Peaks

DNAse I-seq

CHIP-seq

SIGNAL PROCESSING FOR NEXT-GEN SEQUENCING DATA

Gene models

RNA-seq

RIP/CLIP-seq

Binding sites



Transcripts

The Power of Next-Gen Sequencing



For more Seq technologies, see: <u>http://liorpachter.wordpress.com/seq/</u>

Next-Gen Sequencing as Signal Data



- ✓ Map reads (red) to the genome. Whole pieces of DNA are black.
- ✓ Count # of reads mapping to each DNA base → signal

Outline

- Read mapping: Creating signal map
- Finding enriched regions
 - CHIP-seq: peaks of protein binding



• **RNA-seq**: from enrichment to transcript quantification



 Application: Predicting gene expression from transcription factor and histone modification binding

Read mapping

- Problem: match up to a billion short sequence reads to the genome
- Need sequence alignment algorithm faster than BLAST



Read mapping (sequence alignment)

- Dynamic programming
 - Optimal, but SLOW
- BLAST
 - Searches primarily for close matches, still too slow for high throughput sequence read mapping
- Read mapping
 - Only want very close matches, must be super fast

Index-based short read mappers

- Similar to BLAST
- Map all genomic locations of all possible short sequences in a hash table
- Check if read subsequences map to adjacent locations in the genome, allowing for up to 1 or 2 mismatches.
- Very memory intensive!



Trapnell and Salzberg 2009, Slide adapted from Ray Auerbach

Read Alignment using Burrows-Wheeler Transform



- Used in Bowtie, the current most widely used read aligner
- Described in Coursera course: Bioinformatics Algorithms (part 1, week 10)

Trapnell and Salzberg 2009, Slide adapted from Ray Auerbach

Read mapping issues

- Multiple mapping
- Unmapped reads due to sequencing errors
- VERY computationally expensive
 - Remapping data from The Cancer Genome Atlas consortium would take 6 CPU years¹
- Current methods use heuristics, and are not 100% accurate
- These are open problems



FINDING ENRICHED REGIONS: CHIP-SEQ DATA ANALYSIS



- Determine locations of transcription factors and histone modifications.
- The binding of these factors is what regulates whether genes get transcribed.

CHIP-seq protocol

DNA bound by histones and transcription factors



CHIP-seq Data



Basic interpretation: Signal map to represents binding profile of protein to DNA

How do we identify binding sites from CHIP-seq signal "peaks"?

Park 2009 Nature Reviews Genetics



- Background assumption: all sequence reads map to random locations within the genome
- Divide genome into bins, distribution of expected frequencies of reads/bin is described by the Poisson distribution.
- Assign p-value based on Poisson distribution for each bin based on # of reads

Is a Poisson background reasonable for CHIP-seq data?



 "Input" is from a CHIP-seq experiment using an antibody for a non-DNA binding protein

ENCODE NF-Kb CHIP-seq data

Is a Poisson background reasonable for CHIP-seq data?

 "Input" experiment: Do CHIP-seq using an antibody for a protein that doesn't bind DNA



• There are also "peaks" in the input!

Peakseq

Rozowsky *et al.* 2009 *Nature Biotech* Gerstein Lab

Determining protein binding sites by comparing CHIP-seq data with input



Candidate binding site identification



- Use Poisson distribution as background, as in the "naïve" analysis discussed earlier
- Normalize read counts for mappability (uniqueness) of genomic regions
- Use large bin size, finer resolution analysis later

Input normalization



✓ Normalize based on slope of least squares regression line.
Normalized reads = CHIP-seq reads/(slope*input reads)

Input normalization

All data points

Candidate peaks removed



Using regression based on all data points (including candidate peaks) is overly conservative.



Enriched target sites

- Binomial distribution
 - Each genomic region is like a coin
 - The combined number of reads is the # of times that the coin is flipped
 - Look for regions that are "weighted" toward sample, not input

ENCODE NF-Kb CHIP-seq data

Multiple Hypothesis Correction

- Millions of genomic bins → expect many bins with p-value < 0.05!
- How do we correct for this?

