

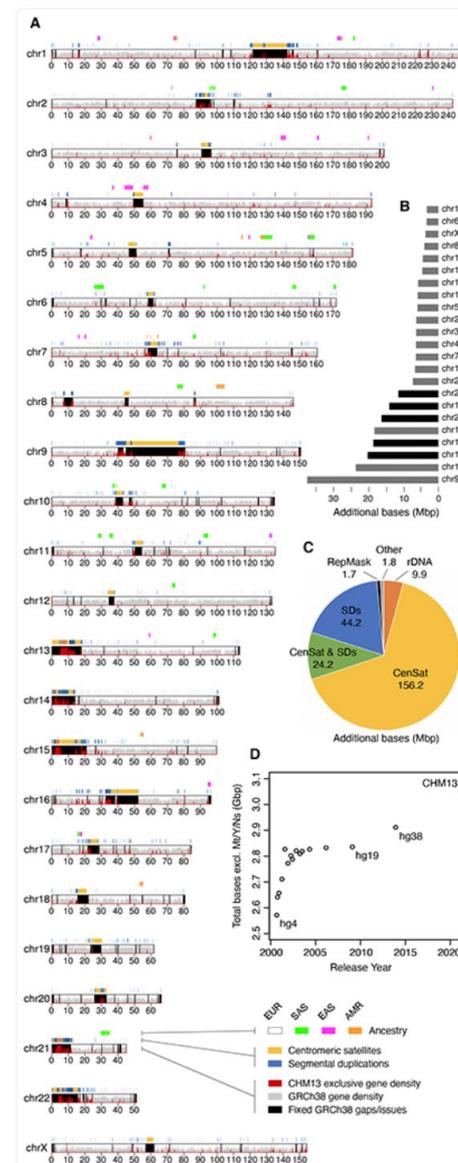
Syllabus and grading

#	Date	Short title	Lecturer	Subject
1	10/102024	introduction	MR	Overview of Bioinformatics, sequence alignment
2	17/102024	Linux/shell/ssh	AD	Introduction to Linux and the command line, bash scripting and ssh
3	24/102024	R (1)	AD	Introduction to the R programming language and Rstudio usage
4	31/102024	R (2)	AD	Advances R subjects, introduction to Bioconductor
5	07/112024	QC+RNASeq	MR	Next generation sequencing: introduction, quality control and gene expression analysis for RNAseq
6	14/112024	bedtools/vcftools/samtools	AD	Command line tool usage: bedtools, vcftools, samtools etc.
7	21/112024	Denovo	MR	NGS for denovo genome and transcriptome assembly
8	28/112024	Exome/SNP calling	AD	Pipelines for SNP calling, especially for exome sequencing using the GATK pipeline
9	05/122024	ChipSeq/chirp	MR	NGS analysis for molecular interactions (ChipSeq, (Par-)Clip, structural sequencing, chromosome conformation capture (3C))
10	12/122024	presentations	MR+AD	Pipelines for SNP calling, especially for exome sequencing using the GATK pipeline
11	19/122024	presentations	MR+AD	Paper presentations by students
12	09/012025	metabolomics	MR	Genome-scale models of metabolism and macromolecular expression, Biological applications of Transformers
13	16/012025	final projects support	MR+AD	Support for the final project

Grade	100%
Presentation	30%
Exercises	20%
Final Project	50%

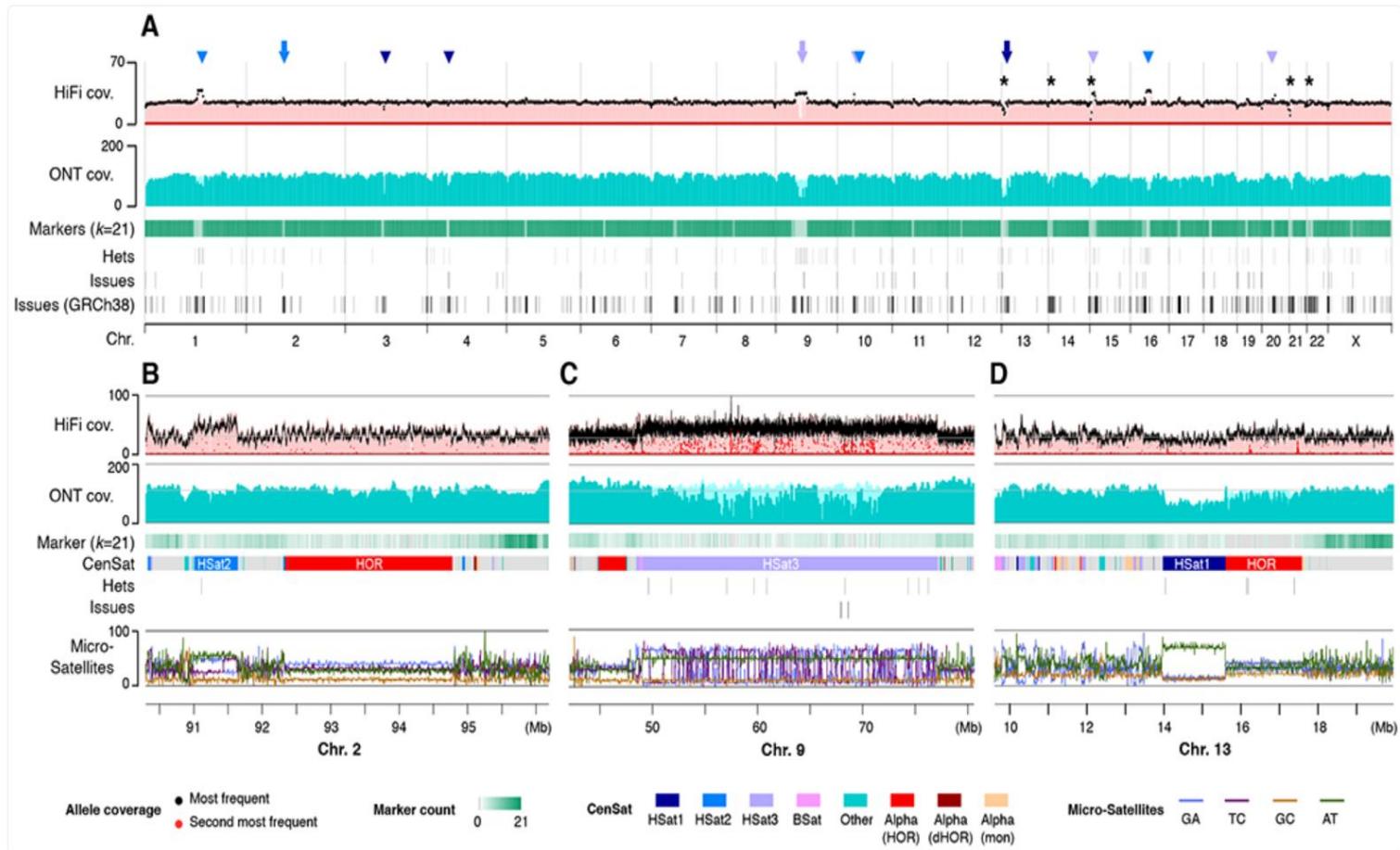
Science. 2022 April ; 376(6588): 44–53. doi:10.1126/science.abj6987.

The complete sequence of a human genome



(A) Ideogram of T2T-CHM13v1.1 assembly features. Bottom to top: gaps/issues in GRCh38 fixed by CHM13 overlaid with the density of g CHM13 in red; segmental duplications (SDs) (42) and centromeric satellites (CenSat) (30); and CHM13 ancestry predictions (EUR, European; Asian; EAS, East Asian; AMR, Ad Mixed American). (B) Additional (non-syntenic) bases in the CHM13 assembly relative to GRCh38 per c the acrocentrics highlighted in black, and (C) by sequence type (note that the CenSat and SD annotations overlap). (D) Total non-gap bases i genome releases dating back to September 2000 (hg4) and ending with T2T-CHM13 in 2021.

Fig. 3. Sequencing coverage and assembly validation.



(A) Uniform whole-genome coverage of mapped HiFi and ONT reads is shown with primary alignments in light shades and marker-assisted alignments overlaid in dark shades. Large HSat arrays (30) are noted by triangles, with inset regions are marked by arrowheads and the location of the rDNA arrays marked with asterisks. Regions with low unique marker frequency (light green) correspond to drops in unique marker density, but are recovered by the lower-confidence primary alignments. Annotated assembly issues are compared for T2T-CHM13 and GRCh38. (B–D) Enlargements corresponding to regions of the genome featured in Fig. 2. Uniform coverage changes within certain satellites are reproducible and likely caused by sequencing bias. Identified heterozygous variants and assembly issues are marked below and typically correspond with low coverage of the primary allele (black) and elevated coverage of the secondary allele (red). % microsatellite repeats for every 128 bp window is shown at the bottom.

Single Cell Sequencing

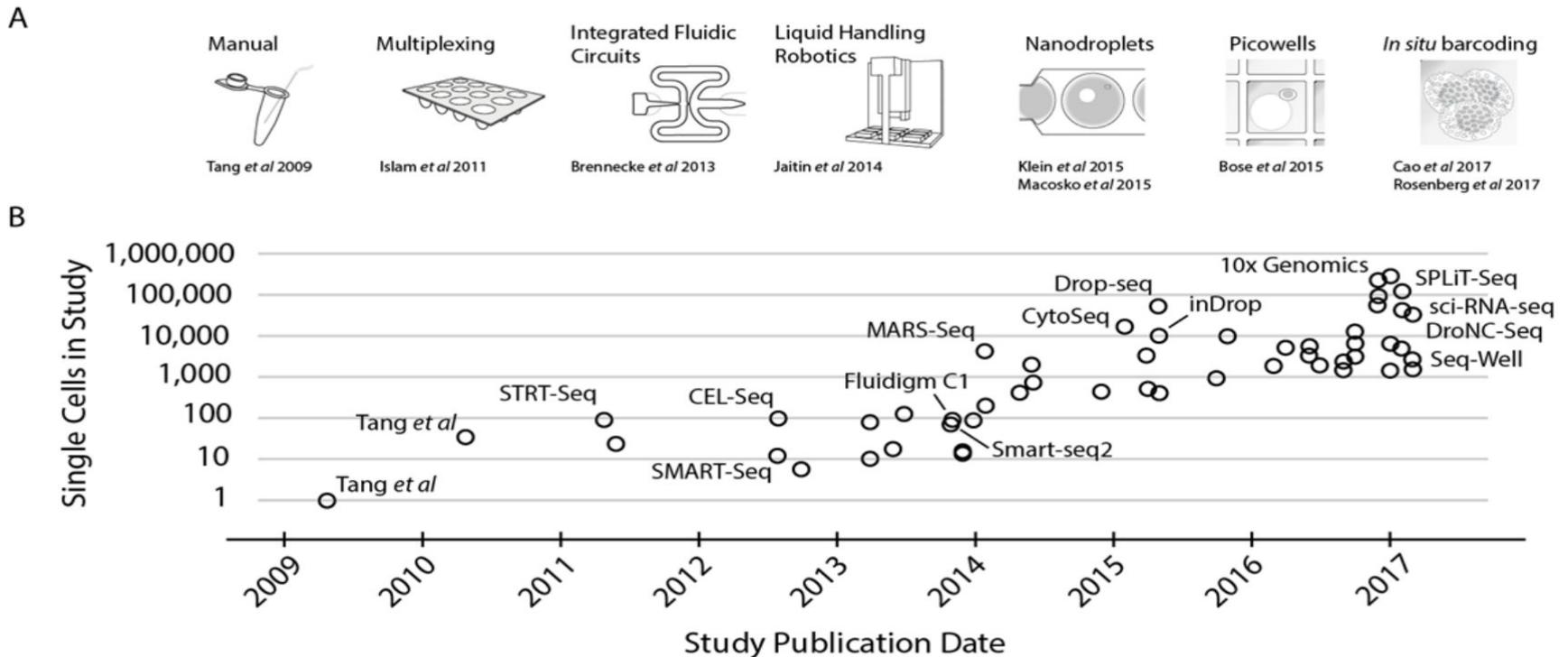


Figure 1: Scaling of scRNA-seq experiments (A) Key technologies allowing jumps in experimental scale. A jump to ~100 cells was enabled by sample multiplexing, a jump to ~1,000 cells by large scale studies using integrated fluidic circuits (IFCs), followed by a jump to several thousands using liquid handling robotics. Further order of magnitude jumps were enabled by random capture technologies through nanodroplets and picowell technologies. Recent studies have employed *in situ* barcoding to reach the next order of magnitude. **(B)** Cell numbers reported in representative publications by publication date. Key technologies and protocols are marked, and a full table with corresponding numbers is available in **Supplementary Table 1**.

