

CHIP-seq and ATAC-seq for Epigenomics Profiling and Functional Analysis

Yurii Chinenov, PhD

David Oliver, PhD

Hospital for Special Surgery
Genomics Center



Presentation Breakdown

- **Part 1: Sequencing Technology**
 - Sequencing
 - ChIP-seq/ATAC-seq
- **Part 2: Pipelines and Downstream Analysis**
 - Sequencer to Data
 - Data to Meaning

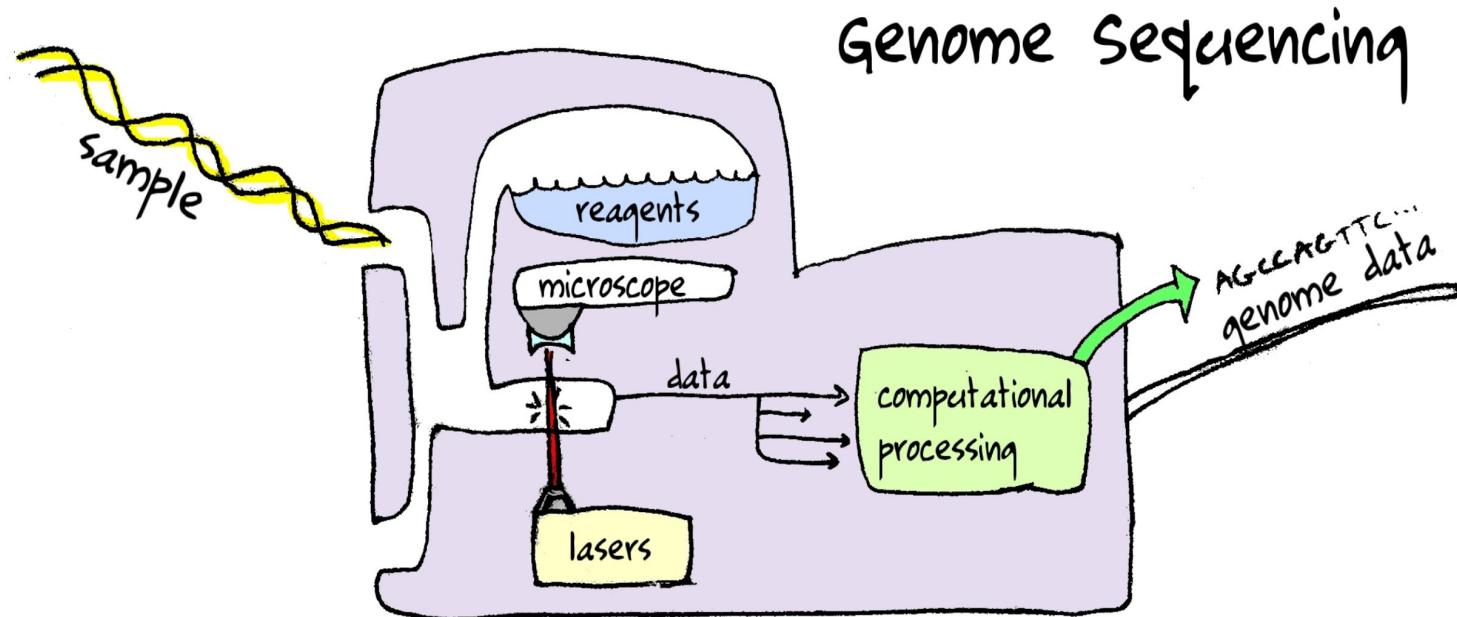
Sequencing Technology



Sequencing Terminology

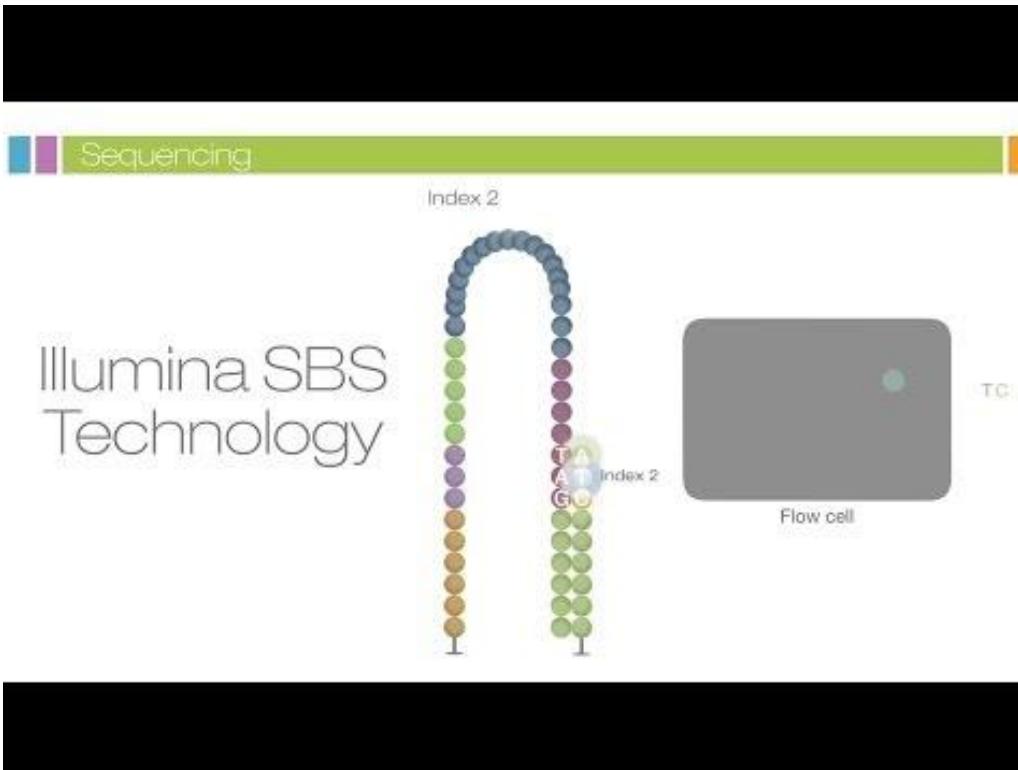
- **Sample** - The raw RNA or DNA to be sequenced
- **Library** - The **sample** after it has been prepared for sequencing (a selected portion of your sample)
- **Flow Cell** - The disposable portion of the sequencer where the **library** is sequenced
- **Cluster** - A single molecule from the **library** amplified 1000 times to create one ‘spot’ on the **flow cell**
- **Reads** - The base calls for each **cluster**.
- **Quality** - The per-base probability that each base in the **read** was called correctly

Short Read Sequencing





Short Read Sequencing By Synthesis





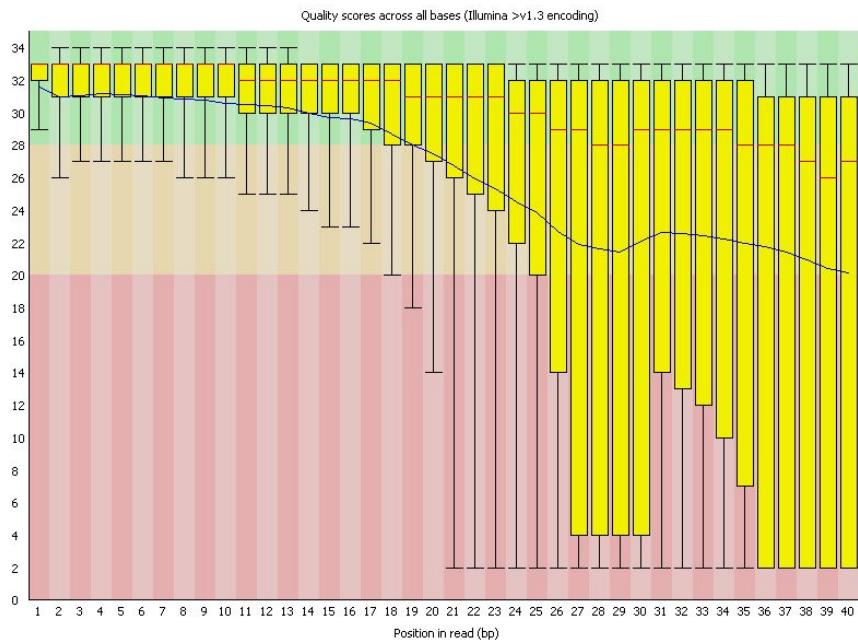
The Cost of Sequencing (Illumina)

- Quality and length of sequencing has increased
- Sequencer chemistry and technology has improved
- Cost of sequencing has decreased dramatically

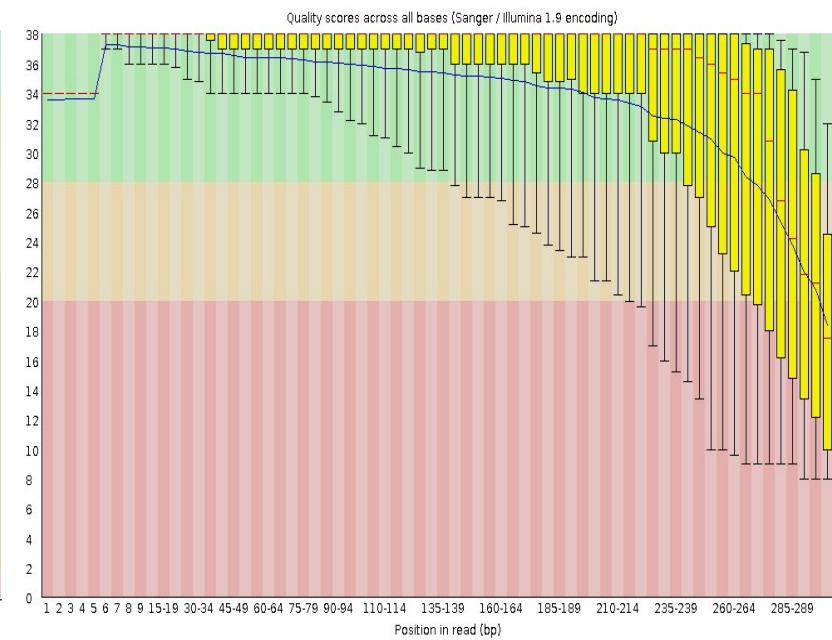


Quality of sequencing has increased

Circa GAllx

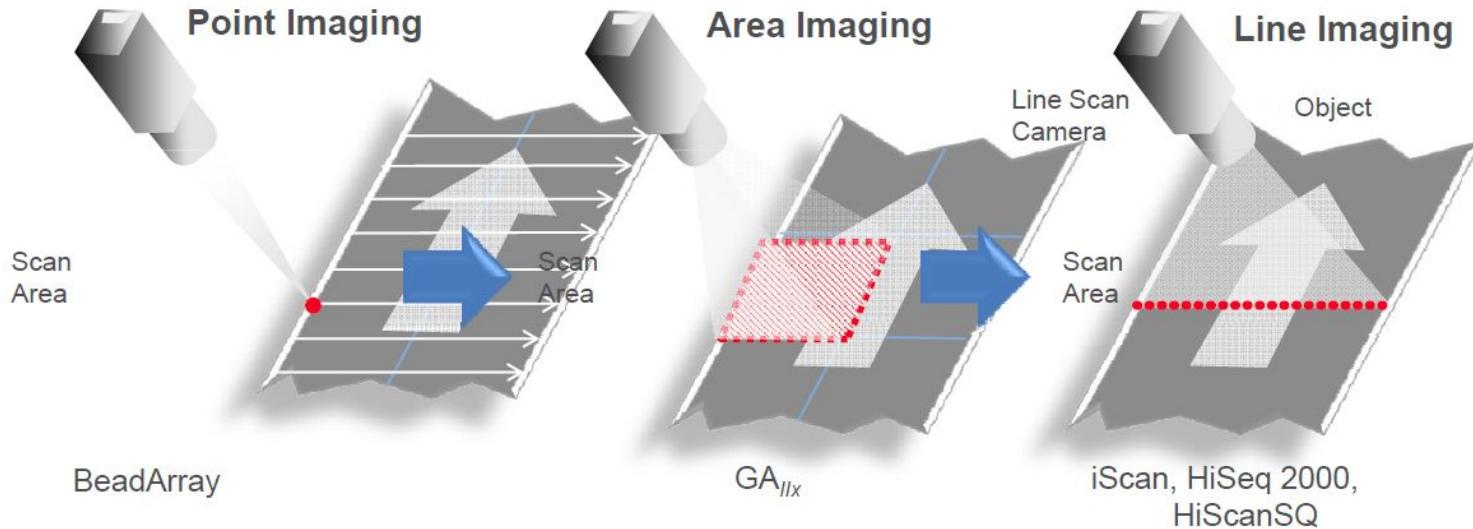


Circa MiSeq



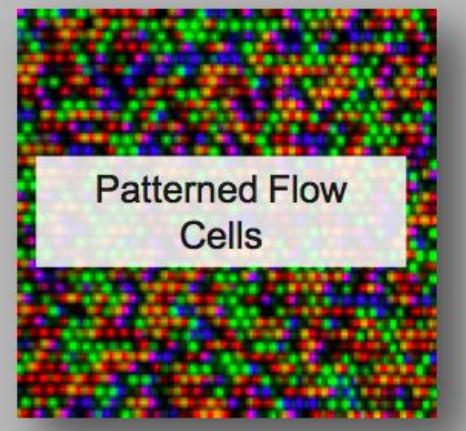
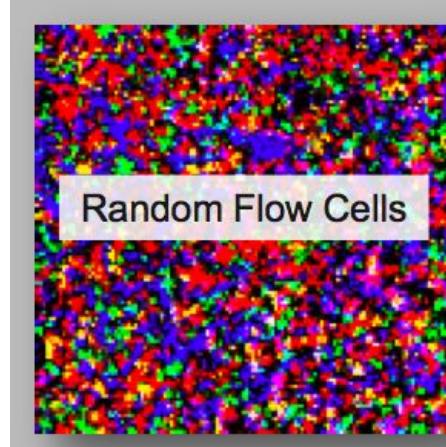
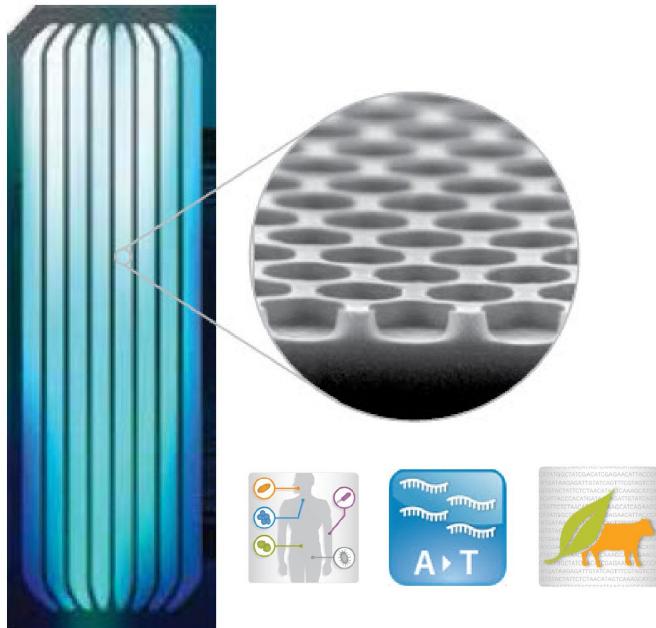


Sequencer technology has improved



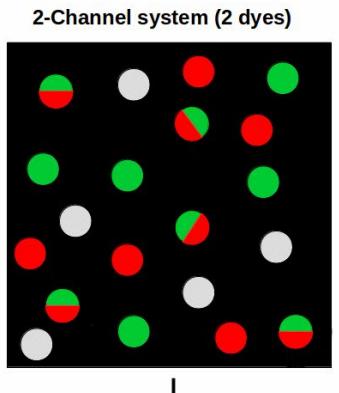
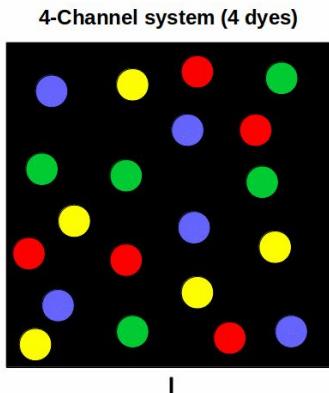


Sequencer technology has improved





Sequencer Chemistry has improved



4 Filter channels

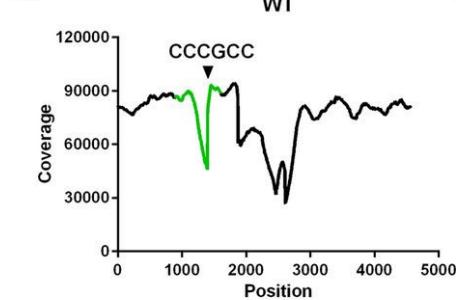
2 Filter channels

T C A G

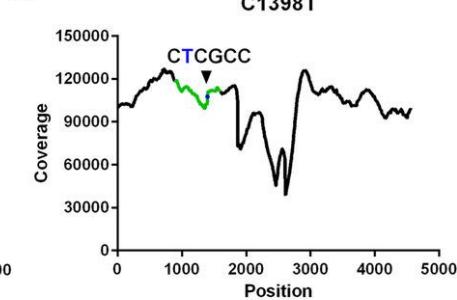
T C A G*

*No detected dye

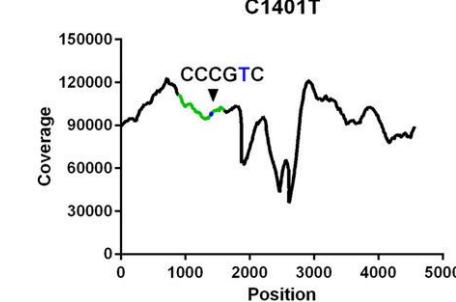
a.



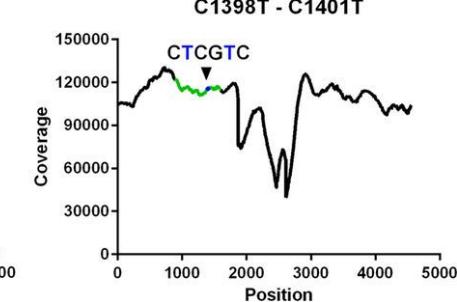
b.



