

Computational Genomics

- suggesting biological functions of predicted genes, through homology search
- suggesting possible genes associated with a particular disease, and hence reducing the search space for relevant genes
- suggesting an organism's biology through genome comparison,
- suggesting component-candidate list and their possible interaction relationships in a biological pathway/network
- providing powerful tools for studies of biological evolution
 - sequence/genome comparison
 - phylogenetic profile analysis
- have played key roles in the human and other genome
 - genome assembly
 - protein-coding gene prediction
 - genome annotation



Challenges in Computational Genomics

- One challenge comes directly from the sheer amount of sequence data and the rate at which the data is being generated
 - Thousands of genomes have been (re)sequenced
 - Thousands of genomes are being sequenced
 - prokaryotic genomes / eukaryotic genomes
- The amount of information potentially drivable through comparative genome analysis could be enormous knowing that functional elements are often conserved among "related" genomes
 - how to effectively derive them?!

BIG DATA



Challenges in Computational Genomics

- Prediction of protein-coding genes still represents a challenging problem
 - accurate prediction of exon/intron boundaries
 - prediction of alternatively spliced gene forms
- Protein-coding genes account for ~3% of the human genome. What and where are the other "functional elements (ncRNAs?)" in the rest of the genome?
 - how to identify them?
 - how to (help to) predict their functions?

Challenges in Computational Genomics

- Identification of RNA-coding genes
 - what are the identifiable characteristics of RNA genes?
- Particularly, identification of small regulatory RNA
 - short interference RNAs (siRNA)
 - microRNA (miRNA)
 - Small RNA (smRNA)
- Identification of regulatory elements/binding sites
 - transcription regulatory binding sites
 - splice factor binding sites
 - other classes of regulatory elements?

Challenges in Computational Genomics

- Identification of other types of functional elements
 - transposons
 -
- Identification of genome variations – polymorphisms
 - identification of SNPs
 - prediction of haplotype blocks
- Recognition of genome structures
 - operons, regulons in microbes
 - genomic structures in eukaryotic genomes

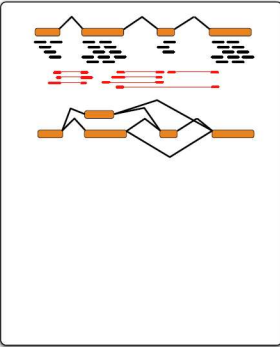
Challenges in Computational Genomics

- Genome is not a linear sequence; It is a 3D structure! **3D Genome**
 - accurate identification and characterization of functional elements by looking at the genome as a 3D DNA structure

.... and many other outstanding challenges!

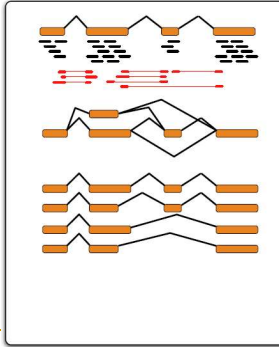
Infer alternative isoforms

- ✓ Map reads
- ✓ Build splicing graph using spliced reads
- ✓ Generate transcripts from graph
- ✓ explain coverage by sparse, weighted sum of transcripts



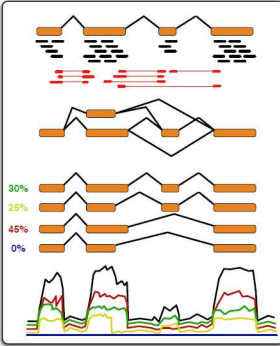
Infer alternative isoforms

- ✓ Map reads
- ✓ Build splicing graph using spliced reads
- ✓ Generate transcripts from graph
- ✓ explain coverage by sparse, weighted sum of transcripts



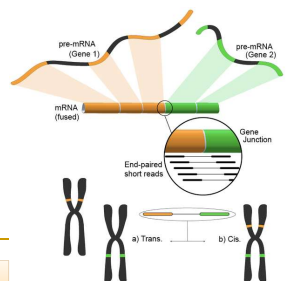
Infer alternative isoforms

- ✓ Map reads
- ✓ Build splicing graph using spliced reads
- ✓ Generate transcripts from graph
- ✓ explain coverage by sparse, weighted sum of transcripts



Gene fusion

- ✓ fusion genes → cancer → “cancer genes”
- ✓ resulting from chromosomal rearrangements in cancer
- ✓ *trans*-splicing




[Wikipedia]

Discovery of ncRNAs

- ✓ Small ncRNAs, Long ncRNAs (lncRNAs)
- ✓ New putative lncRNAs generally identified by RNA-seq → RNA immunoprecipitation followed by sequencing (RIP-seq), parallel analysis of RNA structure (PARS), fragmentation sequencing (Frag-seq) ...

