# Systematic Normalization with Multiple Housekeeping Genes for the Discovery of Genetic Dependencies in Cancer

Oliver Bonham-Carter[a] and Yee Mon Thu[b]

Depts of Computer Science[a] and Biology[b], Allegheny College

Meadville, PA

ALLEGHENY COLLEGE

https://www.cs.allegheny.edu
obonhamcarter@allegheny.edu
ythu@allegheny.edu

## PROJECT OBJECTIVES

We analyze gene expression data to discover pairs of genes whose expressions demonstrate patterns of correlation using a computational approach. *This project presents:*

► A focus on genes suppressing genome instability (GIS genes) since function or expression may often be altered in cancer.

► A computational method to determine normalizing factors that make it possible to discover pairs of GIS genes which show consistent correlation.

► Normalizing factors, created by a selection of cancer-specific housekeeping genes providing ability to compare gene expressions data and to treat these values by linear regression.
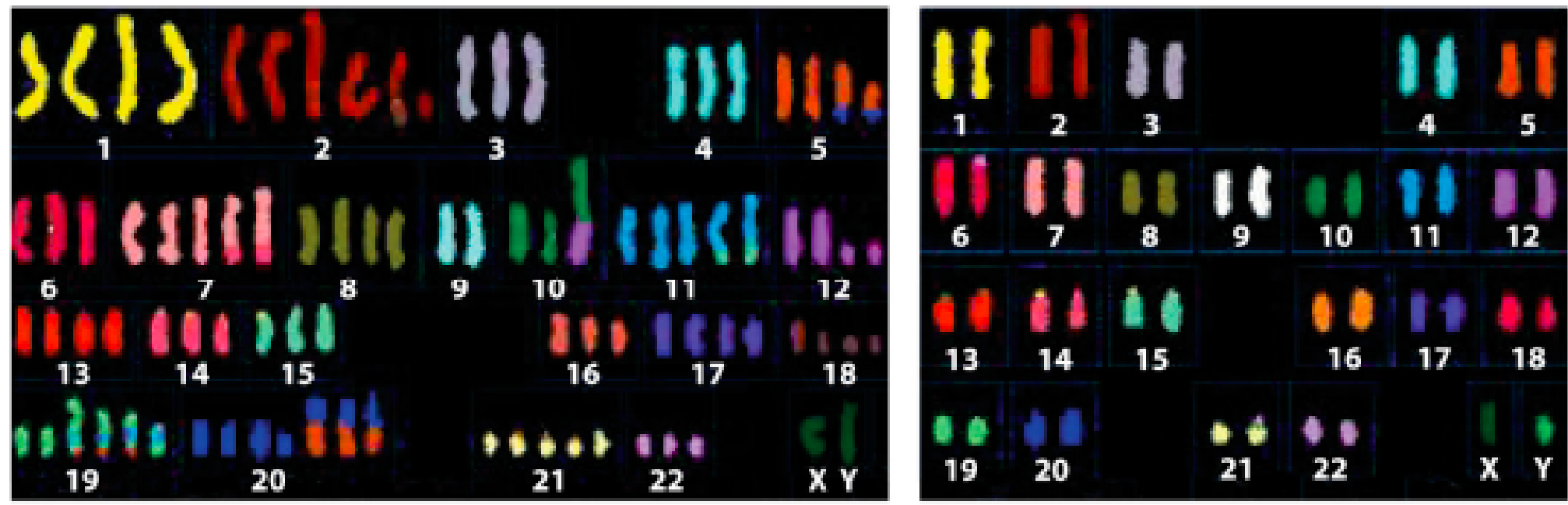


