# Meta-analysis of gene expression and survival data using the GXA framework: a new prognosis tool for breast cancer

Markus Schröder[1,2], Daniel Gusenleitner[1], Matthew Schwede[1], Alexander Goesmann[2], Aedín C. Culhane[1], John Quackenbush[1] and Benjamin Haibe-Kains[1]

[1]Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, USA
[2]Computational Genomics, Centrum for Biotechnology (CeBiTec), University of Bielefeld, Bielefeld, Germany

## Introduction to *Survcomp*

Gene expression profiling has generated unprecedented insight into our molecular understanding of cancer. In breast cancer, gene expression profiling studies have been widely employed and have not only advanced our understanding of disease, but have provided multi-gene predictive and prognostic tests including Oncotype DX, mammaprint, Veridex GGI and the Breast BioClassifier for breast cancer molecular subtypes.

To identify new prognostic genes and gene signatures, several risk prediction models have been introduced recently.

Our SurvComp package is providing functions to assess and to statistically compare the performance of these risk prediction (survival) models. It includes:
i. Implementation of state-of-the-art statistics developed to measure the performance of risk prediction models
ii. Combining these statistics estimated from multiple datasets using a meta-analytical framework
iii. Functions to visualize those measurements in a clear and compact way

*SurvComp* is available on Bioconductor.org

## A new Prognosis Tool for Breast Cancer

Several excellent online resources for mining of gene expression data exist, including Oncomine, NextBio and the Gene Expression Atlas (GXA). However, none provide survival data in addition to significant gene rankings. We are building an online resource combining gene expression and survival data from multiple datasets to enable clinicians and biologists to assess the prognostic values of genes of interest.

We extend the existing framework of the GXA in our pipeline. The GXA, maintained by the European Bioinformatics Institute, is an added-value database of gene expression sequencing data for different cell types, organism parts, developmental stages, disease states, sample treatments and other biological/experimental conditions.

The structure of the GXA has three layers: R Analytics, Database and Front End. It includes a pipeline to port data from ArrayExpress/GEO to GXA. Each dataset is annotated with a standard experimental factor ontology (EFO). We describe our new gene mining approach below.

