# CAP 5510: Introduction to Bioinformatics

## Giri Narasimhan

ECS 254; Phone: x3748
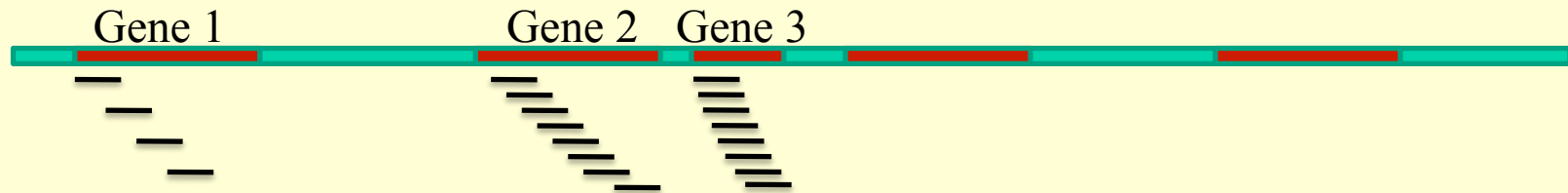
giri@cis.fiu.edu

www.cis.fiu.edu/~giri/teach/BioinfS11.html

# Applications of NGS

- RNA-Seq
- ChIP-Seq
- SNP-Seq
- Metagenomics
- Alternative Splicing
- Copy Number Variations (CNV)
- …

# RNA-Seq



- ## Align reads to genes and count
- ## Assume uniform sampling
  - Count of number of reads mapped per gene is a measure of its expression level
  - Expression of Gene 2 is twice that of Gene 1
  - Expression of Gene 3 is twice that of Gene 2

# Expression Level of Gene

❑ RPKM = Ng / (N X L)

