



ENHANCING BREAST CANCER RISK PREDICTION IN HISPANIC WOMEN THROUGH TRANSFER LEARNING

Zhengwei Song

Supervisor: Tian Gu, Ph.D.

April 26, 2024

Introduction

Breast Cancer in Hispanic Women

- Second leading cause of cancer death.
- 30% lower diagnosis rate compared to non-Hispanic Whites.
- Increased risk of developing aggressive tumors: HER2+ and triple-negative breast cancer.
- Lack of effective prediction tools for early prevention

Current Studies:

- Genetic and clinical factors integration has improved risk stratification in White women, while underrepresentation in studies leads to suboptimal risk prediction in Hispanic women.

Objective

Implement novel **transfer learning** models to improve risk prediction by utilizing estimates from other populations (White and Asian) to benefit Hispanic women, that is, to **enhance the accuracy of breast cancer risk prediction for Hispanic women**.

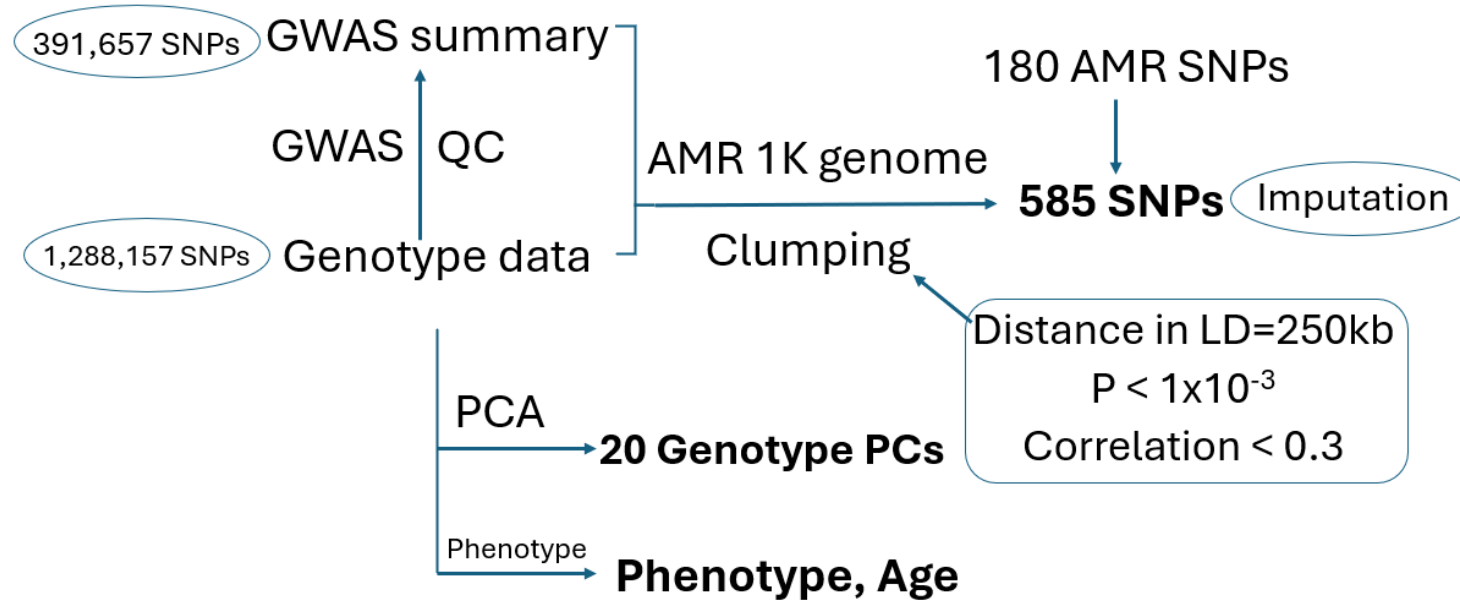
Data (Genotype & Phenotype)

- **The Multiethnic Cohort Study:** a prospective cohort study that was initiated between 1993 and 1996 and includes subjects from various ethnic groups including **Hispanics** and **Japanese-Americans** primarily from greater Los Angeles, CA
- **The Cancer Genetic Markers of Susceptibility (CGEMS) study:** a prospective cohort study by National Cancer Institute (NCI) in **European** cohort that was provided a blood sample between 1989 and 1990 and were free of diagnosed breast cancer at blood collection and followed for incident disease until May 2004

