# Scalable
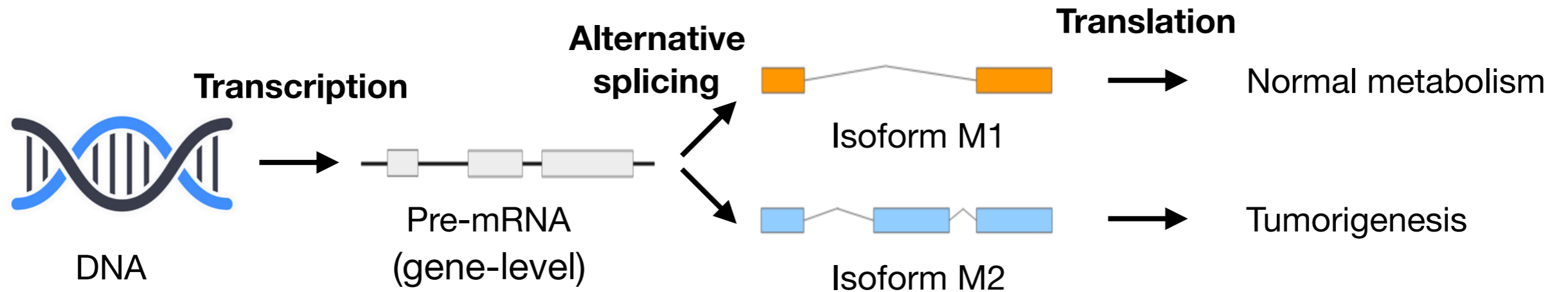# differential transcript usage analysis
# for single-cell applications

JEROEN GILIS

EuroBioc2019 presentation
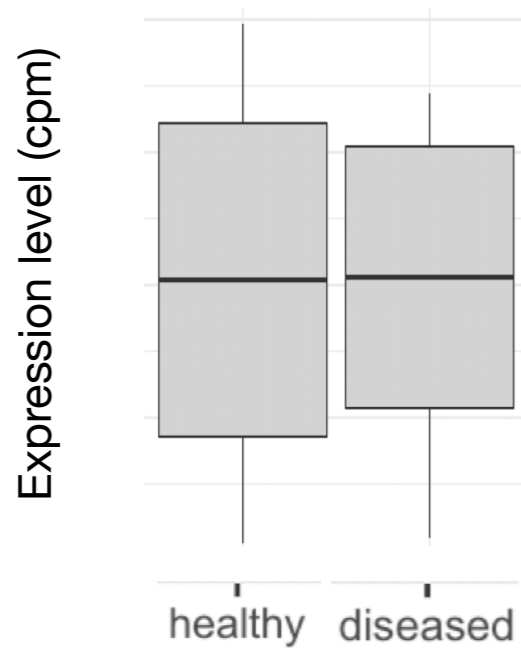
Promotor: Prof. Lieven Clement
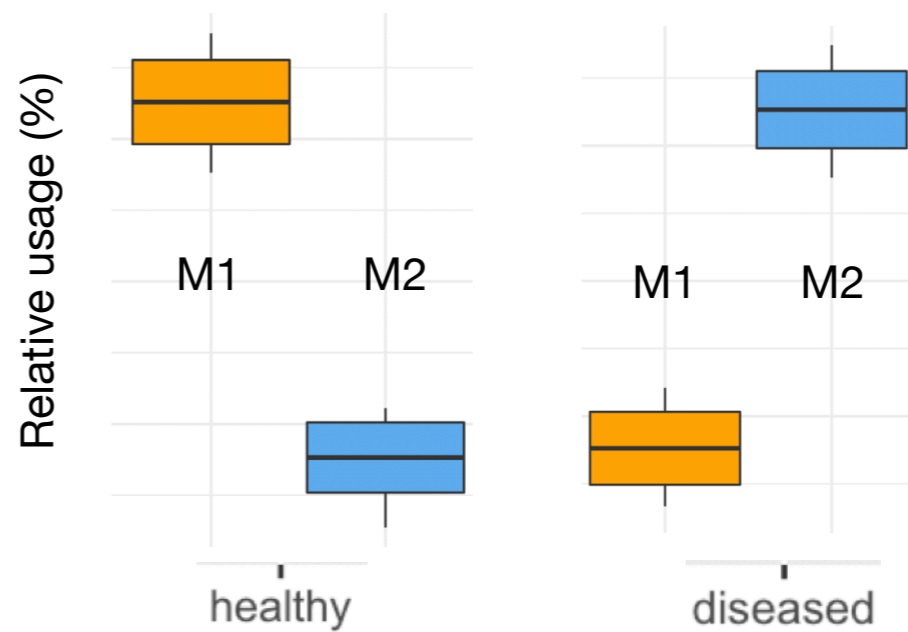
Supervisor: Dr. Koen Van den Berge

# Differential Transcript Usage (DTU)

**Method development**

• Our workflow unlocks edgeR for DTU analysis

**DGE**
$$Y_{gi} \sim NB\,(\mu_{gi},\, \varphi_g)$$
$$\log\,(\mu_{gi}) = \eta_{gi}$$
$$\eta_{gi} = \beta_0 + \beta_{gc}^{C} + \log\,(S_i)$$

**Method development**

• Our workflow unlocks edgeR for DTU analysis

**DTE**
$$Y_{ti} \sim NB\,(\mu_{ti},\,\varphi_t)$$
$$\log\,(\mu_{ti})\,=\,\eta_{ti}$$
$$\eta_{ti}\,=\,\beta_0\,+\,\beta_{tc}^{\,c}+\log\,(S_i)$$

**Method development**

- Our workflow unlocks edgeR for DTU analysis

$$\mathbf{DTU} \left\{ \begin{array}{l} \mathrm{Y}_{ti} \sim NB\,(\mu_{ti},\,\varphi_t) \\[1em] \log\,(\mu_{ti}) \;=\; \eta_{ti} \\[1em] \eta_{ti} \;=\; \beta_0 \;+\; \beta_{tc}^{\,c} + \log\,(\textcolor{red}{\mathrm{T}_{ti}}) \end{array} \right.$$

- Our workflow takes the **gene-level counts (total counts, $T_{ti}$) as offsets** to the GLM framework ⟶ edgeR-total

5

**Method development**

- Our workflow unlocks edgeR for DTU analysis

**DTU**
$$Y_{ti} \sim NB(\mu_{ti}, \varphi_t)$$
$$\log(\mu_{ti}) = \eta_{ti}$$
$$\eta_{ti} = \beta_0 + \beta_{tc}^C + \log(T_{ti})$$

- Our workflow takes the **gene-level counts (total counts, $T_{ti}$) as offsets** to the GLM framework $\longrightarrow$ edgeR-total

- DEXSeq

| Counts | | Sample 1 | … | Sample m | | Sample 1 | … | Sample m | 'other' counts |
|---|---|---|---|---|---|---|---|---|---|
| | Tx 1 | 112 | … | 15 | Tx 1 | 25 | … | 3 | |
| | Tx t | … | … | … | Tx t | … | … | … | |
| | Tx n | 62 | … | 348 | Tx n | 88 | … | 212 | |

