

VUmc  Cancer Center Amsterdam

Statistical Modeling of HPV-induced Transformation

Viktorian Miok

Overview

- **Introduction**
- **Method**
- **Comparison**
- **Application**

Introduction

Cervical Cancer

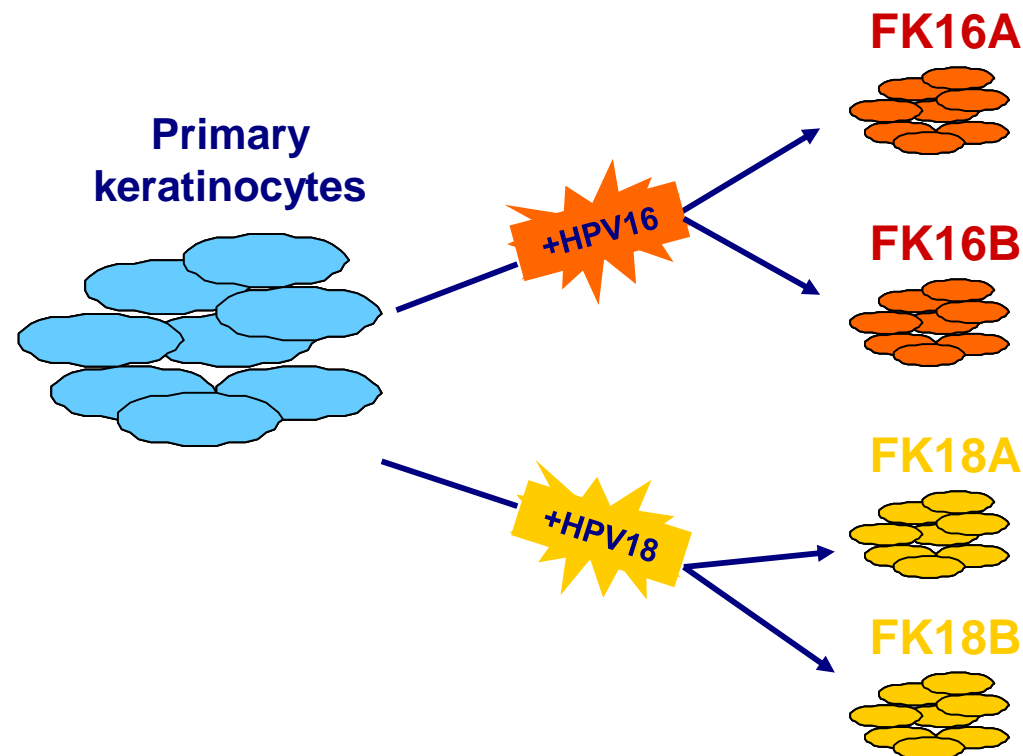
- Second most common cancer in women worldwide.
- The best known example how virus causes cancer.
- Caused by HPV virus, in 80% cases HPV16 and HPV18.
- Genetic disease.

Our approach to investigate HPV-induced transformation

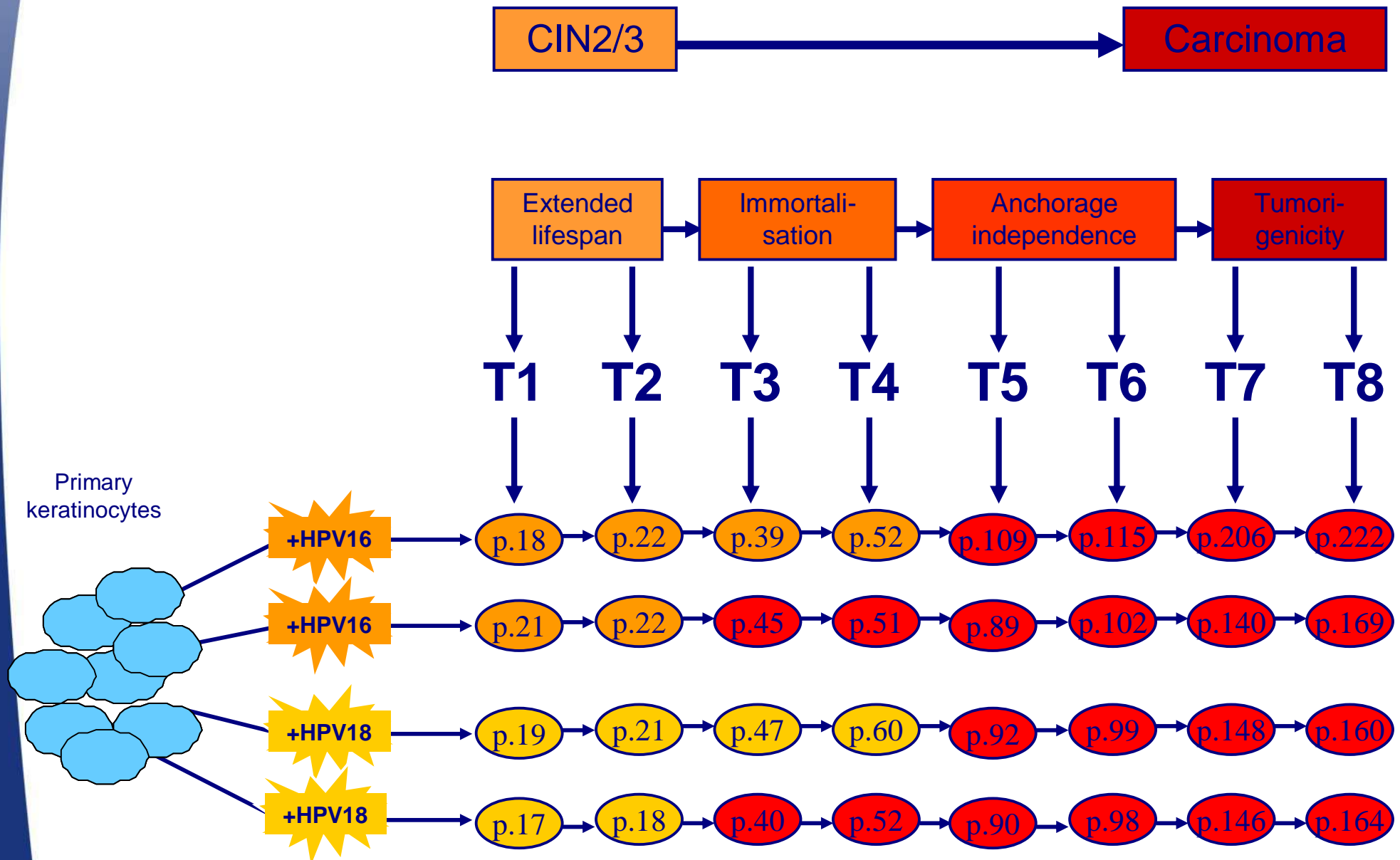
- Cell line model – in vitro model system of HPV-induced transformation
- Time structure – profiling at different moments in time for multiple molecular levels
- Integration – high-throughput multi level molecular data sets
- Aim: Insights in necessary (epi) genetic changes

In vitro model system of HPV-induced transformation

- Suitable model to investigate HPV-induced cancer development.
- Originate from same parental cells.



