

De Novo Transcriptome Assembly

Dr. Jung Soh
657.001 Transcriptomics (WS 2018/2019)

What is transcriptomics?

- ▶ Transcriptomics is
 - the study of the complete set of **RNA transcripts** produced by a given cell, organ, or living organism.
- ▶ The main focus of transcriptomics is
 - to evaluate **differential expression** of genes between conditions.
 - High-throughput technologies such as microarray and **RNA-sequencing** have become the standard for such experiments.
- ▶ Transcriptomics software tools are also used
 - to study RNA modifications, RNA-protein interactions, non-coding RNA and RNA structure.

What is RNA-seq?

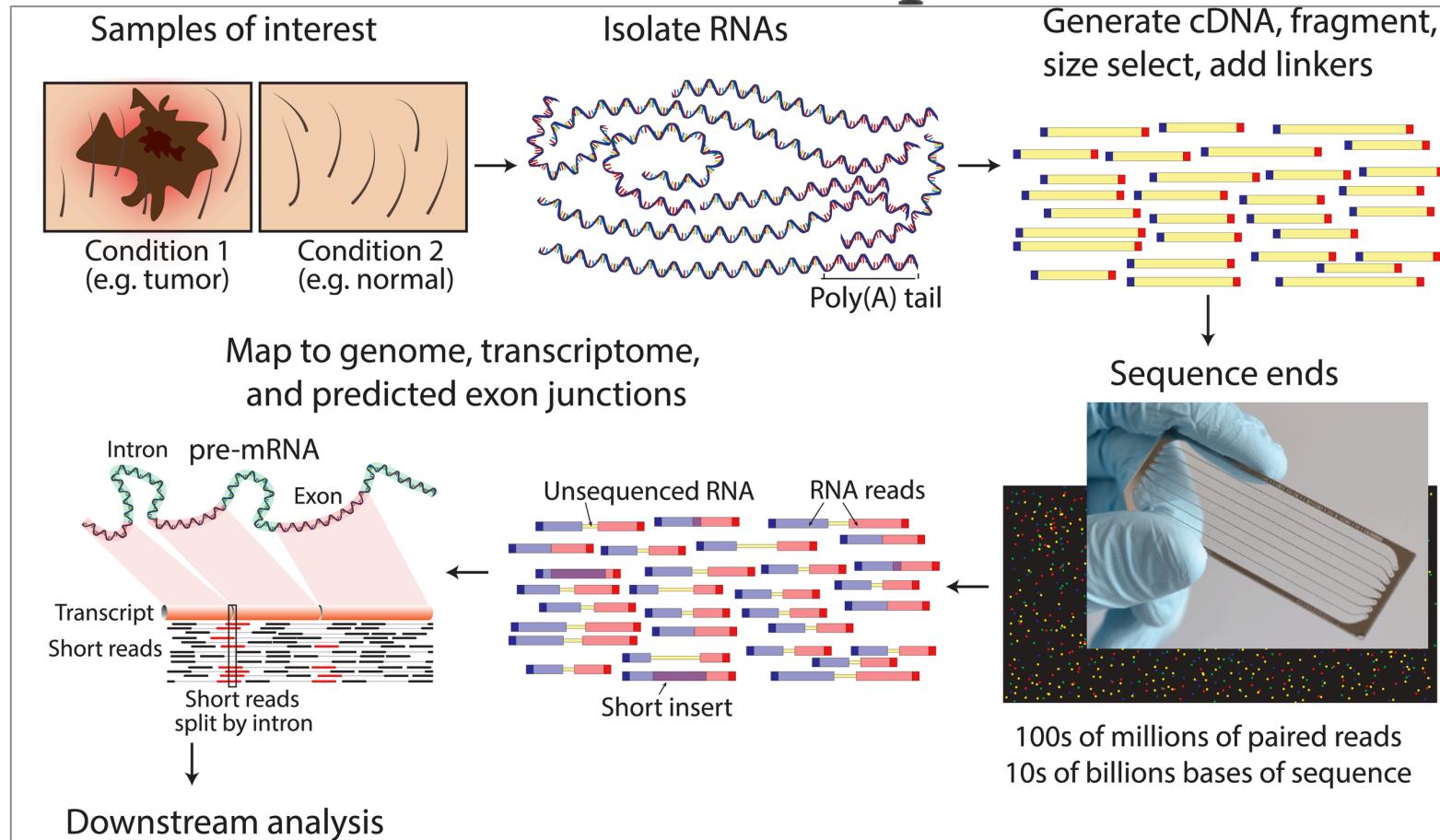
- ▶ An experimental protocol
 - that uses **next-generation sequencing** technologies
 - to sequence the **RNA** molecules within a biological sample
 - in an effort to determine the **primary sequence** and **relative abundance** of each RNA

Martin JA, Wang Z (2011) Next-generation transcriptome assembly. Nat Rev Genet. 12(10):671-682

Why RNA-seq?

