



Genomics NGS Service

Bioinformatics Analysis of RNA-seq de-novo transcriptome by Trinity

Help manual

2017

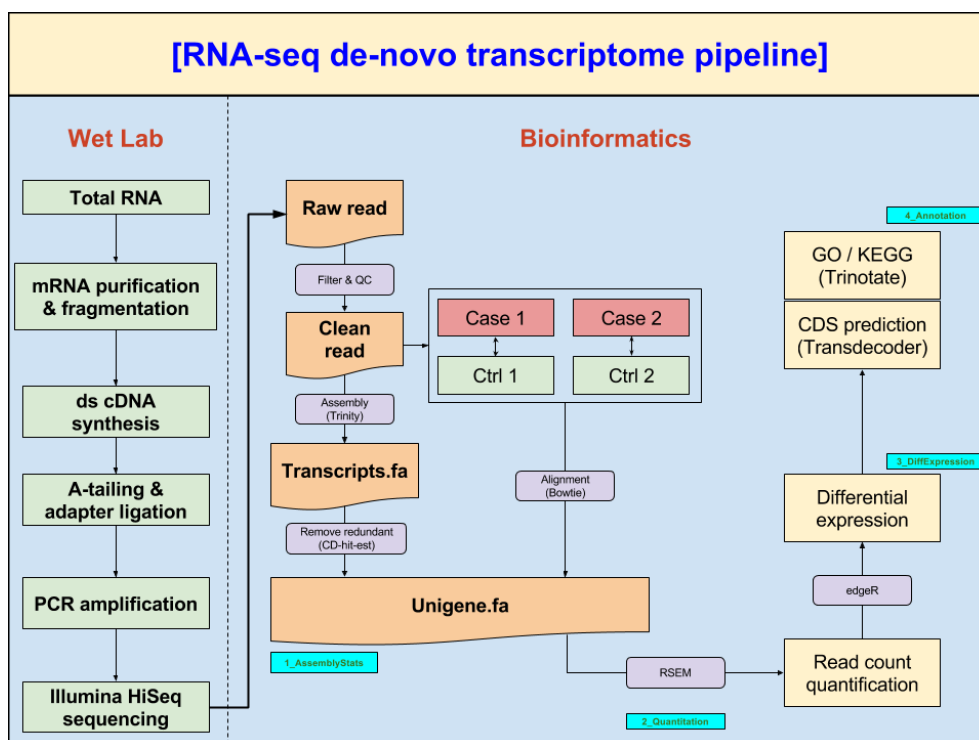
Genomics NGS Analysis Team





Table of Contents

Experiment Process	3
Bioinformatics analysis	4
1. Assembly Stats (1_AssemblyStats)	4
2. Read count quantification (2_Quantitation).....	5
[Alignment stats]:.....	5
[RSEM output]:	6
3. DGE comparisons (3_DiffExpression)	7
[DE output]:	7
4. Annotation (4_Annotation)	9
[GO annotation of transcripts]:	10
[KEGG pathway annotation]:	11
[GO enrichment analysis]:.....	13
5. Reference	15



Experiment Process

- a.) Purify and fragment mRNA: Using poly-T oligo-attached beads to purify mRNA, which is also fragmented for cDNA synthesis.
- b.) Double strand cDNA synthesis: Using reverse transcriptase and random primer to synthesize first strand cDNA, and using dUTP in place of dTTP to generate double-strand cDNA.
- c.) A-tailing and Adaptor Ligation: A single 'A' nucleotide is added to 3' end of ds cDNAs. Then, multiple indexing adapters are ligated to 5' and 3' of the ends of the ds cDNA.
- d.) PCR amplification Using PCR to selectively amplify those DNA fragments that have adapters on both ends.
- e.) Library quality validating: Library was validated on Agilent 2100 Bio-analyzer and Real-Time PCR System.
- f.) Sequencing by Illumina HiSeq platform



Bioinformatics analysis

1. Assembly Stats (1_AssemblyStats)

“Trinity v2.3.2” is a well-known transcriptome de-novo assembly tool. It combines three independent software modules: Inchworm, Chrysalis, and Butterfly, applied sequentially to process large volumes of RNA-seq reads. Trinity partitions the sequence data into many individual de Bruijn graphs, each representing the transcriptional complexity at a given gene or locus, and then processes each graph independently to extract full-length splicing isoforms and to tease apart transcripts derived from paralogous genes [1].

While Trinity job has been completed, it might usually contain lots of duplicate transcripts existed in data. Thus, we commonly use another clustering tool: CD-HIT-EST [2], for processing redundant transcripts removal and try to get more specific unigenes.