**Figure:** 1. A karyotype of reduced stability of typical cancer cells, left. Increased stability in cells, right.
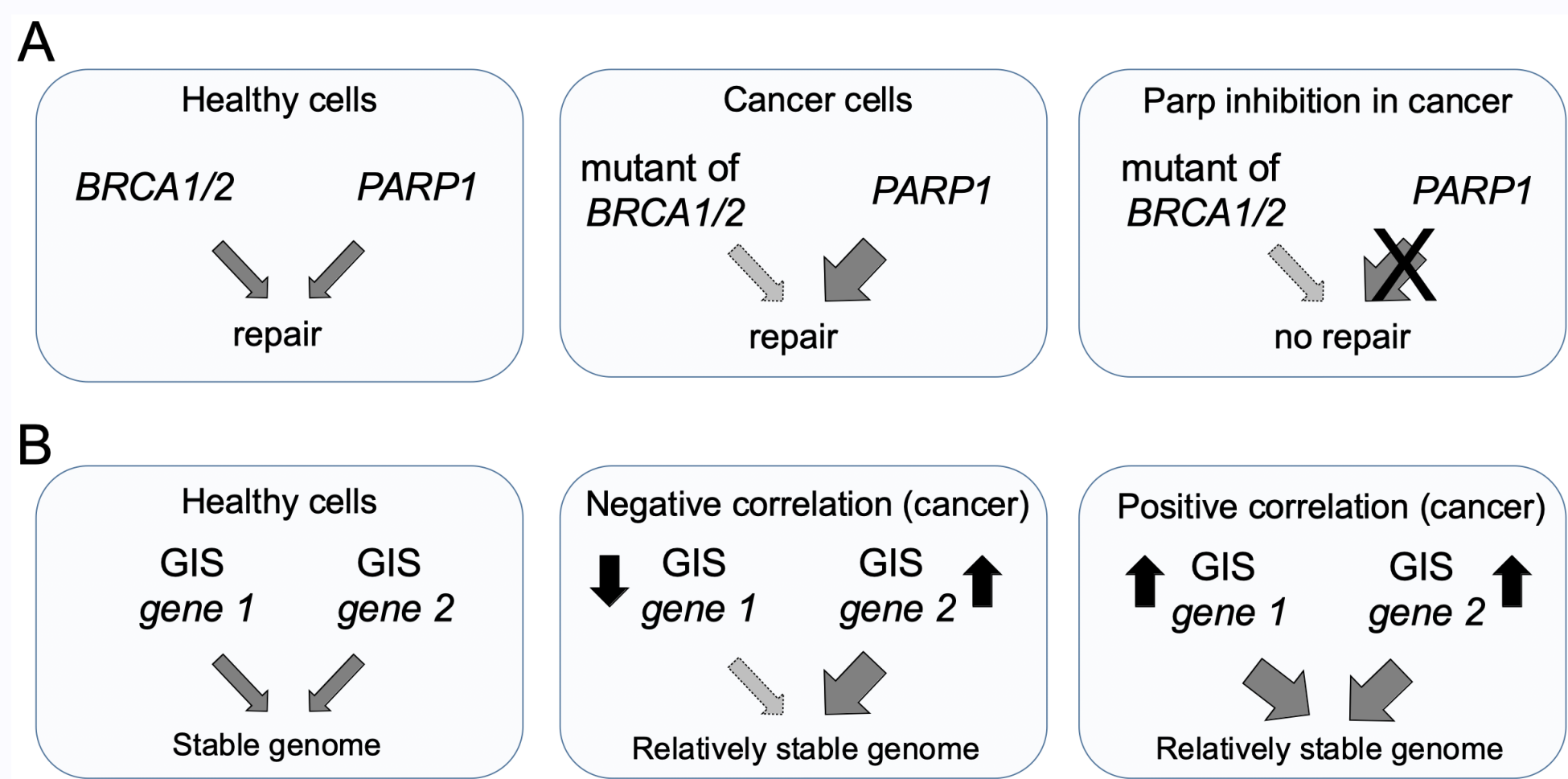
## GENE CORRELATION



**Figure:** 2. Co-Dependance by genes

Genetic interactions between genes are often responsible for DNA repair or genome stability.

► FPKM datasets are used to determine the existence of a positive or negative correlation between the expression of two given GIS (suppress genome instability) genes in cancer, Shown in Figure 2.

► Correlations can reveal if two GIS genes coordinate or if an alteration in the expression of one GIS gene increases dependency of cancer cells on another GIS gene.