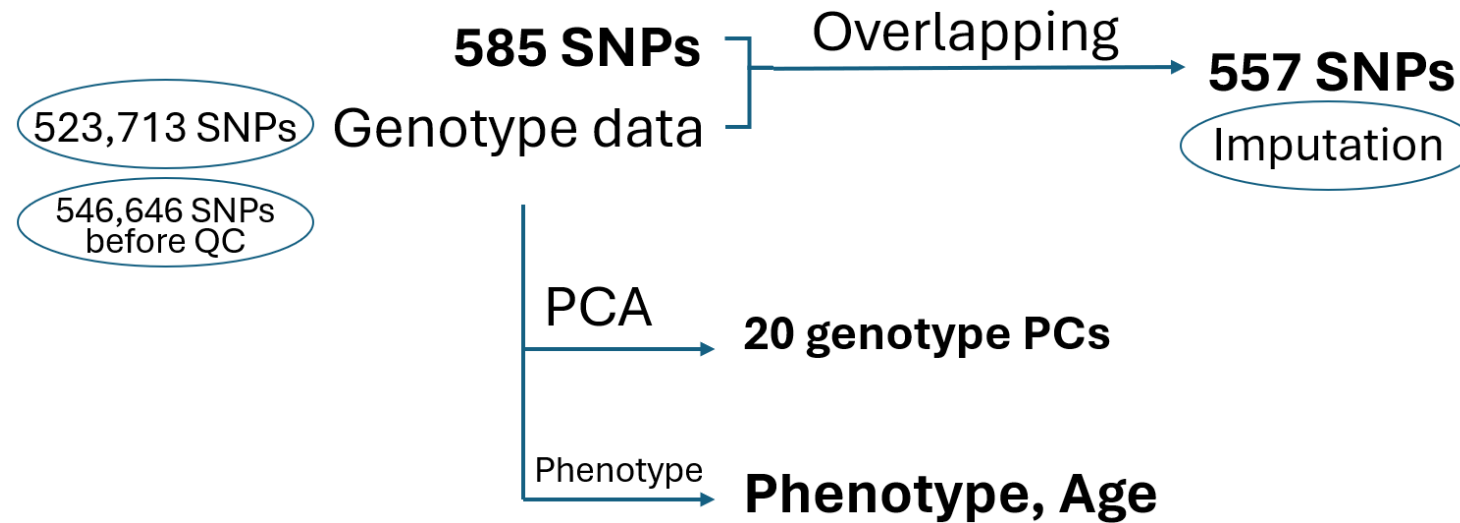


- Target Data: **n=1064 Hispanics (AMR, 520 cases and 544 controls)** descent populations from The Multiethnic Cohort Study
- Two source data: **n₁=2287 European (EUR, 1145 cases & 1142 controls)** and **n₂=1707 Japanese (JPN, 885 cases & 822 controls)** descent populations, from CGEMS and the Multiethnic Cohort Study, respectively
- Both datasets were conducted another preprocessing steps, including outcome (i.e. Breast Cancer) normalization, missing data imputation, and Linkage Disequilibrium (LD) Clumping

Data



(A) The Multiethnic Cohort Study



(B) The Nurses' Health Study

Method

Ridge

For each sample in each ethnic group, the linear model takes the form as

- Residualized Breast Cancer (y) $\sim \beta_0 + \beta_1 * \text{SNP}_1 + \beta_2 * \text{SNP}_2 + \dots + \beta_{557} * \text{SNP}_{557} + \beta_{558} * \text{PC}_1 + \beta_{559} * \text{PC}_2 + \dots + \beta_{577} * \text{PC}_{20} + \varepsilon$

Additionally, with Ridge, the final parameter takes the form as

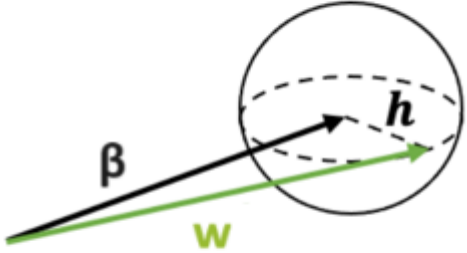
- $\hat{\beta}(\lambda) = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^{577} \beta_j x_{ij})^2 + \lambda \sum_{j=1}^{577} \beta_j^2 \right\}$, where $n=1707$ for JPN, $n=2287$ for EUR, $n=1064$ for AMR



