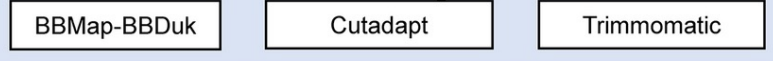


RNA-seq: Getting counts

RNA-seq data analysis workflow

(1) Raw gene expression quantification

Trimming



Alignment against genome



Hybrid alignment (genome + transcriptome)



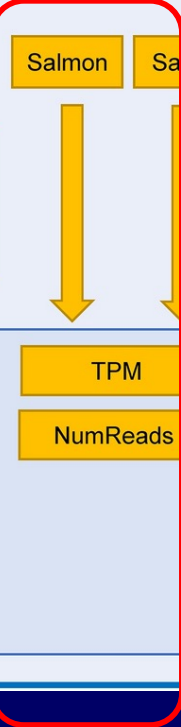
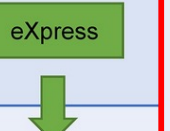
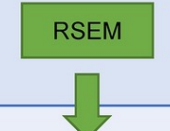
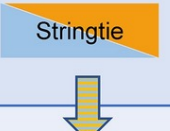
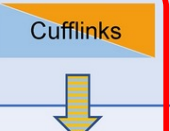
Alignment against transcriptome



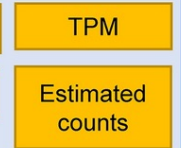
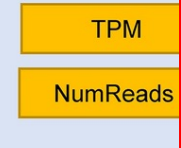
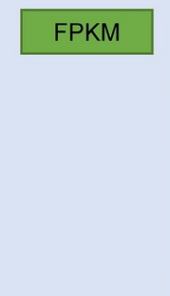
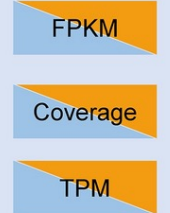
Pseudoalignment



