## Lecture 12. Gene Finding & RNAseq

Michael Schatz

March 4, 2020 Applied Comparative Genomics



# Assignment 4: Due Wed Mar 4

#### Assignment 4: Bedtools and Intro to Machine Learning

Assignment Date: Wednesday Feb 19, 2020 Due Date: Wednesday, March 4, 2020 (9 11:59pm

#### Assignment Overview

In this assignment, you will analyze variant data and make different visualization in the language of your choice. (We suggest Python, R, or perhaps Excel.) Make sure to show your work/code in your writeup! As before, any questions about the assignment should be posted to Plazza.

#### Question 1. De novo mutation analysis [20 pts]

For this question, we will be focusing on the de novo variants identified in this paper: http://www.nature.com/articles/hp/genmed201627

Download the de novo variant positions from here (Supplementary Table 54): http://www.nature.com/article-assets/npg/np/genmed/2016/np/genmed/201627/estref/np/genmed/201627-s3.shs

Download the gene annotation of the human genome here: htp://ttp.ensembi.org/pub/velease-87/gff3/homo\_sapiens/Homo\_sapiens.GRCh38.87.gff3.gz

Download the annotation of regulatory variants from here: htp://tip.ensembi.org/pub/release-87/regulation/homo\_sapiens/homo\_sapiens.GRCh38.Regulatory\_Build.regulatory\_features.20161111.gft.gz

Download chromosome 22 from build 38 of the human genome from here: http://hgdownload.cse.ucsc.edu/goldenPath/hg38/chromosomes/chr22.fa.gz

NOTE The variants are reported using version 37 of the reference genome, but the annotation is for version 38. Fortunately, you can 'lift-over' the variants to the coordinates on the new reference genome using several availate tools. I recommend the USCS liftover tool that can do this in batch by converting the variants into BED format. Note, some variants may not successfully lift over, especially if they become repetitive and/or missing in the new reference, so please make a note of how many variants fail liftover.

· Question 1a. How much of the genome is annotated as a gene?



# Assignment 5: Due Wed Mar II

#### Assignment 5: Annotations and RNA-seq

Assignment Date: Wednesday, March 4, 2020 Due Date: Wednesday, March 11, 2020 @ 11:59pm

#### Assignment Overview

In this assignment, you will analyze gene expression data and learn how to make several kinds of plots in the environment of your choice. (We suggest Python or R.) Make sure to show your work/code in your writeup! As before, any questions about the assignment should be posted to Plazza.

#### Question 1. Gene Annotation Preliminaries [10 pts]

Download the annotation of build 38 of the human genome from here: ftp://ftp.ensembi.org/pub/release-87/gtl/homo\_sapiens/Homo\_sapiens.GRCh38.87.gtf.gz

- Question 1a. How many annotated protein coding genes are on each autosome of the human genome? [Hint: Protein coding genes will have "gene" in the 3rd column, and contain the following text: gene, biotype "protein, coding"]
- Question 1b. What is the maximum, minimum, mean, and standard deviation of the span of protein coding genes? [Hint: use the genes identified in 1b]
- Question 1c. What is the maximum, minimum, mean, and standard deviation in the number of exons for protein coding genes? (Hint: you should separately consider each isoform for each protein coding gene)

#### Question 2. Sampling Simulation [10 pts]

A typical human cell has -250,000 transcripts, and a typical bulk RNA-seq experiment may involve millions of cells. Consequently in an RNAseq experiment you may start with trillions of RNA molecules, although your sequencer will only give a few million to billions of reads. Therefore your RNAseq experiment will be a small sampling of the full composition. We hope the sequences will be a representative sample of the total population, but if your sample is very unlucky or biased it may not represent the true distribution. We will explore this concept by sampling a small subset of transcripts (500 to 50000) out of a much larger set (1M) so that you can evaluate this bias.

In data1.bit with 1,000,000 lines we provide an abstraction of RNA-seq data where normalization has been performed and the number of times a gene name occurs corresponds to the number of transcripts in the sample.

### How is Structural Variation Classified?



- Just as with small variant calling, we typically classify structural variation by comparing two
  individuals to each other.
  - There are many ways to do this "comparison" which will be covered later.
- One individual is designated as the "reference" and the other the "query".
- We then define query structural variants "with respect to" the reference
  - Mike has an insertion with respect to the reference
  - Bob has a deletion with respect to the reference
- SVs are usually defined as being longer than 50 bp.

