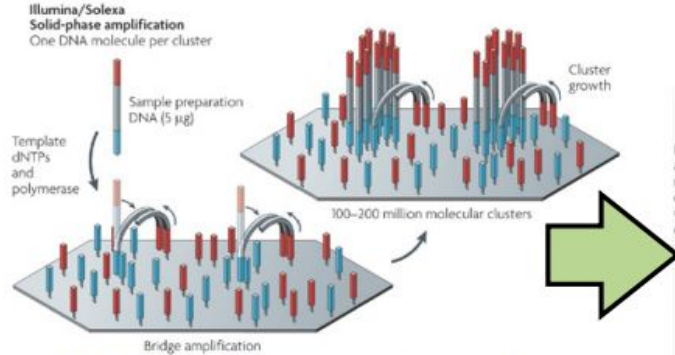


Pre-processing and normalization

Jeff Leek

@jtleek

www.jtleek.com

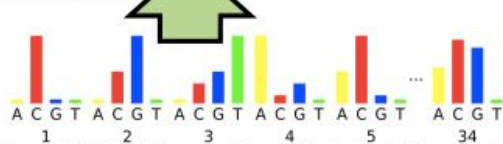


Source: Metzker ML. Sequencing technologies - the next generation. Nat Rev Genet. 2010

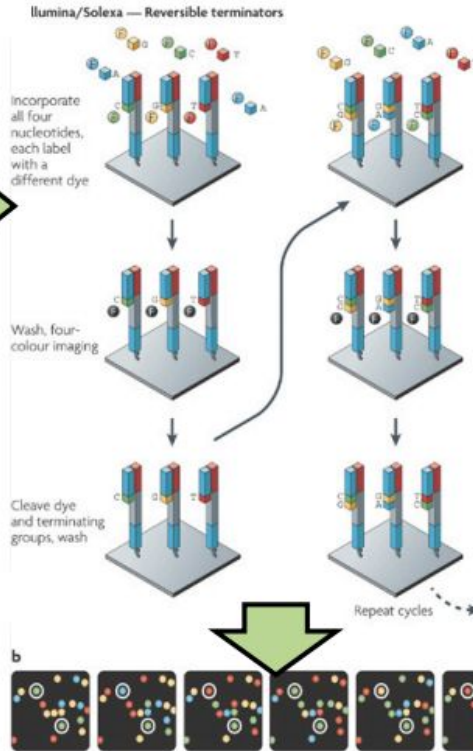


name
sequence
quality scores

x 100s of millions



Source: Whiteford et al. Swift: primary data analysis for the Illumina Solexa sequencing platform. Bioinformatics. 2009

