§ ゲノム勉強会

2020-08-04 - 2020-08-07

孫 建強

RNA-Seq

- Introduction
- **Quality control**
- O Mapping and quantification
- **O** Assembly

Central dogma



RNA-Seq





RNA-Seq



- Illumina HiSeq / NextSeq / MiSeq
- IonTorrent
- PacBio
- Oxford nanopore

RNA-Seq library preparation



RNA-Seq / single-end approach



RNA-Seq / paired-end approach



RNA-Seq

- Introduction
- **Quality control**
- O Mapping and quantification
- **O** Assembly

RNA-Seq



```
@HWI:2:2102:31421:12730
read sequence
                         TAAGGAGATCAAGAACGGGCGACTTGCGCTGTTGGCGTT
                          +
 Quality score
                         @@1BDDDDB??D:CGFIFEBGGABGFGFFFBABBBB?<>
                         @HWT:2:1102:21126:37630
                         GCCCCAGCCCAGTCCTAGCCCCAGCTCCAGTTCCGGCT
                          +
                         @HWI:2:1102:17396:37651
                         AAAAAAGGAACAAAAGAGGACAAGACCCAATCCACAAT
                          +
                         @@CFDFDFFBHHAGGFHIIGCHJJIGHIEGIJJJEECD:
                         @HWI:2:1102:17975:37502
                         TTGATGTTAGGCAAAGTCAAGAAGTTCTTGGTGATGTGA
                          +
                         CCCFFFFFHHHHHJJJJJJJJJJJJJJJJJJJJJJ
```

		51	3	71	G	91	Ε	111	0
		52	4	72	Н	92	\mathbf{i}	112	р
33	!	53	5	73	I	93]	113	q
34	**	54	6	74	J	94	٨	114	r
35	#	55	7	75	Κ	95	_	115	S
36	\$	56	8	76	L	96		116	t
37	%	57	9	77	Μ	97	a	117	u
38	&	58	:	78	Ν	98	b	118	V
39	T	59	,	79	0	99	С	119	W
40	(60	<	80	Ρ	100	d	120	Χ
41)	61	=	81	Q	101	e	121	У
42	*	62	>	82	R	102	f	122	Ζ
43	+	63	?	83	S	103	g	123	{
44	,	64	@	84	Т	104	h	124	
45	-	65	Α	85	U	105	i	125	}
46		66	В	86	V	106	j	126	~
47	/	67	С	87	W	107	k		
48	0	68	D	88	Х	108	1		
49	1	69	Е	89	Υ	109	m		
50	2	70	F	90	Ζ	110	n		

@HWI:2:2102:31421:12730 TAAGGAGATCAAGAACGGGGCGACTTGCGCTGTTGGCGTT + @@1BDDDDB??D:CGFIFEBGGABGFGFFFBABBBBB?<>

error rate p

$$p = 10^{\frac{Q}{-10}} = 10^{\frac{64-33}{-10}} = 7.9 \times 10^{-4}$$

13

		51	3	71	G	91	Ε	111	0
		52	4	72	Н	92	$\mathbf{\lambda}$	112	р
33	!	53	5	73	I	93]	113	q
34	**	54	6	74	J	94	٨	114	r
35	#	55	7	75	Κ	95	_	115	S
36	\$	56	8	76	L	96		116	t
37	%	57	9	77	Μ	97	a	117	u
38	&	58	:	78	Ν	98	b	118	V
39	T	59	,	79	0	99	С	119	W
40	(60	<	80	Ρ	100	d	120	X
41)	61	=	81	Q	101	е	121	У
42	*	62	>	82	R	102	f	122	Ζ
43	+	63	?	83	S	103	g	123	{
44	,	64	Q	84	Т	104	h	124	
45	-	65	Α	85	U	105	i	125	}
46		66	В	86	V	106	j	126	~
47	/	67	С	87	W	107	k		
48	0	68	D	88	Х	108	٦		
49	1	69	Е	89	Υ	109	m		
50	2	70	F	90	Ζ	110	n		

@HWI:2:2102:31421:12730 TAAGGAGATCAAGAACGGGGCGACTTGCGCTGTTGGCGTT + @@1BDDDDB??D:CGFIFEBGGABGFGFFFBABBBBB?<> Quality score = (ASCII code) - 33 = 16

error rate p

$$p = 10^{\frac{Q}{-10}} = 10^{\frac{49-33}{-10}} = 2.5 \times 10^{-2}$$

14

		51	3	71	G	91	Ε	111	0
		52	4	72	н	92	\mathbf{i}	112	р
33	!	53	5	73	I	93]	113	q
34	**	54	6	74	J	94	٨	114	r
35	#	55	7	75	Κ	95	_	115	S
36	\$	56	8	76	L	96		116	t
37	%	57	9	77	Μ	97	a	117	u
38	&	58	:	78	Ν	98	b	118	V
39	¥	59	,	79	0	99	С	119	W
40	(60	<	80	Ρ	100	d	120	Χ
41)	61	=	81	Q	101	е	121	У
42	*	62	>	82	R	102	f	122	Ζ
43	+	63	?	83	S	103	g	123	{
44	,	64	Q.	84	Т	104	h	124	
45	-	65	Α	85	U	105	i	125	}
46		66	В	86	V	106	j	126	~
47	/	67	С	87	W	107	k		
48	0	68	D	88	Χ	108	1		
49	1	69	Е	89	Υ	109	m		
50	2	70	F	90	Ζ	110	n		