Ho, Steve S., Alexander E. Urban, and Ryan E. Mills. "Structural variation in the sequencing era." Nature Reviews Genetics (2019): 1-19.

### Modeling the function



# Oxford Nanopore MinION





# Assignment 5: Due March 8

#### Assignment 5: Genome Arithmetic

Assignment Date: Thursday, March 1, 2018 Due Date: Thursday, March 8, 2018 @ 11:59pm

#### Assignment Overview

In this assignment, you will call structural variants and analyze the properties of variants in the human genome. Make sure to show your work in your writeup! As before, any questions about the assignment should be posted to Piazza.

#### Question 1. Gene Annotation Preliminaries [10 pts]

Download the annotation of build 38 of the human genome from here: ftp://ftp.ensembl.org/pub/release-87/gtf/homo\_sapiens/Homo\_sapiens.GRCh38.87.gtf.gz

- Question 1a. How many many GTF data lines are in this file? [Hint: The first few lines in the file beginning with "#" are so-called "header" lines describing thing like the creation date, the genome version (more on that later in the course), etc. Header lines should not be counted as data lines.]
- Question 1b. How many annotated protein coding genes are on each autosome of the human genome? [Hint: Protein coding genes will have "gene" in the 3rd column, and contain the following text: gene\_biotype "protein\_coding"]
- Question 1c. What is the maximum, minimum, mean, and standard deviation of the span of protein coding genes? [Hint: use the genes identified in 1b]
- Question 1d. What is the maximum, minimum, mean, and standard deviation in the number of exons for protein coding genes? [Hint: you should separately consider each isoform for each protein coding gene]

### **Goal: Genome Annotations**

atgactatgctaagctgcggctatgctaatgcatgcggctatgctaagctcatgcggctatgctaagctgggaat cgatgacaatgcatgcggctatgctaatgcatgcggctatgcaagctgggatccgatgactatgctaagctgcg gctatgctaatgcatgcggctatgctaagctcatgcgg

## **Goal: Genome Annotations**

gcggctatgctaatgaatggtctt<del>gggatttaccttggaatgctaagctg</del>ggatccgatgacaatgcatgcggct atgctaatgaatggtcttgggatt Gene! gctatgctaagctgggatccgat atgcggctatgcaagctgggatccg atgactatgctaagctgcggctatgctaatgcatgcggctatgctaagctcatgcggctatgctaagctgggaatcgatgacaatgcatgcggctatgctaatgcatgcggctatgcaagctgggatccgatgactatgctaagctgcg

gctatgctaatgcatgcggctatgctaagctcatgcgg

gctatgctaagctgggaatgcatgcg

## **Genetic Code**

#### 1st base

	U.,		c		A		6	6			
U		Phenylalanine Phenylalanine Leucine Leucine	UCU UCC UCA UCG	Serine Serine Serine Serine	UAU UAC UAA UAG	Tyrosine Tyrosine Stop Stop	UGU UGC UGA UGG	Cysteine Cysteine Stop Tryptophan	DCAG		
c	CUU CUC CUA CUG	Leucine Leucine Leucine Leucine	CCU CCC CCA CCG	Proine Proine Proine Proine	CAU CAC CAA CAG	Histidine Histidine Glutamine Glutamine	CGU CGC CGA CGG	Arginine Arginine Arginine Arginine	UCAG		
A	AUU AUC AUA AUG	isoleucine isoleucine Methionine (Start)	ACU ACC ACA ACG	Threonine Threonine Threonine Threonine	AAU AAC AAA AAG	Asparagine Asparagine Lysine Lysine	AGU AGC AGA AGG	Serine Serine Arginine Arginine	UCAG		
G	GUU GUC GUA GUG	Valine Valine Valine	GCU GCC GCA GCG	Alanine Alanine Alanine	GAU GAC GAA GAG	Aspertic Acid Aspertic Acid Glutamic Acid Glutamic Acid	GGU GGC GGA GGG	Glycine Glycine Glycine Glycine	UCAG		



# Outline

- I. Alignment to other genomes
- 2. Prediction aka "Gene Finding"
- 3. Experimental & Functional Assays



# Outline

- I. Alignment to other genomes
- 2. Prediction aka "Gene Finding"
- 3. Experimental & Functional Assays