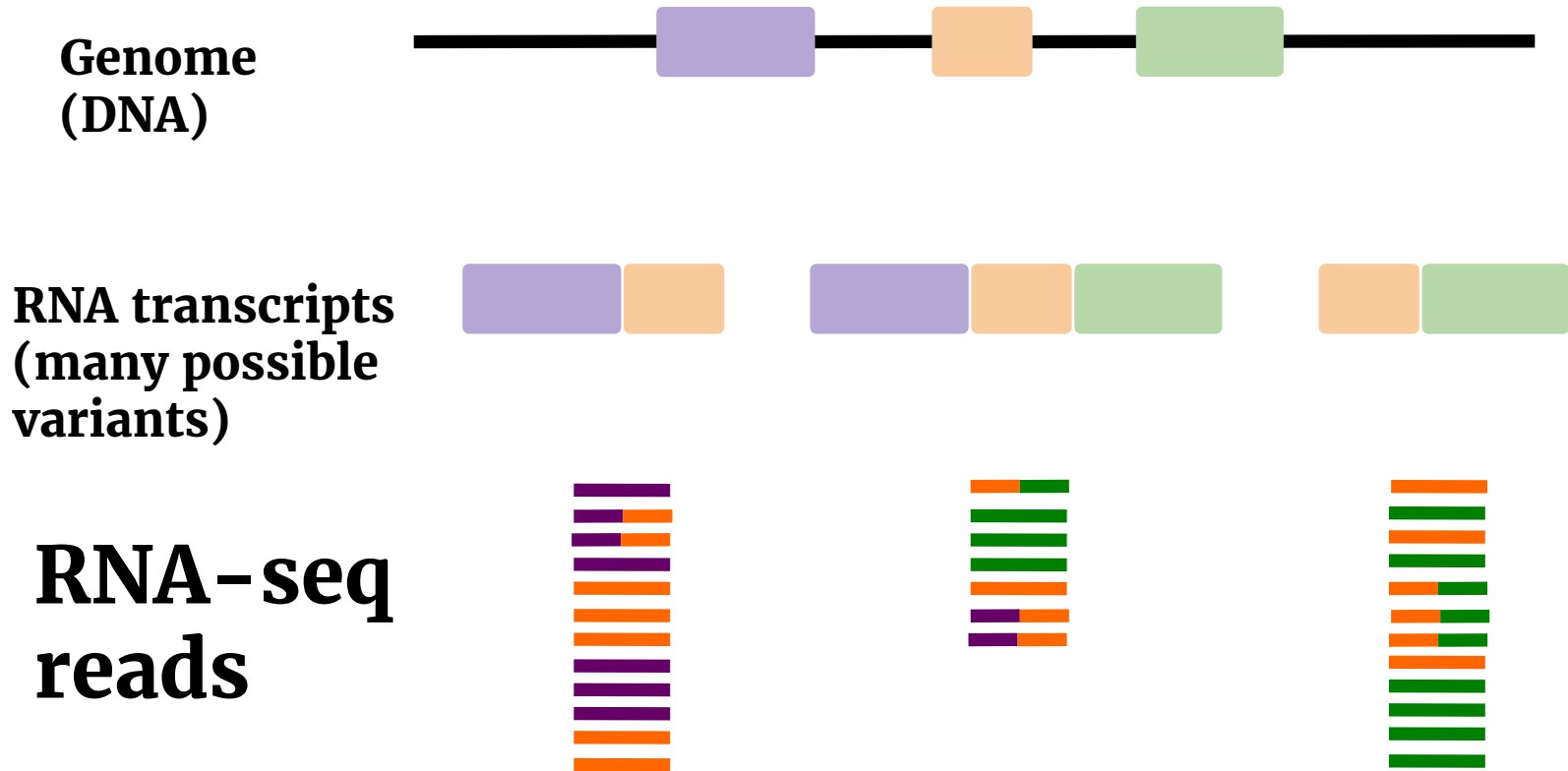


Source: Metzker ML. Sequencing technologies - the next generation. Nat Rev Genet. 2010

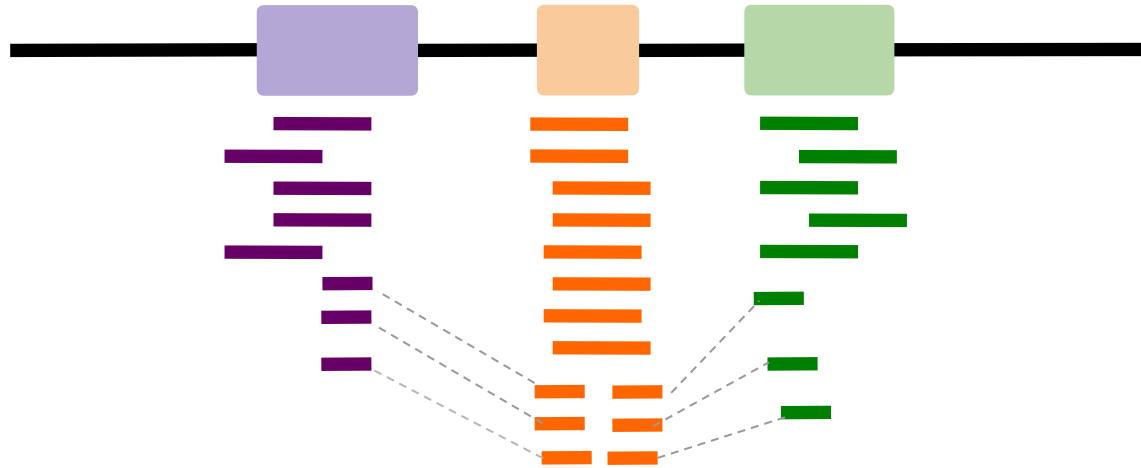
Preprocessing

Convert raw data to “processed”

