https://doi.org/10.1038/s41588-025-02153-x

Mapping genetic diversity with the GenomeIndia project

Chandrika Bhattacharyya, Krithika Subramanian, Bharathram Uppili, Nidhan K. Biswas, Shweta Ramdas, Karthik Bharadwaj Tallapaka, Prathima Arvind, Khader Valli Rupanagudi, Arindam Maitra, Tulasi Nagabandi, Tiyasha De, KuldeepSingh, Praveen Sharma, Nanaocha Sharma, Sunil K. Raghav, Punit Prasad, E. V. Soniya, Abdul Jaleel, Shijulal Nelson Sathi, Madhvi Joshi, Chaitanya Joshi, Mayurika Lahiri, Santosh Dixit, L. S. Shashidhara, Nachimuthu Senthil Kumar, H. Lalhruaitluanga, Lal Nundanga, Venkataram Shivakumar, Ganesan Venkatasubramanian, Naren P. Rao, Mohd Ashraf Ganie, Imtiyaz Ahmad Wani, Ganganath Jha, Ashwin Dalal, Murali Dharan Bashyam, Pritish Kumar Varadwaj, Sanjeev BS, Yogesh Simmhan, Chirag Jain, Durai Sundar, Ishaan Gupta, Pankaj Yadav, Himanshu Sinha, Manikandan Narayanan, Karthik Raman, Raghu Padinjat, Radhakrishnan Sabarinathan, GenomeIndia Consortium, Yadati Narahari, Vijayalakshmi Ravindranath, Thangaraj Kumarasamy, Divya Tej Sowpati, Mohammed Faruq, Analabha Basu & Bratati Kahali

The rich ethnolinguistic and sociocultural differences that exist in India offers a unique opportunity to study human diversity. With the whole genomes of 10,000 healthy and unrelated Indians from 83 populations, the GenomeIndia project captures the genetic diversity of one of the highly underrepresented populations in the global genomics landscape.

India, the most populous nation in the world with 1.44 billion people from over 4,600 distinct endogamous groups, represents an invaluable source of ethnolinguistic, sociocultural and genetic diversity. Deciphering the genetic makeup of human populations helps in understanding their ancestry, evolutionary history, admixture patterns, disease susceptibility and drug response. To this end, world over, genomics consortia have been established. Nevertheless, the global genomic landscape is predominantly Euro-centric¹; although some projects have documented worldwide genetic diversity by discovering DNA sequence variation in several human populations, for example the 1000 Genomes project², followed by others across the world, recent inclusions being at biobank scale^{3–7}. Unfortunately, India has been severely underrepresented in these global studies. To address this lacuna, a few studies have been initiated to understand the human genetic diversity of Asia (GenomeAsia)⁸ and India (IndiGen and LASI-DAD)^{9,10}.

India experienced several waves of migration^{11,12} with evidence of population expansion during the Paleolithic and Neolithic periods. The demographic history of India is also unique, with exposure to different environments and widespread admixture, followed by formation of numerous endogamous groups¹³. Many of the contemporary Indian populations have originated from a few founding groups and have maintained distinct identities through centuries of endogamy. These

endogamous groups^{13,14}, apart from sharing genetic similarities, possess unique variations, including distinct disease-causing mutations with amplified frequencies within specific groups.

Check for updates

Overview of the project

In light of the above, the GenomeIndia national consortium was constituted in late 2017, comprising 20 institutions, with the Centre for Brain Research (CBR) at the Indian Institute of Science in Bangalore as the coordinating center. The GenomeIndia project was launched in January 2020, with funding from the Department of Biotechnology, Ministry of Science and Technology, Government of India.

The GenomeIndia Project is the largest, as well as the most comprehensive and well-designed study of India's genetic diversity. Compared with previous studies, the sampling strategy of GenomeIndia is extensive, nuanced and balanced, with respect to ethnic, socio-cultural, geographic, biogeographic and linguistic diversity of India. An emphasis on isolated populations from all corners of the country provided a robust representation of India's genetic landscape. Blood samples and associated phenotype data were collected from over 20,000 individuals, representing 83 population groups. These comprised 30 tribal and 53 non-tribal populations spread across India (Fig. 1). The tribal populations tend to be smaller groups (<1 million), some of which are isolated and restricted to certain geographical regions. Our sampling strategy was engineered at balancing between sampling from many populations and comprehensive sampling from fewer populations, yet with appropriate representation of population diversity. With approximately 100 samples per group, we aimed to estimate relatively rare alleles (at least 1% minor allele frequencies, MAFs), which are important to understand complex diseases. We focused on genetic diversity by sampling unrelated individuals, limiting the inclusion of first- and second-degree relatives, and emphasized trios (parent-child pairs), thereby ensuring accurate estimation of allele frequencies across groups. In addition, samples are available in the established biobank along with blood biochemistry, and anthropological, and socio-economic data for each participant, setting it apart from earlier studies. The objective was to

create a compendium of genomic variations, with regard to common and rare variants to gain insights into the genetic diversity, draw pharmacogenetic inferences, and construct a reference haplotype panel to facilitate imputation and therefore genotype–phenotype association studies for the Indian population.