Single Cell Sequencing

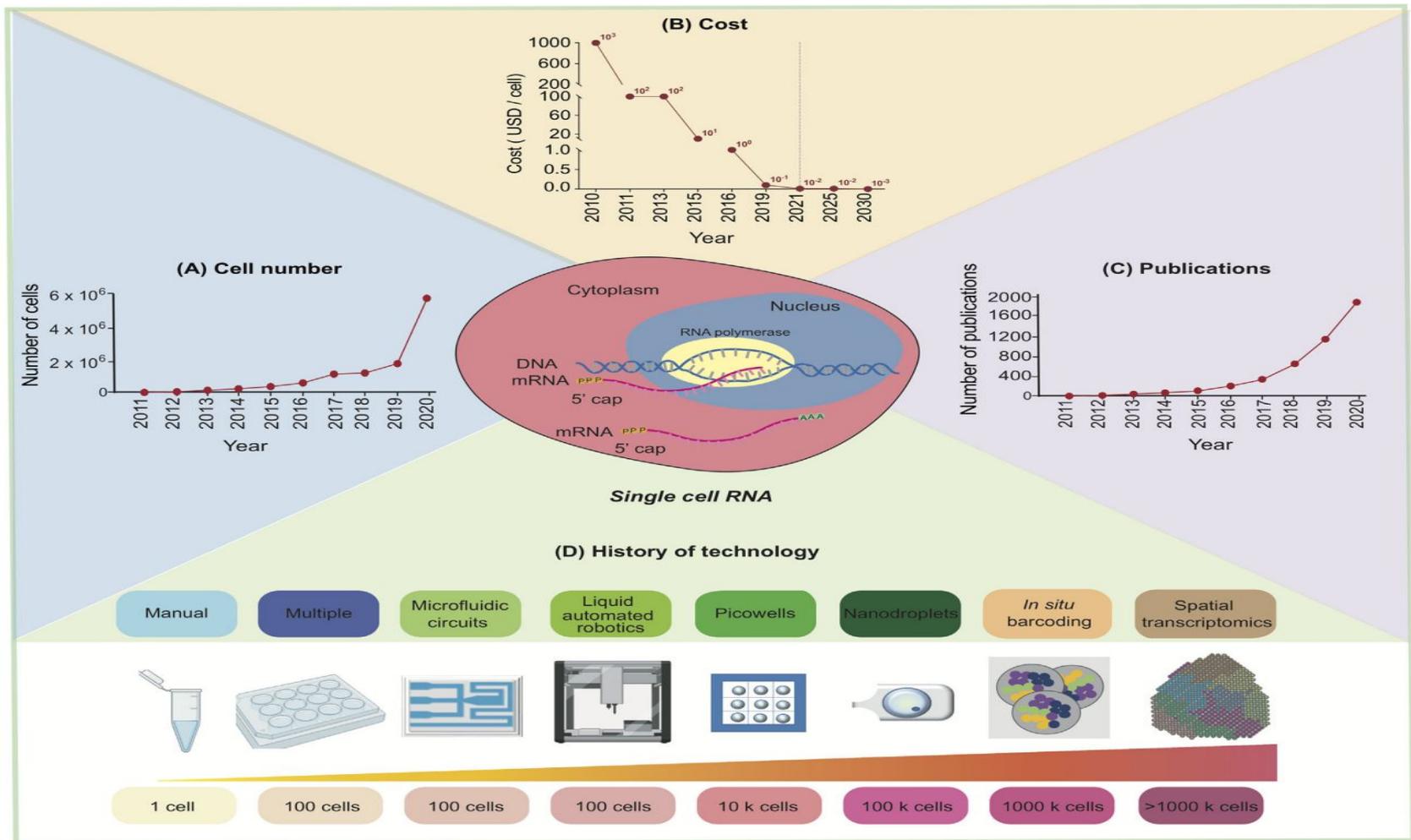
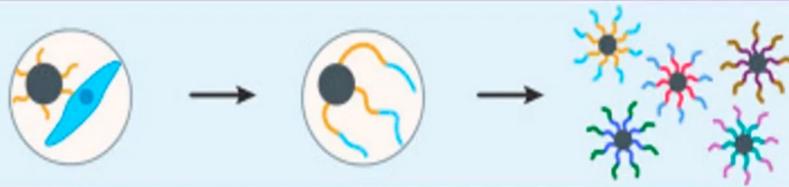
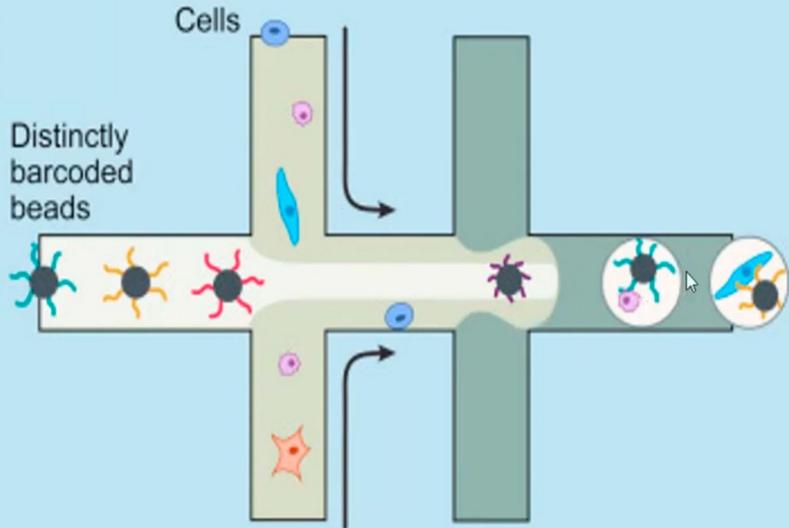


FIGURE 1 Development of single-cell RNA sequencing technology. With the technological advances in single-cell RNA sequencing (scRNA)-seq, (A) the number of analyzed cells increased, (B) the cost (in US dollar) was exponentially reduced, (C) the number of published papers increased and (D) the history of technology evolution in the last decade using more sophisticated, accurate, high throughput analysis was achieved. Part (D) is created with icons from BioRender with license for publication

Bead: Cell barcode and unique molecular identifiers (UMIs)

Drop-seq single cell analysis



1000s of DNA-barcoded single-cell transcriptomes

- Cell barcode: which cell the read comes from
- UMI: which mRNA molecule the read comes from (helps to detect PCR duplicates)

B Barcoded primer bead

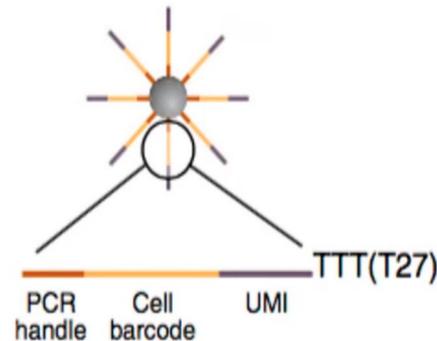


Figure by Macosko et al, *Cell*, 161:1202-1214, 2015

From reads to digital gene expression matrix (DGE)



Overview of DGE extraction

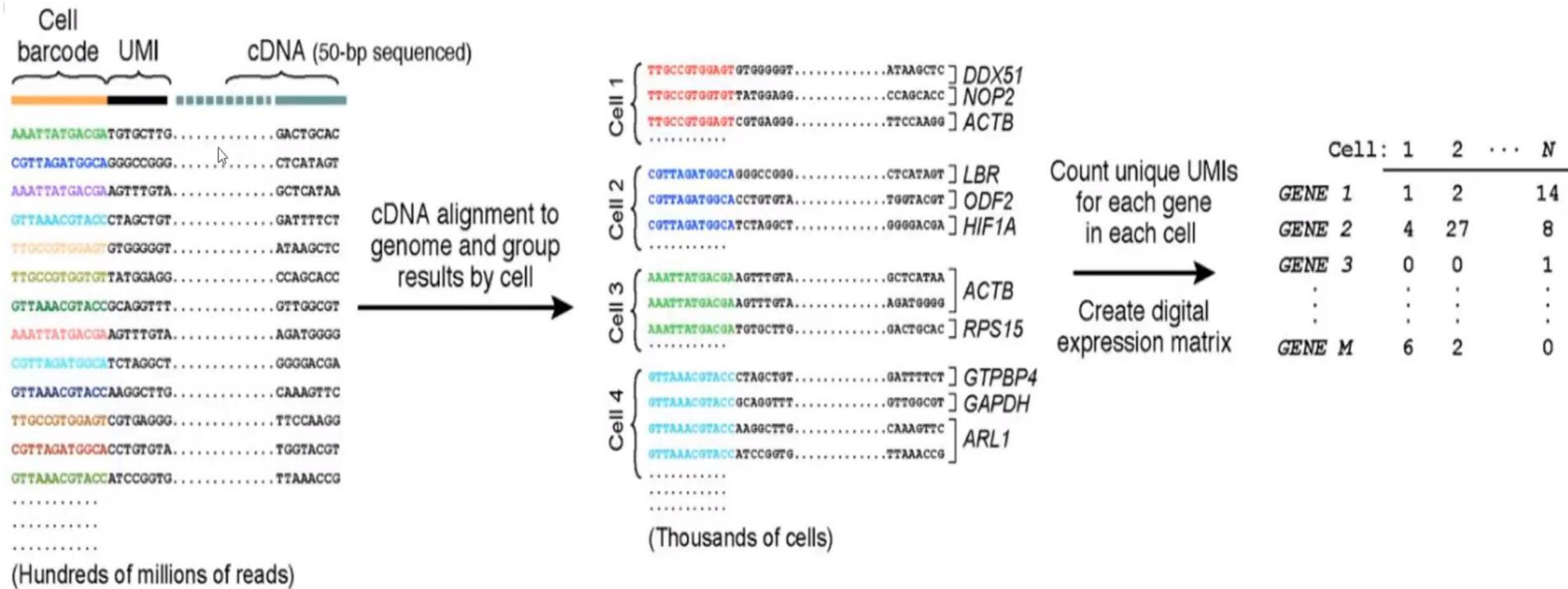
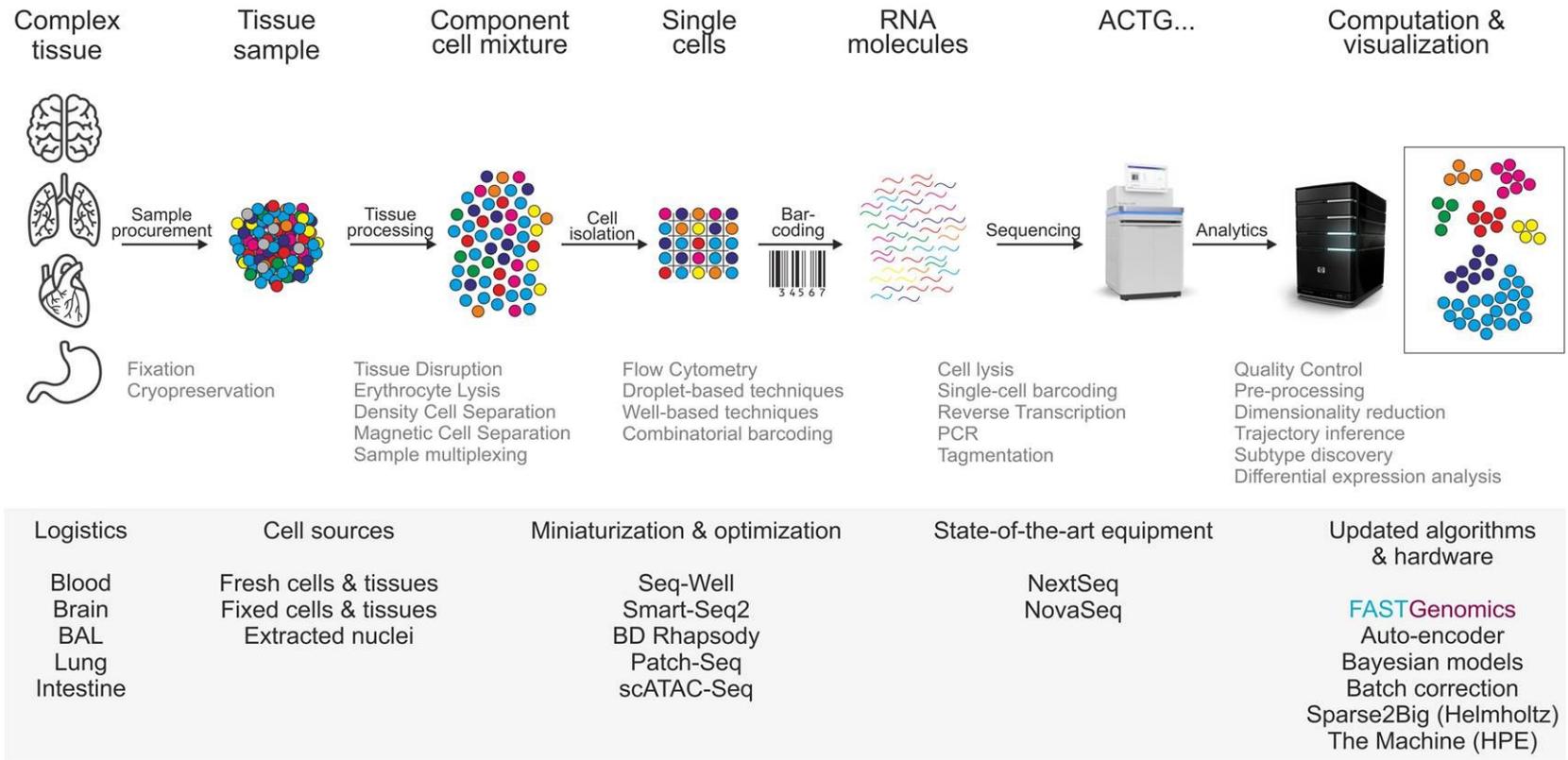


