# RNA-Seq Data Analysis

## Marta R. Hidalgo

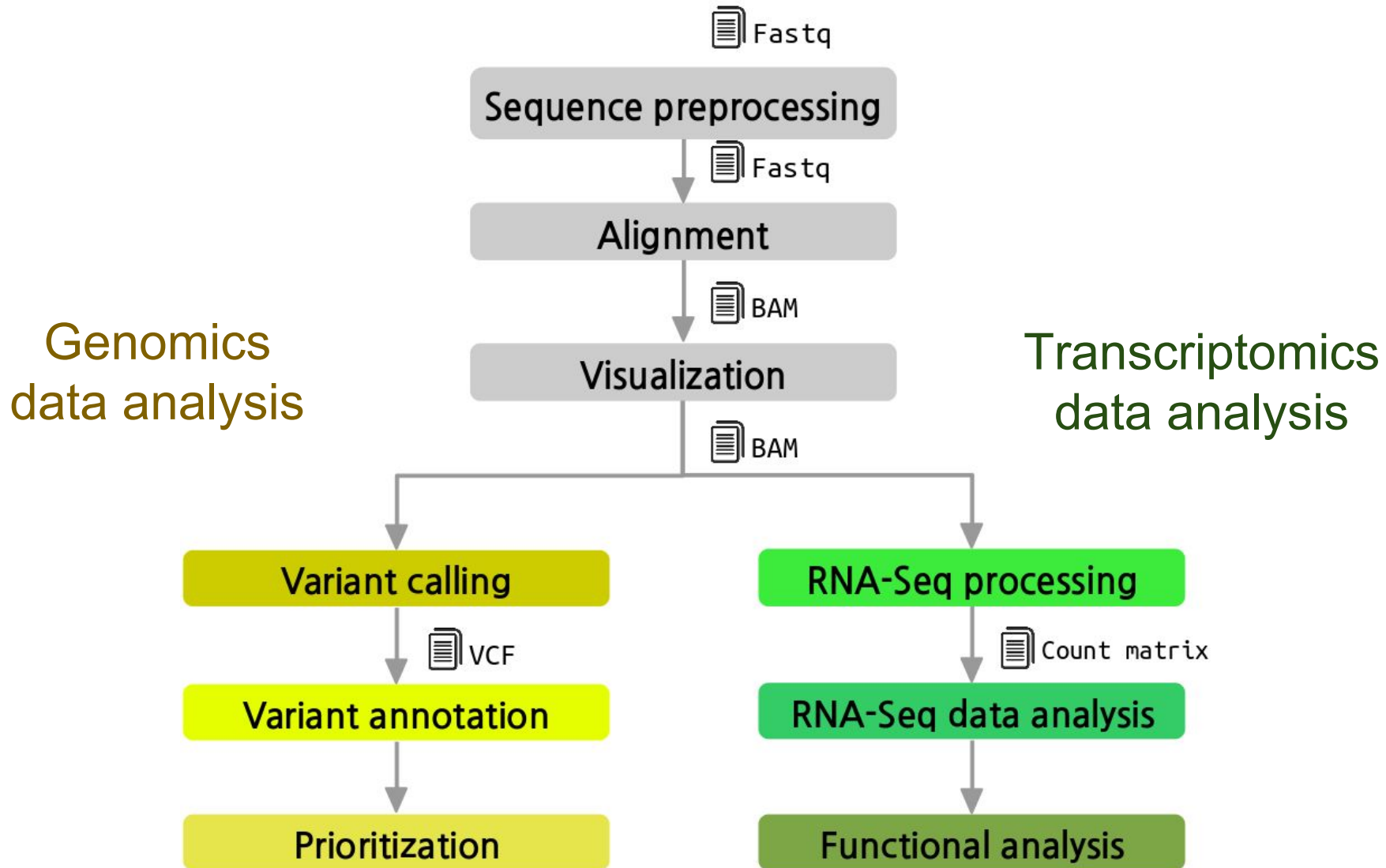Máster en Biotecnología Biomédica
Universidad Politécnica de Valencia