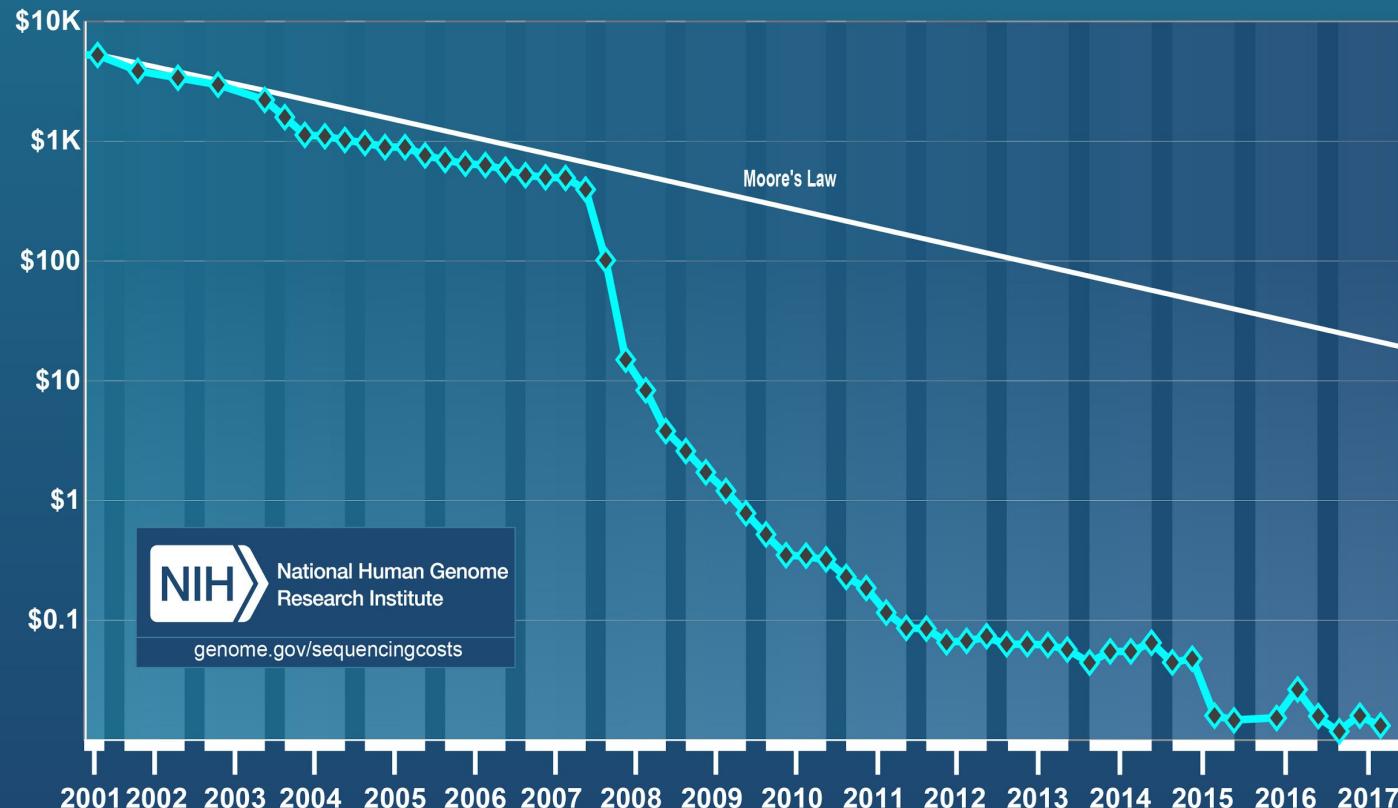
c.



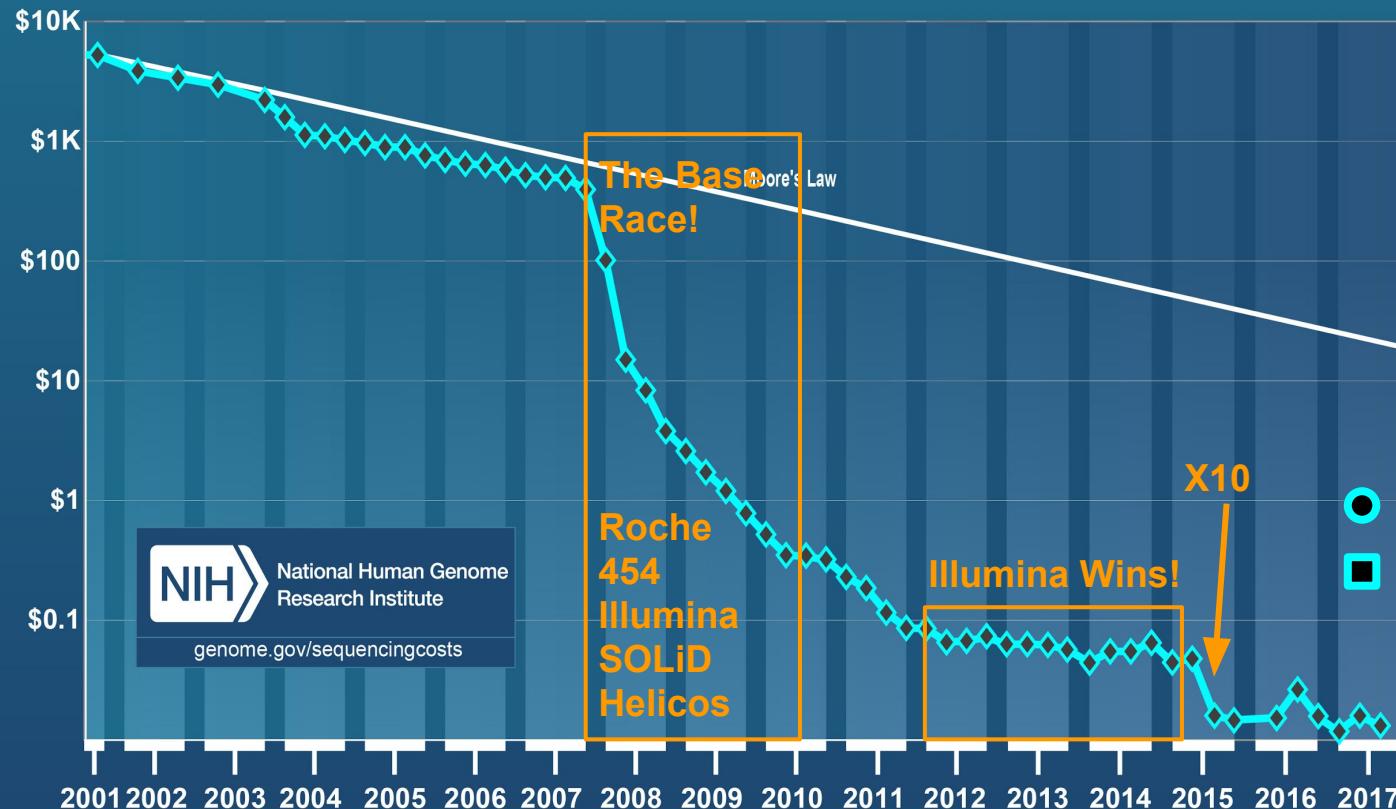
d.



Cost per Raw Megabase of DNA Sequence

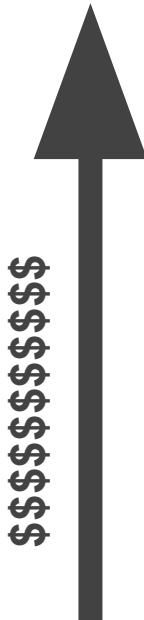


Cost per Raw Megabase of DNA Sequence

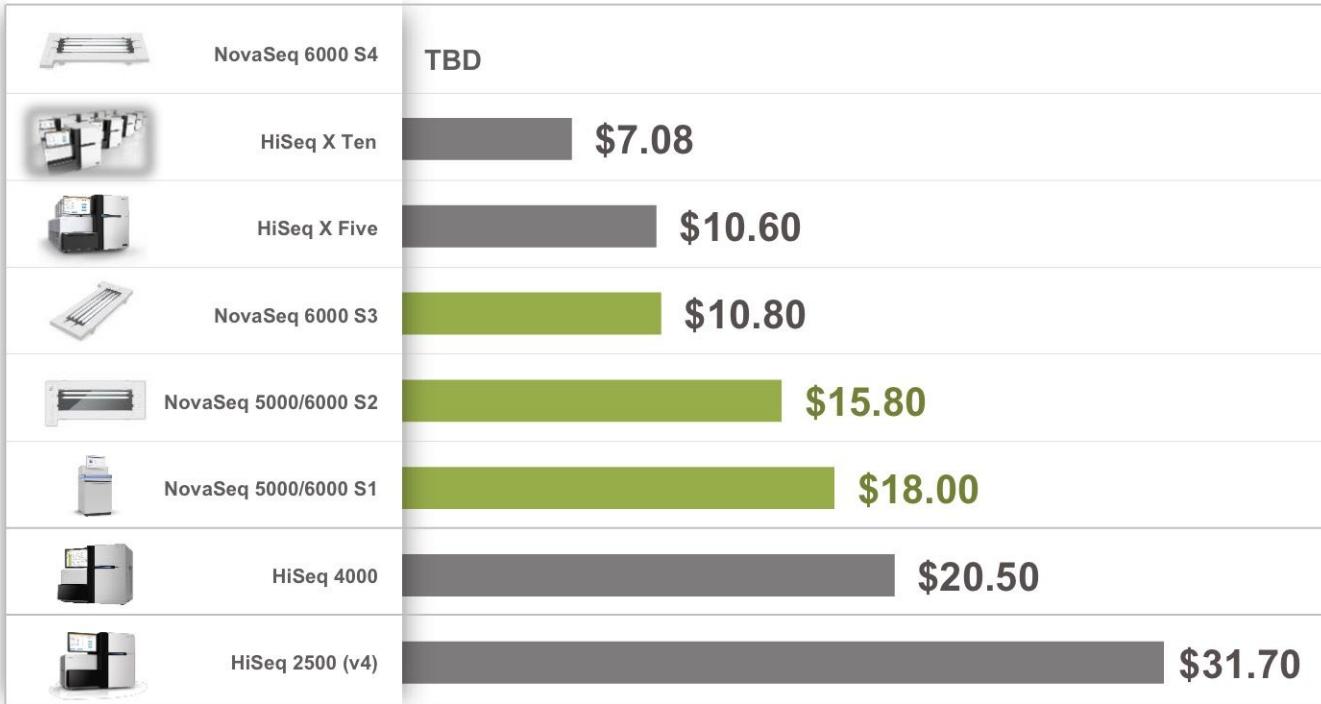




With Great Quality Comes Great Cost

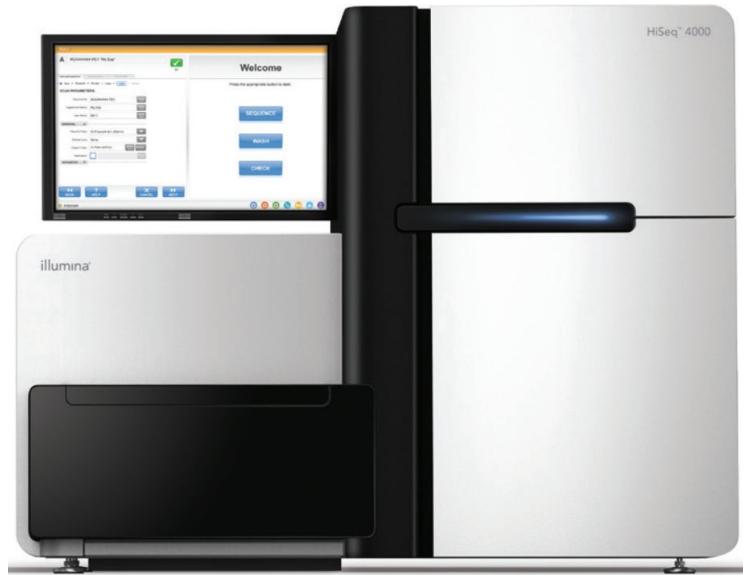


List Price per Gb





Current Illumina Machines





Non-Illumina Sequencing for ChIP/ATAC





Sequencing Experiment Design

- Experimental Considerations
 - Special needs
 - Is an input control necessary
 - How many replicates?
- Library Preparation
 - How long should fragments be?
 - Paired or single end?
 - How long should my reads be?
- Sequencing Considerations
 - What platform should be used
 - How deep should I sequence?



Experimental Considerations: Special Needs

- Are there special adapters?
 - Non-standard adapters require special library prep and a sequencing facility that's willing to accommodate your needs.
- Is the sample non- "standard" material?
 - Too short/Too long ($50 > \text{length} < 500$)
 - Low Quality (e.g. ancient DNA)
 - GC rich/Highly repetitive
- How much starting material is there?
 - Is there enough material for the assay?



Experimental Considerations: Library Preparation

- How long should fragments be?
 - ChIP-seq -- 150-200bp
 - ATAC-seq -- 100-800bp
- Paired or single end?
 - ChIP-seq -- Single End
 - ATAC-seq -- Paired End
- How long should my reads be?
 - > 50bp

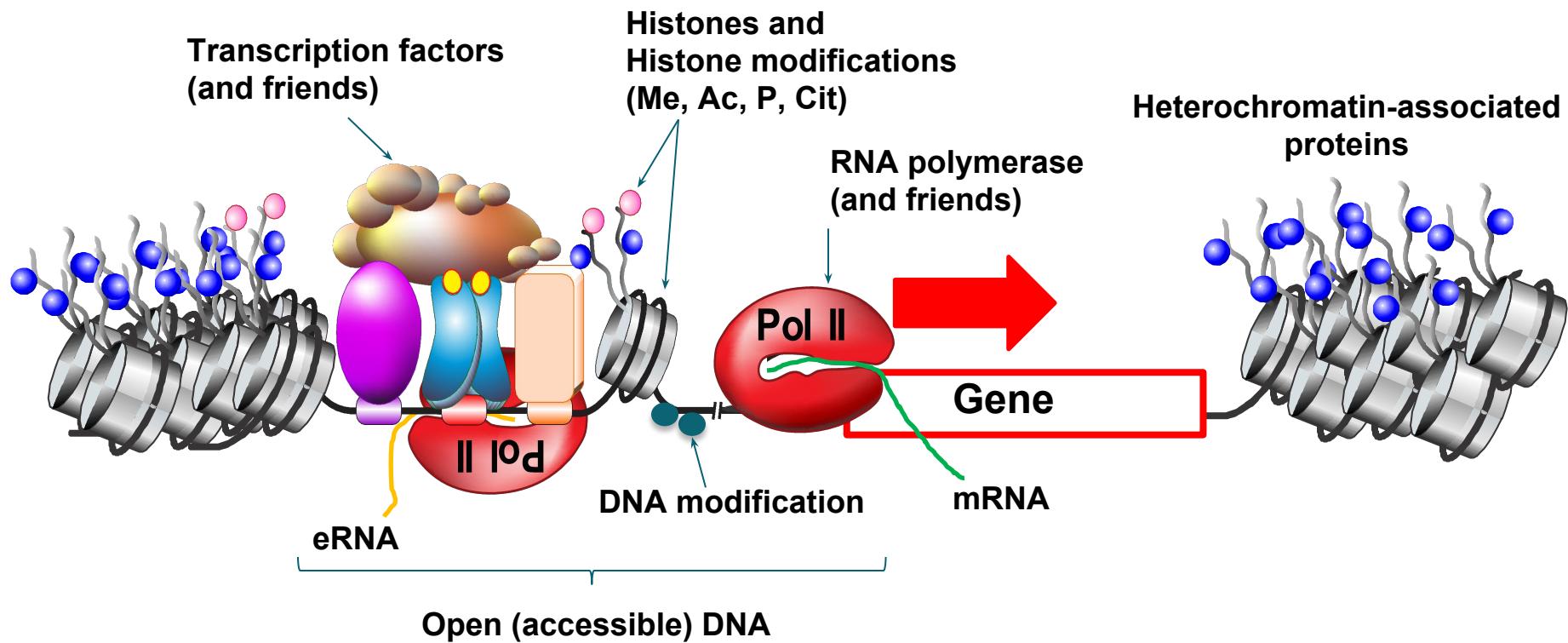


Experimental Considerations: Sequencing

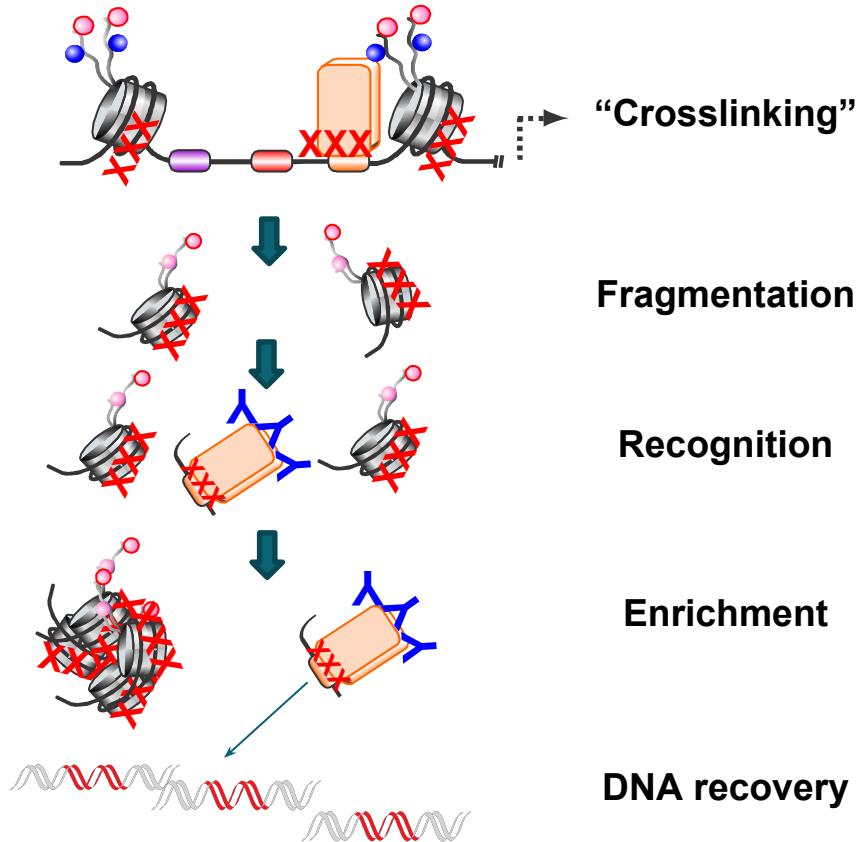
- What platform should I use?
 - For ChIP/ATAC-seq there is no need to use anything other than Illumina.
 - Other NGS based questions may require a specific platform
- How deep should I sequence?
 - ChIP - 15-25M single end reads
 - ATAC - 50-200M paired end reads
 - Use QC steps to determine if you need more.

