

Estimating cell type composition in whole blood using differentially methylated regions

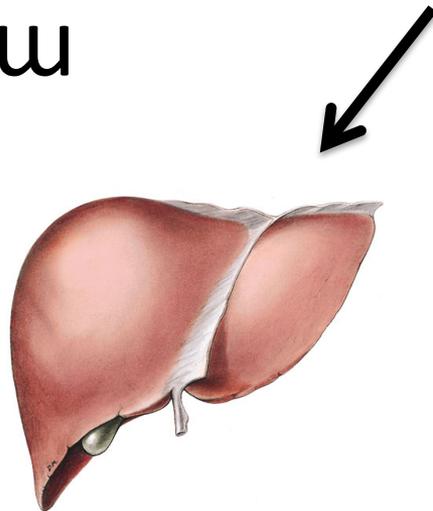
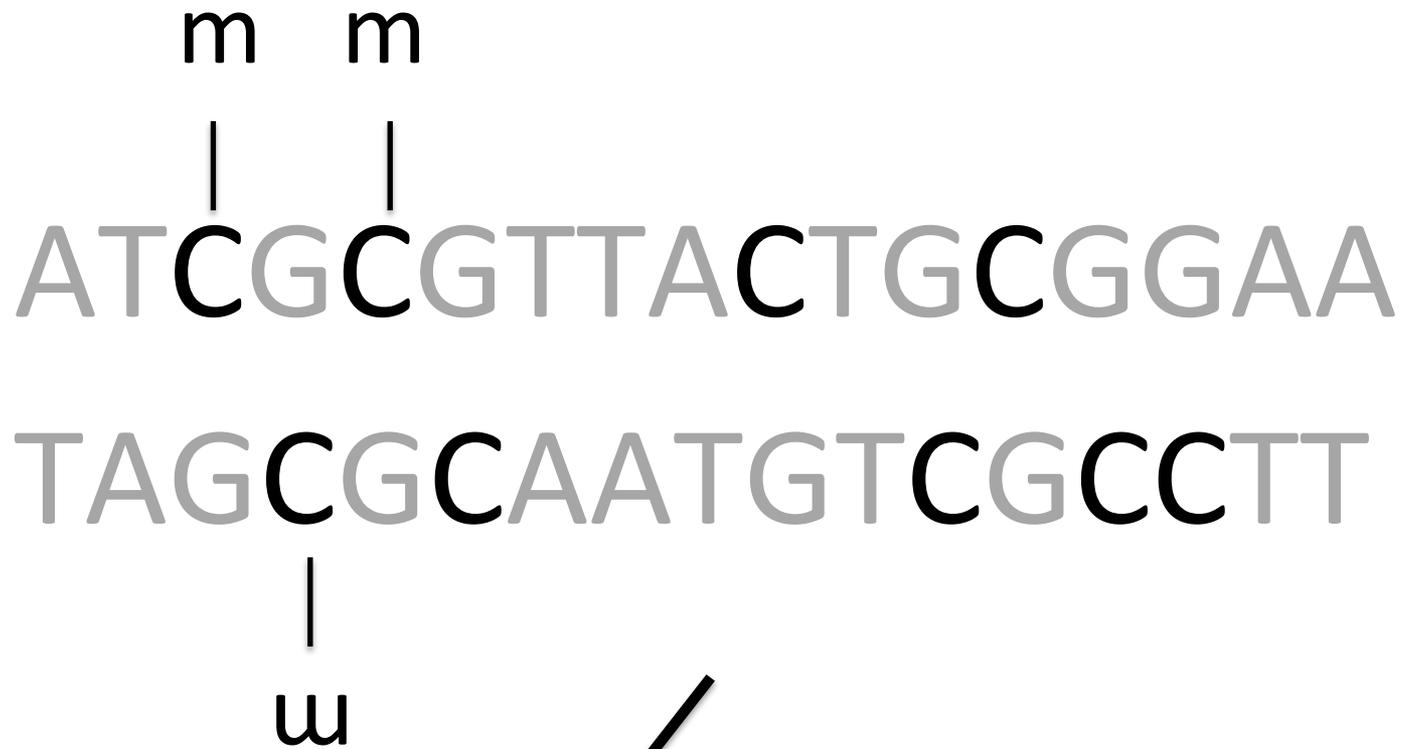
Stephanie Hicks

Bioconductor 2017

What is DNA Methylation?



What is DNA Methylation?



What is DNA Methylation?