Enrolment of subjects and sample collection

The GenomeIndia consortium followed the Helsinki Declaration for research protocols, sample collection, and ethical conduct, which was adapted across all the sampling centers (Fig. 1) and approved by their respective institutional human ethics committees. The spread of the language families in India overlaps with geography, and language is also an established proxy for genetic diversity in the Indian population^{12,13}. We thus ensured that the four large major language families of Indo-European, Dravidian, Austro-Asiatic and Tibeto-Burman are appropriately represented in our sampling (Fig. 1), while accounting for the size of these ethnolinguistic groups. Moreover, within a broad geographic region, we sampled populations that belong to the distinct bio-geographies. To ensure that we do not miss out on the large array of rare variants in each population group, we sequenced a median of 159 unrelated individuals from each non-tribal group and 75 from each tribal group. We also included 3-6 trios in each group, for haplotype phasing and imputation, inference of de novo variations, and to assess the quality of called variants.

A total of 20,459 self-declared healthy adults (>18 years) without any diagnosed monogenic diseases or chromosomal abnormalities consented to participate in the study (Fig. 2a). Informed written consent, in English and/or their native language, was obtained from all participants. The consenting individuals underwent detailed anthropometric assessments, such as height, weight, hip circumference, waist circumference and blood pressure, after assigning 12-digit IDs as part of the de-identification process. Blood samples (10 ml each) were collected from these individuals belonging to 83 populations, which inhabit over 100 distinct geographical locations (Fig. 1). Aliquots of samples were subjected to complete blood counts, as well as biochemical investigations that encompass glucose measures, lipid profiles, and liver and kidney function tests (Fig. 2b). This repository of samples with associated biomedical information distinguishes GenomeIndia from other population-scale projects. An important feature of this project is thus the biobank of over 20,000 samples that are available for future research.

Genotyping and WGS

Initially, 13,242 samples were genotyped on genome-wide SNP array to identify the subset used for whole-genome sequencing (WGS), primarily by checking for relatedness within a specific population. To maximize diversity, in addition to the selected trios, individuals in a population group who are unrelated to each other beyond first cousins were

chosen for WGS (Fig. 2a). In total, DNA samples from 10,074 individuals were subjected to WGS. The four primary sequencing centers (CBR, CSIR-CCMB, CSIR-IGIB and BRIC-NIBMG) and one satellite sequencing center (GBRC) performed WGS, following uniform protocols in benchmarking and analysis, to ensure consistency and reliability. Subsequent quality checks were performed at the four primary sequencing centers (Fig. 2c). Variant quality was benchmarked with Genome in a Bottle (GIAB) samples, showing recall above 0.97, precision above 0.99, and an F1 score above 0.98 for single nucleotide variants (SNVs) and insertions and deletions (indels) for all sequencing centers. Five DNA samples, comprising one trio and two unrelated samples from the same population group, were sequenced and analyzed in a pairwise manner between two institutes. This process was repeated cyclically and cross-validated across four primary sequencing centers and achieved precision and recall of 0.99 and 0.97, respectively. The high concordance of genotype calls examined across sequencing centers suggests negligible batch effects. As the sequencing was conducted in five different centers, to eliminate any remaining possible batch effects, despite the above precautionary measures, we used additional tests on individual gVCF files before they were considered for joint genotyping (Fig. 2c). Our current genetic variant call-set is derived from 9,772 individuals - 4,696 male participants and 5,076 female participants.

Data handling and processing

To fasttrack a project of this magnitude, we distributed uniform protocols for raw data quality control, alignment and variant calling across the four primary sequencing centers (Fig. 1). In a first for the nation, each of these centers processed approximately 2,500 samples (totaling over 1 petabyte of raw data per center) from fastq to single-sample gVCF. Overall, this process consumed upwards of 0.7 million CPU hours and 4.5 petabytes of storage at the centers. These data are housed at a dedicated server at the Indian Biological Data Centre (IBDC), Regional Centre for Biotechnology (RCB), Faridabad, India.