Applications of Sequencing Technology: ChIP-seq/ATAC-seq

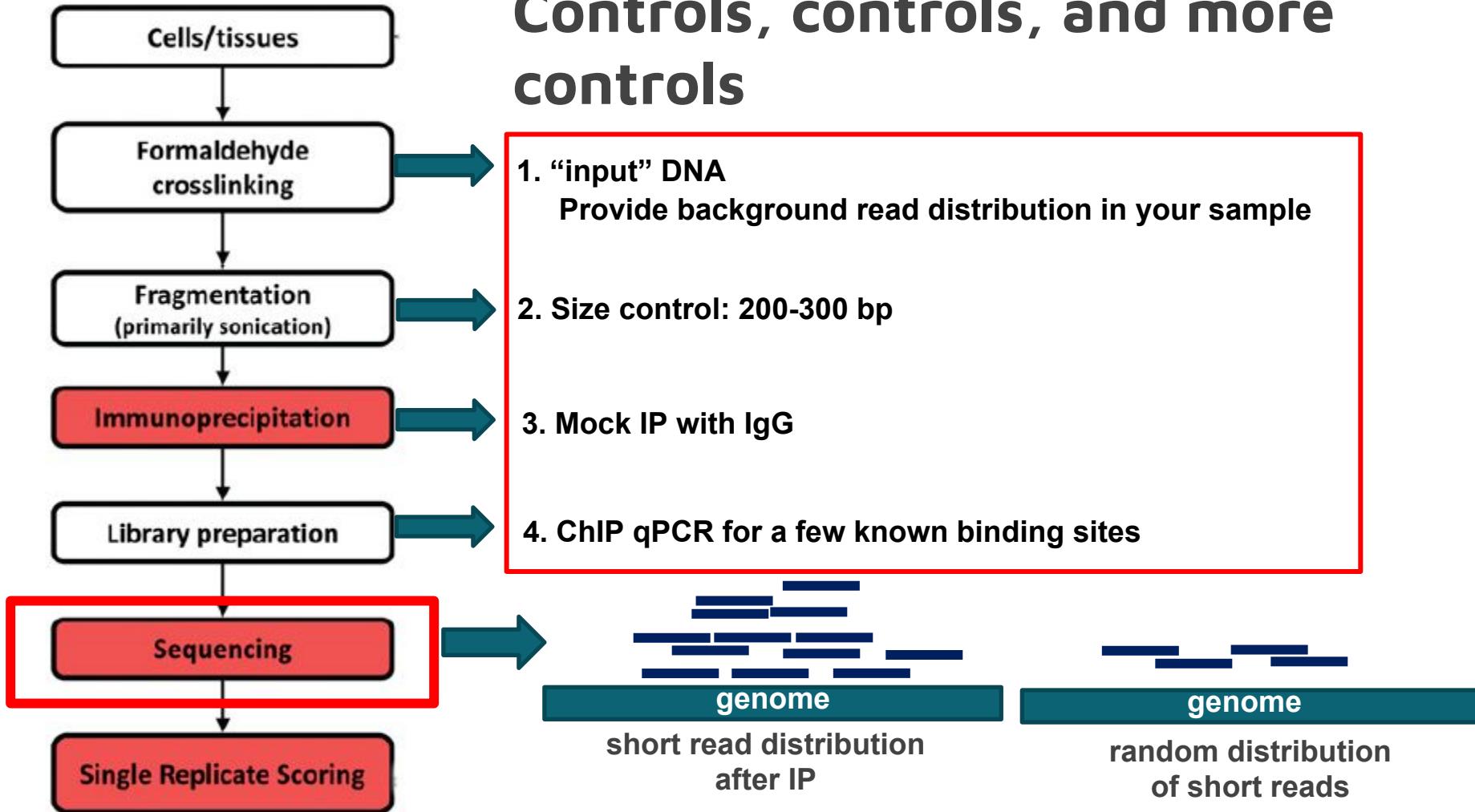
Identifying signals in chromatin



ChIP-seq: Identifying signals in chromatin



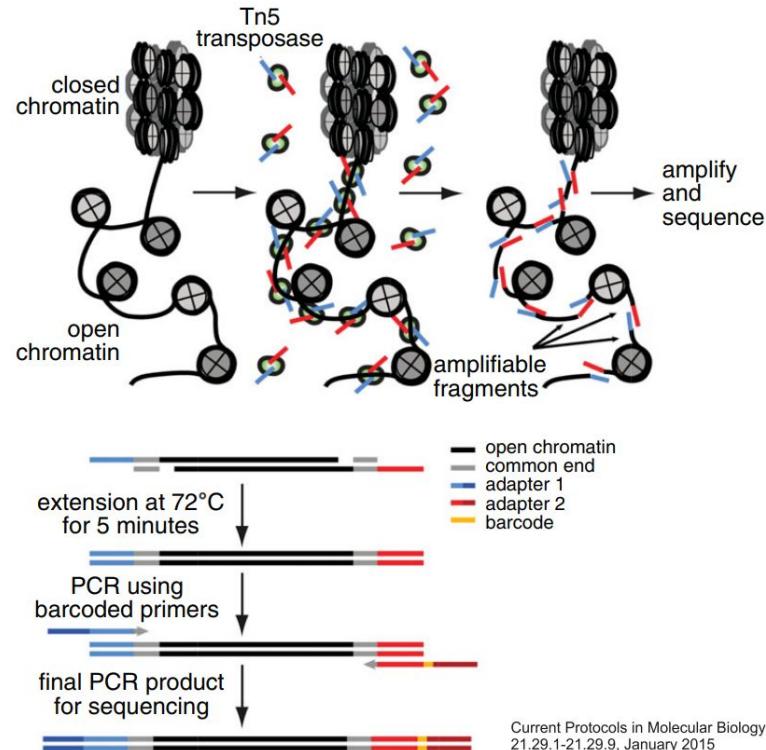
Controls, controls, and more controls





ATAC-seq

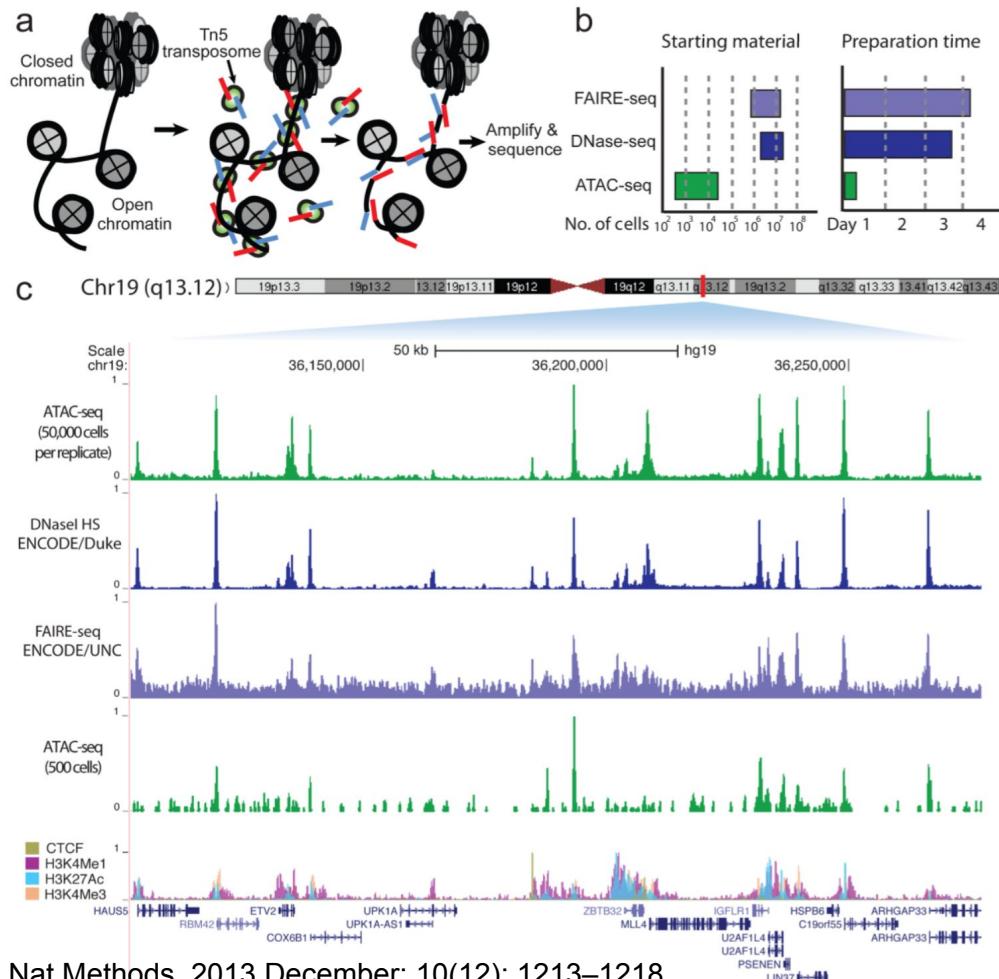
- Maps “open” chromatin genome-wide
- DNA Fragmentation / incorporation of sequencing adaptors in one step





Why ATAC-seq?

- Advantages:
 - More information
 - Less time
 - Less material
 - No sonication!
 - No need for “good quality antibody”
- Disadvantages
 - More Information





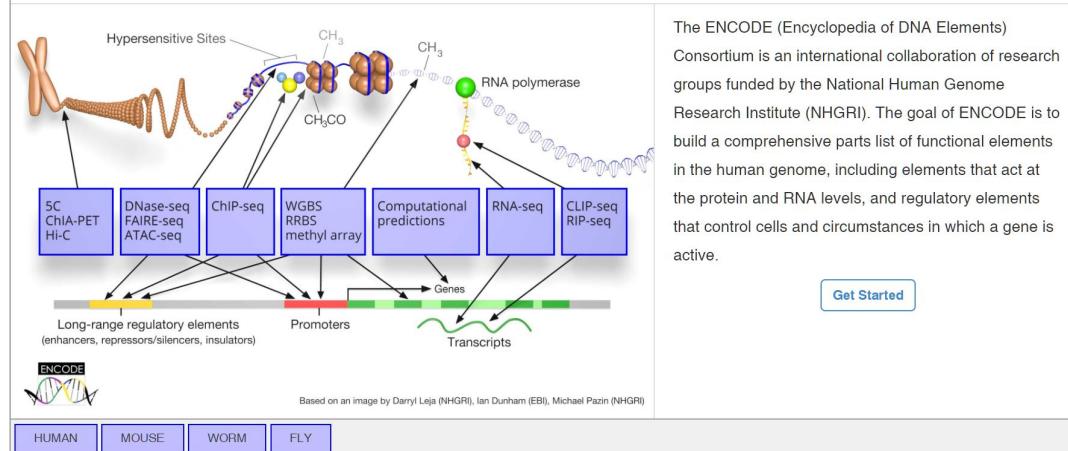
Caveats to ChIP/ATAC-seq

- False positive rate for peak detection is high
- Variation between peak detection methods is HIGH
- Association between binding and function is LOW
- The cost of the experiment is HIGHER than you plan
- Optimization requires MUCH MORE time than you plan
- The bioinformatics takes more time than you expect

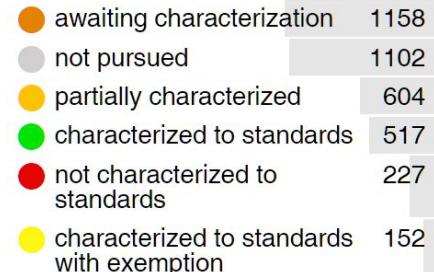


Encode: Resource for Protocols, Workflows, and Data

ENCODE: Encyclopedia of DNA Elements



Eligibility status



<https://www.encodeproject.org/search/?type=AntibodyLot>

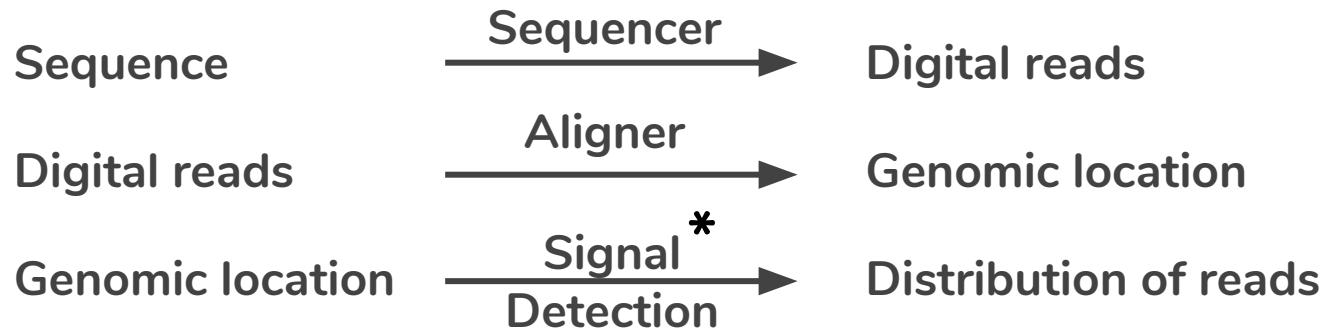
End Part 1: Break



Sequencing Pipeline: Reads to Data



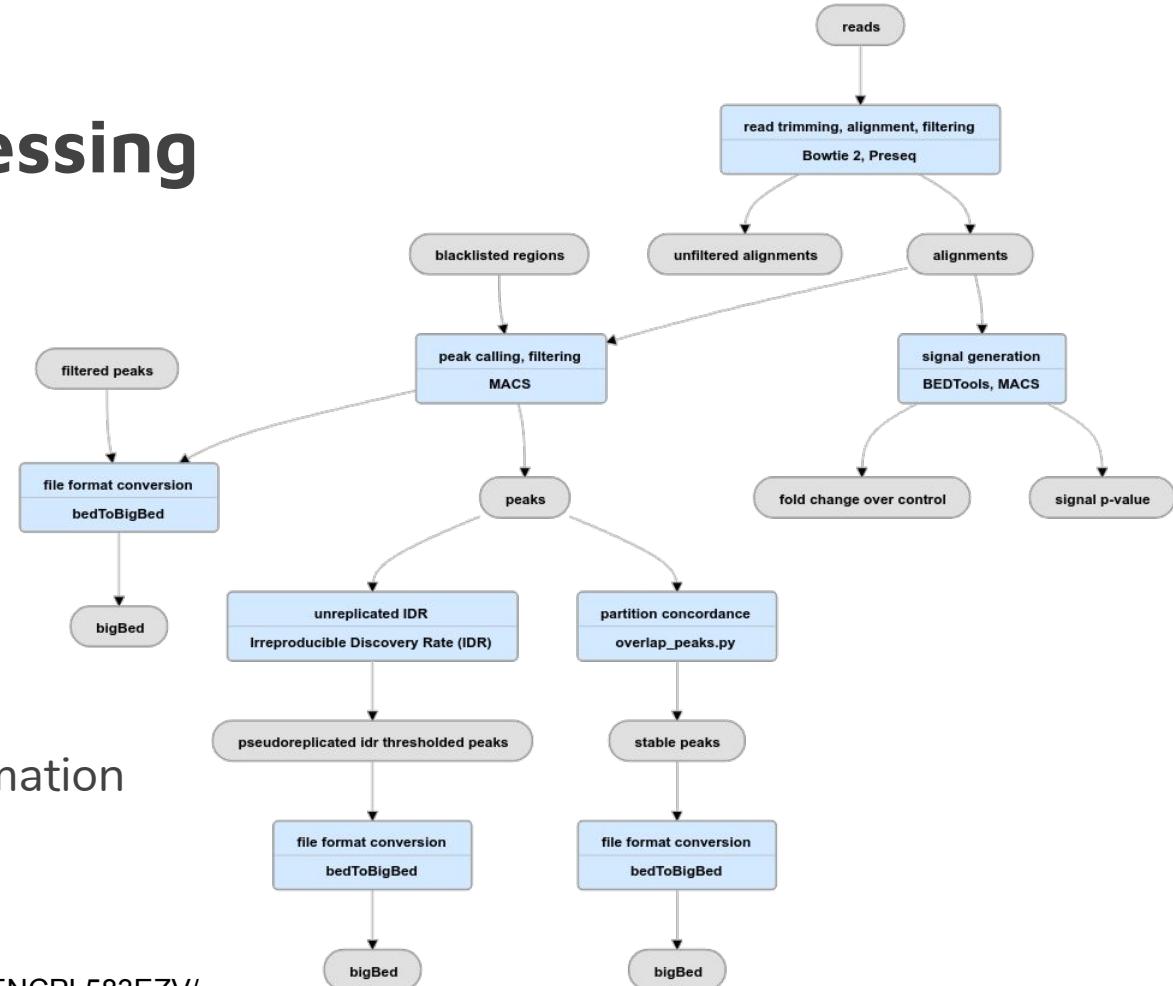
Downstream Terminology



* Specifically, in ChIP/ATAC-seq this is called peak calling and the read distribution is called a peak.

Read Processing

1. Adapter Trimming
2. Read QC
3. Read Alignment
4. Alignment QC
5. Peak Calling
6. Peak QC
7. Peak Confidence Estimation
8. Peak Visualization





Read Processing: Filtering



NOTE

Be aware that the CASAVA 1.8 FASTQ files contain all reads, both the reads that passed filtering, as well as the reads that did not pass filtering. If you use third-party software that cannot handle this, we recommend that clean up the FASTQ files first using the <is filtered> field described above. You can use the

```
cd /path/to/project/sample  
mkdir filtered  
for fastq in *.fastq.gz; do zcat $fastq | grep  
    -A 4 '^@.* [^:]*:N:[^:]*:' > filtered/$fastq  
; done
```

WRONG!

1. grep:

```
$ cat input.fq | grep -A 3 '^@.* [^:]*:N:[^:]*:' | grep -v "^--$"  
> output.fq
```

2. Fastq_illumina_filter

- http://cancan.cshl.edu/labmembers/gordon/fastq_illumina_filter/



Read Processing: Adapter Trimming

1. Cutadapt

- a. Auto-detect adapters
- b. Handles short adapters well
- c. Doesn't keep searching

2. Trimmomatic

- a. No auto-detection
- b. Keeps trimming as long as it finds adapter sequence

3. Trim Galore! (recommended)

- a. Wrapper for Cutadapt and FastQC
- b. https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/



Read Processing: FastQC

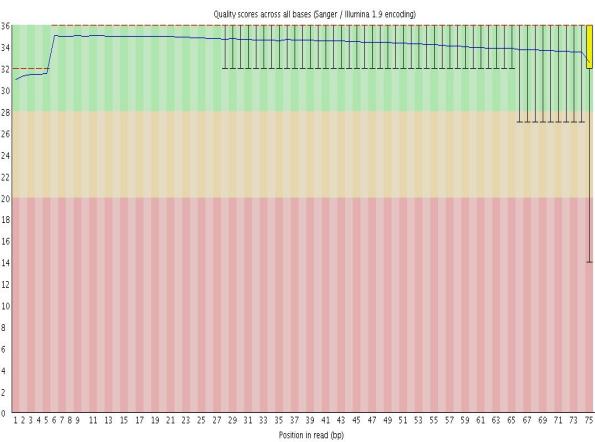
1. FastQC

- a. May be the only thing that (almost) everyone agrees on (<https://www.biostars.org/p/204261/>).