ATCGCGTTACTGCGGAA

TAGCGCAATGTCGCCTT

m

|

|

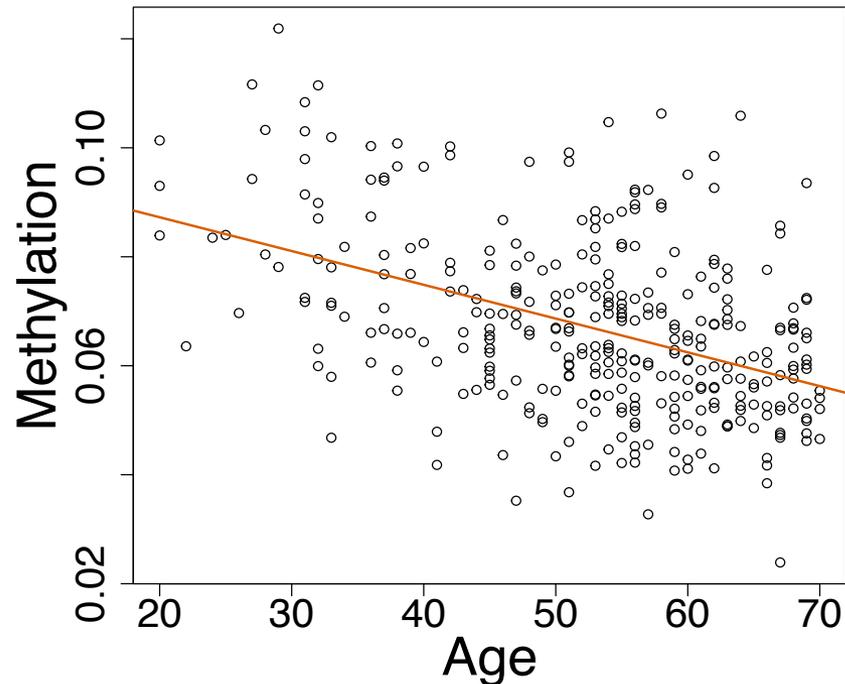
m

|

m

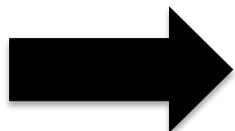
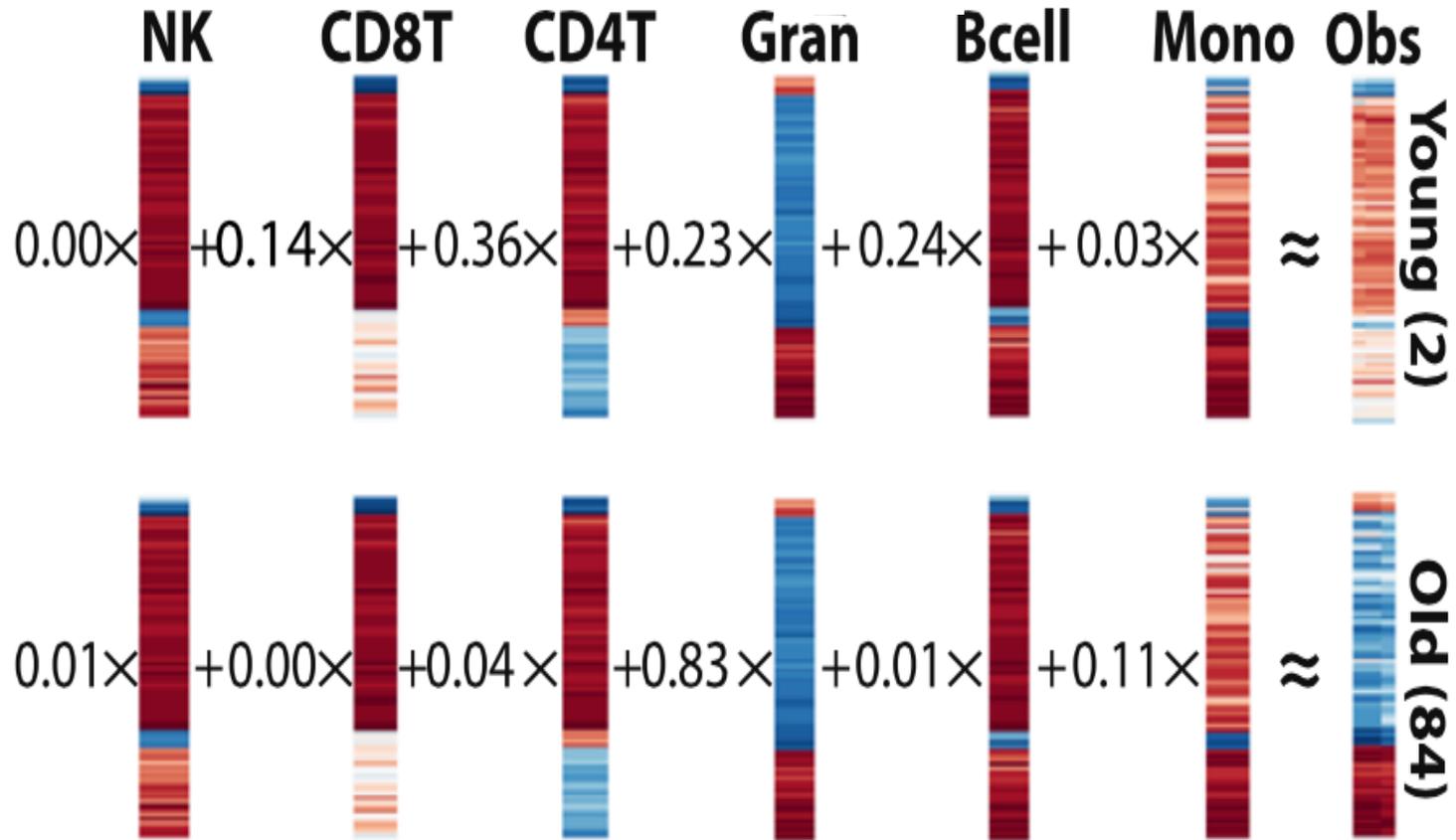


DNA methylation in whole blood correlates with age at this one CpG



Data from GSE32148

Cell composition changes with age



- Different cell compositions in whole blood imply different observed whole blood DNA methylation profiles
- Important to estimate differences in cell composition