# Basic Local Alignment Search Tool

- Rapidly compare a sequence Q to a database to find all sequences in the database with an score above some cutoff S.
  - Which protein is most similar to a newly sequenced one?
  - Where does this sequence of DNA originate?
- Speed achieved by using a procedure that typically finds "most" matches with scores > S.
  - Tradeoff between sensitivity and specificity/speed
    - Sensitivity ability to find all related sequences
    - Specificity ability to reject unrelated sequences

(Altschul et al. 1990)

## Seed and Extend

FAKDFLAGGVAAAISKTAVAPIERVKLLLQVQHASKQITADKQYKGIIDCVVRIPKEQGV FLIDLASGGTAAAVSKTAVAPIERVKLLLQVQDASKAIAVDKRYKGIMDVLIRVPKEQGV

- Homologous sequences are likely to contain a short high scoring word pair, a seed.
  - Smaller seed sizes make the sense more sensitive, but also (much) slower
  - Typically do a fast search for prototypes, but then most sensitive for final result
- BLAST then tries to extend high scoring word pairs to compute high scoring segment pairs (HSPs).
  - Significance of the alignment reported via an e-value

## Seed and Extend

- Homologous sequences are likely to contain a short high scoring word pair, a seed.
  - Smaller seed sizes make the sense more sensitive, but also (much) slower
  - Typically do a fast search for prototypes, but then most sensitive for final result
- BLAST then tries to extend high scoring word pairs to compute high scoring segment pairs (HSPs).
  - Significance of the alignment reported via an e-value

## **BLAST E-values**

E-value = the number of HSPs having alignment score S (or higher) expected to occur by chance.

- $\rightarrow$  Smaller E-value, more significant in statistics
- $\rightarrow$  Bigger E-value, less significant
- → Over I means expect this totally by chance (not significant at all!)

The expected number of HSPs with the score at least S is :

 $E = K^* n^* m^* e^{-\lambda S}$ 

K, λ are constant depending on model
n, m are the length of query and sequence
E-values quickly drop off for better alignment bits scores

## Very Similar Sequences

Query: HBA\_HUMAN Hemoglobin alpha subunit Sbjct: HBB HUMAN Hemoglobin beta subunit

Score = 114 bits (285), Expect = 1e-26Identities = 61/145 (42%), Positives = 86/145 (59%), Gaps = 8/145 (5%)

Query 2 LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-----DLSHGSAQV 55 L+P +K+ V A WGKV + E G EAL R+ + +P T+ +F F D G+ +V

Sbjct 3 LTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKV 60

Query 56 KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPA 115 K HGKKV A ++ +AH+D++ + LS+LH KL VDP NF+LL + L+ LA H Sbjct 61 KAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGK 120

Query 116 EFTPAVHASLDKFLASVSTVLTSKY 140 EFTP V A+ K +A V+ L KY

Sbjct 121 EFTPPVQAAYQKVVAGVANALAHKY 145

## Quite Similar Sequences

```
Query: HBA HUMAN Hemoglobin alpha subunit
Sbjct: MYG HUMAN Myoglobin
Score = 51.2 bits (121), Expect = 1e-07,
Identities = 38/146 (26%), Positives = 58/146 (39%), Gaps = 6/146 (4%)
Query 2 LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-----DLSHGSAQV
                                                                       55
                               +G E L R+F
         LS +
                 V
                     WGKV A
                                            РТ
                                                  F
                                                     F
                                                            D
                                                                 S +
Sbjct 3 LSDGEWOLVLNVWGKVEADIPGHGOEVLIRLFKGHPETLEKFDKFKHLKSEDEMKASEDL
                                                                       62
Query 56 KGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPA
                                                                       115
         K HG V AL
                                 + L+ HA K ++
                       +
                                                    + + S C + L + P
Sbjct
      63 KKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKHPG
                                                                       122
Ouery 116 EFTPAVHASLDKFLASVSTVLTSKYR
                                      141
                              + S Y+
           +F
                  +++K L
Sbjct 123 DFGADAQGAMNKALELFRKDMASNYK
                                      148
```

## Not similar sequences

Query: HBA\_HUMAN Hemoglobin alpha subunit Sbjct: SPAC869.02c [Schizosaccharomyces pombe]