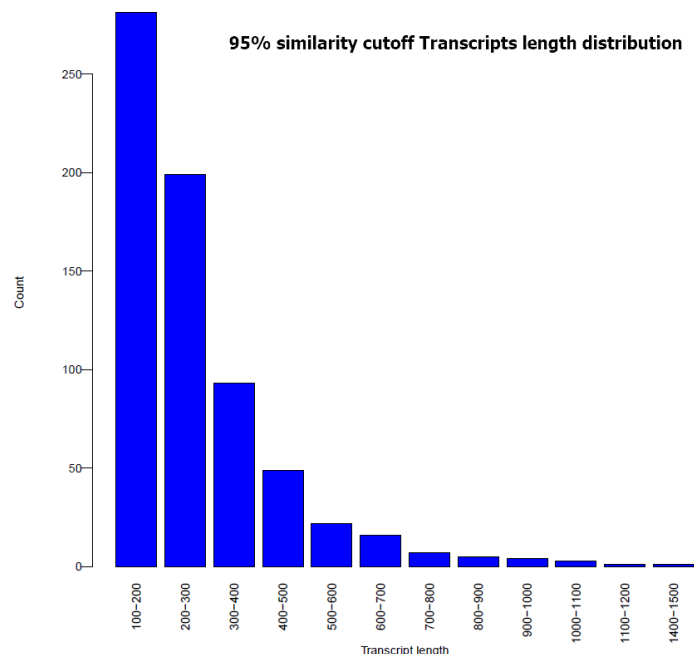
- Trinity parameters:
 - Minimum contig length => 150 bp
- CD-HIT-EST parameters:
 - sequence identity threshold => 95%

Trinity_assembled.final.stats.txt

```
## Counts of transcripts, etc.
Total trinity 'genes': 682
Total trinity transcripts: 686
Percent GC: 44.38
Contig N10: 742
Contig N20: 525
Contig N30: 425
Contig N40: 346
Contig N50: 300

Median contig length: 216
Average contig: 280.01
Total assembled bases: 192085
```

Trinity.95.dist.pdf



2. Read count quantification (2_Quantitation)

In this stage, the de-novo assembled transcriptome will be regarded as backbone reference. All of the samples are going to be aligned for calculating the abundance of read count.

The alignment tool we used is “bowtie v1.1.2” [3], and the read count quantification tool we used is “RSEM v1.2.31” [4].

[Alignment stats]:

Sample 1	
Reads	8,168
Mapped reads	8,168
Pct align	100.0000
...	
Pct mismatch	15.9525
...	
Mapq mean	255.0000
...	
Insert mean	133.0756
...	
Num ref seqs	681
Num ref aligned	403

Note:

- **Pct align:** percent of reads that aligned.
- **Pct mismatch:** percent of reads that have mismatches
- **Mapq:** stats for mapping quality
- **Insert:** stats for insert size
- **Num ref aligned:** number of transcripts aligned by reads

[RSEM output]:

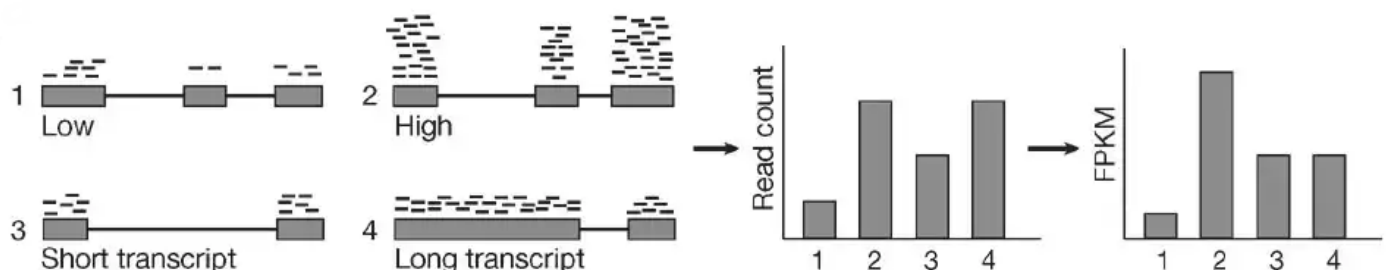
- RSEM.isoforms.results: EM read counts per Trinity transcript (e.g. TRINITY_DN100_c0_g1_i1)
 - RSEM.genes.results: EM read counts per Trinity gene (e.g. TRINITY_DN100_c0_g1)
- * **Basically, we are using “RSEM.isoforms.results” for the downstream jobs.**

transcript_id	gene_id	length	effective_length	expected_count	TPM	FPKM	IsoPct
TRINITY_DN0_c0_g1_i1	TRINITY_DN0_c0_g1	253	117.66	0	0	0	0
TRINITY_DN102_c0_g1_i1	TRINITY_DN102_c0_g1	214	79.28	7	4704.5	21970.57	100
TRINITY_DN107_c0_g1_i1	TRINITY_DN107_c0_g1	214	79.28	2	1344.14	6277.31	100
TRINITY_DN107_c0_g2_i1	TRINITY_DN107_c0_g2	346	210.35	2	506.58	2365.78	100
TRINITY_DN108_c0_g1_i1	TRINITY_DN108_c0_g1	261	125.6	1	424.19	1981.04	100
TRINITY_DN108_c0_g2_i1	TRINITY_DN108_c0_g2	272	136.53	1	390.23	1822.43	100
TRINITY_DN10_c0_g1_i1	TRINITY_DN10_c0_g1	568	432.34	64	7886.96	36833.05	100
TRINITY_DN10_c0_g2_i1	TRINITY_DN10_c0_g2	194	60.1	0	0	0	0
TRINITY_DN110_c0_g1_i1	TRINITY_DN110_c0_g1	211	76.37	1	697.62	3257.99	100

Note:

- **effective_length**: counts only the positions that can generate a valid fragment.
- **expected_count**: sum of the posterior probability of each read comes from this transcripts over all reads.
- **TPM**: Transcripts Per Million. It is a relative measure of transcript abundance. The sum of all transcripts' TPM is 1 million.
- **FPKM**: Fragment Per Kilobase of transcript per Million mapped reads. If reads are paired-end, each R1 or R2 mapped to transcript will be counted 1.
- **IsoPct**: isoform percentage. It is the percentage of expression for a given transcript compared with all expression from that Trinity component. If its parent gene has only one isoform or the gene information is not provided, this field will be set to 100.