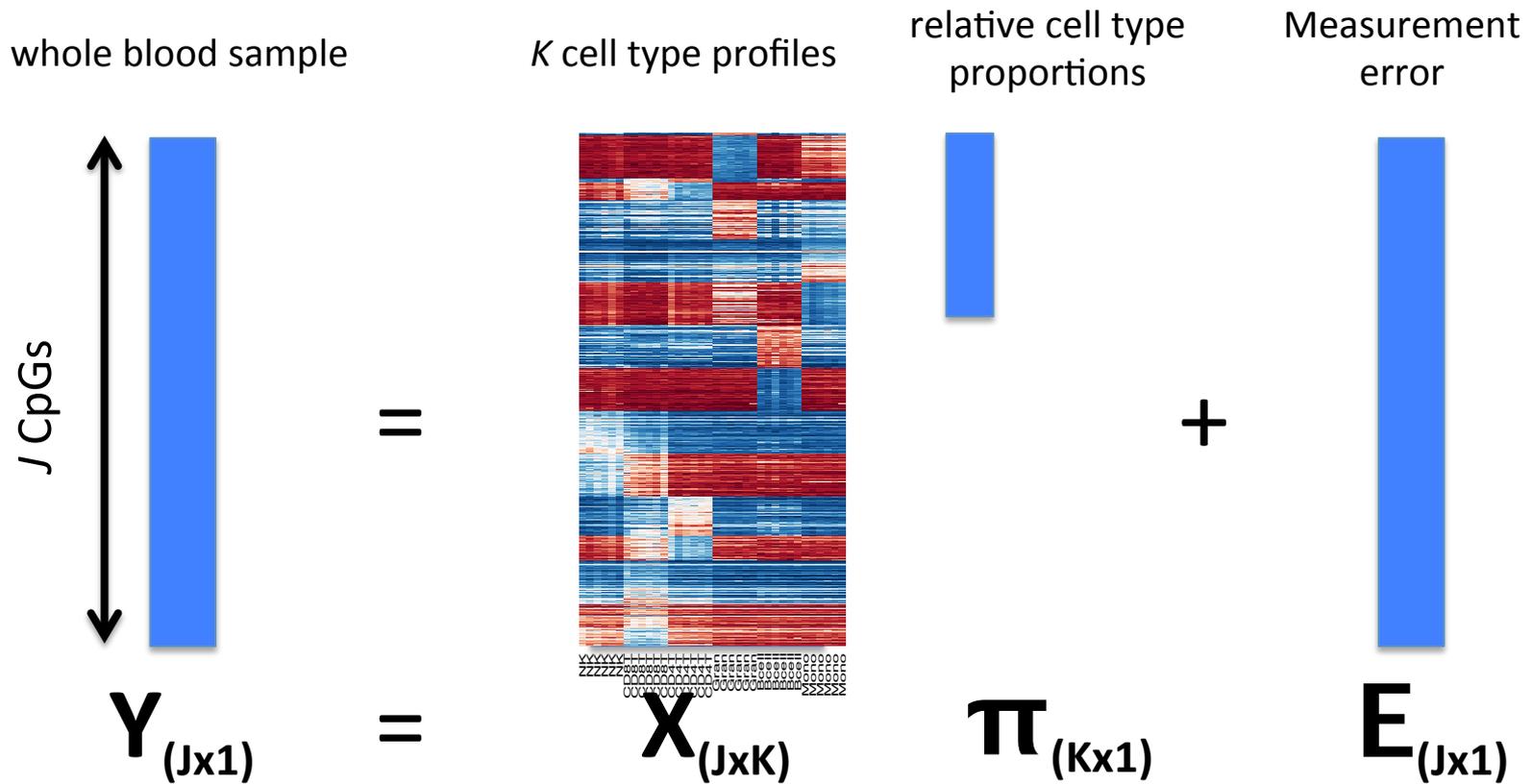
Statistical Model: Houseman et al. (2012)