## Workflow



Gene Expression Datasets

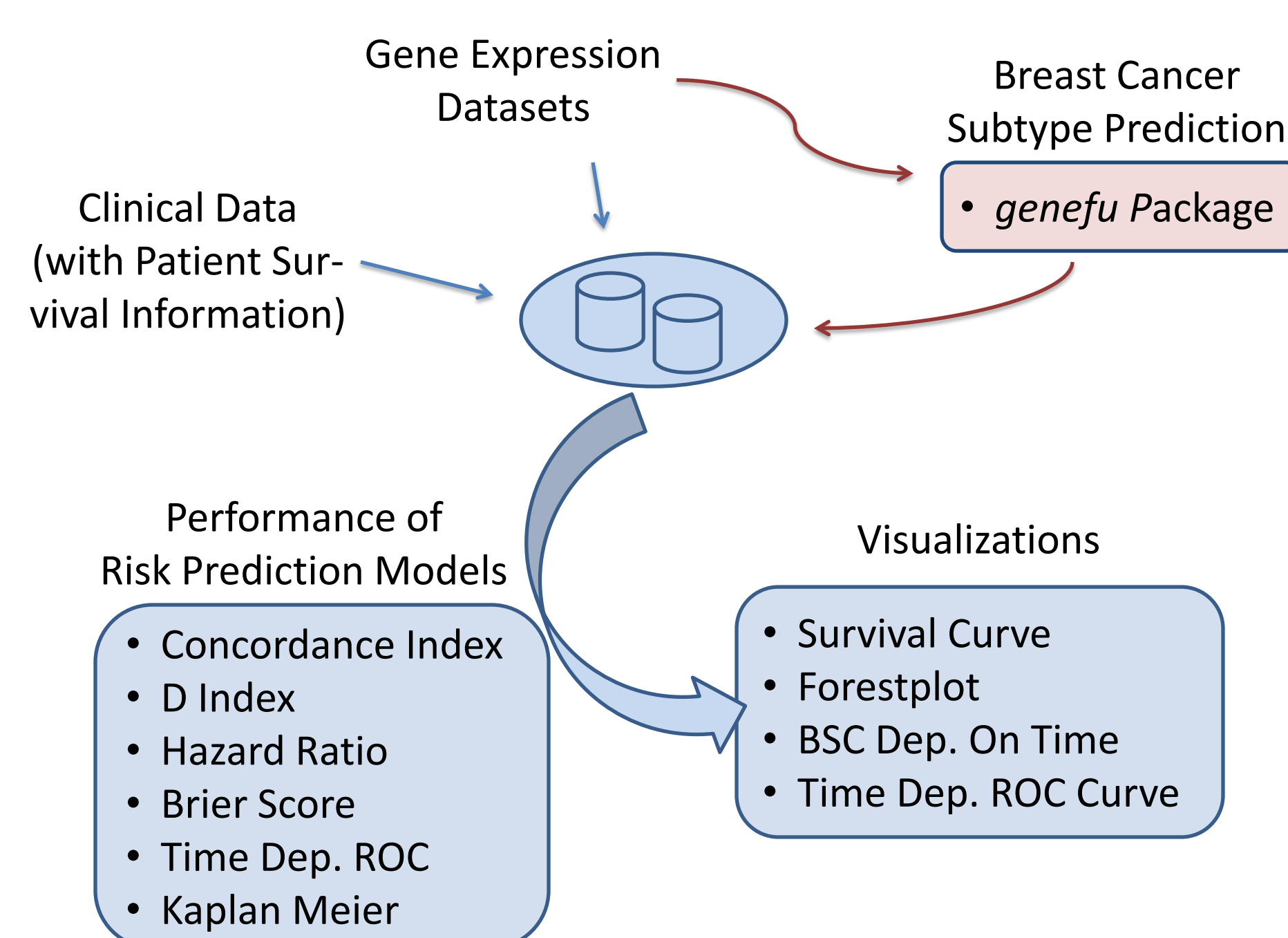Clinical Data (with Patient Survival Information)
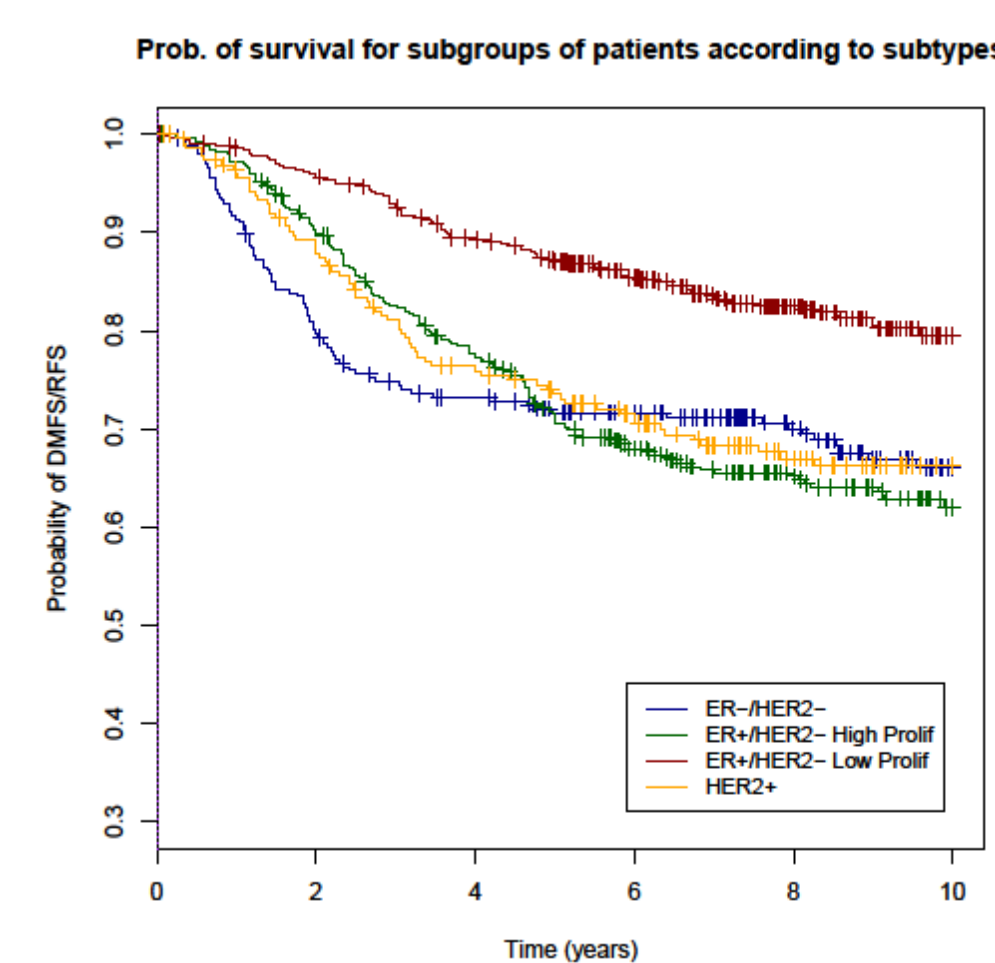
GXA

