

RNAseq analysis

Overview

- Data QC
- File formats
- Aligning sequence reads
- Viewing data on a genome
- Calculate expression levels

Pipeline RNAseq

Tool	Function(s)
FastQC	QC
Bowtie/TopHat	Splice-aware alignment
Samtools	Manipulate alignment (SAM/BAM) files
IGV/UCSC	Data viewing
htseq-count	Calculate gene expression

Fastq - sequence format

Fastq format – fasta with qualities

- p = the probability that the corresponding base call is wrong
- Qualities
 - $p = 0.1 \quad Q = 10$
 - $p = 0.01 \quad Q = 20$
 - $P = 0.001 \quad Q = 30$
- Encoding: Sanger/Phred format can encode a quality score from 0 to 93
- using ASCII 33 to 126: $Q + 33$ - ASCII code

$$Q_{\text{sanger}} = -10 \log_{10} p$$

Dec	Hx	Oct	Hex	Chr	Dec	Hx	Oct	Hex	Chr
32	20	040	4#32;	Space	64	40	100	c#64;	0
33	21	041	4#33;	!	65	41	101	c#65;	A
34	22	042	4#34;	"	66	42	102	c#66;	B
35	23	043	4#35;	#	67	43	103	c#67;	C
36	24	044	4#36;	\$	68	44	104	c#68;	D
37	25	045	4#37;	%	69	45	105	c#69;	E
38	26	046	4#38;	&	70	46	106	c#70;	F
39	27	047	4#39;	*	71	47	107	c#71;	G
40	28	050	4#40;	{	72	48	110	c#72;	H
41	29	051	4#41;	}	73	49	111	c#73;	I
42	2A	052	4#42;	*	74	4A	112	c#74;	J
43	2B	053	4#43;	+	75	4B	113	c#75;	K
44	2C	054	4#44;	,	76	4C	114	c#76;	L
45	2D	055	4#45;	-	77	4D	115	c#77;	M
46	2E	056	4#46;	.	78	4E	116	c#78;	N
47	2F	057	4#47;	/	79	4F	117	c#79;	O
48	30	060	4#48;	0	80	50	120	c#80;	P
49	31	061	4#49;	1	81	51	121	c#81;	Q
50	32	062	4#50;	2	82	52	122	c#82;	R
51	33	063	4#51;	3	83	53	123	c#83;	S
52	34	064	4#52;	4	84	54	124	c#84;	T
53	35	065	4#53;	5	85	55	125	c#85;	U
54	36	066	4#54;	6	86	56	126	c#86;	V
55	37	067	4#55;	7	87	57	127	c#87;	W
56	38	070	4#56;	8	88	58	130	c#88;	X
57	39	071	4#57;	9	89	59	131	c#89;	Y
58	3A	072	4#58;	:	90	5A	132	c#90;	Z
59	3B	073	4#59;	:	91	5B	133	c#91;	[
60	3C	074	4#60;	<	92	5C	134	c#92;	\
61	3D	075	4#61;	=	93	5D	135	c#93;]
62	3E	076	4#62;	>	94	5E	136	c#94;	^
63	3F	077	4#63;	?	95	5F	137	c#95;	_