Try to remove technological
artifacts

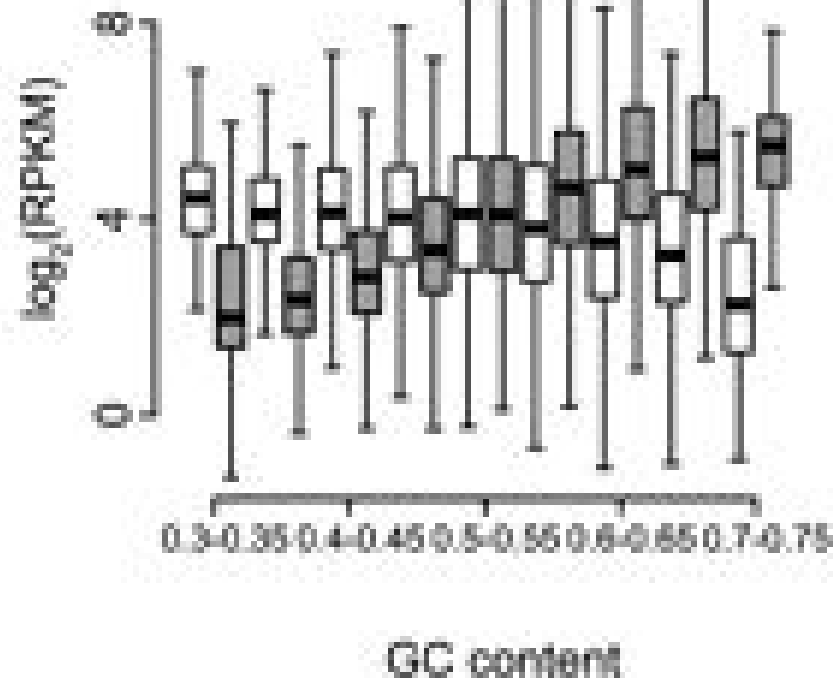


**Genome
(DNA)**