Preliminary findings after joint genotyping

The total number of raw genetic variants obtained were approximately 180 million, comprising biallelic SNVs and short indels, constituting biallelic singletons, doubletons and a small proportion of multi-allelic variants over autosomes, sex chromosomes and mitochondrial genome. After another round of quality checks for variants and genotypes at the sample level, after joint genotyping, 9,772 samples were retained for further consideration. This resulted in around 130 million autosomal variants. Excluding singletons and doubletons, as expected, most variants (65%) are ultra-rare, with a MAF of less than 0.1% in the overall population. Together, the identified genetic variants decode extensive genetic diversity that has been hitherto uncaptured in the Indian population.

Fig. 1| Participating centers and populations sampled in the GenomeIndia project. The background colors of the map represent the distribution of the four major language families. Locations, language families and social groups associated with the samples collected are depicted on the map. Map source: https://surveyofindia.gov.in/pages/downloads, downloaded in December 2024. Twenty participating institutes across the nation, along with their responsibilities highlighted by a color code, are as follows. AIIMS-J, All India Institute of Medical Sciences, Jodhpur; BRIC-CDFD, Centre for DNA Fingerprinting and Diagnostics; BRIC-IBSD, Institute of Bioresources and Sustainable Development; BRIC-ILS, Institute of Life Sciences; BRIC-NIBMG, National Institute of Biomedical Genomics; BRIC-RGCB, Rajiv Gandhi Centre for Biotechnology; CBR, Centre for Brain Research (coordinating center); CSIR-CCMB, Centre for Cellular and Molecular Biology; CSIR-IGIB, Institute of Genomics and Integrative Biology; GBRC, Gujarat Biotechnology Research Centre (A Satellite Sequencing Centre); IISER-Pune, Indian Institute of Science Education and Research, Pune; IIITA, Indian Institute of Information Technology, Allahabad; IISc, Indian Institute of Science; IITD, Indian Institute of Technology, Delhi; IITJ, Indian Institute of Technology, Jodhpur; IITM, Indian Institute of Technology, Madras; MZU, Mizoram University; NCBS, National Centre for Biological Sciences; NIMHANS, National Institute of Mental Health and Neurosciences; SKIMS, Sher-i-Kashmir Institute of Medical Science.

The future ahead

Capacity-building and sharing digital public data. The successful execution of the GenomeIndia project has helped us to develop a robust ethical and legal framework to collect and govern the use of genomic data. A highlight of this project is the creation of digital public data, which will be available at IBDC and accessible for approved research in the foreseeable future. Similarly, the biobank containing over 20,000

samples will be a rich resource for multidisciplinary research. The project witnessed a synergy among 20 institutions with rigorous and uniform practices for consent, sampling, large-scale data generation and analysis. The capacity built in informatics infrastructure has provided a sustainable strategy that can form the basis for future genomics consortia in India, and contribute to South Asian specific genetic variations in global databases.



nature genetics



Fig. 2| Sample processing workflow and quality checks. a, Volunteer recruitment, sample collection and processing workflow in the GenomeIndia project.
b, Phenotype data collected from participants in the GenomeIndia project. Details included anthropometry, and various blood tests to study their hematological parameters, glucose and lipid metabolism markers, kidney and liver function indices. BP, blood pressure; DLC, differential leukocyte count; SGOT, serum glutamic oxaloacetic transaminase; SGPT, serum glutamic pyruvic transaminase.

c, Quality checks (QCs) undertaken to ensure uniformity and reliability across the sequencing centers. Top, four (GIAB) cell lines were sequenced at each center and compared with the Genome in a Bottle (GIAB) truth set provided by National Institute of Standards and Technology (NIST). Bottom, five DNA samples (1 trio, 2 unrelated) from a randomly chosen population group were selected and sent to another sequencing center in a cyclical manner. The variants for each of these five samples were compared for concordance. Pr, precision; re, recall.

