From gene expression to networks: R workflow development for analysis and visualization of RNA-seq data



Environnemen

Sciences Végétales

### Medicago truncatula

### Legume plants :

- Atmospheric N2 → a non-limiting source of nitrogen
- Symbiosis with Rhizobia
- Nodules
- Symbiotic nitrogen fixation (SNF)





### Medicago truncatula

### Legume plants :

- Atmospheric N2 → a non-limiting source of nitrogen
- Symbiosis with Rhizobia
- Nodules
- Symbiotic nitrogen fixation (SNF)



Development and function of root and nodule

 $\rightarrow$  sensitive to environmental stresses



- Reduction of nitrogen fixation efficiency
- Unstable yields in legume crops



**Biotic** stresses



### Systemic response

Split root experimental system















Ruffel et al, 2008 Plant Physiol. Jeudy et al, 2010, New phytol. Laguerre et al, 2012, New phytol.



What are the molecular mechanisms underline the systemic signaling in response to abiotic and biotic stresses in *Medicago truncatula*?



What are the molecular mechanisms underline the systemic signaling in response to abiotic and biotic stresses in *Medicago truncatula*?





**Tissu factor** 

## Time-line treatment factor

Control - 6h - 1d - 3d - 5d

**Replicate factor (3)** 



What are the molecular mechanisms underline the systemic signaling in response to abiotic and biotic stresses in *Medicago truncatula*?





Are these mechanisms, common or specific to different stresses ?



### **RNAseq data analysis**









## R workflow





# Filtering & Normalization

### Normalization<sup>2</sup>

#### Method

Trimmed Mean of M-values (TMM)

#### Tool

calcNormFactors function edgeR package

#### Input

Filtered counts table

#### Outputs

Normalized counts Mean;STD Normalized counts Log2 Normalized counts Log2 Mean; STD Normalized counts



Netbio

Robinson et al., 2010, Bioinformatics <sup>2</sup> Dillies et al., 2013, Briefing in Bioinformatics







## **WORKFLOW** Filtering

**AUTOMATED** 

**RNA-seq ANALYSIS** 

Normalization

Data quality control

**Differential expression** analysis

> **Co-expression** analysis

**Co-expression** Networks



## **Filtering**

#### Method

Counts per million (CPM)

#### Tool

*cpm* function edgeR package<sup>1</sup> (v.3.22.3)

#### Input

Table of counts

### **Outputs**

Non expressed genes Low expressed genes

# Data quality control



13/12/18

#### Tool

#### ggplot2, stats, FactoMineR packages









#### PCA on Normalized Counts



Boxplot of normalized counts

















#### AUTOMATED RNA-seq ANALYSIS WORKFLOW

Filtering

Normalization

Data quality control

Differential expression analysis

> Co-expression analysis

Co-expression Networks





#### <u>Outputs</u>

#### For each contrast



**DEGs Results** 





Gene ID mRNA ID Gene Name Medtr0007s0390 Medtr0007s0390.1 receptor-like Serine/Threonine-kinase plant Medtr0184s0010 Medtr0184s0010.1 UDP-glucosyltransferase Medtr0221s0020 Medtr0221s0020.2 transmembrane protein, putative Medtr0221s0020 Medtr0221s0020.1 transmembrane protein, putative Medtr0384s0020 Medtr0384s0020.1 hypothetical protein Medtr0508s0010 Medtr0508s0010.1 ribosomal protein S8e family protein logFC logCPM LR **PValue** FDR 5.55391970508026 4.81041290555745 16.1086222041737 5.98111063390226e-05 0.0108642973940508 4.58308025736963 5.32108313676033 16.1233960404379 5.93463462776849e-05 0.0108642973940508 -4.58621127038228 5.95209903996132e-05 0.0108642973940508 6.11277538970362 16.1178308342364 -4.58621127038228 6.11277538970362 5.95209903996132e-05 0.0108642973940508 16.1178308342364 -5.45355905337383 5,690285243976 19.7243656112101 8.94535369589355e-06 0.00266942133362172 4.85866597399211 5.95949005692566 11.8634299454211 0.000572470120976213 0.0488095034982535

MA-plot



plotSmear function
 (edgeR package)



Gene clustering

heatmap.2 function (gplots package) (v.3.0.1)

