Breast Cancer microarray datasets, ER+ individuals and DMFS response **Datasets** 8 datasets, with 1555 patients Linear scaling of expression values Cleaning Missing values in gene expressions filled with KNN imputation Remove 60% low variance genes (Dataset section) Remove duplicate samples or samples with missing values Remove studies with too few samples Curation of covariates: size, age, grade (Supplementary table 1 for availability) 7 datasets, with 1021 patients (Table 1) Step 1: Resample studies with replacement **Simulate** Non-parametric Independent Bootstrap Step 2: Resample individuals with replacement **Datasets** (Simulation Parametric Step 3: Use generative model to simulate response Approach section) Bootstrap True Model: Fit CoxBoost Randomly draw samples from U(0,1) and -> true coefficients Compute T and C (function 2-4) True cumulative hazards 1 CV and 1 Repeat for 100 iterations: Within- and A matrix of C Index: CSV value: across study Simulate a Average of Train on dataset i, validation list of (i, j) element For all diagonal and Test on dataset j datasets (Simulation off-diagonal with the Approach section) i, j elements 4-fold CV on set i (i, i) element method above (Supplementary figures 1) **Proportions of** Re-weight Improves CV and CSV, the Alterations of covariates: probabilities of difference is not equalized sources of size, age, grade resampling (Supplementary figure 1, panel b) heterogeneity Increases the variances of Gene Filter genes with CSV, does not eliminate the high Integrative Expression difference Correlation Covariance (Supplementary figure 1, panel c) Use identical Eliminate the difference Difference in model from all (Supplementary figure 1, panel d) true models sets combined The gap is greatly reduced and almost eliminated All sources combined (Supplementary figure 1, panel f) Only the absolute value of the gaps are changed, the Remove outlier CAL trends stay the same. (Supplementary figure 7)