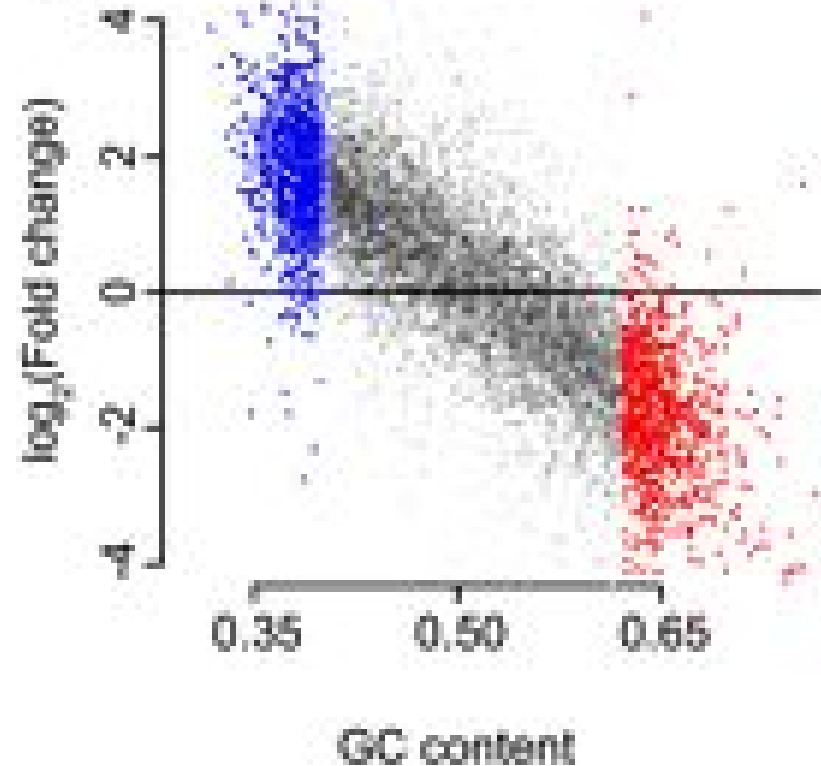


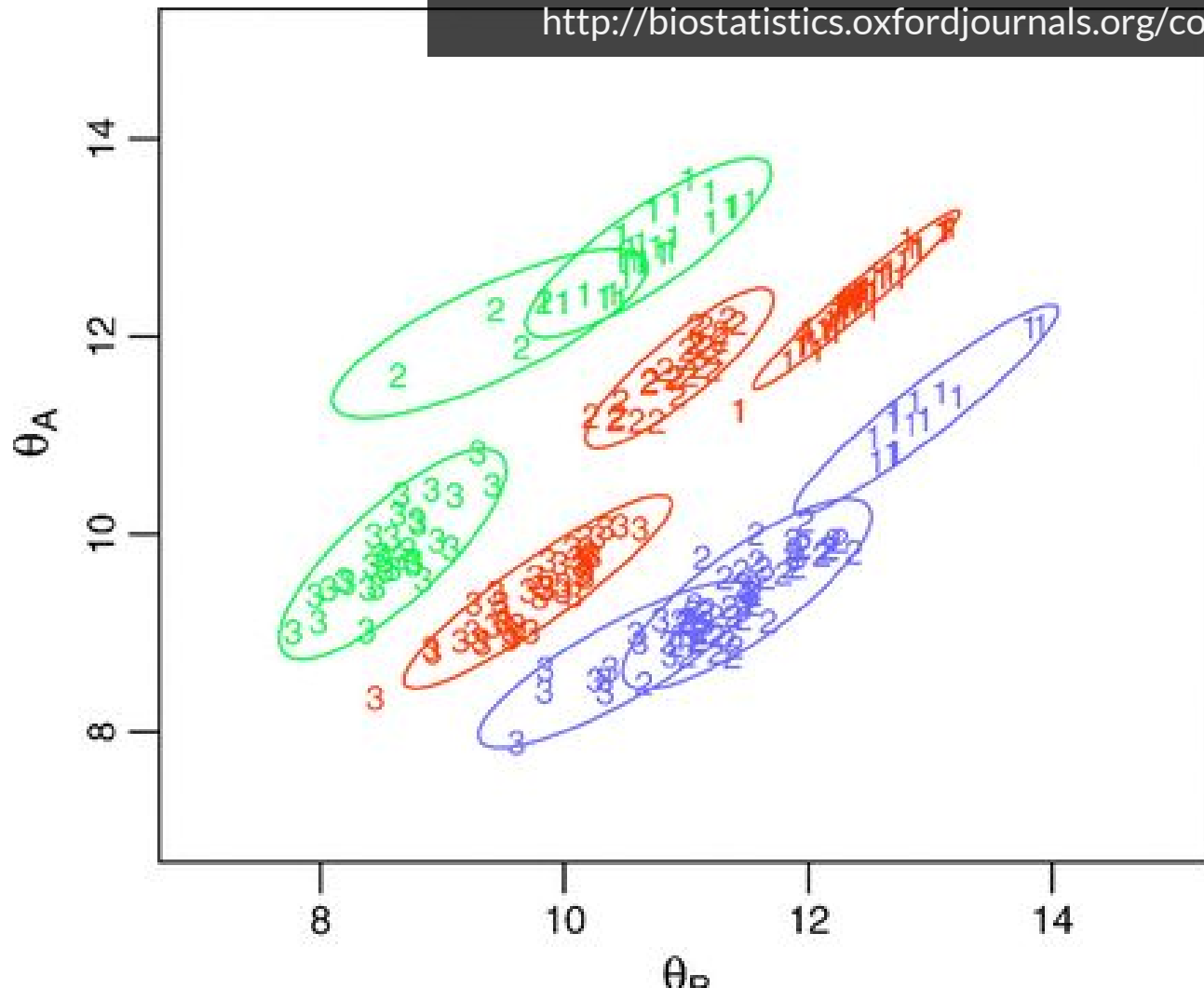
expression = 24

(b)



(c)