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



Ref: (<http://dx.doi.org/10.1038/nmeth.1613>)

3. DGE comparisons (3_DiffExpression)

As we got the read quantification data, various user-provided different comparisons are going to be calculated by “edgeR v3.5” [5], an R package which could process multiple differential expression analysis of RNA-seq expression profile with biological replication.

[DE output]:

- {comparisons}.edgeR.DE_results

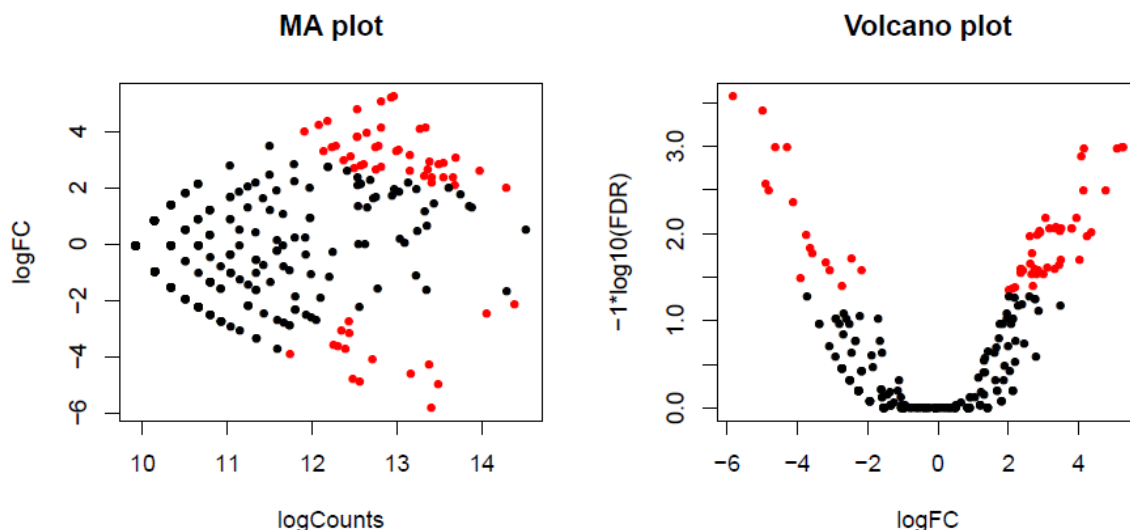
transcript_id	sampleA	sampleB	logFC	logCPM	PValue	FDR
TRINITY_DN265_c0_g2_i1	GSNO_1	wt_1	-5.8337	13.40097	8.42E-07	0.000266
TRINITY_DN386_c0_g1_i1	GSNO_1	wt_1	-4.98723	13.48578	2.43E-06	0.000384
TRINITY_DN121_c0_g1_i1	GSNO_1	wt_1	5.256032	12.96129	1.14E-05	0.001021
TRINITY_DN594_c0_g1_i1	GSNO_1	wt_1	5.223703	12.93235	1.34E-05	0.001021
TRINITY_DN93_c0_g1_i1	GSNO_1	wt_1	-4.63185	13.16146	1.71E-05	0.001021
TRINITY_DN318_c0_g1_i1	GSNO_1	wt_1	-4.28979	13.37738	1.94E-05	0.001021
TRINITY_DN185_c0_g2_i1	GSNO_1	wt_1	4.144669	13.33855	2.50E-05	0.001064

Note:

- **logFC**: log difference between sampleA and sampleB.
- **logCPM**: log counts per million, which is as similar as measuring expression level
- **FDR**: false discovery rate, which could help for validating the false positives in p-value result

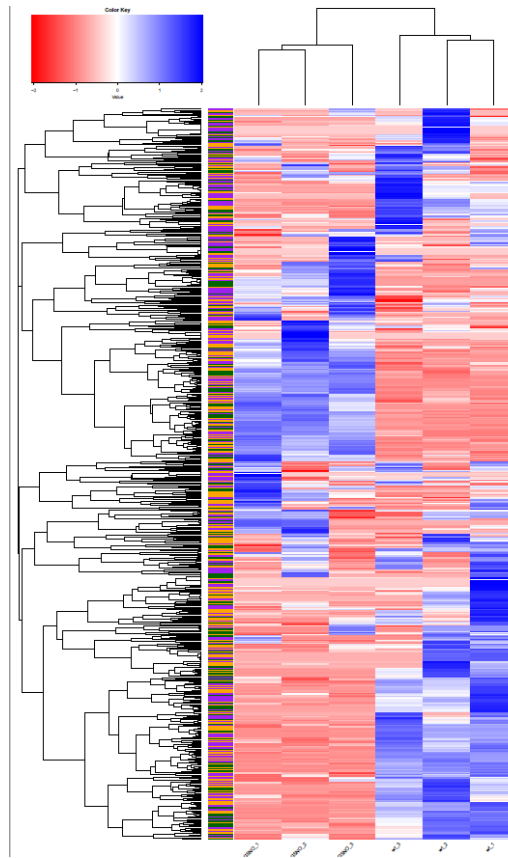
- {comparisons}.edgeR.DE_results.MA_n_Volcano

