# Ethics in AI-Driven Personalized Medicine

## Potential Biases in Treatment Recommendations

AI systems for precision oncology frequently rely on The Cancer Genome Atlas (TCGA), which molecularly characterized over 20,000 primary cancer and matched normal samples spanning 33 cancer types. Initiated in 2006 as a joint effort between the National Cancer Institute and the National Human Genome Research Institute, TCGA generated more than 2.5 petabytes of genomic, epigenomic, transcriptomic, and proteomic data over its dozen-year run. Despite its breadth, the dataset skews toward patients treated at major U.S. academic centers and lacks granular representation from many global populations. Consequently, AI models trained on TCGA may underperform for ethnic groups or regions that did not substantially contribute samples, leading to misestimated variant frequencies and suboptimal treatment suggestions.

Social determinants of health such as socioeconomic status, environmental exposures, and access to care are not captured in the molecular assays archived by TCGA. Omitting these real world confounders forces models to attribute outcome differences solely to molecular features, amplifying risks of biased predictions. Technical artifacts also arise from heterogeneous sequencing platforms and varying annotation completeness across the 33 cancer types, potentially skewing AI's learned patterns.

## Fairness Strategies

### Diversify Recruitment

Partner with hospitals in underrepresented regions to supplement TCGA's existing 33-cancer-type cohort.

Encourage community clinics to contribute de-identified samples alongside socioeconomic metadata.

### Leverage Public Data Access Tools

Utilize the Genomic Data Commons Data Portal to harmonize new datasets with TCGA's 2.5 PB repository, ensuring consistent processing pipelines.

Develop federated-learning frameworks that respect patient privacy while enabling model training across disparate institutions.

### Augment Minority Profiles

Employ generative adversarial networks to simulate rare-variant genomic profiles, maintaining privacy and enriching minority representation.

Use oversampling techniques (e.g., SMOTE) to balance class distributions during classifier training.

Integrate Socioeconomic Features

Link genomic entries with geocoded indices of environmental risk and insurance status.

Model interactions between molecular and social determinants to reflect multifactorial disease drivers.

Conduct Algorithmic Audits and Transparency

 Regularly evaluate predictive performance metrics (sensitivity, specificity, calibration) across ethnic, age, and gender subgroups.

Publicly disclose dataset composition, known blind spots, and Pan-Cancer Atlas findings (published 2018) to inform stakeholders.