Normalization

Remove technological biases

Make samples comparable

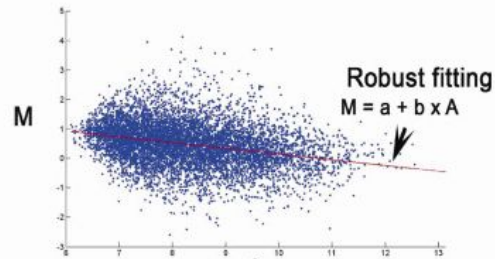
Sample 1

Peak Coordinates
Read Coordinates

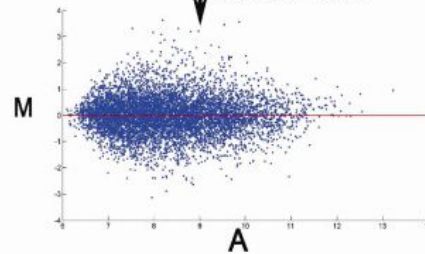
Sample 2

Peak Coordinates
Read Coordinates

MA plot of common peaks



A
Normalization



Quantile normalization

Most common technique

Bulk distributions exactly the same

Raw data

2	4	4	5
5	14	4	7
4	8	6	9
3	8	5	8
3	9	3	5

**Order values
within each sample
(or column)**

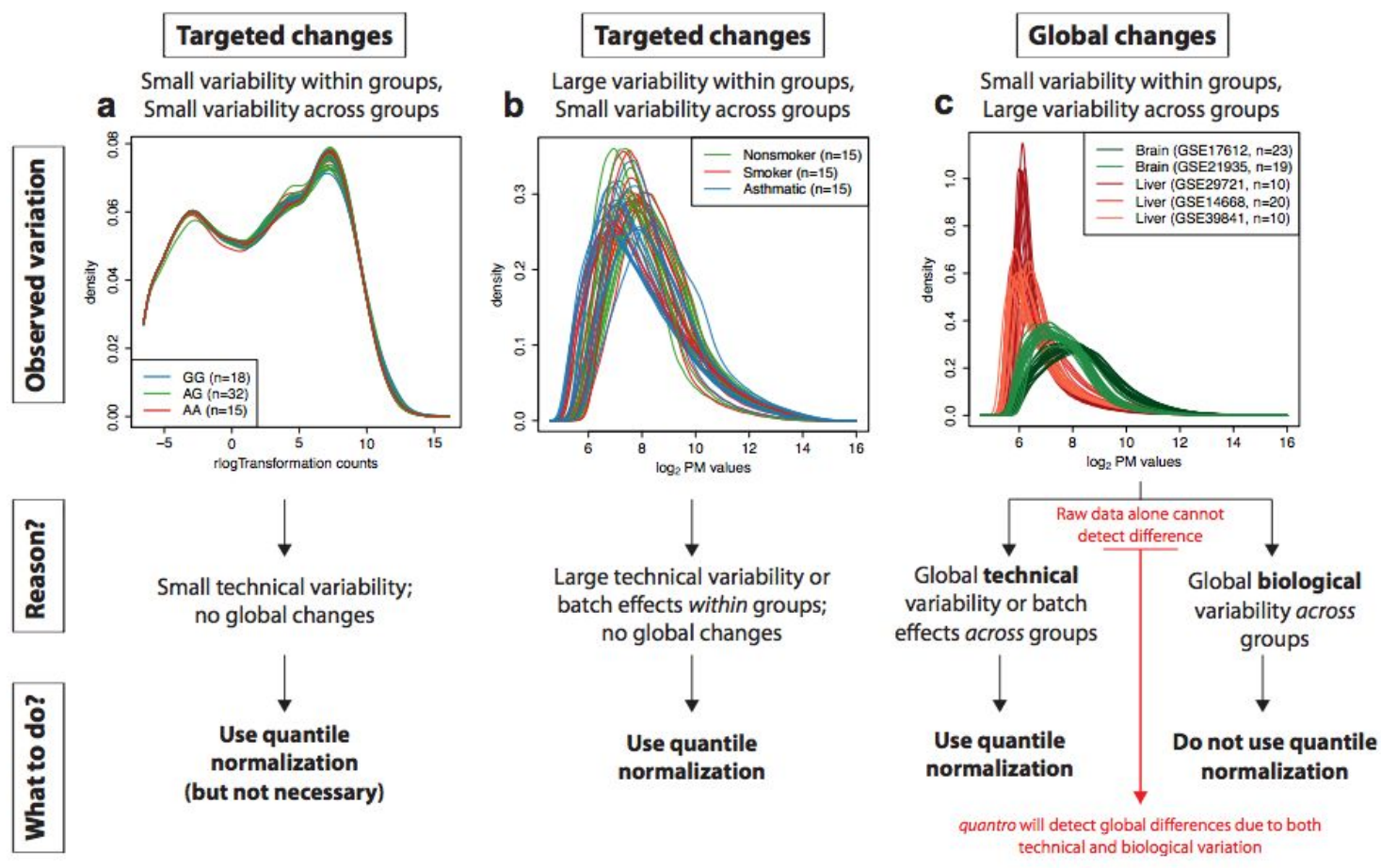
2	4	3	5
3	8	4	5
3	8	4	7
4	9	5	8
5	14	6	9