Website

(Network Common Data Form)

netCDF

*SurvComp* (Survival Analysis)
*genefu* (Breast Cancer Molecular Subtype Prediction)
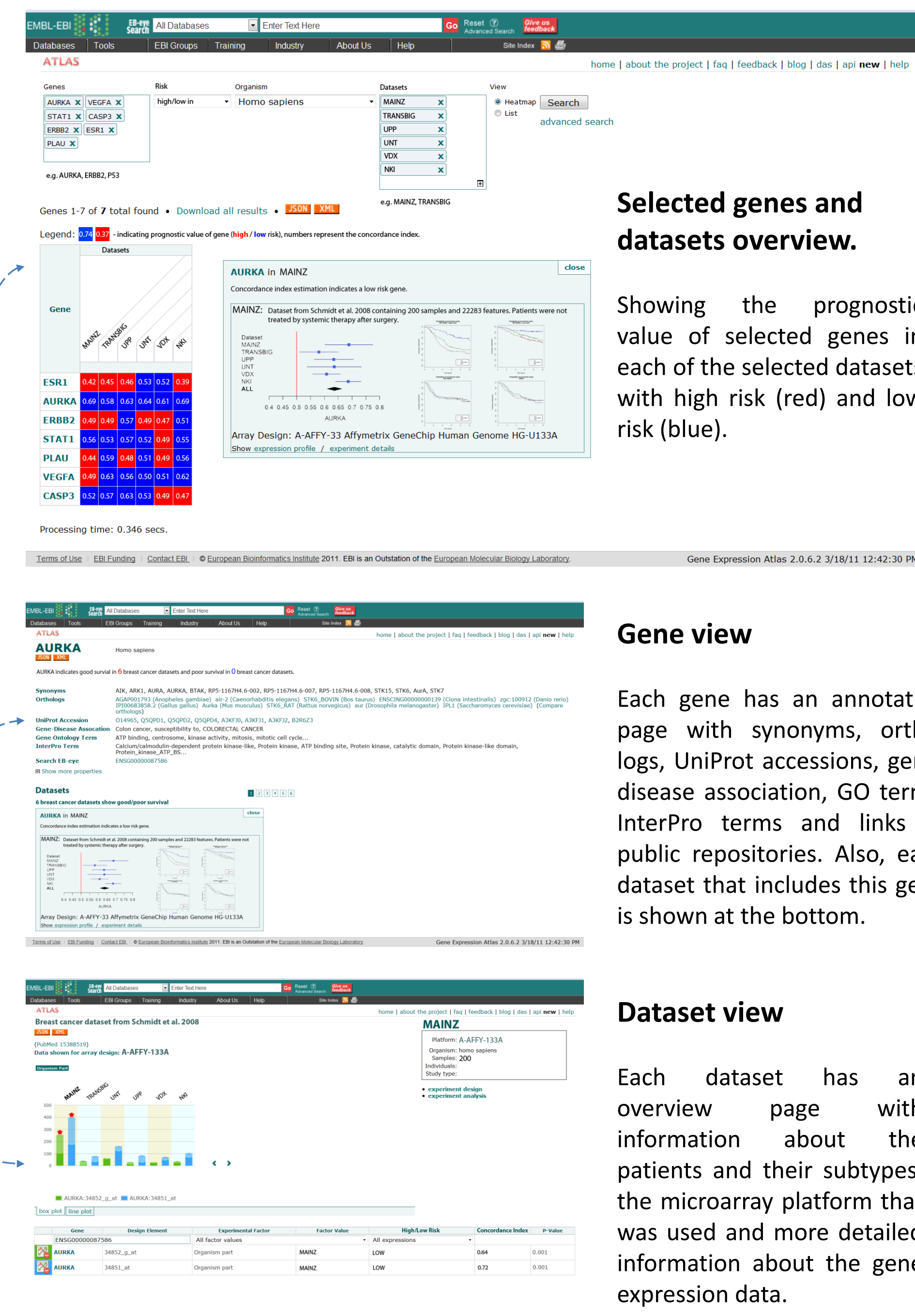
## Survival Analysis Workflow