Red dot: p-value < 0.05



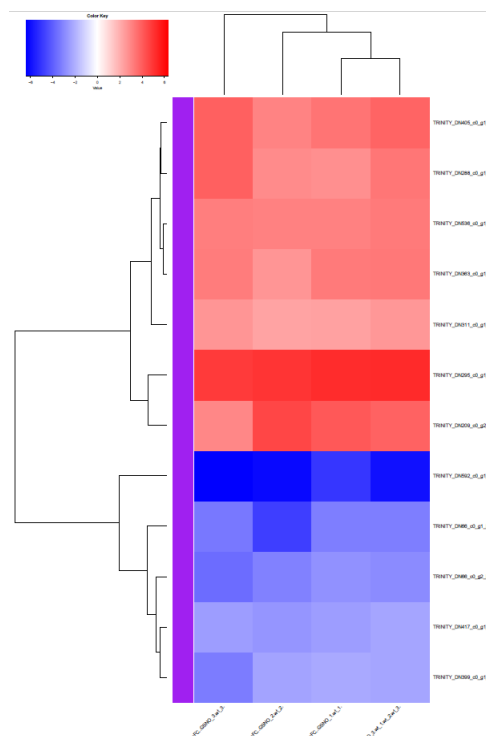
- All_samples_heatmap.pdf

Select **FPKM** value to compare DE by heatmap in each comparison.



- all_groups_heatmap.pdf (only intersection genes within groups will be shown)

Select **p-value<0.05** and **logFC>1** data to compare DE by heatmap in all comparisons, and normalized by **z-score**.



4. Annotation (4_Annotation)

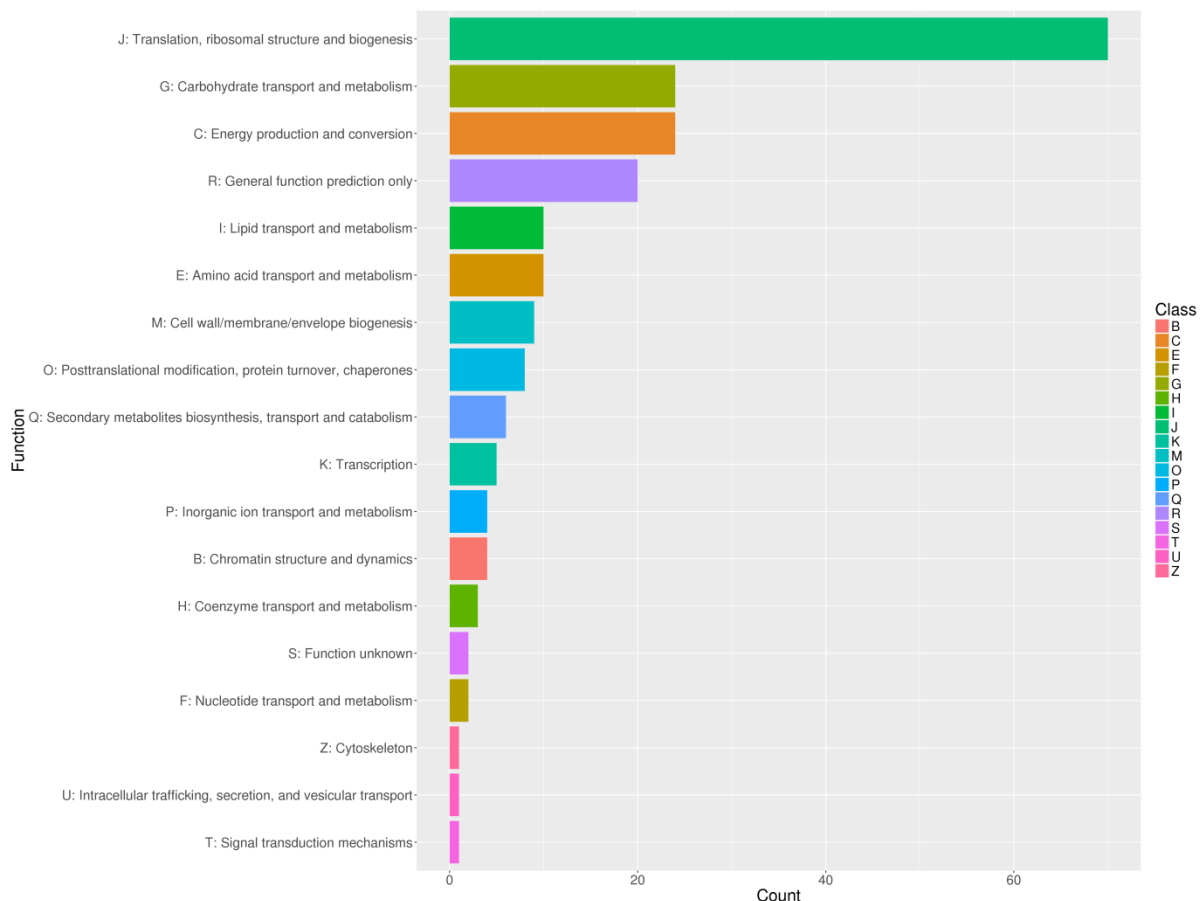
Before annotation work start, we need to parse coding regions within transcripts by gene prediction tool – “Transdecoder v3.0.1” [6] and retrieve protein sequences in the meanwhile.