Time Effect

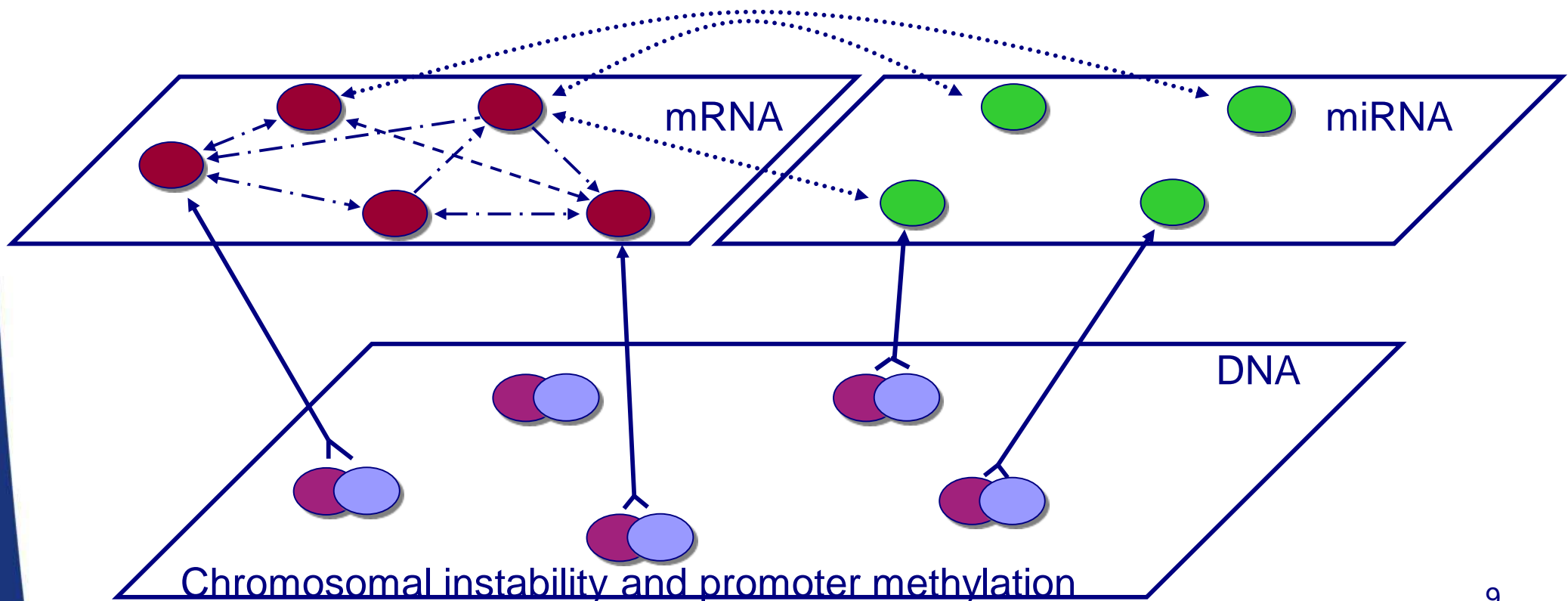


DNA Methylation

- Mechanism of epigenetic deregulation – gene silencing.
- Known to be involved in HPV-induced transformation.
- Measured in mRNA and miRNA gene expression experiments
- Comparing mRNA and miRNA gene expression in the samples treated with DAC and in samples without treatment.
- Results are partially validated using Illumina human methylation arrays.

Multilevel Integrative Analysis

- Goal of project is to perform multilevel integrative analysis of longitudinal data sets
- Necessity of novel integrative statistical methods



Method

Integration

$$GE = \text{Cell Line} + \text{Time}$$



$$GE = \text{Cell Line} + \text{CN} + \text{Time}$$



$$GE = \text{Cell Line} + \text{CN} + \text{Time} + \text{Methylation}$$



Temporal Differential Gene Expression

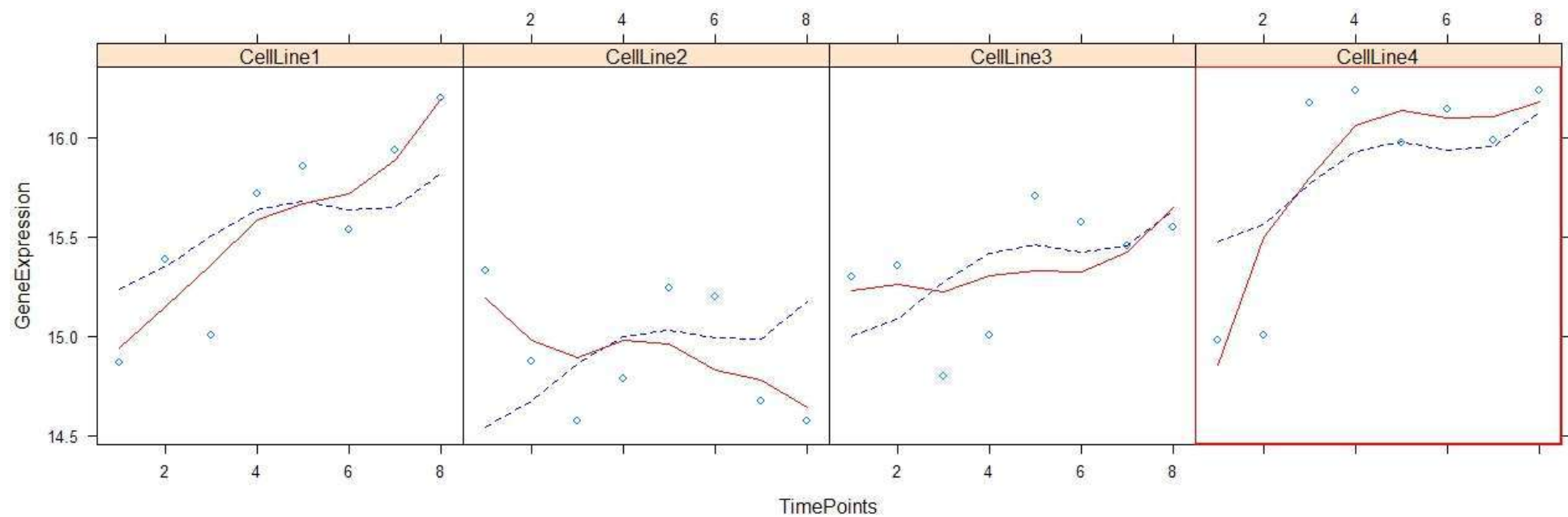
- Generalized linear mixed model

$$GE = \underbrace{\text{Cell Line}}_{\text{fixed effect}} + \underbrace{CN_i}_{\text{random effect}} + \text{Time}$$

- Penalized splines
- Parameters and hyper-parameters estimated using empirical Bayes procedure employing INLA
- Shrinkage of dispersion-related parameters
- The likelihood ratio test