Gene Expression Datasets

Clinical Data (with Patient Survival Information)

Breast Cancer Subtype Prediction
- *genefu* Package

Performance of Risk Prediction Models
- Concordance Index
- D Index
- Hazard Ratio
- Brier Score
- Time Dep. ROC
- Kaplan Meier

Visualizations
- Survival Curve
- Forestplot
- BSC Dep. On Time
- Time Dep. ROC Curve

## Kaplan Meier Survival Curves



Prob. of survival for subgroups of patients according to subtypes

**The Kaplan Meier survival curve**: estimation of the survival expectancy for a group of patients with respect to time.

Kaplan Meier survival curve for a combination of subtype information and survival data for patients from six breast cancer datasets with a total of 1467 patients. The molecular subtypes are ER-/HER2-, ER+/HER2- High Proliferation, ER+/HER2- Low Proliferation and HER2+.

## Website Hierarchy



**Selected genes and datasets overview.**

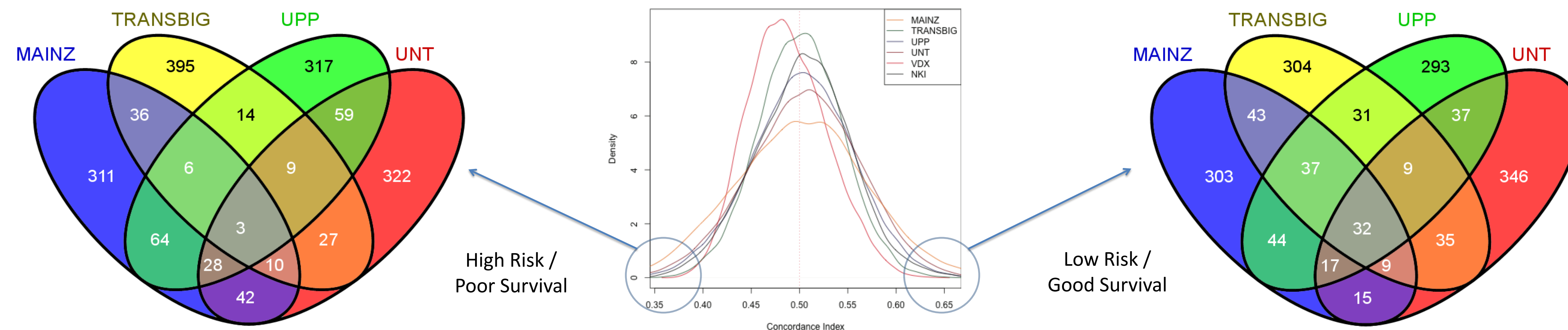Showing the prognostic value of selected genes in each of the selected datasets with high risk (red) and low risk (blue).