@HWI:2:2102:31421:12730 TAAGGAGATCAAGAACGGGGCGACTTGCGCTGTTGGCGTT + @@1BDDDDB??D:CGFIFEBGGABGFGFFFBABBBBB?<> Quality score = (ASCII code) - 33 = 30

error rate p

$$p = 10^{\frac{Q}{-10}} = 10^{\frac{63-33}{-10}} = 1.0 \times 10^{-3}$$

15

Quality control (QC)



QC report



QC report / NG



QC report / before QC





QC report / after QC

paired-end reads



reversed reads



Position in read (bp)



MINLEN: 20 Drop the read if it is below a specified length

http://www.usadellab.org/cms/?page=trimmomatic

RNA-Seq

- Introduction
- **Quality control**
- O Mapping and quantification
- **O** Assembly



Genome sequence

CEnsembl

http://plants.ensembl.org/Arabidopsis_thaliana/Info/Index

Gene annotation

What can I find? Protein-coding and non-coding genes, splice variants, cDNA and protein sequences, non-coding RNAs.



More about this genebuild

Download genes, cDNAs, ncRNA, proteins - FASTA - GFF3

Update your old Ensembl IDs



Example gene



Arabidopsis_thaliana.TAIR10.dna.toplevel.fa.gz

Arabidopsis_thaliana.TAIR10.47.gff3.gz





reads

TCCCTTTTAGCCCCTTAG
CGGGGCTATCGAAATCGA
CCCCAGGCGGGGCTATCG
CCAGGCGGGGCTATCGAA
CCCTTTTAGCCCCTTAGT
GGCGGGGCTATCGAAATC
CCTTAGTAAGGGTCGCGC
GATCCCTTTTAGCCCCTT

mapping

reference GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCTTAGTAAGGGTCGCGCGAAAA

reads

TCCCTTTTAGCCCCTTAG CGGGGGCTATCGAAATCGA CCCCAGGCGGGGGCTATCG CCAGGCGGGGGCTATCGAA CCCTTTTAGCCCCTTAGT GGCGGGGGCTATCGAAATC CCTTAGTAAGGGTCGCGC

mapping

GATCCCTTTTAGCCCCTT

reference GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCCTTAGTAAGGGTCGCGCGAAAA

reads

TCCCTTTTAGCCCCTTAG CGGGGCTATCGAAATCGA CCCCAGGCGGGGGCTATCG CCAGGCGGGGGCTATCGAA CCCTTTTAGCCCCTTAGT GGCGGGGGCTATCGAAATC

mapping

CCTTAGTAAGGGTCGCGC GATCCCTTTTAGCCCCTT

reference GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCTTAGTAAGGGTCGCGCGAAAA

reads mapping CCCCAGGCGGGGCTATCG TCCCTTTTAGCCCCTTAG GGCGGGGCTATCGAAATC CCCTTTTAGCCCCTTAGT CGGGGCTATCGAAATCGA **CCTTAGTAAGGGTCGCGC** CCAGGCGGGGCTATCGAA GATCCCTTTTAGCCCCTT reference GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCTTAGTAAGGGTCGCGCGAAAA



Mapping software



(Bowtie) Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biology, 10:R25, 2009. Fast gapped-read alignment with Bowtie 2. Nature Methods, 9:357–359, 2012.

TopHat: Discovering Splice Junctions With RNA-Seq. Bioinformatics, 25(9):1105-11, 2009.

TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. Genome Biology, 14: R36, 2013.

HISAT: a fast spliced aligner with low memory requirements. Nature Methods, 12:357–360, 2015.

Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. Nature Biotechnology, 37:907–915, 2019.

Fast and accurate short read alignment with Burrows-Wheeler Transform. Bioinformatics, 25:1754-60, 2009.

Fast and accurate long-read alignment with Burrows-Wheeler Transform. Bioinformatics, 26:589-595, 2010.

STAR: ultrafast universal RNA-seq aligner. Bioinformatics, 29:1, 2013.

RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. BMC Bioinformatics, 12:323, 2011. (kallisto) Near-optimal probabilistic RNA-seq quantification. Nature Biotechnology, 34:525–527, 2016.

Salmon provides fast and bias-aware quantification of transcript expression. Nature Methods, 14,:417–419, 2017.

Mapping software

alignment-based



Alignment-based software align reads based on

- text search by FM-index using Burrows-Wheeler transformation text
- local alignment by Smith-Waterman algorithm

BWT

The ability to resolve the inherent complexity of gene families related to important agronomic traits demonstr ates the impact of IWGSC RefSeq v1.0 on dissecting quan titative traits genetically and implementing modern bre eding strategies for future wheat improvement.