## Single gene expression profile



ggline function (ggpubr package) (v.0.1.8)

#### AUTOMATED RNA-seq ANALYSIS WORKFLOW





### <u>Outputs</u>

#### Summary results





# Histogram of Up/Down DGE (ggplot2 package)

Comparison table

(data.table package) (v.1.11.4)

Contrast1	Contrast2	Contrast3	Contrast4	Contrast5	Contrast6
0	0	0	1	0	0
1	0	0	1	1	0
1	0	0	1	1	0
0	0	0	1	0	0
0	0	0	0	0	0
0	0	0	0	0	0
	Contrast1 0 1 0 0 0	Contrast1 Contrast2 0 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0	Contrast1 Contrast2 Contrast3 0 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0	Contrast1         Contrast2         Contrast3         Contrast4           0         0         0         1           1         0         0         1           1         0         0         1           0         0         0         1           0         0         0         1           0         0         0         0           0         0         0         0	Contrast1         Contrast2         Contrast3         Contrast4         Contrast5           0         0         0         1         0           1         0         0         1         1           1         0         0         1         1           0         0         0         1         1           0         0         0         1         0           0         0         0         0         0         0           0         0         0         0         0         0         0

Venn diagram – Intersection – Union function





of raw counts

 $p_{ij} = \frac{y_{ij}/s_j + 1}{\sum_{\ell} y_{i\ell}/s_{\ell} + 1}$ 

Sj: scaling normalization factors



13/12/18

#### **AUTOMATED RNA-seq ANALYSIS** WORKFLOW

Filtering

Normalization

Data quality control

Differential expression analysis

Co-expression analysis

Co-expression Networks





### **Outputs**

Summary coseq: ************************************	****	* * *		
<pre>Model: Gaussian_pk_Lk_Ck Transformation: arcsin ************************************</pre>	****	* * *		
Clusters fit: 10,11,12,13,14,15 Clusters with errors: Selected number of clusters via ICL of selected model: -349190.	5,16,17,18,19,20 ICL: 12 3	,21,22,23,24,25,	26,27,28,29,30	
Number of clusters = 12 ICL = -349190.3 ******	*****	* * *		
Cluster sizes: Cluster 1 Cluster 2 Cluster 818 213 79 Cluster 9 Cluster 10 Cluster 1 393 243 42	3 Cluster 4 C 7 445 1 Cluster 12 0 392	luster 5 Cluste 280	er 6 Cluster 7 Clu 438 602	ster 8 313
Number of observations with MAP 4310 (80.5%)	9 > 0.90 (% of t	otal):		
Number of observations with MAP Cluster 1 Cluster 2 Cluster 3 710 186 663 (86.8%) (87.32%) (83.19%) Cluster 10 Cluster 11 Cluster 201 327 294 (82.72%) (77.86%) (75%)	<pre>P &gt; 0.90 per clu Cluster 4 Clust 354 248 (79.55%) (88.5 12</pre>	ster (% of total er 5 Cluster 6 C 349 4 7%) (79.68%) (	L per cluster): Cluster 7 Cluster 8 430 269 (71.43%) (85.94%)	Cluster 9 279 (70.99%)
18800-	-340000 -		800 -	
18000-	-342500 -		600 - Seg	May
0000 H 186000 -	<u>전</u> -345000-		o too - oo	ordtienal probability 0.8 <0.8
182000 -	-3/7500-		200-	
Log Likelihood	10 15	luser L 25 30	Max condition	al probability

15

**Outputs** 



DN2\_3d SN2\_1d

Normalized expression profiles

### Clusters table

Gene_ID	Clusters
Medtr0001s0660	2
Medtr0002s0020	1
Medtr0002s1060	10
Medtr0003s0560	3
Medtr0003s0580	4
Medtr0003s0590	3