Unidad de
Bioinformática y
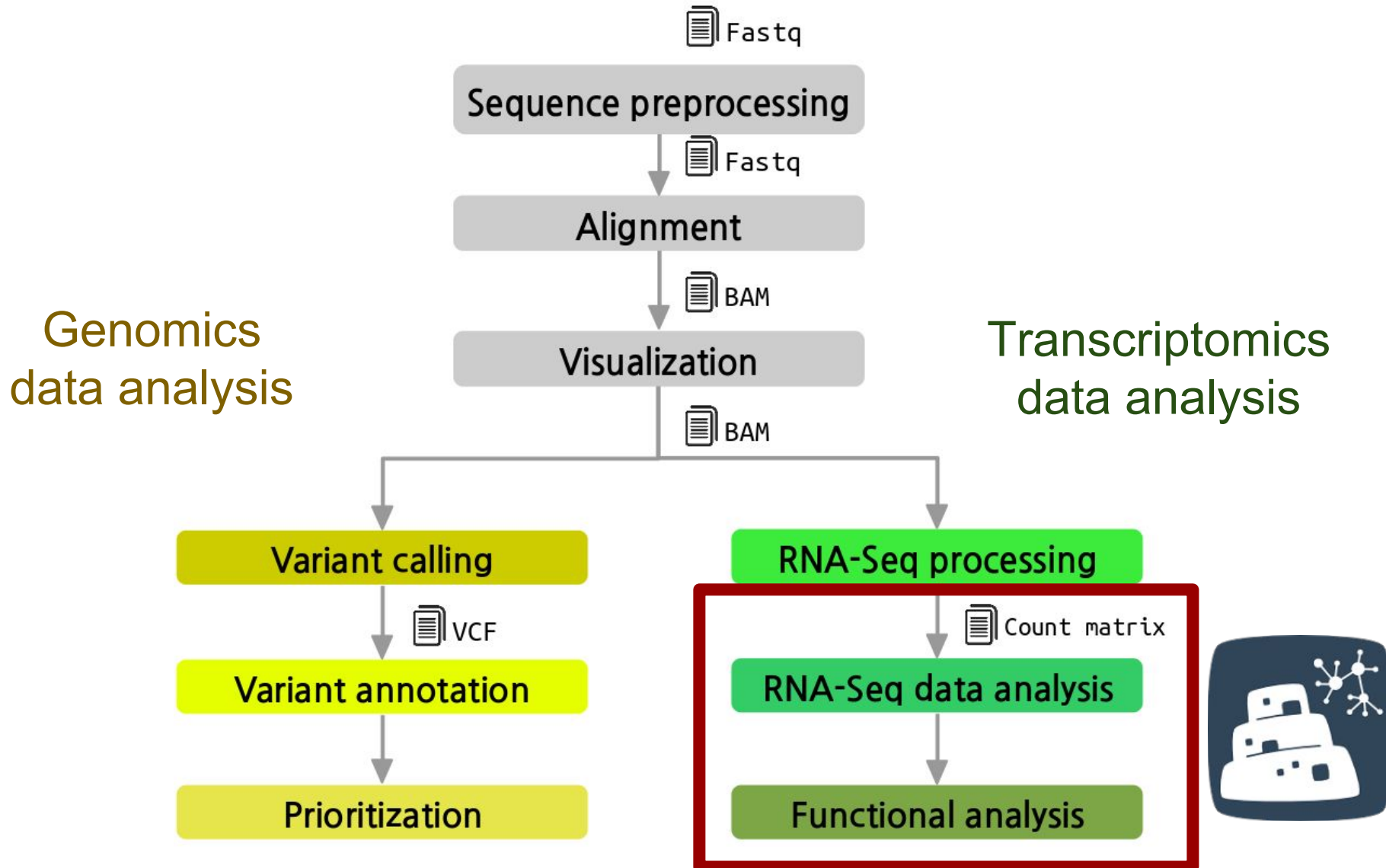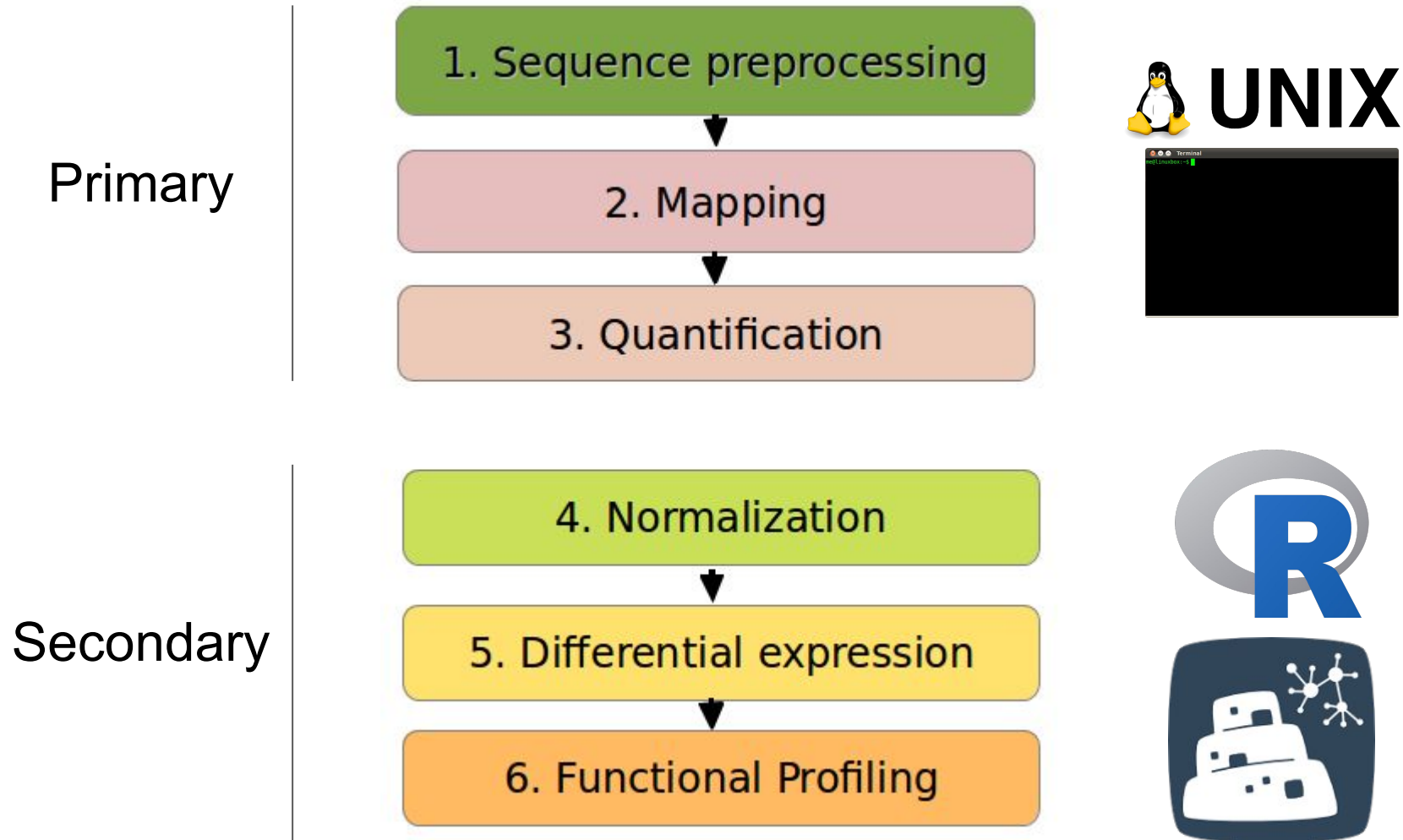Bioestadística