- ▶ To study functions based on gene expression changes
 - Drug treatment vs. no treatment
 - Patients vs. healthy people
 - Wild type vs. knock-out
 - Expressions in different tissues
- ▶ Some features available only at RNA level
 - Alternative isoforms, RNA editing, transcript fusion

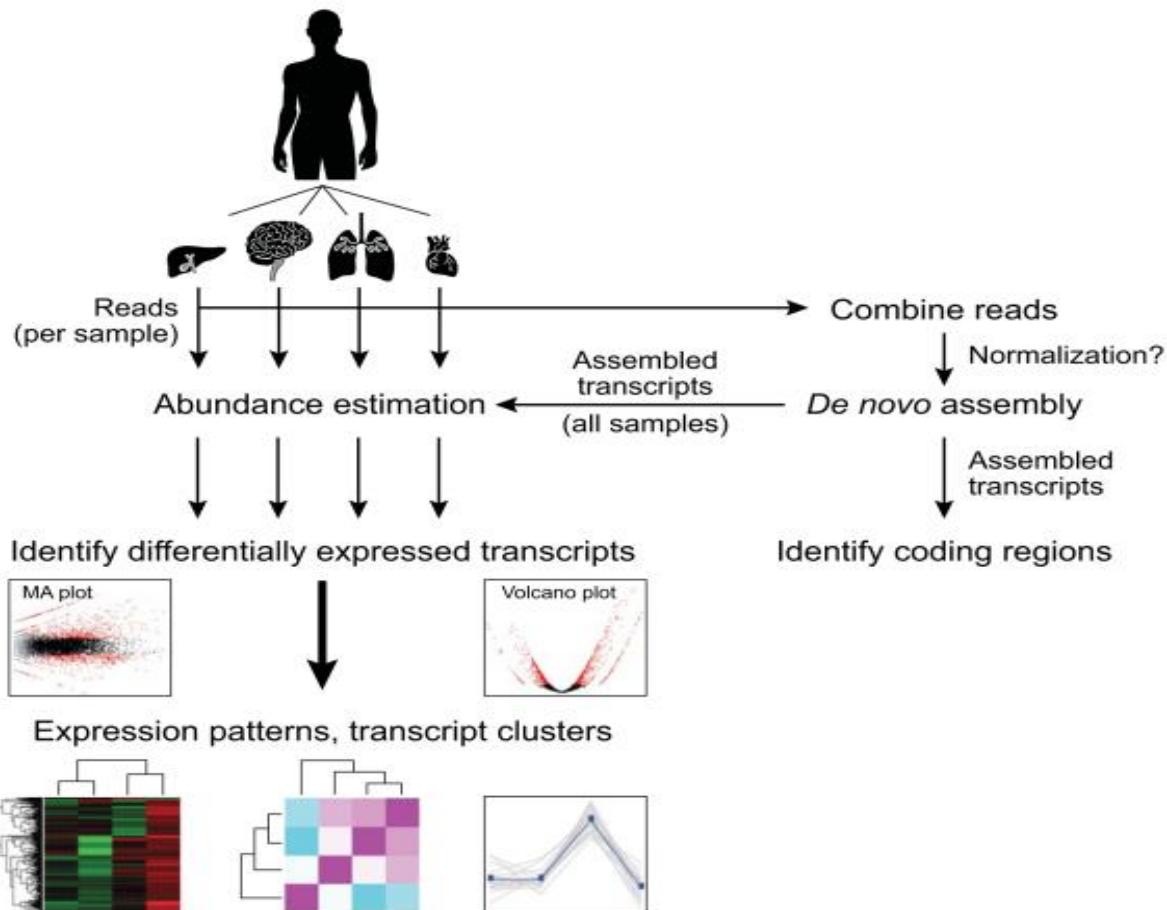
Overview of RNA-seq



RNA-seq challenges

- ▶ Difficulty with sampling
 - More fragile than DNA
 - Different sizes of RNA
- ▶ Relative abundance of RNAs hard to control
 - Can vary by orders of magnitude
 - Uneven coverage
 - Many reads from a small number of highly expressed genes
- ▶ Computational analysis challenges
 - Assembly, (spliced) alignment, quantification, normalization, visualization, ...