RNA-DNA Differences

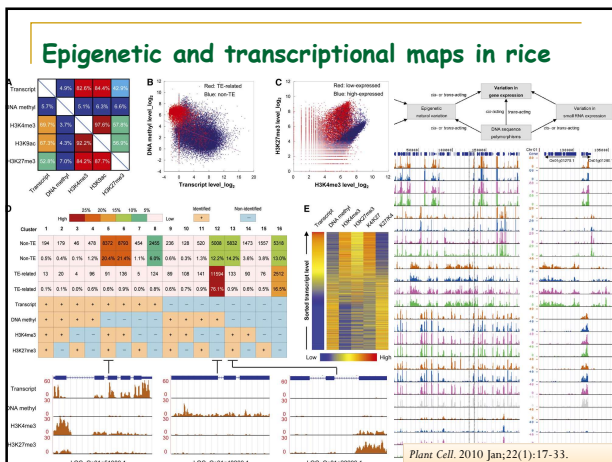
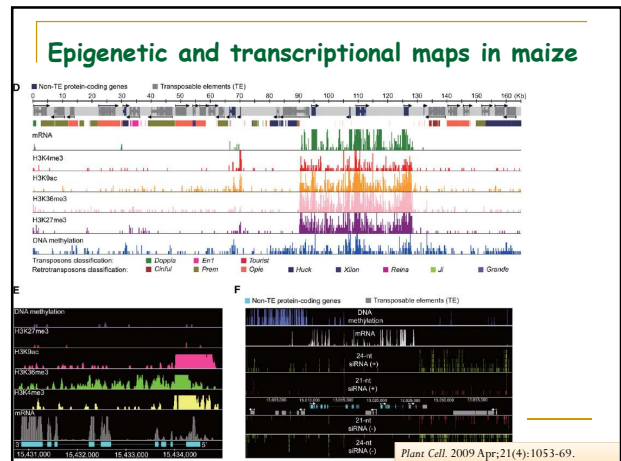
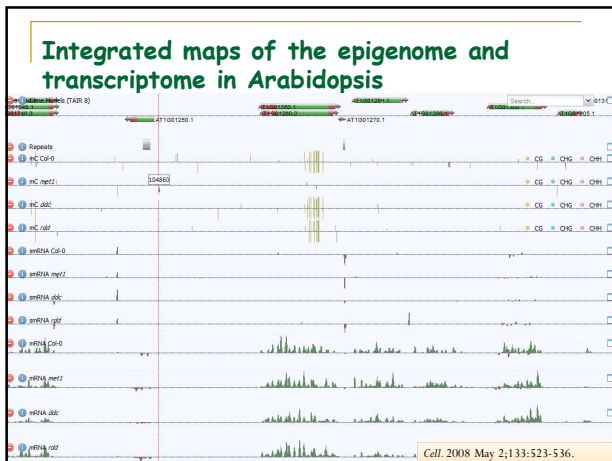
- ✓ SNPs/CNVs
- ✓ RNA editing



Widespread RNA and DNA Sequence Differences in the Human Transcriptome
 Mingyao Li,^{1*} Isabel X. Wang,^{2*} Yuo Li,^{3,7} Alan Bruzel,⁴ Allison L. Richards,⁴ Jonathan M. Toung,⁵ Vivian G. Cheung^{3,4,7}

¹Departments of Biostatistics and Epidemiology, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA. ²Department of Genetics, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA. ³Department of Pediatrics, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA. ⁴Cell and Molecular Biology Graduate Program, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA. ⁵Genomics and Computational Biology Graduate Program, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA. ⁶Departments of Genetics, University of North Carolina School of Medicine, Chapel Hill, NC 27599, USA. ⁷Department of Biostatistics, University of North Carolina School of Medicine, Chapel Hill, NC 27599, USA. ⁸Howard Hughes Medical Institute, Chevy Chase, MD 20815, USA.

Science, 2011 May 19.



Biocomputing