Burrows-Wheeler transformation (BWT)

fCetyntsnesrfsedotegty0gsogsedyceqe1t.vSW Gf.I p rrcf tuelrrta ii aene oe hnhhvvrhsterRtrvldggrmmShdiitrnlooe nnn eaTttwnmtglmb ttd dtaaxlteppiiaoleeoa eiioii roaeeiiiiooaeeeattm snc rmfprmmmme ottttub e egpotete esei nnnacnriaaa encna ssiiuiiqtf lio eltt

- high compression ratio
- high-performance text searching

BWT



acaacg\$ 1234560 BWT



	F	L			F L			F	L
	\$ac??	'?g	\$ac????	g\$a????	\$a????g	\$a?????	g\$?????	\$?????	'g g
	aac??	?c	aac????	caa????	aa????c	aa?????	ca?????	a?????	°C C
	aca??	?\$	aca????	\$ac????	ac????\$	ac?????	\$a?????	a?????	\$\$
-	acg??	'?a ←	acg???? ←	_aac???? ←	- ac????a 🛶	_ac????? ←	- aa????? 🛶	-a?????	'a ← a
	caa??	?a	caa????	aca????	ca????a	ca?????	ac?????	c?????	'a a
	cg\$??	?a	cg\$????	acg????	cg????a	cg?????	ac?????	c?????	'a a
	g\$a??	?c	g\$a????	cg\$????	g\$????c	g\$?????	cg?????	g?????	°C C

- F-column equals to L-column sorted in ascended rank
- L-column equals to BW translated string when F-column is sorted in ascending rank
LF mapping is an algorithm mapping from the last column of the BWT to the first column.



LF mapping is an algorithm mapping from the last column of the BWT to the first column.



BWT



LF mapping is an algorithm mapping from the last column of the BWT to the first column.



LF mapping is an algorithm mapping from the last column of the BWT to the first column.



LF mapping is an algorithm mapping from the last column of the BWT to the first column.

LF(i) = C[L[i]] + ooc(L[i], 0, i)

C[s] index into F column where c begin

ooc(s, from, to)a function that calculates the
number of occurrences of
character s between from and to



LF mapping is an algorithm mapping from the last column of the BWT to the first column.

LF(i) = C[L[i]] + ooc(L[i], 0, i)

Compute LF(i) when i = 0, 1, 2, ..., 6

$$LF(0) = C[g] + ooc(g, 0, 0) = 6 + 1 = 7$$

$$LF(1) = C[c] + ooc(c, 0, 1) = 4 + 1 = 5$$

$$LF(2) = C[\$] + ooc(\$, 0, 2) = 0 + 1 = 1$$

$$LF(3) = C[a] + ooc(a, 0, 3) = 1 + 1 = 2$$

$$LF(4) = C[a] + ooc(a, 0, 4) = 1 + 2 = 3$$

$$LF(5) = C[a] + ooc(a, 0, 5) = 1 + 3 = 4$$

$$LF(6) = C[c] + ooc(c, 0, 6) = 4 + 2 = 6$$





LF mapping



\$



\$g







\$gca



\$gcaa



\$gcaac



\$gcaaca





suffix array





acaacg\$ BWT F \$acaacg aacg\$ac acaacg\$ acg\$aca caacg\$a cg\$acaa g\$acaac



















Query = aac



Query = aac



Query = aac

62

a atcg a С 0 0 0 0 0 0 0 0 0 a 0 a 0 C

Reference

$$F(i,j) = \max \begin{cases} 0 \\ F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$$

s(C, C) = 2 $s(C, ^C) = -2$ d = 1

63

		a	С	a	a	t	С	g
	0	0	0	0	0	0	0	0
a	0 -	12						
a	0							
С	0							

Reference

$$F(i,j) = \max \begin{cases} 0 \\ F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$$

s(C,C) = 2 $s(C,^C) = -2$ d = 1

64

Query



Reference

$$F(i,j) = \max \begin{cases} 0 \\ F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$$

s(C,C) = 2 $s(C,^C) = -2$ d = 1

65



Reference

$$F(i,j) = \max \begin{cases} 0 \\ F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$$

s(C,C) = 2 $s(C,^C) = -2$ d = 1

66

Query



Reference

$$F(i,j) = \max \begin{cases} 0 \\ F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$$

s(C,C) = 2 $s(C,^C) = -2$ d = 1

67



Reference

$$F(i,j) = \max \begin{cases} 0 \\ F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$$

s(C, C) = 2 $s(C, ^C) = -2$ d = 1

68



Reference

s(C, C) = 2 $s(C, ^C) = -2$ d = 1

69

Query



Reference

acaatcg || | aa-c

$$F(i,j) = \max \begin{cases} 0 \\ F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$$

s(C,C) = 2 $s(C,^C) = -2$ d = 1

70

Bowtie

backtracking search

double index





Bowtie2

- Similar to Bowtie, Bowtie2 uses BWT index and use 16 bases as a seed
- Extend seed to both ends with Smith-Waterman algorithm
 - allow gaps


HISAT2

- Alignment with FM index algorithm
 - global FM index
 - local FM index
 - generated from 64,000 bp subsequences of genome
 - overlap is 1024 bp between two subsequences

а

b

- covered whole genome
- Find alignemnt candidates with global FM index and build a correct alignment with local FM index



STAR

- STAR consists of two major steps: seed searching and clustering/stitching/scoring
- seed searching
 - search MMP (Maximal Mappable Prefix) from 5'-end
 - then, search MMP from the unmapped position of the given read
 - trim poly-A tails and poor quality tails
- clustering/stitching/scoring
 - intron size
 - mismatches
 - gaps