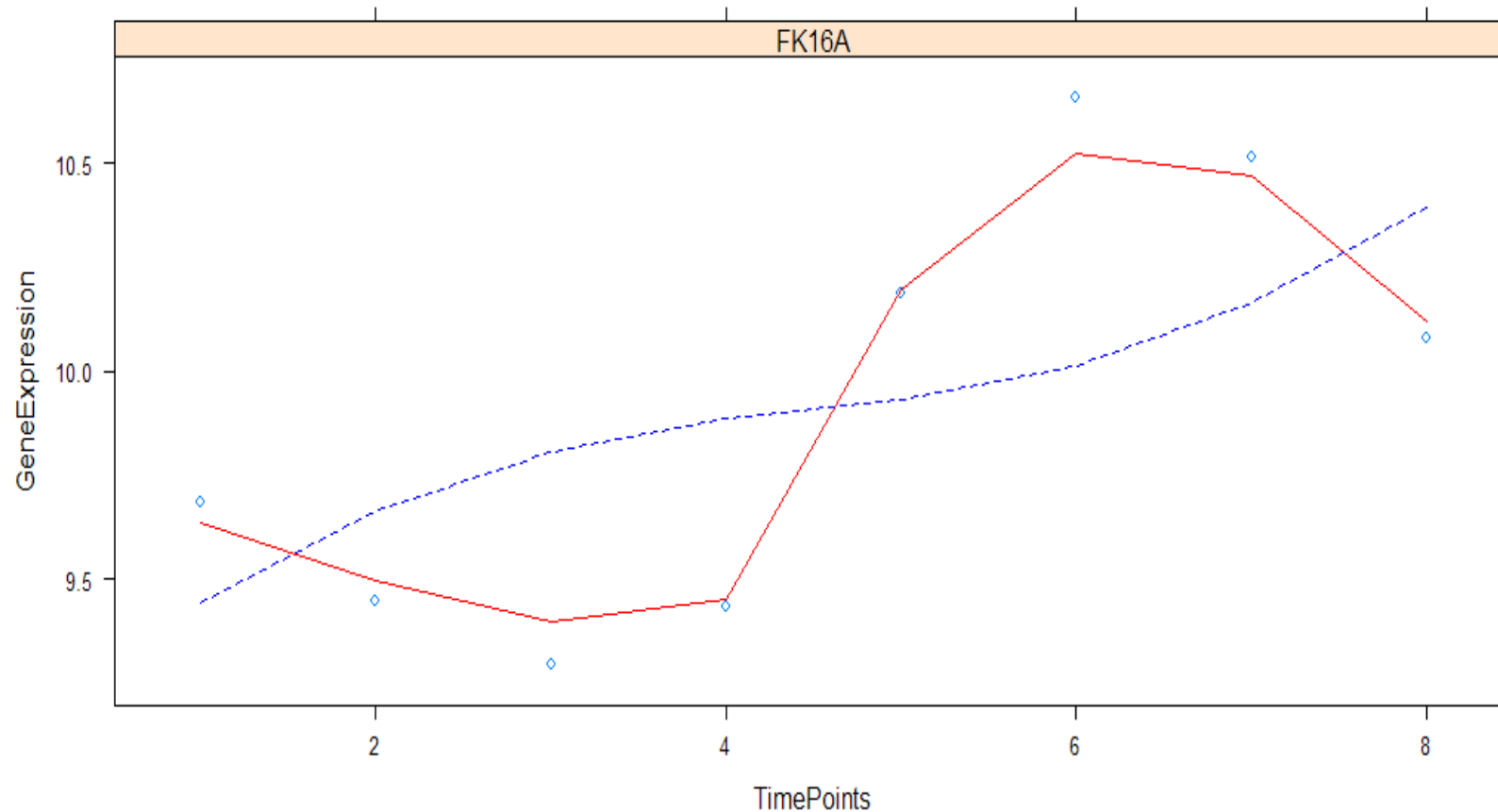
Time effect

-  Different spline per cell line
-  Same spline for all cell lines



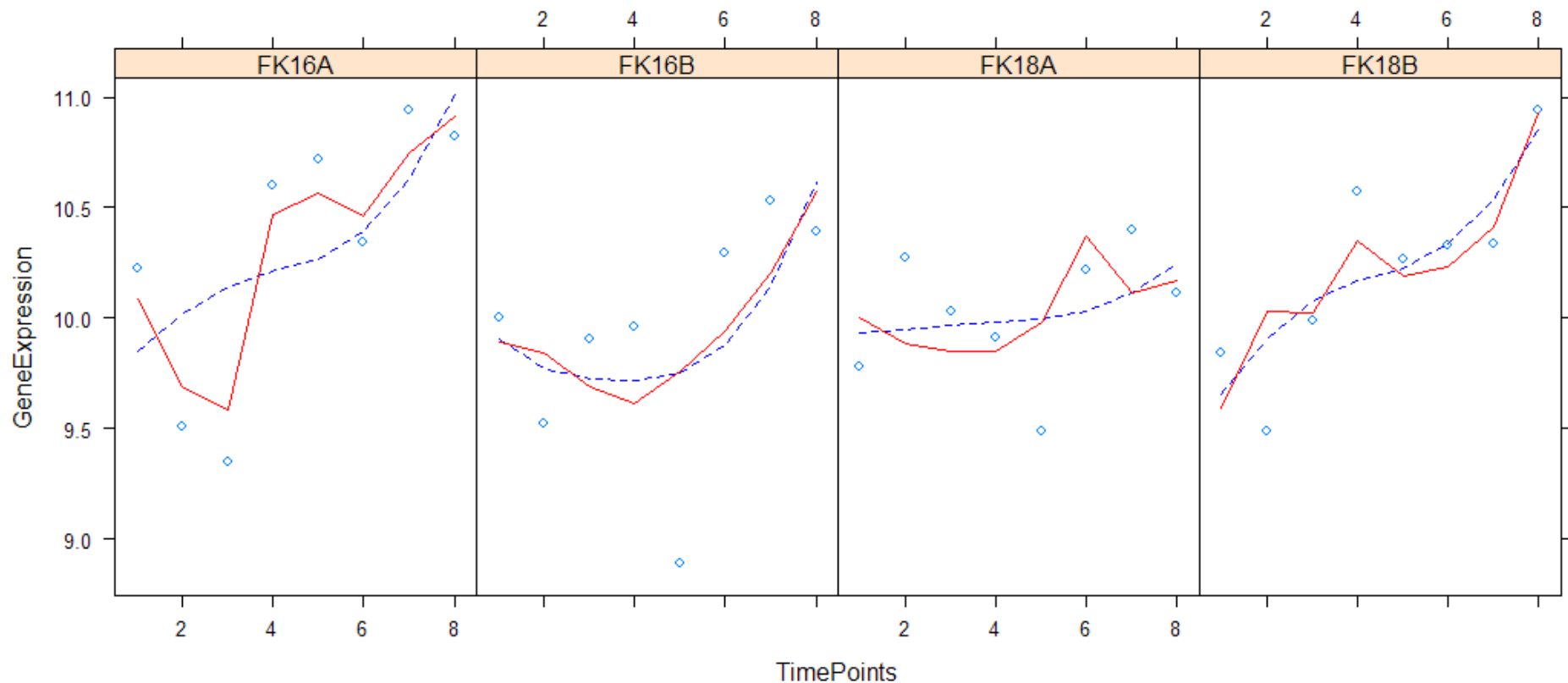
Copy number Effect

- — Cell line + CN + Time
- - - - Cell line + Time



Gene - SLC25A36

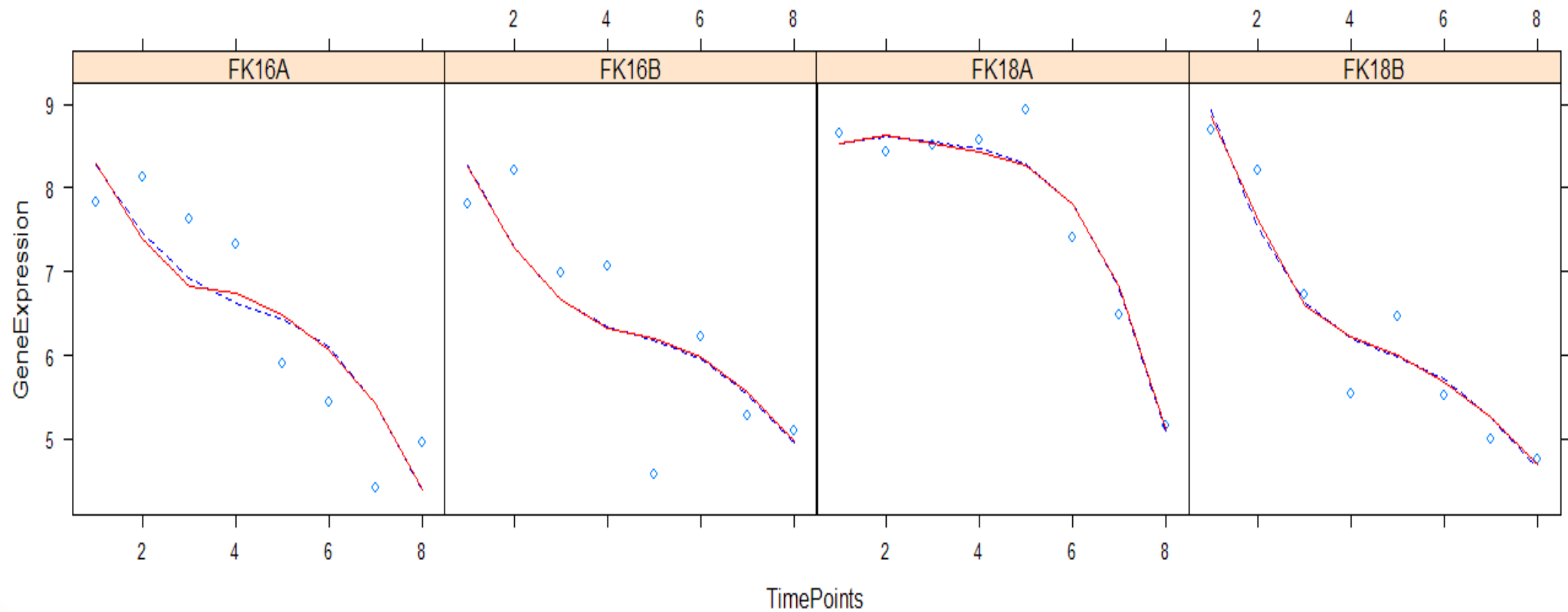
- Cell line + CN + Time
- Cell line + Time