PRINCIPE FELIPE
CENTRO DE INVESTIGACION

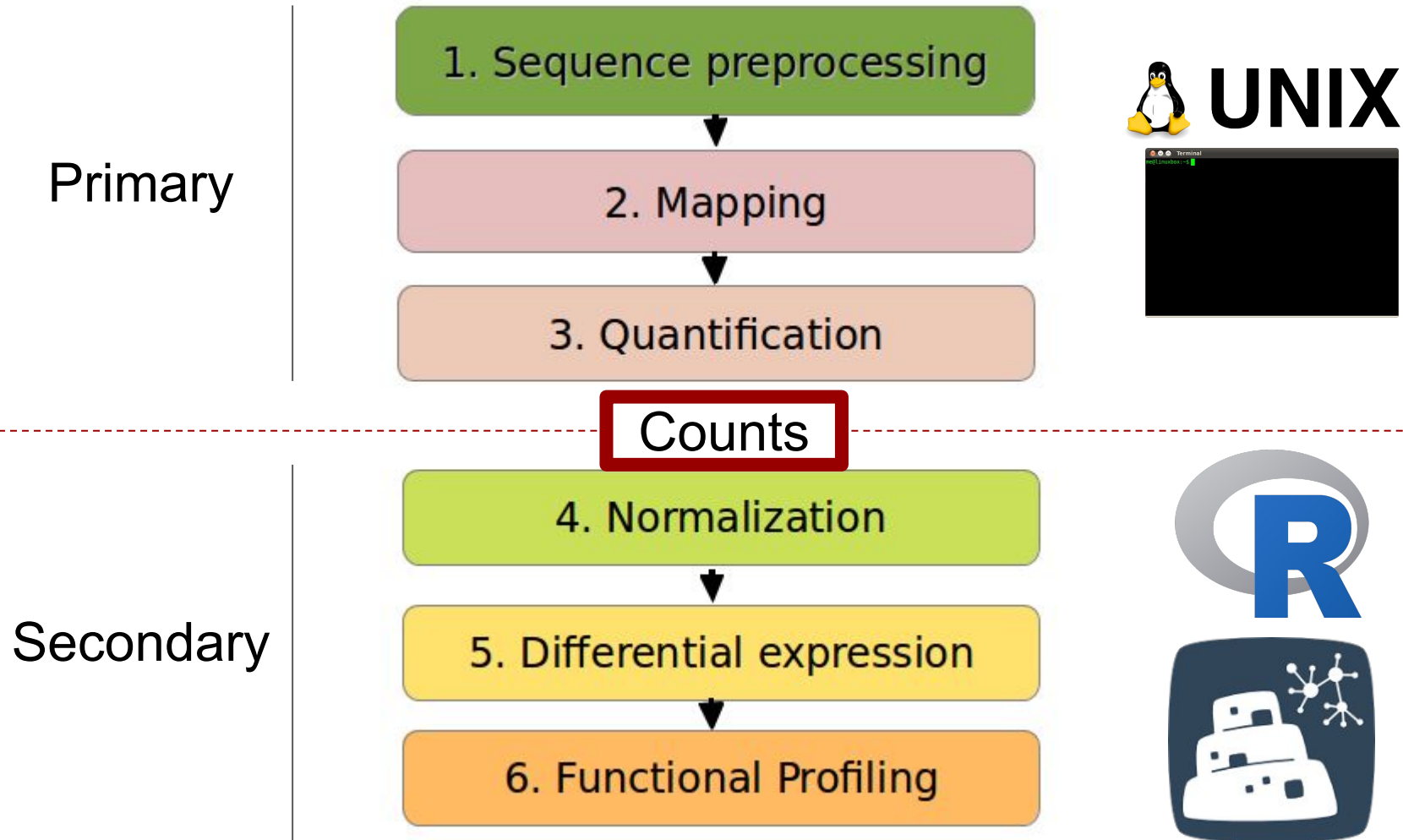# NGS data analysis pipeline

# NGS data analysis pipeline

# RNA-seq data analysis pipeline
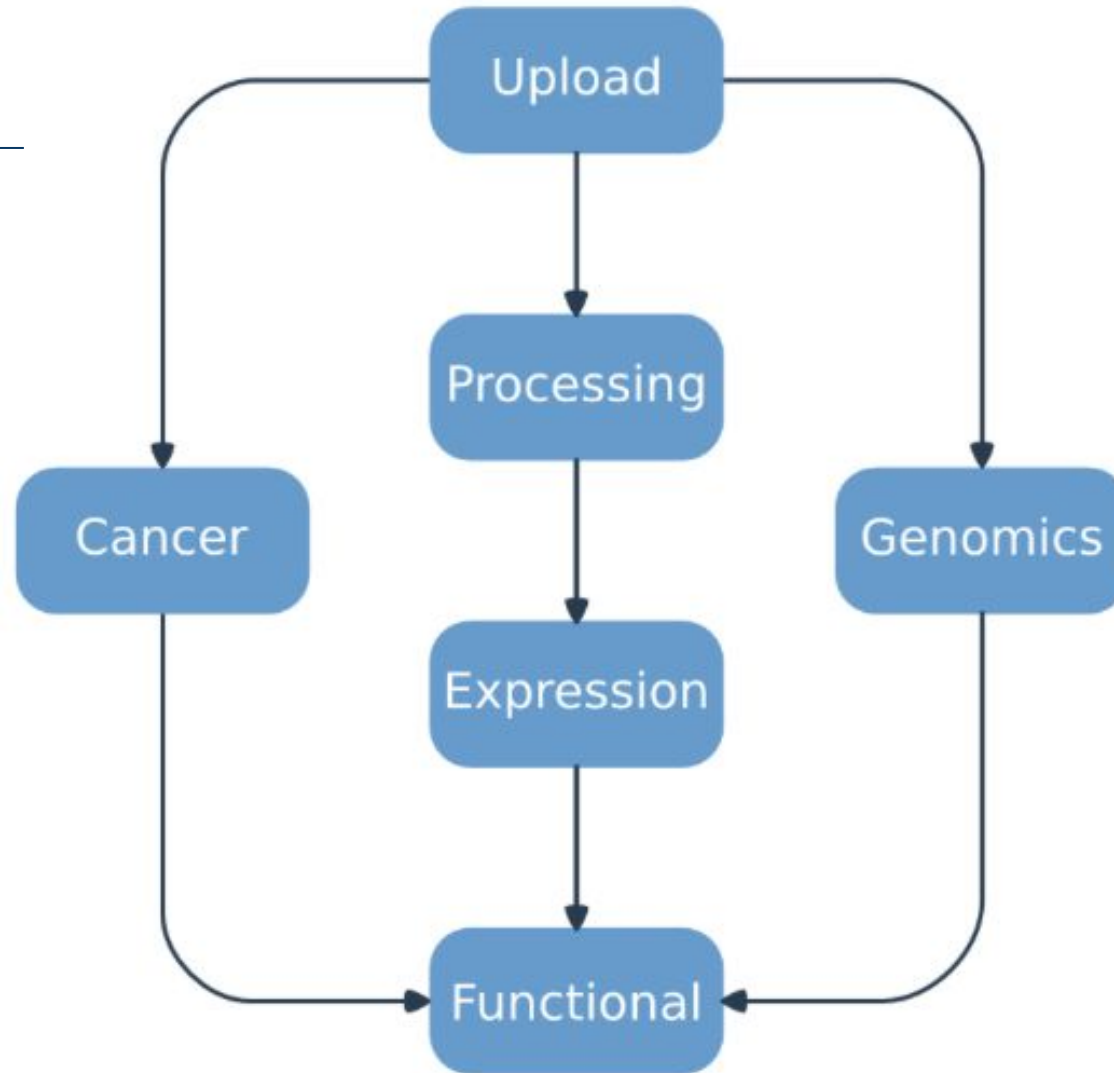
# RNA-seq data analysis pipeline

# Counts

# Counts

Sample

Gene

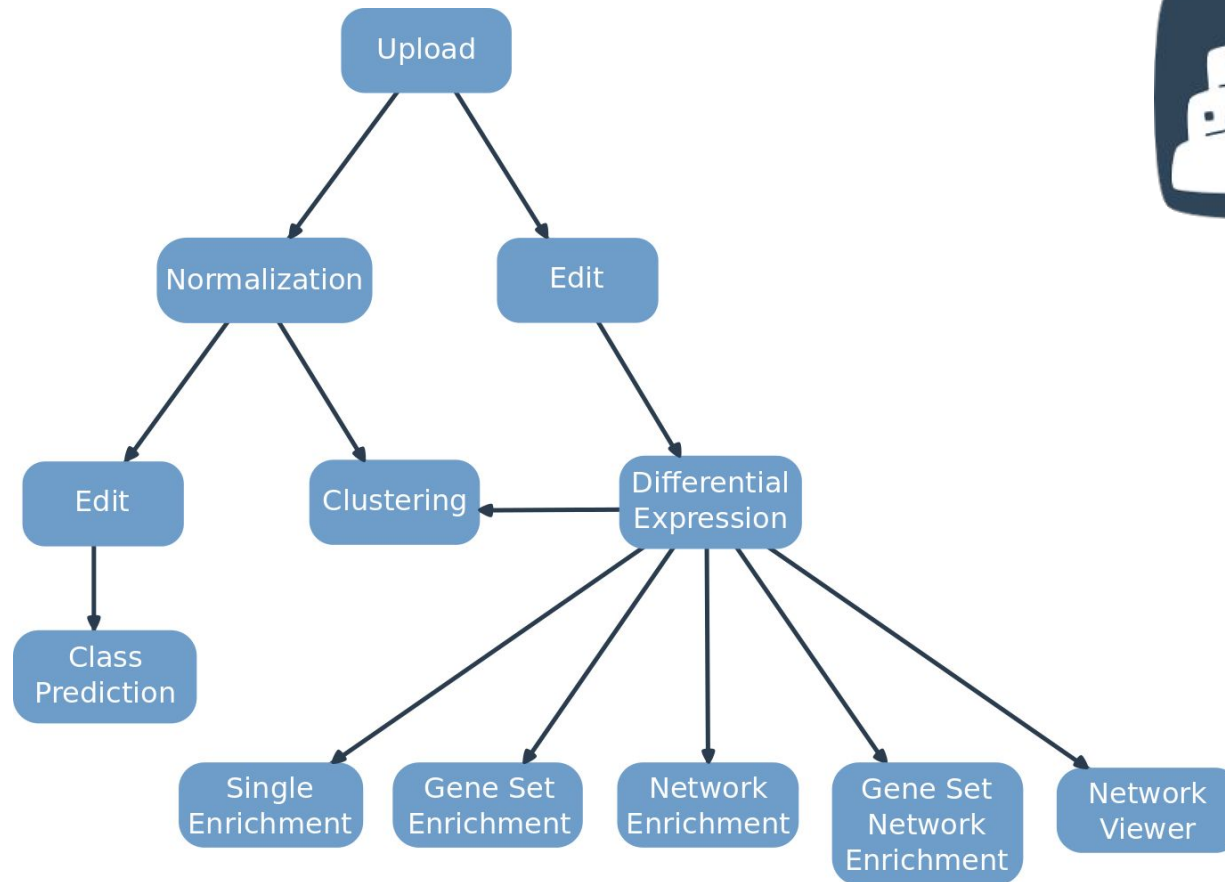| Ensembl | Gene.Name | T1 | T2 | T3 | T4 | T5 | WT1 | WT2 | WT3 | WT4 | WT5 | WT6 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ENSMUSG00000000134 | Tfe3 | 312 | 295 | 333 | 258 | 392 | 257 | 344 | 223 | 423 | 277 | 389 |
| ENSMUSG00000000142 | Axin2 | 165 | 171 | 138 | 166 | 203 | 170 | 172 | 119 | 203 | 147 | 178 |