De novo transcriptome assembly and analysis



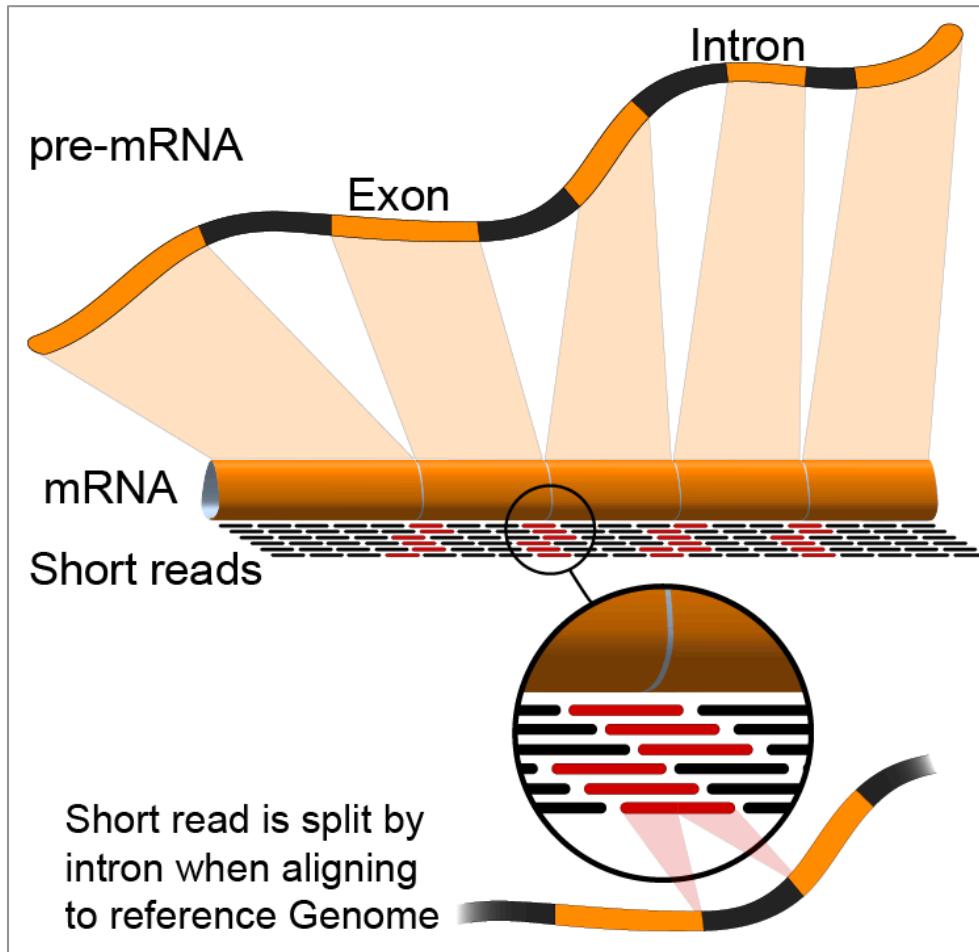
Major transcriptomic analysis tasks

- ▶ *De novo* transcriptome assembly
- ▶ Functional annotation of transcriptome
- ▶ Differential expression analysis
- ▶ Novel transcript discovery

Read alignment

- ▶ Essential method for
 - Transcript abundance estimation
 - De novo assembly quality assessment
- ▶ Can use reference genome or *de novo* assembled transcriptome
- ▶ Issues with aligning to reference genome
 - Reads spanning across exon junction
 - Alternative splicing
- ▶ Most common alignment results format
 - SAM: sequence alignment/map (or BAM: binary version)
- ▶ Many alignment tools
 - Bowtie 2, BWA, SOAP, SHRiMP, mrFAST, mrsFAST, ZOOM, SSAHA2, ...