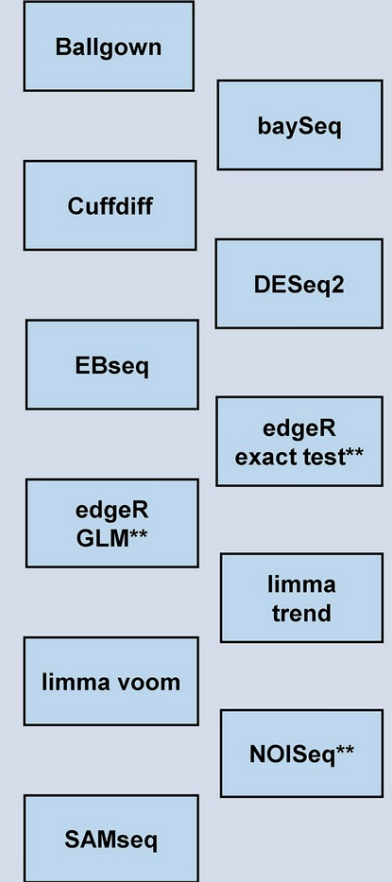
Counting



Normalization



(2) Differential gene expression



RNA-Seq: Getting counts

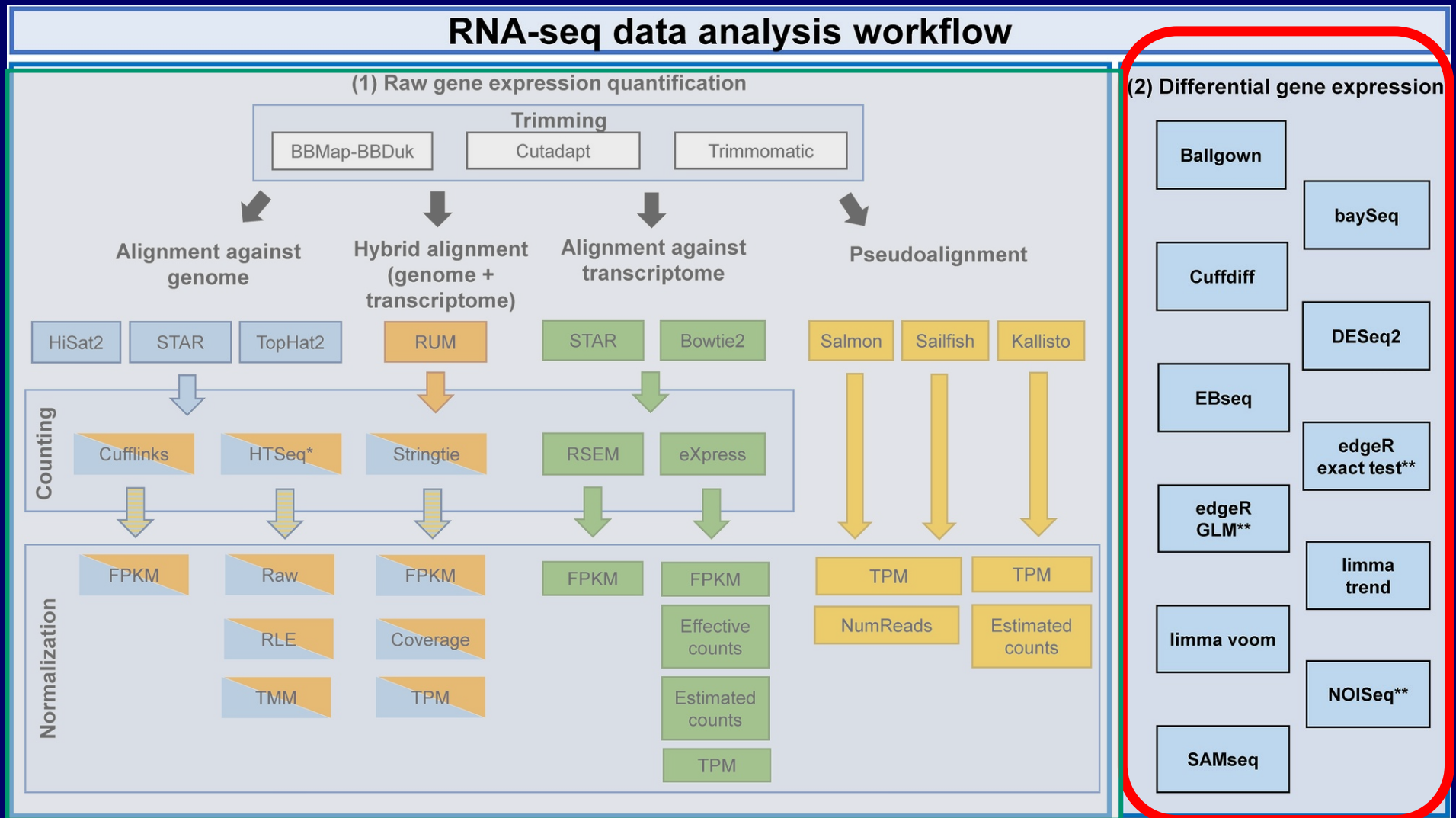
- ❑ Raw – counts (reads) per gene.
- ❑ Normalized
 - ❑ FPKM (Fragments Per Kilobase gene length and per Million reads)
 - ❑ TPM (Transcripts Per Million)
- ❑ Depending on the which program will be used for identifying DEGs.
 - ❑ DESeq (DESeq2) requires raw counts
 - ❑ CuffLinks generated normalized counts as well as models for CuffDiff.

RNA-Seq Overview

Four major steps, **semi-independent** of each other.

- I. Mapping → produce SAM/BAM or counts data.
- II. Quantification → produce RPKM for each gene/transcript.
- III. Identifying DEG (Differentially expressed genes) → gene list.

RNA-seq: Identify DEGs



Many options at this stage. Personal favorites –
Cuffdiff and **DESeq2**

Identification of Differentially Expressed Genes (DEGs)

```
module load cufflinks
```

```
## Frist merge the gtf files for samples to be compared.
```

```
In /ufrc/gms6014/share/genome/dm6/annotation/genes.gtf dm6.gtf
```

```
In /ufrc/gms6014/share/genome/dm6/sequence/genome.fa dm6.fa
```

```
cuffmerge -g dm6.gtf -s dm6.fa -p 2 WG_assemblies.txt
```

```
./WG_young_1.clout/transcripts.gtf  
./WG_young_2.clout/transcripts.gtf  
./WG_old_1.clout/transcripts.gtf  
./WG_old_2.clout/transcripts.gtf
```