Wilting et al., Genes,Cromosomes&Cancer, 2008

Gene - CADM1

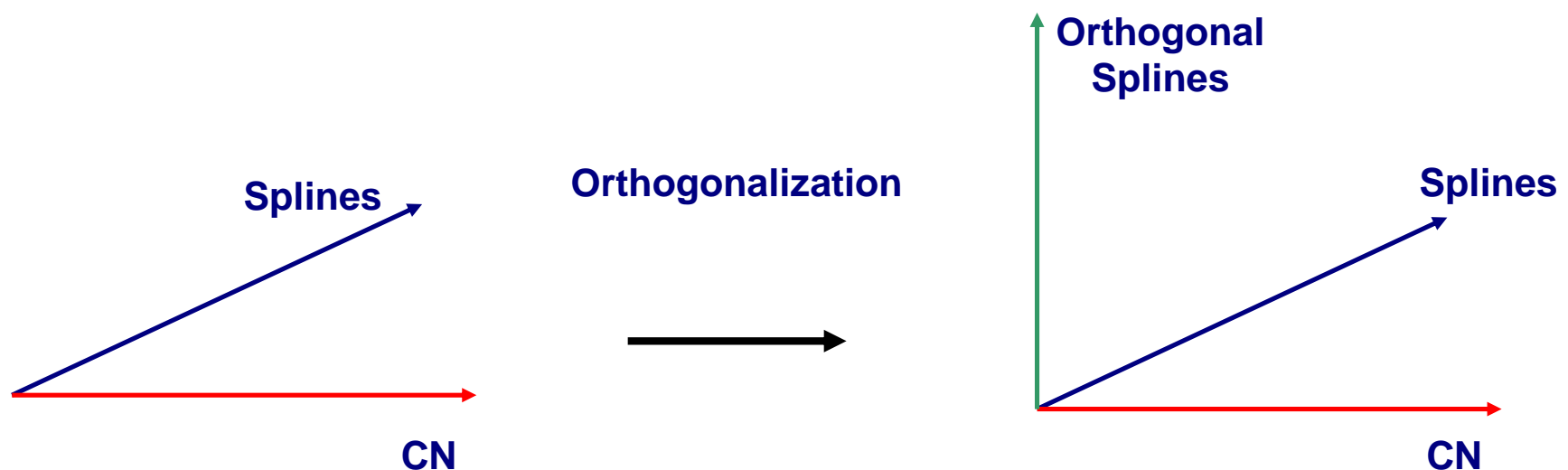
- — Cell line + CN + Time
- ⋯ Cell line + Time





Steenbergen et al., JNCI, 2004
Overmeer et al., J Pathol., 2008

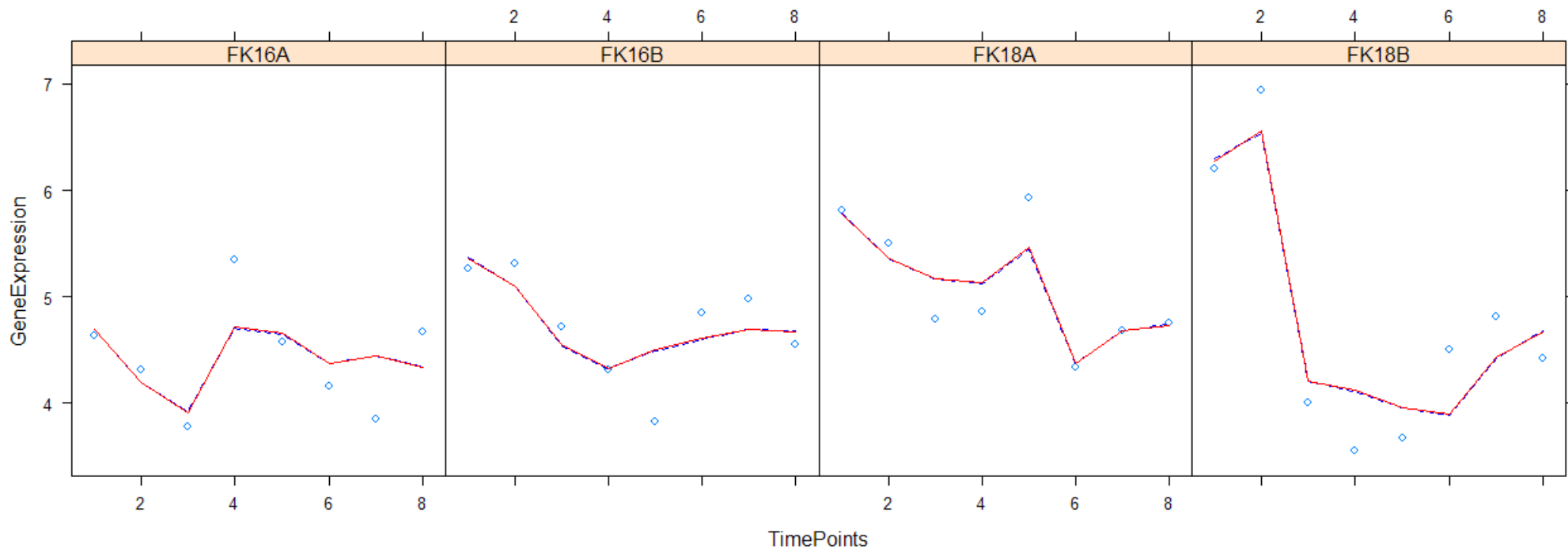
Competition between time and CN effect

- Time effect – flexibility of splines consume effect of copy number
- Potential solution
 - Orthogonalization of splines onto copy number design matrix