Aligning RNA-seq reads to genome



SAM/BAM format

► Header section

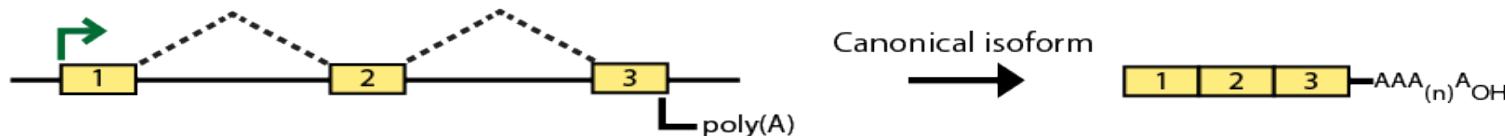
```
mgriffit@linus270 ~> samtools view -H /gscmnt/gc13001/info/model_data/2891632684/build136494552	alignments/136080019.bam | grep -P "SN\\:22|HD|RG|PG"
@HD VN:1.4 SO:coordinate
@SQ SN:22 LN:51304566 UR:ftp://ftp.ncbi.nih.gov/genbank/genomes/Eukaryotes/vertebrates_mammals/Homo_sapiens/GRCh37/special_requests/GRCh37-lite.fa.gz AS:GRCh37-lite M5:a718aca6135fdca8357d5bfe9
4211dd SP:Homo sapiens
@RG ID:2888721359 PL:illumina PU:D1BA4ACXX.3 LB:H_KA-452198-0817007-cDNA-3-lib1 PI:365 DS:paired end DT:2012-10-03T19:00:00-05000 SM:H_KA-452198-0817007 CN:WUGSC
@PG ID:2888721359 VN:2.0.8 CL:tophat --library-type fr-seconstrand --bowtie-version=2.1.0
@PG ID:MarkDuplicates PN:MarkDuplicates PP:2888721359 VN:1.85(exported) CL:net.sf.picard.sam.MarkDuplicates INPUT=[/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-1543 e10-2-5.gsc.wustl.edu-jwalker-1543-136080019/scratch-Ilg6Y/H_KA-452198-0817007-cDNA-3-lib1-2888360300.post_dup.bam METRICS_FILE=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-1543 4-136080019/staging-liuJS/H_KA-452198-0817007-cDNA-3-lib1-2888360300.metrics REMOVE_DUPLICATES=false ASSUME_SORTED=true MAX_FILE_HANDLES_FOR_READ_ENDS_MAP=9500 TMP_DIR=[/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-1543-136080019/scratch-Ilg6Y] VALIDATION_STRINGENCY=SILENT MAX_RECORDS_IN_RAM=500000 PROGRAM_RECORD_ID=MarkDuplicates PROGRAM_GROUP_NAME=MarkDuplicates MAX_SEQUENCES_FOR_DISK_READ_ENDS_MAP=50000 SORTING_COLLECTION_SIZE_RATIO=0.25 READ_NAME_REGEX=[a-zA-Z0-9]+:[0-9]+:[0-9]+:[0-9]+.* OPTICAL_DUPLICATE_PIXEL_DISTANCE=100 VERBOSITY=INFO QUIET=false COMPRESSION_LEVEL=5 CREATE_INDEX=false CREATE_M5_FILE=false
mgriffit@linus270 ~>
```

► Alignment section (10 alignments shown)