Source: http://en.wikipedia.org/wiki/FASTQ_format

Illumina sequence identifiers

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
```

EAS139	the unique instrument name
136	the run id
FC706VJ	the flowcell id
2	flowcell lane
2104	tile number within the flowcell lane
15343	'x'-coordinate of the cluster within the tile
197393	'y'-coordinate of the cluster within the tile
1	the member of a pair, 1 or 2 (<i>paired-end or mate-pair reads only</i>)
Y	Y if the read is filtered, N otherwise
18	0 when none of the control bits are on, otherwise it is an even number
ATCACG	index sequence

Run/read QC

- First thing when receive a FASTQ file to check:
 - data amount
 - data quality
 - get first idea of whether sequenced what were aiming to sequencing
- Influences downstream e.g. checking that there are no non-genomic sequences in reads that would prevent mapping to reference
- Extremely important for feedback upstream to wet-lab

Potential QC issues

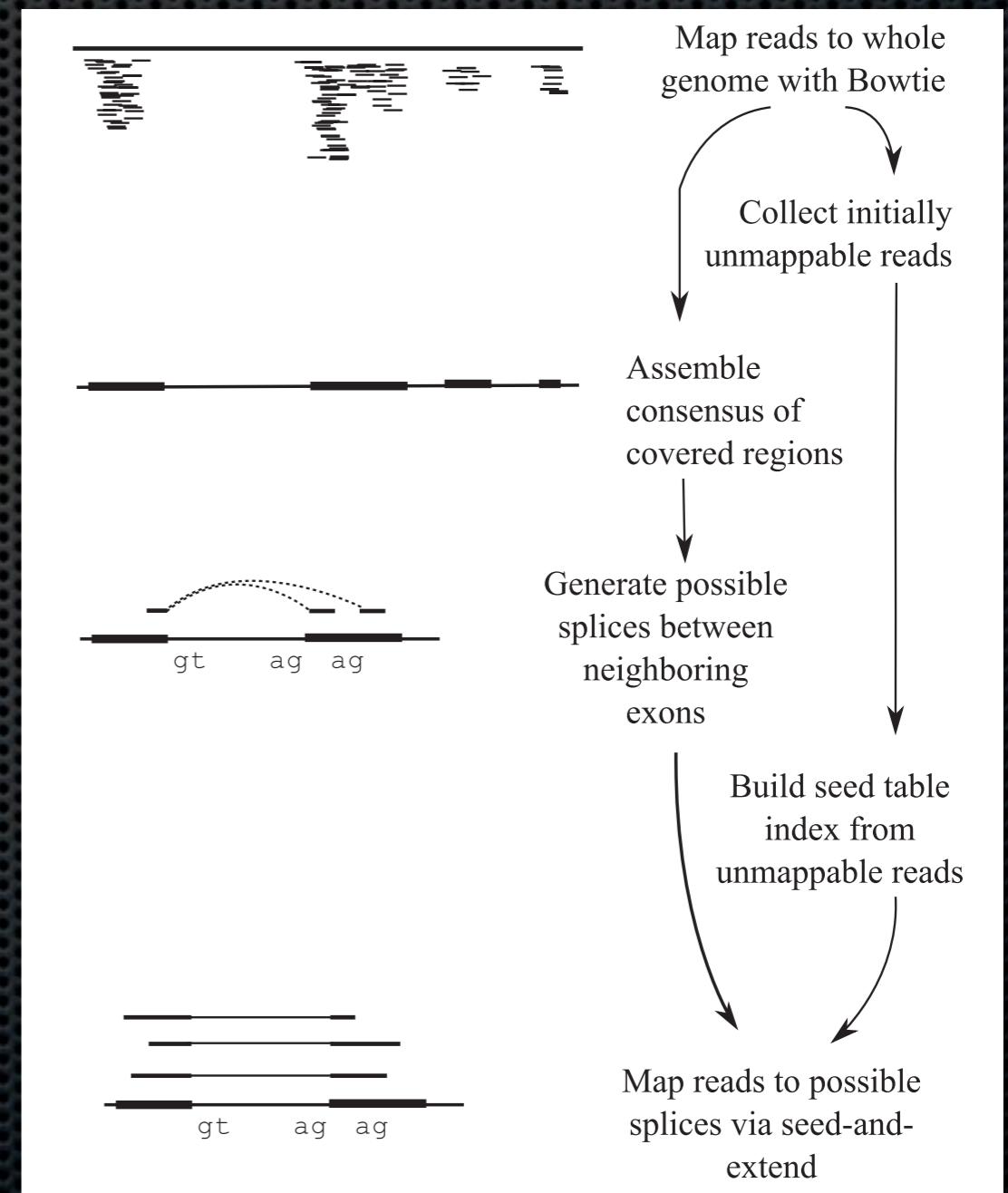
- Numbers of reads (no adaptors?, too little/much DNA loaded?, imaging problem?)
- Qualities (reagent problems?, imaging?)
- by cycle: fall with increasing cycle (Illumina)
- by read: usually bad bases distributed across reads
- Base distribution
- strong skewness at particular cycles (indicative of adaptors present)
- skewness across all cycles (may be indicative of not sequencing targeted region)
- Levels of duplication

FastQC

- <http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/>
- Don't take all warnings at face value
- Has both a GUI and command line interface.
- Used by NSC to generate our standard QC report
- Example...

Aligning RNAseq reads to genome

- Bowtie - fast short read aligner
- TopHat
 - Split reads
 - Align
 - Predict splicing
 - Realign



Commands I

```
# Make directories for course
```

```
mkdir ref
```

```
mkdir reads
```

```
mkdir results
```

```
# Files
```

```
reads/sample_1_hsa21_R1.fastq, reads/sample_1_hsa21_R2.fastq
```

```
ref/chr21.fa
```

```
ref/Homo_sapiens_chr21.GRCh37.64.gtf
```

```
# Prepare bowtie index
```

```
bowtie-build chr21.fa chr21
```

```
# Run tophat to do alignment
```

```
tophat -r 50 -G ref/Homo_sapiens_chr21.GRCh37.64.gtf --no-novel-juncs \  
--library-type=fr-unstranded -p 1 -o results \  
ref/chr21 reads/sample_1_hsa21_R1.fastq reads/sample_1_hsa21_R2.fastq
```

Commands II