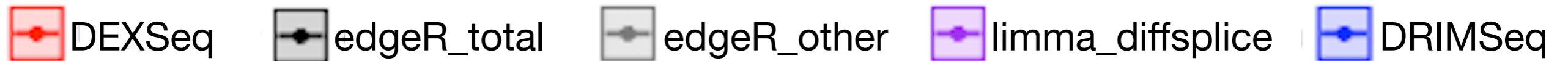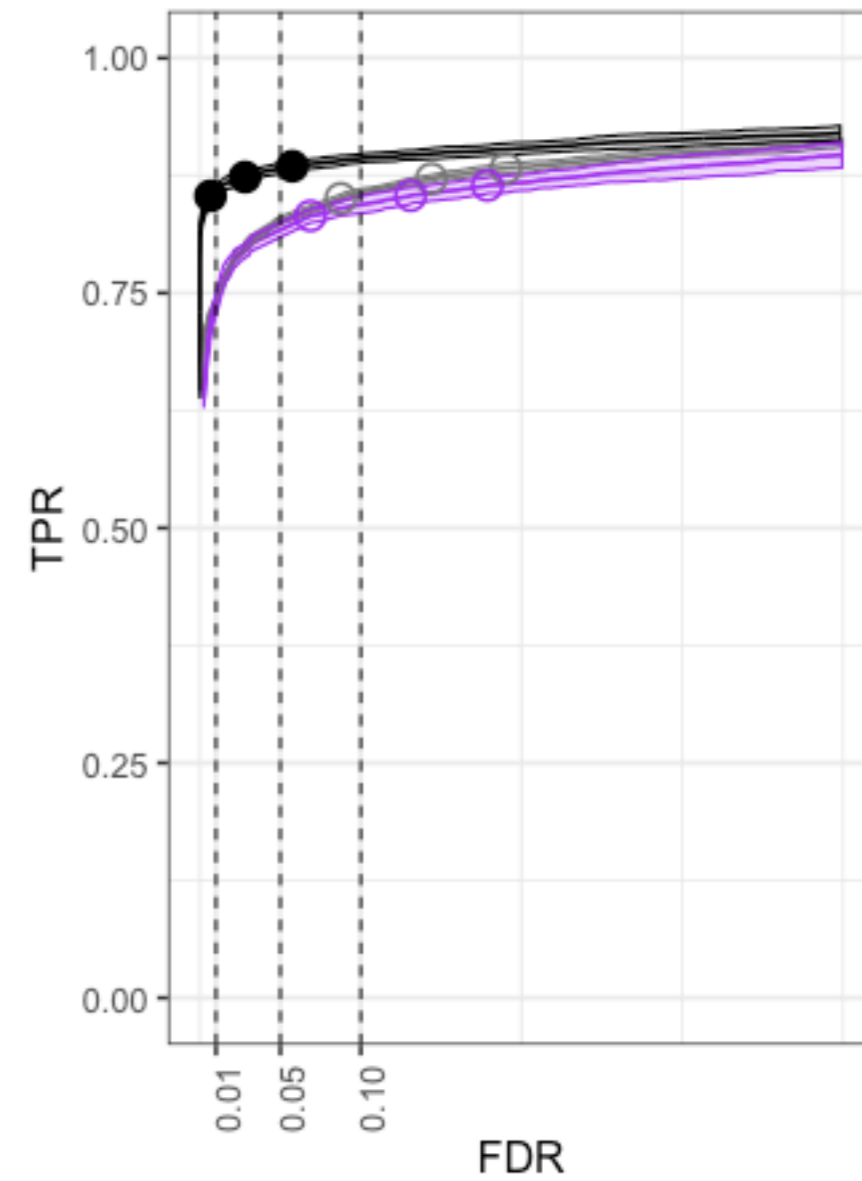- Our second workflow takes the **other counts as offsets** $\longrightarrow$ edgeR-other