$$Y_{ij} = \sum_{k=1}^K \pi_{ik} X_{jk} + \epsilon_{ij}$$

$i = (1, \dots, N)$ = whole blood samples

$j = (1, \dots, J)$ = CpGs

$k = (1, \dots, K)$ = cell type profiles



New platform technologies emerging

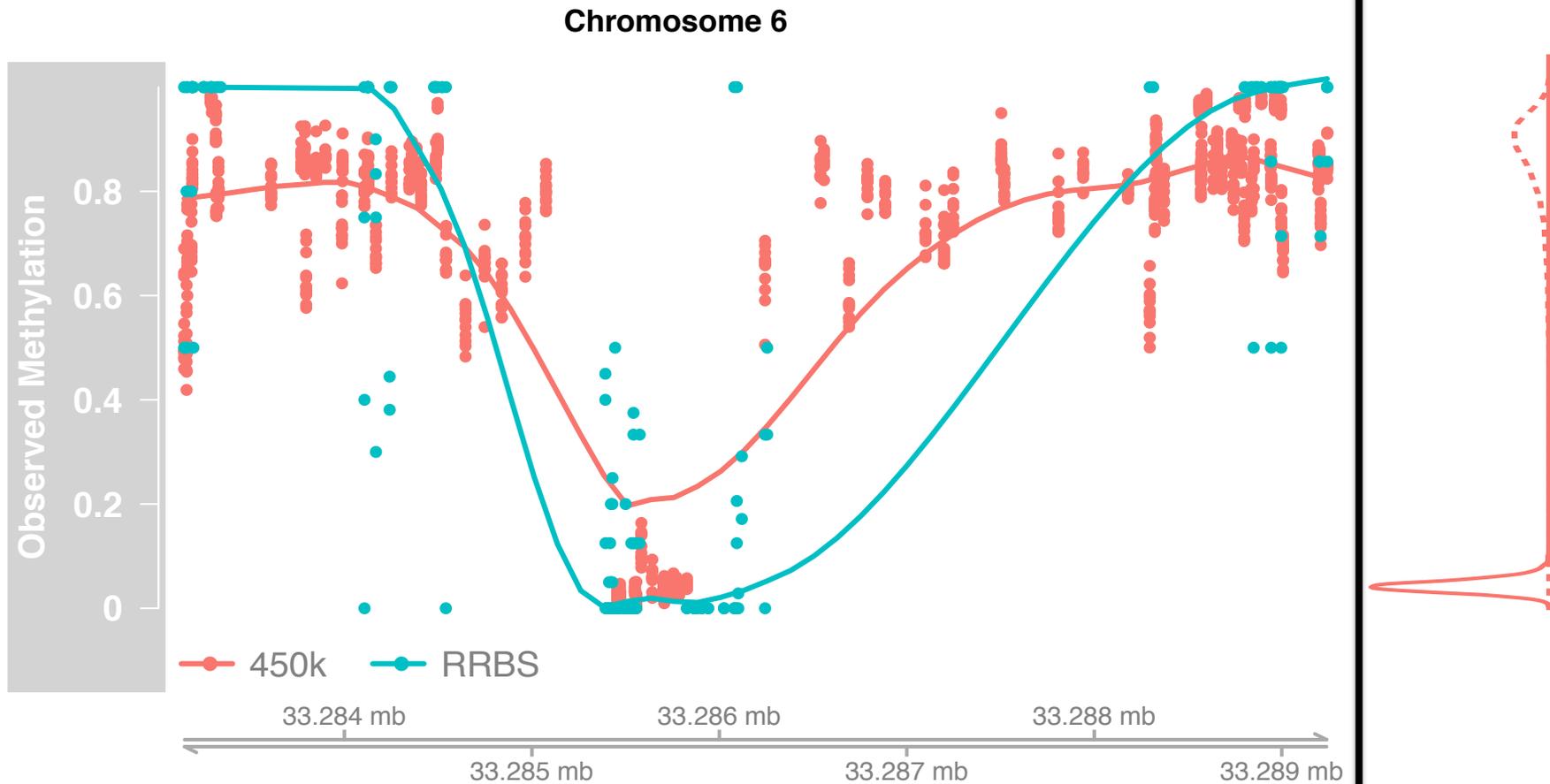
First approach

- Apply Houseman method using new platform technology

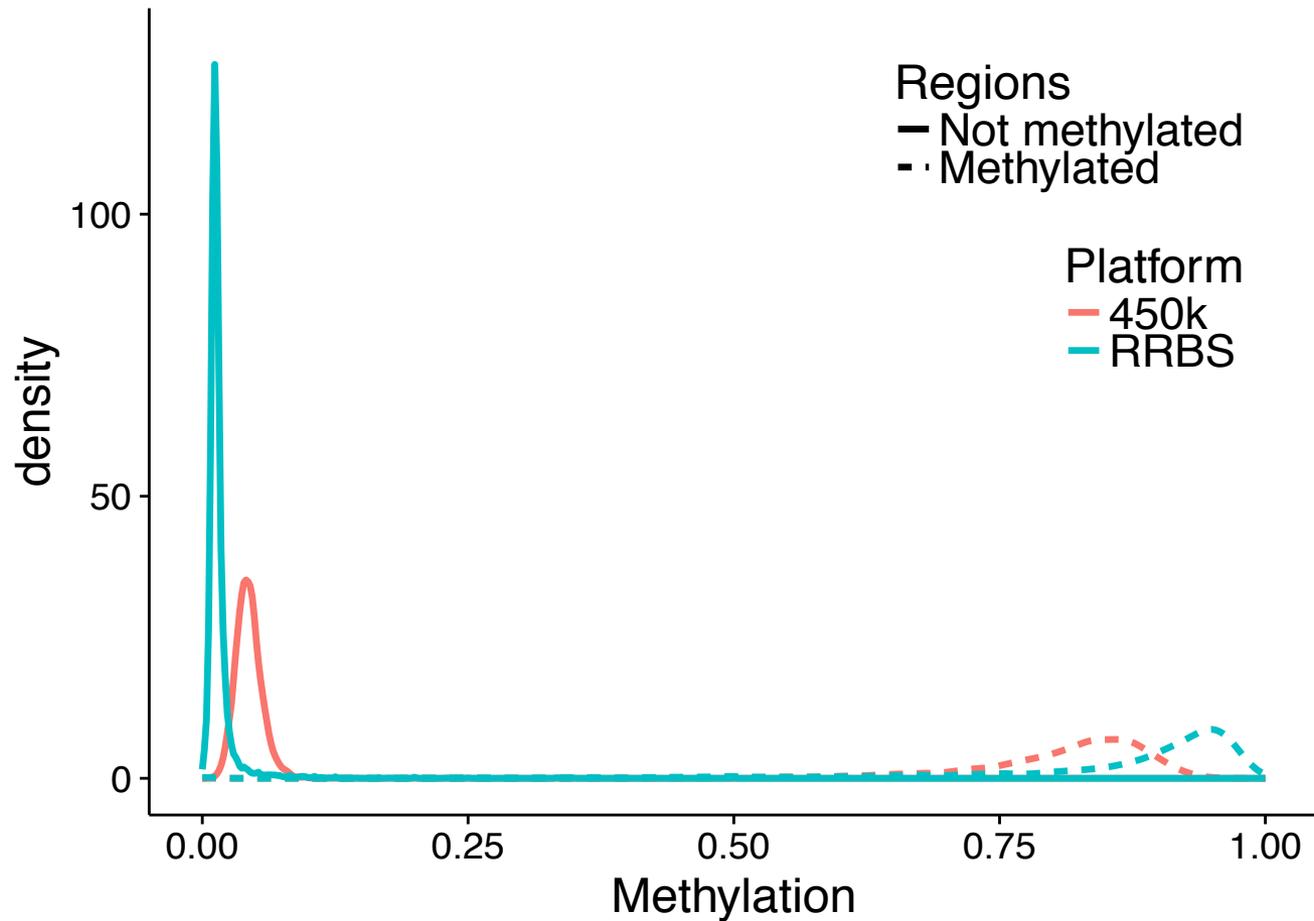
Problems with this approach

1. Observed methylation levels depend on platform used
2. Not all CpGs are included in new platforms