Identification of Differentially Expressed Genes (DEGs)

```
module load cufflinks
```

```
## Frist merge the gtf files for samples to be compared.
```

```
In /ufrc/gms6014/share/genome/dm6/annotation/genes.gtf dm6.gtf
```

```
In /ufrc/gms6014/share/genome/dm6/sequence/genome.fa dm6.fa
```

```
cuffmerge -g dm6.gtf -s dm6.fa -p 2 WG_assemblies.txt
```

```
./WG_young_1.clout/transcripts.gtf  
./WG_young_2.clout/transcripts.gtf  
./WG_old_1.clout/transcripts.gtf  
./WG_old_2.clout/transcripts.gtf
```

Identification of differentially expressed genes (DEGs)

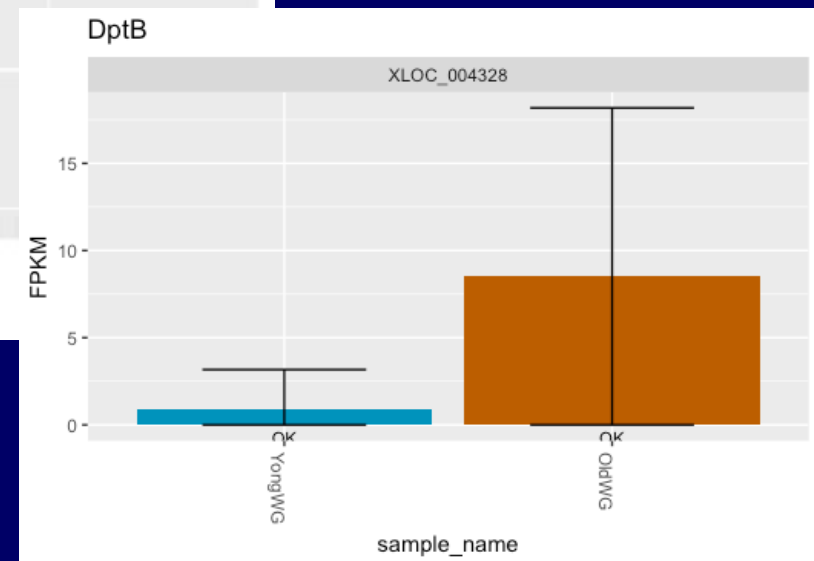
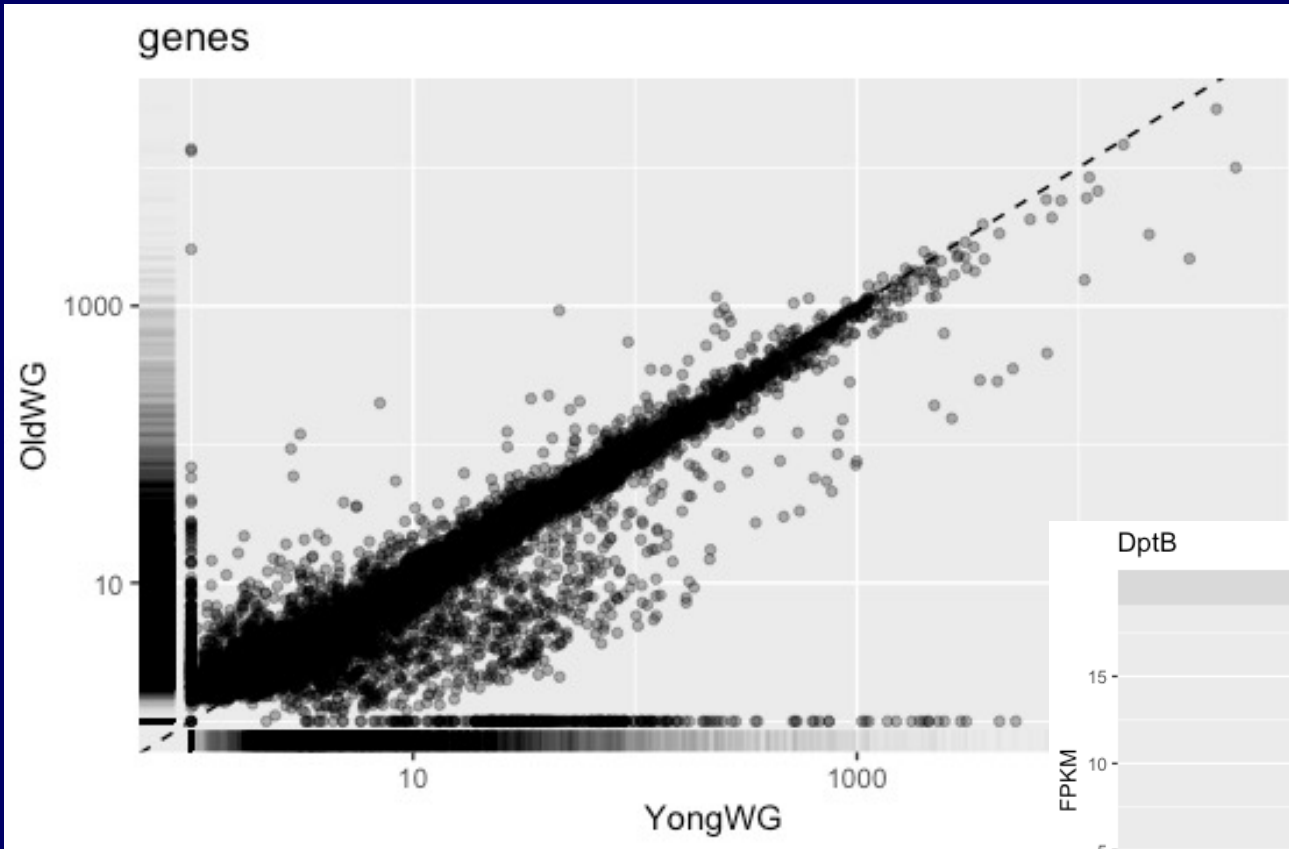
```
module load cufflinks
```

```
cuffdiff -o Old_v_Young -b ./index/Dm6.44.fa -u Merged/merged.gtf -p 2 -L youngWG,oldWG \  
./starMap/WG_young_1Aligned.sortedByCoord.out.bam,./starMap/WG_young_2Aligned.sortedByCoord.out.bam \  
./starMap/WG_old_1Aligned.sortedByCoord.out.bam,./starMap/WG_old_2Aligned.sortedByCoord.out.bam
```