| ENSMUSG00000000148 | Brat1 | 213 | 196 | 207 | 224 | 350 | 204 | 268 | 143 | 300 | 177 | 288 |
| ENSMUSG00000000149 | Gna12 | 684 | 684 | 613 | 545 | 900 | 496 | 672 | 426 | 1023 | 583 | 797 |
| ENSMUSG00000000154 | Slc22a18 | 3 | 2 | 3 | 2 | 2 | 3 | 3 | 2 | 1 | 1 | 3 |
| ENSMUSG00000000157 | Itgb2l | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ENSMUSG00000000159 | Igsf5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ENSMUSG00000000167 | Pih1d2 | 15 | 19 | 6 | 10 | 9 | 5 | 5 | 5 | 7 | 6 | 6 |
| ENSMUSG00000000168 | Dlat | 899 | 777 | 967 | 756 | 1116 | 777 | 1047 | 614 | 1155 | 894 | 1126 |
| ENSMUSG00000000171 | Sdhd | 1055 | 1003 | 1047 | 914 | 1430 | 939 | 1192 | 766 | 1390 | 916 | 1412 |
| ENSMUSG00000000182 | Fgf23 | 1 | 0 | 3 | 1 | 0 | 2 | 0 | 2 | 2 | 0 | 0 |
| ENSMUSG00000000183 | Fgf6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| ENSMUSG00000000184 | Ccnd2 | 1961 | 1978 | 1804 | 1779 | 2090 | 1655 | 2148 | 1585 | 2504 | 1895 | 2274 |
| ENSMUSG00000000194 | Gpr107 | 784 | 733 | 667 | 615 | 889 | 654 | 818 | 483 | 1034 | 627 | 1015 |
| ENSMUSG00000000197 | Nalcn | 1120 | 1009 | 1047 | 917 | 1356 | 1129 | 1202 | 758 | 1625 | 1127 | 1044 |