Fitting - orthogonal vs. standard

-  Cell Line + CN + orthogonal (Time)
-  Cell Line + CN + Time

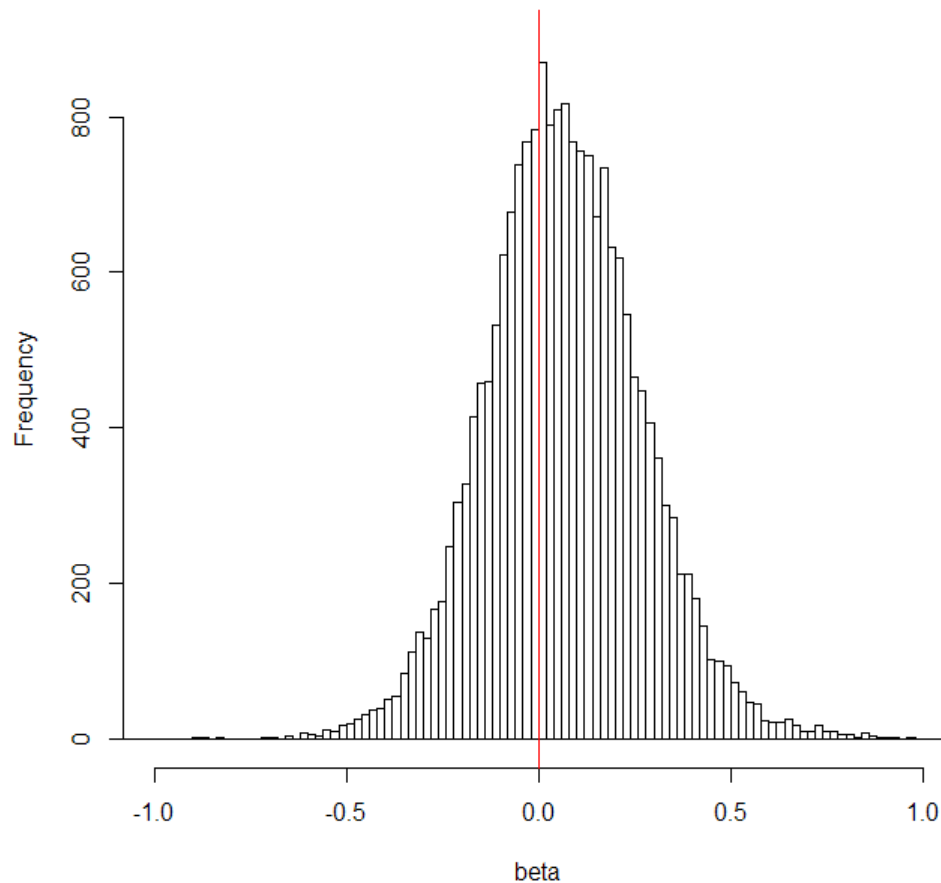


Parameter of CN effect

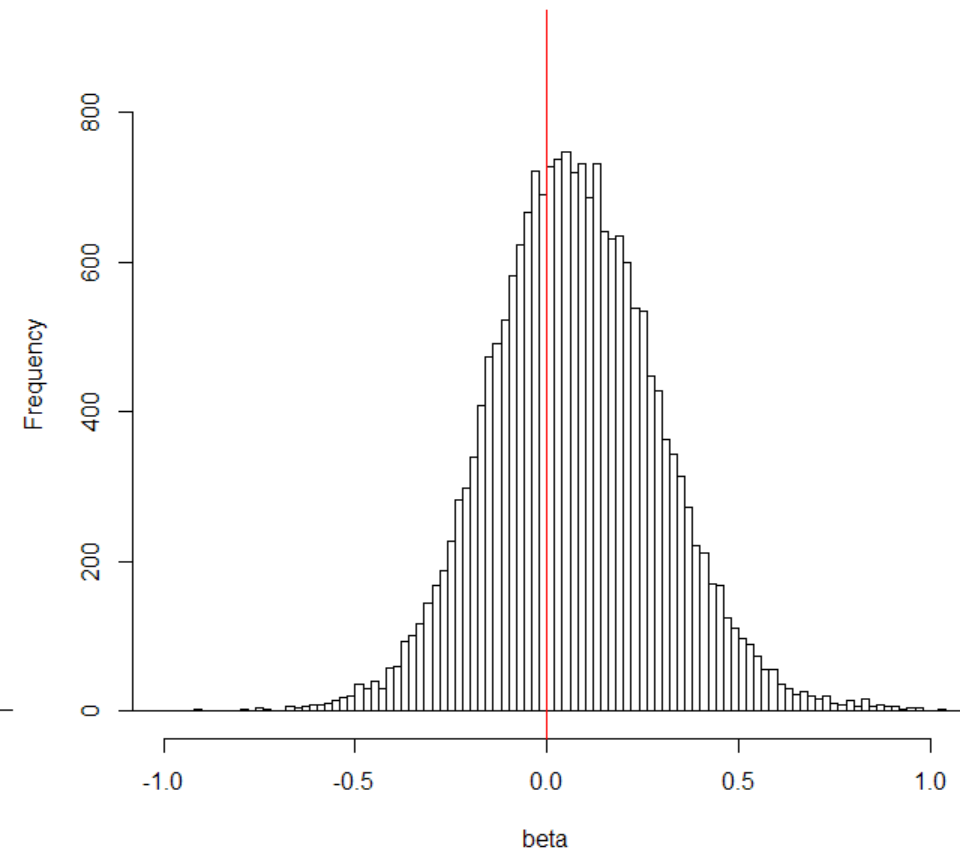
Standard

Orthogonalized





CN parameter - Standard

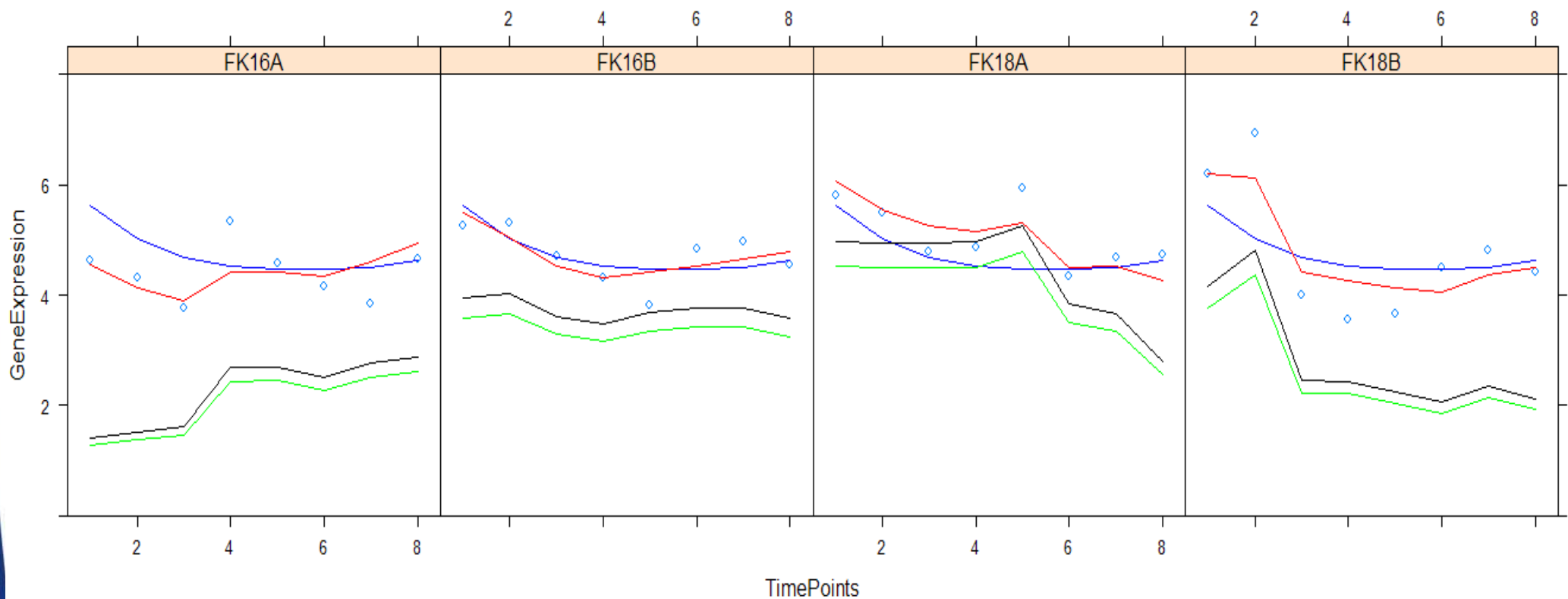


CN parameter - Orthogonalized



Fitting - orthogonal vs. standard design

-  Full model
-  Time
-  Orthogonal (CN)
-  CN



Comparison

Comparison

- Comparison of our method with
 - **timecourse** – Tai and Speed, Annals of Statistics, 2006.
 - **EDGE** – Storey et al., PNAS, 2005.
- Fair comparison – model without copy number
- Method is applied on two data sets
 - Data from our experiment
 - Data from Storey et al., PNAS, 2005.