Figure by Macosko et al, Cell, 161:1202-1214, 2015

Single Cell RNA Sequencing and its main applications

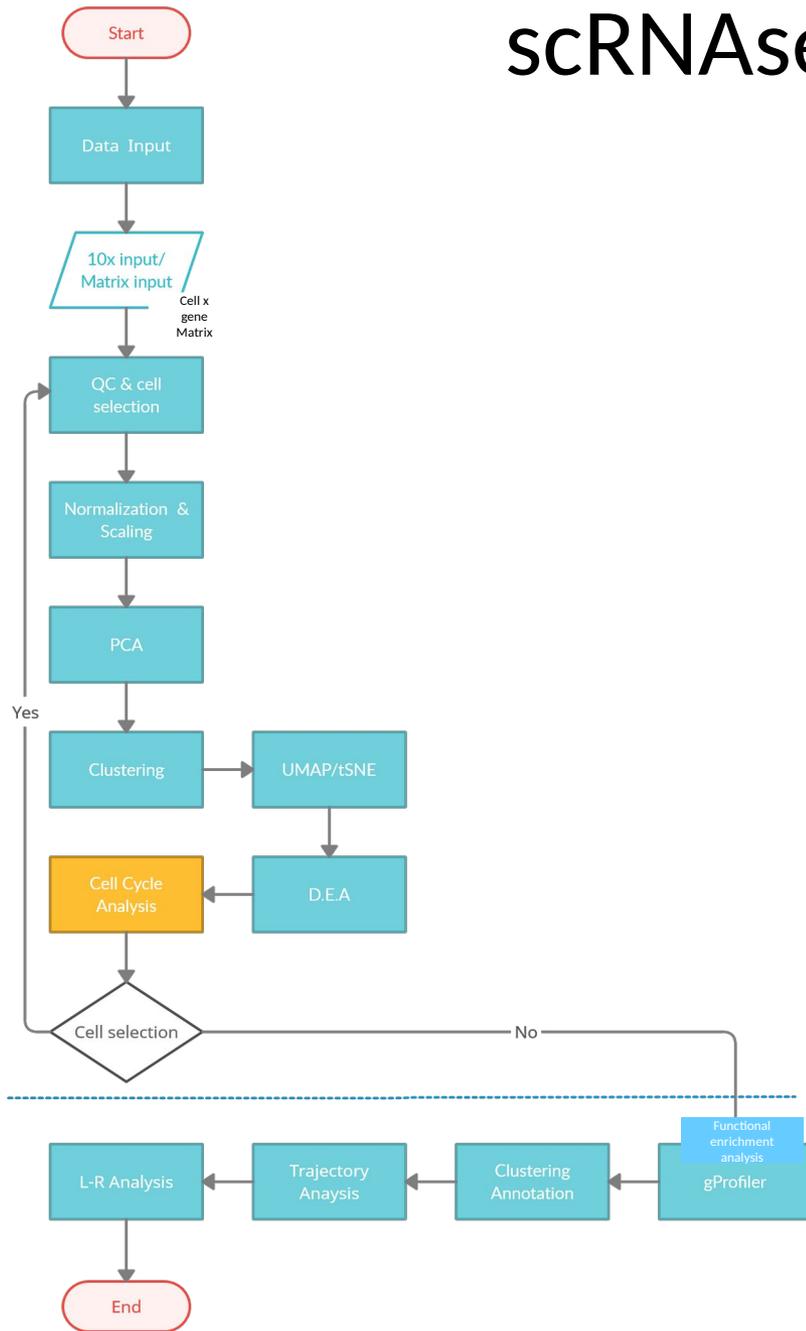
Single Cell RNA Sequencing and its main applications

- ❑ Identification of new cell populations and subpopulations in complex tissues
- ❑ Studying gene dynamics in developmental studies
- ❑ Immune cell profiling
- ❑ Cancer research
- ❑ Personalized medicine
- ❑ Cell atlases



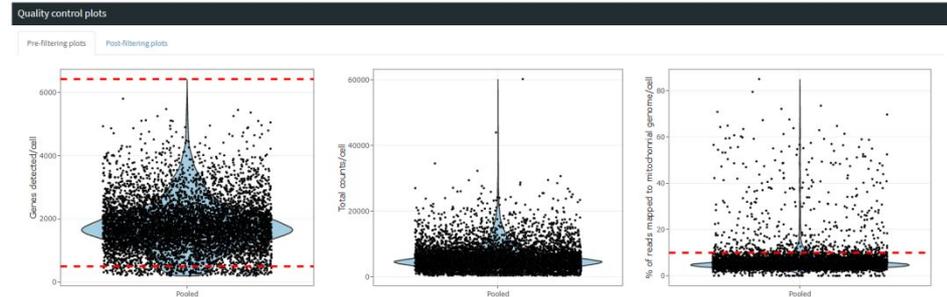
(slides by ITBI student Dimitra Panou)

scRNAseq pipeline



1. Quality Control & Cell selection

- ☐ Detect + remove low quality cells from downstream analysis
 - Genes detected/cell
 - Total reads/cell
 - % of reads in mitochondrial genome/cell



Red lines show the filtering points

scRNAseq pipeline

2. Normalization & Scaling

- **Global normalization**
 - Correcting for sequencing depth differences between cells
 - Log transformation
- **Detection of Highly variable genes**
 - Mean.var.plot method (**mvp**) highly variable genes
 - Scaling transformation
- **Calculation of scaled values for all genes**
 - Scales + centers the genes in dataset

3. Dimensional reduction

PCA analysis

Detection of most informative principal components

- Moving to PCA space can help reducing runtime of cell clustering
- May fail to capture local patterns in scRNA data

Non linear methods

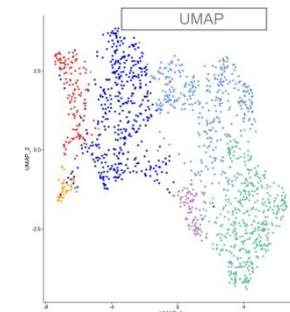
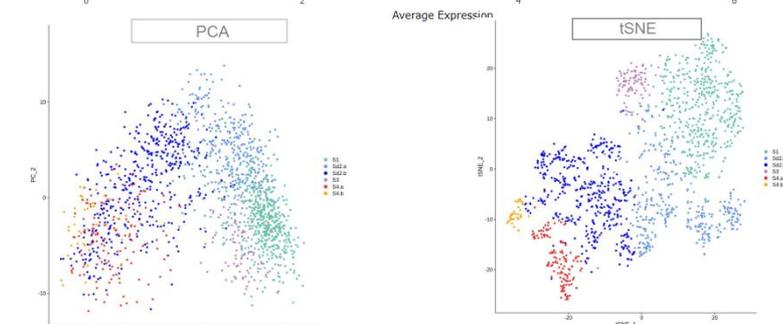
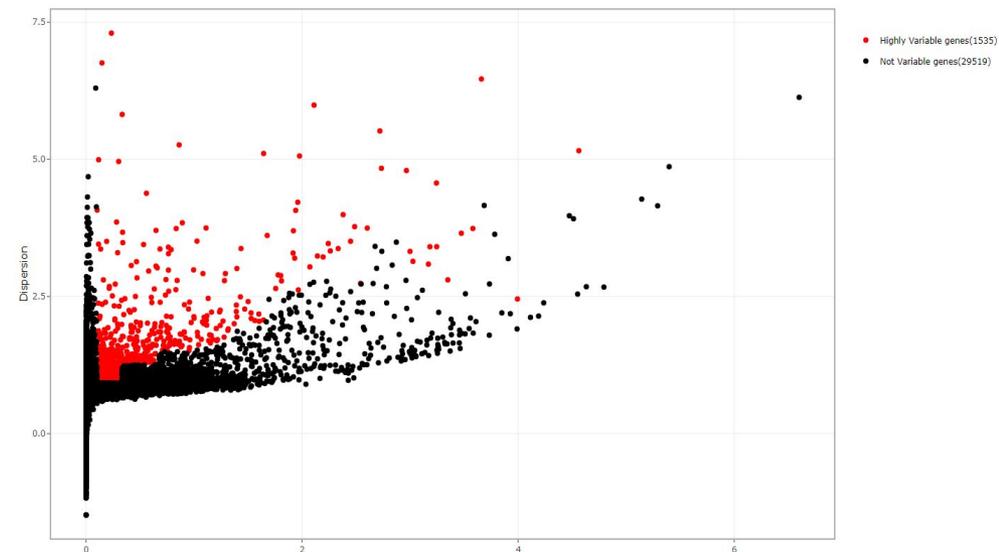
t-Distributed Stochastic Neighbor Embedding (t-SNE)

- Can capture subtle local patterns of expression in the data
- Places cells with similar local neighborhoods in high dimensional space together in low dimensional space
- It may fail to give a precise representation of clusters' size and distances

Uniform Manifold Approximation and Projection (UMAP)

- Preserves better the global structure of the data
- Faster runtime than tSNE
- It may fail to illuminate the lineage structure of the data

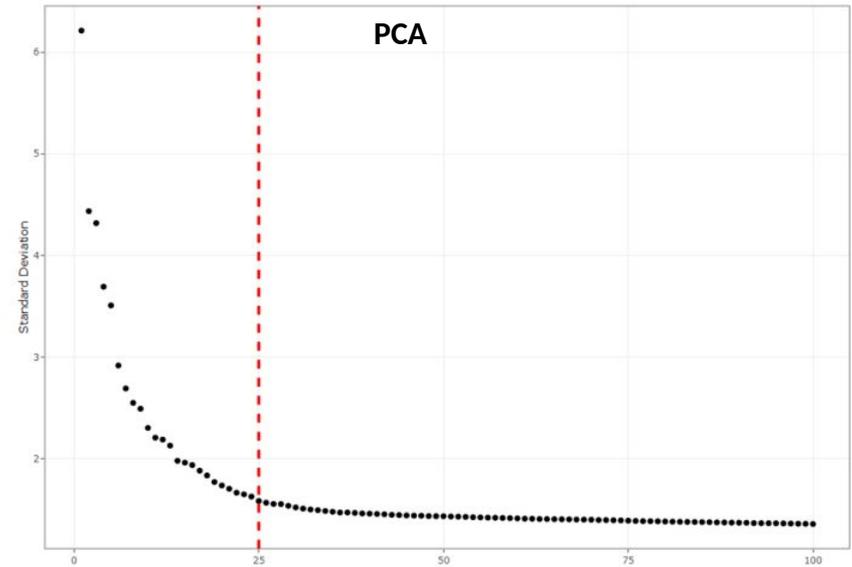
Highly variable genes



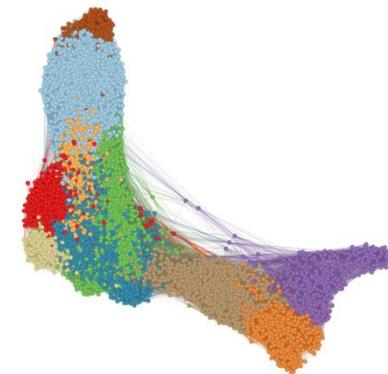
scRNAseq pipeline

4. Clustering analysis

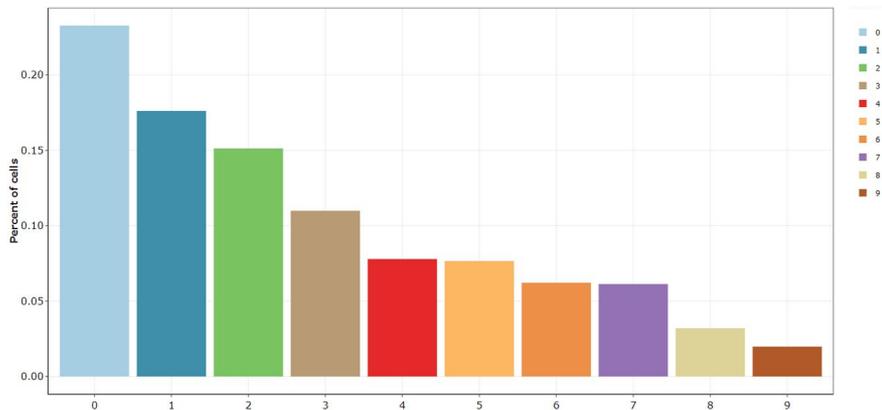
- ❑ Creation of a Shared nearest neighbor (SNN) graph
- ❑ Clusters represent
 - cell population
 - cell sub-population
 - cell state