Platform-dependent differences between 450k array and RRBS platforms



Platform-dependent differences between 450k array and RRBS platforms



New platform technologies emerging

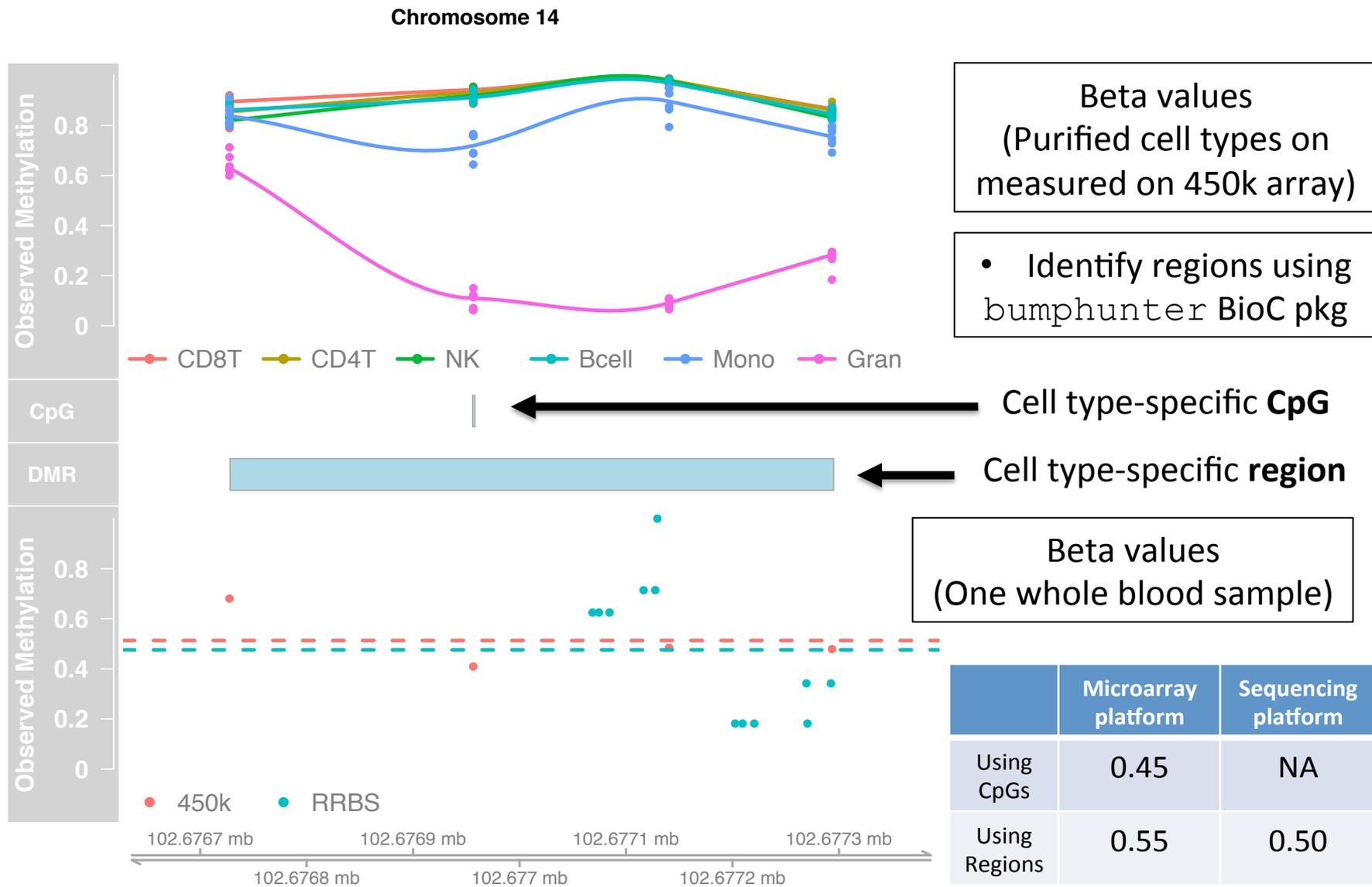
First approach

- Apply Houseman method using new platform technology

Problems with this approach

1. Observed methylation levels depend on platform
2. Not all CpGs are included in new platforms

Cell types preserve their methylation state across regions



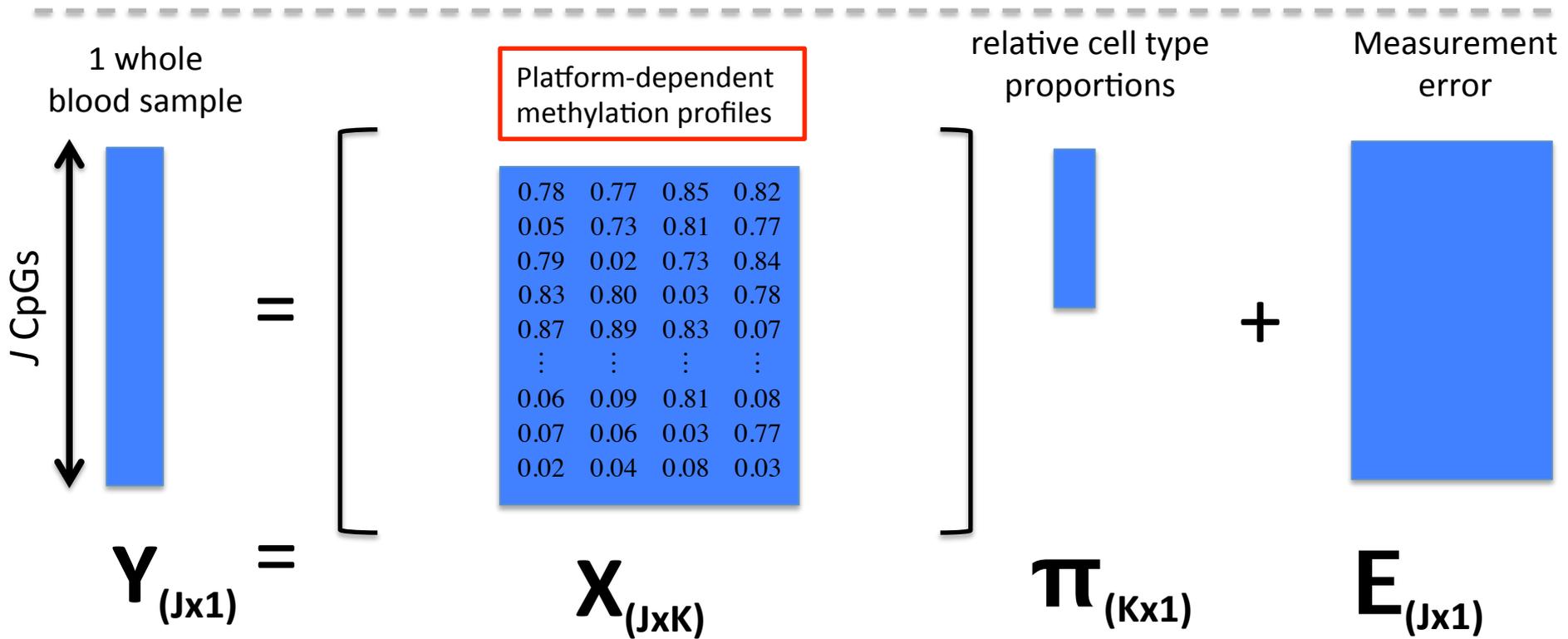
Recall Houseman Model:

$$Y_{ij} = \sum_{k=1}^K \pi_{ik} X_{jk} + \varepsilon_{ij}$$

$i = (1, \dots, N) =$ whole blood samples

$j = (1, \dots, J) =$ CpGs

$k = (1, \dots, K) =$ cell type profiles



Our proposed model:

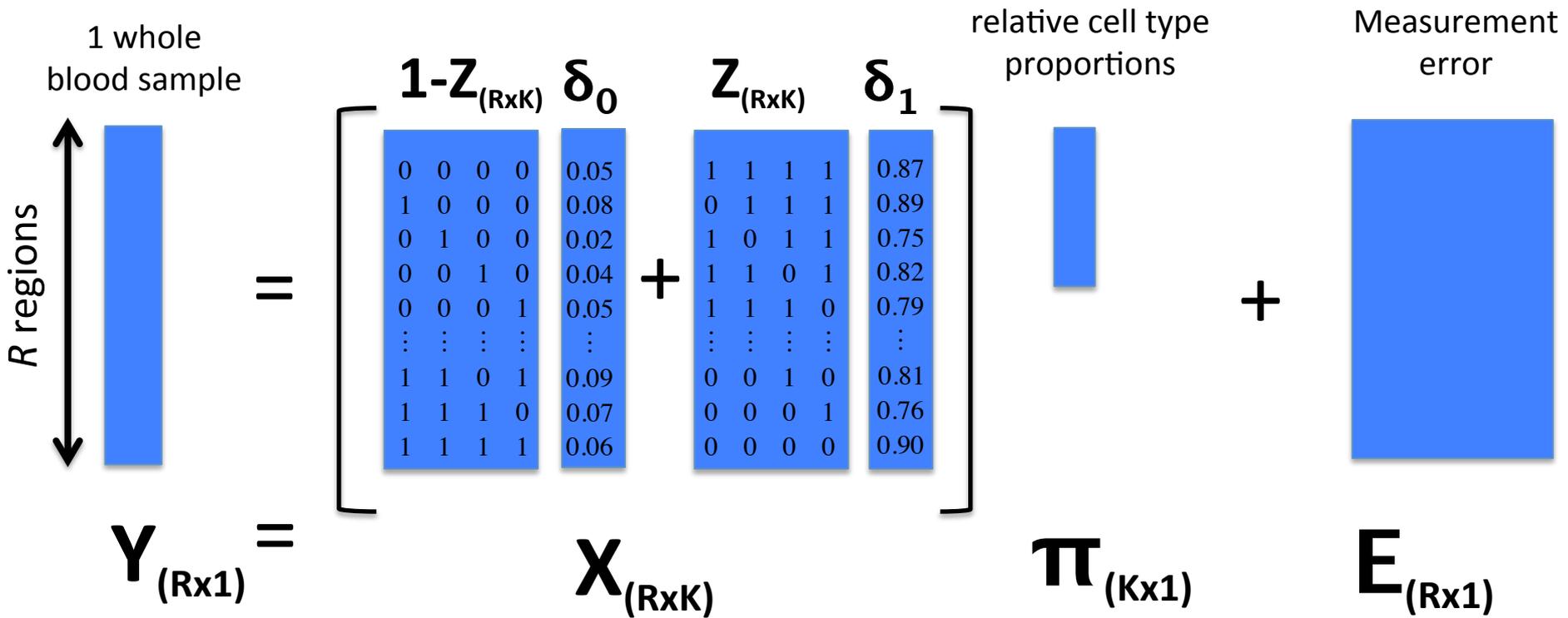
$$Y_r = \sum_{k=1}^K \pi_k \left[(1 - Z_{rk}) \delta_{0,r} + Z_{rk} \delta_{1,r} \right] + \varepsilon_r$$