**Gene view**

Each gene has an annotation page with synonyms, orthologs, UniProt accessions, gene-disease association, GO terms, InterPro terms and links to public repositories. Also, each dataset that includes this gene is shown at the bottom.

**Dataset view**

Each dataset has an overview page with information about the patients and their subtypes, the microarray platform that was used and more detailed information about the gene expression data.

## Genome Scale Overview of Gene (n >20.000) Concordance Indices from Multiple Datasets



High Risk / Poor Survival

Low Risk / Good Survival

Representations of the 500 lowest (**left**) and 500 highest (**right**) concordance indices from 4 different datasets. High concordance indices for genes indicate a good survival / low risk for patients, low concordance indices for genes indicate poor survival / high risk for patients. In the **middle** figure, the concordance indices for genes from six different datasets are shown as a genome scale density plot, providing an overview of the concordance index distribution in each dataset.

**The Concordance Index**: probability that, for a pair of randomly chosen comparable samples, the patient with the higher risk prediction will experience an event before the other patient.

## Analysis of Concordance Indices from Six Datasets

Low Concordance Indices / High Risk

| Gene Symbol | Gene Description |
|---|---|
| TXNIP | thioredoxin interacting protein |
| PKP2 | plakophilin 2 |
| PDLIM5 | PDZ and LIM domain 5 |
| LMO4 | LIM domain only 4 |
| USP34 | ubiquitin specific peptidase 34 |
| HPS1 | Hermansky-Pudlak syndrome 1 |
| MUC5AC | mucin 5AC, oligomeric mucus/gel-forming |
| SIVA1 | SIVA1, apoptosis-inducing factor |
| CCND2 | cyclin D2 |
| SERPINA5 | serpin peptidase inhibitor, clade A member 5 |

We took the union of the genes from six datasets, which resulted in 19768 unique genes. Those genes were ranked according to their concordance index and summed over the six datasets with leaving the highest rank out, e.g. a gene with the lowest concordance index in all datasets would have score five since it has rank one in all datasets. We removed the highest rank for each gene over the six datasets to be more sensitive to outliers.

The **left table** shows the genes with the lowest scores strongly related to high risk patients.