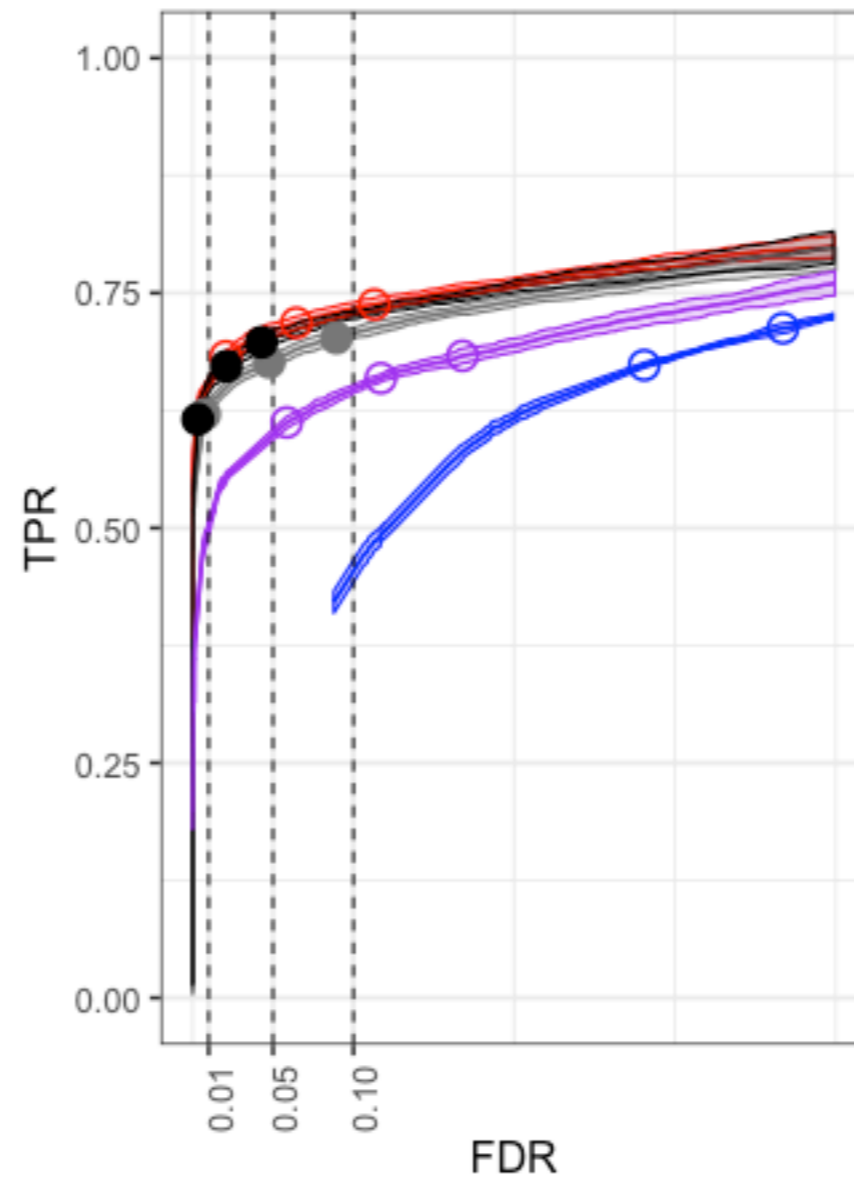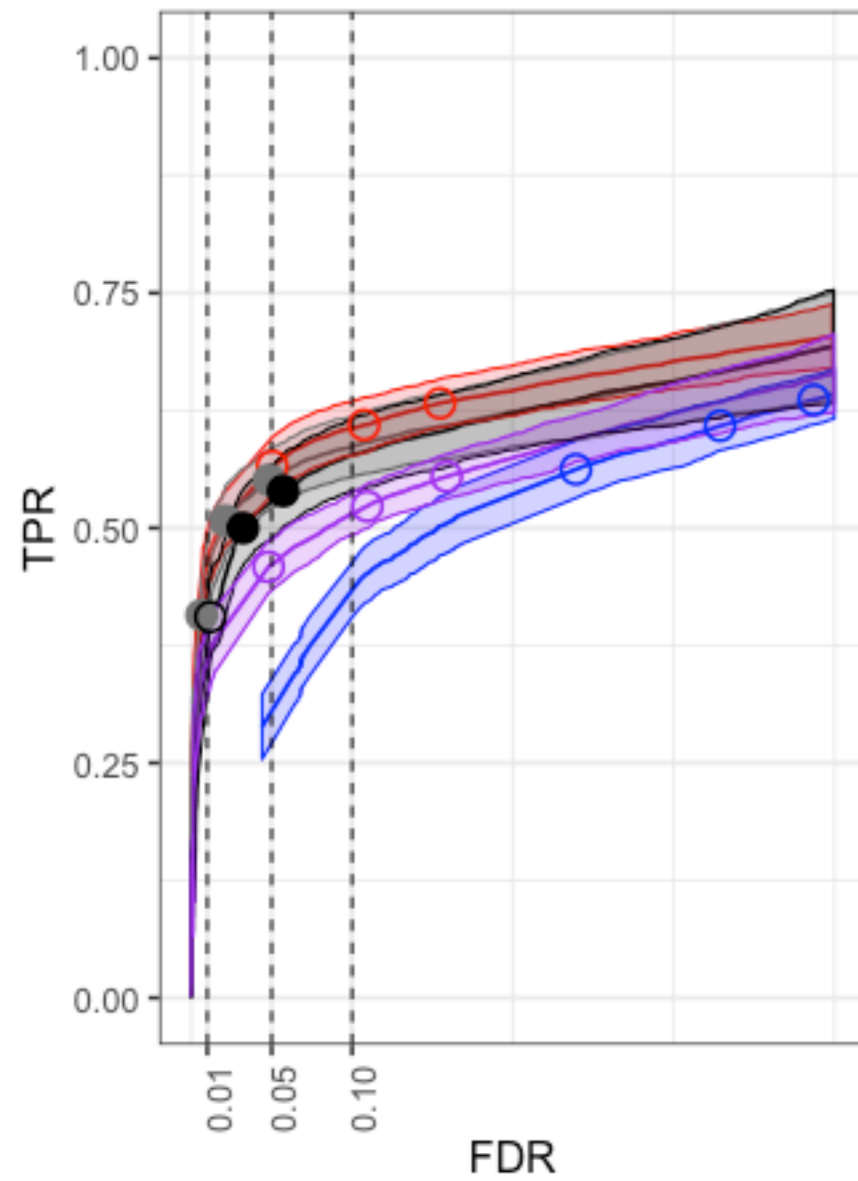# Performance evaluation on real bulk data
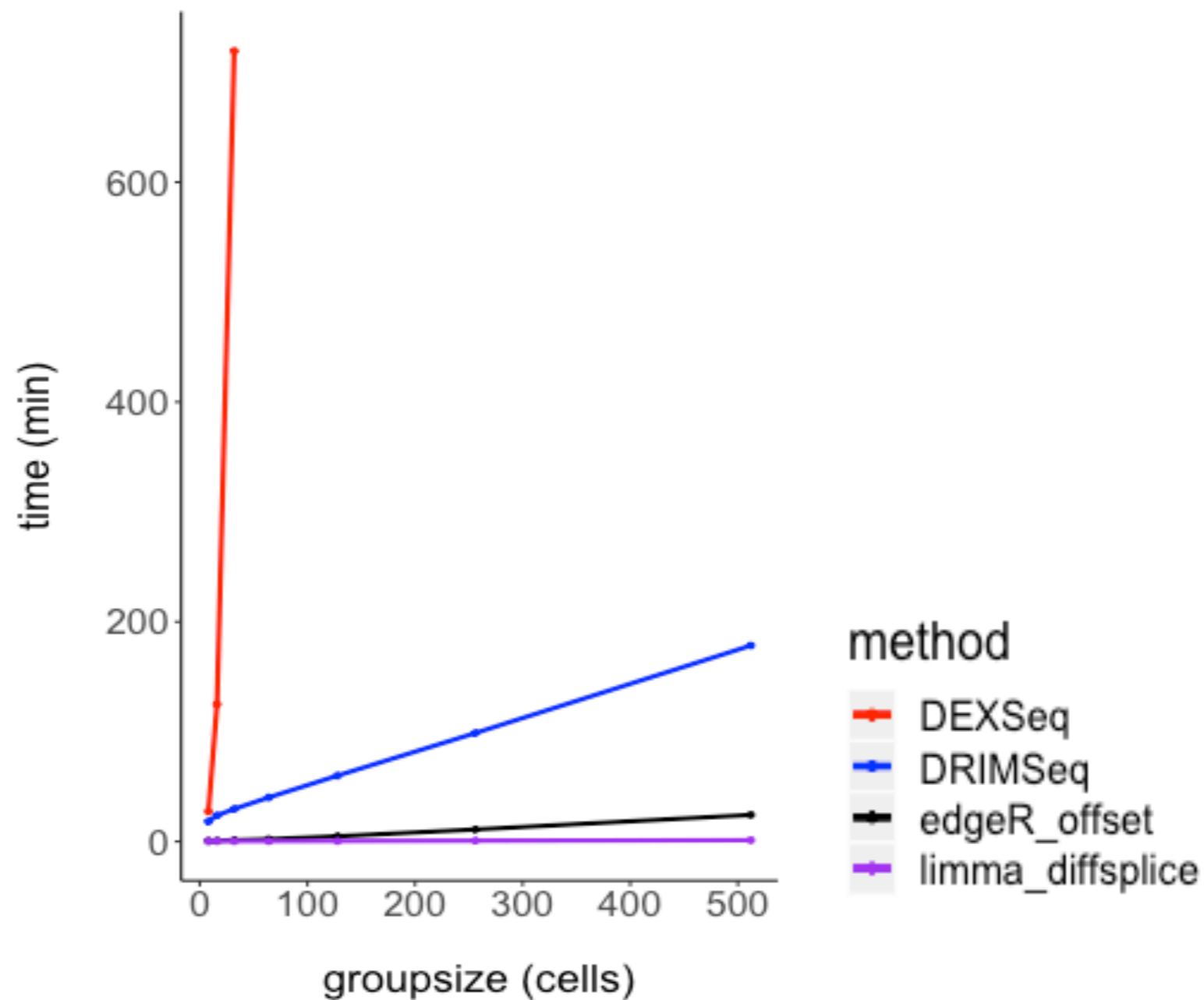
*Gtex dataset, Nature Genetics 45, 580-585 (2013)*