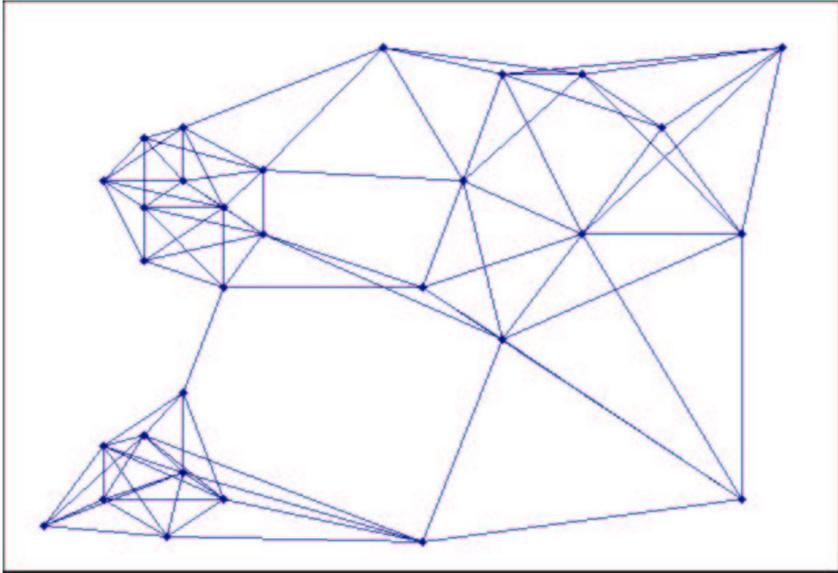
Shared Nearest Neighbors graph



Clustering

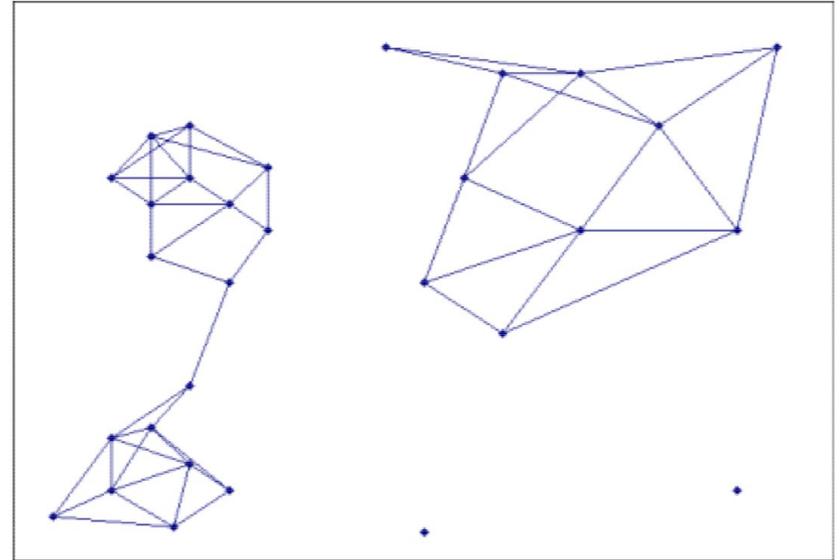


Shared Nearest Neighbors graph



(a) Near Neighbor Graph.

$k=5$



(b) Unweighted Shared Nearest Neighbor.

link if p_1 and p_2 have each other in their nearest neighbor lists

scRNAseq pipeline

5. Differential Expression Analysis

Design of the analysis

- Cells belonging to one cluster VS Cells belonging to another
- Cells belonging to one cluster VS Cells belonging to the rest of the clusters

Selection of D.E.A test

- Wilcoxon test, Student's t-test, Poisson, MAST *

Marker gene analysis

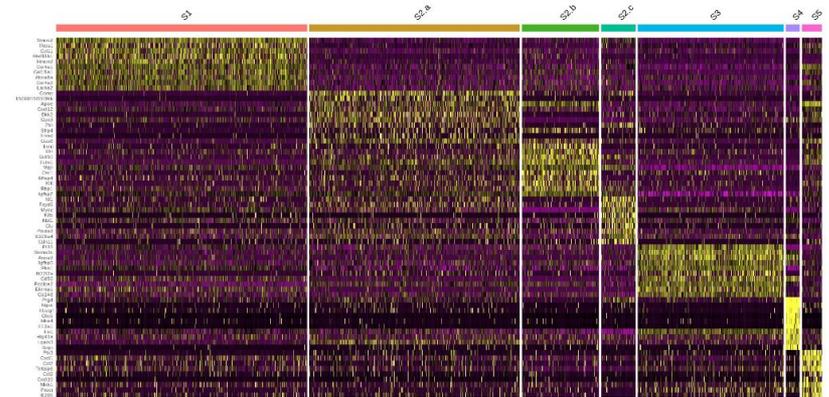
- Identify marker genes per cluster
- Those genes can distinguish one cluster from the rest
- High average expression in cells of the cluster, low in the other cells

- Wilcoxon test
- $\log_{2}FC \geq 0.25$
- $Pval < 0.01$
- Percentage of expression $\geq 25\%$

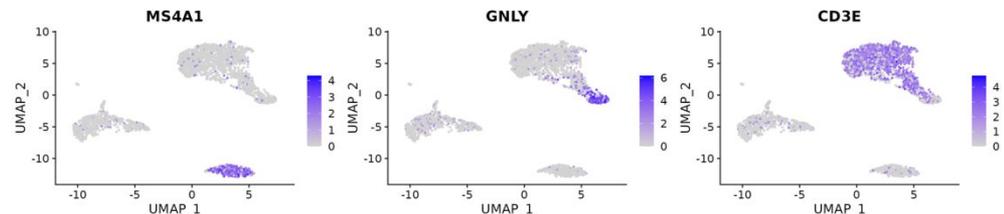
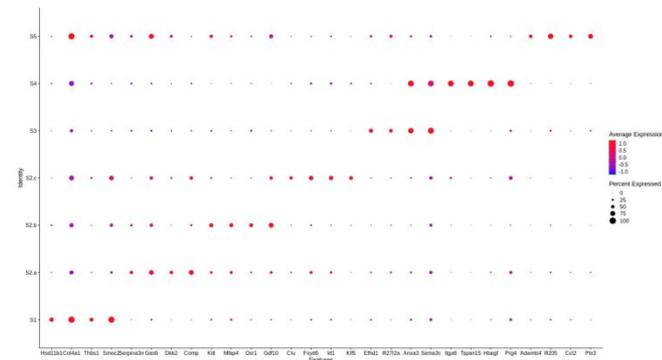
Inspection of top marker genes

- Feature plots in UMAP space
- Color denotes normalized expression

Differential expression analysis



Marker genes for each cluster

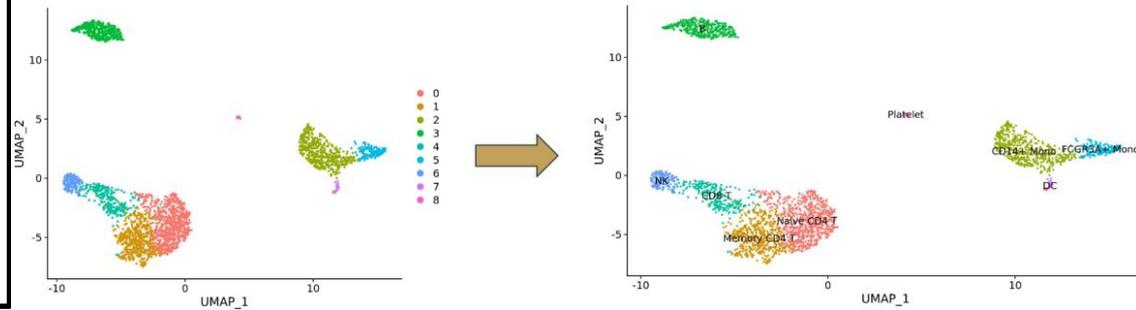


* MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data

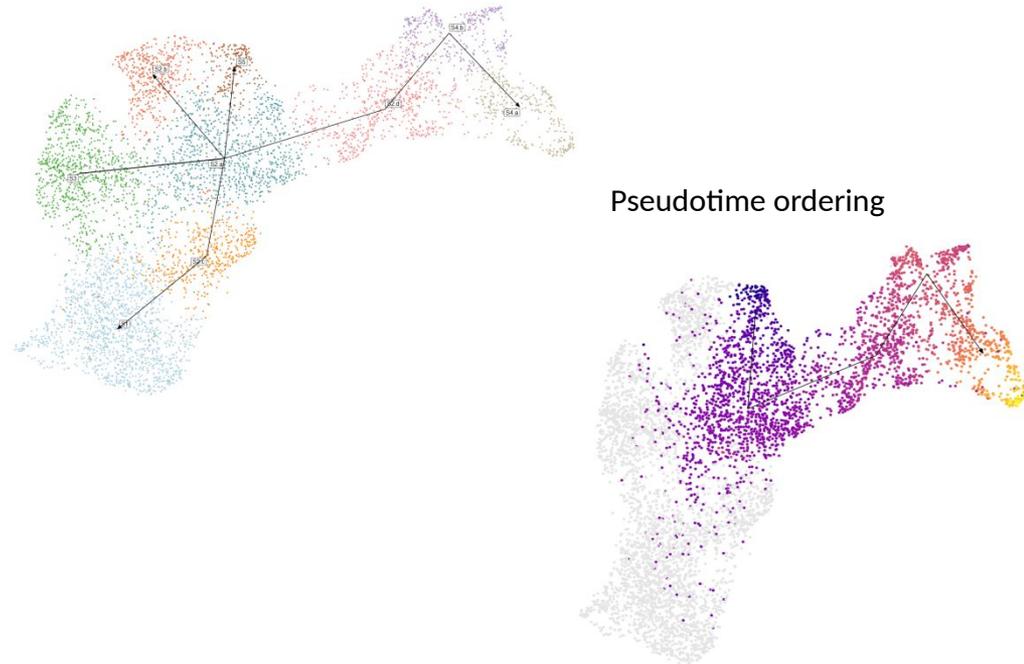
scRNAseq pipeline

6. Cluster annotation

- Compare cluster marker genes to canonical markers for different cell types from the literature
- Using computational methods, match cluster labels from a different dataset (e.g. a cell atlas of the studied organism) to your own clusters



Minimum spanning tree



7. Trajectory-Pseudotime analysis

- Infer the lineage structure of the dataset
- Order the cells along the predicted topology
- PCs as input
- Output in UMAP plot

- Useful links

- [Seurat - Guided Clustering Tutorial](#)

- https://satijalab.org/seurat/articles/pbmc3k_tutorial.html

Online scRNAseq analysis

- <https://singlecell.usegalaxy.eu/>

- <http://scala.fleming.gr/app/scala>

- <https://crescent.cloud/>

Isoform quantitation tools in the literature

- 26 tools found in literature that support transcript DE
 - 10 still active
 - 6 user friendly enough for being used (!)
 - open-source with source code released under a license

	Name	Since	Citations
1	Tuxedo Suite	2012	5390
2	RSEM	2011	4068
3	New Tuxedo Suite	2016	215
4	sleuth	2017	169
5	BitSeq	2012	164
6	EBSeq	2015	4

(slide by A. Dimopoulos)

De-novo genome sequence assembly, Genome-Based and Genome-Free Transcript Reconstruction and Analysis Using RNA-Seq Data

based on material from Mathias Haimel, EBI

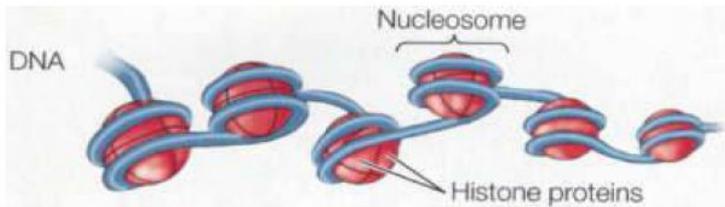
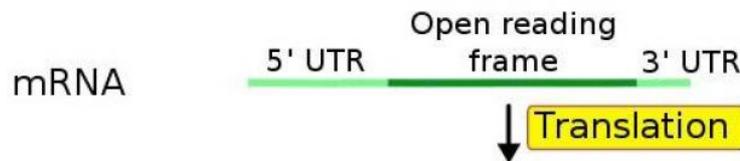
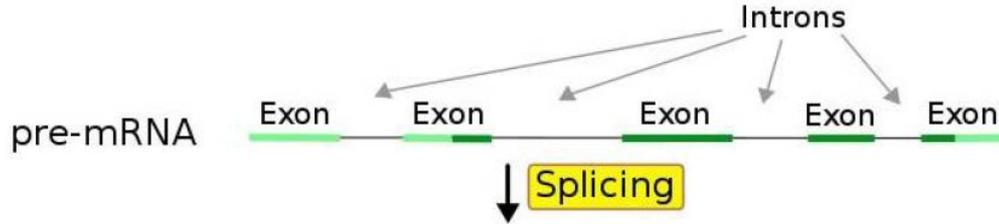
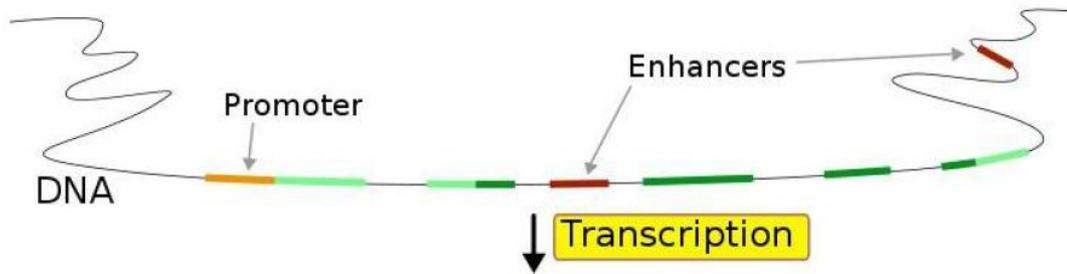
https://www.ebi.ac.uk/training/online/sites/ebi.ac.uk.training.online/files/user/18/private/velvet_1.pdf