## METHOD: PROOF OF CONCEPT

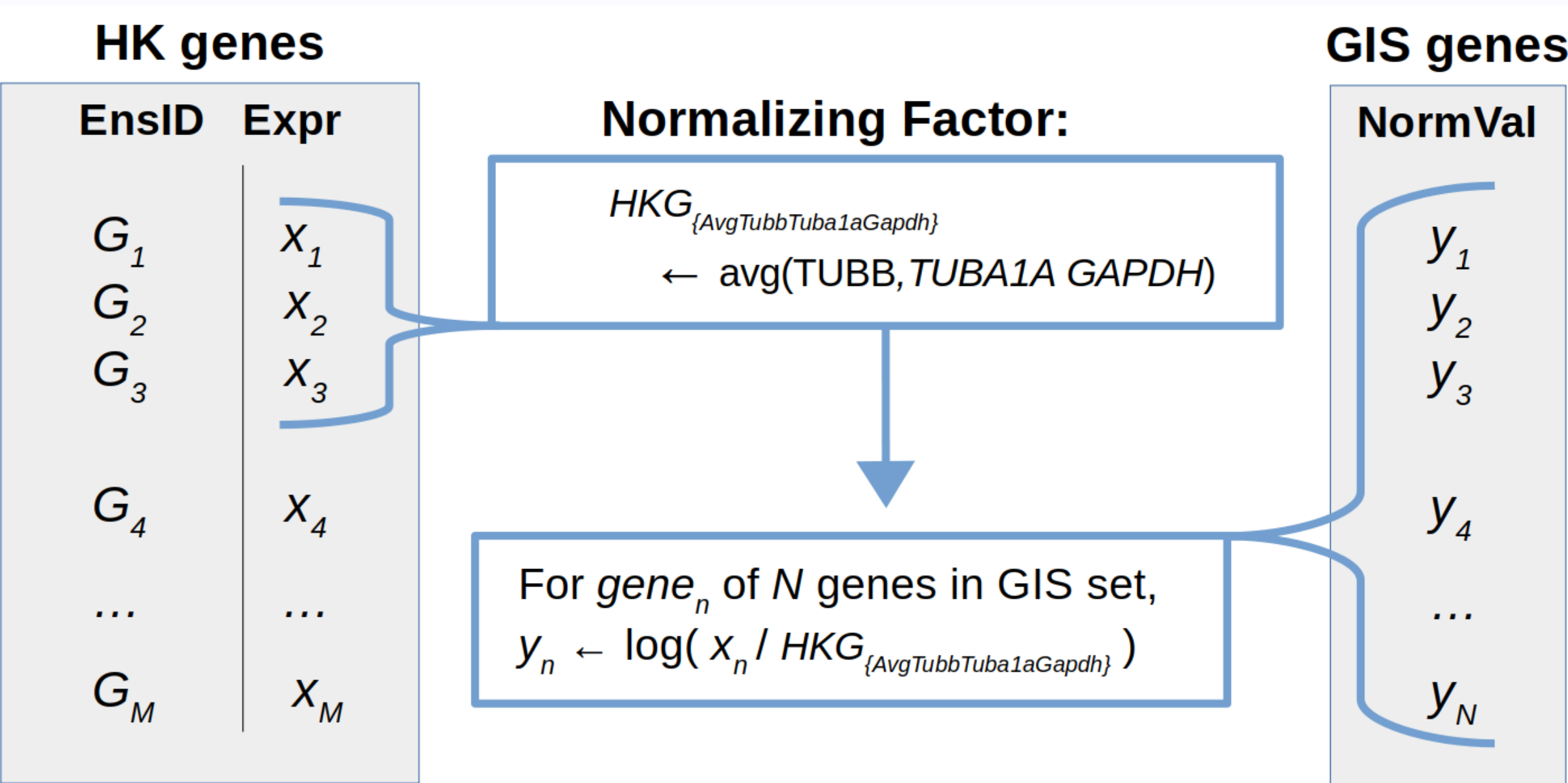Random selection of