Identifying genetic diversity and uniqueness in Indian context. A direct outcome of the GenomeIndia data is decoding extensive genetic diversity from Indian populations that has been uncaptured until now. However, these data will also propel several ongoing and future compelling research questions and hypotheses. These questions include understanding the contribution of genetic variations to drug metabolism and modulation of drug responses, delving into patterns of natural selection that act on the genes and regulatory regions of the genome, as well as deciphering the unique regions of tandem repeats and transposable elements in the Indian population. The identified variants will enable us to create a larger set of benign innocuous variants that can be eliminated from candidate variant lists in clinical cases. Models that assess polygenic score (PGS) - an estimate of individual level genetic risk for susceptibility to a disease based on their genetic makeup - currently rely on allele frequencies of Euro-centric data and hence suffer from accuracy and portability issues in Indian population. The population-specific linkage disequilibrium structure unraveled by GenomeIndia data will enable training of better PGS models applicable to Indians, and perhaps South Asians in general. Overall, such research, which is being addressed by several working groups, will help to deepen our understanding of the genetic basis of different diseases. This will help to address health disparities that arise owing to a combination of sociocultural and environmental variations in various Indian population groups. Analysis towards these goals is ongoing, and we are actively preparing a detailed manuscript based on the results obtained so far and further insights anticipated in the near future.

Genome-wide arrays for South Asian ancestry. Our preliminary data indicate that imputing genotypes from Indian or South Asian ancestry, leveraging the reference haplotype panel constructed with rare and common variants obtained from GenomeIndia, will yield higher imputation power and accuracy, as well as allelic concordance compared with other widely used worldwide datasets, such as the HRC (Haplotype Reference Consortium) and TOPMed (Trans-Omics for Precision Medicine). This is owing to the benefits of ancestry-matched genomic linkage disequilibrium patterns, especially for rare and low-frequency variants in the Indian population. We are embarking on designing a genotype array for the Indian population, which we expect will notably improve on the recently developed SARGAM (South Asian Research Genotyping Array for Medicine)¹⁵ array. Our focus is on the array content being tailored for individuals of Indian ancestry and will be centered around the discovery of common and rare genetic variants that is possible owing to the nuanced sampling in the GenomeIndia project. We foresee that this effort will successfully capture genetic variation across diverse population groups and facilitate future large-scale genetic association studies in the country. Importantly, many of these common variants are rare or non-existent in global variant databases, thereby highlighting the utility of this variant resource for clinical diagnostics and related applications in Indian and South Asian populations.

Concluding remarks

The GenomeIndia project is poised to provide valuable resources to advance medical genetics, enabling further clinical research. In-depth analysis of 9,772 diverse genomes along with the blood biochemistry and anthropometry data will improve disease diagnostics, predict the genetic basis of drug responses, and kickstart precision medicine efforts in India. This reference dataset could aid in the diagnostic process of rare diseases, potentially saving time and money, helping patients to circumvent an arduous road to healthcare. The findings from the project are also expected to inspire future functional studies that aim to determine the possible roles of population-specific genetic variants in various diseases and aiding in interventional and translational research in India.

Chandrika Bhattacharyya D^{1,24}, Krithika Subramanian^{2,3,24}, Bharathram Uppili⁴, Nidhan K. Biswas ¹, Shweta Ramdas², Karthik Bharadwaj Tallapaka⁵, Prathima Arvind², Khader Valli Rupanagudi², Arindam Maitra ¹, Tulasi Nagabandi⁵, Tiyasha De⁴, Kuldeep Singh⁶, Praveen Sharma⁶, Nanaocha Sharma⁷, Sunil K. Raghav⁸, Punit Prasad⁸, E. V. Soniya⁹, Abdul Jaleel⁹, Shijulal Nelson Sathi⁹, Madhvi Joshi¹⁰, Chaitanya Joshi¹⁰, Mayurika Lahiri¹¹, Santosh Dixit¹¹, L. S. Shashidhara¹¹, Nachimuthu Senthil Kumar **D**¹², H. Lalhruaitluanga 🕲 ¹², Lal Nundanga¹², Venkataram Shivakumar¹³, Ganesan Venkatasubramanian¹³, Naren P. Rao¹³, Mohd Ashraf Ganie¹⁴, Imtiyaz Ahmad Wani¹⁴, Ganganath Jha¹⁵, Ashwin Dalal¹⁶, Murali Dharan Bashyam¹⁶, Pritish Kumar Varadwaj¹⁷, Sanjeev BS¹⁷, Yogesh Simmhan¹⁸, Chirag Jain¹⁸, Durai Sundar ¹⁹, Ishaan Gupta¹⁹, Pankaj Yadav²⁰, Himanshu Sinha ¹, Manikandan Narayanan²¹, Karthik Raman²¹, Raghu Padinjat **D**²², Radhakrishnan Sabarinathan 10²², GenomeIndia Consortium*, Yadati Narahari^{2,18,25}, Vijayalakshmi Ravindranath^{2,25}