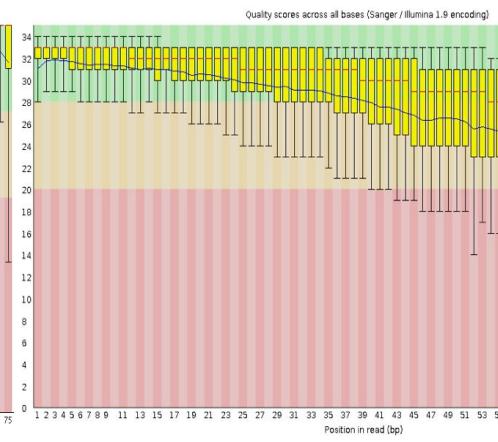


FastQC: Per-Base Quality Score

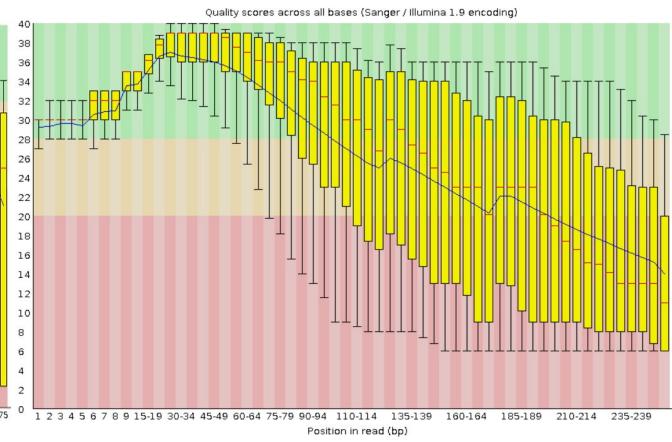
Good



Bad



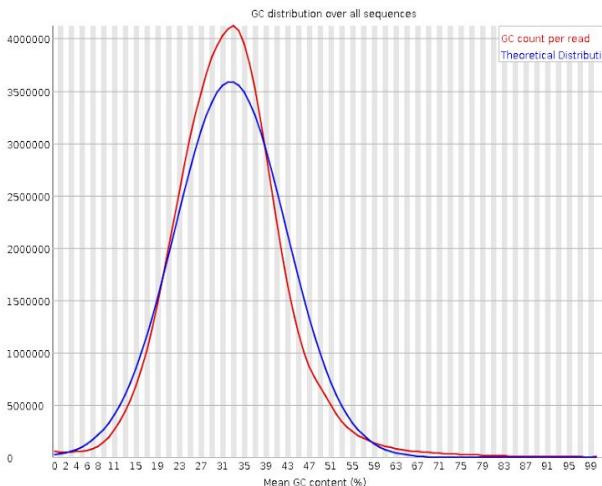
Ugly



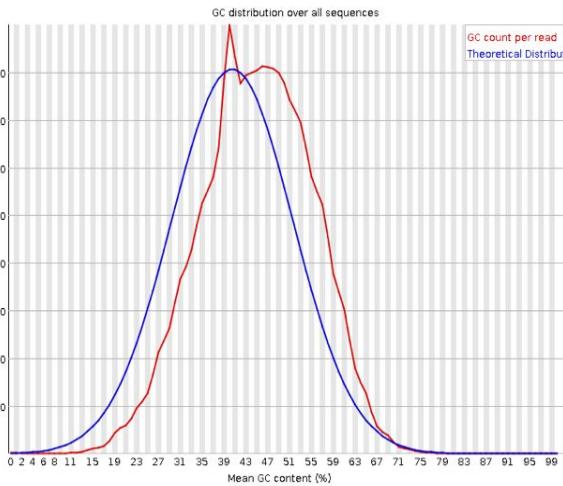


FastQC: Average GC Content

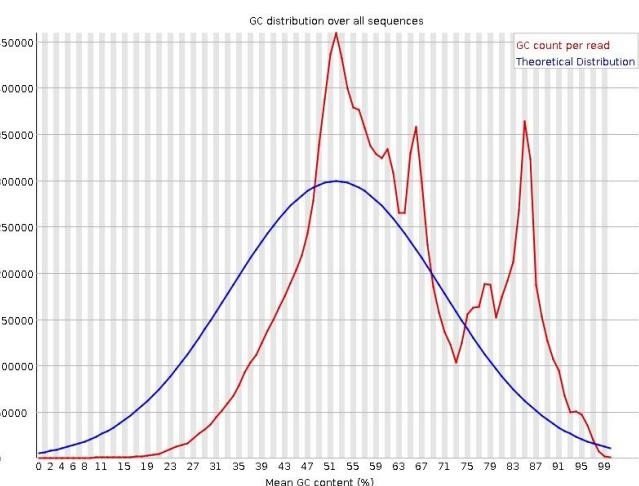
Good



Bad



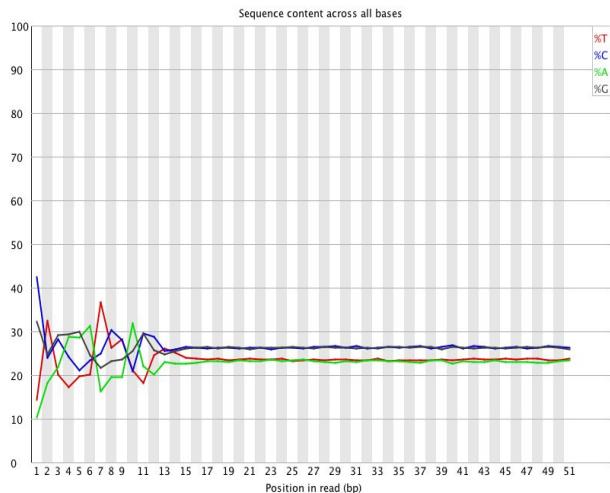
Ugly



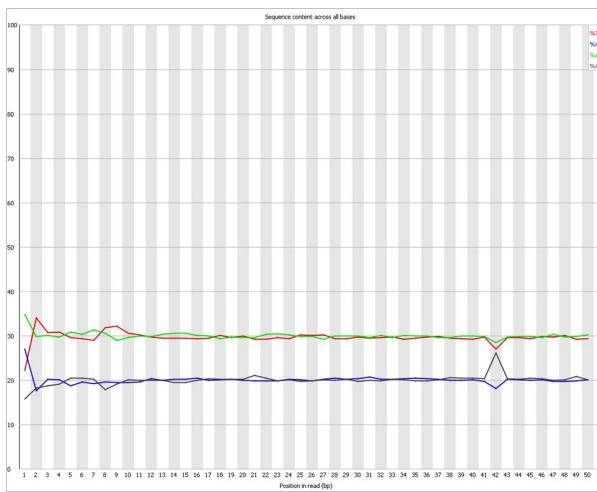


FastQC: Per-Base Nucleotide Content

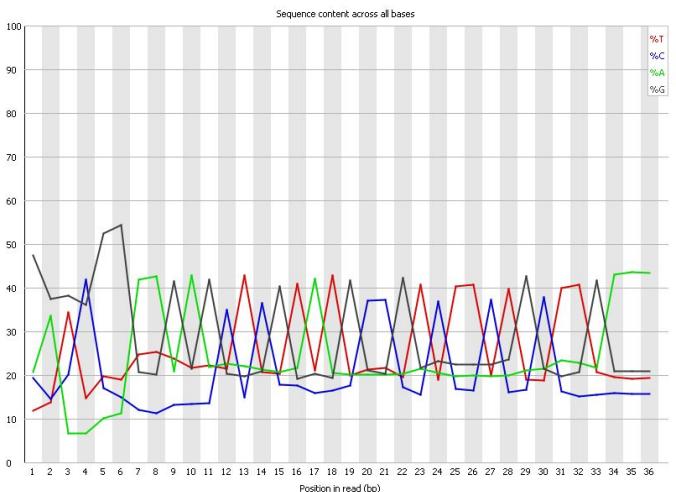
Good



Bad



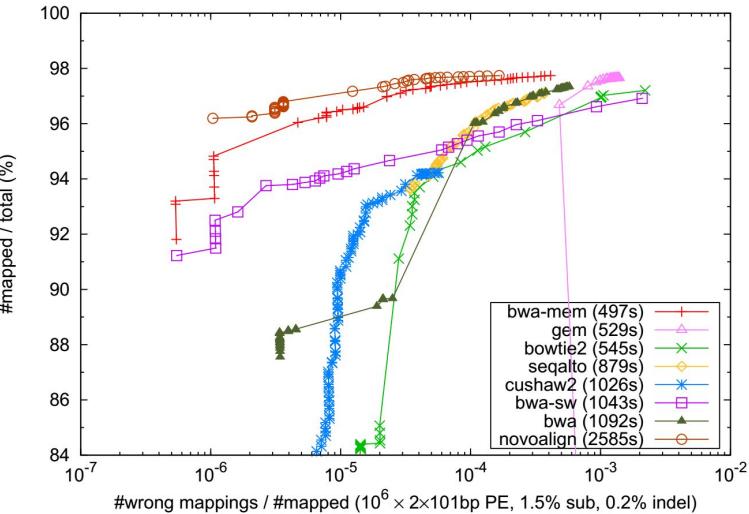
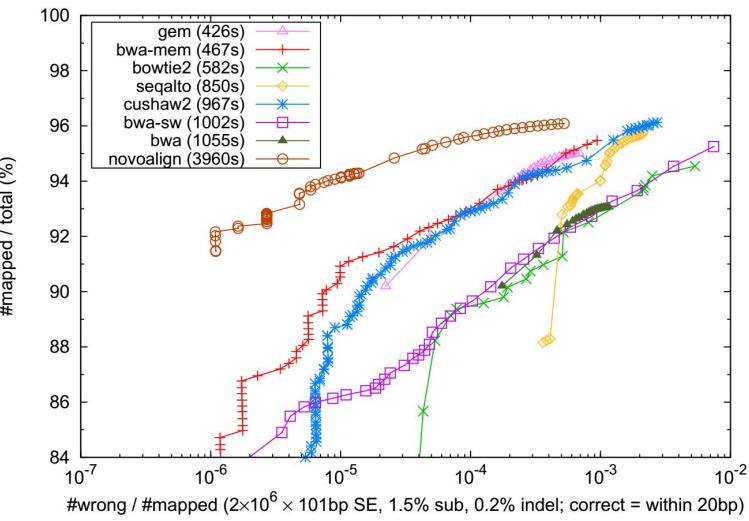
Ugly





Read Alignment

1. Bowtie2
 - a. Most popular
2. BWA
 - a. Sensitive but SLOW
3. ELAND (Illumina)
 - a. Core's often use this
4. Novoalign
 - a. Proprietary
5. STAR
 - a. Careful!
 - b. Need to prevent spliced alignment



Read Alignment: QC

alignment summary

18200175 (100.00%) were paired; of these:

497844 (2.74%) aligned **concordantly** 0 times

11291185 (62.04%) **aligned concordantly exactly 1 time**

6411146 (35.23%) **aligned concordantly >1 times**

497844 pairs aligned concordantly 0 times; of these:

139262 (27.97%) **aligned discordantly 1 time**

358582 pairs aligned 0 times concordantly or discordantly; of these:

717164 **mates make up the pairs**; of these:

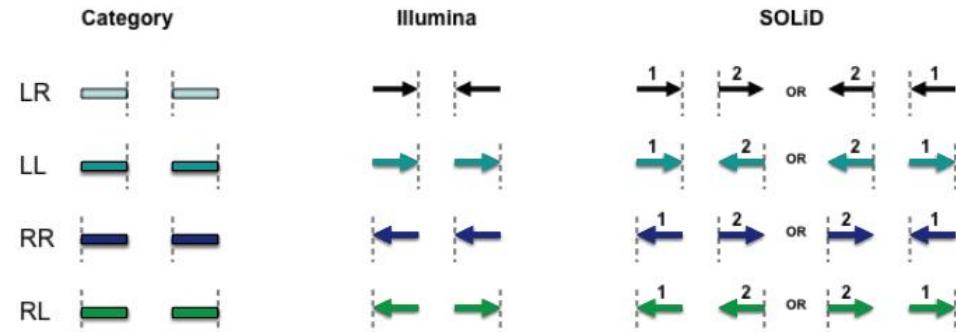
162541 (22.66%) aligned 0 times

174764 (24.37%) aligned exactly 1 time

379859 (52.97%) aligned >1 times

99.55% overall alignment rate

Interpretation of read pair orientations



LR

Normal reads.

The reads are left and right (respectively) of the unsequenced part of the sequenced DNA fragment when aligned back to the reference genome.

LL,RR

Implies inversion in sequenced DNA with respect to reference.

RL

Implies duplication or translocation with respect to reference.



Alignments: SAM Format

Col	Field	Type	Regexp/Range	Brief description
1	QNAME	String	[!-?A-~]{1,254}	Query template NAME
2	FLAG	Int	[0,2 ¹⁶ -1]	bitwise FLAG
3	RNAME	String	* [!-()+-<>-~] [!-~]*	Reference sequence NAME
4	POS	Int	[0,2 ³¹ -1]	1-based leftmost mapping POSition
5	MAPQ	Int	[0,2 ⁸ -1]	MAPping Quality
6	CIGAR	String	* ([0-9]+[MIDNSHPX=])+	CIGAR string
7	RNEXT	String	* = [!-()+-<>-~] [!-~]*	Ref. name of the mate/next read
8	PNEXT	Int	[0,2 ³¹ -1]	Position of the mate/next read
9	TLEN	Int	[-2 ³¹ +1,2 ³¹ -1]	observed Template LENgth
10	SEQ	String	* [A-Za-z.=.]+	segment SEQuence
11	QUAL	String	[!-~]+	ASCII of Phred-scaled base QUALity+33



Filtering Alignments

- Samtools

- -f keeps all alignments with specific flag bits
- -F excludes all alignments with specific flags
- Decoding SAM flags:
 - <https://broadinstitute.github.io/picard/explain-flags.html>



Removing PCR Duplicates

- Samtools rmdup
 - Does not remove PCR duplicates for chromosomal translocation reads
- Picard Tools MarkDuplicates
 - Marks PCR duplicates and will remove them if requested
 - Removes translocation duplicates missed by samtools rmdup command.



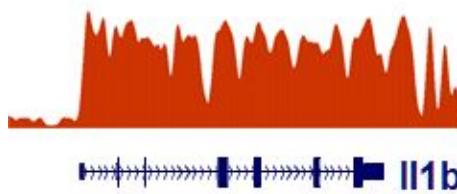
Peak calling: genomic read distribution

TF ChIP



Narrow Peaks

Histone Mod
ChIP



Broad Peaks

ATAC



Mixed

Peak Callers

- MACS
 - ENCODE
 - Most widely used
 - Stable performance
- SPP
 - R implementation
- FSEQ
 - Kernel density based
 - ENCODE
 - Can be convinced to call peaks.

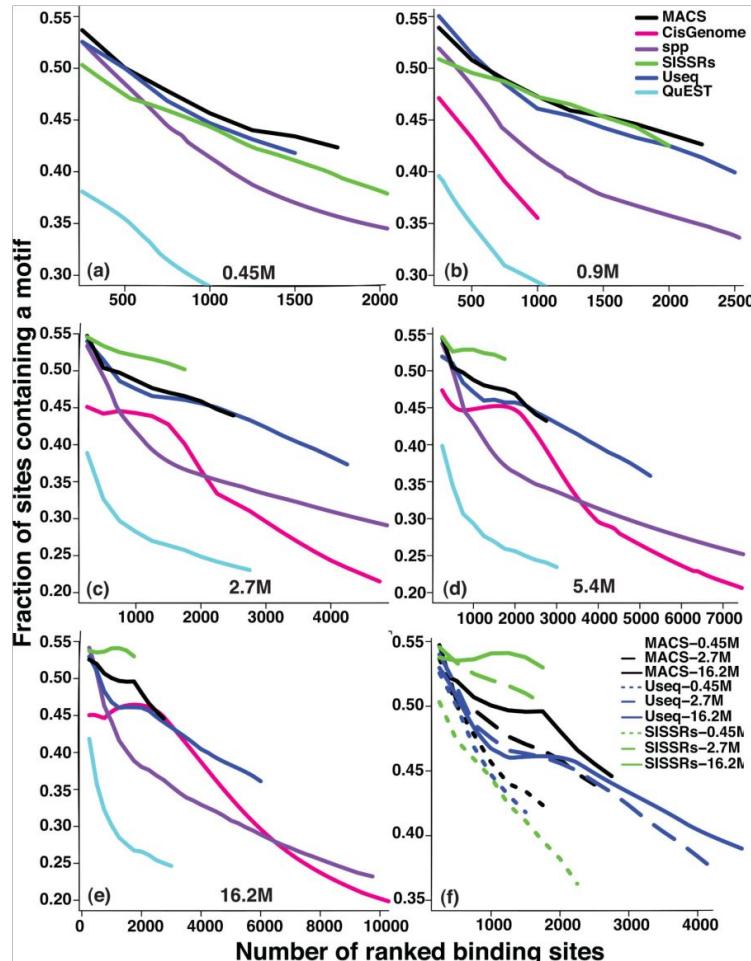


Figure 3. Quality of the Su(Hw) peaks

The fraction Su(Hw) peaks, identified by the indicated peak callers, that contains a Su(Hw) binding motif is plotted as a function of the number of top-ranked binding sites at the