ten data sets of breast cancer gene expression.
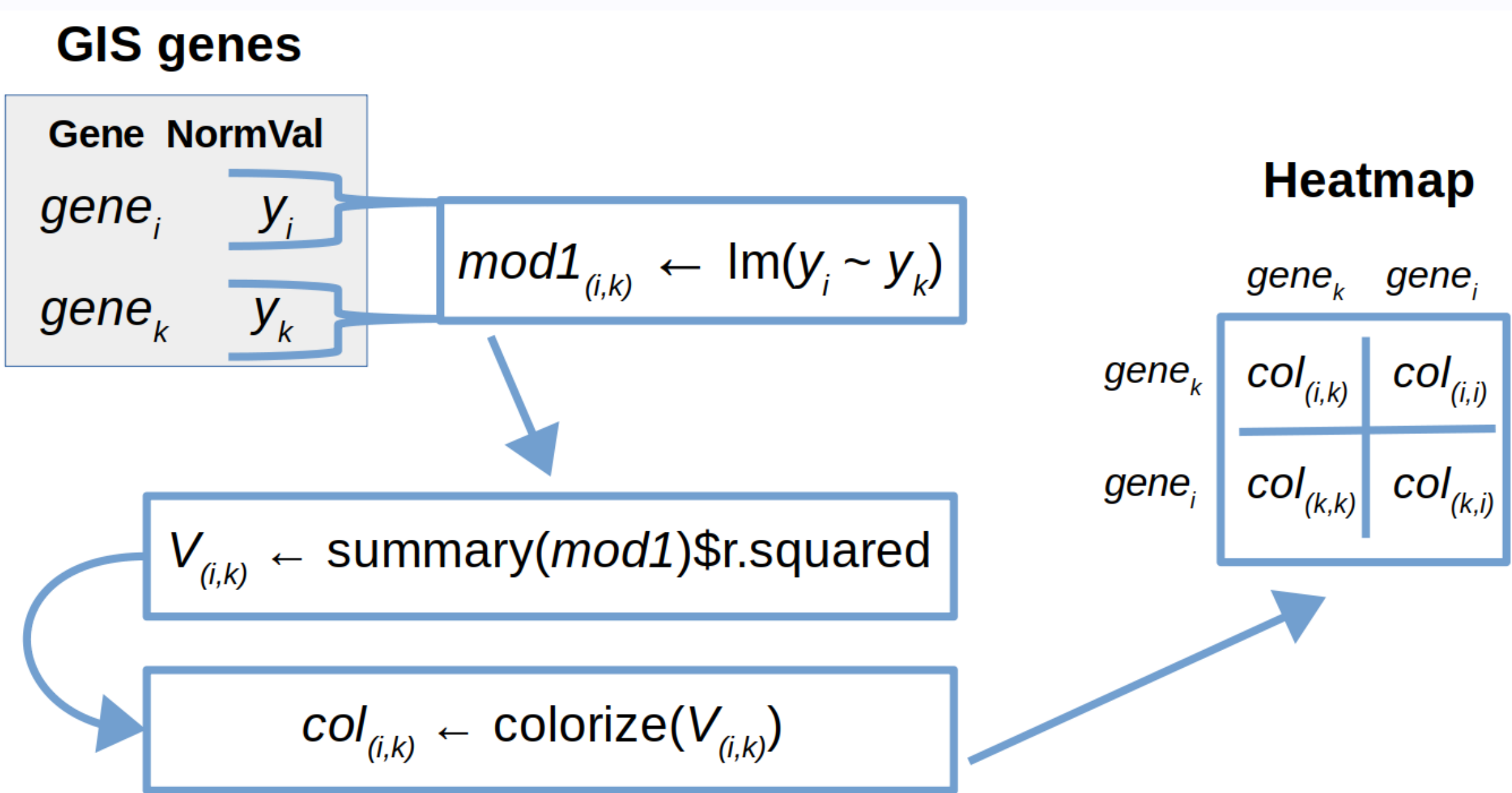