$$\delta_{0,r} \sim N(\alpha_0, \sigma_0^2)$$

$$\delta_{1,r} \sim N(\alpha_1, \sigma_1^2)$$

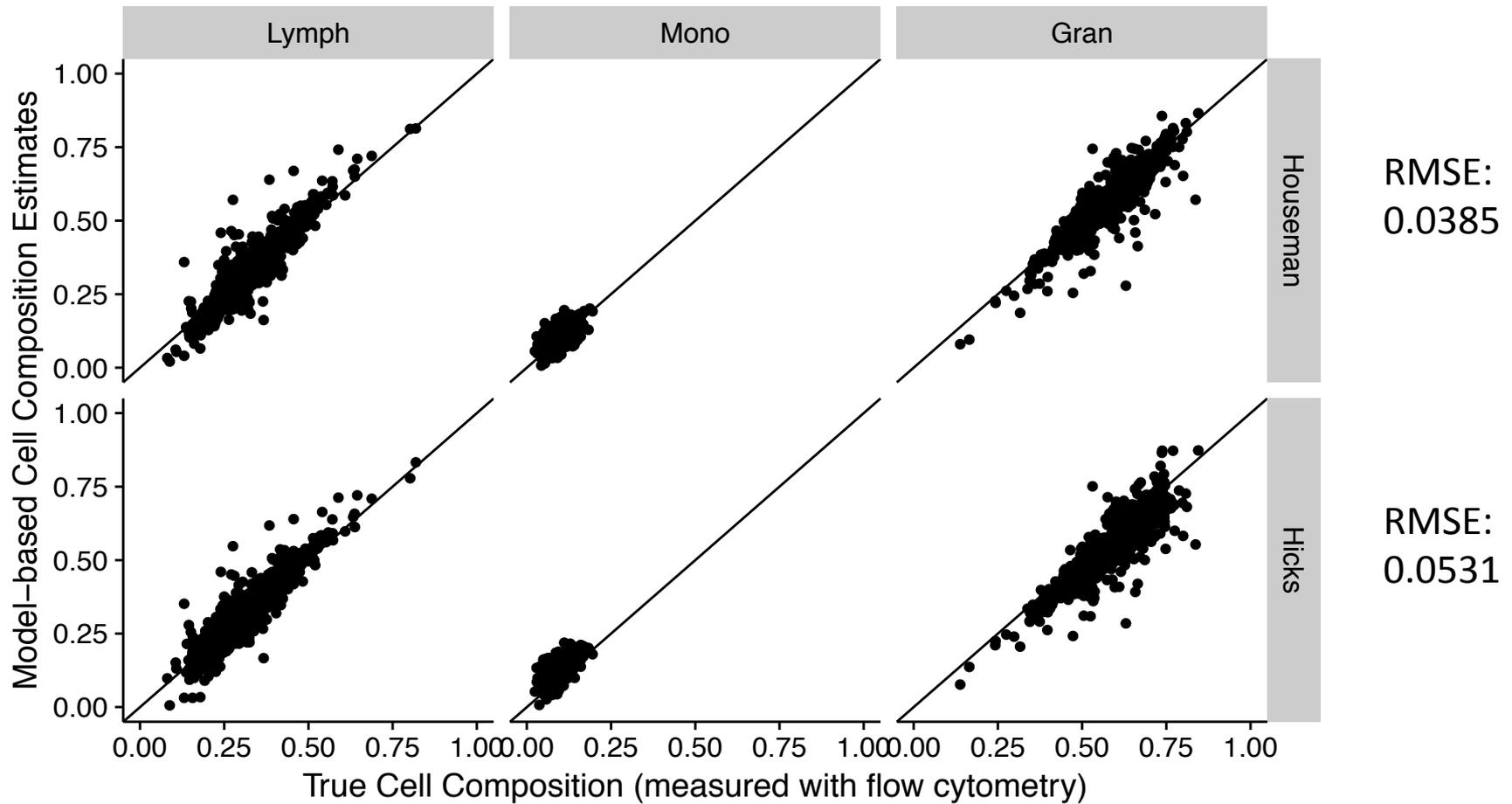
$$\varepsilon_r \sim N(0, \sigma^2)$$

$Z_{rk} = \begin{cases} 1 & \text{if region } r \text{ and cell type } k \text{ is methylated} \\ 0 & \text{otherwise} \end{cases}$
 $r = (1, \dots, R) = \text{differentially methylated regions}$
 $k = (1, \dots, K) = \text{cell types}$



How does our model perform?

$N = 800$ whole blood samples run on 450k microarray platform



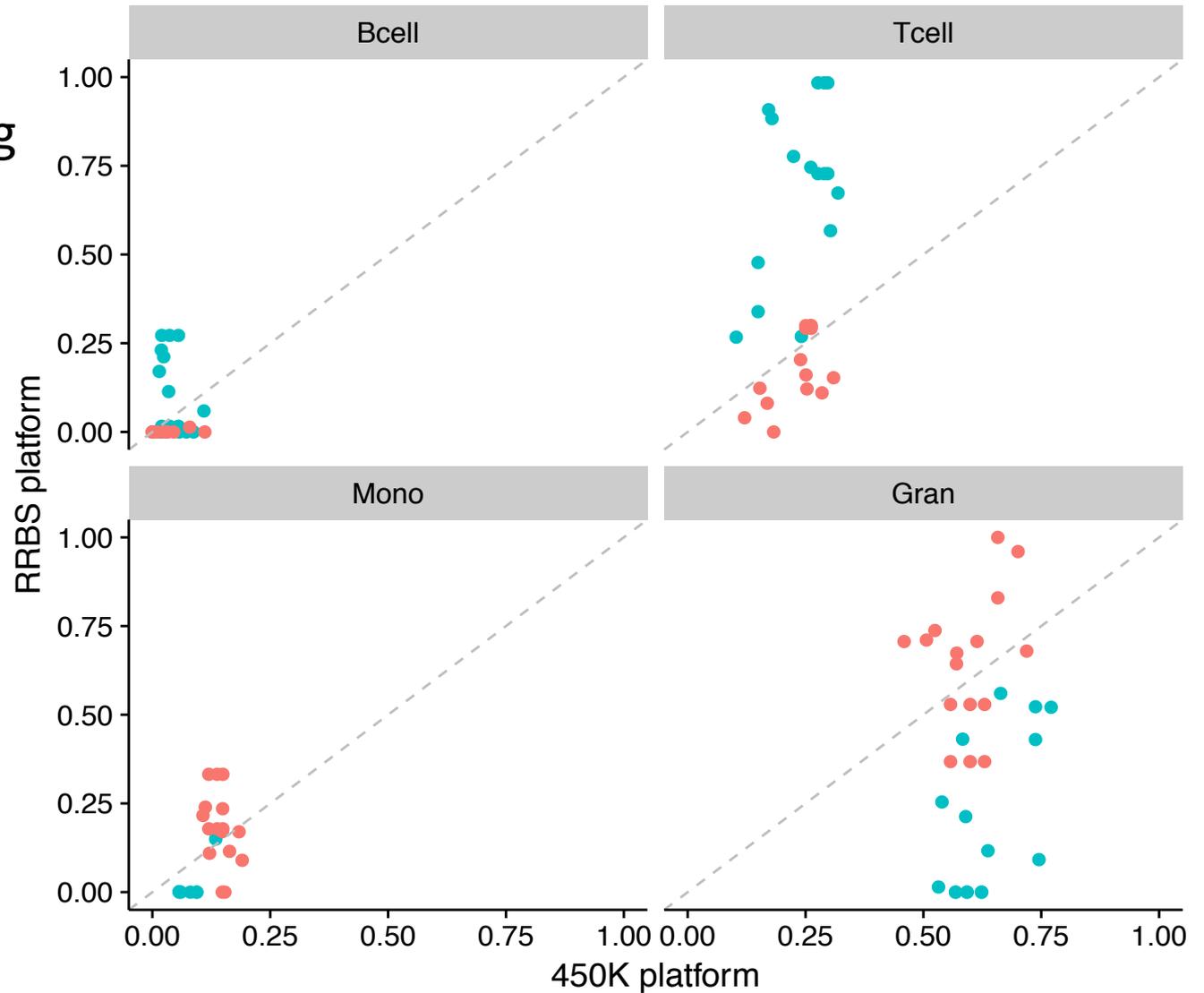
Cell composition estimates from whole blood samples measured on two platforms