SAM format



RECORDS



SAM format

- FLAG 1 read paired 2 read mapped in proper pair 4 read unmapped 8 mate unmapped 16 read being reverse complemented 32 mate being reverse complemented 64 first in the template 128 last in the template 256 not primary alignment 512 not passing platform/vendor quality checks 1024 read is PCR or optical duplicate 2048 supplementary alignment
- I insertion to the reference D deletion from the reference N skipped region from the reference S soft clipping (clipped sequences present in SEQ) H hard clipping (clipped sequences NOT present in SEQ)

TAG AS alignment score X0 number of best hits XM number of mismatches in the alignment NH number of reported alignments that contain the query in the current record BAM



IGV



kallisto



Quantification





gene	count
LFY	2
STM	1
FT	2
EIN3	3

Gene annotation

- Gene annotations are recorded by GTF or GFF3 format
 - tab delimited text file
 - nine columns records chromosome name, feature, position, strand, gene ID, exon ID, and so on

C.*Ensembl*

http://plants.ensembl.org/Arabidopsis_thaliana/Info/Index

Gene annotation

What can I find? Protein-coding and non-coding genes, splice variants, cDNA and protein sequences, non-coding RNAs.



Download genes, cDNAs, ncRNA, proteins - FASTA - GFF3

Update your old Ensembl IDs



Example gene



Example transcript

Arabidopsis_thaliana.TAIR10.dna.toplevel.fa.gz

GTF (Gene Transfer Format)

chr1 at exon 9873504 9874841 . + . gene_id "G01"; transcript_id "T01"; gene_name "ZNF366"; chr1 at CDS 9873504 9874841 . + 0 gene_id "G01"; transcript_id "T01"; gene_name "ZNF366"; chr1 at exon 9877488 9877679 . + . gene_id "G01"; transcript_id "T01"; gene_name "ZNF366"; chr1 at CDS 9877488 9877679 . + 0 gene_id "G01"; transcript_id "T01"; gene_name "ZNF366"; chr1 at exon 9888412 9888586 . + . gene_id "G01"; transcript_id "T01"; gene_name "ZNF366"; chr1 at CDS 9888412 9888586 . + 0 gene_id "G01"; transcript_id "T01"; gene_name "ZNF366"; chr1 at CDS 9888412 9888586 . + 0 gene_id "G01"; transcript_id "T01"; gene_name "ZNF366"; chr1 at exon 9891475 9891998 . + . gene_id "G01"; transcript_id "T01"; gene_name "ZNF366"; chr1 at CDS 9891475 9891995 . + 2 gene_id "G01"; transcript_id "T01"; gene_name "ZNF366";

GTF (General Feature Format)

ctg123	at	mRNA	1300	9950	+		<pre>ID=t_012143;gene_name=EDEN</pre>
ctg123	at	exon	1300	1500	+		Parent=t_012143
ctg123	at	exon	3000	3902	+		Parent=t_012143
ctg123	at	CDS	3301	3902	+	0	Parent=t_012143
ctg123	at	exon	5000	5500	+		Parent=t_012143
ctg123	at	CDS	5000	5500	+	1	Parent=t_012143
ctg123	at	exon	7000	9000	+		Parent=t_012143
ctg123	at	CDS	7000	7600	+	1	Parent=t_012143
ctg123	at	exon	9400	9950	+		Parent=t_012143



RNA-Seq

- Introduction
- **Quality control**
- O Mapping and quantification
- **O** Assembly

Assembly

GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCTTAGTAAGGGTCGCGCGAAAA



sequencing

CCCCAGGCGGGGCTATCG GGCGGGGCTATCGAAATC CGGGGCTATCGAAATCGA CCAGGCGGGGCTATCGAA TCCCTTTTAGCCCCTTAG CCCTTTAGCCCCTTAGT CCTTAGTAAGGGTCGCGC GATCCCTTTTAGCCCCTT

GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCCTTAGTAAGGGTCGCGCGAAAA



sequencing

CCCCAGGCGGGGCTATCG GGCGGGGCTATCGAAATC CGGGGCTATCGAAATCGA CCAGGCGGGGCTATCGAA TCCCTTTTAGCCCCTTAG CCCTTTAGCCCCTTAGT CCTTAGTAAGGGTCGCGC GATCCCTTTTAGCCCCTT

assembly

CCCCAGGCGGGGCTATCG TCCCTTTAGCCCCTTAG GGCGGGGCTATCGAAATC CCCTTTAGCCCCTTAGT CGGGGCTATCGAAATCGA CCTTAGTAAGGGTCGCGC CCAGGCGGGGCTATCGAA GATCCCTTTTAGCCCCTT

GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCCTTAGTAAGGGTCGCGCGAAAA

Hamiltonian path problem

ATGATCCTAGACCCTGAT



ATGATCC CCTGAT CCTAGA GACCC





(c) wikipedia

de novo genome assembly

• Eulerian path









scaffold

NNNNNN 🛛

assembly quality

Table 1. Assembly statistics of IWGSC RefSeq v1.0.

