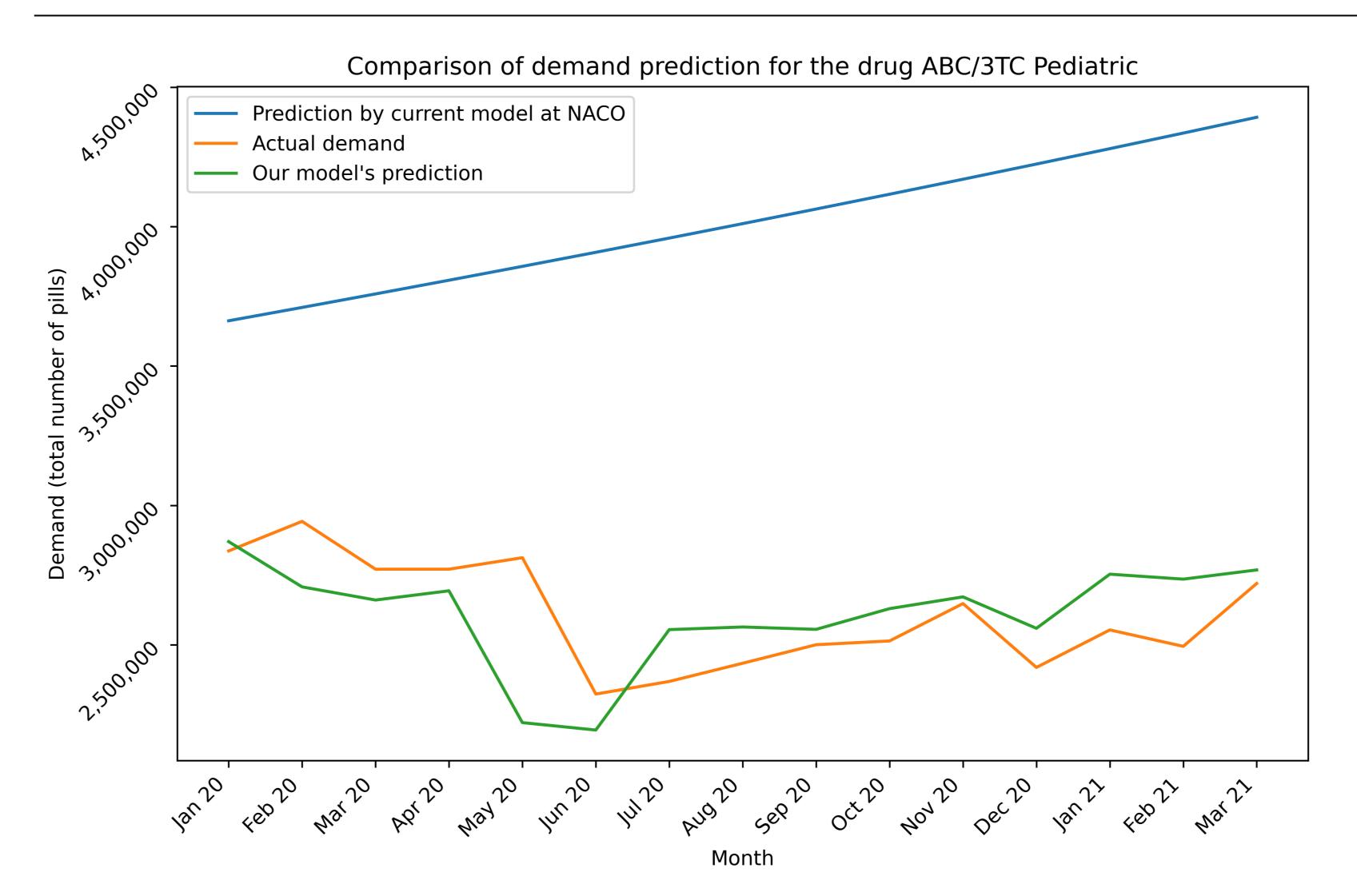
Figure 5. ABC/3TC Pediatric Consumption, January 2016