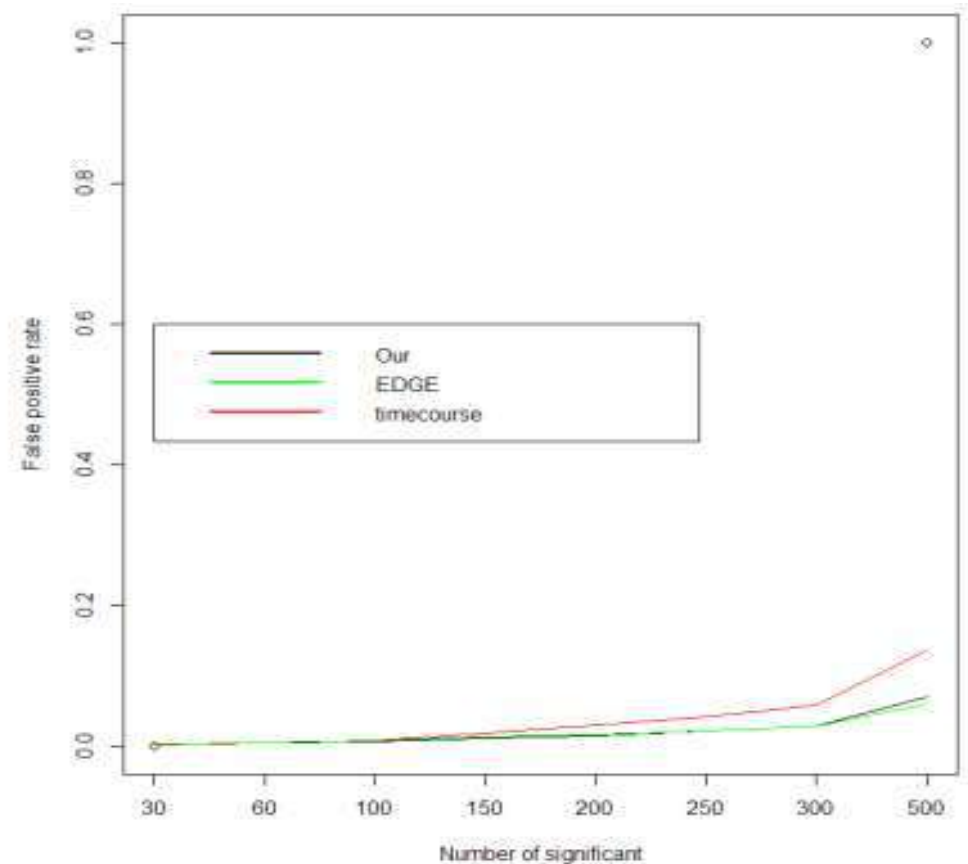
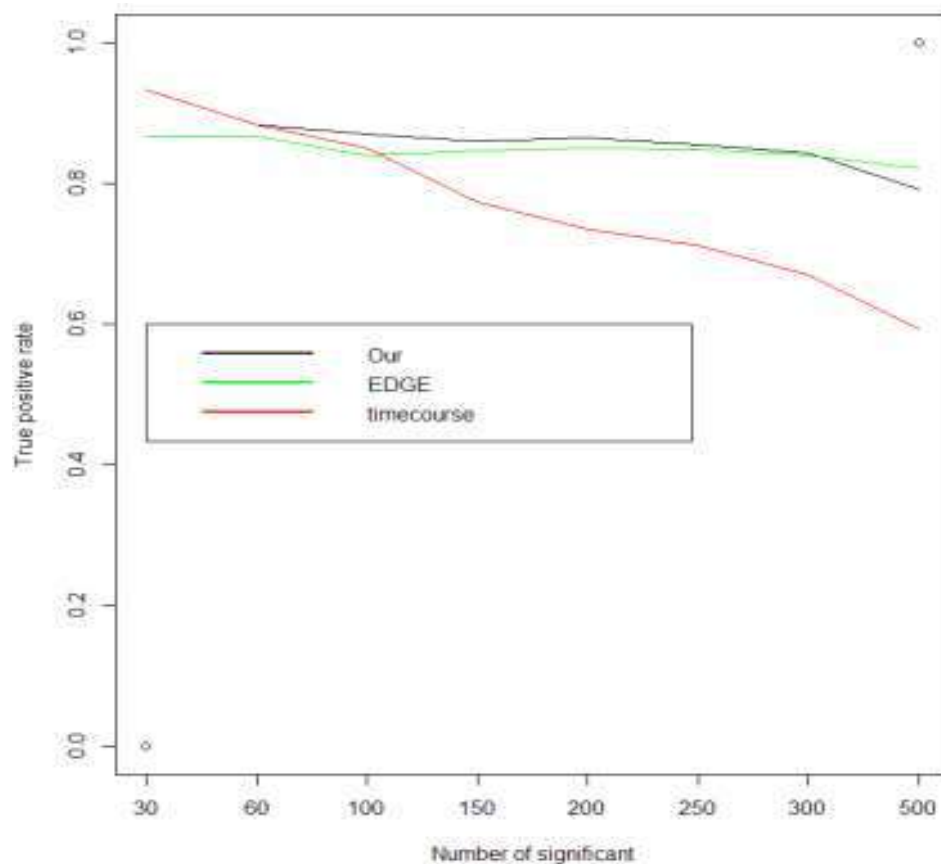
Comparison with **timecourse** and **EDGE**

- Sensitivity and specificity
 - Truth – significant genes among methods.
 - For different number of significant genes calculate true positive and false positive rate.
- Reproducibility
 - Equally divided data set in two groups
 - For different number of significant genes calculate number of overlap genes.