Assembly characteristics	Values
Assembly size	14.5 Gb
Number of scaffolds	138,665
Size of assembly in scaffolds ≥ 100 kb	14.2 Gb
Number of scaffolds ≥ 100 kb	4,443
N50 contig length	51.8 kb
Contig L50 number	81,427
N90 contig length	11.7 kb
Contig L90 number	294,934
Largest contig	580.5 kb
Ns in contigs	0
N50 scaffold length	7.0 Mb
Scaffold L50 number	571
N90 scaffold length	1.2 Mb
Scaffold L90 number	2,390
Largest scaffold	45.8 Mb
Ns in scaffolds	261.9 Mb
Gaps filled with BAC sequences	183 (1.7 Mb)
Average size of inserted BAC sequence	9.5 kb
N50 superscaffold length	22.8 Mb
Superscaffold L50 number	166
N90 superscaffold length	4.1 Mb
Superscaffold L90 number	718
Largest superscaffold	165.9 Mb
Sequence assigned to chromosomes	14.1 Gb (96.8%)
Sequence \geq 100 kb assigned to chromosomes	14.1 Gb (99.1%)
Number of superscaffolds on chromosomes	1,601
Number of oriented superscaffolds	1,243
Length of oriented sequence	13.8 Gb (95%)
Length of oriented sequence \geq 100 kb	13.8 Gb (97.3%)
Smallest number of superscaffolds per subgenome chromosome	35 (7A), 68 (2B), 36 (1D)
Largest number of superscaffolds per subgenome chromosome	111 (4A), 176 (3B), 90 (3D)
Average number of superscaffolds per chromosome	76

genome -				
contias -	90k	60k	40k 38k	
oonago				
			N50 = 40k	
			L50 = 3	

95

reference-guided assembly



Expression analysis

O DE analysis

O Functional analysis

DE analysis





DE analysis



StringTie-Ballgown



99

https://doi.org/10.1038/nprot.2016.095

DE analysis



DE analysis / normalization



gene	ctrl 1	ctrl 2	ctrl 3	cold 1	cold 2	cold 3
FT	12	11	6	29	26	19
LFY	42	36	18	84	58	52
STM	3	4	1	17	33	23
PIN	32	23	14	42	38	25
EIN3	231	364	107	553	371	429
ERS1	96	103	45	182	91	85

Most genes are up-regulated by cold treatment?

DE analysis / normalization



gene	ctrl 1	ctrl 2	ctrl 3	cold 1	cold 2	cold 3
FT	12	11	6	29	26	19
LFY	42	36	18	84	58	52
STM	3	4	1	17	33	23
PIN	32	23	14	42	38	25
EIN3	231	364	107	553	371	429
ERS1	96	103	45	182	91	85
lib size	416	541	191	907	617	633

normalization

gene	ctrl 1	ctrl 2	ctrl 3	cold 1	cold 2	cold 3
FT	28.8	20.4	31.4	32.0	42.1	30.0
LFY	101.0	66.5	94.2	93.0	94.0	82.1
STM	7.2	7.4	5.3	18.7	53.5	36.3
PIN	76.9	42.5	73.3	45.3	61.6	39.5
EIN3	555.3	672.8	560.2	609.3	601.3	677.8
ERS1	230.8	190.4	235.6	201.7	147.5	134.3
lib size	1000.0	1000.0	1000.0	1000.0	1000.0	10000

MA plot

gene	ctrl 1	ctrl 2	ctrl 3	cold 1	cold 2	cold 3
FT	28.8	20.4	31.4	32.0	42.1	30.0
LFY	101.0	66.5	94.2	93.0	94.0	82.1
STM	7.2	7.4	5.3	18.7	53.5	36.3
PIN	76.9	42.5	73.3	45.3	61.6	39.5
EIN3	555.3	672.8	560.2	609.3	601.3	677.8
ERS1	230.8	190.4	235.6	201.7	147.5	134.3





gene	ctrl	cold
FT	26.9	34.0
LFY	87.2	89.7
STM	6.6	36.2
PIN	64.2	48.8
EIN3	642.3	629.5
ERS1	218.9	161.2

gene	mean	log ₂ FC
FT	30.5	0.34
LFY	88.5	0.04
STM	21.4	2.46
PIN	56.5	-0.396
EIN3	635.9	-0.029
ERS1	190.1	-0.441

DE analysis / count-based normalization



DE analysis / count-based normalization



DE analysis / variance



DE analysis / variance



DE analysis / variance


DE analysis / variance



DE analysis / variance



DE analysis / variance



DE analysis / dispersion



DE analysis / dispersion



DE analysis / statistical test



How do we compare the expression between groups?





 $E[ctrl] = \beta_1$ $E[cold] = \beta_2$

then, test the null hypothesis $\beta_1 = \beta_2$.

How do we compare the expression between groups?





 $E[ctrl] = \beta_1$ $E[cold] = \beta_1 + \beta_2$

then, test the null hypothesis $\beta_2 = 0$.

How do we compare the expression between groups?