“Trinotate v3.0.2” is a comprehensive annotation suite designed for functional annotation of de novo assembled transcriptomes, from model or non-model organisms [7]. Our functional annotation works including:

- blastx / blastp: homology search to known & reviewed database (UniprotKB/Swiss-Prot)
- PFAM: protein domain identification
- signalP / TmHMM protein signal peptide and transmembrane domain prediction
- COG / GO / KEGG: functional & pathway annotation

[Protein group function annotation by COG/eggNOG]

In order to extract the maximum amount of information from the rapidly accumulating genome sequences, all conserved genes need to be classified according to their homologous relationships. Each COG consists of individual orthologous proteins or orthologous sets of paralogs from at least three lineages. Orthologs typically have the same function, allowing transfer of functional information from one member to an entire COG.



[GO annotation of transcripts]:

All transcripts are searched to GO slim database which contain a subset of the terms in the whole GO. GO slims are particularly useful for giving a summary of the results of GO annotation of a genome, microarray, or cDNA collection when broad classification of gene product function is required. Once the GO terms have been corresponded to the transcripts, Map2Slim could help us to dig out more informative annotation of transcripts' function.

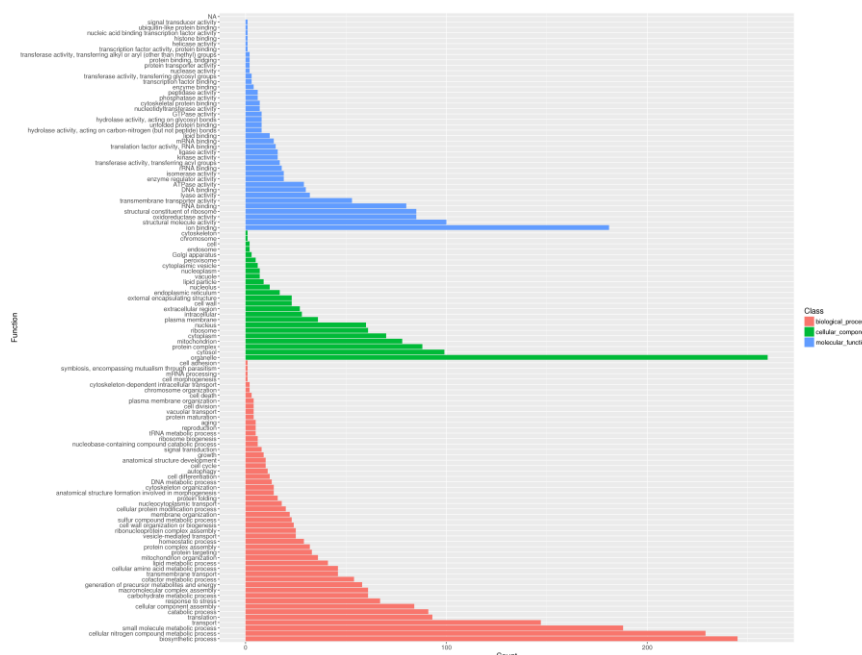
Trinotate_report.xls.gene_ontology (GO terms extraction)

TRINITY_DN0_c0_g1_i1	GO:0003674,GO:0003735,GO:0005198,GO:0005575,GO:0005622
TRINITY_DN0_c0_g2_i1	GO:0003674,GO:0003735,GO:0005198,GO:0005575,GO:0006412
TRINITY_DN102_c0_g1_i1	GO:0000166,GO:0003674,GO:0003824,GO:0004550,GO:0005488
TRINITY_DN105_c0_g1_i1	GO:0003674,GO:0003735,GO:0005198,GO:0005575,GO:0005840
TRINITY_DN105_c0_g2_i1	GO:0003674,GO:0003735,GO:0005198,GO:0005575,GO:0005840
TRINITY_DN109_c0_g1_i1	GO:0000139,GO:0002790,GO:0005575,GO:0005789,GO:0006810
TRINITY_DN10_c0_g1_i1	GO:0003674,GO:0003824,GO:0004092,GO:0005575,GO:0005739

GO_mapping.txt (informative GO annotation)

biological_process	GO:0009058	biosynthetic process	245	The chemical reactions and pathways resulting in the formation of substances; typically the energy-requiring part of metabolism in which simpler substances are transformed into more complex ones. [GOC:cm]
biological_process	GO:0034641	cellular nitrogen compound metabolic process	229	The chemical reactions and pathways involving various organic and inorganic nitrogenous compounds, as carried out by individual cells. [GOC:mah]
biological_process	GO:0044281	small molecule metabolic process	188	The chemical reactions and pathways involving small molecules, any low molecular weight, monomeric, non-encoded molecule. [GOC:curators, GOC:pde, GOC:rw]
biological_process	GO:0006910	transport	147	The directed movement of substances (such as macromolecules, small molecules, ions) or cellular components (such as complexes and organelles) into, out of or within a cell, or between cells, or within a multicellular organism. [GOC:cm]
biological_process	GO:0006412	translation	93	The cellular metabolic process in which a protein is formed, using the sequence of a mature mRNA molecule to specify the sequence of amino acids in a polypeptide chain. Translation is mediated by the ribosome. [GOC:cm]
biological_process	GO:0009056	catabolic process	91	The chemical reactions and pathways resulting in the breakdown of substances, including the breakdown of carbon compounds with the liberation of energy for use by the cell or organism. [ISBN:019547688]
biological_process	GO:0022607	cellular component assembly	84	The aggregation, arrangement and bonding together of a cellular component. [GOC:cm, complete]
biological_process	GO:0006950	response to stress	67	Any process that results in a change in state or activity of a cell or an organism (in terms of movement, secretion, enzyme production, gene expression, etc.) as a result of a disturbance in organismal or cellular homeostasis. [GOC:cm]
biological_process	GO:0005975	carbohydrate metabolic process	61	The chemical reactions and pathways involving carbohydrates, any of a group of organic compounds based on the general formula C _x (H ₂ O) _y . Includes the formation of carbohydrate derivatives by the addition of other groups. [GOC:cm]
biological_process	GO:0065003	macromolecular complex assembly	61	The aggregation, arrangement and bonding together of a set of macromolecules to form a complex. [GOC:cm]
biological_process	GO:0006091	generation of precursor metabolites and cofactors	58	The chemical reactions and pathways resulting in the formation of precursor metabolites, substances from which energy is derived, and any process involved in the liberation of energy from these substances. [GOC:cm]
biological_process	GO:0051186	cofactor metabolic process	54	The chemical reactions and pathways involving a cofactor, a substance that is required for the activity of an enzyme or other protein. Cofactors may be inorganic, such as the metal atoms zinc, iron, and copper. [GOC:cm]