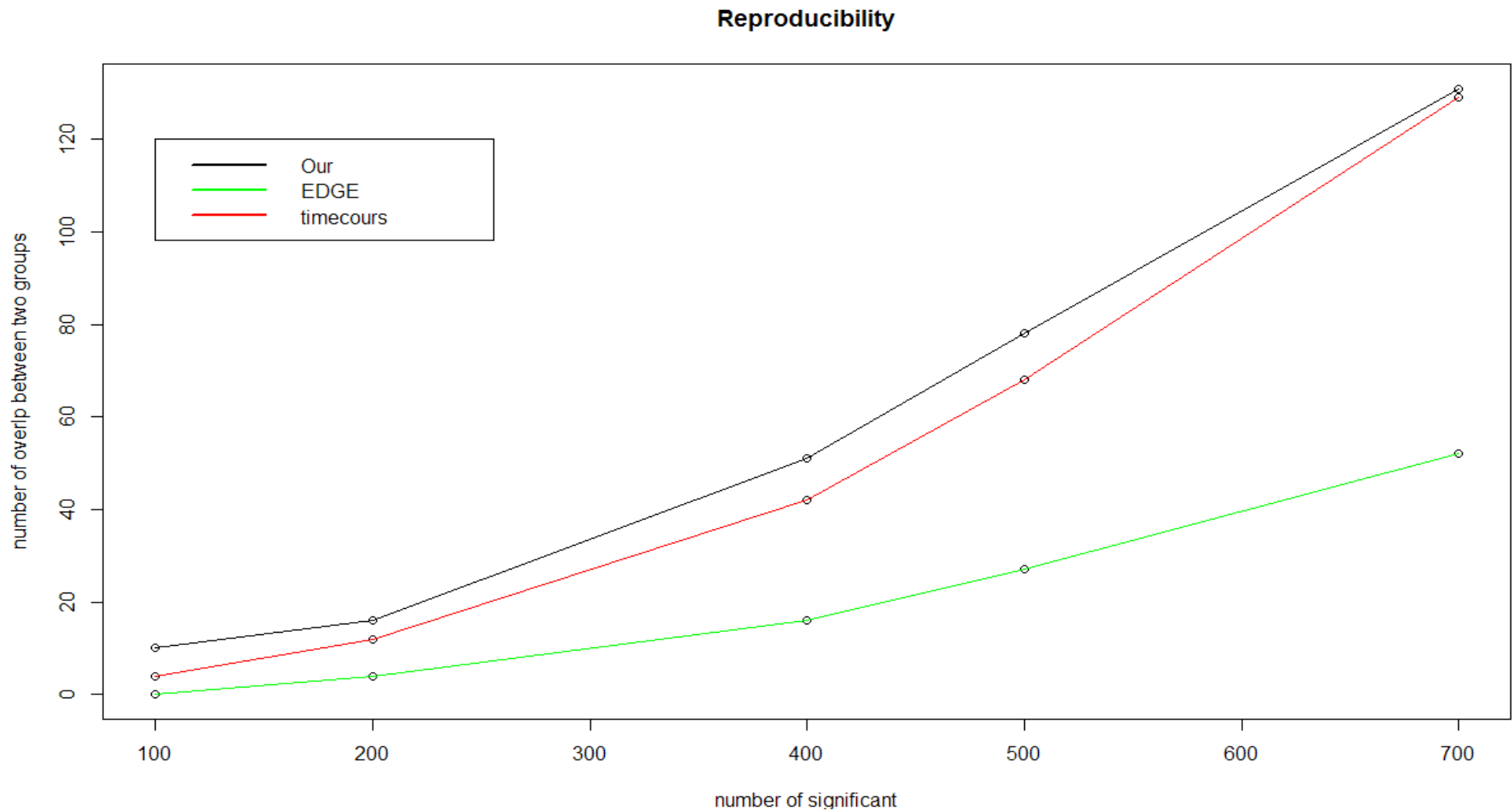
Sensitivity and specificity – our data

- Sensitivity

- Specificity

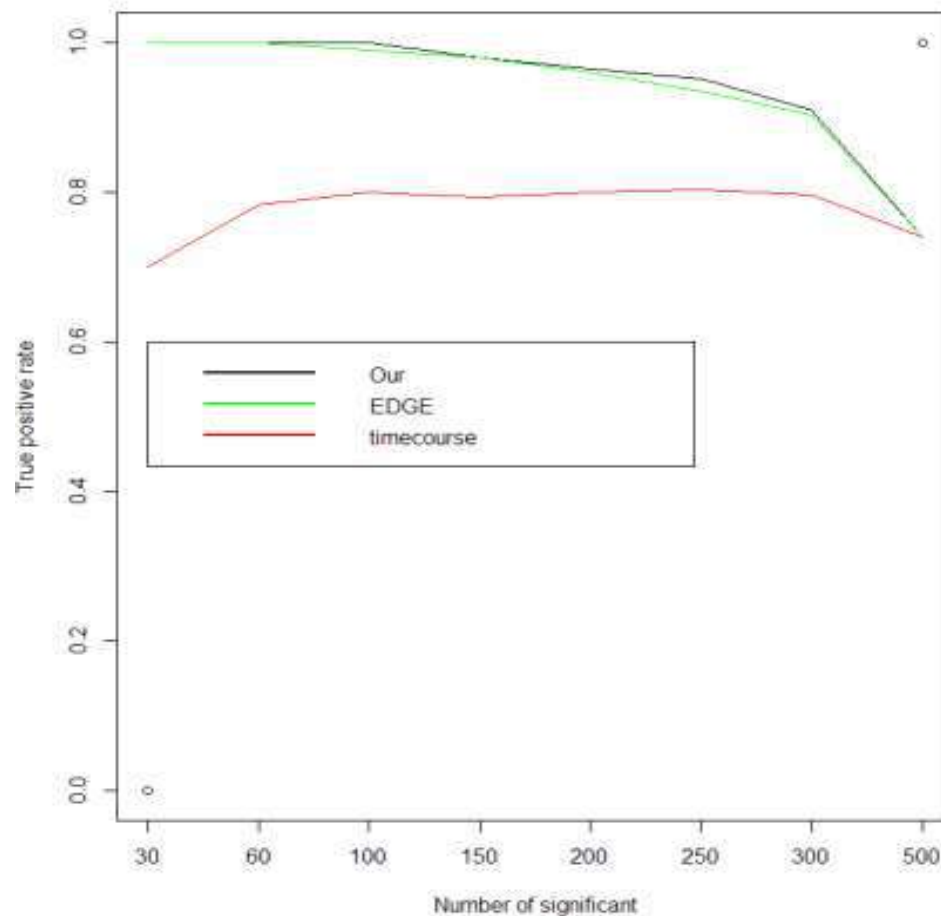


Reproducibility – our data

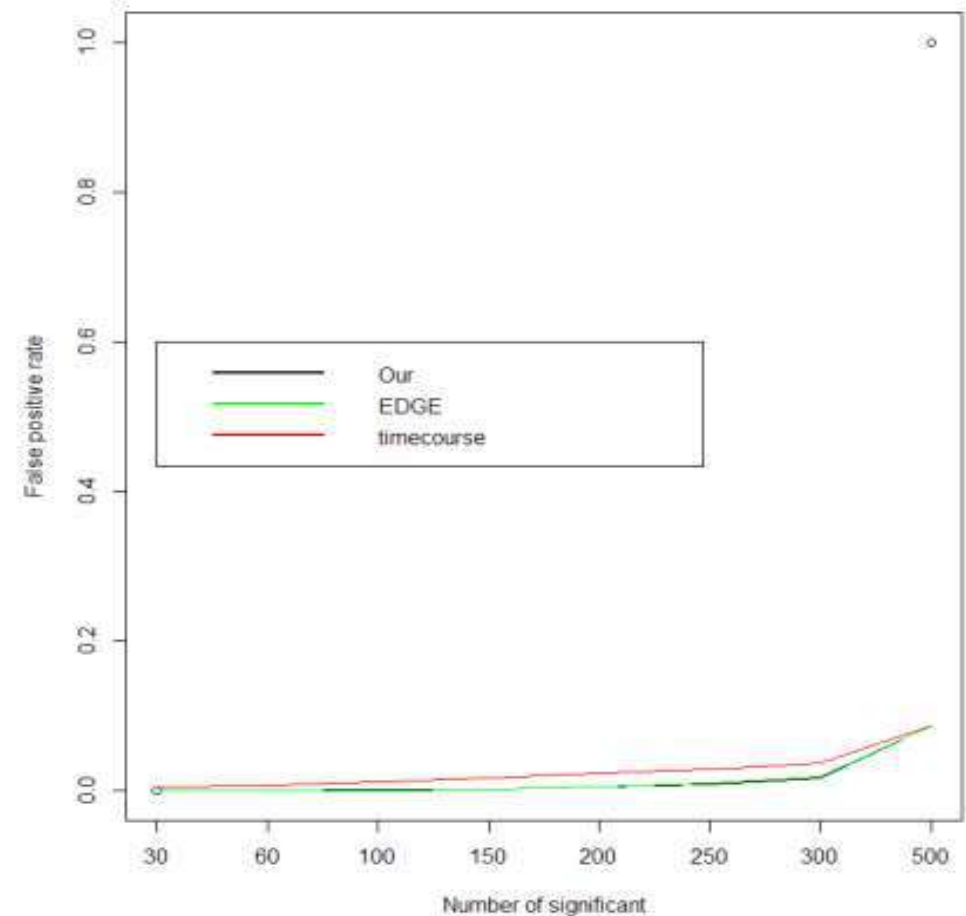


Sensitivity and specificity - EDGE data

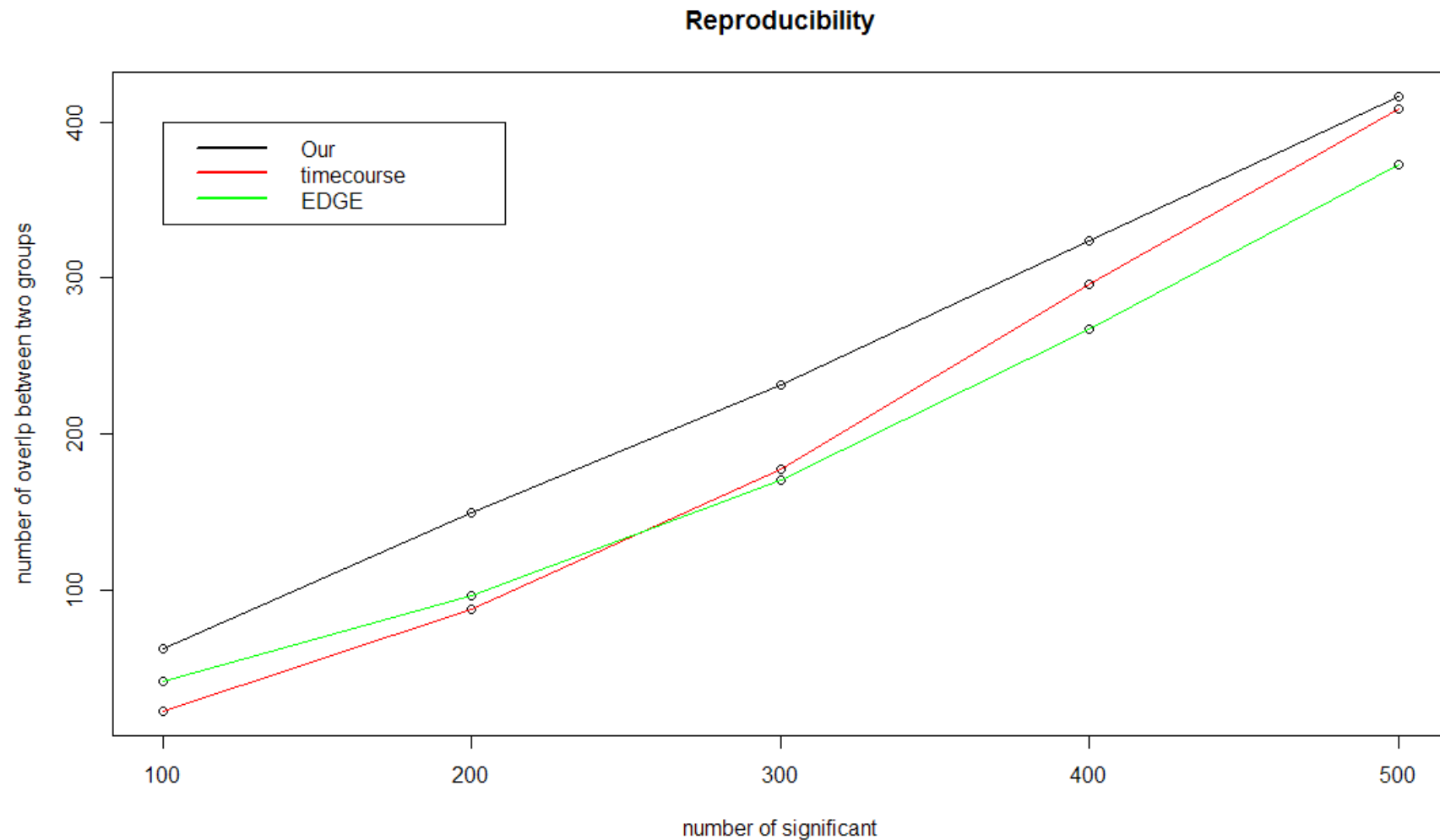
- Sensitivity



- Specificity



Reproducibility – EDGE data



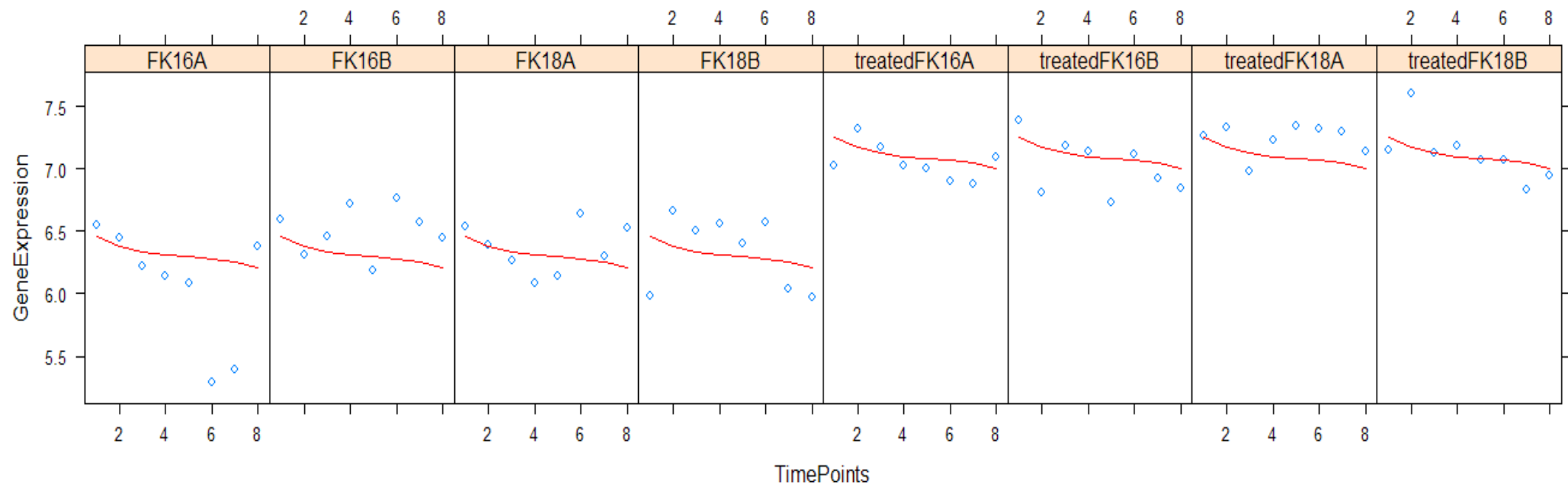
Application

Testing for differential gene expression

Number of significant genes at 5% FDR level	Same Spline		Different Spline	
	Standard	Orthogonal	Standard	Orthogonal
CN Effect	524	544	923	936
Time Effect	4959	5011	10487	10567

Testing for treatment effect

- Treatment effect - 16976 genes significant at 5% FDR level



Summary