Assume that

 $E[ctrl] = \beta_1$ $E[cold] = \beta_1 + \beta_2$

then, test the null hypothesis $\beta_2 = 0$. In this experiment, we have 6 replicates,

 $E[r1] = \beta_1$ $E[r2] = \beta_1$ $E[r3] = \beta_1$ $E[r4] = \beta_1 + \beta_2$ $E[r5] = \beta_1 + \beta_2$ $E[r6] = \beta_1 + \beta_2$

How do we compare the expression between groups?



Assume that

 $E[ctrl] = \beta_1$ $E[cold] = \beta_1 + \beta_2$

then, test the null hypothesis $\beta_2 = 0$. In this experiment, we have 6 replicates,

 $(E[r1] \quad E[r2] \quad E[r3] \quad E[r4] \quad E[r5] \quad E[r6]) = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$ design matrix

118

How do we compare the expression between groups?



To test null hypothesis $\beta_2 = 0$, we create two models and decide which one is better.



reduced model $(E[r1] \ E[r2] \ E[r3] \ E[r4] \ E[r5] \ E[r6]) = \begin{pmatrix} 1 & 0 \\ 1$

How do we compare the expression between groups?



Likelihood ratio test (LRT) and Wald test can

be used for comparing the two models. 120

Multiple testing corrections

If we test a null hypothesis with cutoff 0.05, the type I error rate should be less than

0.05

If we test two null hypotheses with cutoff 0.05 independently, the type I error rate should be less than

$$1 - (1 - 0.05)^2 = 0.0975$$



FDR



DE analysis



¹²³



- long genes are expected to have a large number of reads
- normalize read counts by both gene length and library size
 - 1. normalize gene length to 1 kbp
 - 2. normalize library size to 1 million

Normalization / count-based normalization



- gene length between control and cold treatment are same
- to identify differentially expressed genes between groups
 - normalize by library size is enough

note

Expression analysis

O DE analysis

O Functional analysis

AT1G62360 (STM)
gene ontology

- BP carpel development, cytokinin biosynthetic process, floral meristem determinacy, plasmodesmata-mediated intercellular transport, regulation of meristem structural organization, regulation of transcription by RNA polymerase II, stem cell population maintenance
- CC cytoplasm, cytoplasm, endosome, microtubule cytoskeleton, nucleus, plasma membrane, plasmodesma
- MF DNA-binding transcription factor activity, DNA-binding transcription factor activity, RNA polymerase II-specific, RNA binding, RNA polymerase II cis-regulatory region sequence-specific DNA binding, protein binding, protein homodimerization activity



QuickGO - https://www.ebi.ac.uk/QuickGO

в

в

в

в

B

в

в

в

What functions are these DEGs related to?



What functions are these DEGs related to?



GO:0010582 GO:0009723 GO:0009415 GO:0008150

What functions are these DEGs related to?



What functions are these DEGs related to?



	DEGs	nonDEGs	total
with GO:0010582	k	m - k	m
without GO:0010582	n - k	N - m - n + k	N - m
total	n	N - n	Ν

Fisher's exact test

$$P(X = k) = \frac{\binom{m}{k}\binom{N-m}{n-k}}{\binom{N}{n}}$$
 132

GSEA (gene set enrichment analysis)



Focus on DEGs annotated with specific GO terms.

Focus on all genes ranked by the statistics, and check whether the gene is annotated with specific GO terms or not from the top untill the bottom of the rank.

GWAS

- Introduction
- O Mapping
- O Variant calling

Genomic variants



Ancestor REF GGGAACCCCAGGGCTATCGAAATCGATCCCAA

Cardamine amara inhabits in water



- A1 GGGAACCGCAGGGCAATCGAATCGACCCCAA
- A2 GGGAACCGCAGGGCAATCGAAATCGACCCCAA
- A3 GGGAACCCCAGGGCAATCGAAATCGACCCCAA
- A4 GGGAACCCCAGGGCAATCGAAATCGACCCCAA
- A5 GGGAACCCCAGGGCAATCGAAATCGACCCCAA

Cardamine rivularis inhabits in glassland



R1 GGGAACCCCAAGGCTATCGTAATCGACCCCAA R2 GGGAACCCCAAGGCTATCGTAATCGACCCCAA R3 GGGAACCCCAAGGCTATCGTAATCGACCCCAA R4 GGGAACCCCAGGGCTATCGTAATCGACCCCAA A5 GGGAACCCCAGGGCTATCGTAATCGACCCCAA

contribute water-stress tolerance?

GWAS



REF	GGGAACCCCAGGGCTATCGAAATCGATCCCAA
XI-1A.1	GGGAACCGCAGGGCAATCGAAATCGACCCCAA
XI-1A.2	GGGAACCGCAGGGCAATCGAAATCGACCCCAA
XI-18.1	GGCAACCCCAGGGCAATCGAAATCGACCCCAA
XI-18.2	GGCAACCCCAGGGCAATCGAAATCGACCCCAA
XI-2.1	GGGAACCCCAAGGGTATCGAAA-CGATCGCAA
XI-2.2	GGGAACCCCAAGGCTATCGTAA-CGATCGCAA
XI-2.3	GGGAACCCCAAGGCTATCGTAA-CGATCGCAA
:	GGGAACCCCAAGGCTATCGTAATCGATCGCAA
:	GGGAACCCCAGGGCTATCGTAATCGATCCCAA
Admix.3	GGGAACCCCAGGGCTATCGTAATCGATCCCAA