- Ng = Number of reads mapped to gene
- N = Total number of mapped reads (in millions)
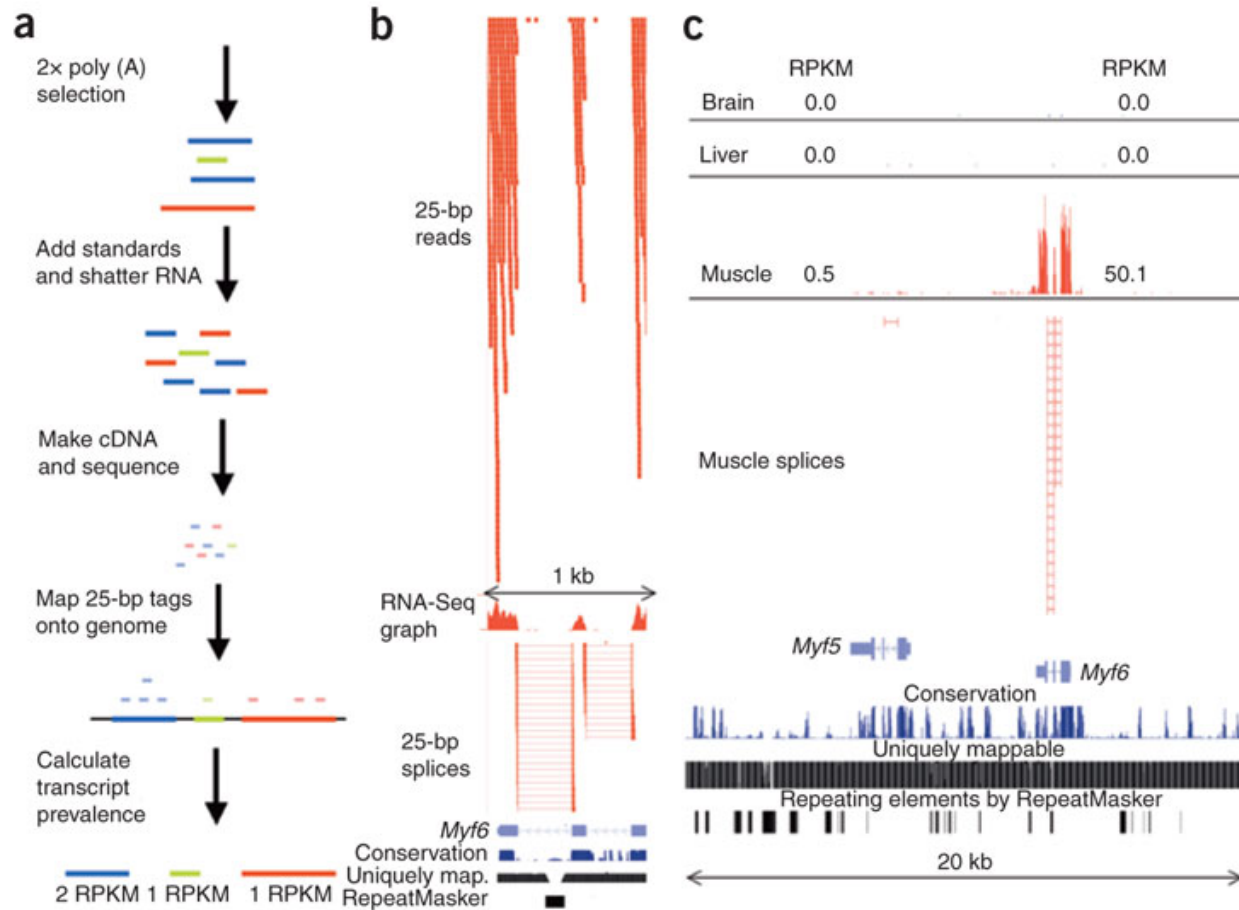- L = Length of gene in KB
- [Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B., Nat Methods. 2008 Jul;5(7):621-8. **Mapping and quantifying mammalian transcriptomes by RNA-Seq.**]
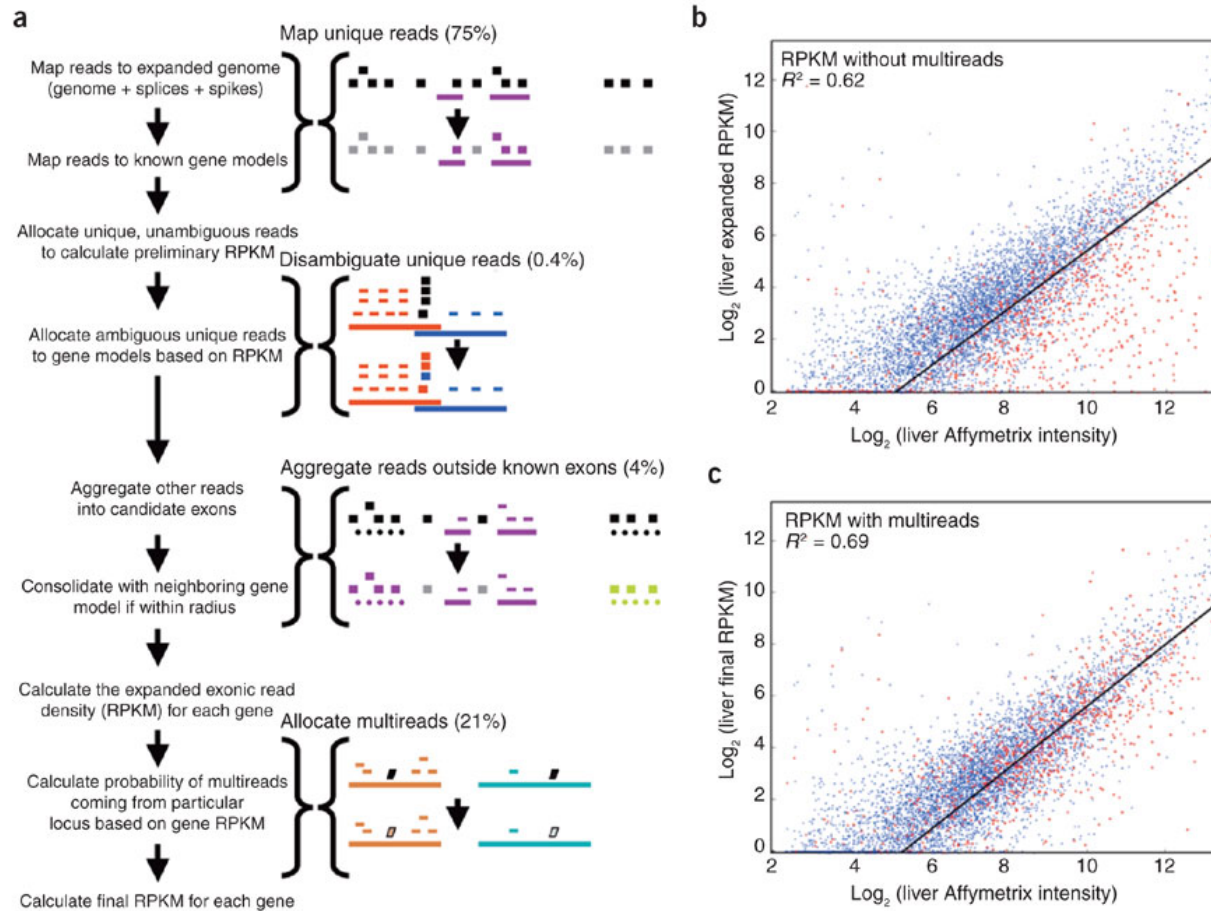
# Complications