Practice – observe CuffDiff results

- Transfer the CuffDiff result folder (“Old_v_Young”) from HiPerGator to your own computer.
- Observe (you may force it to be opened by Excel by adding .xlsx extension):
 - Gene_exp.diff
 - Gene_exp.tracking.
- Further explore of results with the R package “cummeRbund” .

Observe CuffDiff results with cummeRbund



RNA-Seq Overview

Four major steps, **semi-independent** of each other.

- I. Mapping → produce SAM/BAM or counts data.
- II. Quantification → produce RPKM for each gene/transcript.
- III. Identifying DEG (Differentially expressed genes) → gene list.
- IV. Identifying affected biological processes/pathways.

Functional Analysis of HTS data

- Gene Ontology –
<http://www.geneontology.org/>
- Regulatory pathways.
- Modeling & Systems Biology.

Gene Ontology – hierarchical framework of terms / concepts

AmiGO : Tree View - Microsoft Internet Explorer

File Edit View Favorites Tools Help

[Top Docs](#) [Gene Ontology](#) [GO Links](#) [GO Summary](#)

- GO:0003673 : Gene Ontology (46199)
 - GO:0008150 : biological_process (30188)
 - GO:0016265 : death (525)
 - GO:0008219 : cell_death (484)
 - GO:0012501 : programmed_cell_death (447)
 - GO:0006915 : apoptosis (419)
 - GO:0006916 : anti-apoptosis (111)
 - GO:0008632 : apoptotic_program (51)
 - GO:0008637 : apoptotic_mitochondrial_changes (11)
 - GO:0030262 : apoptotic_nuclear_changes (10)
 - GO:0030263 : apoptotic_chromosome_condensation (1)
 - GO:0006309 : DNA_fragmentation (9)
 - GO:0030264 : nuclear_fragmentation (0)
 - GO:0006919 : caspase_activation (16)
 - GO:0006921 : disassembly_of_cell_structures (10)
 - GO:0008633 : induction_of_proapoptotic_gene_products (0)
 - GO:0045884 : regulation_of_survival_gene_products (7)
 - GO:0006917 : induction_of_apoptosis (148)
 - GO:0008624 : induction_of_apoptosis_by_extracellular_signals (46)
 - GO:0008629 : induction_of_apoptosis_by_intracellular_signals (23)
 - GO:0019051 : induction_of_apoptosis_by_virus (0)
 - GO:0006925 : killing_of_inflammatory_cells (0)
 - GO:0006927 : killing_transformed_cells (3)
 - GO:0006926 : killing_virus-infected_cells (1)
 - GO:0045476 : nurse_cell_apoptosis (1)
 - GO:0006924 : peripheral_killing_of_activated_T-cells (0)
 - GO:0012502 : induction_of_programmed_cell_death (148)
 - GO:0006917 : induction_of_apoptosis (148)
 - GO:0008624 : induction_of_apoptosis_by_extracellular_signals (46)
 - GO:0008629 : induction_of_apoptosis_by_intracellular_signals (23)
 - GO:0019051 : induction_of_apoptosis_by_virus (0)
 - GO:0005575 : cellular_component (22371)
 - GO:0003674 : molecular_function (37018)

