# Enhancing Diversity in Imaging Data: A Framework for Inclusive AI Development

**Lawrence Guan, MD,** Resident Physician, Yale School of Medicine
Sophie Chheang, MD, MBA; Irene Dixe De Oliveira Santo, MD, CIIP

## Background/Problem Being Solved

Artificial intelligence (AI) has rapidly transformed radiology, with algorithms excelling in tasks from pathology detection to guiding therapeutic interventions. However, the effectiveness, fairness, and generalizability of these tools hinge on the diversity of training datasets. Despite advances, radiology datasets often lack sufficient demographic, geographic, and disease-specific diversity, reinforcing biases and limiting applicability, particularly in underrepresented populations. This study aims are to: (1) identify gaps in publicly available radiology datasets by analyzing their demographic composition, geographic distribution, and disease variability; (2) propose actionable strategies for recruiting diverse patient populations to improve dataset inclusivity.

## Intervention(s)

We analyzed 15 publicly available radiology datasets (e.g., NIH ChestX-ray14, UK Biobank) to evaluate demographic and geographic representation. Metadata, including patient age, gender, race/ethnicity, and geographic origin, were extracted and compared against global population distributions using the Dataset Representation Index (DRI). To address identified gaps, we developed recruitment strategies focused on underrepresented populations and evaluated the role of federated learning for secure international data sharing.

## Barriers/Challenges

The analysis uncovered significant demographic imbalances: rural populations, women, and racial minorities were underrepresented. Data primarily originated from North America and Europe, with limited contributions from Africa and South Asia, restricting applicability to low- and middle-income countries (LMICs).
Strategies that yield potential for mitigate these imbalances include: (1) implementation of targeted recruitment strategies and global collaborations to increase representation across key metrics, improving dataset diversity, and (2) federated learning enabled secure and privacy-preserving data sharing between international institutions, enhancing inclusion without compromising legal or ethical standards.

## Outcome

Our findings reveal the critical need to address disparities in radiology datasets to mitigate risks of biased AI systems. Collaborative efforts with global health institutions, community-based recruitment initiatives, and the adoption of federated learning frameworks are vital for achieving equitable representation. Standardized demographic reporting and bias monitoring are essential to ensure ongoing fairness and transparency in AI development.

Scientific Research & Applied Informatics Posters and Demonstrations

## Conclusion/Statement of Impact/Lessons Learned

This study highlights the importance of fostering diversity in imaging datasets to enhance the reliability and fairness of radiology AI tools. By adopting our proposed framework, stakeholders can build inclusive systems that better serve global populations and contribute to more equitable healthcare delivery.

## Keywords

Artificial Intelligence/Machine Learning; Patient/Family Experience; Quality Improvement & Quality Assurance