Thangaraj Kumarasamy 5,25 , Divya Tej Sowpati 5,25 , Mohammed Faruq 4,25 , Analabha Basu 1,25 & Bratati Kahali 2,25

¹BRIC - National Institute of Biomedical Genomics (BRIC-NIBMG), Kolkata, India. ²Centre for Brain Research (CBR), IISc Campus, Bengaluru, India. ³Manipal Academy of Higher Education, Karnataka, India. ⁴CSIR - Institute of Genomics & Integrative Biology (CSIR-IGIB), New Delhi, India. ⁵CSIR - Centre for Cellular and Molecular Biology (CSIR-CCMB), Hyderabad, India. ⁶All India Institute of Medical Sciences (AIIMSJ), Jodhpur, India. 7BRIC - Institute of Bioresources & Sustainable Development (BRIC-IBSD), Imphal, India, ⁸BRIC - Institute of Life Sciences (BRIC-ILS), Bhubaneswar, India, ⁹BRIC - Raiiv Gandhi Centre for Biotechnology (BRIC-RGCB), Thiruvananthapuram, India. ¹⁰Gujarat Biotechnology Research Centre (GBRC), Gandhinagar, India. ¹¹Indian Institute of Science Education and Research (IISER), Pune, India. ¹²Mizoram University (MZU), Aizawl, India. ¹³National Institute of Mental Health & Neurosciences (NIMHANS), Bengaluru, India. ¹⁴Sher-i-Kashmir Institute of Medical Sciences (SKIMS), Srinagar, India.¹⁵Vinoba Bhave University, Hazaribagh, India.¹⁶BRIC - Centre for DNA Fingerprinting and Diagnostics (BRIC-CDFD), Hyderabad, India.¹⁷Indian Institute of Information Technology (IIITA), Allahabad, India.¹⁸Indian Institute of Science (IISc), Bengaluru, India.¹⁹Indian Institute of Technology Delhi (IITD), New Delhi, India.²⁰Indian Institute of Technology Jodhpur (IITJ), Jodhpur, India.²¹Indian Institute of Technology Madras (IITM), Chennai, India.²²National Centre for Biological Sciences (NCBS), Bengaluru, India.²⁴These authors contributed equally: Chandrika Bhattacharyya, Krithika Subramanian. ²⁵These authors jointly supervised this work: Yadati Narahari, Vijayalakshmi Ravindranath, Kumarasamy Thangaraj, Divya Tej Sowpati, Mohammed Faruq, Analabha Basu, Bratati Kahali. *A list of authors and their affiliations appears at the end of the paper. e-mail: thangs@ccmb.res.in; tej@ccmb.res.in; faruq.mohd@igib.in; ab1@nibmg.ac.in; bratati@cbr-iisc.ac.in

Published online: 08 April 2025

References

- 1. Popejoy, A. B. & Fullerton, S. M. Nature **538**, 161–164 (2016).
- 2. The 1000 Genomes Project Consortium. Nature **467**, 1061–1073 (2010).
- 3. Okada, Y. et al. Nat. Commun. 9, 1631 (2018).
- Choudhury, A. et al. *Nature* 586, 741–748 (2020).
 Halldorsson, B. V. et al. *Nature* 607, 732–740 (2022).
- Halldorsson, B. V. et al. Nature 607, 732–740 (2022)
 Walters, R. G. et al. Cell Genom. 3, 100361 (2023).
- Walters, R. G. et al. Cell Genom. 3, 100361 (2023).
 Bianchi, D. W. et al. Nat. Med. 30, 330–333 (2024)
- Blanchi, D. W. et al. Nat. Med. **30**, 330–333 (2024).
 GenomeAsia100K Consortium. Nature **576**, 106–111 (2019).
- Jain, A. et al. Nucleic Acids Res. 49. D1225–D1232 (2021).
- Jain, A. et al. Nucleic Acids Res. 49, D1223–D1232 (2021).
 Kerdoncuff, E. et al. Preprint at *bioRxiv* https://doi.org/10.1101/2024.02.15.580575
- (2024). 1. The second of the second second
- 11. Thangaraj, K. et al. Science **308**, 996 (2005).
- 12. Kumar, S. et al. BMC Evol. Biol. 8, 230 (2008).
- Basu, A. et al. Proc. Natl Acad. Sci. USA 113, 1594–1599 (2016).
 Nakatsuka, N. et al. Nat. Genet. 49, 1403–1407 (2017).
- Nakatsuka, N. et al. Nat. Genet. 49, 1403–1407 (
 Wall, J. D. et al. Nat. Commun. 14, 3377 (2023).