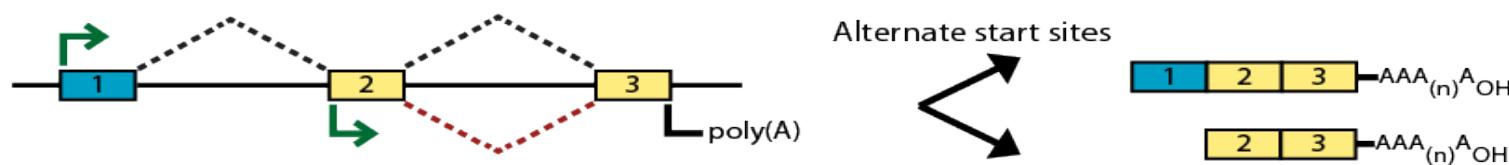
```
mgriffit@linus270 ~> samtools view -f 3 -F 1804 /gscmnt/gc13001/info/model_data/2891632684/build136494552	alignments/136080019.bam | head
HWI-ST495_129147882:3:2114:15769:38646 99 1 11306 3 100M = 11508 302 ACTGCGGGGCCCTCTGGCTACTGTATAGGGGACATTGCAGGGCTCTTGCTCAAGGTGAGTGGCACACGC
CCCCFFHHHHJJJJJJJJJJHJJHJJJJHIIJJJJJHFD#####DDDDDDDDDDDDDDCCCC->@CDDDDDDDD?1+B
1 XN:i:0 X0:i:0 CP:i:102519765 AS:i:-5 XS:A:- YT:Z:UU
HWI-ST495_129147882:3:2114:15769:38646 147 1 11508 3 100M = 11306 -302 ACTCTCAAATGGGATTCTGGTTAAAGTATAAAAATGTTAATGGGAACTGATTACCATCAGAACATTGTACTGTATCCCACAG5
;5:CDCDCDECFD@H=9E?7EEIIIIHCEGGIJJJJJJJHIF@?00IHHFFGG?*JJJJGHGEJJJJJJJHJCJJHFFHGHFFEDFCB
1 XN:i:0 X0:i:0 CP:i:102519563 AS:i:-6 XS:A:+ YT:Z:UU
HWI-ST495_129147882:3:1210:1257:16203 163 1 11810 3 100M = 12055 345 CCTGCGTAGTTAACAGGAGTTGCCAGCACGGGATCATTCAACCATTCTTCTGTTACTTGCAGCCTTTTGTGACCTCTTCTTC
CCCCFFHHFHAGGGIIJJJJJEHGIGGGJJJJG@?EHIGIJGDGHIIHGIGGJJJJJJHGGHHFFCDDDDCDCCCCA;>@@@@:AA>AA
0 XN:i:0 X0:i:0 CP:i:102519261 AS:i:0 XS:A:- YT:Z:UU
HWI-ST495_129147882:3:1210:1257:16203 83 1 12055 3 100M = 11810 -345 GAGCACTGGAGTGGAGTTCTCTGGAGAGGAGCCATGCCAGTGGGATGGGCCATTGTTCATCTCTGGCCCTGTGCTCATGTAACCTAAC
CC>4>CCACACDCCB2BDCCE@ECCFFHHHHJJJJHIIJJJJHHHEIIHGJJJJHIIJJJJJJJJJJJJJJJJHGHHHDFEFFCCCC
0 XN:i:0 X0:i:0 CP:i:102519016 AS:i:0 XS:A:+ YT:Z:UU
HWI-ST495_129147882:3:2111:3117:78828 163 1 12634 3 100M = 12746 212 GCCCTCCCCAGCATCAGGTCTCCAGAGCTCGAGAACGAGGCCGACTTGGATCACACTTGTGAGTGTGCCAGTGTGACAGGTGAGGAGA<
@FFFFFHRRHHFHGIIFGA#DHEGII=GHIJIIIIIIIIIIIIIIIIIIIIHDDFFEEEEECCACCCCCC:AADCCBCC=CAC<CCCCC:@C@#B@B#>