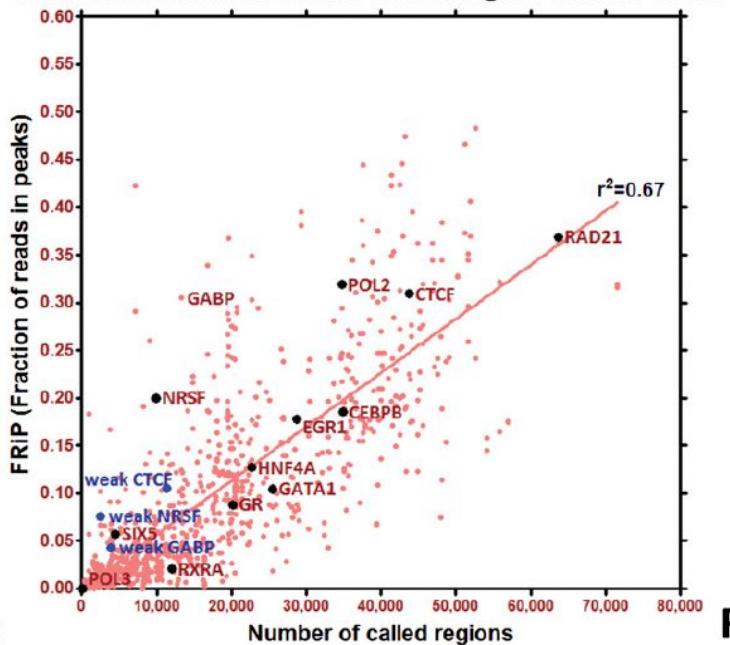
Summary QC Reports with MultiQC



[MultiQC Example](#)

Quality Control: FRiP

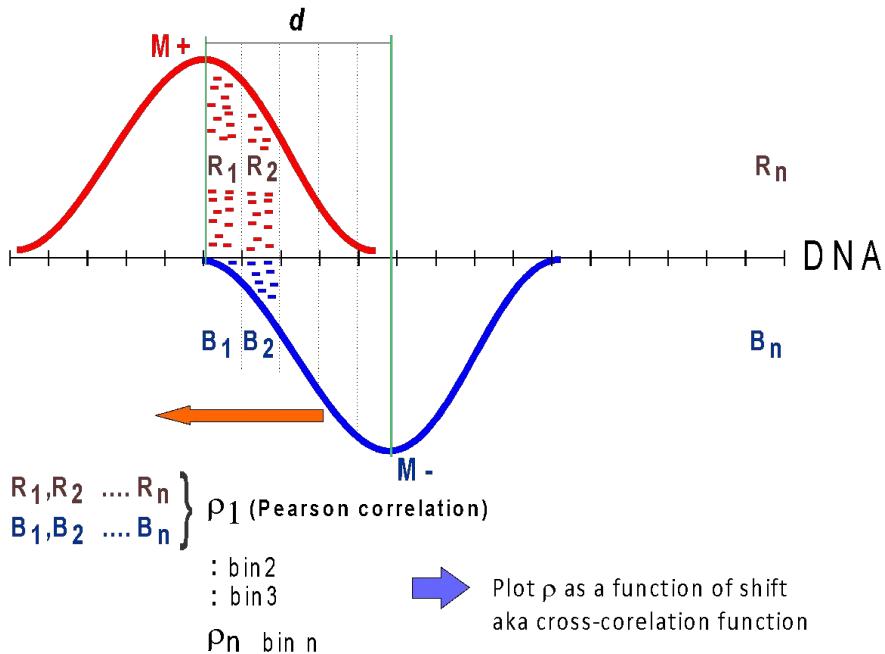
Correlation between the number of called regions and FRiP scores



- FRiP = Reads in Peaks / Total Mapped Reads
- FRiP < 1% may be a cause for concern
- Correlated with peak number

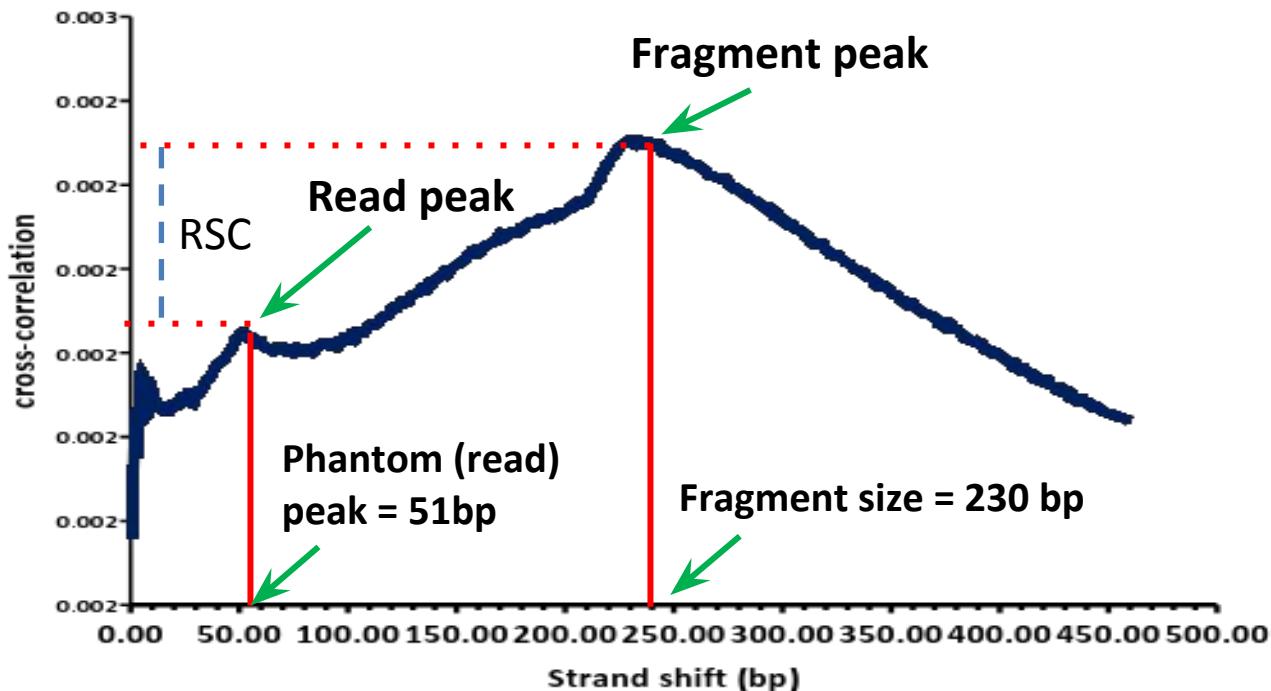


Quality Control: Cross Correlation Analysis

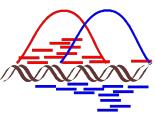


- When the signal from the **positive strand (red)** matches the signal from the **negative strand (blue)** the correlation is maximized

Cross Correlation Analysis Example

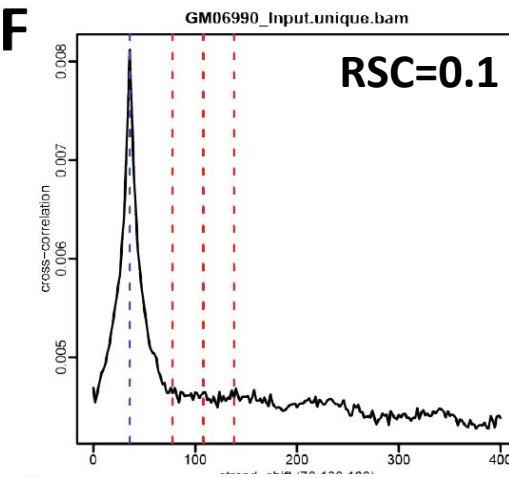
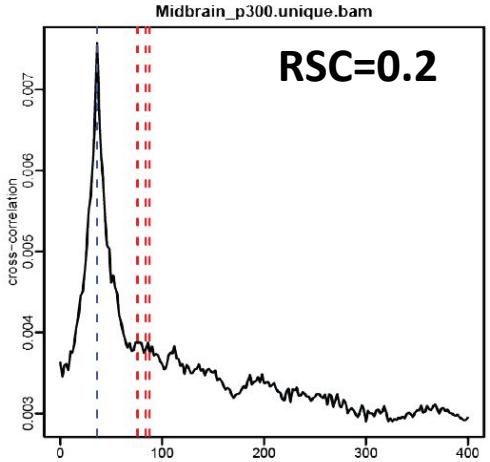


Relative strand correlation (RSC) = Fragment Peak/Read Peak = 1.617 (should be >0.8)

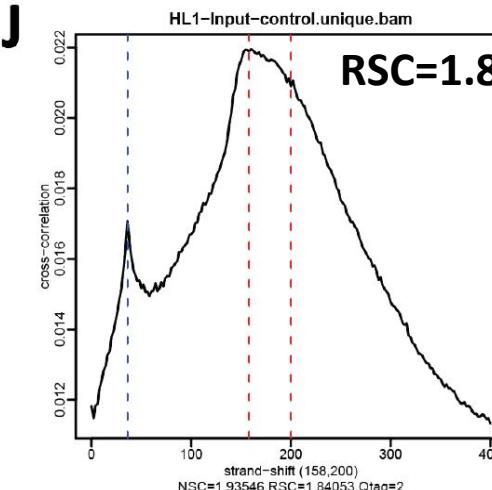
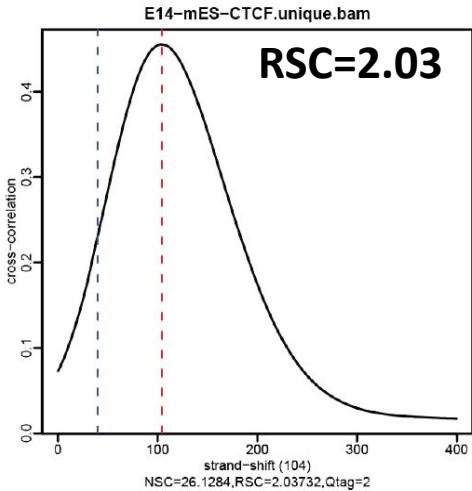


Cross Correlation analysis currently in repository

Bad datasets



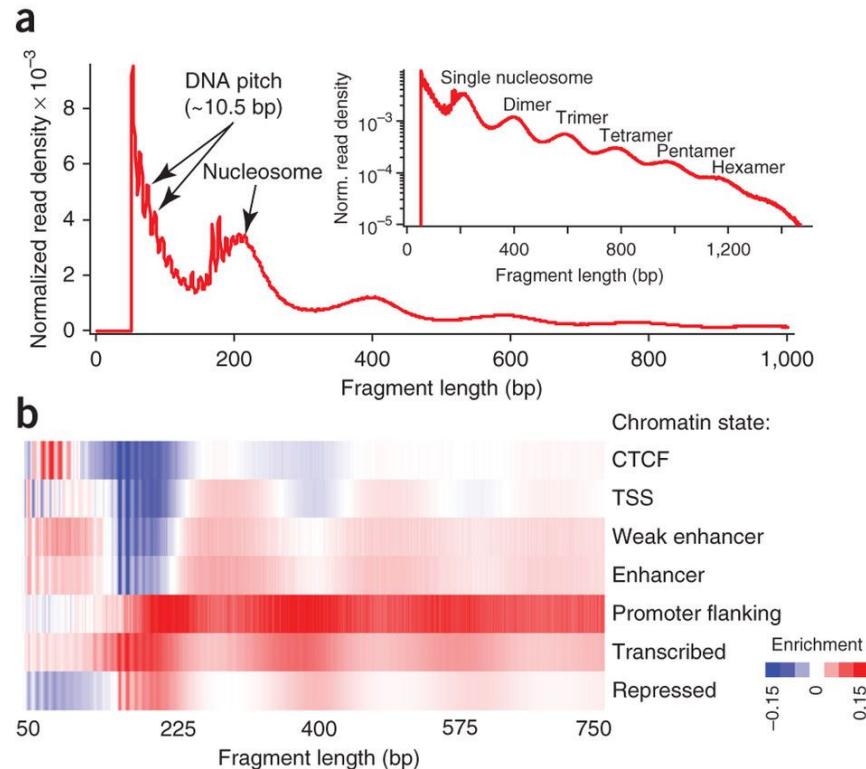
Good datasets





ATAC-seq Fragment Size Distribution

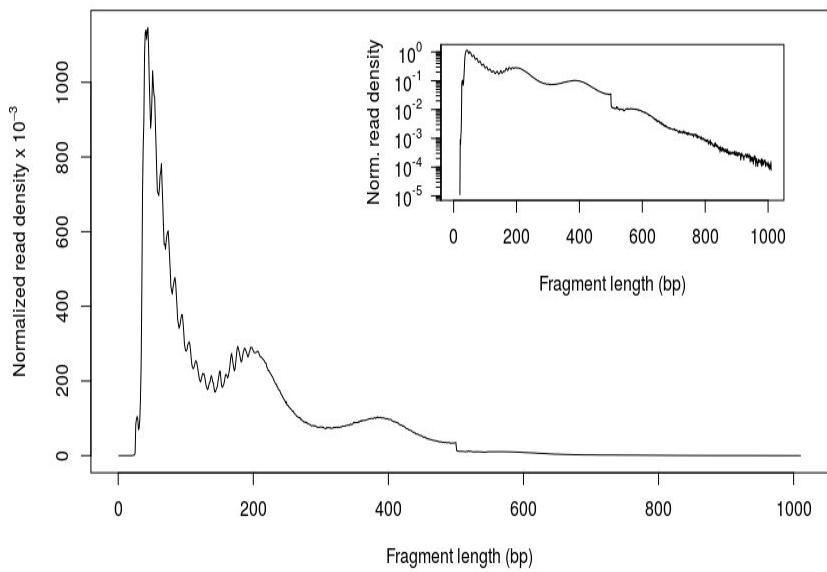
- Looking for distinct nucleosome profile.



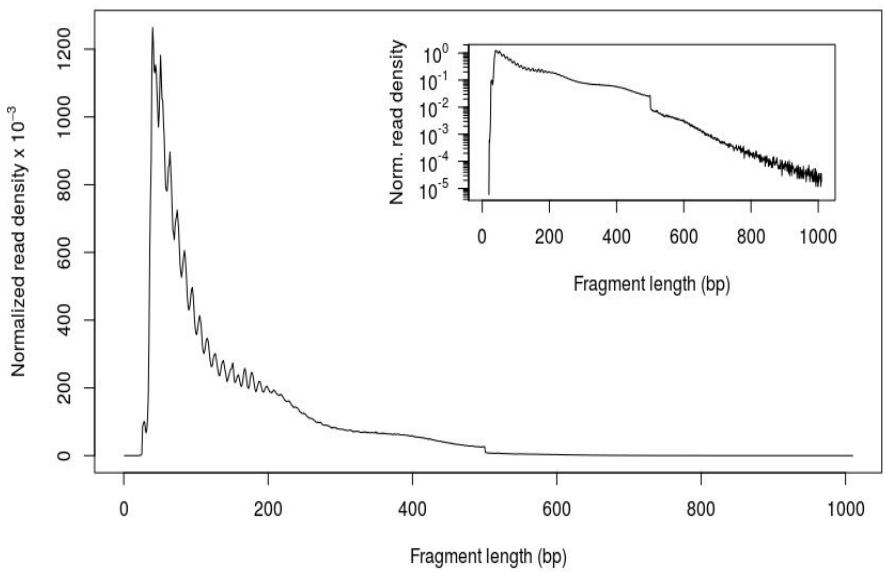


ATAC-seq Fragment Size Distribution

Good



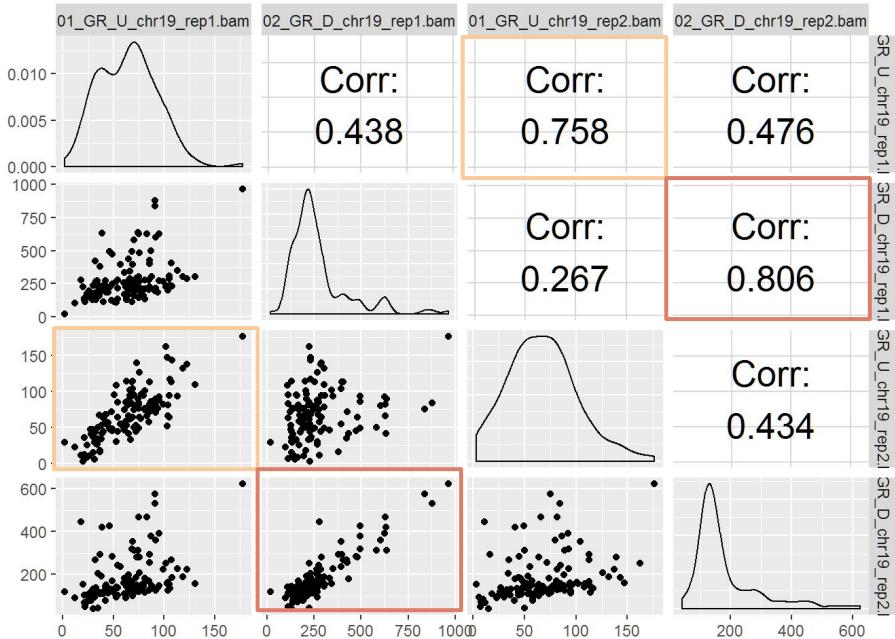
Poor





Correlation Between Read in Peaks

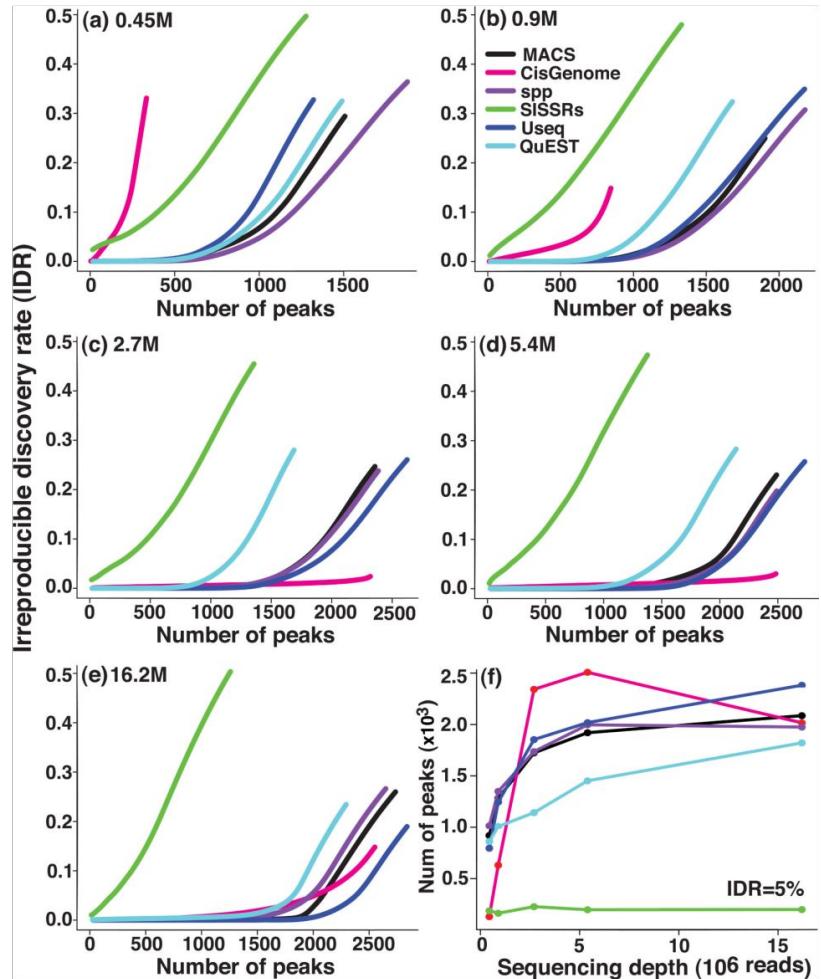
1. Union of peaks
between replicates
2. Count reads in peaks
3. Calculate correlation