GO_barchart.png (according to GO_mapping.txt)



KEGG PATHWAY is a collection of manually drawn pathway maps representing our knowledge on the molecular interaction, reaction and relation networks for:

1. Metabolism

Global/overview, Carbohydrate, Energy, Lipid, Nucleotide, Amino acid, Other amino, Glycan, Cofactor/vitamin, Terpenoid/PK, Other secondary metabolite, Xenobiotics, Chemical structure

2. Genetic Information Processing

3. Environmental Information Processing

4. Cellular Processes

5. Organismal Systems

6. Human Diseases

7. Drug Development

Our KEGG result is generated from ec number data which is come from annotated GO terms.

ec2kegg.xls

[illegible]

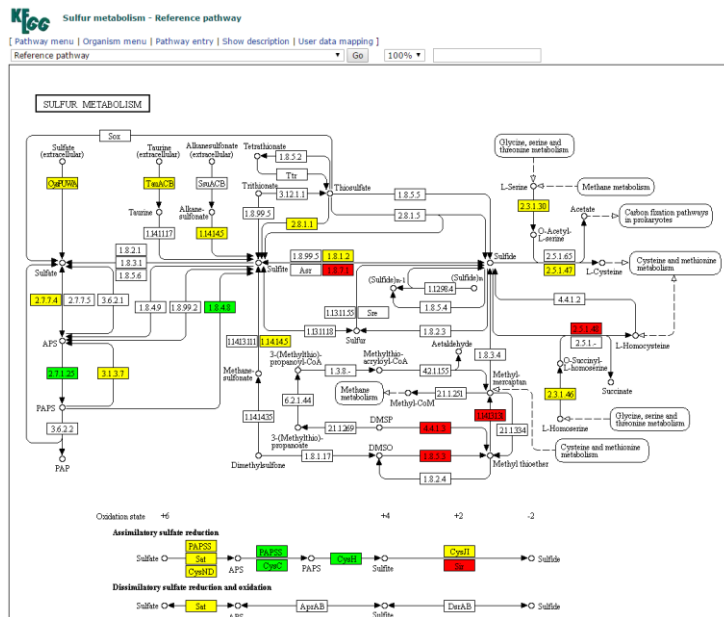
[Column definition]

- Total(EC_All) = number of ECs associated with the KEGG pathway;
- Total(EC_Ref(ead)) = number of ECs in reference genome ead (*E. adhaerens* OV14) associated with the KEGG pathway;
- Total(EC_Given) = number of tested ECs found to be associated with the KEGG pathway;
- Total(EC_Shared) = number of tested ECs that are shared with reference genome;
- Total(EC_Unique_Ref) = number of ECs that are unique to the reference genome;
- Total(EC_Unique_Given) = number of ECs that are unique to the tested genome.

$$(EC_Shared) + (EC_Unique_Given) = (EC_Given)$$

$$(EC_Shared) + (EC_Unique_Ref) = (EC_Ref(ead))$$

Click URL and get the pathway information



[Pathway map color definition]

green – an enzyme unique to a reference organism, (EC_Unique_Ref)

red – an enzyme unique to a given list, (EC_Unique_Given)

yellow – a shared enzyme. (EC_Shared)




[GO enrichment analysis]

One of the main uses of the GO is to perform enrichment analysis on gene sets. For example, given a set of genes that are up-regulated under certain conditions, an enrichment analysis will find which GO terms are over-represented (or under-represented) using annotations for that gene set.

In go enrichment analysis, we are using “Goseq v3.6” to finished this work. [8]

[3 GO enrichment dataset]:

e.g. <GSNO_1> v.s. <wt_1>: GSNO_1 is treatment & wt_1 is control

 GSNO_1_vs_wt_1.edgeR.DE_results.P0.05_C1.GSNO_1-UP.subset.GOseq.enriched ← up-regulated
 GSNO_1_vs_wt_1.edgeR.DE_results.P0.05_C1.wt_1-UP.subset.GOseq.enriched ← down-regulat
 GSNO_1_vs_wt_1.edgeR.DE_results.P0.05_C1.DE.subset.GOseq.enriched

[GSNO_1.UP.subset.GOseq.enrichment]