and Brian Haas

Broad Institute, modified by M. Reczko



Next Generation Sequencing

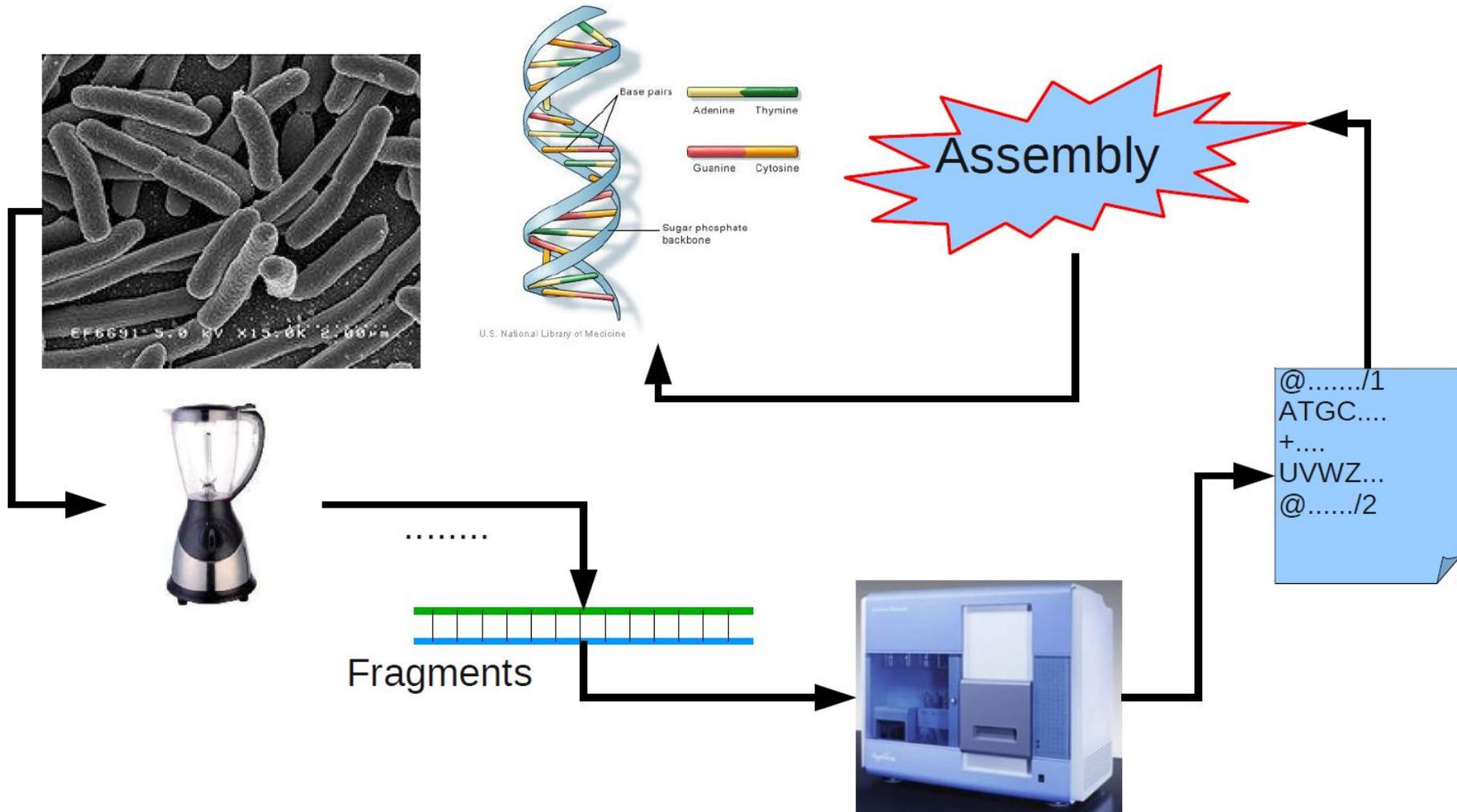


Whole Genome sequencing

RNA-Seq
Whole Transcriptome sequencing

ChIP-Seq
Chromatin Immunoprecipitation with DNA sequencing

Next Generation Sequencing



De novo transcriptome assembly

No genome required

Empower studies of non-model organisms

- expressed gene content
- transcript abundance
- differential expression

Shortest Superstring Problem

- Problem: Given a set of strings, find a shortest string that contains all of them
- Input: Strings s_1, s_2, \dots, s_n
- Output: A string s that contains all strings s_1, s_2, \dots, s_n as substrings, such that the length of s is minimized
- **Complexity**: NP – complete
- **Note**: this formulation does not take into account sequencing errors

Shortest Superstring Problem: Example

The Shortest Superstring problem

Set of strings: {000, 001, 010, 011, 100, 101, 110, 111}

Concatenation

Superstring

000 001 010 011 100 101 110 111

010

110

011

Shortest

superstring

000

0 0 0 1 1 1 0 1 0 0

001

111

101

100

PHASE TWO: INTERPRETATION



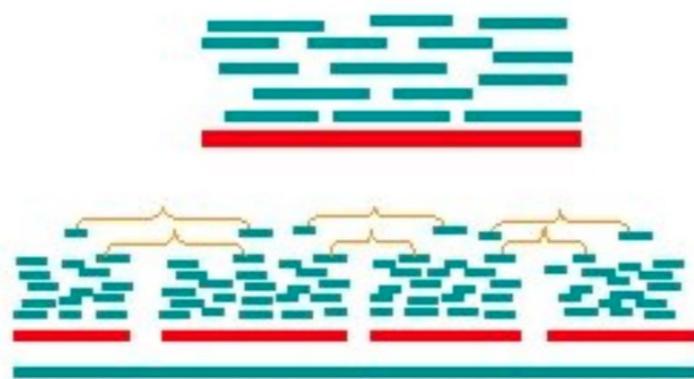
Overlap-Layout-Consensus

Assemblers: ARACHNE, PHRAP, CAP, TIGR, CELERA

Overlap: find potentially overlapping reads



Layout: merge reads into contigs and contigs into supercontigs



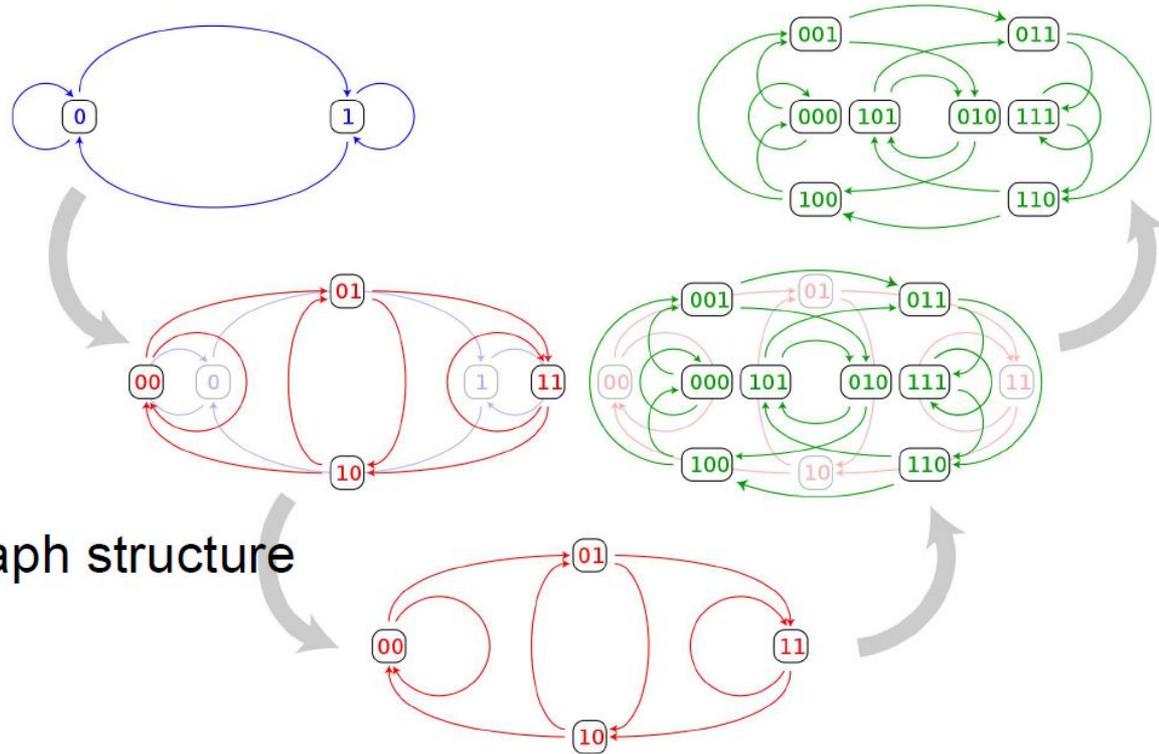
Consensus: derive the DNA sequence and correct read errors

..ACGATTACAATAGGTT..

The General Approach to
De novo DNA/RNA-Seq Assembly
Using De Bruijn Graphs

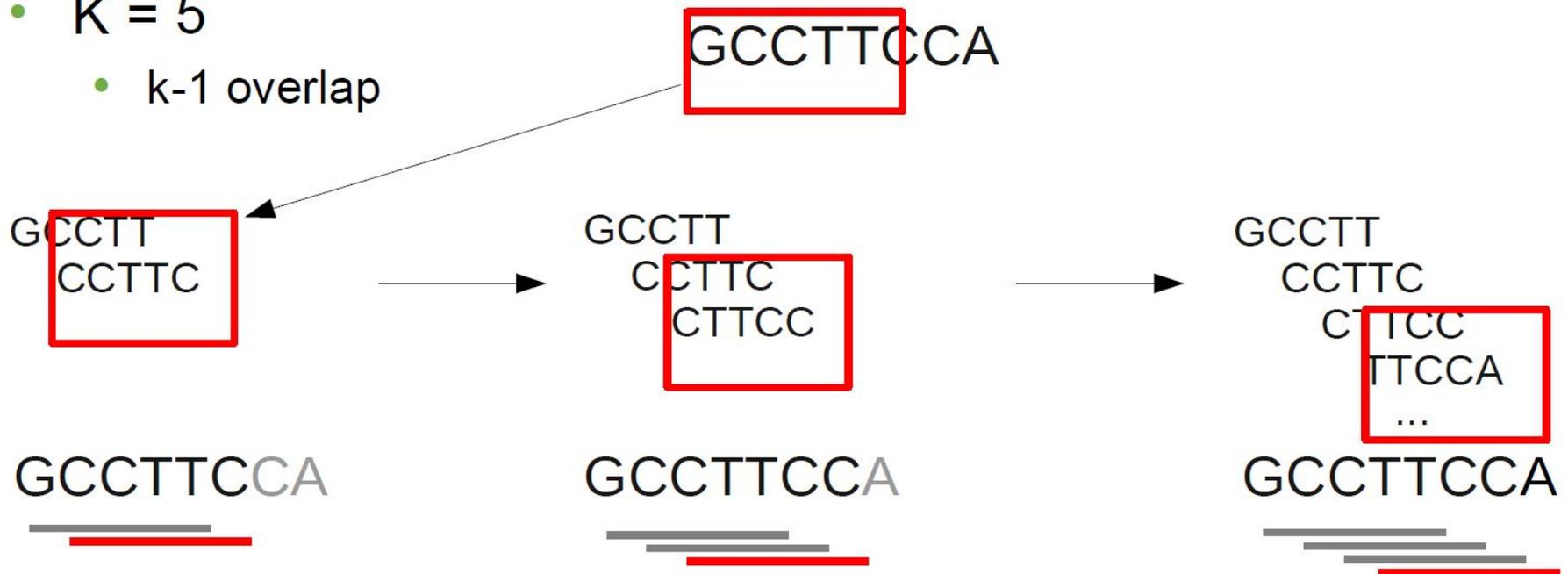
De Bruijn graph

- A concept in combinatorial mathematics
 - In combinatorics, de bruijn graph is usually fully connected
 - http://en.wikipedia.org/wiki/De_Bruijn_graph
- de bruijn sequence
 - Related concept
 - Path through graph
- Velvet
 - de Bruijn inspired graph structure