Score = 33.1 bits (74), Expect = 0.24 Identities = 27/95 (28%), Positives = 50/95 (52%), Gaps = 10/95 (10%)

Query	30	ERMF	LSFPT	TKTYI	FPHFD	LSH	GSAQ	VKGHGI	KVA	DAL	TNA	VAHVDDM	PNALSAI	SDLHAH	89
		++M	++P		P+F+	+H	+	-	- +A	AL	N	++DD+	+LSA	D	
Sbjct	59	QKML	GNYPE	V]	LPYFN	KAH	QISL	SQPI	RILA	FAL	LNY	AKNIDDL	-TSLSAF	MDQIVV	112
Query	90	K	LRVDP	VNFKI	LLSHC	LLV	TLAA	HLPAEI	-TP	A	120				
		K	L++	++ -	++ HC	LL '	T+	LP++	TP.	A					
Sbjct	113	KHVG	LQIKA	EHYP	IVGHC	LLS	TMQE	LLPSD	ATP.	A	147				

## **Blast Versions**

Program	Database	Query		
BLASTN	Nucleotide	Nucleotide		
BLASTP	Protein	Protein		
BLASTX	Protein	Nucleotide translated into protein		
TBLASTN	Nucleotide translated into protein	Protein		
TBLASTX	Nucleotide translated into protein	Nucleotide translated into protein		

## **NCBI** Blast



- Nucleotide Databases
  - nr:All Genbank
  - refseq: Reference organisms
  - wgs:All reads

- Protein Databases
  - nr:All non-redundant sequences
  - Refseq: Reference proteins



# Outline

- I. Alignment to other genomes
- 2. Prediction aka "Gene Finding"
- 3. Experimental & Functional Assays





## Bacterial Gene Finding and Glimmer (also Archaeal and viral gene finding)

Arthur L. Delcher and Steven Salzberg Center for Bioinformatics and Computational Biology Johns Hopkins University School of Medicine

Step One

• Find open reading frames (ORFs).



Step One

• Find open reading frames (ORFs).



• But ORFs generally overlap ...



All ORFs longer than 100bp on both strands shown - color indicates reading frame Longest ORFs likely to be protein-coding genes

Note the low GC content

All genes are ORFs but not all ORFs are genes



#### Campylobacter jejuni RM1221 30.3%GC

	_			 	 			
		-	-	 		_	_	

#### Mycobacterium smegmatis MC2 67.4%GC

_

### Note what happens in a high-GC genome

#### Mycobacterium smegmatis MC2 67.4%GC


#### *Mycobacterium smegmatis* MC2 67.4%GC

-									
			_		 	_	_		
	· ;	··	· 2 . <del></del>	· · · · · · · · · · · · · · · · · · ·	 .ह. <mark>ज्ञान्त</mark> ्र	· e	yy 17 <mark>.2000</mark> 0	्र ::: <del></del>	ko <u>ko noka</u> ko

P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?

\$ ./coinflip.pl 0.046875 1000

- 0: ННННННННННН 15
- 1: HHHHHHT 7
- 2: HHHHHHHHHH 12
- 3: НННННННННННННННННН 24
- 4: HT 2
- 5: НННННННННН 14
- 6: HHHHHHHHH 10
- 7: НННННННННННТ 14
- 8: HHHHHT 6
- 9: HHHHHHHHH 11
- 11: НННННННННННННННННННННННННННННННННН
- 13: HHHT 4
- 14: ННННННННННН 15
- 15: ННННННННННННННННННННННННННННННННН
- 16: HHHHHT 6
- 17: НННННННННННННННННННННННННННННННН
- 18: ННННННННННННННННННННН 26
- 19: HHHHHHHHHH 12

P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?



P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?

Geometric Distribution:  $P(X=x) = p_{heads}^{x-1}p_{tails}$ 



Flips until heads

P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?

Geometric Distribution:  $P(X=x) = p_{heads}^{x-1}p_{tails}$ 



Flips until heads

P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?

Geometric Distribution:  $P(X=x) = p_{heads}^{x-1}p_{tails}$ 



Flips until heads

## **Stop Codon Frequencies**



If the sequence is mostly A+T, then likely to form stop codons by chance!