```
# View coverage
```

```
genomeCoverageBed -split -bg -ibam accepted_hits_np.bam > accepted_hits_coverage.bed
```

```
# Sort hits, and make sam file
```

```
samtools sort -n accepted_hits.bam accepted_hits_n
```

```
samtools view -h accepted_hits_n.bam > accepted_hits_n.sam
```

```
# View results: UCSC, IGV
```

<http://genome.ucsc.edu/>

<http://www.broadinstitute.org/igv/>

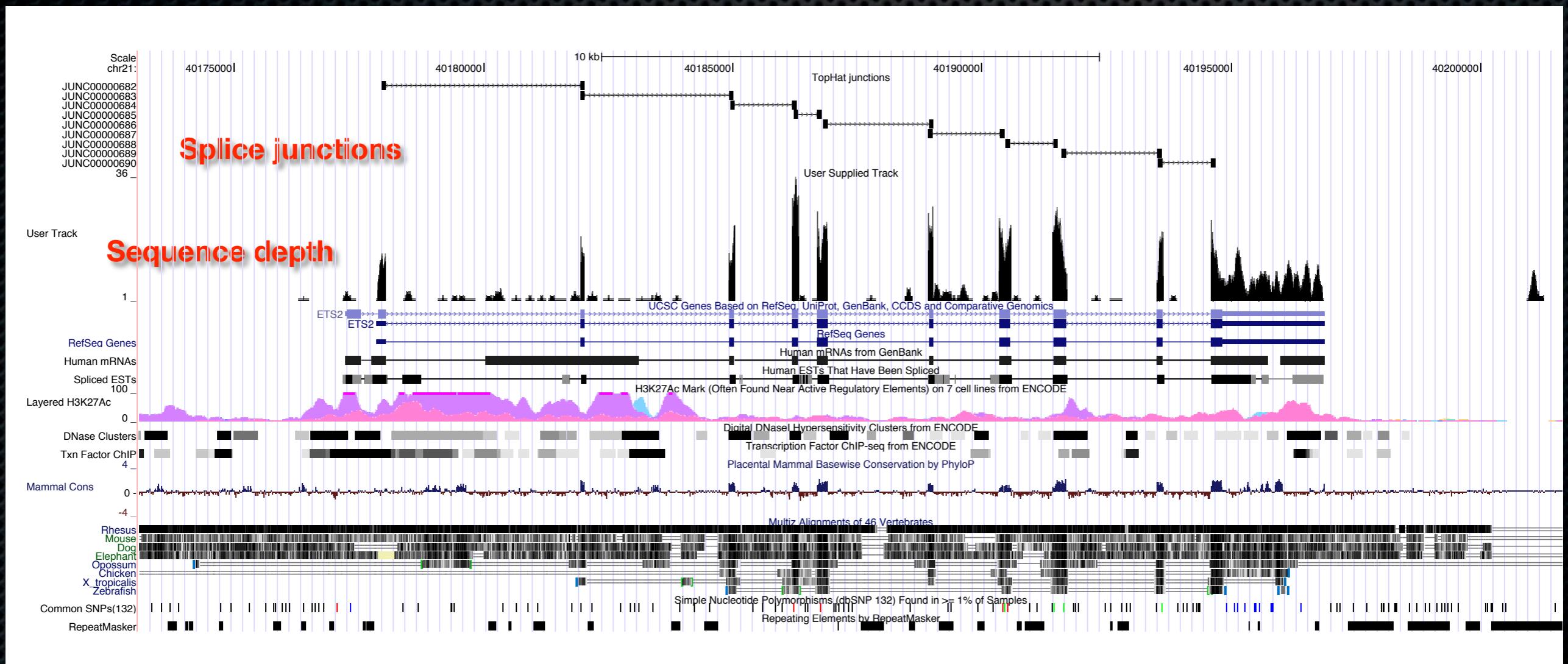
```
# Count expression levels
```

```
htseq-count -m intersection-strict --stranded=no -t exon -i gene_id \  
results/accepted_hits_n.sam ref/Homo_sapiens_chr21.GRCh37.64.gtf \  
> results/sample_1_hsa21.count
```

Viewing data

- UCSC
- IGV

Viewing RNAseq data



Expression levels

- Count expression levels from alignment
 - HTSeq package
- Differential expression
 - R
 - DEseq (and others)
 - Cufflinks

Questions?

11TH-GRADE ACTIVITIES:

USEFULNESS
TO CAREER
SUCCESS