- Source Estimates: $\hat{\beta}_{JPN}, \hat{\beta}_{EUR}$

Transfer Learning Models

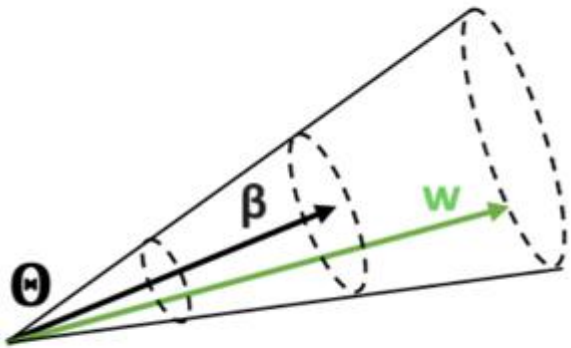
A



$$\hat{w}_{distTL} = \operatorname{argmin}_w \frac{1}{n} \|Y_k - X_k w\|_2^2 + \lambda_d \|w - \hat{\beta}_k\|^2, \text{ where } \lambda_d \text{ is the tuning parameters, } k=1 \text{ or } 2; \hat{\beta}_1 = \hat{\beta}_{JPN}, \hat{\beta}_1 = \hat{\beta}_{EUR}, \text{ or } \hat{\beta}_2 = (\hat{\beta}_{JPN}, \hat{\beta}_{EUR})$$

Distance-based similarity
measure: $\|\beta - w\|_2 \leq h$

B



Geometric illustration of the distance-based similarity characterization (A); the angle-based characterization (B).

$$\hat{w}_{angleTL} = \operatorname{argmin}_w \frac{1}{n} \|Y_k - X_k w\|_2^2 + \lambda_a \|w\|^2 - 2\eta \hat{\beta}_k^T w, \text{ where } \lambda_a \text{ and } \eta \text{ are the tuning parameters, } k=1 \text{ or } 2; \hat{\beta}_1 = \hat{\beta}_{JPN}, \hat{\beta}_1 = \hat{\beta}_{EUR}, \text{ or } \hat{\beta}_2 = (\hat{\beta}_{JPN}, \hat{\beta}_{EUR})$$

Angle-based similarity
measure: $\sin \Theta(\beta, w) \leq d$

Transfer Learning Models

- **Target-Only** Model (baseline model using AMR Only): $\hat{\mathbf{w}}_{target} = \underset{\mathbf{w}}{argmin} \frac{1}{n} \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|^2$;
➡ Testing $\hat{\mathbf{y}}'_{target} = \mathbf{X}'\hat{\mathbf{w}}_{target} + \epsilon$;
- **Source-Only** Model (directly using JPN or EUR): $\hat{\mathbf{w}}_{src} = \underset{\mathbf{w}}{argmin} \frac{1}{N_k} \|\mathbf{Y}_k - \mathbf{X}_k\mathbf{w}\|_2^2 + \lambda_k \|\mathbf{w}\|^2$,
➡ Testing $\hat{\mathbf{y}}'_{src} = \mathbf{X}'\hat{\mathbf{w}}_{src} + \epsilon$;
- Distance-based Transfer Learning (**DistTL**): $\hat{\mathbf{w}}_{DistTL} = \underset{\mathbf{w}}{argmin} \frac{1}{n} \|\mathbf{Y}_k - \mathbf{X}_k\mathbf{w}\|_2^2 + \lambda_d \|\mathbf{w} - \hat{\boldsymbol{\beta}}_k\|^2$,
➡ Testing $\hat{\mathbf{y}}'_{DistTL} = \mathbf{X}'\hat{\mathbf{w}}_{DistTL} + \epsilon$;
- Angle-based Transfer Learning (**AngleTL**): $\hat{\mathbf{w}}_{AngleTL} = \underset{\mathbf{w}}{argmin} \frac{1}{n} \|\mathbf{Y}_k - \mathbf{X}_k\mathbf{w}\|_2^2 + \lambda_a \|\mathbf{w}\|^2 - 2\eta \hat{\boldsymbol{\beta}}_k^T \mathbf{w}$
➡ Testing $\hat{\mathbf{y}}'_{AngleTL} = \mathbf{X}'\hat{\mathbf{w}}_{AngleTL} + \epsilon$;
- **Aggregated**: Testing $\hat{\mathbf{y}}'_{weighted} = w_1 * \hat{\mathbf{y}}'_{target} + w_2 * \hat{\mathbf{y}}'_{src-JPN} + w_3 * \hat{\mathbf{y}}'_{src-EUR} + w_4 * \hat{\mathbf{y}}'_{DistTL-JPN} + w_5 * \hat{\mathbf{y}}'_{DistTL-EUR} + w_6 * \hat{\mathbf{y}}'_{DistTL-J\&E} + w_7 * \hat{\mathbf{y}}'_{AngleTL-JPN} + w_8 * \hat{\mathbf{y}}'_{AngleTL-EUR} + w_9 * \hat{\mathbf{y}}'_{AngleTL-J\&E}$
 where weight w is obtained from independent **training data** from