#### In High A+T (Low G+C):

Frequent stop codons; Short Random ORFs; long ORFs likely to be true genes

#### In High G+C (Low A+T):

Rare stop codons; Long Random ORFs; harder to identify true genes

#### A relationship between GC content and coding-sequence length.

Oliver & Marín (1996) J Mol Evol. 43(3):216-23.

## Probabilistic Methods

- Create models that have a probability of generating any given sequence.
  - Evaluate gene/non-genome models against a sequence
- Train the models using examples of the types of sequences to generate.

- Use RNA sequencing, homology, or "obvious" genes

- The "score" of an orf is the probability of the model generating it.
  - Most basic technique is to count how kmers occur in known genes versus intergenic sequences
  - More sophisticated methods consider variable length contexts, "wobble" bases, other statistical clues



# Overview of Eukaryotic Gene Prediction

CBB 231 / COMPSCI 261

W.H. Majoros



### **Eukaryotic Gene Syntax**



Regions of the gene outside of the CDS are called *UTR*'s (*untranslated regions*), and are mostly ignored by gene finders, though they are important for regulatory functions.



### Representing Gene Syntax with ORF Graphs

After identifying the most promising (i.e., highest-scoring) signals in an input sequence, we can apply the gene syntax rules to connect these into an *ORF graph*:



An ORF graph represents all possible *gene parses* (and their scores) for a given set of putative signals. A *path* through the graph represents a single gene parse.



### **Conceptual Gene-finding Framework**

TATTCCGATCGATCGATCTCTCTAGCGTCTACG CTATCATCGCTCTCTATTATCGCGCGATCGTCG ATCGCGCGAGAGTATGCTACGTCGATCGAATTG

> identify most promising signals, score signals and content regions between them; induce an ORF graph on the signals



find highest-scoring path through ORF graph; interpret path as a gene parse = gene structure







- Similar to Markov models used for prokaryotic gene finding, but system may transition between multiple models called states (gene/non-gene, intergenic/exon/intron)
- Observers can see the emitted symbols of an HMM (i.e., nucleotides) but have no ability to know which state the HMM is currently in.
  - But we can *infer* the most likely hidden states of an HMM based on the given sequence of emitted symbols.



AAAGCATGCATTTAACGTGAGCACAATAGATTACA



## Eukaryotic Gene Finding with GlimmerHMM

Mihaela Pertea Associate Professor JHU

## **HMMs and Gene Structure**

- Nucleotides {A,C,G,T} are the observables
- Different states generate nucleotides at different frequencies

A simple HMM for unspliced genes:



#### AAAGC ATG CAT TTA ACG AGA GCA CAA GGG CTC TAA TGCCG

 The sequence of states is an annotation of the generated string – each nucleotide is generated in intergenic, start/stop, coding state

## **HMMs & Geometric Feature Lengths**

$$P(x_0...x_{d-1} \mid \theta) = \left(\prod_{i=0}^{d-1} P_e(x_i \mid \theta)\right) p^{d-1}(1-p)$$



### **Generalized HMMs Summary**

• GHMMs generalize HMMs by allowing each state to emit a subsequence rather than just a single symbol

• Whereas HMMs model all feature lengths using a geometric distribution, coding features can be modeled using an arbitrary length distribution in a GHMM

• Emission models within a GHMM can be any arbitrary probabilistic model ("submodel abstraction"), such as a neural network or decision tree

• GHMMs tend to have many fewer states => simplicity & modularity

### **GlimmerHMM** architecture



### Coding vs Non-coding

A three-periodic ICM uses three ICMs in succession to evaluate the different codon positions, which have different statistics:

## P[C|M<sub>0</sub>] P[G|M<sub>1</sub>] P[A|M<sub>2</sub>] ICM<sub>1</sub> ICM<sub>2</sub> ATC GAT CGA TCA GCT TAT CGC ATC

The three ICMs correspond to the three phases. Every base is evaluated in every phase, and the score for a given stretch of (putative) coding DNA is obtained by multiplying the phase-specific probabilities in a mod 3 fashion:

$$\prod_{i=0}^{L-1} P_{(f+i)(\mathrm{mod}3)}(x_i)$$

GlimmerHMM uses 3-periodic ICMs for coding and homogeneous (non-periodic) ICMs for noncoding DNA.

## **Signal Sensors**

Signals – short sequence patterns in the genomic DNA that are recognized by the cellular machinery.