1 XN:i:0 X0:i:0 CP:i:102518437 AS:i:-5 XS:A:- YT:Z:UU
HWI-ST495_129147882:3:2111:3117:78828 83 1 12746 3 100M = 12634 -212 GGGACTGGCGTCGCCCTTACGGGCCATCTCTGTCATGGGAGGCGCTTCATGCCCTCAGCATCCCCTGATCTCCCTGTGATGTD
DCABDBDDDDDDDDDDDBDDDBB@BDBD@Q;CCCCCDFE@Q;?<HIGGEIHEGJJJJIGGIIHEGFEHFJIIIIIGJJJJHHHHFFFFFC@Q
1 XN:i:0 X0:i:0 CP:i:102518325 AS:i:-5 XS:A:- YT:Z:UU
HWI-ST495_129147882:3:1102:4242:26638 99 1 13503 3 100M = 13779 376 CGCTTGCCCTTCCTTGCTCTGCCCTGGAGACCGGTGTTGTCTAGGGCTGCTGAGGGATCTGCTACAAAGGTGAAACCCAGGAGTGTGAC
CCCCFFHHHHJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJHGGHFFFDEEEEEECCACDDACDCCDDDB?>B@A@CDC
0 XN:i:0 X0:i:0 CP:i:114357414 AS:i:0 XS:A:+ YT:Z:UU
HWI-ST495_129147882:3:1309:15328:74082 99 1 13534 3 100M = 13780 346 AGACGGTGTGTCATGGGCCCTGGCTCGAGGGATCTGCTACAAGGTGAAACCCAGGAGTGTGAGTCCAGAGTGTGAGTCCAGGACCCAGG@
CCCCFDHHHHJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJHBFHIIJJJJJJJJH=EEEEECCEDCDCD#####DDDBCCD
0 XN:i:0 X0:i:0 CP:i:114357383 AS:i:0 XS:A:- YT:Z:UU
HWI-ST495_129147882:3:1308:10126:19636 99 1 13779 3 100M = 14027 348 CCTCTGAGGGACTGCCATTGGCTGCCACCTCTAGAACGCGAGCGAGGCCACTCTGCTACTGCTCTTCTATAATAAACTAAAGTAGTGC
CCCCFFHHHHJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJHJJHGHFFFDEEEEEECCDCDCDFACCAACCDCCCCCD
0 XN:i:0 X0:i:0 CP:i:114357140 AS:i:0 XS:A:+ YT:Z:UU
HWI-ST495_129147882:3:1102:4242:26638 147 1 13779 3 100M = 13503 -376 CCTCTGAGGGACTGCCATTGGCTGCCACCTCTAGAACGCGAGGCCACTCTGCTACTGCTCTTCTATAATAAACTAAAGTAGTGC#
##DCDDCCBBAACCDDBBDH@C=GIJIIIIIIJJHJJJJJJJJJJJJJJJJHGHGJJIIJJJJJJJJHGGHHHHHHHHHCCCC
0 XN:i:0 X0:i:0 CP:i:114357140 AS:i:0 XS:A:+ YT:Z:UU
mgriffit@linus270 ~>
```