$$\mathbf{y}_{AMR} = w_1 * \hat{\mathbf{y}}_{target} + w_2 * \hat{\mathbf{y}}_{src-JPN} + w_3 * \hat{\mathbf{y}}_{src-EUR} + w_4 * \hat{\mathbf{y}}_{DistTL-JPN} + w_5 * \hat{\mathbf{y}}_{DistTL-EUR} + w_6 * \hat{\mathbf{y}}_{DistTL-J\&E} + w_7 * \hat{\mathbf{y}}_{AngleTL-JPN} + w_8 * \hat{\mathbf{y}}_{AngleTL-EUR} + w_9 * \hat{\mathbf{y}}_{AngleTL-J\&E}$$

Transfer Learning Models

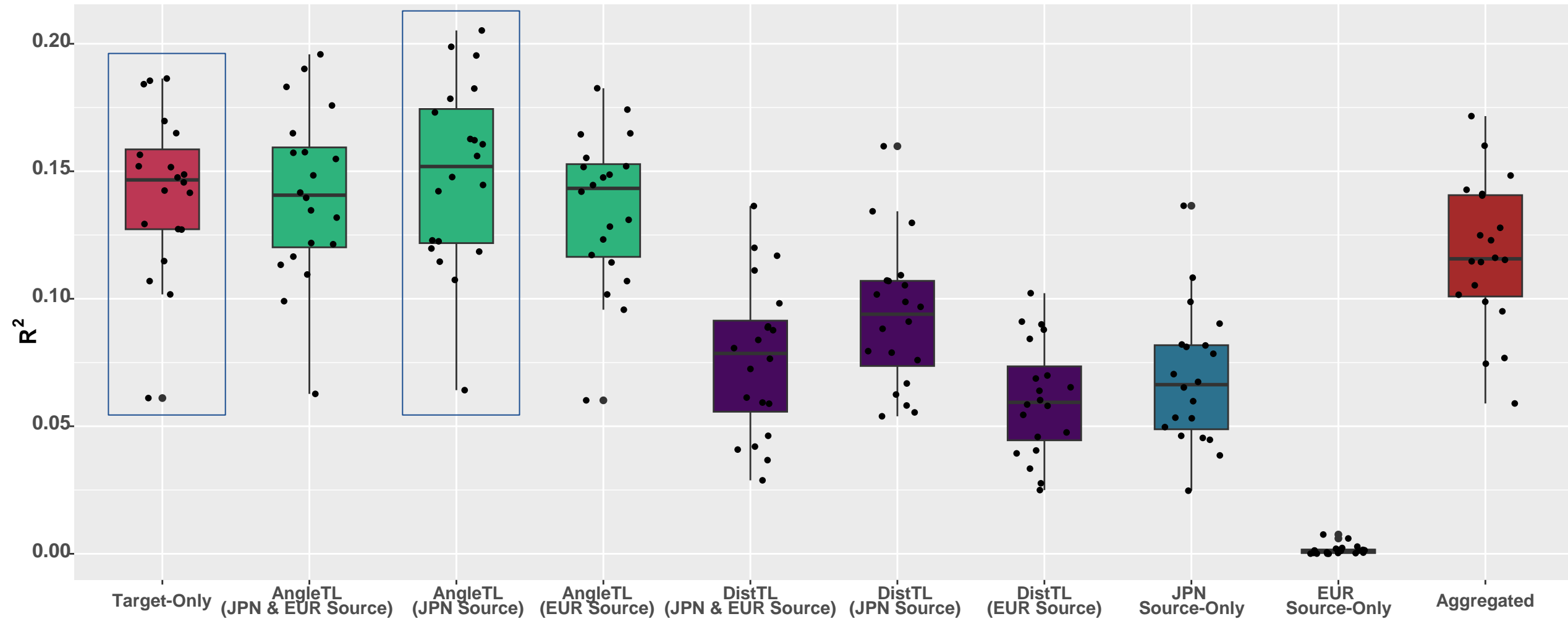
- Comparison

Model Performance Measure:

$$R^2 = (\text{Correlation between predicted } \hat{y}' \text{ and observed } y'_{\text{AMR}})^2$$

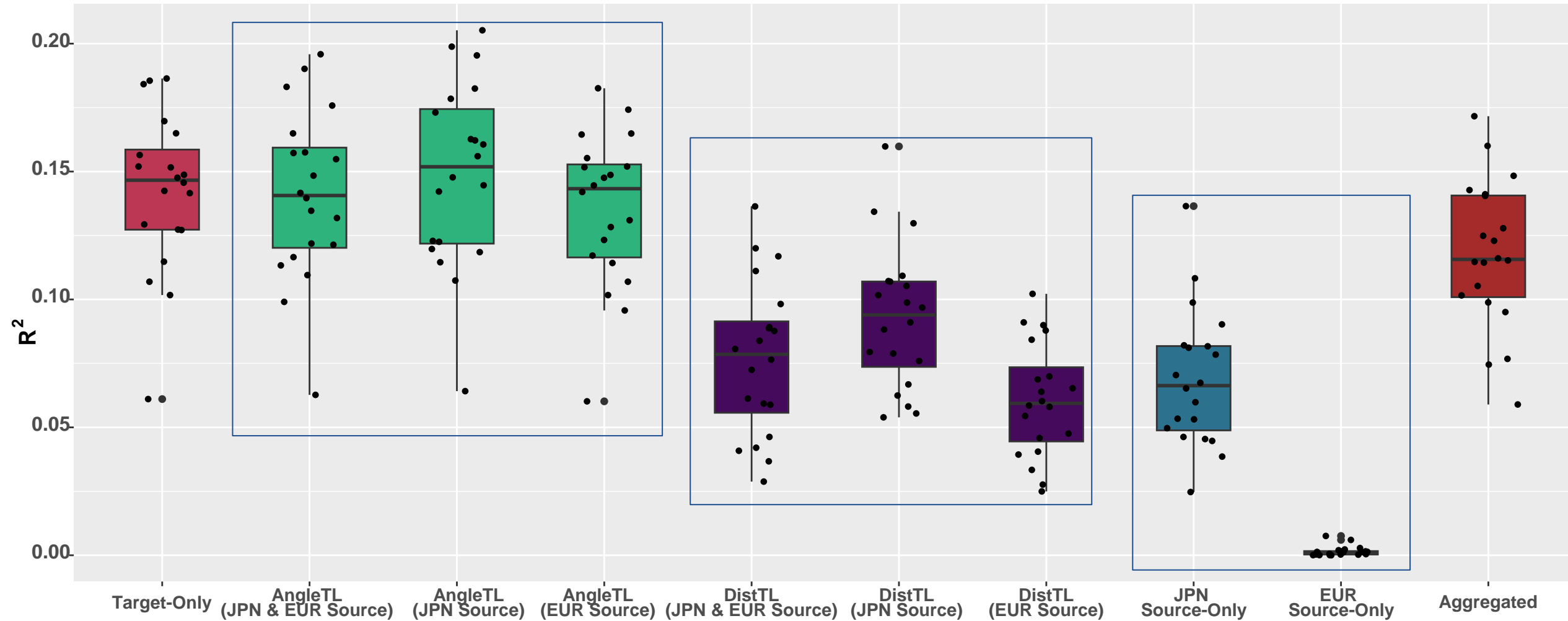
Result

R^2 across models



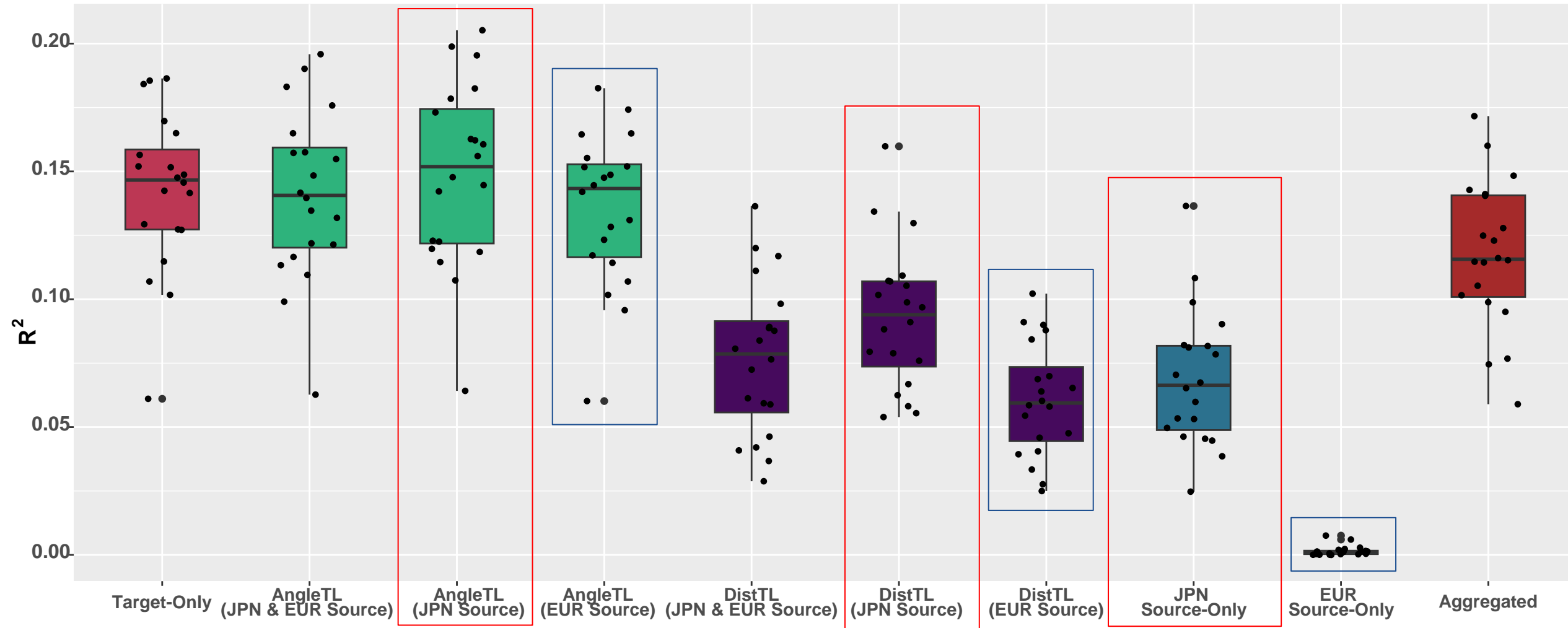
- AngleTL with JPN source showed a >10% improvement in R^2 values over target-only model

R² across models



- AngleTL has superior predictive accuracy compared to source-only model and DistTL
- Under the same source, the performance of transfer learning (AngleTL or DistTL) is consistent with source-only.

R^2 across models



- Models using EUR source all perform poorly, i.e. the similarity between AMR and EUR is low.

Discussion

Discussion

SNP Clumping & Data Quality:

- Reduced genetic marker redundancy
- Isolated most representative SNPs for analysis, which improved reliability of subsequent findings

Transfer Learning Performance:

- Potential to address risk prediction disparities in underrepresented populations, with utilization of data from well-researched populations to benefit underserved groups
- Demonstrated potential of transfer learning in genetic research

Thank you

Appendix

Breast Cancer Genotype Data Quality Control (QC)

- Step1: **Genotype Missing Rate** $< 5\%$ (--geno 0.05), i.e. removing SNPs with NA accounting for more than 5%
- Step2: **Sample Missing Rate** $< 5\%$ (--mind 0.05), i.e. removing samples with NA accounting for more than 5%
- Step3: **Minor Allele Frequency (MAF)** $> 1\%$ (--maf 0.01), i.e. removing SNPs with $MAF < 1\%$
- Step4: **Hardy-Weinberg Equilibrium (HWE)** with **p-value** $< 10^{-6}$ (--hwe 0.000001), i.e. removing SNPs with the p-value $> 10^{-6}$