Fig. 1: Unweighted neighbour-joining tree based on 3,010 samples and computed on a simple matching distance matrix for filtered SNPs. Wang, W., Mauleon, R., Hu, Z. et al. Genomic variation in 3,010 diverse accessions of Asian cultivated rice. Nature 557, 43-49 (2018). https://doi.org/10.1038/s41586-018-0063-9

■ XI-1A ■ XI-1B XI-2 XI-3 ■ XI-adm □ GJ-adm GJ-trp GJ-sbtrp GJ-tmp 🗖 cA ■ cB Admix

GWAS

- Introduction
- O Mapping
- O Variant calling

Variant calling



reads

TCCCTTTTAGCCCCTTAG
CGGGGCTATCGAAATCGA
CCCCAGGCGGGGCTATCG
CCAGGCGGGGCTATCGAA
CCCTTTTAGCCCCTTAGT
GGCGGGGCTATCGAAATC
CCTTAGTAAGGGTCGCGC
GATCCCTTTTAGCCCCTT

mapping

reference GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCTTAGTAAGGGTCGCGCGAAAA

reads

TCCCTTTTAGCCCCTTAG CGGGGGCTATCGAAATCGA CCCCAGGCGGGGGCTATCG CCAGGCGGGGGCTATCGAA CCCTTTTAGCCCCTTAGT GGCGGGGGCTATCGAAATC CCTTAGTAAGGGTCGCGC

mapping

GATCCCTTTTAGCCCCTT

reference GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCCTTAGTAAGGGTCGCGCGAAAA

reads

TCCCTTTTAGCCCCTTAG CGGGGCTATCGAAATCGA CCCCAGGCGGGGGCTATCG CCAGGCGGGGGCTATCGAA CCCTTTTAGCCCCTTAGT GGCGGGGGCTATCGAAATC

mapping

CCTTAGTAAGGGTCGCGC GATCCCTTTTAGCCCCTT

reference GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCTTAGTAAGGGTCGCGCGAAAA

		TCCCTTTTAGCCCCTTAG
		CGGGGCTATCGAAATCGA
		CCCCAGGGGGGCTATCG
reads	ads	CCAGGGGGGCTATCGAA
reado		CCCTTTTAGCCCCTTAGT
		GGCGGGGCTATCGAAATC
		CCTTAGTAAGGGTCGCGC
		GATCCCTTTTAGCCCCTT
	mapping	
		CCCCAGGCGGGGCTATCG TCCCTTTAGCCCCTTAG
		GGCGGGGCTATCGAAATC CCCTTTTAGCCCCTTAGT
_	L	CGGGGCTATCGAAATCGA CCTTAGTAAGGGTCGCGC
	•	CCAGGCGGGGCTATCGAA GATCCCTTTTAGCCCCTT
.		

reference GGGAGCACCCCAGGCGGGGCTATCGAAATCGATCCCTTTTAGCCCCTTAGTAAGGGTCGCGCGAAAA

STAR

- STAR consists of two major steps: seed searching and clustering/stitching/scoring
- seed searching
 - search MMP (Maximal Mappable Prefix) from 5'-end
 - then, search MMP from the unmapped position of the given read
 - trim poly-A tails and poor-quality tails
- clustering/stitching/scoring
 - intron size
 - mismatches
 - gaps


BWA

- There are three alignment modes
 - short reads (< 100 bp)
 - BWA-backtrack
 - long reads (70 < 1M bp)

0 googol\$

1 oogol\$g

2 ogol\$go

3 gol\$goo

4 ol\$goog

5 1\$googo

6 \$googol

X = googol

Pos

- BWA-SW
- BWA-MEM



SAM format

HEADER

@HD VN:1.0
@SD SN:1 LN:30427671
@SD SN:2 LN:19698289
@SD SN:3 LN:23459830
@RG ID:MAK1 PL:ILLUMINA LB:IL200 SM:Col-0
@PG ID:bwa PN:bwa VN:0.6.2

RECORDS



SAM format

- FLAG 1 read paired 2 read mapped in proper pair 4 read unmapped 8 mate unmapped 16 read being reverse complemented 32 mate being reverse complemented 64 first in the template 128 last in the template 256 not primary alignment 512 not passing platform/vendor quality checks 1024 read is PCR or optical duplicate 2048 supplementary alignment
- CIGAR M alignment match (can be a sequence match or mismatch) I insertion to the reference D deletion from the reference N skipped region from the reference S soft clipping (clipped sequences present in SEQ) H hard clipping (clipped sequences NOT present in SEQ)

TAG AS alignment score X0 number of best hits XM number of mismatches in the alignment NH number of reported alignments that contain the query in the current record **147**

GWAS

- Introduction
- O Mapping
- **O** Variant calling
- Ο



pre-processing / mark duplicates



pre-processing / local realignment







.bam mark duplicates local realignment **BQSR**

Base Quality Score Recalibration

- Build a correction model with non-variants loci
- Use @RG metadata to correct the reported quality score
- BQSR can performed iteratively until covergence





We need variants information before BQSR!