### **Identifying Signals In DNA**

We slide a fixed-length model or "window" along the DNA and evaluate score(signal) at each point:



When the score is greater than some threshold (determined empirically to result in a desired sensitivity), we remember this position as being the potential site of a signal.

The most common signal sensor is the Position Weight Matrix:

	<b>A</b> 100%	<b>T</b> 100%	<b>G</b> 100%	A = 19% T = 20% C = 29% G = 32%	A = 24% T = 18% C = 26% G = 32%
--	------------------	------------------	------------------	--	--

## Splice site prediction



The splice site score is a combination of:

- first or second order inhomogeneous Markov models on windows around the acceptor and donor sites
- Maximal dependence decomposition (MDD) decision trees
- longer Markov models to capture difference between coding and noncoding on opposite sides of site (optional)
- maximal splice site score within 60 bp (optional)

## **GlimmerHMM** architecture



### Gene Prediction with a GHMM

Given a sequence S, we would like to determine the parse  $\phi$  of that sequence which segments the DNA into the most likely exon/intron structure:



The parse  $\phi$  consists of the coordinates of the predicted exons, and corresponds to the precise sequence of states during the operation of the GHMM (and their duration, which equals the number of symbols each state emits).

This is the same as in an HMM except that in the HMM each state emits bases with fixed probability, whereas in the GHMM each state emits an entire feature such as an exon or intron.

### **Evaluation of Gene Finding Programs**

### Nucleotide level accuracy



Sensitivity:

Specificity:

 $Sn = \frac{TP}{TP + FN}$  $Sp = \frac{TP}{TP + FP}$ 

What fraction of reality did you predict?

What fraction of your predictions are real?

## **More Measures of Prediction Accuracy**

### Exon level accuracy



$$ExonSn = \frac{TE}{AE} = \frac{\text{number of correct exons}}{\text{number of actual exons}}$$
$$ExonSp = \frac{TE}{PE} = \frac{\text{number of correct exons}}{\text{number of predicted exons}}$$

## GlimmerHMM is a high-performance ab initio gene finder

Arabidopsis thaliana test results

	Nu	clec	otide		Exo	n	Gene			
	Sn	Sp	Acc	Sn	Sp	Acc	Sn	Sp	Acc	
GlimmerHMM	97	99	98	84	89	86.5	60	61	60.5	
SNAP	96	99	97.5	83	85	84	60	57	58.5	
Genscan+	93	99	96	74	81	77.5	35	35	35	

- All three programs were tested on a test data set of 809 genes, which did not overlap with the training data set of GlimmerHMM.
- All genes were confirmed by full-length Arabidopsis cDNAs and carefully inspected to remove homologues.

## **GlimmerHMM on human data**

	Nuc Sens	Nuc Spec	Nuc Acc	Exon Sens	Exon Spec	Exon Acc	Exact Genes
GlimmerHMM	86%	72%	79%	72%	62%	67%	17%
Genscan	86%	68%	77%	69%	60%	65%	13%

GlimmerHMM's performace compared to Genscan on 963 human RefSeq genes selected randomly from all 24 chromosomes, non-overlapping with the training set. The test set contains 1000 bp of untranslated sequence on either side (5' or 3') of the coding portion of each gene.

# Gene Finding Overview

- Prokaryotic gene finding distinguishes real genes and random ORFs
  - Prokaryotic genes have simple structure and are largely homogenous, making it relatively easy to recognize their sequence composition
- Eukaryotic gene finding identifies the genome-wide most probable gene models (set of exons)
  - "Probabilistic Graphical Model" to enforce overall gene structure, separate models to score splicing/transcription signals
  - Accuracy depends to a large extent on the quality of the training data



# Outline

- I. Alignment to other genomes
- 2. Prediction aka "Gene Finding"
- 3. Experimental & Functional Assays