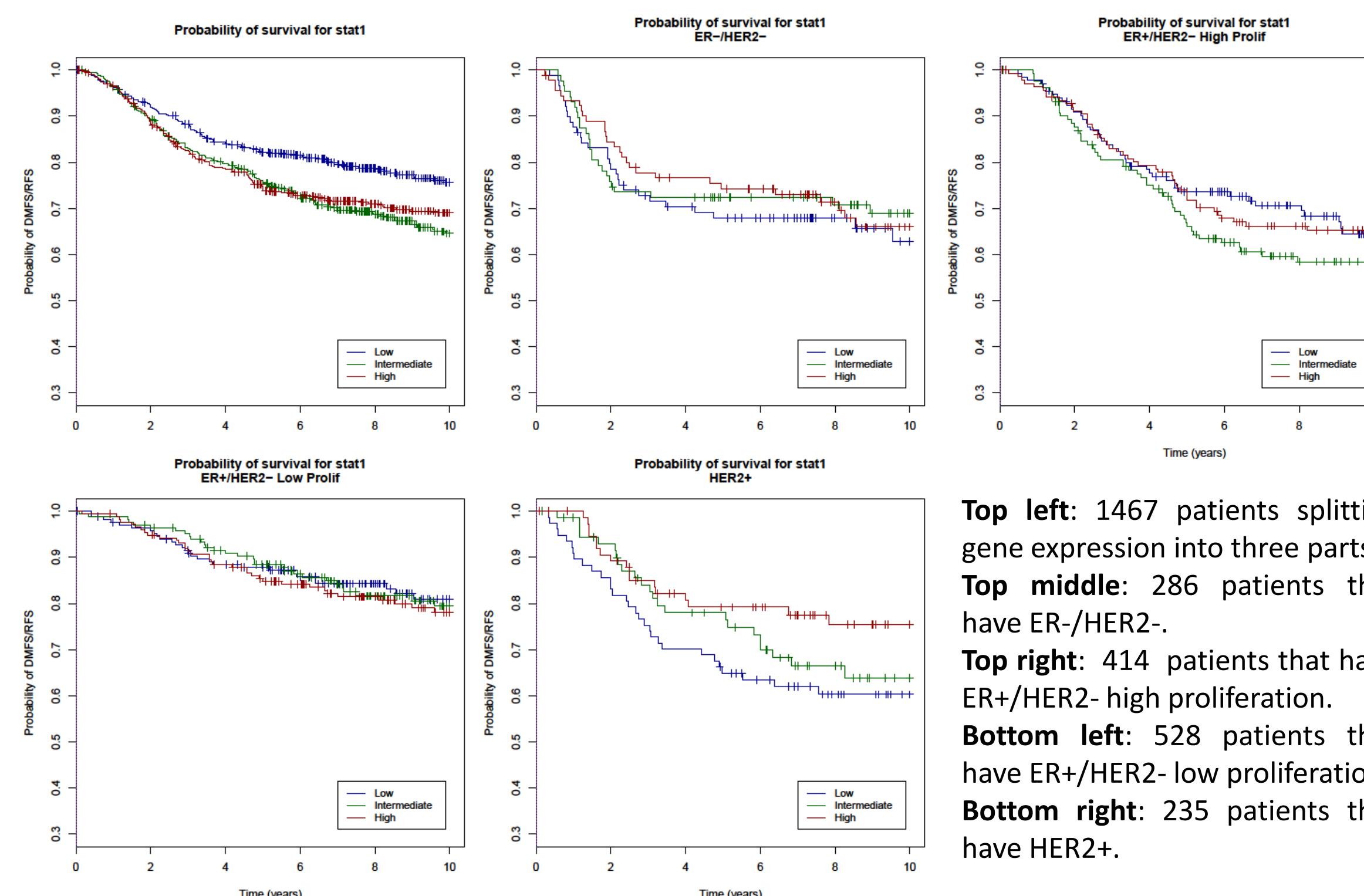The **right table** shows the genes with the lowest scores strongly related to low risk patients.