Maximum conditional probability



**AUTOMATED** 

**RNA-seq ANALYSIS** WORKFLOW

Filtering

Normalization

Data quality control

Differential expression analysis

**Co-expression** 

analysis

**Co-expression** 

Networks

**FUNCTIONAL MODULES** 

**IDENTIFICATION** 

**Enrichment analysis** 

П























AUTOMATED RNA-seq ANALYSIS WORKFLOW		Tal	ble of occurren	се			
Filtering	Source_Gene	Edge_score	Target_Gene	SR1	SR2	SR3	Cluster_group
Normalization	MtrunA17Chr1g0148241 MtrunA17Chr1g0149811 MtrunA17Chr4g0071991	2 2 1	MtrunA17Chr1g0209251 MtrunA17Chr1g0209691 MtrunA17Chr7g0240451	N N N	3 7 8	1 3 N	N31 N73 N8N
Data quality control	MtrunA17Chr5g0398081	3	MtrunA17Chr5g0444331	1	1	2	112
Differential expression analysis							
Co-expression analysis							
Co-expression Networks							
FUNCTIONAL MODULES							
Enrichment analysis							



AUTOMATED RNA-seq ANALYSIS WORKFLOW	Table of occurrence						
Filtering	Source_Gene	Edge_score	Target_Gene	SR1	SR2	<b>SR3</b>	Cluster_group
Normalization	MtrunA17Chr1g0140241 MtrunA17Chr1g0149811 MtrunA17Chr4g0071991	2	MtrunA17Chr1g0209691	N	7 8	3 N	N73
Data quality control	MtrunA17Chr5g0398081	3	MtrunA17Chr5g0444331	1	1	2	112
Differential expression analysis			SR1				SR2
Co-expression analysis	Cluster_grou	Nitroge	Drought 1				
Co-expression Networks	MtrunA17Chr1g018557 AdminA17Chr3g00898 MtrunA17Chr4g0049607 AdminA17Chr3g014552	nA17Chr4g0044901	0.6 -			0.8 -	
FUNCTIONAL MODULES IDENTIFICATION	MtrunA17Chr4g007267 UrrunA17Chr4g0022991 HtrunA17Chr5g0398001 MtrunA17Chr3g01	Mininal/Chr8g035041	351 0.2 -			0.4 <b>-</b> 0.2 <b>-</b>	<b>■</b>
Enrichment analysis	MtrunA17Chr7g0245551 4thunA17Chr4g0065221 MtrunA17Chr2g031522 MtrunA17Chr2g031522	A17Chr1g020930hrmA17Chr5g0444331	0.0 -	•	SR3	0.0 -	
	MtrunA17Chr7g0267511 MtrunA17Chr	r6g048 <del>5333Mfru</del> nA17Chr2g0287841		Apha	anon ²	nyce	S
	Sub-Netw co-expre	work of ession				•	
13/12/18					T <b>T</b>	<b>↓</b> +	

#### AUTOMATED RNA-seq ANALYSIS WORKFLOW

Filtering

Normalization

Data quality control

Differential expression analysis

Co-expression analysis

Co-expression Networks





### <u>Method</u>

Hypergeometric test

$$P(X=k)=rac{\binom{K}{k}\binom{N-K}{n-k}}{\binom{N}{n}}$$



N : population size

K : number of success in the population

n : sample size

k : number of success in the sample

#### AUTOMATED RNA-seq ANALYSIS WORKFLOW

Filtering

Normalization

Data quality control

Differential expression analysis

Co-expression analysis

Co-expression Networks





## <u>Method</u>



### Differential expression analysis



- N : population size
- K : number of success in the population
- n : sample size
- k : number of success in the sample

Co-expression analysis











## **Conclusion & Prospects**





## Acknowledgments



### **MASCE** team



Brigitte Brunel PR1, Montpellier SupAgro Montpellier, France <u>Contact</u> <u>Publications</u>



Karine Heulin IE, Inra Montpellier, France <u>Contact</u>



Antoine Le Quéré CR2, IRD Rabat, Maroc <u>Contact</u> Publications



Marc Tauzin TR, Inra Montpellier, France Contact



Stefano Colella CR1, Inra Montpellier, France <u>Contact</u> <u>Publications</u>





Stephane BOVIN Mathilde TANCELIN



#### Genomic networks team

Marie-Laure Martin-Magniette (DR, INRA) Christine Paysant-Le-Roux (IE, INRA) Véronique Brunaud (IR, INRA) Nathalie Boudet (MdC, UEVE) Cécile Guichard (IE, INRA) Guillem Rigaill (CR, INRA) Jean Philippe Tamby (IE, INRA) Julien Rozière (IE, CDD) Margot Correa (IE, CNRS)



Maxime BONHOMME

## Thank you for your attention !