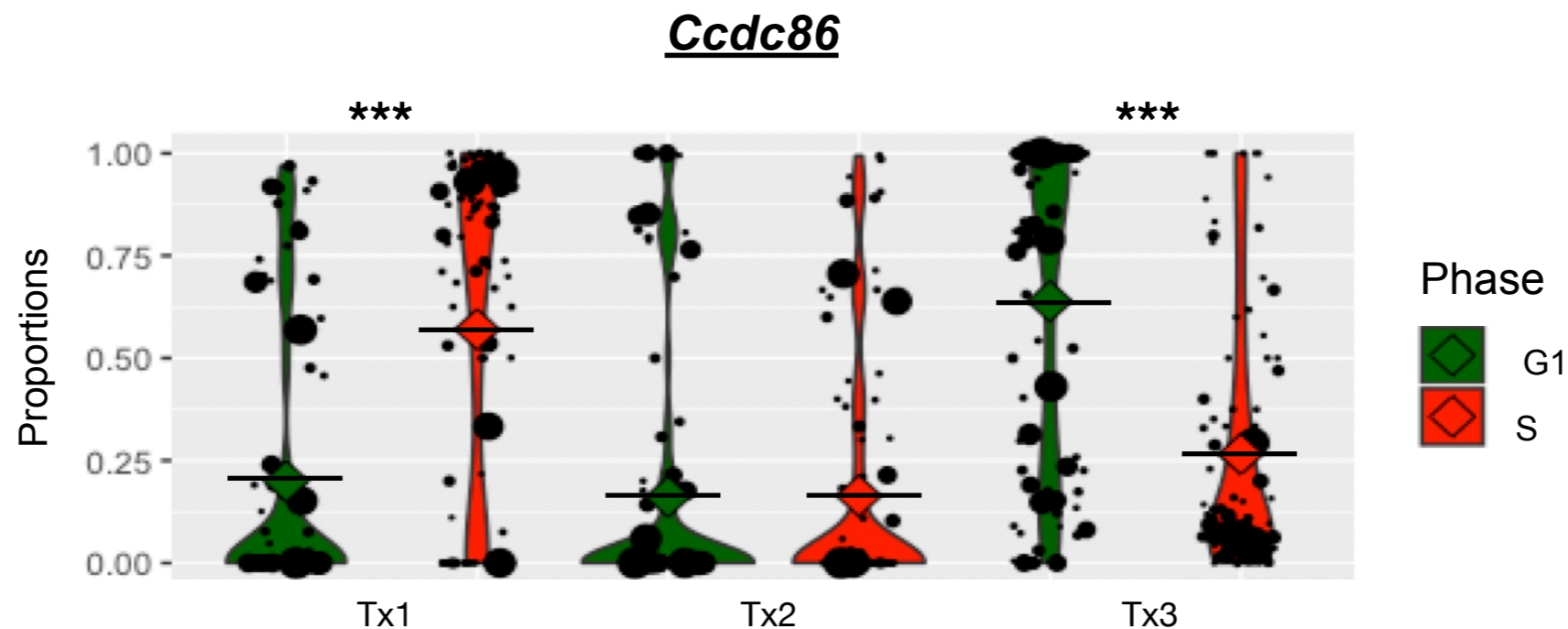# Scalability benchmark on real single-cell data

- Our workflow performs a DTU analysis between two groups of 512 cells in ~20 minutes
- DEXSeq scales quadratically

# Single-cell transcriptomics case study

*Dataset from Buettner et al., Nature Biotechnology 33; 155-160 (2015)*

- Dataset; 288 mouse embryonic stem cells, different cell cycle stages (G1, S and G2M)

- Runtime; < 2 minutes

- Significant enrichment in cell cycle processes

- Several DTU genes are;

    ✦ Biologically relevant

    ✦ Not picked up in a gene-level analysis

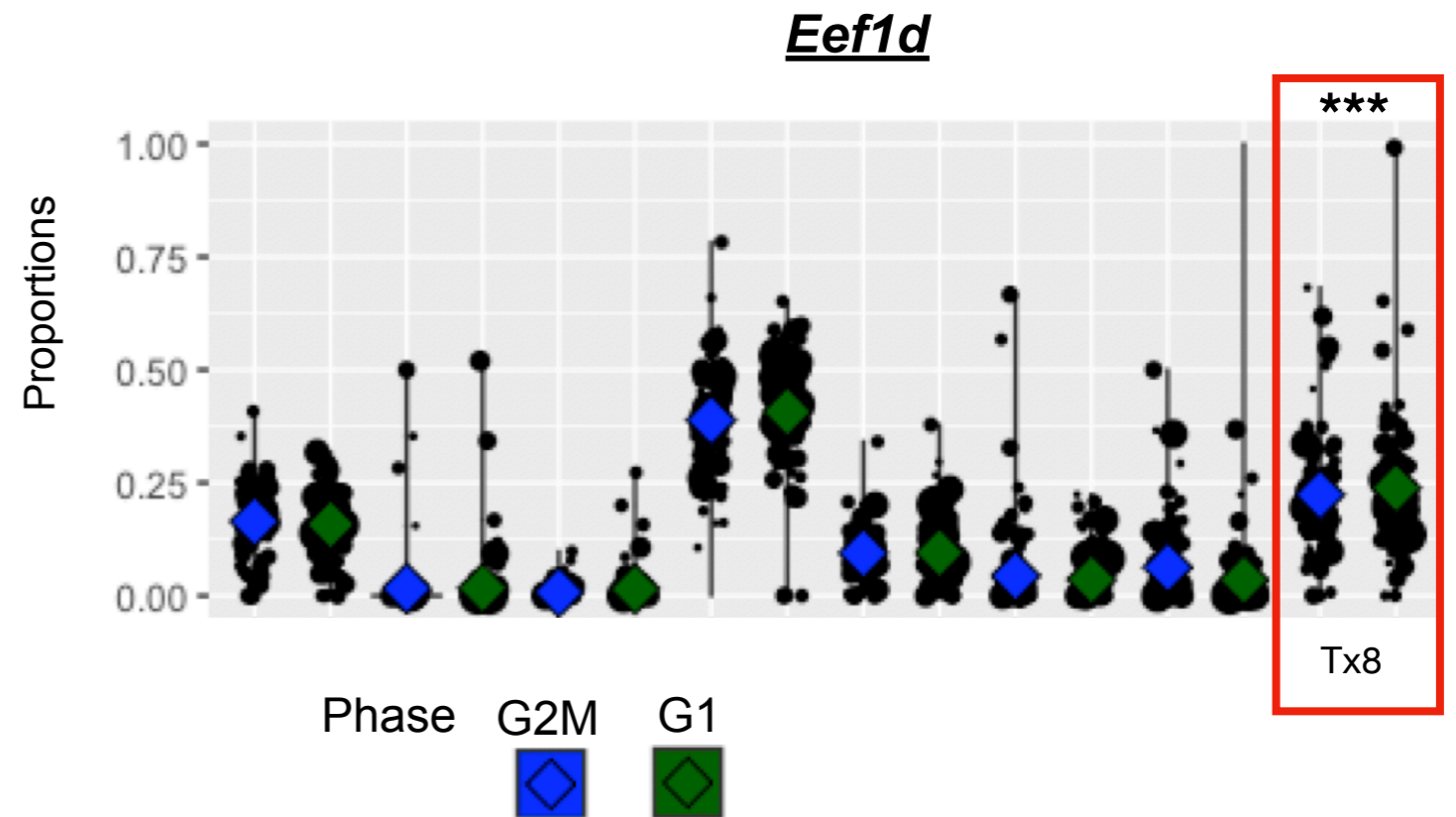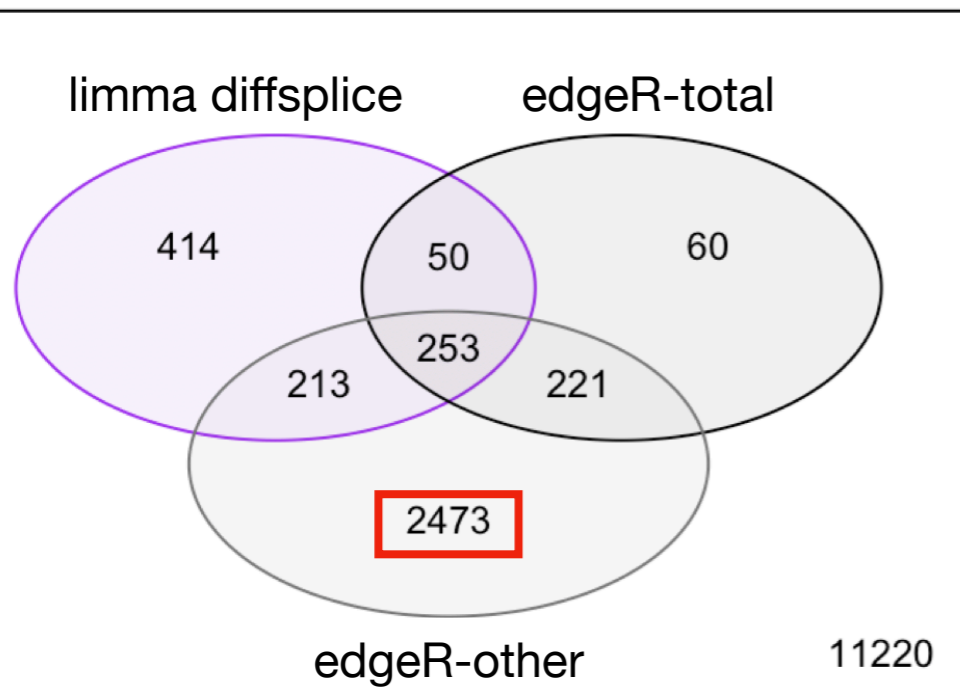    ✦ Clearly differentially used when visualised