Irreproducible discovery rate (IDR)

1. Peak reproducibility between replicates.
2. ENCODE uses it
3. Slightly better than Fisher's and Stouffer's methods.
4. Honestly, use one of these but the IDR program is finicky and the method is not defined for more than 2 replicates



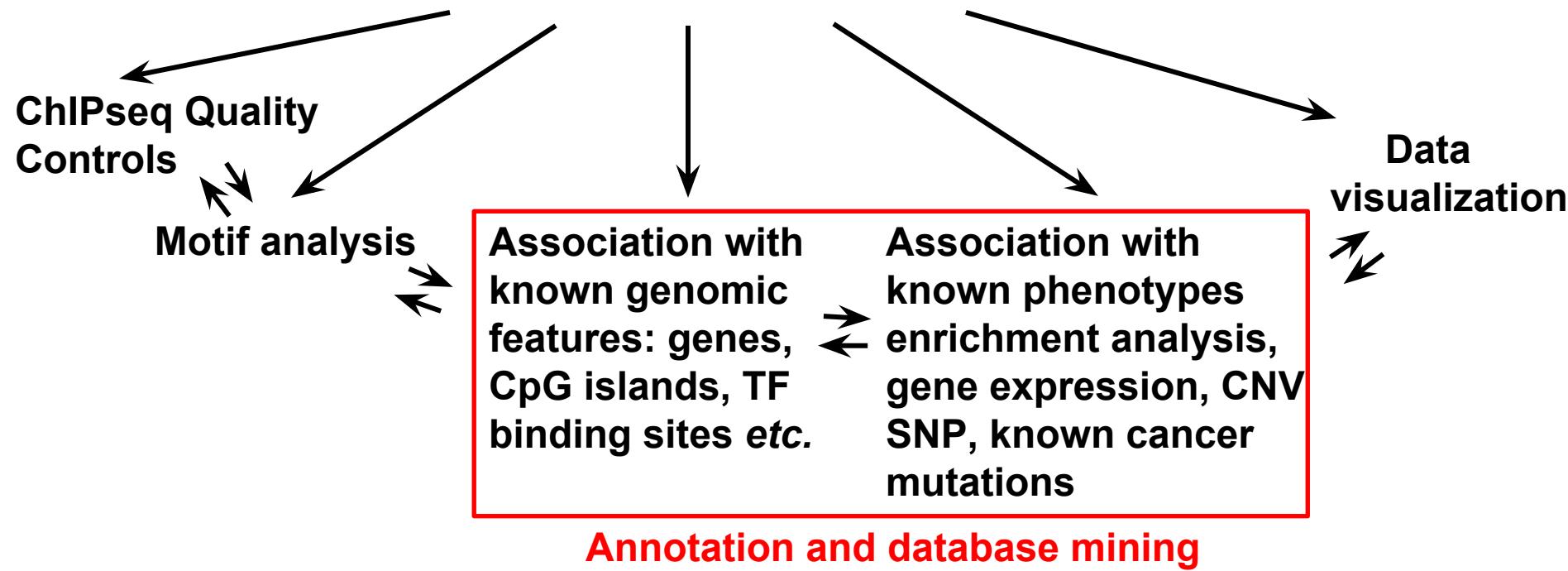
End Part 2: Break



Functional Analysis: Data to Meaning

After peak calling:

BAM file of aligned reads
BED file of peak coordinates



DNA Motif discovery

Supervised

Compares peak sequences to known binding sites databases such as **JASPAR** (<http://jaspar.genereg.net/>) and **TRANSFAC**

- Can only find “known” sites
- Relies on the quality of databases
- Can work with a small number of sequences
- Biased
(Trust, but verify – there are mistakes and circular annotations)

Unsupervised

Doesn’t care about binding sites.
Treats a collection of sequences as a text that might contain repeating words.

- Slow
- Noisy
- Need large number of sequences
- Need to interpret motifs
(which means comparing to databases)
- unbiased

De novo Motif discovery

The MEME Suite

<http://alternate.meme-suite.org/>



<http://xxmotif.genzentrum.lmu.de/>



<http://rsat.sb-roscoff.fr/index.html>



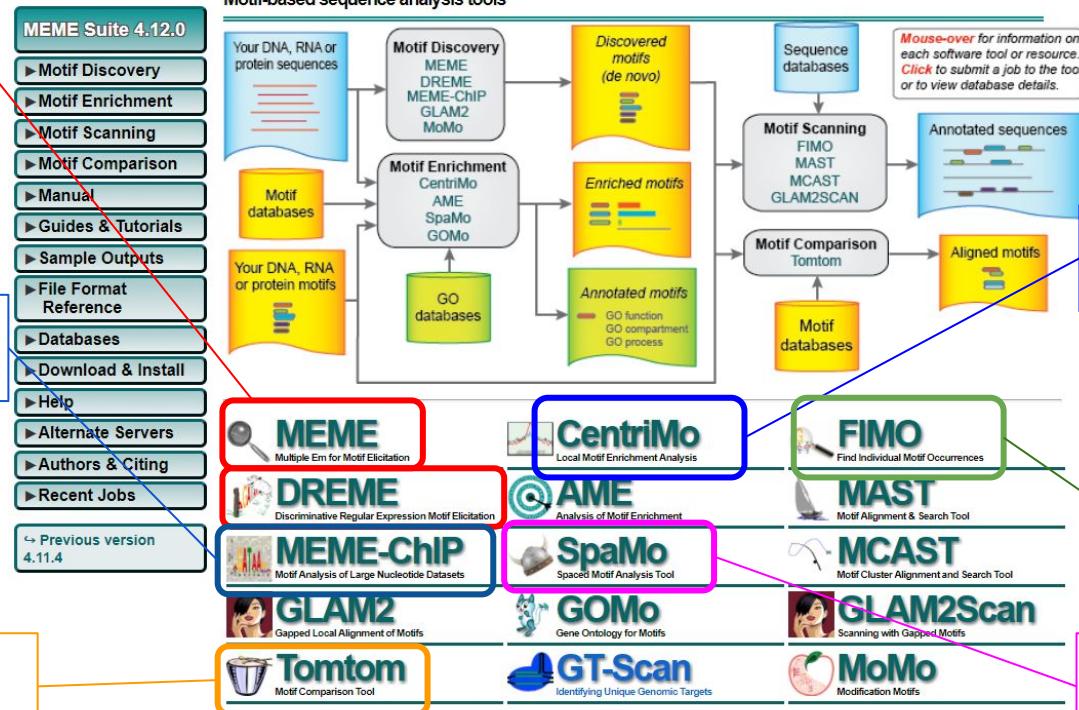
HOMER

<http://homer.ucsd.edu/homer/index.html>

The MEME Suite

Motif-based sequence analysis tools

De novo motif search
overrepresentation
with Gibbs sampling



Performs all analyses
at once

Comparing identified
motifs to known motifs

Motif positional
enrichment

Tabulating motif
occurrence

Co-occurrence of
motifs

De novo motif prediction/overrepresentation

Typical MEME output



If you use MEME-ChIP in your research, please cite the following paper:
Philip Machanick and Timothy L. Bailey, "MEME-ChIP: motif analysis of large DNA datasets", *Bioinformatics*, 27(12), 1696-1697, 2011. [\[full text\]](#)

[MOTIFS](#) | [PROGRAMS](#) | [INPUT FILES](#) | [PROGRAM INFORMATION](#) | [SUMMARY IN TSV FORMAT](#) | [NEW](#) | [MOTIFS IN MEME TEXT FORMAT](#) | [NEW](#)

MOTIFS



The significant motifs ($E\text{-value} \leq 0.05$) found by the programs MEME, DREME and CentriMo; clustered by similarity and ordered by $E\text{-value}$.

Expand All Clusters

Collapse All Clusters



Discovery/Enrichment Program

MEME
1.6e-016

E-value
Known or Similar Motifs

Erg (MA0474.1)
FLI1 (MA0475.1)
Ets1 (MA0098.2)

Distribution
Not Centrally Enriched

SpaMo & FIMO
• Motif Spacing Analysis
• Motif Sites in GFF3

Show 6 More ↴ CentriMo Group ↴



Discovery/Enrichment Program

MEME
4.9e-010

E-value
Known or Similar Motifs

NR3C1 (MA0113.2)
NR3C2_DB3 AR (MA0007.2)

Distribution



SpaMo & FIMO

• Motif Spacing Analysis
• Motif Sites in GFF3

Show 6 More ↴ CentriMo Group ↴

De novo motif prediction/overrepresentation

Typical MEME output

MEME-ChIP
Motif Analysis of Large Nucleotide Datasets

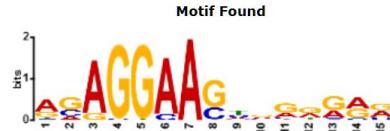
If you use MEME-ChIP in your research, please cite the following paper:
 Philip Machanick and Timothy L. Bailey, "MEME-ChIP: motif analysis of large DNA datasets", *Bioinformatics*, 27(12), 1696-1697, 2011. [\[full text\]](#)

[MOTIFS](#) | [PROGRAMS](#) | [INPUT FILES](#) | [PROGRAM INFORMATION](#) | [SUMMARY IN TSV FORMAT](#) | [MOTIFS IN MEME TEXT FORMAT](#)

MOTIFS

The significant motifs ($E\text{-value} \leq 0.05$) found by the programs MEME, DREME and CentriMo; clustered by similarity and ordered by $E\text{-value}$.

[Expand All Clusters](#) [Collapse All Clusters](#)



Discovery/Enrichment Program	E-value	Known or Similar Motifs	Distribution
MEME	1.6e-016	Erg (MA0474.1) FL1 (MA0475.1) Ets1 (MA0098.2)	Not Centrally Enriched

Reverse Complement Show 6 More

CentriMo Group



Discovery/Enrichment Program	E-value	Known or Similar Motifs	Distribution
MEME	4.9e-010	NR3C1 (MA0113.2) NR3C2 DBD AR (MA0007.2)	

Reverse Complement Show 6 More

CentriMo Group



- Motif Spacing Analysis
- Motif Sites In GFF3

Positions of motifs

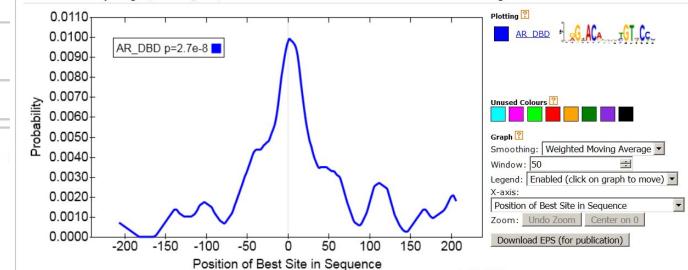
CentriMo
Local Motif Enrichment Analysis

For further information on how to interpret these results or to get a copy of the MEME software please access <http://meme-suite.org>. If you use CentriMo in your research, please cite the following paper:
 Timothy L. Bailey and Philip Machanick, 'Inferring direct DNA binding from ChIP-seq', *Nucleic Acids Research*, 40:e128, 2012. [\[full text\]](#)

[MOTIF PROBABILITY GRAPH](#) | [ENRICHED MOTIFS](#) | [INPUT FILES](#) | [PROGRAM INFORMATION](#)

RESULTS

Motif Probability Graph (score ≥ 5 bits)



Enriched motifs ($E\text{-value} \leq 10$ using the binomial test)

ID	Alt ID	Consensus	E-value	Region Width
NR3C2_DBD	NR3C1	KRGWACAYRTGTWCYH	3.5e-5	133
AR_DBD	AR	DRGWACAYSRTGTWCY	3.9e-5	133

Matching sequences (out of 127)

Intersections: 39 sequences (31%).

Intersection: 39 sequences (31%).

Intersections: 39 sequences (31%).

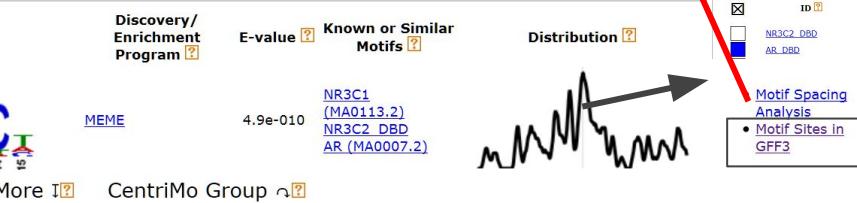
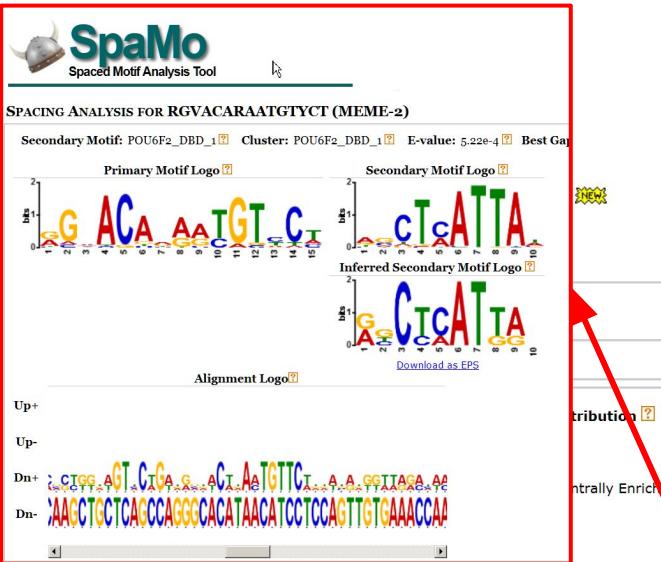
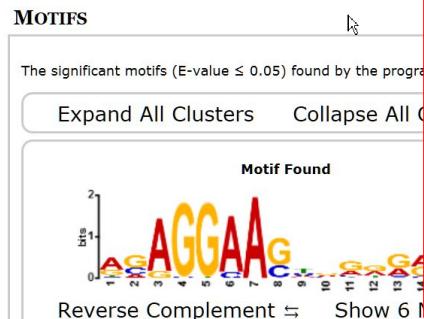
De novo motif prediction/overrepresentation

Typical MEME output

MEME-ChIP
Motif Analysis of Large Nucleotide Datasets

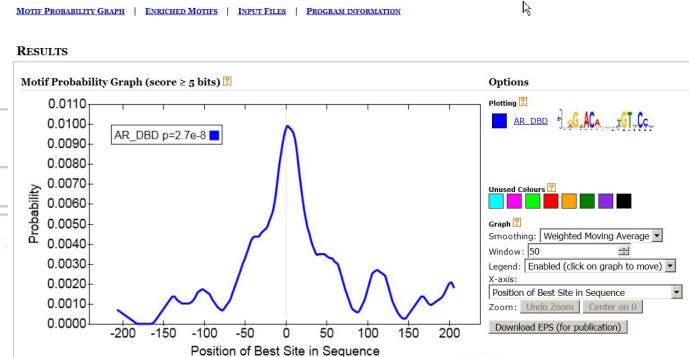
If you use MEME-ChIP in your research, please cite the following paper:
Philip Machanick and Timothy L. Bailey, "MEME-ChIP: motif analysis of ChIP-seq data", *Nucleic Acids Research*, 40:e128, 2012. [full text]

[MOTIFS](#) | [PROGRAMS](#) | [INPUT FILES](#) | [PROGRAM INFORMATION](#)



CentriMo
Local Motif Enrichment Analysis

For further information on how to interpret these results or to get a copy of the MEME software please access <http://meme-suite.org>. If you use CentriMo in your research, please cite the following paper:
Timothy L. Bailey and Philip Machanick, "Inferring direct DNA binding from ChIP-seq", *Nucleic Acids Research*, 40:e128, 2012. [full text]

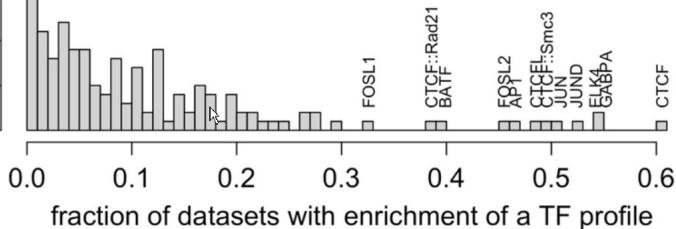


Positions of motifs

My program detected motif X in my peaks. It is a “real” binding site, right?

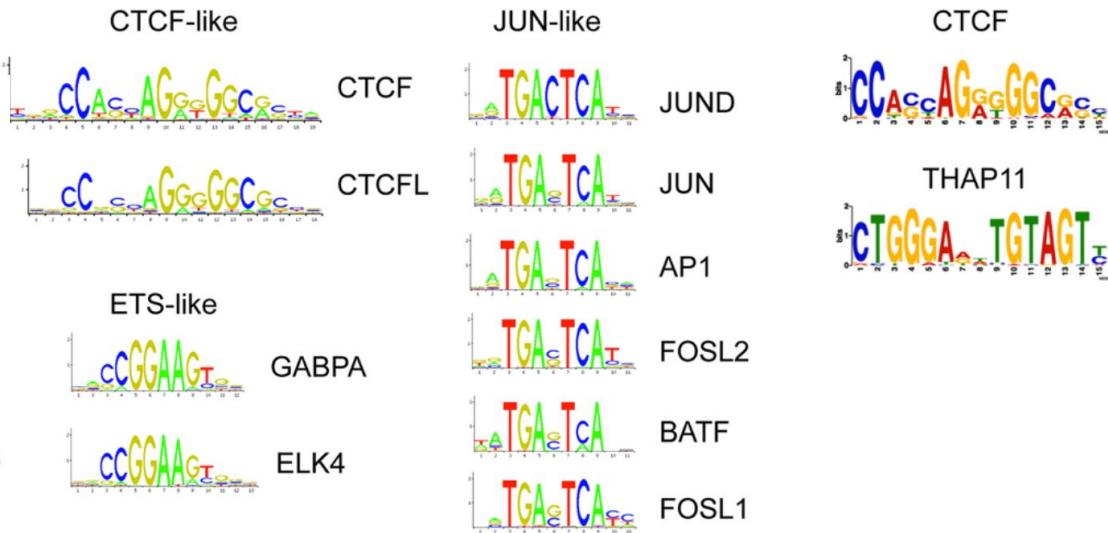


281 human
ChIP-seq datasets



... may be?