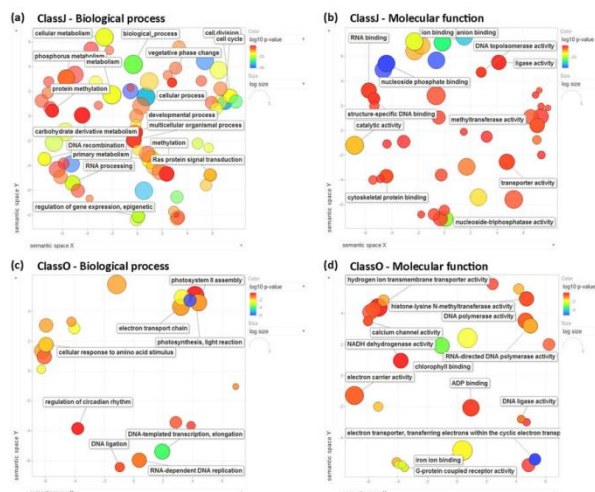
category	over_represented_pvalue	under_represented_pvalue	numDEInCat	numInCat	term	ontology	over_represented_FDR	go_term	gene_ids
GO:0003735	0	1	41	74	structur	MF	0	MF structur	TRINITY_DN105_c0_g2_i1,
GO:0005198	0	1	41	80	structur	MF	0	MF structur	TRINITY_DN105_c0_g2_i1,
GO:0006412	0	1	41	75	translat	BP	0	BP translati	TRINITY_DN105_c0_g2_i1,
GO:0006518	0	1	41	77	peptide	BP	0	BP peptide	TRINITY_DN105_c0_g2_i1,
GO:0009059	0	1	41	90	macron	BP	0	BP macron	TRINITY_DN105_c0_g2_i1,
GO:0019538	0	1	41	89	protein	BP	0	BP protein	TRINITY_DN105_c0_g2_i1,
GO:0030529	0	1	42	89	intracel	CC	0	CC intracel	TRINITY_DN105_c0_g2_i1,

Note:

- **Over-represented (enrichment):** lots of transcripts support certain GO term.
- **Under-represented (depletion):** few of transcripts could be found in certain GO term.
- **NumDEInCat:** number of searched DE transcripts matched with the GO term.
- **NumInCat:** number of total transcripts existed in the GO term.

If user would like to be more visualized your Gene Ontology terms which are derived from gene enrichment analysis, we recommend you this online tool – **REVIGO!** (<http://revigo.irb.hr>) [9]

You just need to copy above table red square columns (“category” and “over_represented_pvale”).



Reference graph:

Forestan C, Aiese Cigliano R, Farinati S, Lunardon A, Sanseverino W, Varotto S. Stress-induced and epigenetic-mediated maize transcriptome regulation by means of transcriptome reannotation and differential expression analysis. *Scientific Reports*. 2016;6:30446. doi:10.1038/srep30446.

Transcripts		Read quantitation										Differential expression									
transcript_id	length_x	Raw count	Raw count	Raw count	Raw count	Raw count	Raw count	FPKM (GS FPKM)	GS FPKM (GS FPKM)	wt_FPKM (wt_FPKM)	logFC (GSIlogFC)	GSIlogFC (GSIlogFC)	GSIlogFC (GSIlogFC)	GSipvalue (GSipvalue)	GSipvalue (GSipvalue)	GSipvalue (GSipvalue)					
TRINITY_253	0	1	2	2	1	2	0	2033.58	4803.43	4798.12	2431.95	6783.29	-	-	-	-					
TRINITY_174	1	0	0	0	0	0	5957.27	0	0	0	0	0	-	-	-	-					
TRINITY_277	0	2	1	2	1	2	0	3395.36	2012.7	4030.29	2049.21	5531.22	-	-	-	-					
TRINITY_568	64	57	42	24	30	17	3683.05	32161.16	28452.74	16392.26	21020.34	14176.85	1.336498	-2.04785	1.150901	2.922137					
TRINITY_194	0	3	2	3	2	2	0	11736.4	9044.33	13414.12	8906.38	14234.47	-	-	-	-					
TRINITY_214	7	7	11	2	1	21970.57	20920.69	38385.04	6931.09	3481.64	5250.01	1.674134	-1.95863	-2.39649	-2.37529	0.266667					
TRINITY_214	2	0	0	2	3	4	6327.31	0	6931.09	10444.92	21004.02	-0.06976	-0.19706	-0.19912	0.072874	1					
TRINITY_346	2	1	7	27	16	27	2365.78	1150.14	9605.29	37220.18	22528.26	48406.53	-3.74586	2.783709	2.674251	0.000921					
TRINITY_261	1	0	0	2	4	3	1981.04	0	4511.76	9158.66	9465.3	-0.98537	-1.12042	0.71261	-0.3346	1					
TRINITY_272	1	1	0	5	2	2	1822.43	1758.26	0	10423.55	4237.67	5753.99	-2.25674	-1.95863	1.978864	-2.37529					
TRINITY_164	0	0	1	1	1	0	0	0	7825.58	7688.88	7479.49	0	-	-	-	-					
TRINITY_211	1	0	0	2	3	1	3257.99	0	0	7174.53	10801.04	5471	-0.98537	-1.12042	-0.21081	0.985615					
TRINITY_175	0	0	0	0	1	2	0	0	0	6011.35	20692.79	-	-	-	-	-					
TRINITY_210	0	0	0	1	0	1	0	0	0	3629.69	0	5548.47	-	-	-	-					
TRINITY_154	0	0	0	0	1	0	0	0	0	9505.25	0	0	-	-	-	-					
TRINITY_226	3	2	3	7	6	8	8204.57	5226.42	9193.4	21353.51	18447.92	36050	-1.26103	-2.08051	1.055444	-1.07267					
TRINITY_193	1	1	2	0	0	0	4206.01	3972.54	9178.02	0	0	0	-	-	-	-					
TRINITY_353	8	3	4	17	16	20	9158.48	3341.06	5317.03	22706.18	21833.82	34608.24	-1.14895	-1.07916	-2.08267	1.816981					
TRINITY_264	2	7	1	0	0	0	3870.24	13052.49	2206.38	0	0	0	-	-	-	-					
TRINITY_185	0	0	0	2	1	2	0	0	10262.41	5082.67	16775.29	-	-	-	-	-					
TRINITY_769	45	43	39	1	1	1	17679.03	16603.66	18113.42	468.82	481.55	561.66	5.256032	-4.88565	4.87162	-5.25578					