# Babelomics 5



Babelomics 5
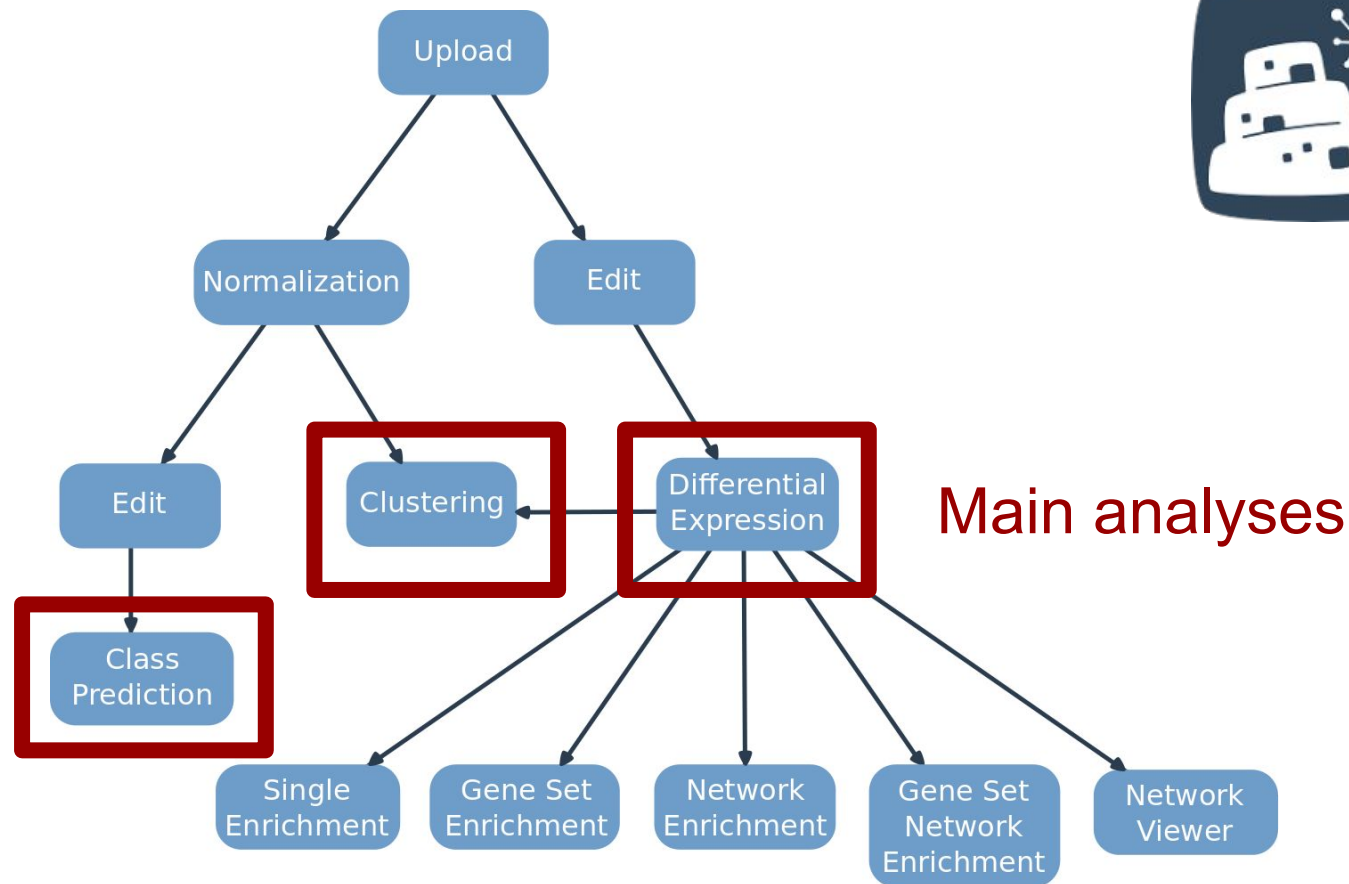GENE EXPRESSION, GENOME VARIATION AND FUNCTIONAL PROFILING ANALYSIS SUITE