**Average across rows
and substitute value
with average**

3.5	3.5	3.5	3.5
5.0	5.0	5.0	5.0
5.5	5.5	5.5	5.5
6.5	6.5	6.5	6.5
8.5	8.5	8.5	8.5

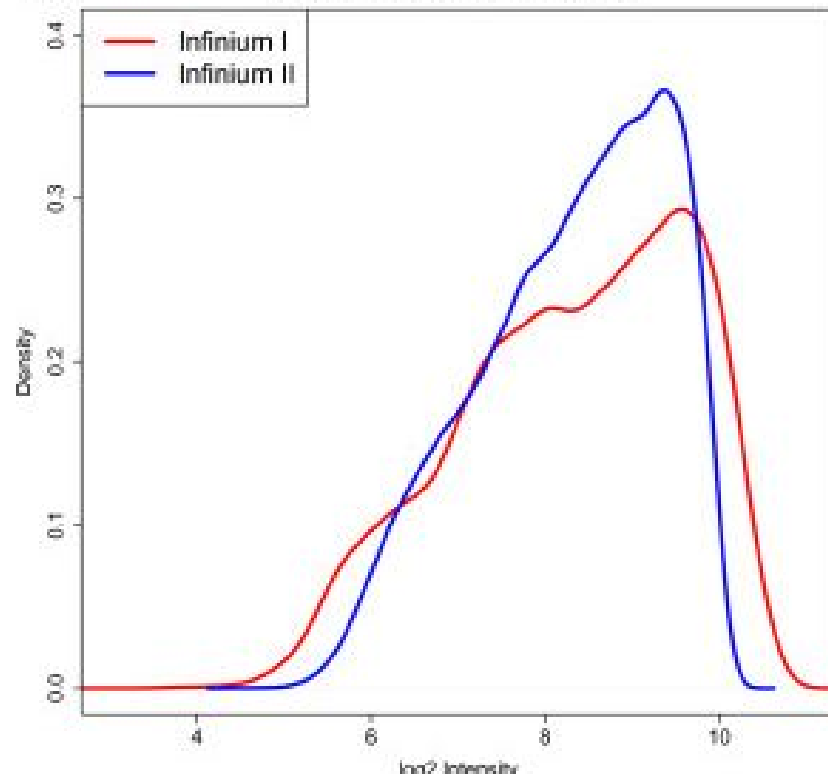
**Re-order averaged
values in original
order**

3.5	3.5	5.0	5.0
8.5	8.5	5.5	5.5
6.5	5.0	8.5	8.5
5.0	5.5	6.5	6.5
5.5	6.5	3.5	3.5

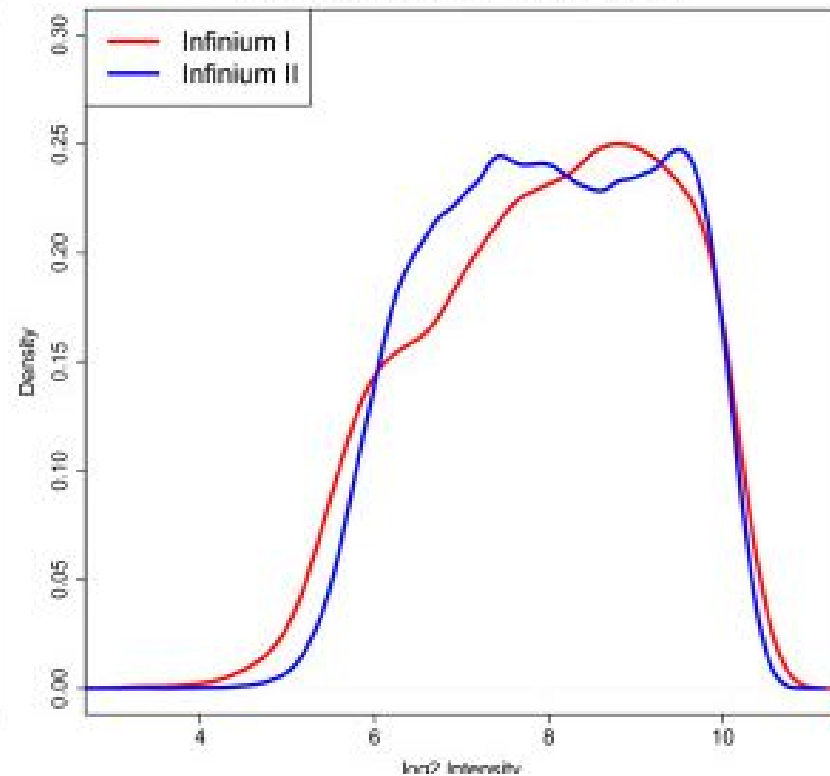


(a)

Methylated Channel



Unmethylated Channel



Notes and further reading

- Preprocessing and normalization are highly platform/problem dependent
- In general check to make sure there aren't bulk differences between samples, especially due to technology
- Bioconductor workflows are a good place to start:

<https://www.bioconductor.org/help/workflows/>

“ First, researchers starting out in genomics must keep in mind that interesting outliers — that is, results that deviate significantly from the sample — will inevitably contain a plethora of experimental or analytical artefacts. ”