High Concordance Indices / Low Risk

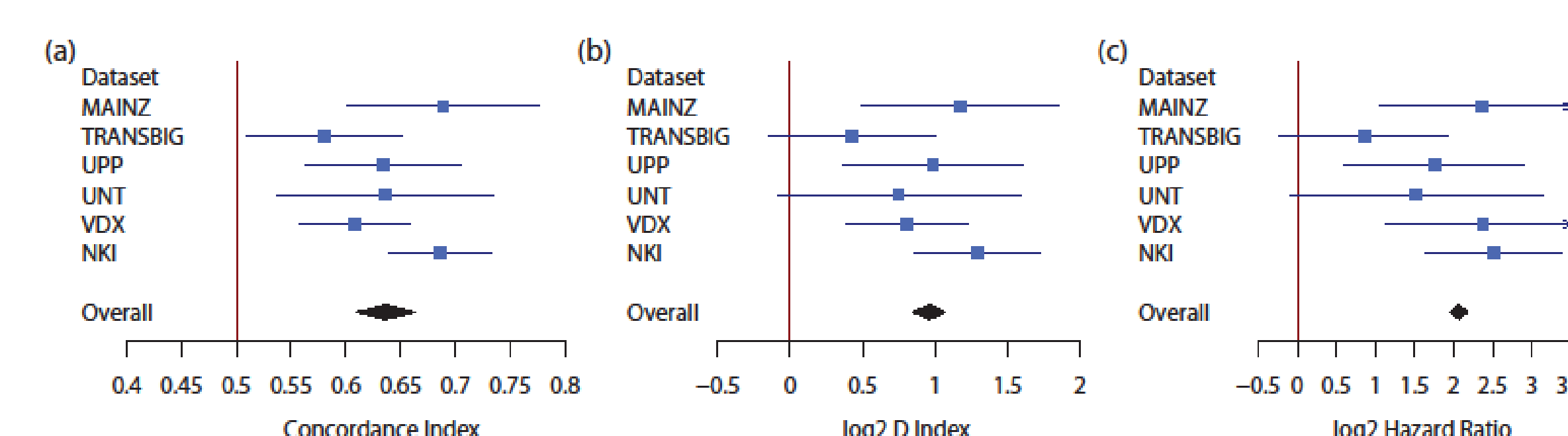| Gene Symbol | Gene Description |
|---|---|
| TRO | trophinin |
| B2M | beta-2-microglobulin |
| DDX47 | DEAD (Asp-Glu-Ala-Asp) box polypeptide 47 |
| TMEM45A | transmembrane protein 45A |
| CCNT2 | cyclin T2 |
| NFIB | nuclear factor I/B |
| TMEM132A | transmembrane protein 132A |
| RAB3IL1 | RAB3A interacting protein (rabin3)-like 1 |
| AGBL2 | ATP/GTP binding protein-like 2 |
| PION | pigeon homolog (Drosophila) |

## Gene Expression Datasets Included in Prototype

| Dataset | Patients [#] | ER+ [#] | HER2+ [#] | Age [years] | Grade [1/2/3] | Platform |
|---|---|---|---|---|---|---|
| MAINZ | 200 | 155 | 23 | 25-90 | 29/136/35 | HGU133A |
| TRANSBIG | 198 | 123 | 35 | 24-60 | 30/83/83 | HGU133A |
| UPP | 251 | 175 | 46 | 28-93 | 67/128/54 | HGU133AB |
| UNT | 137 | 94 | 21 | 24-73 | 32/51/29 | HGU133AB |
| VDX | 344 | 186 | 57 | 26-83 | 7/42/148 | HGU133A |
| NKI | 337 | 212 | 53 | 26-62 | 79/109/149 | Rosetta |
| **Overall** | **1467** | **945** | **235** | **24-93** | **244/549/498** | **Affy/Agilent** |