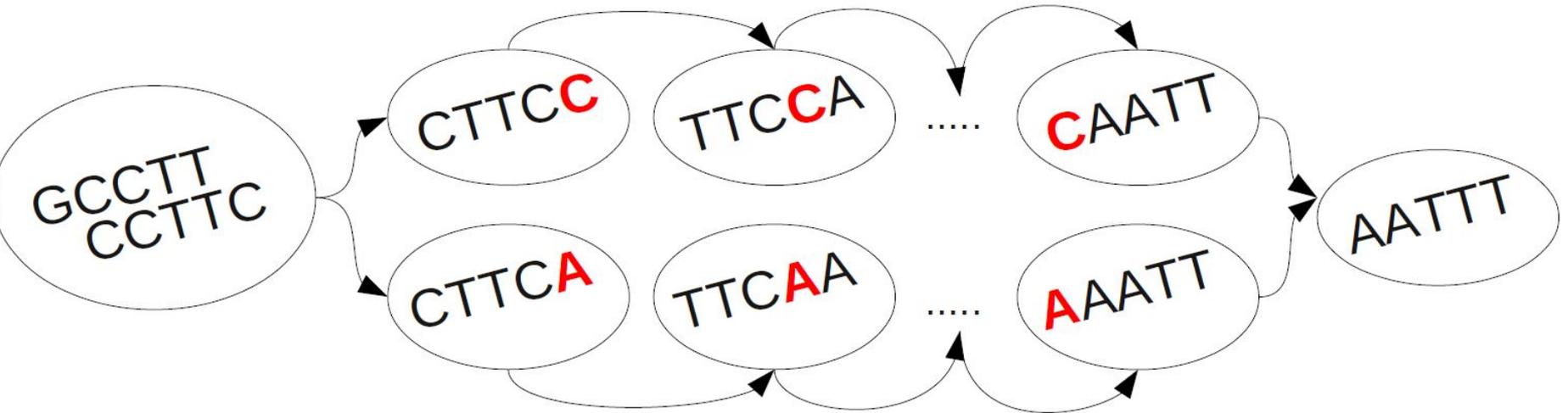
De Bruijn graph (Velvet)

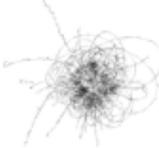
- Representation of
 - a sequence based on short words (k-mers)
 - overlaps between words
- K-mer: word of length k
- K = 5
 - k-1 overlap



De Bruijn graph (Velvet)

GCCTT**C**AATTT
GCCTT**A**AATTT





Example

TAGTCGAGGCTTTAGATCCGATGAGGCTTTAGAGACAG

```
AGTCGAG CTTTAGA  CGATGAG CTTTAGA
GTCGAGG  TTAGATC  ATGAGGC    GAGACAG
          GAGGCTC   ATCCGAT AGGCTTT  GAGACAG
AGTCGAG   TAGATCC ATGAGGC  TAGAGA
TAGTCGA  CTTTAGA CCGATGA    TTAGAGA
          CGAGGCT  AGATCCG TGAGGCT  AGAGACA
TAGTCGA GCTTTAG TCCGATG  GCTCTAG
          TCGACGC    GATCCGA GAGGCTT  AGAGACA
TAGTCGA   TTAGATC GATGAGG TTTAGAG
          GTCGAGG TCTAGAT  ATGAGGC  TAGAGAC
          AGGCTTT  ATCCGAT AGGCTTT  GAGACAG
AGTCGAG   TTAGATT  ATGAGGC  AGAGACA
          GGCTTTA  TCCGATG    TTTAGAG
          CGAGGCT TAGATCC  TGAGGCT  GAGACAG
AGTCGAG  TTTAGATC  ATGAGGC  TTAGAGA
          GAGGCTT  GATCCGA  GAGGCTT  GAGACAG
```



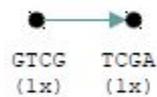
Example

Read: GTCGAGG

●
GTCG
(1x)

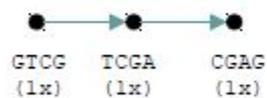
Example

Read: GTCGAGG



Example

Read: GTCGAGG



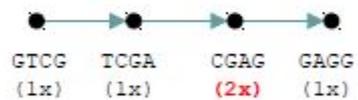
Example

Read: GTCGAGG



Example

New read: CGAGGCT



Example

Read: CGAGGCT



Example

Read: CGAGGGCT



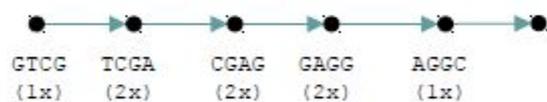
Example

Read: CGAGGGCT



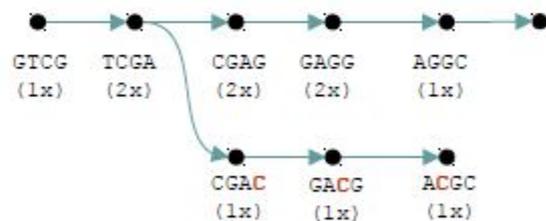
Example

New read: TCGACGC

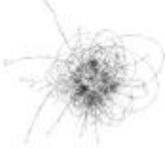


Example

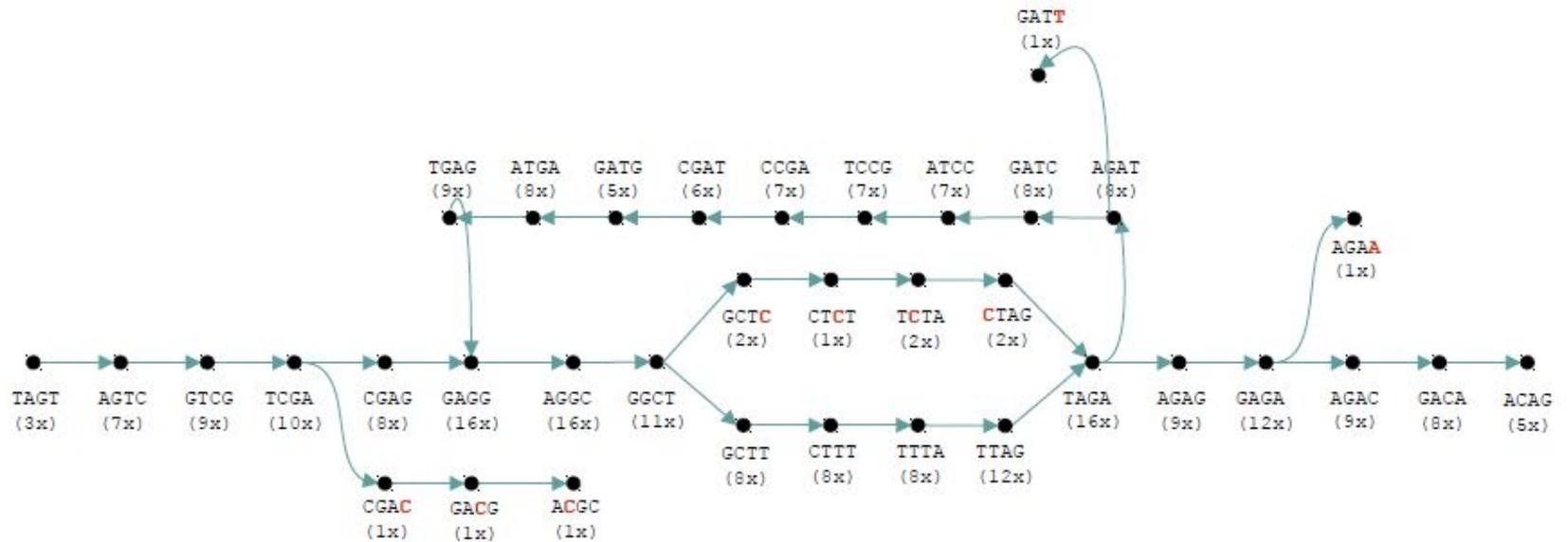
Read: TCGACGC

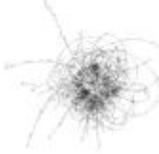


Example



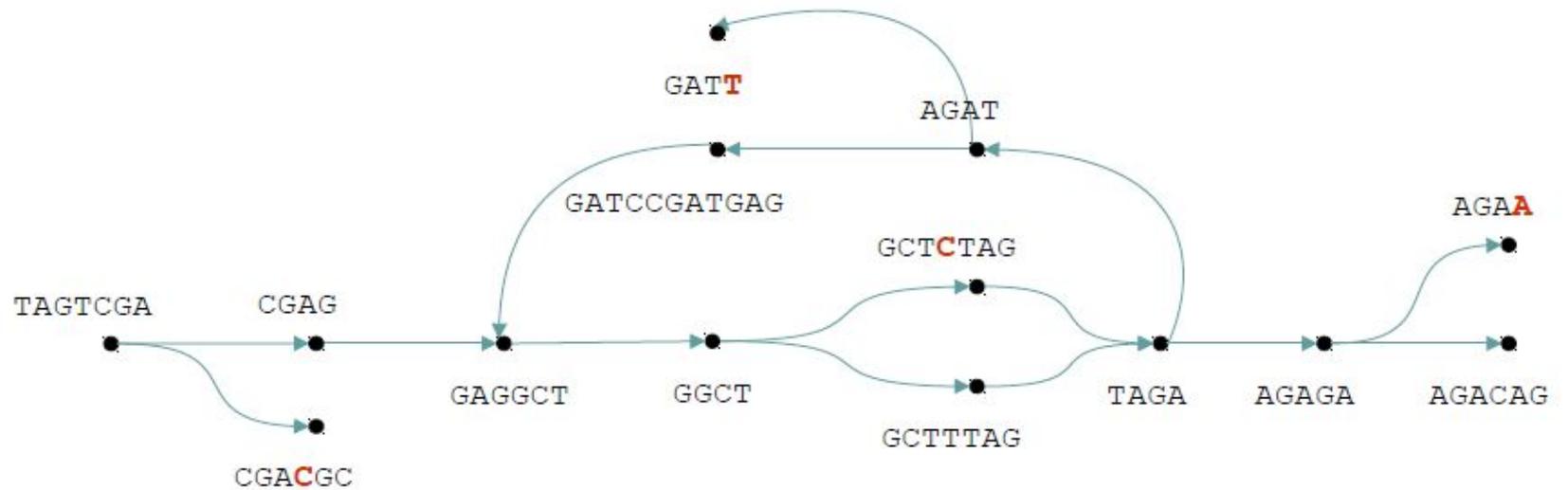
etc...

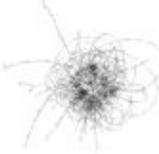




Example

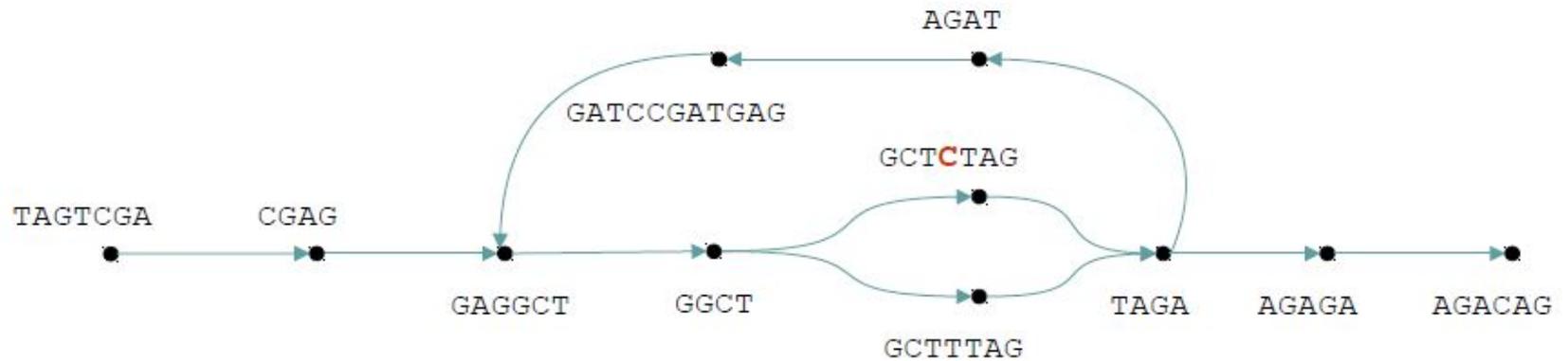
After simplification...

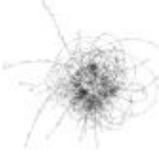




Example

Tips removed...

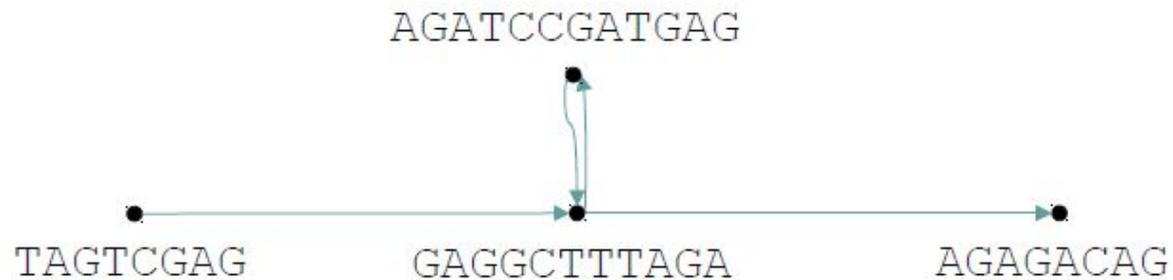




Example

TAGTCGAGGCTTTAGATCCGATGAGGCTTTAGAGACAG

Final simplification...