## Sequencing Assays

#### The \*Seq List (in chronological order)

- 1. Gregory E. Crawford et al., "Genome-wide Mapping of DNase Hypersensitive Sites Using Massively Parallel Signature Sequencing (MPSS)," Genome Research 16, no. 1 (January 1, 2006): 123–131, doi:10.1101/gr.4074106.
- 2. David S. Johnson et al., "Genome-Wide Mapping of in Vivo Protein-DNA Interactions," Science 316, no. 5830 (June 8, 2007): 1497–1502, doi:10.1126/science.1141319.
- 3. Tarjei S. Mikkelsen et al., "Genome-wide Maps of Chromatin State in Pluripotent and Lineage-committed Cells," Nature 448, no. 7153 (August 2, 2007): 553–560, doi:10.1038/nature06008.
- 4. Thomas A. Down et al., "A Bayesian Deconvolution Strategy for Immunoprecipitation-based DNA Methylome Analysis," Nature Biotechnology 26, no. 7 (July 2008): 779–785, doi:10.1038/nbt1414.
- 5. Ali Mortazavi et al., "Mapping and Quantifying Mammalian Transcriptomes by RNA-Seq," Nature Methods 5, no. 7 (July 2008): 621–628, doi:10.1038/nmeth.1226.
- 6. Nathan A. Baird et al., "Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers," PLoS ONE 3, no. 10 (October 13, 2008): e3376, doi:10.1371/journal.pone.0003376.
- 7. Leighton J. Core, Joshua J. Waterfall, and John T. Lis, "Nascent RNA Sequencing Reveals Widespread Pausing and Divergent Initiation at Human Promoters," Science 322, no. 5909 (December 19, 2008): 1845–1848, doi:10.1126/science.1162228.
- 8. Chao Xie and Martti T.Tammi, "CNV-seq, a New Method to Detect Copy Number Variation Using High-throughput Sequencing," BMC Bioinformatics 10, no. 1 (March 6, 2009): 80, doi:10.1186/1471-2105-10-80.
- 9. Jay R. Hesselberth et al., "Global Mapping of protein-DNA Interactions in Vivo by Digital Genomic Footprinting," Nature Methods 6, no. 4 (April 2009): 283–289, doi:10.1038/nmeth.1313.
- Nicholas T. Ingolia et al., "Genome-Wide Analysis in Vivo of Translation with Nucleotide Resolution Using Ribosome Profiling," Science 324, no. 5924 (April 10, 2009): 218–223, doi:10.1126/science.1168978.
- 11. Alayne L. Brunner et al., "Distinct DNA Methylation Patterns Characterize Differentiated Human Embryonic Stem Cells and Developing Human Fetal Liver," Genome Research 19, no. 6 (June 1, 2009): 1044–1056, doi:10.1101/gr.088773.108.
- 12. Mayumi Oda et al., "High-resolution Genome-wide Cytosine Methylation Profiling with Simultaneous Copy Number Analysis and Optimization for Limited Cell Numbers," Nucleic Acids Research 37, no. 12 (July 1, 2009): 3829–3839, doi:10.1093/nar/gkp260.
- 13. Zachary D. Smith et al., "High-throughput Bisulfite Sequencing in Mammalian Genomes," Methods 48, no. 3 (July 2009): 226–232, doi:10.1016/j.ymeth.2009.05.003.
- 14. Andrew M. Smith et al., "Ouantitative Phenotyping via Deep Barcode Sequencing." Genome Research (lulv 21. 2009).



## \*-seq in 4 short vignettes







**Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications.** Sørlie et al (2001) *PNAS*. 98(19):10869-74.

## **RNA-seq** Overview



## **RNA-seq** Overview



## **RNA-seq** Overview



Downstream analysis

## **RNA-seq Challenges**



Challenge I: Eukaryotic genes are spliced

# **RNA-Seq Approaches**



Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the genome with a gapped mapper. Next (novel) transcript discovery and quantification can proceed with or without an annotation file. Novel transcripts are then functionally annotated. **b** If no novel transcript discovery is needed, reads can be mapped to the reference transcriptome using an ungapped aligner. Transcript identification and quantification can occur simultaneously. **c** When no genome is available, reads need to be assembled first into contigs or transcripts. For quantification, reads are mapped back to the novel reference transcriptome and further analysis proceeds as in (**b**) followed by the functional annotation of the novel transcripts as in (**a**). Representative software that can be used at each analysis step are indicated in *bold text*. Abbreviations: *GFF* General Feature Format, *GTF* gene transfer format, *RSEM* RNA-Seq by Expectation Maximization

#### A survey of best practices for RNA-seq data analysis Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8

# **RNA-Seq Approaches**



ated in bold text. Abbreviations: GFF General Feature Format, GTF gene transfer format,

A survey of best practices for RNA-seq data analysis Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8

software that can be used a

RSEM RNA-Seg by Expectation Maximization