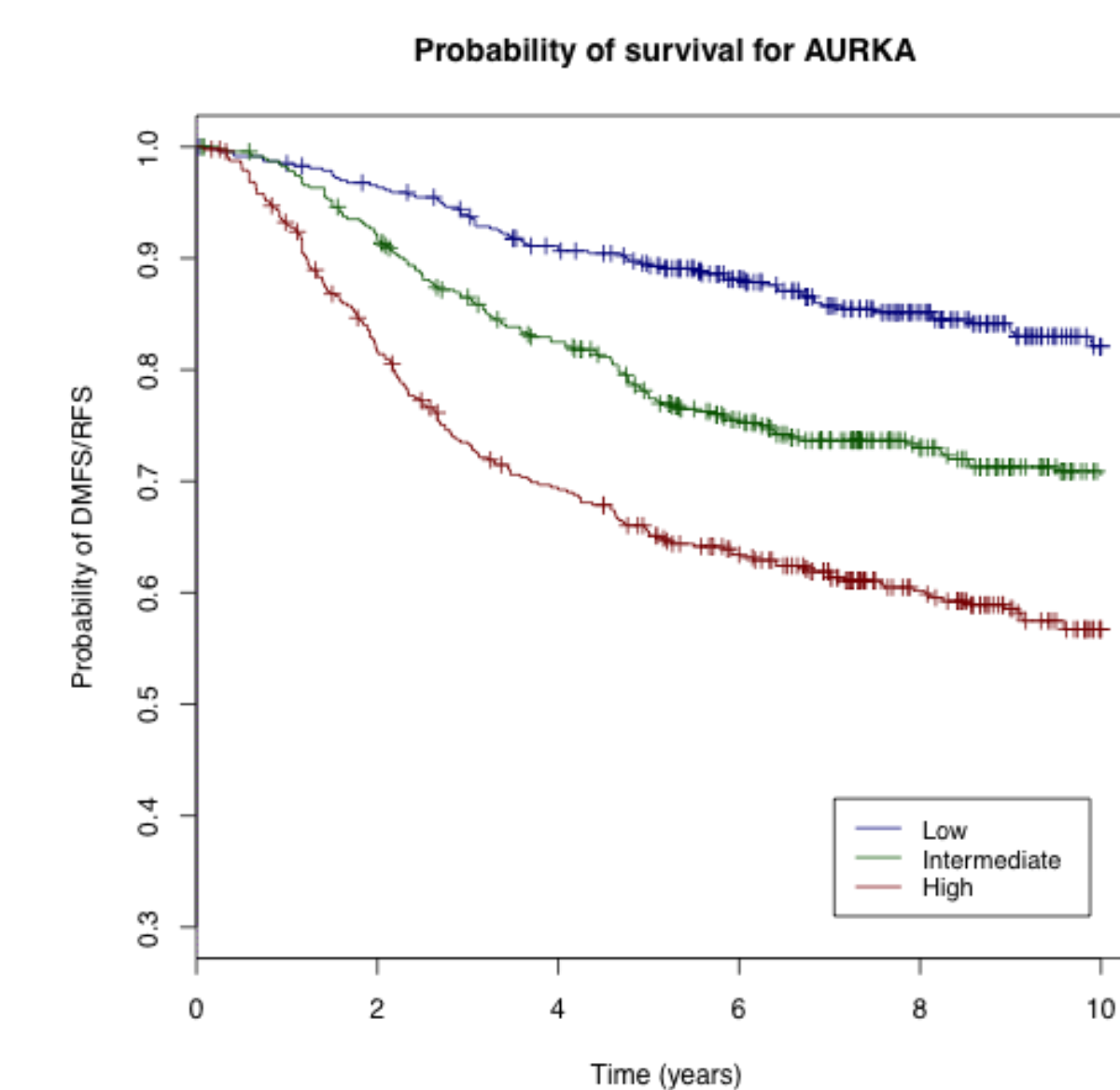## Prognostic Value of STAT1 Using Subtypes



**Top left**: 1467 patients splitting gene expression into three parts.
**Top middle**: 286 patients that have ER-/HER2-.
**Top right**: 414 patients that have ER+/HER2- high proliferation.
**Bottom left**: 528 patients that have ER+/HER2- low proliferation.
**Bottom right**: 235 patients that have HER2+.

## A Meta-Analysis Case Study of AURKA



**Top:** Forestplots show three different measurements for the performance of risk prediction models for the gene AURKA in six different datasets. An overall estimation for each measurement is included. (a) concordance index, (b) log2 D index, (c) log2 Hazard ratio.

**Right:** Kaplan Meier survival curve for the gene AURKA. Gene expression and survival data is a combination of six breast cancer datasets with over 1400 patients. The gene expression is split into 3 parts using quantiles (33% and 66%)

Probability of survival for AURKA

## References

Schröder et al. (2011), SurvComp: an R/Bioconductor package for performance assessment and comparison for survival analysis, in preparation

Haibe-Kains et al. (2008), A comparative study of survival models for breast cancer prognostication based on microarray data: does a single gene beat them all?, *Bioinformatics*, **24**, 2200-2208

Kapushesky et al. (2010), Gene Expression Atlas at the European Bioinformatics Institute, Bioinformatics, **38**, D690-D698