Types of alternative expression

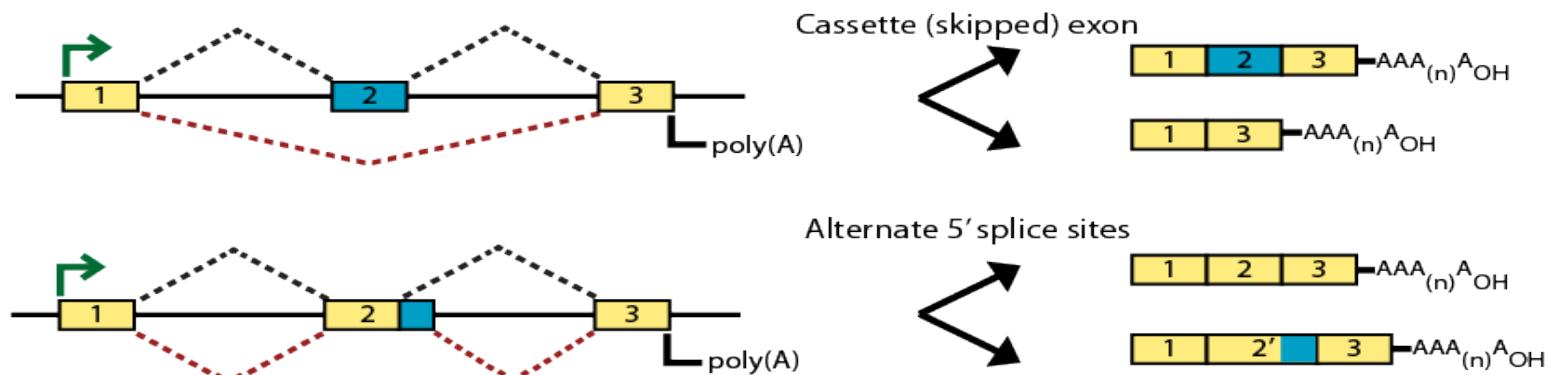
Simple transcription



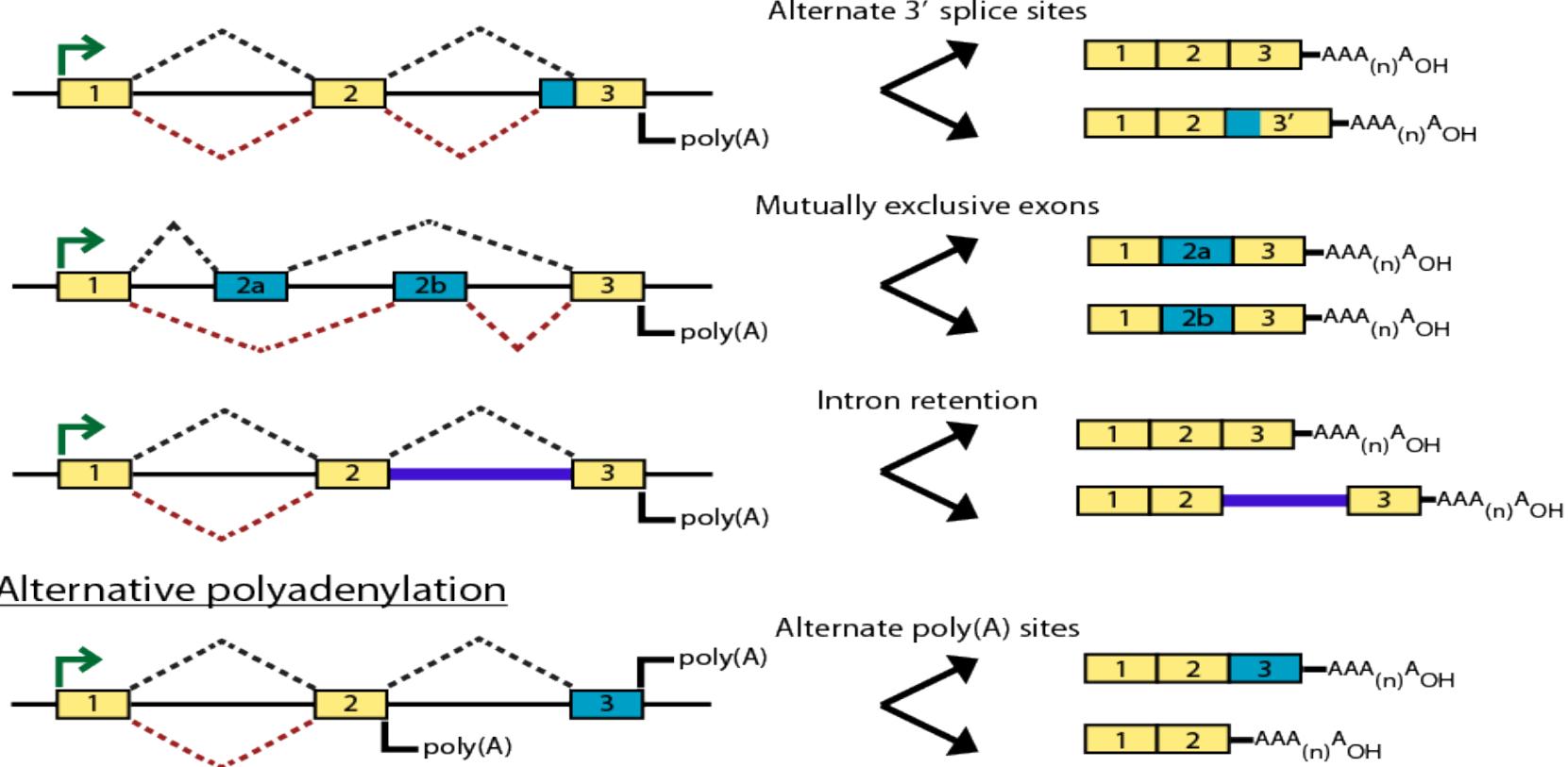
Alternative transcript initiation



Alternative splicing



Types of alternative expression



De novo transcriptome assembly

- ▶ No reference genome, use only RNA-seq reads
- ▶ Required for most non-model organisms
- ▶ Similar to *de novo* genome assembly
 - RNA reads tend to be shorter
 - Uneven coverage depth
- ▶ Many assemblers
 - Trinity, rnaSPAdes, SOAPdenovo-trans, Velvet/Oases, ...