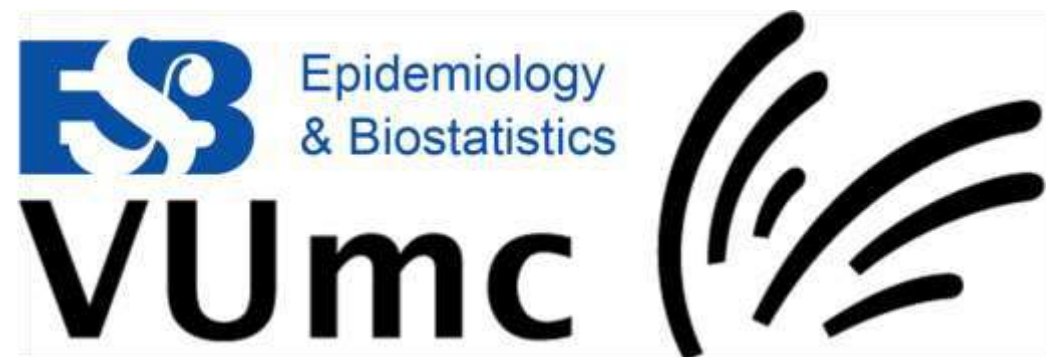
- Generated a reliable multilevel molecular dataset.
- Developed generalized linear mixed model for identification of temporal differential gene expression.
- Improvement in identification of variation in genes due to copy number effect
- Method show improvements in sensitivity, specificity and reproducibility comparing with other methods.

Future plans

- Developing methodology to include spatial effect over genome for copy number
- Including DNA methylation in model
- Integrative analysis for miRNA
 - Applying method on miRNA data
 - Pathway analysis

Acknowledgments

- Wessel van Wieringen
- Saskia Wilting
- Annelieke Jaspers
- Mark van de Wiel
- Peter Snijders
- Renske Steenbergen



**Thank you for your
attention**