The size of the dots (which represent individual cells) are weighted according to the total expression of the gene in that cell.

# Single-cell transcriptomics case study

*Buettner dataset, Nature Biotechnology 33; 155-160 (2015)*

- Dataset; 288 mouse embryonic stem cells, different cell cycle stages (G1, S and G2M)

- Runtime; < 2 minutes for offset-based methods

- Significant enrichment in cell cycle processes

- Some DTU genes display clear DTU in visualisation and are biologically relevant

- edgeR_other method large number of (false) positive results; sensitive to outliers (?)

- Discrepancy between edgeR-total and limma diffsplice; asses formally in single-cell benchmark

# Take-home messages

We are developing a workflow for studying DTU that;

1. Has a performance similar to that of DEXSeq

2. Correctly controls the false discovery rate

3. Scales towards large transcriptomics datasets

# Scalable
# differential transcript usage analysis
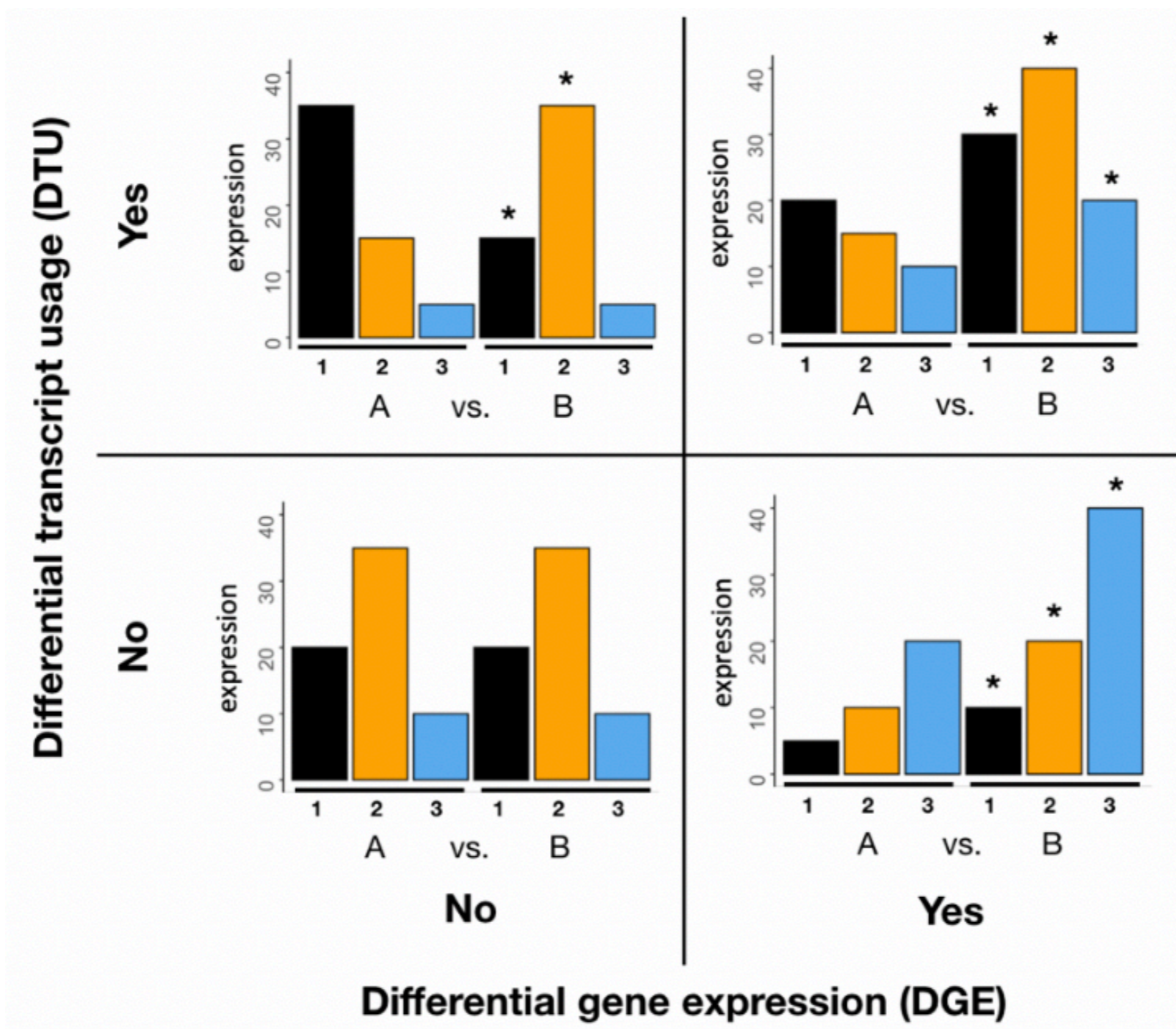# for single-cell applications

JEROEN GILIS

EuroBioc2019 presentation

Promotor: Prof. Lieven Clement

Supervisor: Dr. Koen Van den Berge

# Background - DTU

# Background - DEXSeq

- **Input**: matrix of transcript-level counts (e.g. Salmon or kallisto)

**Transcript-level counts**

|  |  | Sample 1 | Sample 2 | … |
|---|---|---|---|---|
| **Gene A** | Transcript 1 | 20 | 18 | … |
|  | Transcript 2 | 10 | 7 | … |
|  | Transcript 3 | 70 | 45 | … |
| **Gene B** | Transcript 1 | 22 | 0 | … |
|  | Transcript 2 | 3 | 16 | … |
| … | … | … | … | … |

**Complementary counts**

|  |  | Sample 1 | Sample 2 | … |
|---|---|---|---|---|
| **Gene A** | Transcript 1 | 80 | 52 | … |
|  | Transcript 2 | 90 | 63 | … |
|  | Transcript 3 | 30 | 25 | … |
| **Gene B** | Transcript 1 | 3 | 16 | … |
|  | Transcript 2 | 22 | 0 | … |
| … | … | … | … | … |