Acknowledgements The GenomeIndia project was

The GenomeIndia project was funded by the Department of Biotechnology, Ministry of Science and Technology, Government of India (BT/GenomeIndia/2018). The consortium thanks all the participants who volunteered for this study. We extend our sincere gratitude to the dedicated social workers and government agencies at local and regional levels whose unwavering support has been instrumental in the successful collection of samples for the GenomeIndia project. We also thank the Illumina team of joint genotyping for help with the joint genotyper program. We thank the reviewers for their helpful

Competing interest

The authors have no conflict of interest to declare.

Additional information

suggestions for revising the article.

Peer review information Nature Genetics thanks Mark Caulfield and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

GenomeIndia Consortium

Conceptualization and project writing

Analabha Basu^{1,25}, Vijayalakshmi Ravindranath^{2,25}, Bratati Kahali^{2,25}, Kumarasamy Thangaraj⁵, Mitali Mukerji⁴, Sridhar Sivasubbu⁴, Vinod Scaria⁴, Sunil K. Raghav⁸ & Himanshu Sinha²¹

Project coordination and administration group

Suman K. Paine¹, Chandrika Bhattacharyya^{1,24}, Arindam Maitra¹, Nidhan K. Biswas¹, Analabha Basu^{1,25}, Vijayalakshmi Ravindranath^{2,25}, Bratati Kahali^{2,25}, Prathima Arvind², Yadati Narahari², Karthik Bharadwaj Tallapaka⁵, Divya Tej Sowpati^{5,25}, Govindarajan Umapathy⁵, Vinay K. Nandicoori⁵, Rakesh Mishra⁵, Kumarasamy Thangaraj⁵, Mohammed Faruq^{4,25}, Sridhar Sivasubbu⁴, Vinod Scaria⁴, Kuldeep Singh⁶, Praveen Sharma⁶, Nanaocha Sharma⁷, Dinabandhu Sahoo⁷, Sunil K. Raghav⁸, Ajay Parida⁸, E. V. Soniya⁹, Abdul Jaleel⁹, M. Radhakrishna Pillai⁹, Abitha Thomas⁹, Madhvi Joshi¹⁰, Chaitanya Joshi¹⁰, L. S. Shashishara¹¹, Mayurika Lahiri¹¹, Nachimuthu Senthil Kumar¹², H. Lalhruaitluanga¹², Naren P. Rao¹³, Ganesan Venkatasubramanian¹³, Venkataram Shivakumar¹³, Imtiyaz Ahmad Wani¹⁴, Mohd Ashraf Ganie¹⁴, Murali Dharan Bashyam¹⁶, Ashwin Dalal¹⁶, Pritish Kumar Varadwaj¹⁷, Sanjeev BS¹⁷, Y. Narahari^{2,18}, Yogesh Simmhan¹⁸, Arun Kumar¹⁸, Durai Sundar¹⁹, B. Jayaram¹⁹, Ishaan Gupta¹⁹, Pankaj Yadav²⁰, Himanshu Sinha²¹, Karthik Raman²¹, Manikandan Narayanan²¹, Padinjat Raghu²² & Radhakrishnan Sabarinathan²²