- We can download these information from database
- We can provisionally perform variant calling





.bam variant calling genotyping variant selection variant filtration

.vct

- Define active region
- Determine haplotypes by assembly of the active region
 - local realignment of atcive regions with Smith-Waterman algorithm
 - identify possible haplotypes
- Determine likelihoods of the haplotypes given the read data
 - pairwise alignment of each read against each applotype with PairHMM algorithm
 - ouput likelihoods of alleles for each potential variants
- Assign sample genotypes
 - calculate likehoods of each genotype per sample with Bayes 58







.vcf





GATK outputs five major types of variant

- SNP: single nucleotide polymorphism
- INDEL: insetion and deletion
- MIXED: combination of SNPs and indels at a single position
- MNP: multi-nucleotidde polymorphism
- SYMBOLIC



VCF format

##fileformat=VCFv4.0													
##fileDate=20090805													
##source=myImputationProgramV3.1													
##reference=1000GenomesPilot-NCBI36													
##phasing=partial													
##INFO= <id=ns.number=1.type=integer.description="number data"="" of="" samples="" with=""></id=ns.number=1.type=integer.description="number>													
##INFO= <id=dp.number=1.type=integer.description="total depth"=""></id=dp.number=1.type=integer.description="total>													
##INFO= <id=af.number=type=float.description="allele frequency"=""></id=af.number=type=float.description="allele>													
##INFO= <id=aa.number=1.type=string.description="ancestral allele"=""></id=aa.number=1.type=string.description="ancestral>													
##INFO= <id=db.number=0.type=flag.description="dbsnp 129"="" build="" membership.=""></id=db.number=0.type=flag.description="dbsnp>													
##INFO= <id=h2.number=0.tvpe=flag.description="hapmap2 membership"=""></id=h2.number=0.tvpe=flag.description="hapmap2>													
##FILTER= <id=q10.description="ouality 10"="" below=""></id=q10.description="ouality>													
##FILTER= <id=s50.description="less 50%="" data"="" have="" of="" samples="" than=""></id=s50.description="less>													
##FORMAT= <id=gt.number=1.type=string.description="genotype"></id=gt.number=1.type=string.description="genotype">													
##FORMAT= <id=gq,number=1,type=integer,description="genotype quality"=""></id=gq,number=1,type=integer,description="genotype>													
##FORMAT= <id=dp.number=1.type=integer.description="read depth"=""></id=dp.number=1.type=integer.description="read>													
##FORMAT= <id=ho.number=2.type=integer.description="haplotype quality"=""></id=ho.number=2.type=integer.description="haplotype>													
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001	NA00002	NA00003		
20	14370		G	Α	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:51.51	1 0:48:8:51.51	1/1:43:5:		
20	17330		т	Α	3	q10	NS=3;DP=11;AF=0.017	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3	0/0:41:3:.,.		
20	1110696		Α	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	2/2:35:4:		
20	1230237		т		47	PASS	NS=3;DP=13;AA=T	GT:GO:DP:HO	0 0:54:7:56.60	0 0:48:4:51.51	0/0:61:2:		
20	1234567		GTCT	G,GTACT	50	PASS	NS=3; DP=9; AA=G	GT:GQ:DP	0/1:35:4	0/2:17:2	1/1:40:3		
:				-			· ·	•	-	-	-		



IGV



🗯 IGV File	Genomes View Tracks Regions Tools Help						
Human hg19	Load Genome from File Load Genome from URL Load Genome From Server						
1	Remove genomes Create .genome File Create genome File						



	Create .genome file					
Unique identifier	tair10					
Descriptive name	tair10					
FASTA file	/Users/jsun/Desktop/igv_tair10/Arabidopsis_thaliana.TAIR10.dna.toplevel.fa Browse					
Optional						
Cytoband file	Browse					
Gene file	/Users/jsun/Desktop/igv_tair10/Arabidopsis_thaliana.TAIR10.40.gff3 Browse					
Alias file	fill description and set sequence and Browse					
	анносацон mes 2 ок 16е					

IGV



.bam	tair10			
.bai		1		5
	SRR094979.bam Coverage SRR094979.bam		Zoom in to see coverage. Zoom in to see alignments.	
	\$99005608 bar Causage		Zoom in to see severage	
select BAM files.	SRR095698.bam		Zoom in to see alignments.	
IGV File Genomes View Tracks Regions Tools Help Load from File Load from URL				
Load from Server	10 🗘	Q Search	Zoom in to see coverage.	
Favorites Name New Session Sove Session Open Session Arabidopsis_thaliana.TAIR10 Save Session Arabidopsis_thaliana.TAIR10 Reload Session Desktop Severation Severation Save Session Desktop Severation Severation Severation Severation	40.gff3 dna.toplevel.fa dna.toplevel.fa.fai cf.gz f.gz	A Size Kinu 111.2 MB Document 121.7 MB Document 176 bytes Document 21.2 MB gzip coarcl 97.6 MB gzip coarcl 1.18 GB Document 2020 KD Document	Zoom in to see alignments.	
Save Image		332 KB Document 1.38 GB Document 347 KB Document 984.8 MB Document		an an a faith and a second
Exit SRR095786.bam.bai		326 KB Document		1,109M of 3,267M
New Folder		Cancel Open		168

IGV

.bam .bai









The pair reads that are coded by the chromosome on which their mates can be found