- **Statistical model:**

$$
\begin{cases}
Y_{ti} \sim NB\left(\mu_{ti}, \varphi_t\right) \\[2mm]
\log\left(\mu_{ti}\right) = \eta_{ti} \\[2mm]
\eta_{ti} = \boldsymbol{\beta}_{ti}^{S} + \boldsymbol{\beta}_{t}^{T} + \boldsymbol{\beta}_{tc_i}^{TC}
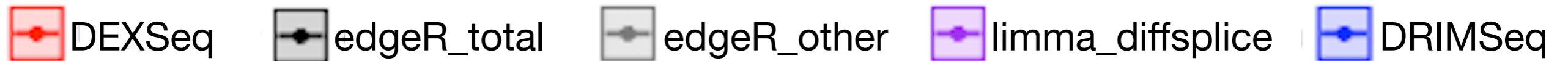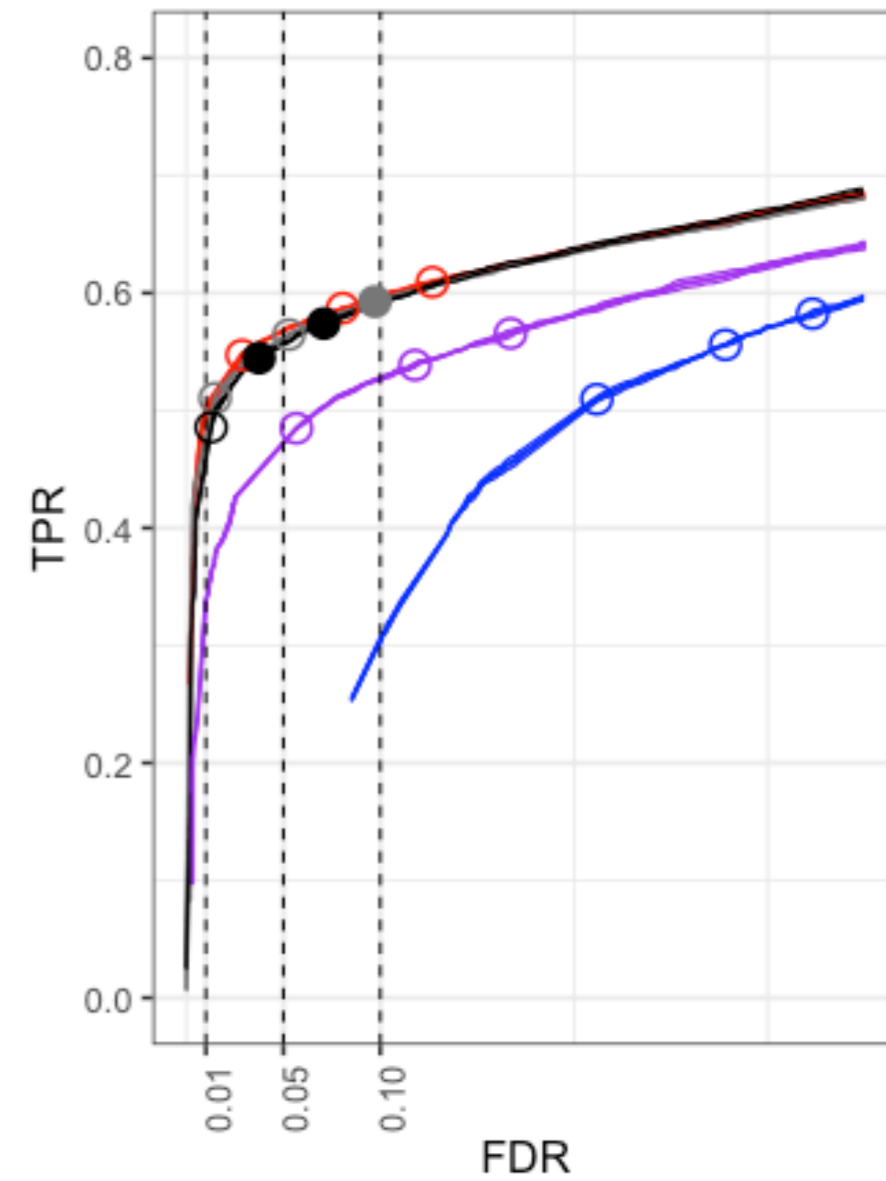\end{cases}
$$

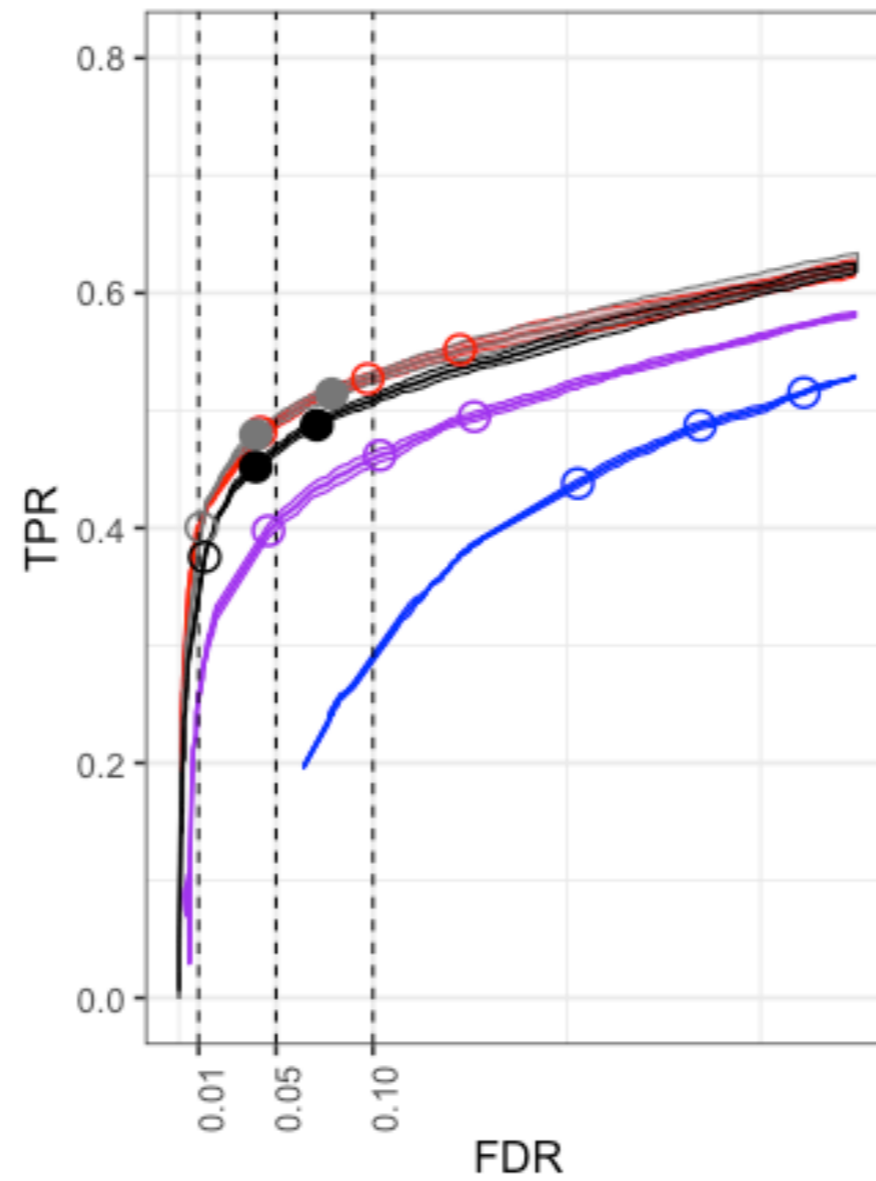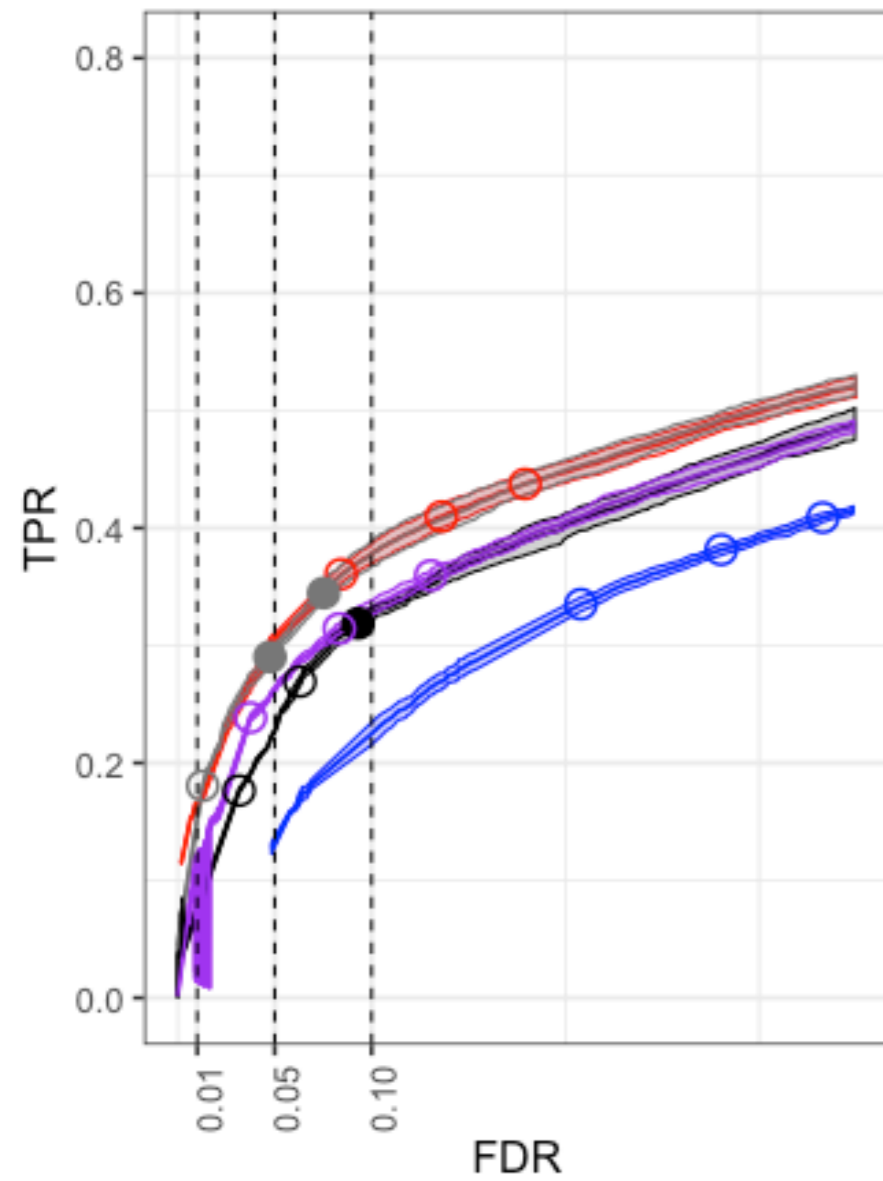# Parametric bulk simulation study