One possible walk through the graph ...

TAGTCGAG
GAGGCTTTAGA
AGATCCGATGAG
GAGGCTTTAGA
AGAGACAG

2. Sequencing, tools and computers.

2.6 Assembly evaluation

During the assembly optimization will be generated several assemblies. The parameters to evaluate the assembly are:

1. Total Assembly Size,

How far is this value from the estimated genome size

2. Total Number of Sequences (Scaffold/Contigs)

How far is this value from the number of chromosomes.

3. Longest scaffold/contig

4. Average scaffold/contig size

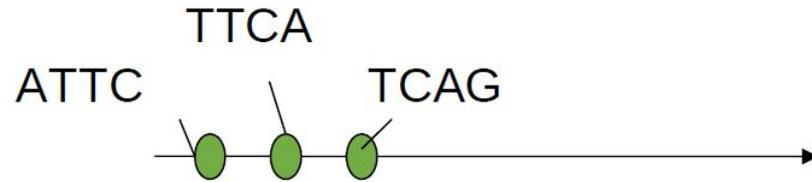
5. N50/L50 (or any other N/L)

Number sequence (N) and minimum size of them (L) that represents the assembly if the sequences are sorted by size, from bigger to small

De Bruijn graph biology extensions (Velvet)

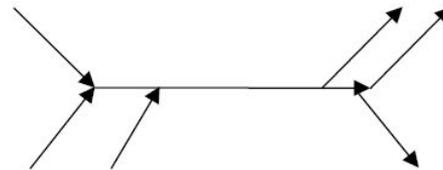
- Handling of reverse strand
 - DNA is read in two directions
 - Paired-end data
- Handling small differences, which are “uninteresting”
 - Errors in sequencing technology
- Memory
 - regularly use 80, 100GB real memory
 - easily get to 1TB real memory requirements

De Bruijn graph representations (Velvet)

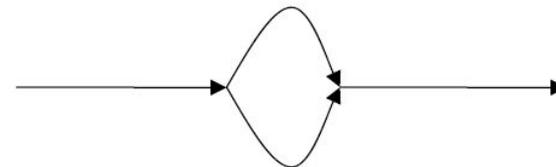


Error free, no repeat,
no polymorphism

Repeat > kmer length

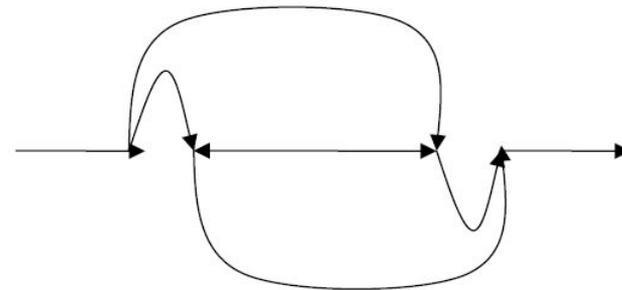


SNP, variant, < kmer length



Structural variant, inversion
Structural variant, deletion...

...



Contrasting Genome and Transcriptome Assembly

Genome Assembly

- Uniform coverage
- Single contig per locus
- Double-stranded

Transcriptome Assembly

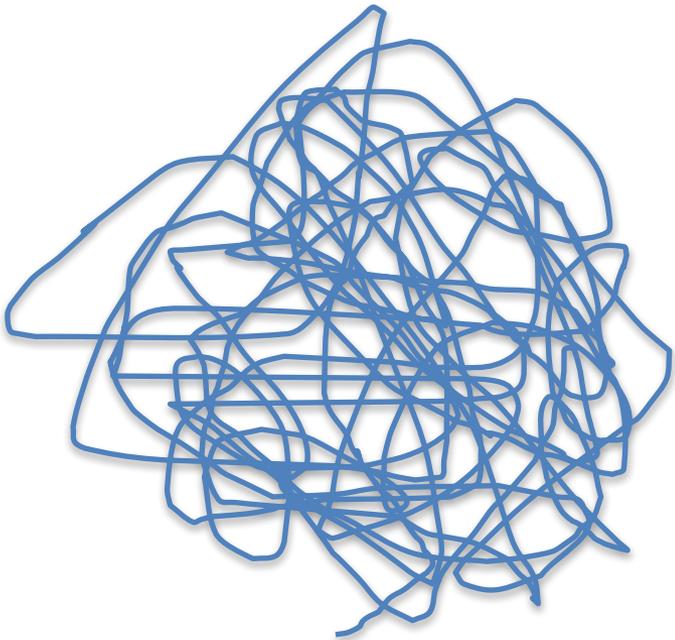
- Exponentially distributed coverage levels
- Multiple contigs per locus (alt splicing)
- Strand-specific



Trinity Aggregates Isolated Transcript Graphs

Genome Assembly

Single Massive Graph



Entire chromosomes represented.

Trinity Transcriptome Assembly

Many Thousands of Small Graphs



Ideally, one graph per expressed gene.

Applied for: Olive fly *Bactrocera oleae* (dakos)



- Ordo: *Diptera*
- Family: *Tephritidae*
- Genus: *Bactrocera*

- Monophagous
- Production losses > 30% possible
- Affects quantity and quality
- Global economic damage estimated: **800.000.000 \$**

Collaborative effort of



Department of Biochemistry and Biotechnology
University of Thessaly

Laboratory of Molecular Biology and Genomics

- K. Mathiopoulos, E. Sagri



ALEXANDER FLEMING
Biomedical Sciences Research Center

- J. Ragoussis, M. Reczko, K. Salpea, V. Harokopos, A. Dimopoulos

Trinity – How it works:



RNA-Seq
reads



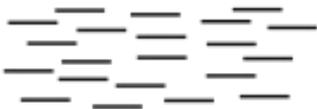
Linear
contigs



de-Bruijn
graphs



Transcripts
+
Isoforms



```
>a121:len=5845  
_____  
>a122:len=2560  
_____  
>a123:len=4443  
_____  
>a124:len=48  
_____  
>a125:len=8878  
_____  
>a126:len=66  
_____
```



...CTTCGCAA...TGATCGGAT...
...ATTTCGCAA...TCATCGGAT...

Thousands of disjoint graphs

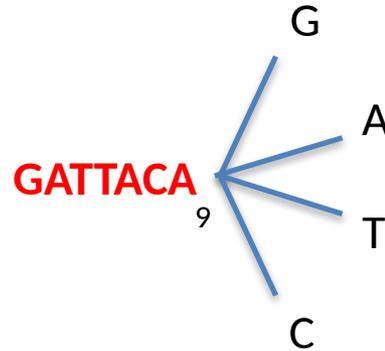


Inchworm Algorithm

Decompose all reads into overlapping Kmers (25-mers)

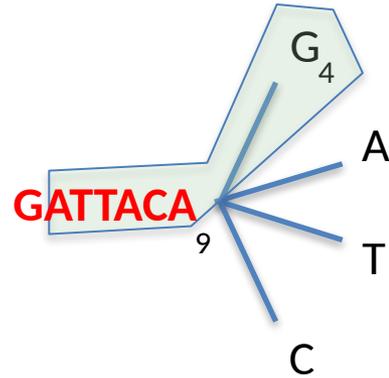
Identify seed kmer as most abundant Kmer, ignoring low-complexity kmers.

Extend kmer at 3' end, guided by coverage.