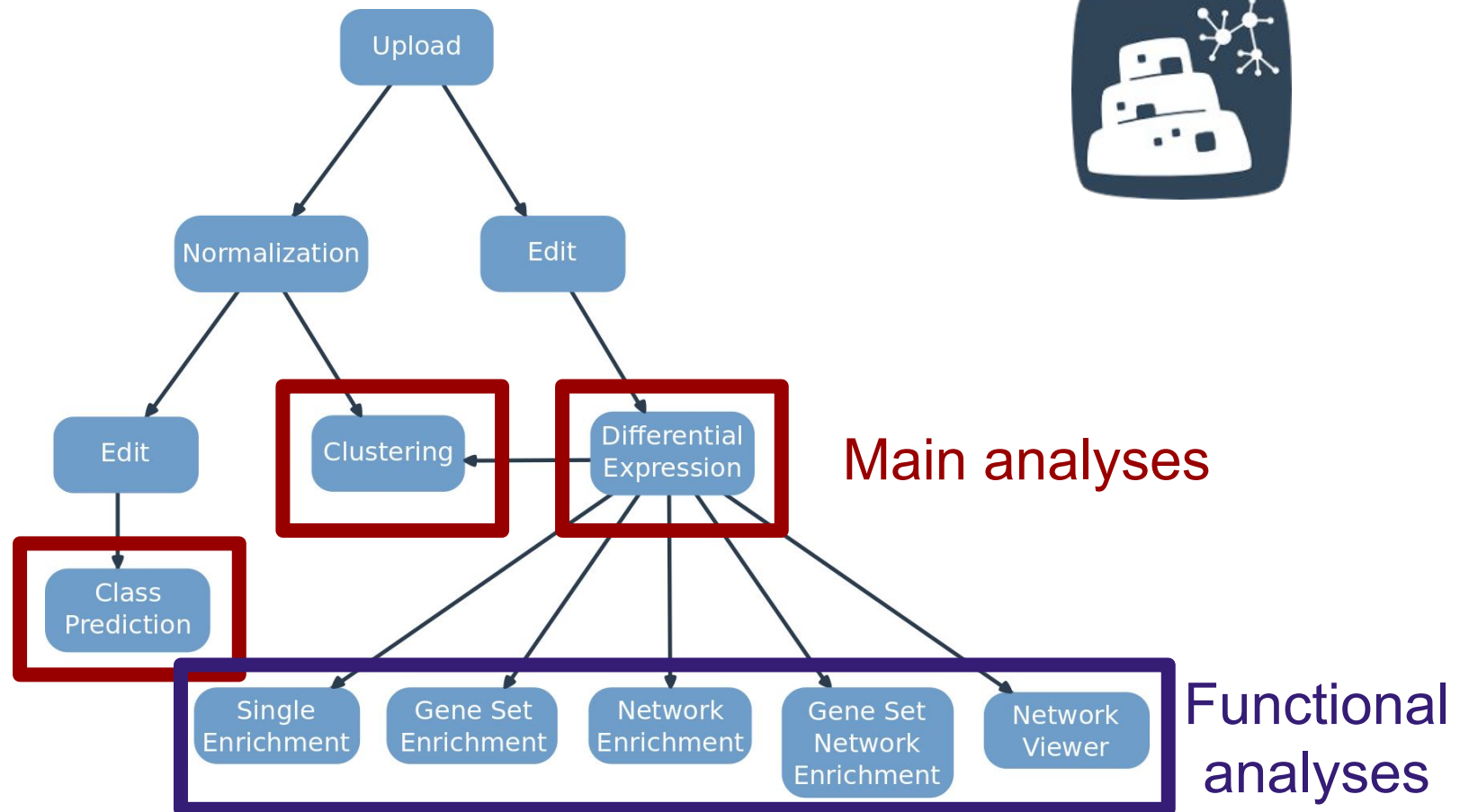
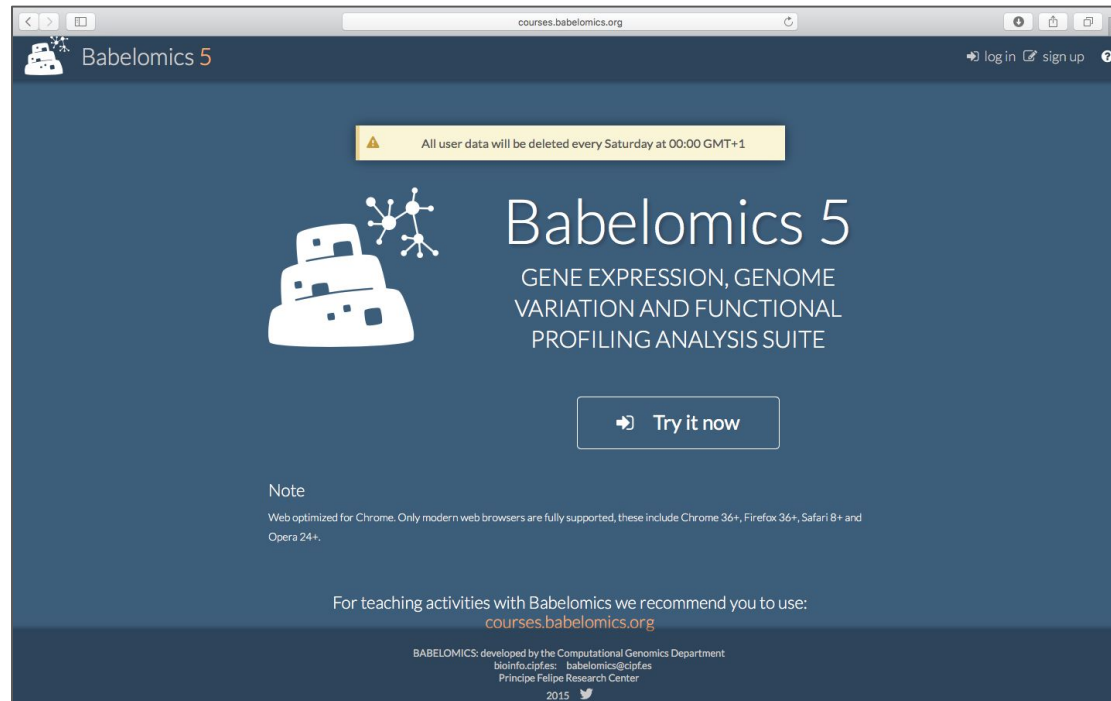http://courses.babelomics.org