Sample collection group

Azad Ali¹, Mahabub Alam¹, Parveena Choudhury¹, Poulomi Ghosh¹, Sukanya Dhar¹, Saurav Roy¹, Nasrin Parvin¹, Rahul Modak¹, Sayan Bhowmick¹, Sourav Gangopadhyay¹, Devashish Tripathi¹, Analabha Basu^{1,25}, Chandrika Bhattacharyya^{1,24}, Suman K. Paine¹, Prathima Arvind², K. S. H. Shafeeq², G. Rajesh², C. Mohana², A. Divakar², Reddy P. Kommaddi², Neha Singh⁵, Priya Pandey⁵, Devavrat Desai⁵, Mahfuj Hassan⁵, Deepak Kumar Kashyap⁵, Vasantha Kumar⁵, Aman Kumar Suryan⁵, Hema Sindhuja Rachiraju⁵, A. Mahesh⁵, Sushmita Nitta⁵, Vijaya Mohan⁵, Karthikeyan Meenakshisundaram⁵, Jagamohan Chhatai⁵, G. Mala⁵, Karthik Bharadwaj Tallapaka⁵, Kumarasamy Thangaraj⁵, Sandeep Kumar Pal⁴, Simmy Kaur⁴, Mahino Fatima⁴, Mohammed Akbar⁴, Rahul C. Bhoyar⁴, Pooja Sharma⁴, Shreya Bari⁴, Tiyasha De⁴, Pratima Pandey⁴, Anushree Mishra⁴, Nishat Ashrafi⁴, Syed Ahmad⁴, Deepak Mudila⁴, Sridhar Sivasubbu⁴, Vinod Scaria⁴, Mohammed Farug^{4,25}, Arun Sree Parameswaran⁶, Dolat Singh Shekhawat⁶, Kuldeep Singh⁶, Nayan Tada⁶, Praveen Sharma⁶, Tanuja Rajial⁶, Varuna Vyas⁶, Arvinda Thoudam⁷, H. Moushmi Sharma⁷, Khuraijam Dolly Devi⁷, Nanaocha Sharma⁷, Teresa Tangpua⁷, Adyasha Mishra⁸, Arup Ghosh⁸, Deepak Jena⁸, Punit Prasad⁸, Soumendu Mahapatra⁸, Sudarshana Jena⁸, Sudeshna Datta⁸, Sunil K. Raghav⁸, E. V. Soniya⁹, Abdul Jaleel⁹, Shijulal Nelson Sathi⁹, M. Radhakrishna Pillai⁹, Abhitha Thomas⁹, Udaya Lekshmi⁹, R. A. Aswanth⁹, Anjana S. Nair⁹, Vasudev Paveri⁹, T. S. Amal⁹, Aman Tripathi¹⁰, Bhagirath Dave¹⁰, Bhumika Prajapati¹⁰, Chaitanya Joshi¹⁰, Madhvi Joshi¹⁰, Ramesh Pandit¹⁰, Sanman Samova¹⁰, Ajay Malik¹¹, Kajal Gaikwad¹¹, L. S. Shashishara¹¹, Mayurika Lahiri¹¹, Santosh Dixit¹¹, Siddharth Gahlaut¹¹, Andrew Vanlallawma¹², H. Lalhruaitluanga¹², John Zohmingthanga¹², Lalawmpuii Pachuau¹², Lalchhandama Chhakchhuak¹², Lalnundanga¹², Nachimuthu Senthil Kumar¹², Ranjan Jyoti Sarma¹², Daddaladka Krishnayya Samartha¹³, Ganesan Venkatasubramanian¹³, Naren P. Rao¹³, Paranthaman V. Kavya¹³,

S. G. Tejaswini¹³, Venkataram Shivakumar¹³, Bashir Ahmad Charoo¹⁴, Imtiyaz Ahmad Wani¹⁴, Mahrukh Hameed Zargar¹⁴, Mohd Ashraf Ganie¹⁴ & Ganganath Jha¹⁵

Biobanking group

C. Mohana², G. Rajesh², A. Divakar², K. H. Rakesh² & Shobha Anilkumar²

Sequencing group

Mahabub Alam¹, Parveena Choudhury¹, Azad Ali¹, Poulomi Ghosh¹, Rahul Modak¹, Sukanya Dhar¹, Nasrin Parvin¹, Sayan Bhowmick¹, Sourav Roy¹, Shouvanik Sengupta¹, Chandrika Bhattacharyya^{1,24}, Indranil Bagchi¹, Subrata Patra¹, Suman K. Paine¹, Arindam Maitra¹, Khader Valli Rupanagudi², M. H. K. Mujawar², Vinayak Hosawad², Tulasi Nagabandi⁵, Valli Undamatla⁵, Neha Singh⁵, Priya Pandey⁵, Mahfuj Hassan⁵, Vasantha Kumar⁵, Devavrat Desai⁵, Pratheusa Maccha⁵, Sushmita Nitta⁵, Aman Kumar Suryan⁵, Hema Sindhuja Rachiraju⁵, Karthik Bharadwaj Tallapaka⁵, Pooja Sharma⁴, Tiyasha De⁴, Shreya Bari⁴, Shahrumi Reza⁴, Divya Goel⁴, Rahul C. Bhoyar⁴, Sandeep Kumar Pal⁴, Bharathram Uppilli⁴, Arushi Batra⁴, Ashvarya Shankar⁴, Gayatri Singh⁴, Suman Mudila⁴, Mahino Fatima⁴, Divya Goel⁴, Saima Iram⁴, Mohamed Imran⁴, Mohit Divakar⁴, Vigneshwar Senthivel⁴, Sridhar Sivasubbu⁴, Vinod Scaria⁴ & Mohammed Faruq^{4,25}