Genome assembly metaphor

DNA clones



Reads



Reconstructed genome

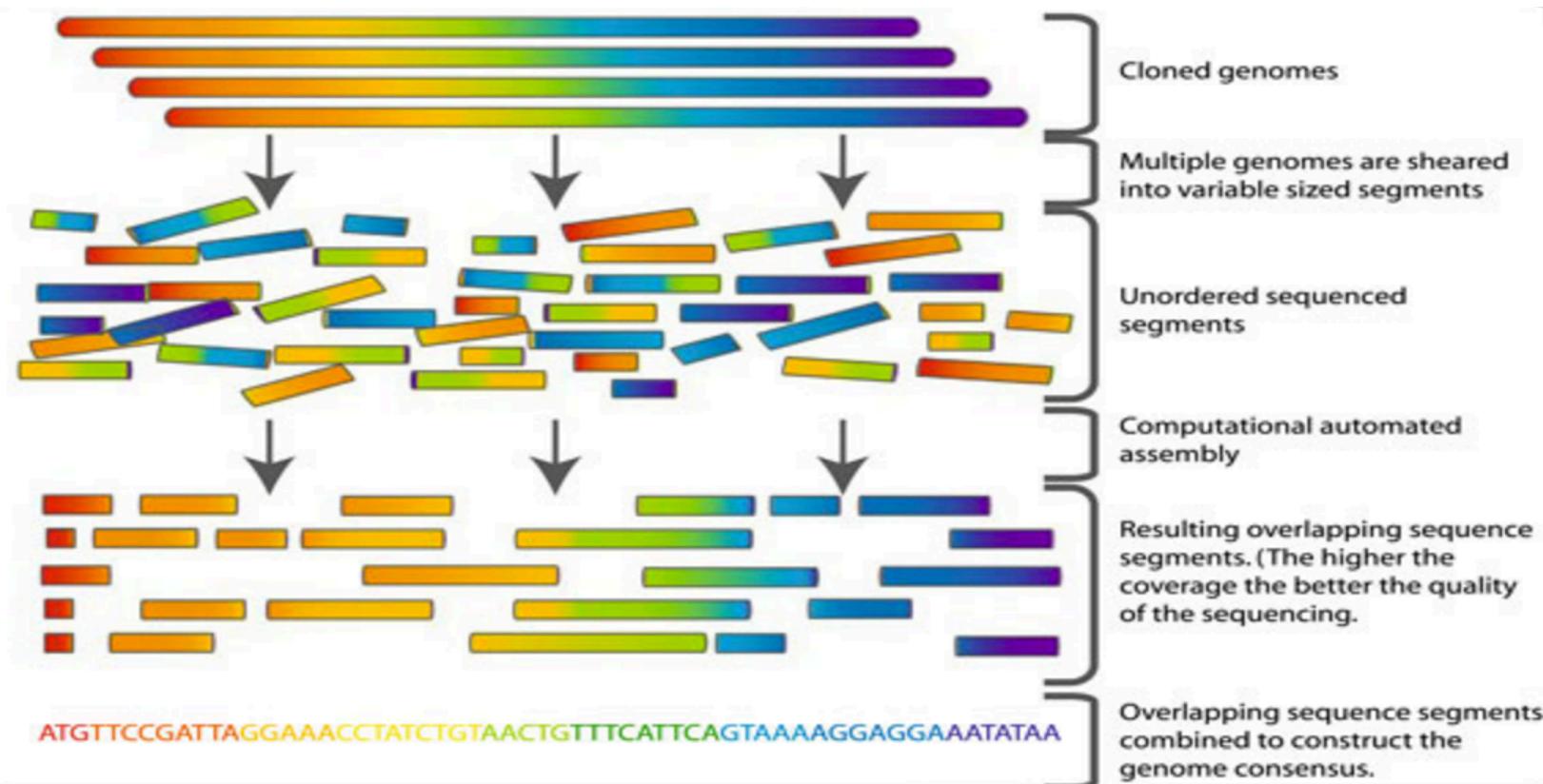
The Call-Chronicle-Examiner
SAN FRANCISCO, THURSDAY, APRIL 10, 1906

EARTHQUAKE AND FIRE: SAN FRANCISCO IN RUINS



<http://www.vicbioinformatics.com/documents/Genome%20Assembly%20Strategies%20-%20Torsten%20Seemann%20-%20IMB%20-%205%20Jul%202010.pdf>

Genome assembly: another view



Commins J, Toft C, Fares MA. (2009) Computational biology methods and their application to the comparative genomics of endocellular symbiotic bacteria of insects. Biol Proced Online. 11:52-78.

Lab overview