$N = 12$ samples measured

on two platforms:

- 450k microarray
- RRBS sequencing

Method • Our method • Houseman



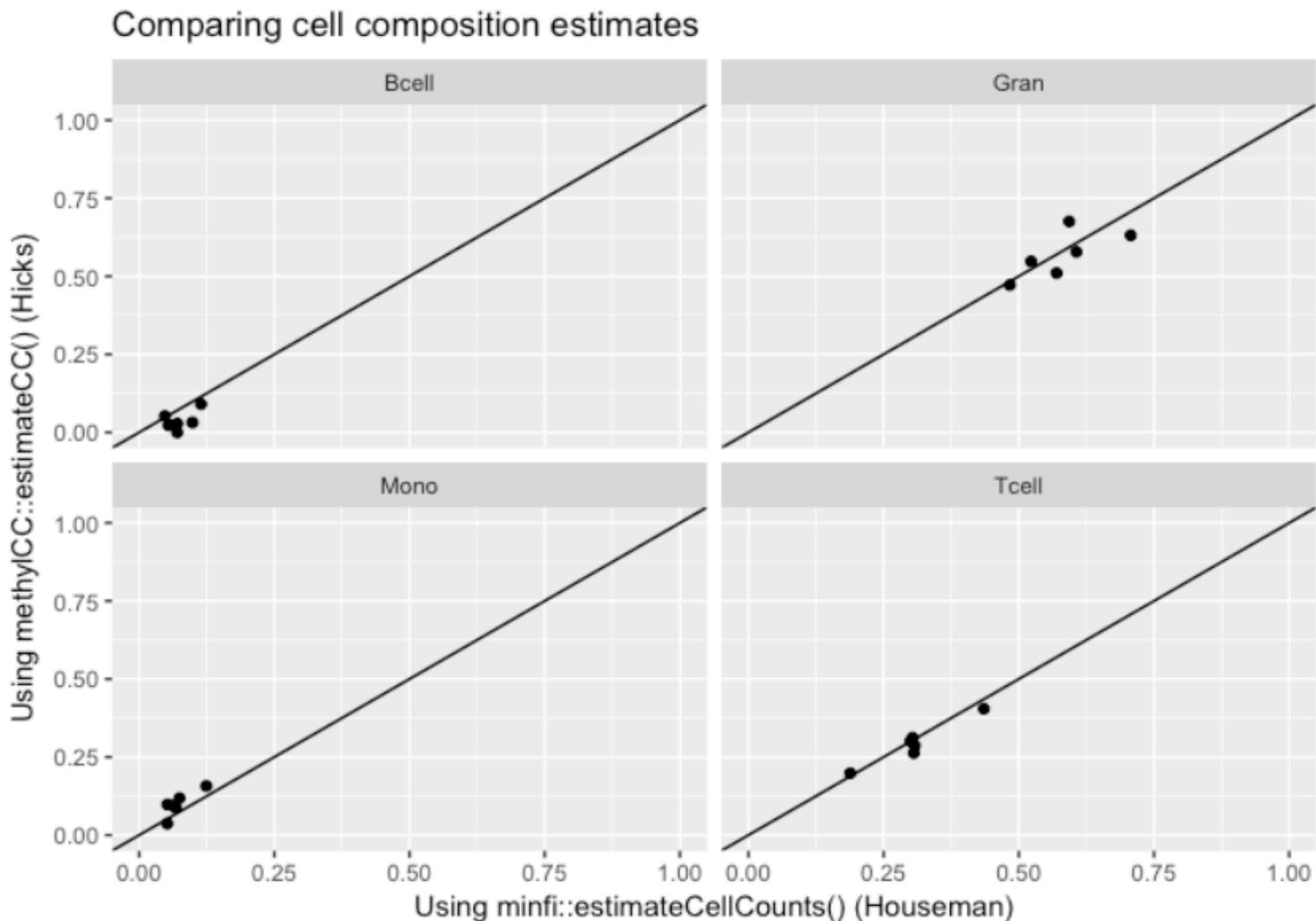
3.1 Demo for BioC 2017

```
library(methylCC)
library(minfi)
library(FlowSorted.Blood.450k)
data(FlowSorted.Blood.450k)
```

```
# Subset RGCharacteristics
rgset <- FlowSorted.Blood.450k
```

```
# Use methylCC:
est.methylCC <- methylCC(rgset)
counts.methylCC <- estimateCellCounts(methylCC)
```

```
# Compare to minfi:
sampleNames(rgset)
counts.minfi <- estimateCellCounts(minfi)
```



For more information

methylCC:

<https://github.com/stephaniehicks/methylCC>

Comments/Suggestions:

email: shicks@jimmy.harvard.edu

GitHub & Twitter: @stephaniehicks

#BioC2017

#RLadies

#dataparasite