- Summit enrichment?
- Site complexity?
- Correlation with peak scores?
- Distance conservation relative to other sites



Motif overrepresentation – data interpretation

Co-occurrence of motifs does not mean co-function

Worsley Hunt R, Wasserman WW. Non-targeted transcription factors motifs are a systemic component of ChIP-seq datasets. *Genome Biol.* 2014 Jul 29;15(7):412. PMID: 25070602;

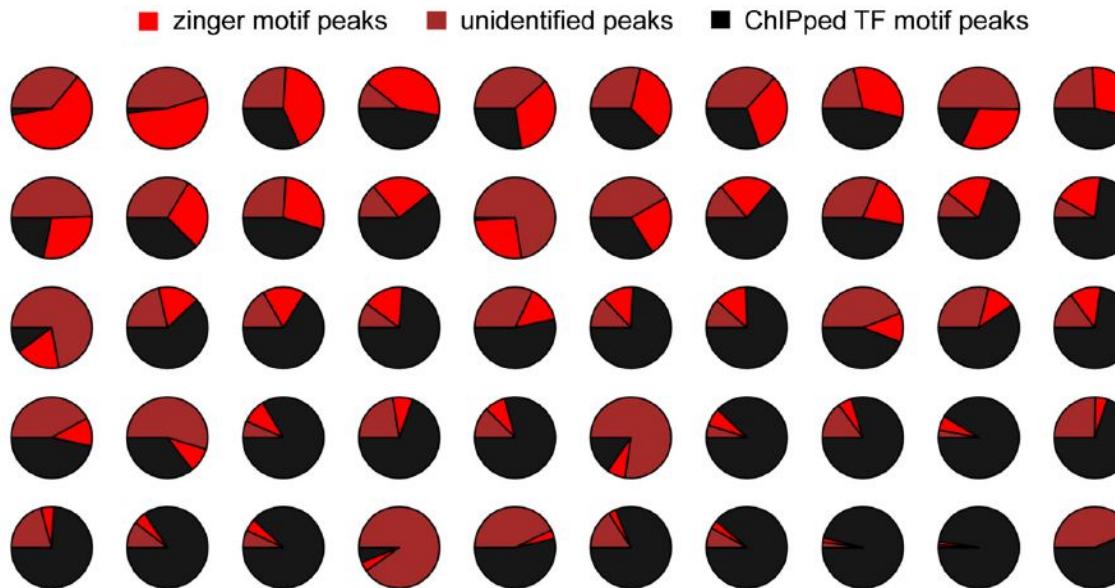


Figure 3 The fraction of zinger motif peaks and ChIPped TF motif peaks varies across ChIP-seq datasets. The pie charts present a random selection of 50 datasets for multiple TFs and cell-lines with zinger motifs present (>1% zinger). The charts are ordered by greatest zinger motif peak enrichment to the least. Black is the portion of peaks with the ChIPped TF's motif, red is the portion of zinger motif peaks, and brown is the remaining portion of peaks that do not contain either the ChIPped TF nor zinger motifs.

My program claims that the site enriched in my peak is a binding site for factor X. Is it?

	ID	Name	Species	Class	Family	Logo
	MA0039.1	Klf4	Mus musculus	C2H2 zinc finger factors	Three-zinc finger Kruppel-related factors	
	MA0039.2	Klf4	Mus musculus	C2H2 zinc finger factors	Three-zinc finger Kruppel-related factors	
	MA0039.3	KLF4	Homo sapiens	C2H2 zinc finger factors	Three-zinc finger Kruppel-related factors	

[Copy](#)[CSV](#)

My program claims that the site enriched in my peak is a binding site for factor X. Is it?

ID	Name	Profile summary
MA0039.1	Klf4	<p>Name: Klf4</p> <p>Matrix ID: MA0039.1</p> <p>Class: C2H2 zinc finger factors</p> <p>Family: Three-zinc finger Kruppel-related factors</p> <p>Collection: CORE</p> <p>Taxon: Vertebrates</p> <p>Species: Mus musculus</p> <p>Data Type: SELEX</p> <p>Validation: 9443972</p> <p>Uniprot ID: Q60793</p>
MA0039.2	Klf4	<p>Three-zinc finger Kruppel-related factors</p>
MA0039.3	KLF	<p>Three-zinc finger Kruppel-related factors</p>
		<p>Three-zinc finger Kruppel-related factors</p>

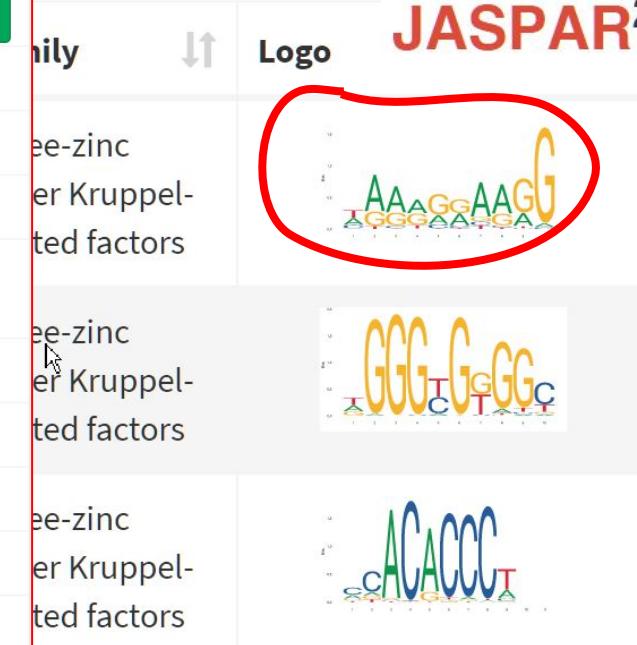
JASPAR 2018

N	5'	- G A T G C	1	2	3	4	5	6	7	8	9	10	11	12	13	14	C G	T A G	- 3'
1	5'	- G A T G C	A	A	A	G	A	A	G	G	G	A	A	G	G	C G	T A G	- 3'	
2	5'	- G A T G C	G	A	A	G	A	A	G	A	A	G	A	T	G	G	C G	T A G	- 3'
3	5'	- G A T G C	A	A	G	T	A	A	T	G	A	G	G	A	G	G	C G	T A G	- 3'
4	5'	- G A T G C	A	A	T	G	A	C	G	A	A	G	G	G	G	G	C G	T A G	- 3'
5	5'	- G A T G C	A	C	C	A	G	T	G	A	A	C	A	G	G	G	C G	T A G	- 3'
6	5'	- G A T G C	A	A	A	A	G	A	A	A	A	T	G	G	G	G	C G	T A G	- 3'
7	5'	- G A T G C	A	C	C	A	G	T	G	A	A	C	A	G	G	G	C G	T A G	- 3'
8	5'	- G A T G C	A	A	A	A	G	A	A	A	A	A	T	G	G	G	C G	T A G	- 3'
9	5'	- G A T G C	G	A	G	T	A	A	T	G	A	A	G	G	G	G	C G	T A G	- 3'
10	5'	- G A T G C	C	G	C	G	G	G	G	G	G	G	A	T	A	G	C G	T A G	- 3'
11	5'	- G A T G C	C	A	A	T	A	A	G	A	G	A	T	A	G	G	C G	T A G	- 3'
12	5'	- G A T G C	C	A	A	T	A	A	G	A	G	A	T	A	G	G	C G	T A G	- 3'
13	5'	- G A T G C	A	A	C	T	G	A	A	G	A	G	G	G	G	G	C G	T A G	- 3'
14	5'	- G A T G C	A	A	C	T	G	A	A	G	A	G	G	G	G	G	C G	T A G	- 3'
15	5'	- G A T G C	A	A	G	G	A	A	G	T	A	A	G	G	G	G	C G	T A G	- 3'
16	5'	- G A T G C	A	T	C	G	T	C	T	G	T	T	A	G	G	G	C G	T A G	- 3'
17	5'	- G A T G C	A	T	A	T	A	G	A	A	A	G	G	G	G	G	C G	T A G	- 3'
18	5'	- G A T G C	A	T	A	T	A	G	A	A	A	G	G	G	G	G	C G	T A G	- 3'
19	5'	- G A T G C	A	T	A	T	A	G	A	A	A	G	G	G	G	G	C G	T A G	- 3'
20	5'	- G A T G C	T	A	G	A	A	G	G	G	G	G	G	G	G	G	C G	T A G	- 3'
21	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
22	5'	- G A T G C	C	C	G	A	A	A	A	A	G	A	A	G	G	G	C G	T A G	- 3'
23	5'	- G A T G C	C	A	T	A	G	G	A	A	C	A	C	A	G	G	C G	T A G	- 3'
24	5'	- G A T G C	A	A	T	A	C	G	G	A	A	G	G	G	G	G	C G	T A G	- 3'
25	5'	- G A T G C	A	A	T	C	A	G	G	A	A	G	G	G	G	G	C G	T A G	- 3'
26	5'	- G A T G C	A	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
27	5'	- G A T G C	A	A	A	A	T	C	G	G	A	A	G	G	G	G	C G	T A G	- 3'
28	5'	- G A T G C	A	A	A	A	T	C	A	G	G	A	A	G	G	G	C G	T A G	- 3'
29	5'	- G A T G C	A	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
30	5'	- G A T G C	A	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
31	5'	- G A T G C	T	A	A	A	G	A	A	G	G	A	A	G	G	G	C G	T A G	- 3'
32	5'	- G A T G C	G	T	T	A	C	G	G	G	G	A	G	G	G	G	C G	T A G	- 3'
33	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
34	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
35	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
36	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
37	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
38	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
39	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
40	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
41	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
42	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
43	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
44	5'	- G A T G C	G	A	A	A	A	A	G	G	G	G	G	G	G	G	C G	T A G	- 3'
45	5'	- G A T G C	T	A	A	A	G	A	A	G	G	A	G	G	G	G	C G	T A G	- 3'
46	5'	- G A T G C	G	G	C	A	C	A	A	G	A	A	G	G	G	G	C G	T A G	- 3'
47	5'	- G A T G C	G	G	G	A	G	A	A	G	A	A	G	G	G	G	C G	T A G	- 3'
48	5'	- G A T G C	G	G	G	A	T	A	G	A	A	G	G	G	G	G	C G	T A G	- 3'
49	5'	- G A T G C	G	G	G	A	T	A	G	A	A	G	G	G	G	G	C G	T A G	- 3'
50	5'	- G A T G C	G	G	G	A	T	A	G	A	A	G	G	G	G	G	C G	T A G	- 3'
51	5'	- G A T G C	G	G	G	A	T	A	G	A	A	G	G	G	G	G	C G	T A G	- 3'
52	5'	- G A T G C	G	G	G	A	T	A	G	A	A	G	G	G	G	G	C G	T A G	- 3'
53	5'	- G A T G C	G	G	G	A	T	A	G	A	A	G	G	G	G	G	C G	T A G	- 3'
54	5'	- G A T G C	G	G	G	A	T	A	G	A	A	G	G	G	G	G	C G	T A G	- 3'

Selected sites were not aligned!

My program claims that the site enriched in my peak is a binding site for factor X. Is it?

Profile summary	
Name:	Klf4
Matrix ID:	MA0039.1
Class:	C2H2 zinc finger factors
Family:	Three-zinc finger Kruppel-related factors
Collection:	CORE
Taxon:	Vertebrates
Species:	Mus musculus
Data Type:	SELEX
Validation:	9443972
Uniprot ID:	Q60793



My program claims that the site enriched in my peak is a binding site for factor X. Is it?

Circular References

Experimental binding site for a gene family member...



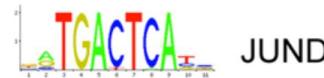
is cited in a review on entire family with a comma-separated list of family members



which creates new entry in a database for a new site for all other family member

Gene family members with similar sites

JUN-like



JUND



JUN



AP1



FOSL2



BATF



FOSL1

Peak annotation



ChIP seq

400-170,000

ATAC seq

100,000-300,000



Peak annotation

1. Descriptive: where binding events occur relative to known genomic features (genes, TSS, transcripts, CpG islands, binding sites of other transcription factors, any other “locations” of your choosing)

Easy – overlay your peak coordinates with those of known genome features

2. Predictive : How does a binding event near feature X affects the function of this feature or any other phenotype of interest?

Difficult

Predictive Peak annotation

- Noise and low-affinity interactions
- Correlation between binding events and transcriptional response is often lacking in a point expression study.

" On average, **14.7%** of genes bound by a factor were differentially expressed following the knockdown of that factor, suggesting that most interactions between TF and chromatin do not result in measurable changes in gene expression levels of putative target genes. "

Cussanovich et all. (2014) PlosGenetics 10(3): e1004226

- Correlation between the number of peaks per gene and expression does exist
- Correlation between specific histone modifications and expression exists, but is it quantitative?

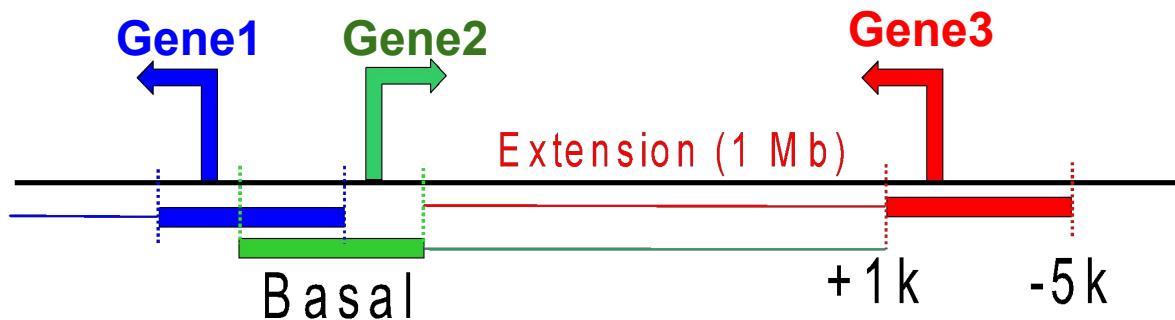
Peak Annotation

- “Naïve” TSS proximity – all peaks within X kb of a TSS
- “Augmented” TSS proximity – all peaks within X kb of TSS and
 - with a binding site of a factor in question
 - ranked by peak scores or p-values
 - overlapping other transcription factors binding sites
 - overlapping open chromatin
 - overlapping peaks of your target from different cells.
- All peak in a “regulatory region” belong to all genes in a regulatory region.
 - experimental topological domain
 - educated guess
- Distance-dependent stratification.

GREAT: Genomic Regions Enrichment of Annotations Tool

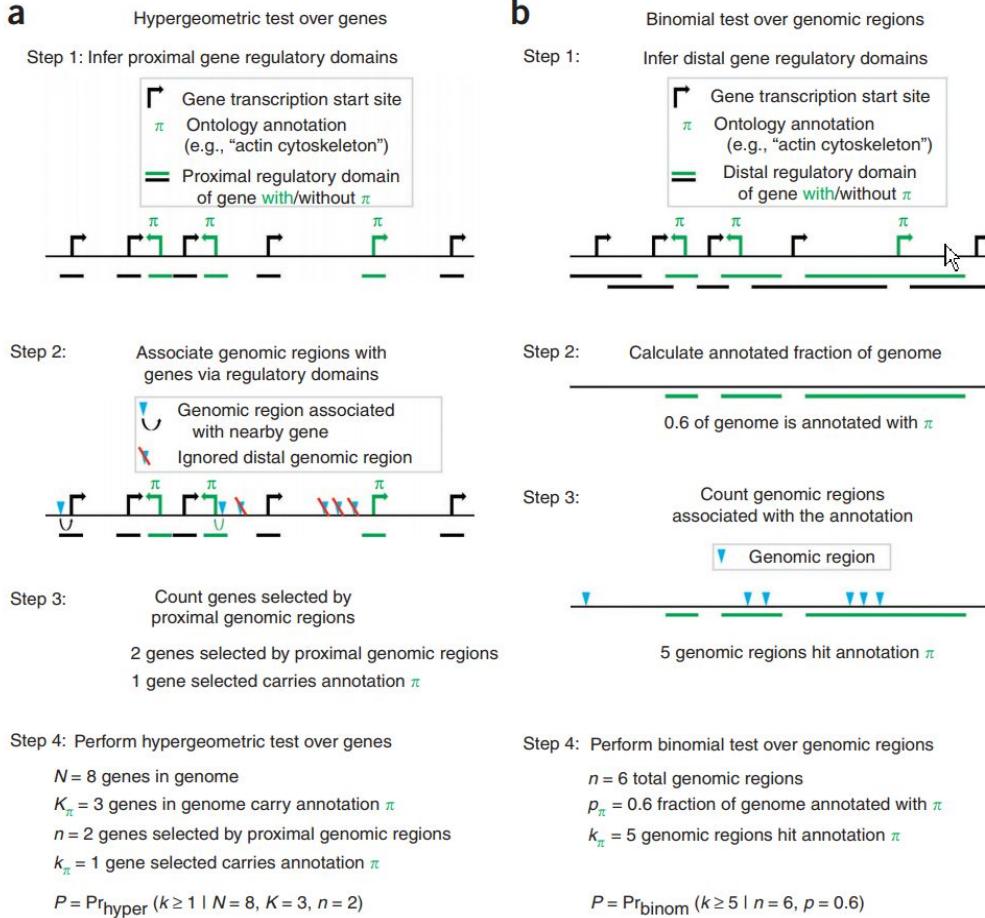


User-defined “Regulatory regions”



Uses experimentally-defined “regulatory regions”, when known

Peak Annotation: GREAT



Nature Biotechnology 28, 495–501 (2010)

Ontology annotation of Regulatory regions

Calculate annotated fraction of genome

Counts peaks that hit “annotated” portion of genome

Peak Annotation: GREAT



GREAT Overview News Use GREAT Demo Video How to Cite Help Forum

GREAT version 3.0.0 current (02/15/2015 to now) ▾

All genomic region-gene association tables (2526 regions, 2944 genes)

Job ID: 20171115-public-3.0.0-FWs9W9

Display name: 02_GR_D_rep2.bed4.bed

What do these tables show?

Genomic region -> gene association table [Download table as text](#).