❑ Repeat regions
  - Paralogs and other homologous regions in genes
  - Ambiguities in maps

❑ Introns and Exons
  - Aligning reads to genome is more complex

❑ Alternative Splicing

❑ Transcription start site is upstream of ORFs

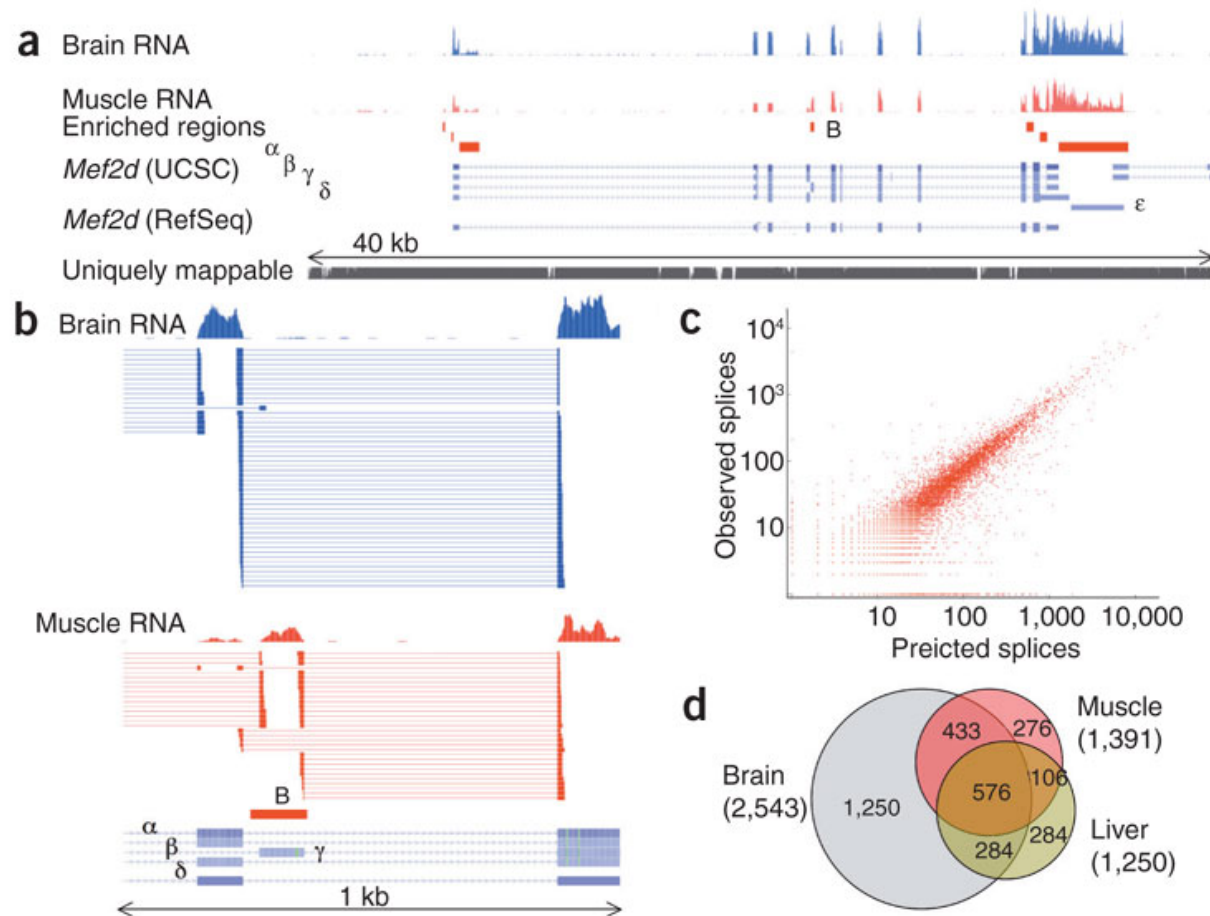❑ Unknown ORFs and Small RNAs
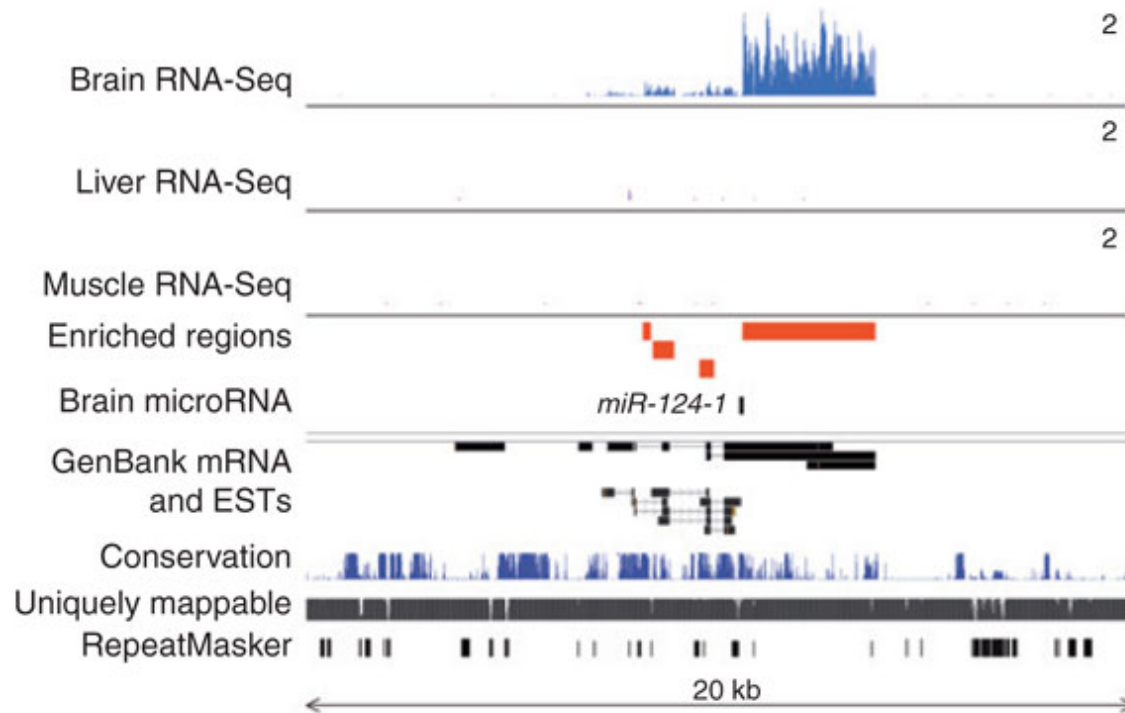
❑ Other transcripts
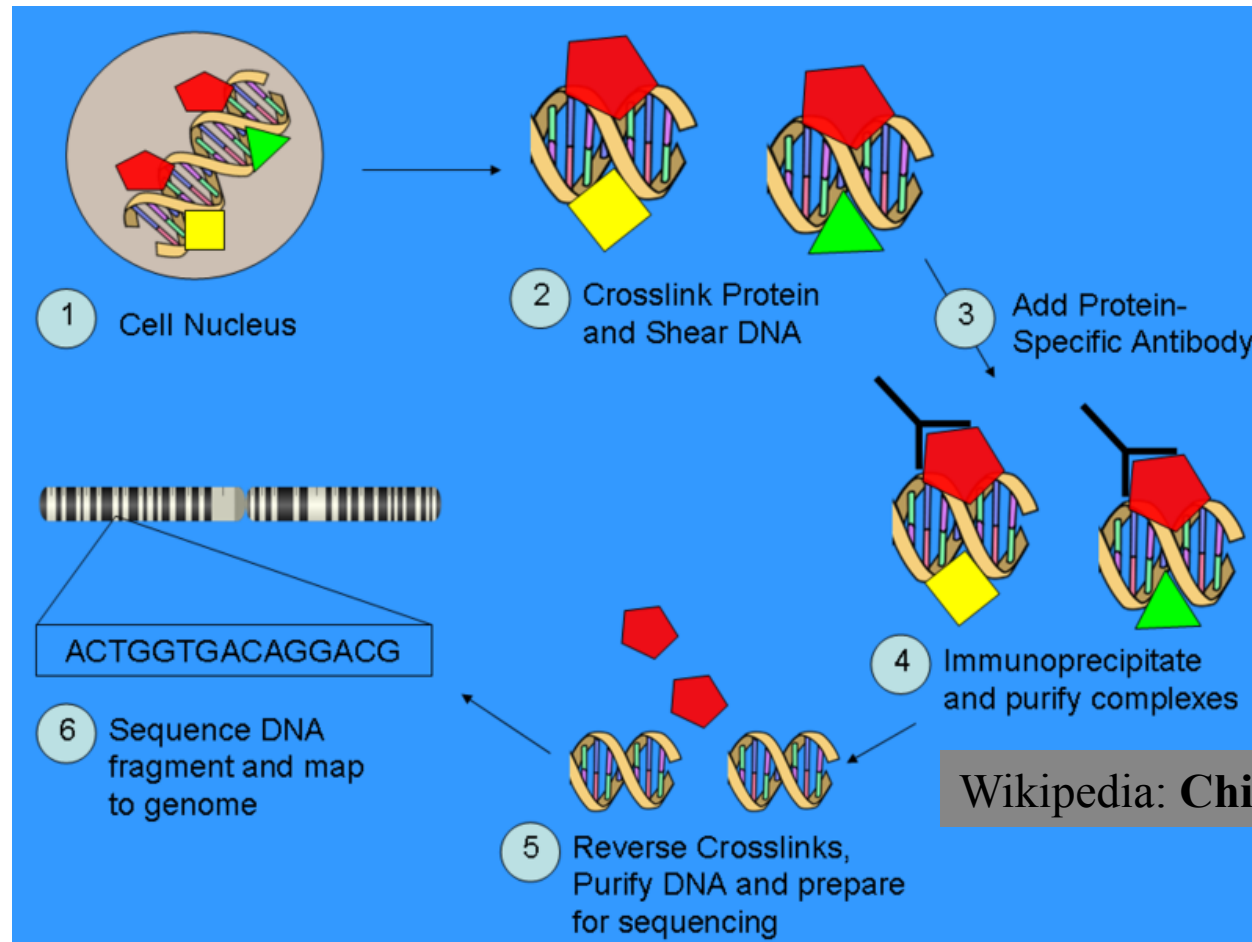
# Mapping Reads to Reference