AL-paed consumption in Jan 16

20000 40000 80000 160000 320000 640000

Figure 6. ABC/3TC Pediatric Consumption, January 2017

## **Initial Results**



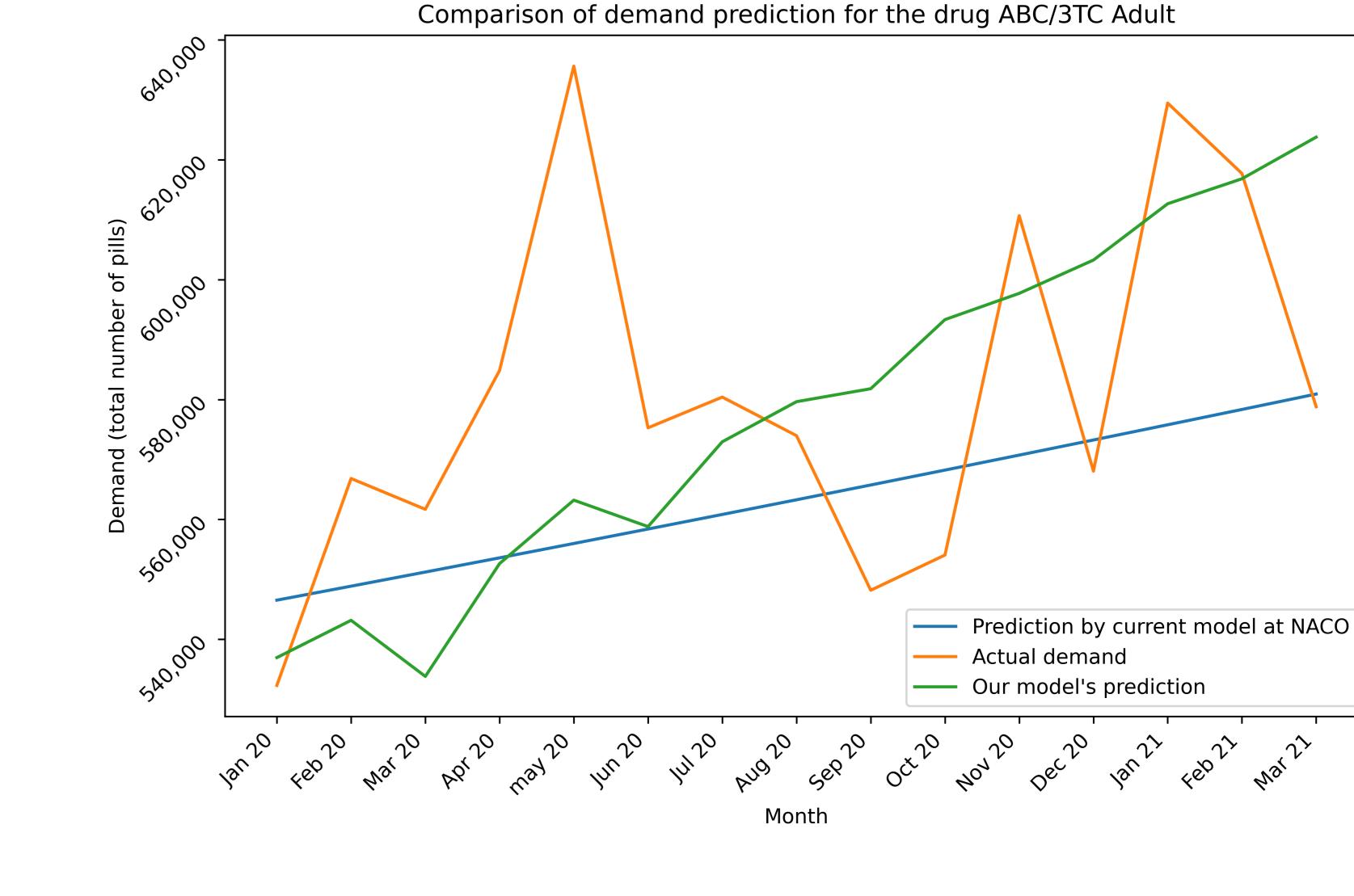


Figure 7. For ABC/3TC Pediatric: Comparison of actual demand, our model's prediction and current model at NACO

Figure 8. For ABC/3TC Adult: Comparison of actual demand, our model's prediction and current model at NACO

| Drug Name            | Actual   | Our Model | NACO     | MAPE<br>(Actual, NACO) | MAPE<br>(Actual, Our mode | % Change<br>I) (Actual, NACO) ( | % Change<br>(Actual, Our model) |
|----------------------|----------|-----------|----------|------------------------|---------------------------|---------------------------------|---------------------------------|
| ABC/3TC<br>Adult     | 8717970  | 8670610   | 8451718  | 4.15%                  | 4.28%                     | 3.05%                           | 0.54%                           |
| ABC/3TC<br>Pediatric | 39109899 | 39140357  | 60255456 | 55.12%                 | 5.90%                     | -54.07%                         | -0.08%                          |

. Model prediction error comparison: MAPE is calculated for 15 months period starting from January 2020 till March 2021

## **Conclusion and Future Directions**

Our results for ABC/3TC (both adult and pediatric) show better predictions than the current model at NACO. Specifically, for ABC/3TC pediatric, our predictions are within 0.08% while the current model over predicts by 54.07%. Our initial results suggest that to model drug demand for different drugs, we can not use the same model for all the drugs. We plan to extend our work in following ways:

- Model development and validation for all drugs.
- User interface and design and implementation, actively incorporating feedback from potential users.
- Applying our learning methodologies to enhance HIV care in India for targeted interventions including improving patient retention in care, personalized follow-up schedules and adherence support.

## Acknowledgments / Partnerships

We thank the National AIDS Control Organization (NACO) for providing us access to the necessary data and to the Johns Hopkins Gupta-Kilinsky India Institute for providing the funding to make this project possible. We also appreciate the efforts of Jaee and Saumya, who helped us clean and preprocess the data.