DAG view

Gene Ontology

Goal – “produce a dynamic controlled vocabulary that can be applied to all organisms even as knowledge of gene and protein roles in cells is accumulating and changing” – GO consortium (~2001)

Ontology:

“ The branch of metaphysics that deals with the nature of being” – The American Heritage Dictionary

Implications of Gene Ontology (I)

Monitoring biological processes or molecular functions beyond individual gene.

Example:

1.) Which biological process (mol. Function) is activated/suppressed following a treatment?

Gene Expression Profile Differences between the two lung cancer cell lines A549 and H23

extracellular (GO:0005576)	1.91E-08	169
Cell Communication	1.32E-07	690
plasma membrane (GO:0005886)	1.34E-07	511
Complement and coagulation cascades - Homo sapiens	1.73E-07	20
Metabolism	2.10E-06	174
carbohydrate metabolism (GO:0005975)	2.45E-06	207
cell adhesion molecule activity (GO:0005194)	0.000102	113
Structural Protein	0.000231	271
extracellular matrix (GO:0005578)	0.000235	53
Cell Growth and Maintenance	0.000569	590
Cell Adhesion	0.000917	100

development	1.40E-07	596
cell differentiation (GO:0030154)	6.60E-05	186
regulation of gene expression, epigenetic (GO:0040029)	7.71E-05	442
cell growth (GO:0016049)	8.37E-05	307
transcription regulator activity (GO:0030528)	0.000307	319
extracellular (GO:0005576)	0.000515	153

Implications of Gene Ontology (II)

Basis for cross genome comparison and integrating knowledge from different model systems.

Term	Fly Genes	Worm Genes	Mouse Genes	Human Genes	Sacc. Yeast Genes	Pombe Yeast Genes	Weed Genes
<input checked="" type="checkbox"/> cell cycle	<u>265</u>	<u>182</u>	<u>294</u>	<u>717</u>	<u>424</u>	<u>622</u>	<u>181</u>
<input type="checkbox"/> cell cycle dependent actin filament reorganization	<u>2</u>	0	0	0	<u>4</u>	0	0
<input checked="" type="checkbox"/> DNA replication and chromosome cycle	<u>134</u>	<u>128</u>	<u>67</u>	<u>175</u>	<u>172</u>	<u>73</u>	<u>146</u>
<input type="checkbox"/> endomitotic cell cycle	0	0	<u>1</u>	<u>1</u>	0	<u>3</u>	0
<input checked="" type="checkbox"/> M phase	<u>171</u>	<u>39</u>	<u>69</u>	<u>181</u>	<u>213</u>	<u>253</u>	<u>3</u>
<input checked="" type="checkbox"/> mitotic cell cycle	<u>133</u>	<u>140</u>	<u>102</u>	<u>314</u>	<u>239</u>	<u>202</u>	<u>141</u>
<input checked="" type="checkbox"/> nuclear migration	<u>1</u>	0	0	0	<u>13</u>	0	0
<input checked="" type="checkbox"/> regulation of cell cycle	<u>42</u>	<u>4</u>	<u>136</u>	<u>383</u>	<u>87</u>	<u>65</u>	<u>3</u>
<input type="checkbox"/> schizogony	0	0	0	0	0	0	0
<input type="checkbox"/> second mitotic wave (sensu Drosophila)	<u>1</u>	0	0	0	0	0	0

Tools associated with GO

- A comprehensive [list](#) at GO web site.
- Tools for browsing, AmiGO, QuickGO at EBI, etc.
- Tools for identifying over represented GOs/pathways, etc.

Using GO to gain comprehensive understanding of cellular differences

Practice: Load a gene list to identify over-represented GO