**Figure:** 3. Determining the normalizing factors



**Figure:** 4. HouseKeeping genes to normalize data sets

► Genomic Data Commons Data Portal (National Cancer Institute).
► Selected subset of genes limited search space, were specific to breast cancer research and reduced noise in results.
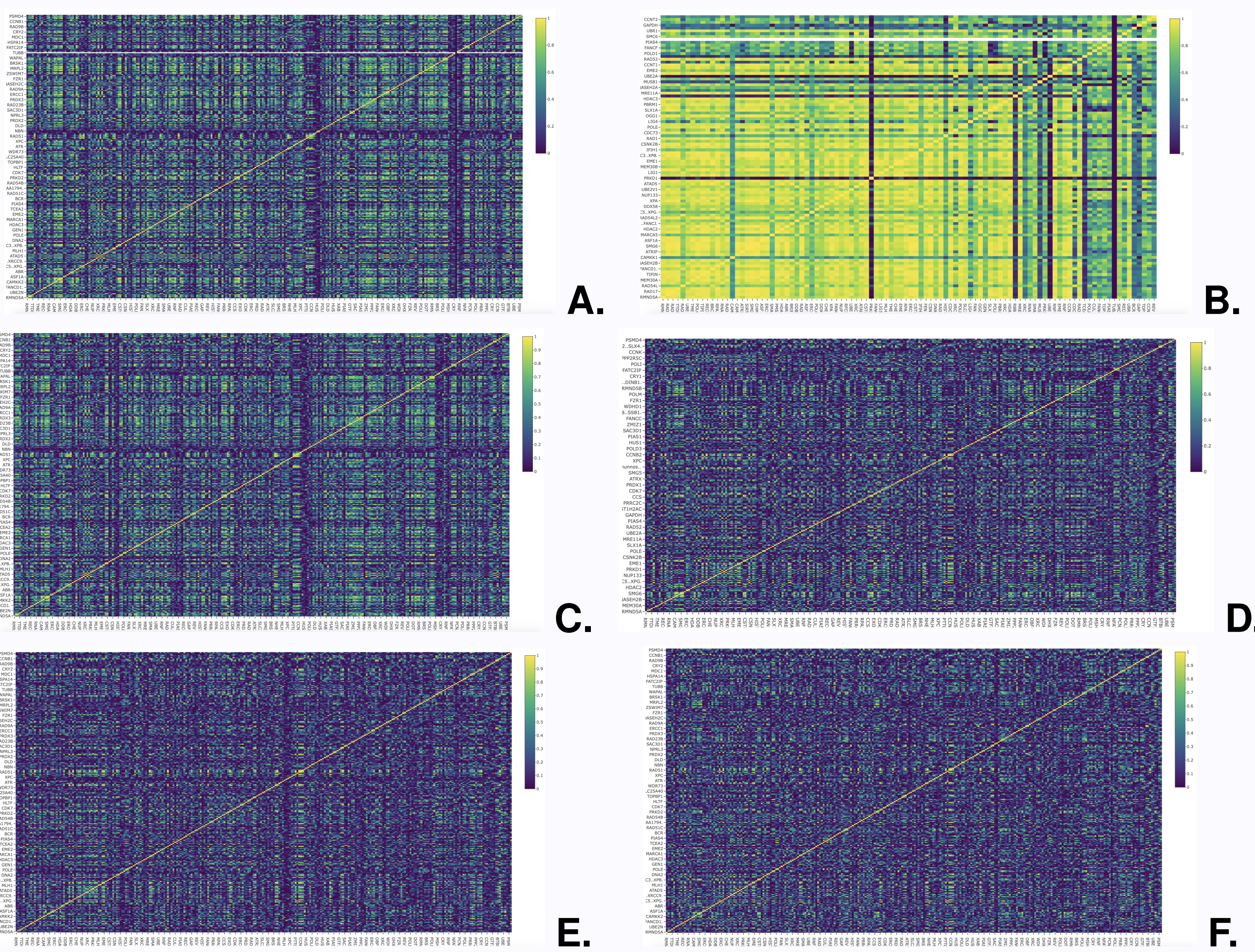
## R-SQUARED VALUES FROM LINEAR MODELS



**Figure:** 5. Heatmaps of $R^2$ values, derived from normalizing factors

► $R^2$ values in heatmaps from linear regression models, *all-against-all* regressions of GIS genes.

► Left to right, $R^2$ values, derived from single housekeeping genes to create normalizing factors; $HKG_{Tubb}$, $HKG_{Tuba1a}$, $HKG_{Tubb}$, $HKG_{AvgTubbTuba1aGapdh}$, see Figure 5,{**A** - **C**}, resp.

► Ten housekeeping genes: $HKG_{AvgG1}$, $HKG_{AvgG3}$, $HKG_{AvgG3}$, See Figure 5,{**D** - **F**}, resp.

## RESULTS

► $HKG_{Tuba1a}$ (**B**): high diversity of $R^2$ values indicated poor normalizing, biologically improbable.

► According to regression model results, normalizing factors that were created from larger groups of housekeeping gene produced correlations implying more biological relevance.

## CONCLUSIONS

► Our results (examples shown in heatmaps of Figure 5) indicate normalized data and contain biologically probable findings.

► **Single Expression Normalization**: Normalizing factors derived from single housekeeping genes did not provide generally consistent correlations that were biologically probable.

► **Multiple Expression Normalization**: Using the averaged expression values of multiple housekeeping genes was an effective approach to finding biologically relevant consistency across our data sets.

► Our method enabled us to identify the co-expression of gene pairs in breast cancer tissues and compared diverse normalization factors.

► Our study also allows reproducibility across data sets and allows for scalability in gene correlation research.