Region	Gene (distance to TSS)
GR_D_Peak1	Arfgef1 (-22,746), Cpa6 (+464,529)
GR_D_Peak2	Prdm14 (-184,573), Ncoa2 (+62,347)
GR_D_Peak3	Ncoa2 (-80,580), Tram1 (+135,201)
GR_D_Peak4	Ube2w (-22,863), Tceb1 (+14,716)

Gene -> genomic region association table [Download t](#)

Gene	Region (distance to TSS)
0610010F05Rik	GR_D_Peak1541 (-28,961)
0610039K10Rik	GR_D_Peak376 (-2,001)
0610040J01Rik	GR_D_Peak750 (+126,350)
1110004E09Rik	GR_D_Peak2192 (-66,452)

GeneOntology

Mouse Phenotype (MGI)

Human Phenotype (OMIM)

Disease Ontology

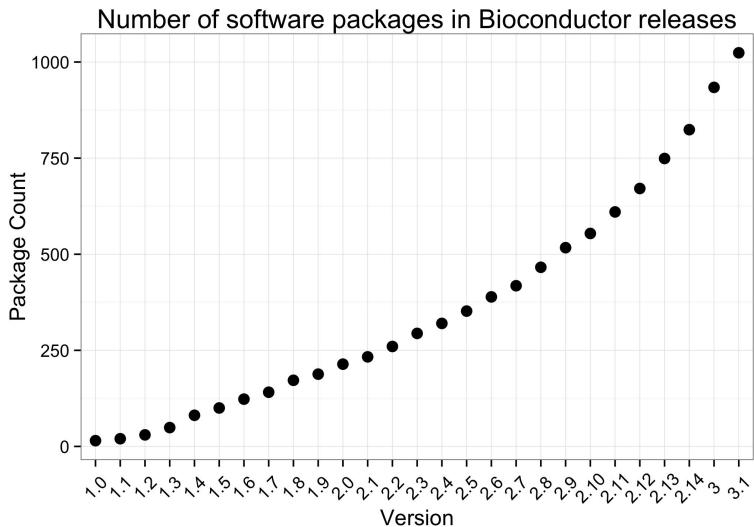
MSigDb (Broad Institute)

PANTHER (UCLA)

InterPro (EMBL)

TreeFAM(EMBL)

Bioconductor: a repository of R tools for the analysis of genomic data.



Current version is 3.6 >1400 packages
<http://bioconductor.org/packages/release/BiocViews.html#Software>

Is updated twice a year

<http://bioconductor.org/>

R environment

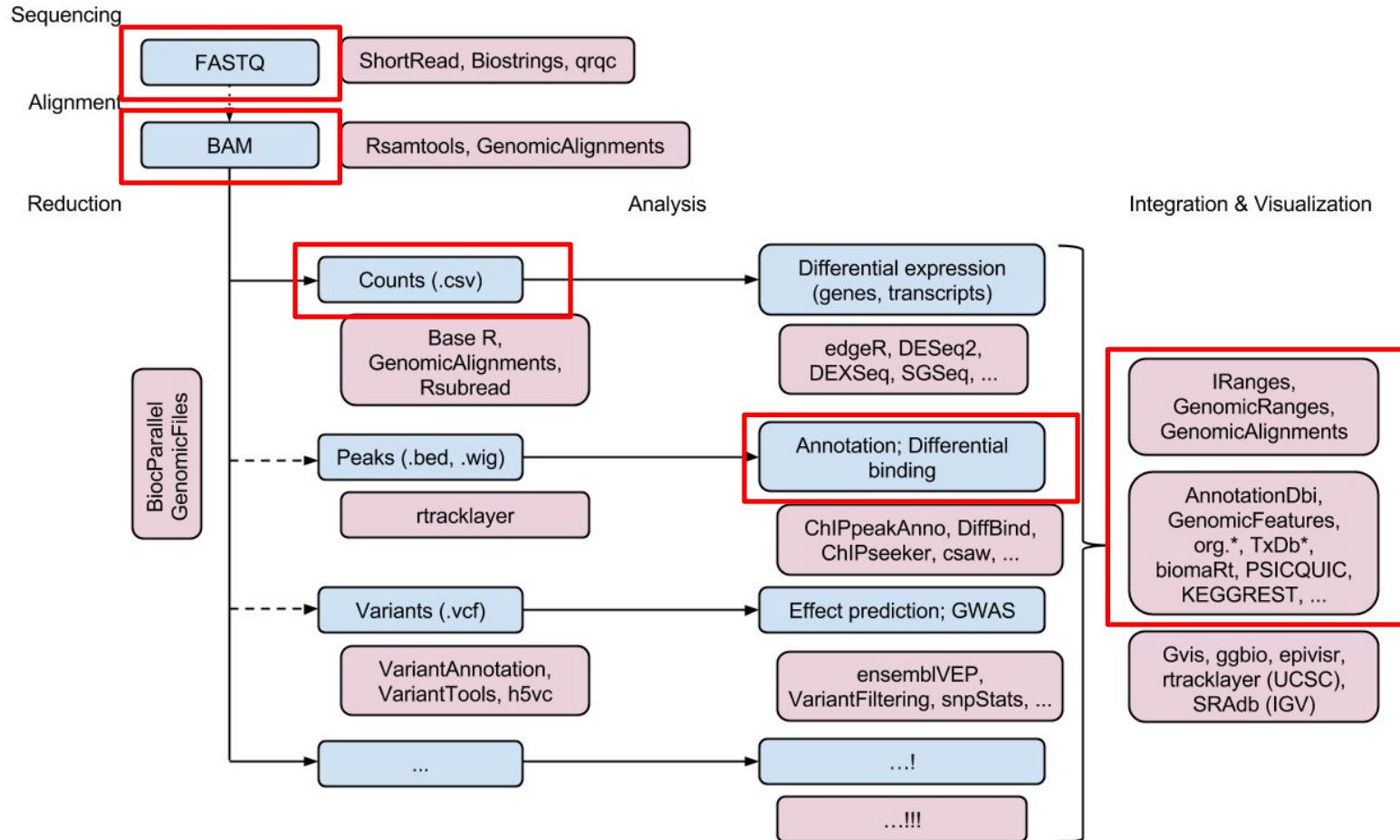
“Compatibility” with ~7000 packages in CRAN

Consistent data structure and access methods

Large user base and effective support

“Good” documentation (but not for all packages)

Bioconductor – (very) general overview



Explore Bioconductor: BioViews

https://bioconductor.org/packages/release/BiocViews.html#_Software

All Packages

Bioconductor version 3.6 (Release)

Autocomplete biocViews search:

▼ Software (1476)

- ▶ AssayDomain (573)
- ▶ BiologicalQuestion (559)
- ▶ Infrastructure (322)
- ▶ ResearchField (413)
- ▶ StatisticalMethod (487)
- ▶ Technology (933)
- ▶ WorkflowStep (774)

▼ AnnotationData (909)

- ▶ ChipManufacturer (388)
- ▶ ChipName (195)
- ▶ CustomArray (2)
- ▶ CustomDBSchema (3)
- ▶ FunctionalAnnotation (14)
- ▶ Organism (595)
- ▶ PackageType (642)
- ▶ SequenceAnnotation (1)



Packages found under Software:

Show All ▾ entries

Package	Maintainer
a4	Tobias Verbeke, Willem Ligtenberg
a4Base	Tobias Verbeke, Willem Ligtenberg
a4Classif	Tobias Verbeke, Willem Ligtenberg
a4Core	Tobias Verbeke, Willem Ligtenberg
a4Preproc	Tobias Verbeke, Willem Ligtenberg
a4Reporting	Tobias Verbeke, Willem Ligtenberg
ABAEnrichment	Steffi Grote

Annotation packages in Bioconductor

TxDb family (eg: `TxDb.Hsapiens.UCSC.hg19.knownGene`) –
Transcriptome coordinates for the known genes (introns, exons, UTRs)

BSgenome family (eg: `BSgenome.Hsapiens.UCSC.hg19`) - complete genome sequence

orgdb family (eg: `org.Hs.eg.db`) gene based gene ID mapper, plus location information

biomaRt – ENSEMBL-centric annotation package for both model and non-model organisms.

ChIP-Seq/Atac-Seq - related packages

ChIPQC - quality controls for Chip-seq experiments

QC Summary

Table 1. Summary of ChIP-seq filtering and quality metrics.

ID	Tissue	Factor	Condition	Replicate	Reads	Dup%	ReadL	FragL	RelCC	SSD	RiP%	RiBL%
BRD4_L_rep1	BRD4		1	3340300	0	51	175	10	0.99	3.2	0.32	
BRD4_L_rep2	BRD4		2	3723935	0	51	160	13	0.93	1.7	0.68	
BRD4_LD_rep1	BRD4		1	3290166	0	51	178	8	1	3.9	0.29	
BRD4_LD_rep2	BRD4		2	3390004	0	51	142	11	0.83	0.55	0.75	
BRD4_U_rep1	BRD4		1	3249236	0	51	163	12	1.1	2.3	0.3	
BRD4_U_rep2	BRD4		2	3685080	0	51	142	9.2	0.92	0.65	0.74	
GR_D_rep1	GR		1	6219988	0	51	267	3	0.77	1	0.55	
GR_D_rep3	GR		2	2814984	0	51	172	5.1	0.82	1.2	0.53	
GR_L_rep1	GR		1	5861427	0	51	267	4.7	0.71	0.12	0.53	
GR_L_rep3	GR		2	2648221	0	51	175	4	0.64	0.067	0.58	
GR_LD_rep1	GR		1	5870854	0	51	26					
GR_LD_rep3	GR		2	3669465	0	51	173					
GR_U_rep1	GR		1	6221929	0	51	26					
GR_U_rep3	GR		2	3065535	0	51	160					
Nelf1_LPSoh	NELF		1	2612013	0	49	194					
Nelf1_LPSih	NELF		1	2605162	0	49	185					
Nelf1_LPS3h	NELF		1	2876635	0	49	192					
Nelf3_LD_1h	NELF		1	2257041	0	49	218					
p65_L_rep1	p65		1	3546168	0	51	170					
p65_L_rep2	p65		2	4770660	0	51	18					
p65_LD_rep1	p65		1	3546954	0	51	168					
p65_LD_rep2	p65		2	3847697	0	51	16					

Summary table

Number of Reads

% Duplicates

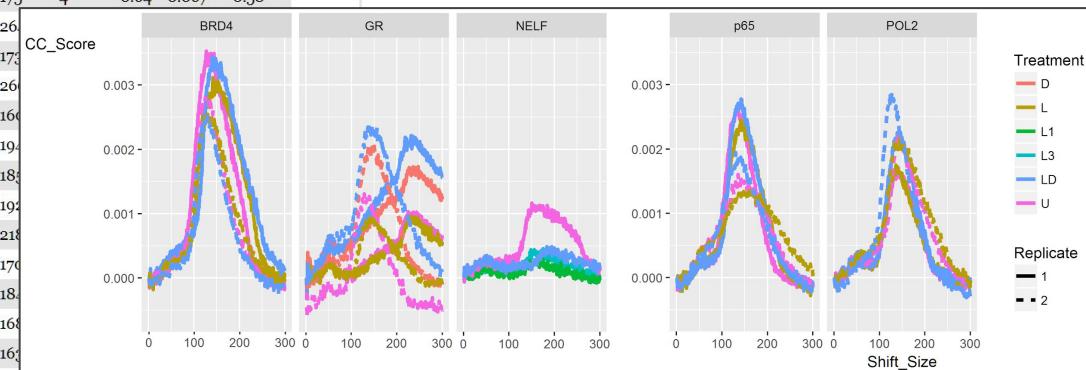
Read Length

Fragment Length

Relative Cross Coverage

Frip%

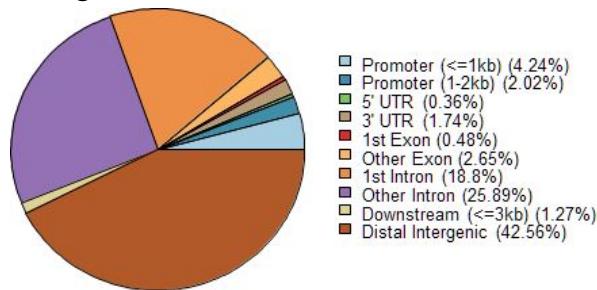
Calculates Cross-Correlation but slow
Provides IRD filtering (very conservative)



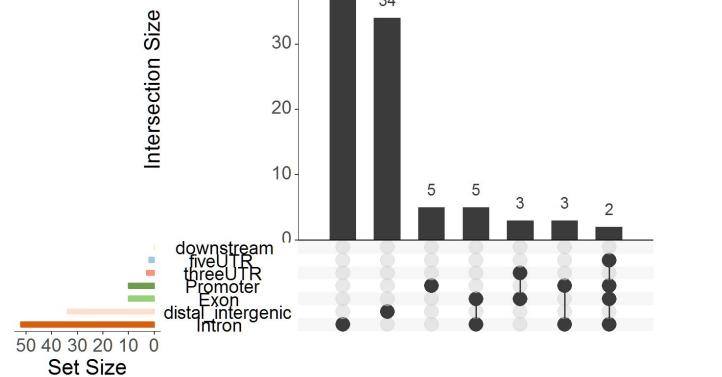
ChIPSeeker , ChIPpeakAnno peak summarization and annotation



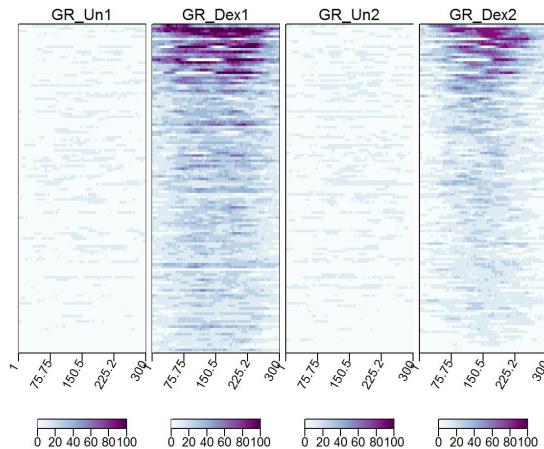
Annotation relative to genomic features



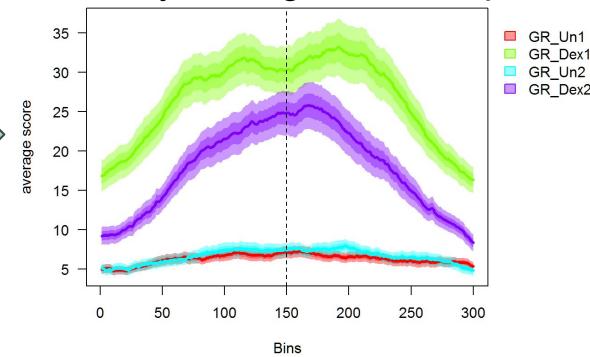
Heatmap of read densities in each peak



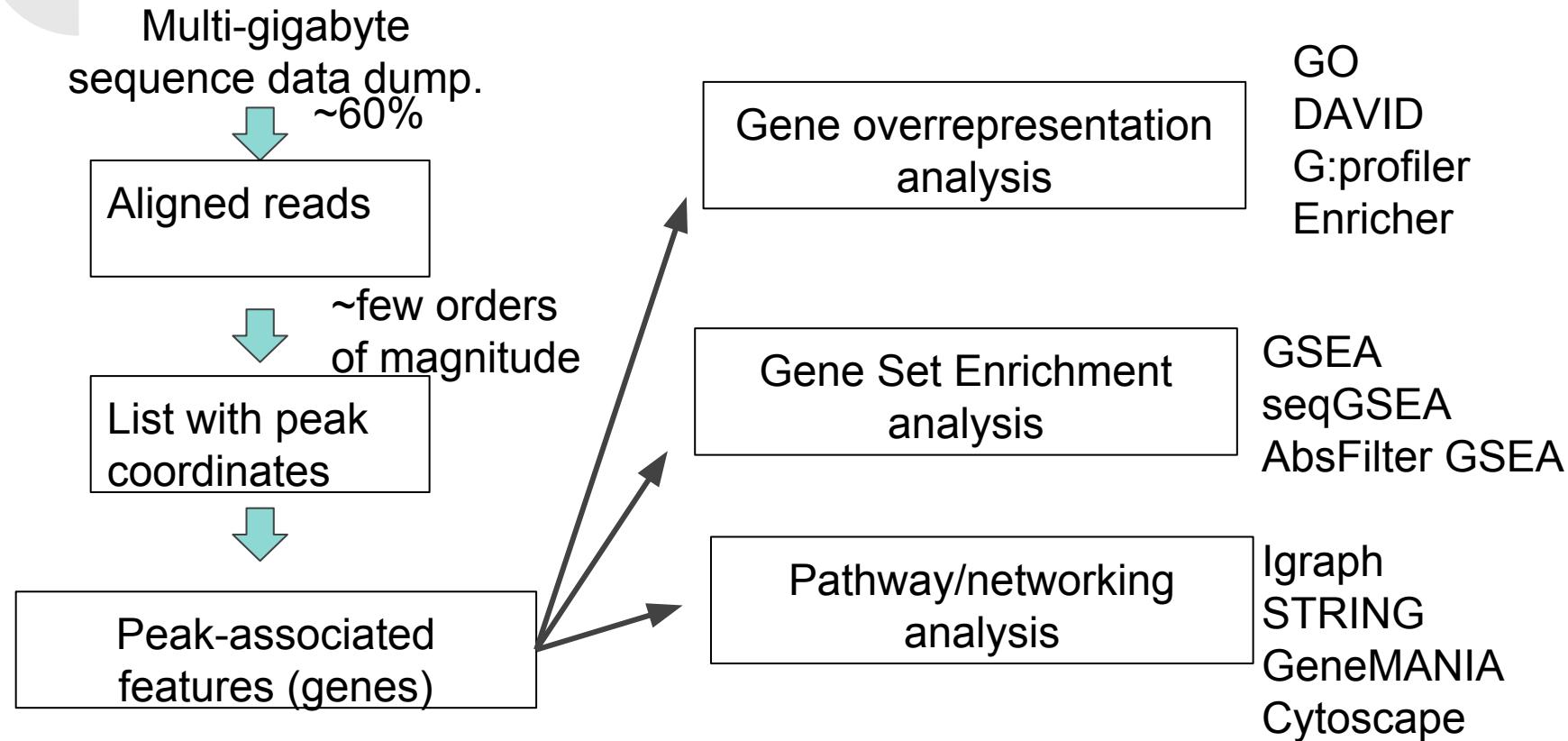
Genomation - peak summarization and annotation



Read density averaged over all peaks



From peak annotation to gene list - data reduction and summarization



Acknowledgements:



Hospital for Special Surgery

Funding:
The Tow Foundation
American Heart Association
NIH