Multiple Hypothesis Correction

- Bonferroni Correction
 - Multiply p-value by number of observations
 - Adjusts p-values \rightarrow expect up to 1 false positive
 - Very conservative

Multiple Hypothesis Correction

- False discovery rate (FDR)
 - Expected number of false positives as a percentage of the total rejected null hypotheses
 - Expectation[false positives/(false positives+true postives)]
- q-value: maximum FDR at which null hypothesis is rejected.
- Benjamini-Hochberg Correction

– q-value = p-value*# of tests/rank

Is PeakSeq an optimal algorithm?

Many other CHIP-seq "peak"-callers

Program	1	aterence ve	Color Co	South States	use in the second	A Contract	n warned	Service and servic	and a top and	Consensation	ston to to	Disconstructures	10 Contractor	a substant rade substant and substant and	orvest
CisGenome	28	1.1	X*	х				х	х		x		х	conditional binomial model	
Minimal ChipSeq Peak Finder	16	2.0.1			x			х				х			
E-RANGE	27	3.1			x			х				х	х	chromsome scale Poisson dist.	
MACS	13	1.3.5		X				х			X		х	local Poisson dist.	
QuEST	14	2.3				x		х			X**		х	chromsome scale Poisson dist.	
HPeak	29	1.1		X				X					х	Hidden Markov Model	
Sole-Search	23	1	X	X				х		X			х	One sample t-test	
PeakSeq	21	1.01			x			х					х	conditional binomial model	
SISSRS	32	1.4		X			х					X			
spp package (wtd & mtc)	31	1.7		х			х		х	X.	х				
				Generating density profiles			Peak assignment		Adjustments w. control data		Significance relative to control data				

X* = Windows-only GUI or cross-platform command line interface

X** = optional if sufficient data is available to split control data

X' = method exludes putative duplicated regions, no treatment of deletions

Wilbanks EG, Facciotti MT (2010) Evaluation of Algorithm Performance in ChIP-Seq Peak Detection. PLoS ONE 5(7): e11471. doi:10.1371/journal.pone.0011471

http://www.plosone.org/article/info:doi/10.1371/journal.pone.0011471



CHIP-seq summary

- Method to determine DNA binding sites of transcription factors or locations of histone modifications
- Must normalize sequence reads to experimental input
- Search for signal enrichment to find **peaks**
 - Peakseq: binomial test + Benjamini-Hochberg correction
 - Many other methods

RNA-SEQ: GOING BEYOND ENRICHMENT

RNA-seq

• Searching for "peaks" not enough:



Nucleotide position

- Are these "peaks" part of the same RNA molecule?
- How much of the RNA is really there?

Wang et al Nature Reviews Genetics 2009

Background: RNA splicing



- pre-mRNA must have introns spliced out before being translated into protein.
- The components that are retained in the mature **mRNA** are called **exons**

Background: alternative splicing



- *Alternative splicing* leads to creation of multiple RNA **isoforms**, with different component exons.
- Sometimes, exons can be retained, or introns can be skipped.

Simple quantification

- Count reads overlapping annotations of known genes
- Simplest method: Make composite model of all isoforms of gene



• Quantification: Reads per kilobase per million reads (RPKM)

Isoform Quantification

- Map reads to genome
- How do we assign reads to overlapping transcripts?



Isoform Quantification

• Simple method: only consider unique reads



Isoform Quantification

- Simple method: only consider unique reads
- Problem: what about the rest of the data?



Expectation Maximization Algorithm

- Assign reads to isoforms to maximize likelihood of generating total pattern of observed reads.
- O. Initialize (expectation): Assign reads randomly to isoforms based on naïve (length normalized) probability of the read coming from that isoform (as opposed to other overlapping isoforms)



Expectation Maximization Algorithm

• 1. Maximization: Choose transcript abundances that maximize likelihood of the read distribution (Maximization).



 2. Expectation: Reassign reads based on the new values for the relative quantities of the isoforms.



Expectation Maximization Algorithm

• 3. Continue **expectation** and **maximization** steps until isoform quantifications converge (it is a mathematical fact that this will happen).



RNA-Seq conclusions

- RNA-Seq is a powerful tool to identify new transcribed regions of the genome and compare the RNA complements of different tissues.
- Quantification harder than CHIP-seq because of RNA splicing
- Expectation maximization algorithm can be useful for quantifying overlapping transcripts