# RNA-seq pipeline in Babelomics 5



http://courses.babelomics.org

# RNA-seq pipeline in Babelomics 5



http://courses.babelomics.org

# RNA-seq pipeline in Babelomics 5



http://courses.babelomics.org

# Babelomics 5: web structure



http://courses.babelomics.org

Practical exercises: http://bioinfo.cipf.es/mbb

# Normalization

## Why normalizing?

- The technology introduces different biases
- We need to remove them to compare
  - Among genes in a sample
  - Among samples

## Biases

- Gene length: larger genes get more reads
- Library depth: deeper libraries get more reads
- RNA composition: some genes steal reads from other genes
- Others

# Normalization methods

- Reads per kilobase per million (**RPKM**): Removes gene length and library depth biases.

$$RPKM = \frac{total\ exon\ reads}{mapped\ reads\ (millions) * exon\ length\ (KB)}$$

- Trimmed means of M-values (**TMM**): assumes only a few genes are differentially expressed and changes library depth.

- **Quantiles**: makes all sample distributions the same.

# Normalization methods in Babelomics 5

- RPKM (gene length required)

- TMM

- TMM with gene length correction (gene length required)

- Automatic selection of the method based on the diagnostic test for differences in RNA composition from NOISeq.

# Babelomics 5: normalization



http://courses.babelomics.org

Normalization exercise: http://bioinfo.cipf.es/mbb

# Differential expression

# Differential expression

Genes that show statistically significant differences in expression level between conditions

# Babelomics 5: differential expression



http://courses.babelomics.org

Normalization exercise: http://bioinfo.cipf.es/mbb

# Functional profiling: single enrichment



Association test

# Functional profiling: gene set enrichment



Block of genes enriched in the annotation **A**

Annotation **C** is homogeneously distributed along the list

Block of genes enriched in the annotation **B**

# Babelomics 5: functional enrichment



http://courses.babelomics.org

Single enrichment exercise: http://bioinfo.cipf.es/mbb

# Classification methods: unsupervised and supervised

# Babelomics 5: classification



http://courses.babelomics.org

Clustering exercise: http://bioinfo.cipf.es/mbb

# Babelomics 5: RNA-seq pipeline

http://courses.babelomics.org



Pipeline exercise: http://bioinfo.cipf.es/mbb