# Alternative Splicing

# microRNA

# Chromatin Immunoprecipitation

❑ Useful for pinpointing location of TFBS for TF

❑ High-throughput method to find all binding sites for a specific TF under specific conditions

❑ Identify sites using

- ChIP-on-chip (Microarray technique)
- ChIP-Seq (Sequencing technique)

❑ Problems: TFs bind to specific TFBS only under specific conditions – hard to predict

# ChIP-Seq

# SNP-Seq

❑ Align reads and look for differences
- 🔴 Differences to reference
  - ➤ Align reads to reference sequence first
- 🔴 Differences within reads
- 🔴 Differences between samples or sets of reads

# Environmental Microbiology

❑ **Conventional methods**
  - 🔴 Culture, then identify
    - ➤ Slow, expensive, labor intensive, unculturable microbes
  - 🔴 PCR-based length heterogeneity studies

❑ **Microarray-based methods**
  - 🔴 Unique probes for organisms (e.g., Virochip)
    - ➤ Only works for sequenced regions of known organisms

❑ **NGS-based methods**

# Metagenomics

❑ Detect known pathogens
❑ Diversity
  🔴 Identity of individual species not needed
❑ Functional profile of community

# NGS-based method

- ❑ Map reads against appropriate database
- ❑ Identify closest hits for each read
- ❑ Generate contigs
- ❑ Generate abundance information
- ❑ Clustering of reads can be beneficial to estimate abundance