BLASTP										BLASTX										Annotation									
UniprotKB	peptide	length	mismatch	gppow	qstaty	qslen	sstart	sends	evalue	klocscore	UniprotKB	peptide	length	mismatch	gppow	qstaty	qslen	sstart	sends	evalue	klocscore	Plan	SignalP	TmHMM	COGs (egg COs)	KEGGS	EC number		
											SDBB_CA	95.238	84	4	0	2	253	111	194	7.84E-55	172				COC3479	GO:005747	RO E0023	1.3.9.1,1.3.9.1	
											VATB_YE	96.491	57	2	0	3	173	430	436	7.47E-32	117						GO:001045	RO E0214	3.6.1,3.6.6,3.6.6
											VATB_YE	96.604	91	4	0	3	275	328	418	2.76E-56	184								
TCTP_CA	100	167	0	0	1	167	1	167	1.51E-120	339	TCTP_CA	100	167	0	0	59	559	1	167	2.82E-103	297				BNDG411	GO:001045			
											RLA3_YE	88.71	62	7	0	194	9	1	62	9.64E-33	111						GO:002265	RO E0294	
											COX12_YI	80.769	78	15	0	279	46	6	83	3.61E-47	149								
											GFP2_YE	80	85	17	0	4	258	95	179	2.15E-44	146								
											GFP1_YE	87.778	90	11	0	1	270	10	99	4.66E-53	168								
											FMP41_YI	58.571	70	29	0	2	211	27	96	2.41E-23	91.7						GO:008577		3
											BDH1_YE	51.724	58	28	0	1	174	76	133	1.60E-12	63.2						GO:008577	RO E0000	1.1.1.4.1.1,1.1.1.4.1.1
											BDH1_YE	69.565	69	21	0	3	209	11	79	3.32E-28	106								
											VATH_YE	74.51	51	13	0	2	154	309	359	4.41E-13	64.3								
											DIFL_ZYC	43.82	89	28	4	19	225	7	93	1.07E-12	61.2						GO:008577		
											MDM3S_Y	83.794	37	6	0	191	81	49	85	4.80E-15	65.9						GO:008577	RO E0227	1.9.3.11.1,1.9.3.11.1,1.9.3.11.1
											COX3_YE	53.866	78	27	2	111	323	1	76	2.50E-08	50.6						GO:000435		1.4.1.2.1.4,1.4.1.2.1.4,1.4.1.2.1.4
											DMBA_ZW	82.799	87	15	0	3	263	236	322	2.01E-44	151						GO:000435		
											LCPN_XY	60.465	43	17	0	57	185	29	121	8.80E-14	66.6						GO:001045	RO K0189	6.2.1.1.6.2,6.2.1.1.6.2,6.2.1.1.6.2

5. Reference

1. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D, Li B, Lieber M, Macmanes MD, Ott M, Orvis J, Pochet N, Strozzi F, Weeks N, Westerman R, William T, Dewey CN, Henschel R, Leduc RD, Friedman N, Regev A. De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. Nat Protoc. 2013 Aug;8(8):1494-512. Open Access in PMC doi: 10.1038/nprot.2013.084.
2. Limin Fu, Beifang Niu, Zhengwei Zhu, Sitao Wu and Weizhong Li, CD-HIT: accelerated for clustering the next generation sequencing data. Bioinformatics, (2012), 28 (23): 3150-3152. doi: 10.1093/bioinformatics/bts565.
3. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biol 10:R25.Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biol 10:R25.
4. Li, Bo & Dewey, Colin N. (2011). *RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome*, BioMed Central
5. McCarthy, J. D, Chen, Yunshun, Smyth and K. G (2012). "Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation." Nucleic Acids Research, 40(10), pp. -9.
6. Haas BJ, Papanicolaou A, Yassour M, et al. De novo transcript sequence reconstruction from RNA-Seq: reference generation and analysis with Trinity. Nature protocols. 2013;8(8):10.1038/nprot.2013.084. doi:10.1038/nprot.2013.084.
7. <https://trinitate.github.io/>
8. Gene ontology analysis for RNA-seq: accounting for selection bias Matthew D. Young, Matthew J. Wakefield, Gordon K. Smyth, Alicia Oshlack Genome Biology 2010, 11:R14 (4 February 2010)
9. Supek, F., Bosnjak, M., Skunca, N. & Smuc, T. REVIGO summarizes and visualizes long lists of gene ontology terms. PLoS One 6, e21800 (2011).