*Dataset from Love et al., F1000Research, 7:952 (2018)*
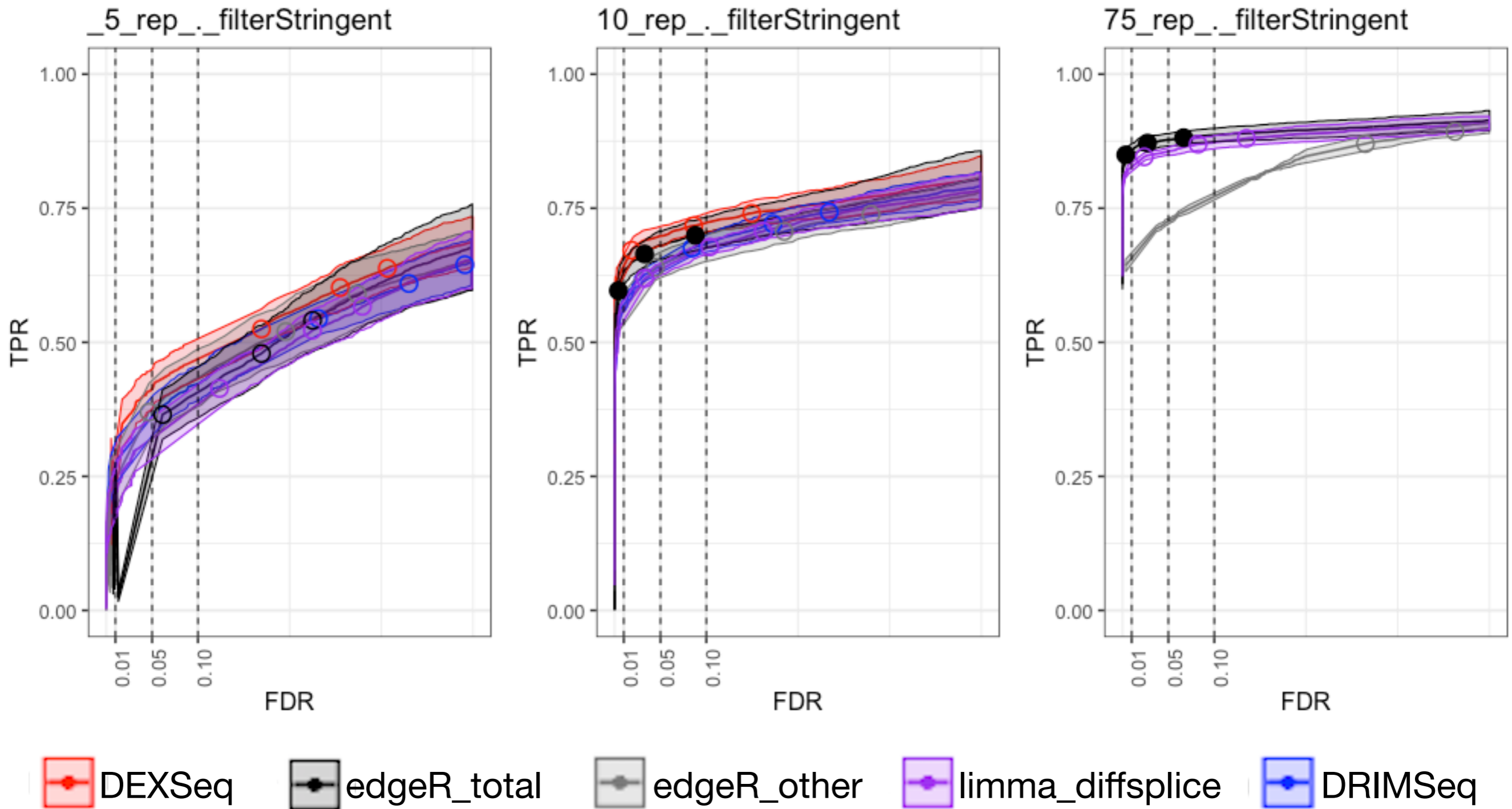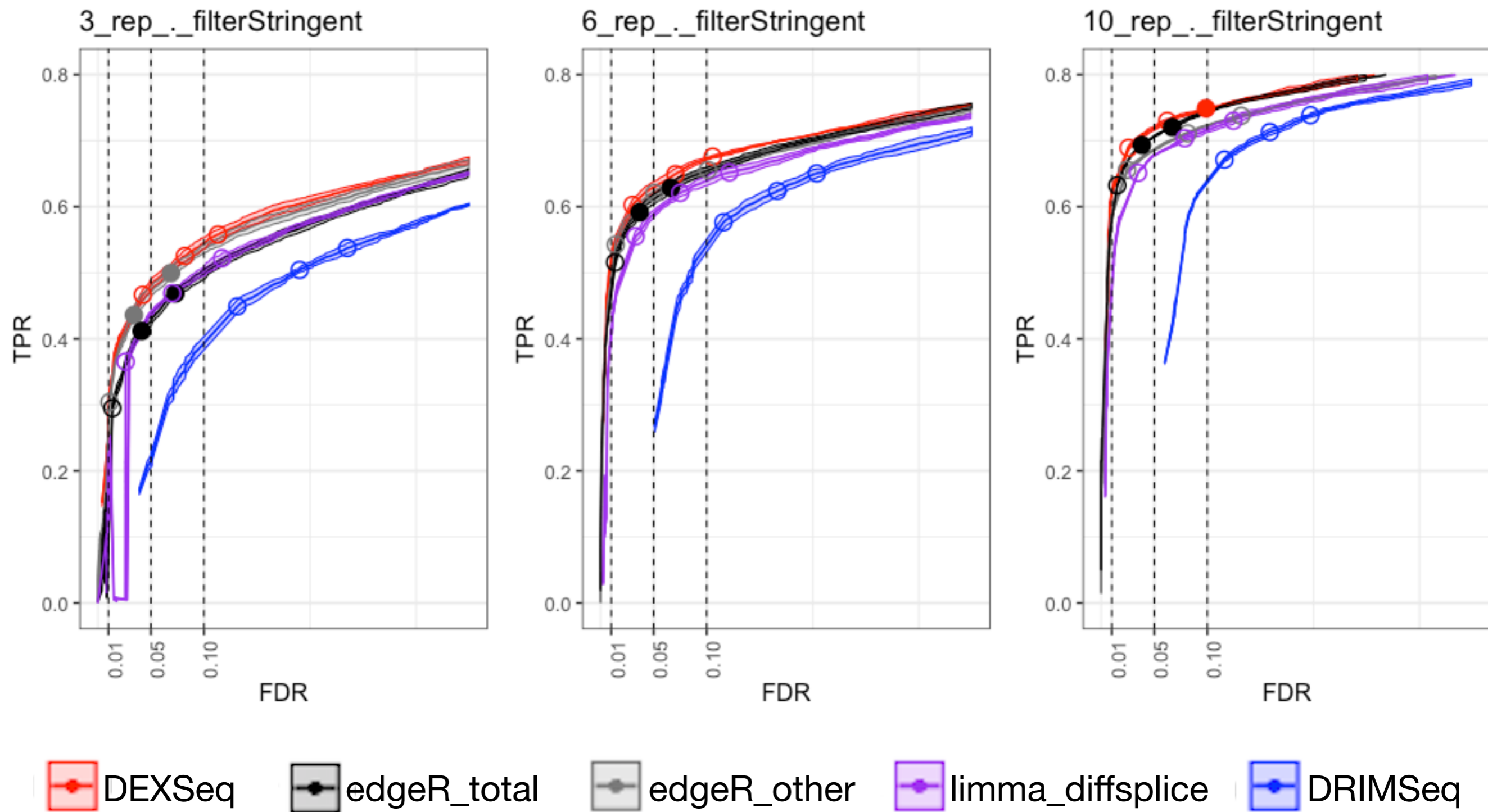
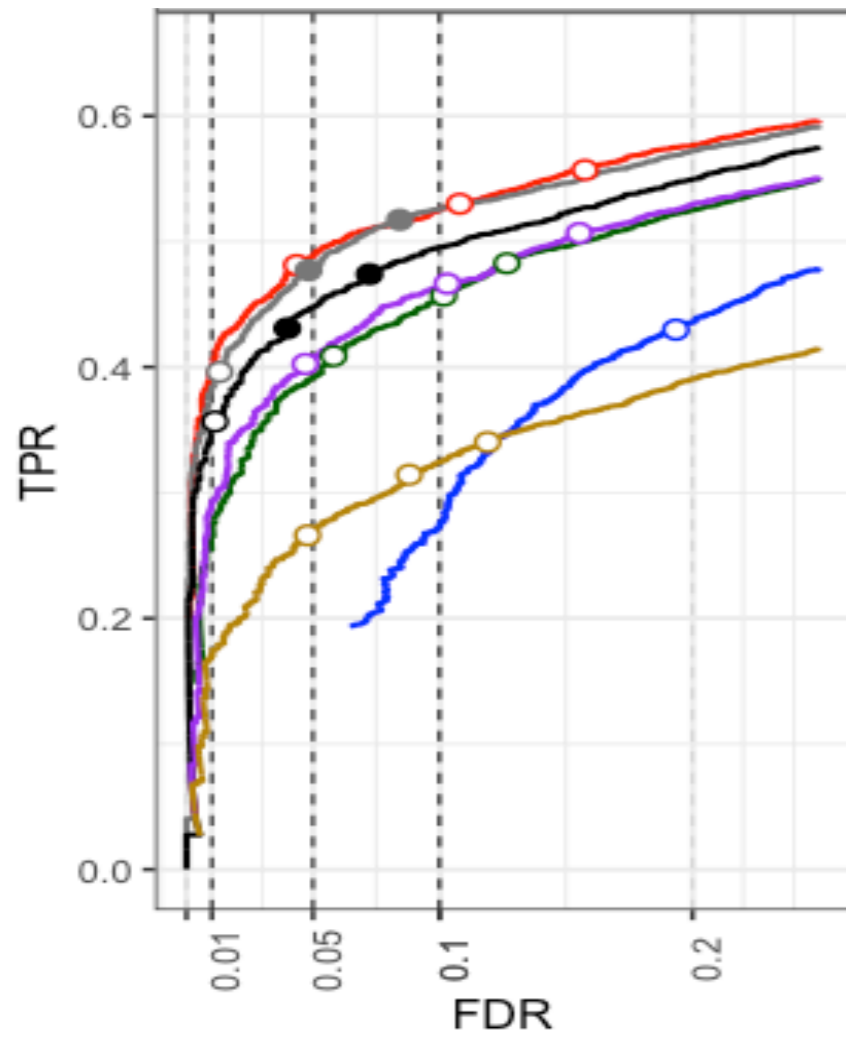# Gtex dataset stringent filtering

# Love dataset stringent filtering

# Other parametric bulk simulations and additional methods



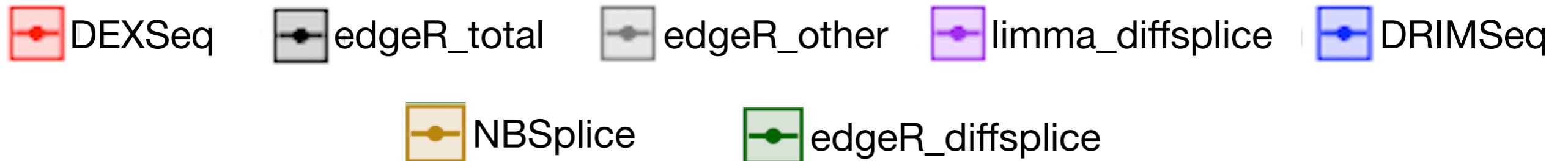DEXSeq · edgeR_total · edgeR_other · limma_diffsplice · DRIMSeq · NBSplice · edgeR_diffsplice

# Results - Scalability

- Methods that require sample-level intercepts scale quadratically with the number of cells
- edgeR one order of magnitude faster than DESeq2
- All methods scale linearly with the number of transcripts