Analysis working group

Chandrika Bhattacharyya^{1,24}, Devashish Tripathi¹, Vinay More¹, Arghya Dey¹, Shouvanik Sengupta¹, Haya Afreen¹, Mahabub Alam¹, Parveena Choudhury¹, Saurav Roy¹, Animesh Kumar Singh¹, Arnab Ghosh¹, Chitrarpita Das¹, Debashree Tagore¹, Subrata Das¹, Suman K. Paine¹, Nidhan K. Biswas¹, Analabha Basu^{1,25}, Krithika Subramanian², Shreya Chakraborty², Raghvendra Agrawal², Sauma Suvra Majumdar², Siddhi Jani², Akkshaya Rajesh², Debasrija Mondal², Anand Kumar², Debdutta Chatterjee², Priyanka Singh², A. Sohan Angelo², Tanmay Panigrahi², Eric Macwan², Rupanwita Majumder², S. Sagar², Samarpita Saha², Shweta Ramdas², Bratati Kahali^{2,25}, Payel Mukherjee⁵, Pratheusa Maccha⁵, Sreelekshmi MS⁵, Jayesh Jain⁵, Sofia Banu⁵, Malini Nemalikanti⁵, Sriram Sudarsanam⁵, Divya Tej Sowpati^{5,25}, Shreya Bari⁴, Bharathram Uppilli⁴, Ankit Mukerji⁴, Bani Jolly⁴, Jupita Handique⁴, Tiyasha De⁴, Pooja Sharma⁴, Sridhar Sivasubbu⁴, Vinod Scaria⁴ & Mohammed Faruq^{4,25}

Joint genotyping group

Chandrika Bhattacharyya^{1,24}, Saurav Roy¹, Arnab Ghosh¹, Analabha Basu^{1,25}, Nidhan K. Biswas¹, Krithika Subramanian², Raghvendra Agrawal², Shreya Chakraborty², Siddhi Jani², Anand Kumar² & Bratati Kahali^{2,25}

Data management group (data organization, data quality control, data storage, data archival, data security, data curation)

Chandrika Bhattacharyya^{1,24}, Saurav Roy¹, Shouvanik Sengupta¹, Animesh Kumar Singh¹, Arghya Dey¹, Arnab Ghosh¹, Azad Ali¹, Devashish Tripathi¹, Haya Afreen¹, Mahabub Alam¹, Parveena Choudhury¹, Vinay More¹, Suman K. Paine¹, Analabha Basu^{1,25}, Nidhan K. Biswas¹, Krithika Subramanian², Anand Kumar², Siddhi Jani², Shreya Chakraborty², Debasrija Mondal², V. Jothibasu², S. Karthik², Shweta Ramdas², Bratati Kahali^{2,25}, Sreelekshmi MS⁵, Payel Mukherjee⁵, Sriram Sudarsanam⁵, Pratheusa Maccha⁵, Jayesh Jain⁵, Divya Tej Sowpati⁵, Bharathram Uppilli⁴, Shreya Bari⁴, Pooja Sharma⁴, Tiyasha De⁴, Shahrumi Reza⁴, Mohammed Faruq^{4,25}, Sanjay Deshpande²³, Deepak T. Nair²³ & Saurabh Raghuvanshi²³

Method development group

Murali Dharan Bashyam¹⁶, Ashwin Dalal¹⁶, Asmita Gupta¹⁶, Sumedha Avadhanula¹⁶, Imlimaong Aier¹⁷, Pritish Kumar Varadwaj¹⁷, Rahul Semwal¹⁷, B. S. Sanjeev¹⁷, Ajeya Bhat¹⁸, Arun Kumar¹⁸, Chirag Jain¹⁸, Nagakishore Jammula¹⁸, Sai Manasa Chadalavada¹⁸, Yogesh Simmhan¹⁸, Nirmal Singh Mahar¹⁹, Ishaan Gupta¹⁹, Durai Sundar¹⁹, Jyoti Sharma²⁰, Pankaj Yadav²⁰, Rajveer Singh Shekhawat²⁰, Soham Biswas²⁰, Ayam Gupta²¹, Harshita Agarwal²¹, Himanshu Sinha²¹, Karthik Raman²¹, Manikandan Narayanan²¹, Venkatesh Kamaraj²¹, Agastya Singh²², Raghu Padinjat²² & Radhakrishnan Sabarinathan²²

²³Indian Biological Data Centre, Haryana, India.