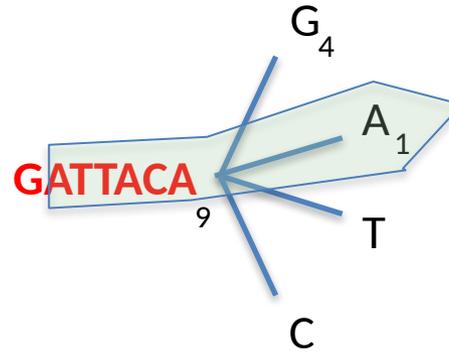


Inchworm Algorithm



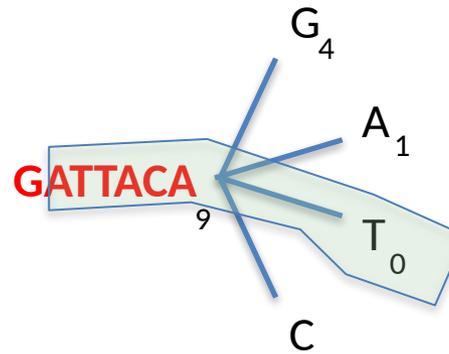


Inchworm Algorithm



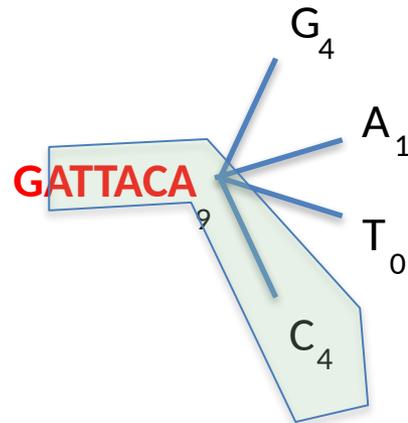


Inchworm Algorithm



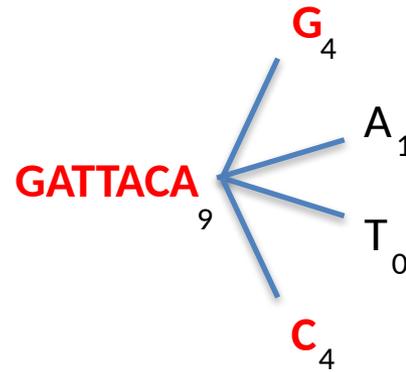


Inchworm Algorithm



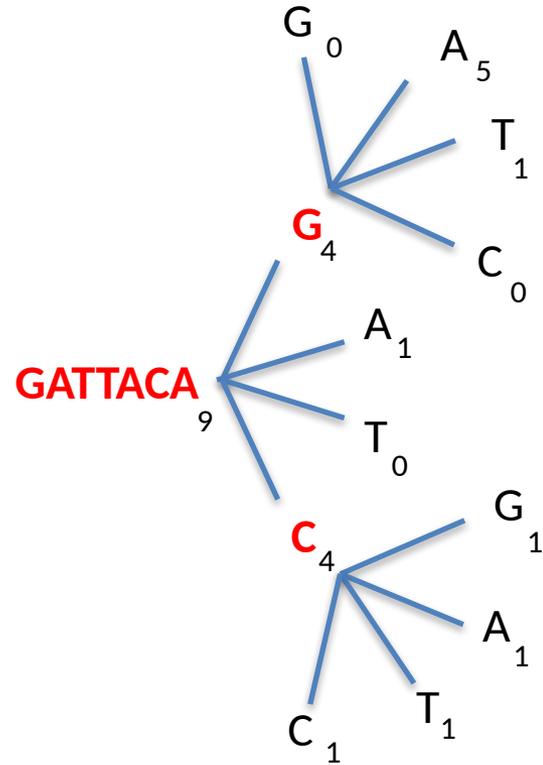


Inchworm Algorithm



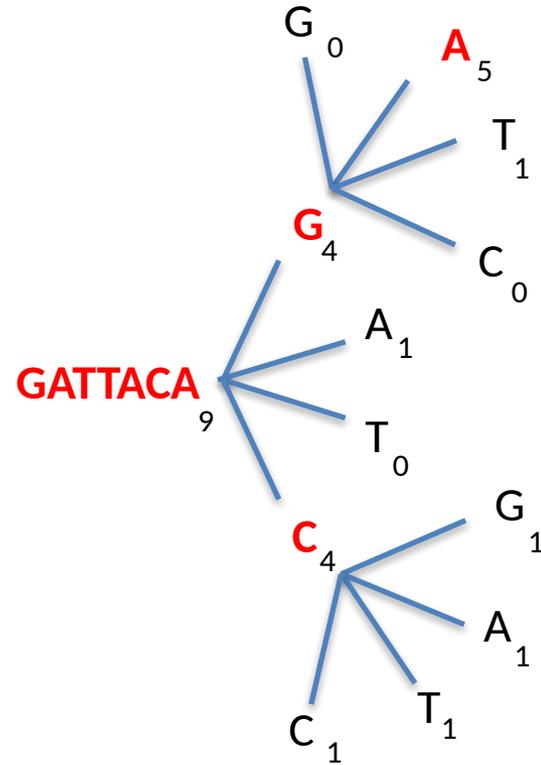


Inchworm Algorithm



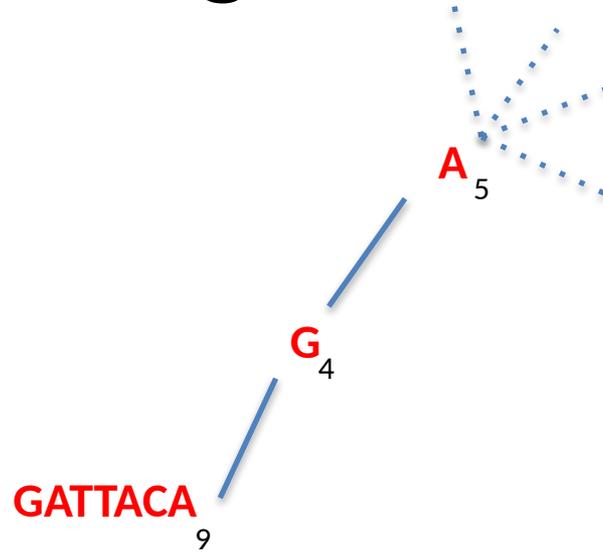


Inchworm Algorithm



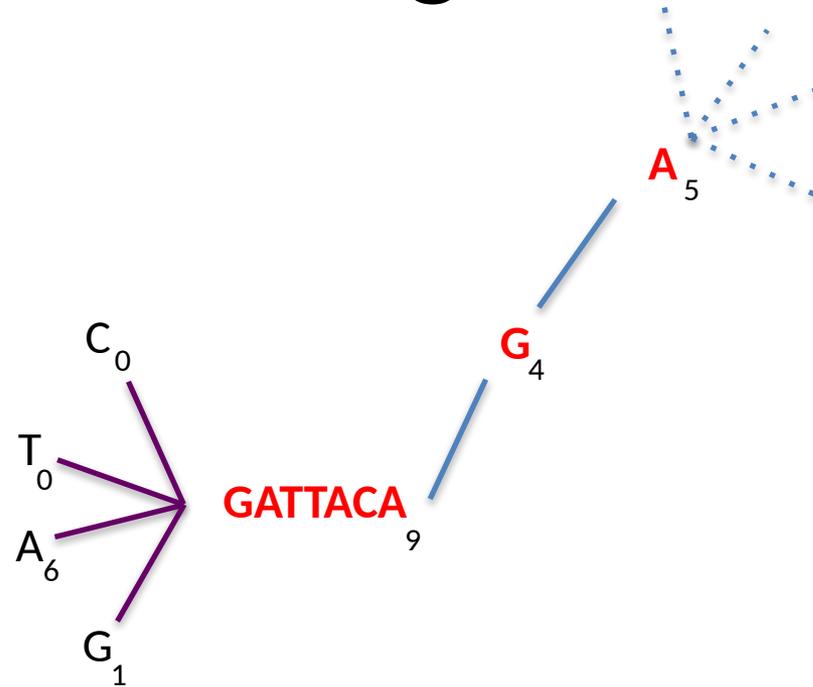


Inchworm Algorithm



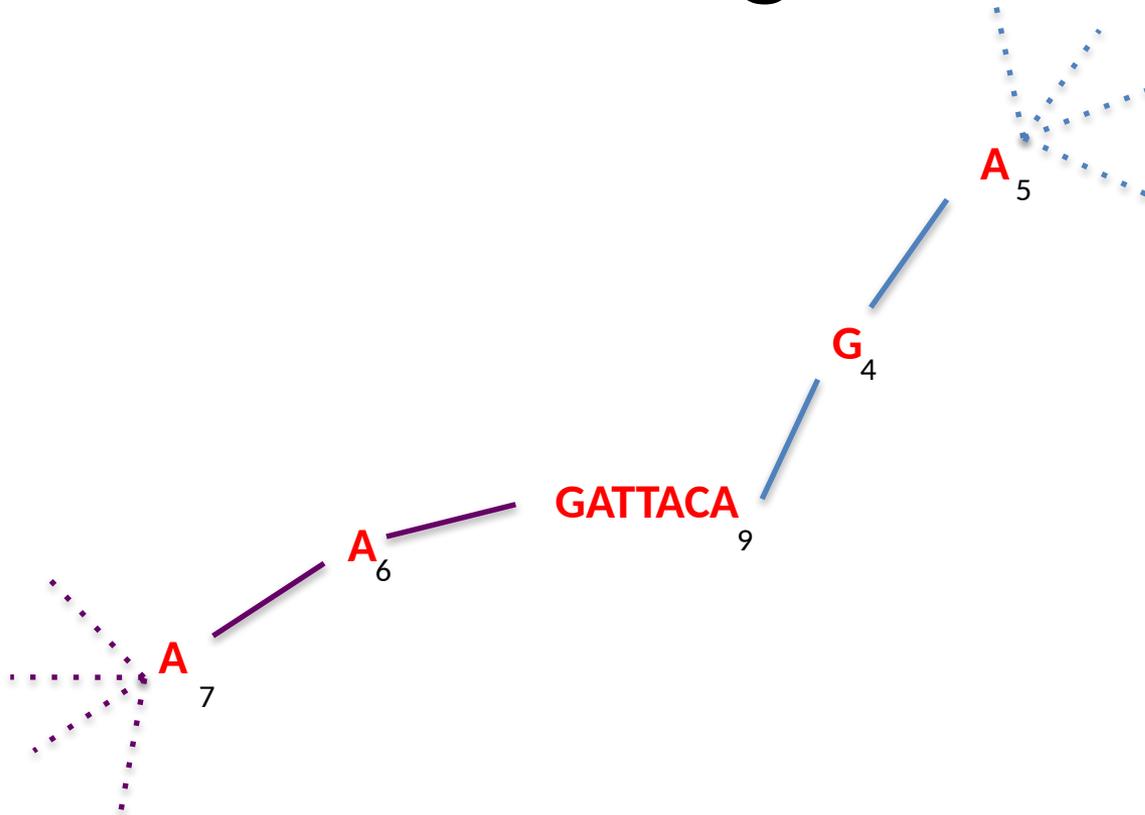


Inchworm Algorithm





Inchworm Algorithm



Report contig:**AAGATTACAGA**....

Remove assembled kmers from catalog, then repeat the entire process.



Inchworm Contigs from Alt-Spliced Transcripts

Expressed isoforms





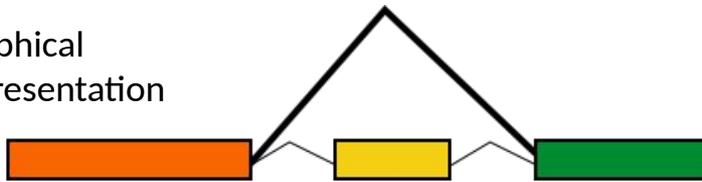
Inchworm Contigs from Alt-Spliced Transcripts

Expressed isoforms

Expression

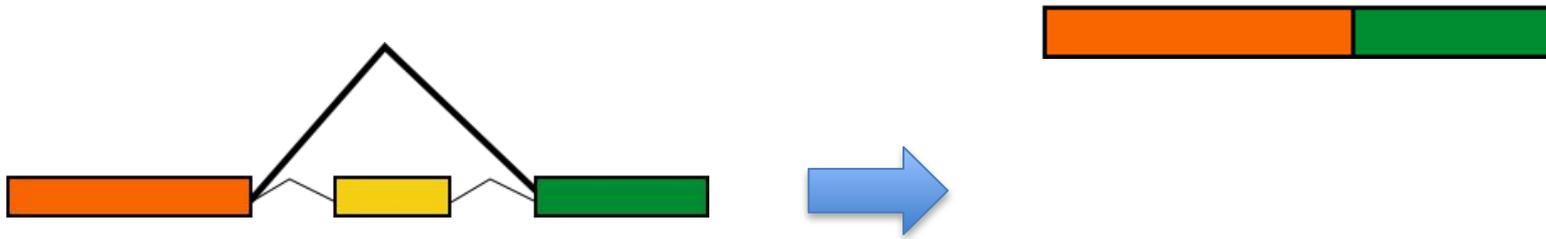


Graphical
representation



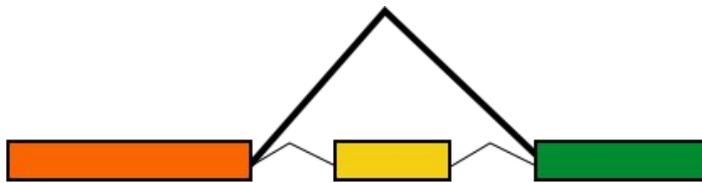


Inchworm Contigs from Alt-Spliced Transcripts





Inchworm Contigs from Alt-Spliced Transcripts



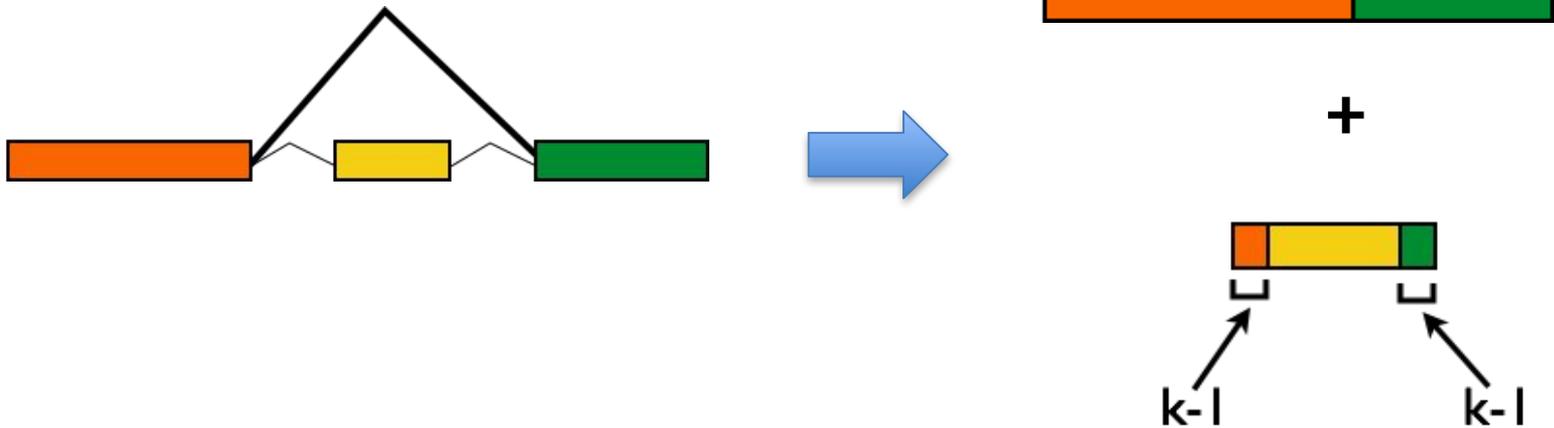
+

No k-mers
in common

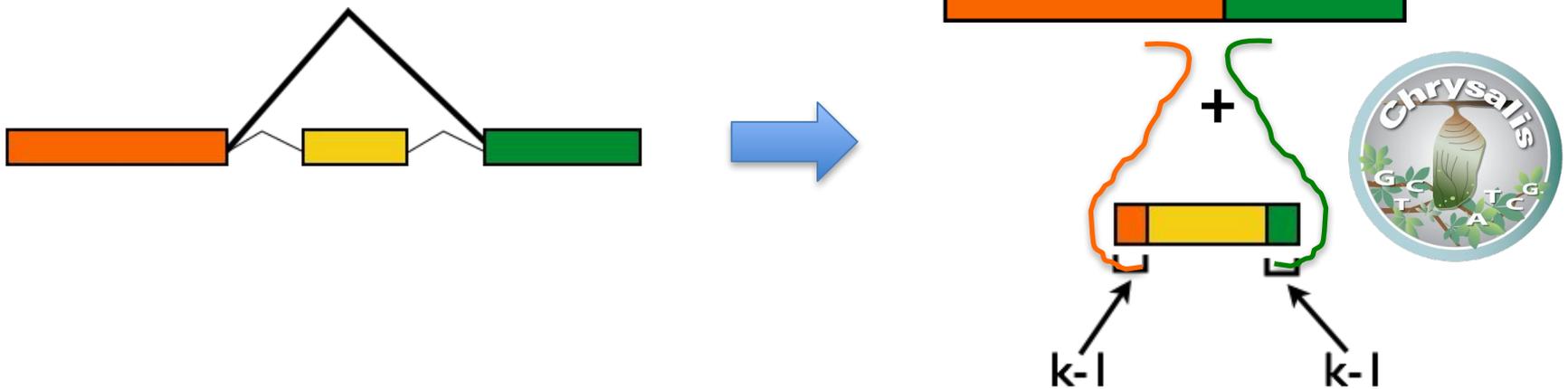




Inchworm Contigs from Alt-Spliced Transcripts



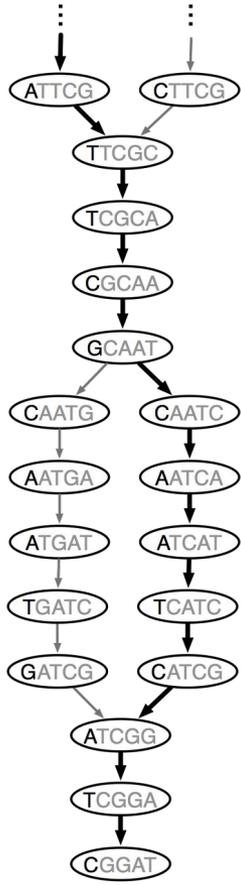
Chrysalis Re-groups Related Inchworm Contigs



Chrysalis uses $(k-1)$ overlaps and read support to link related Inchworm contigs

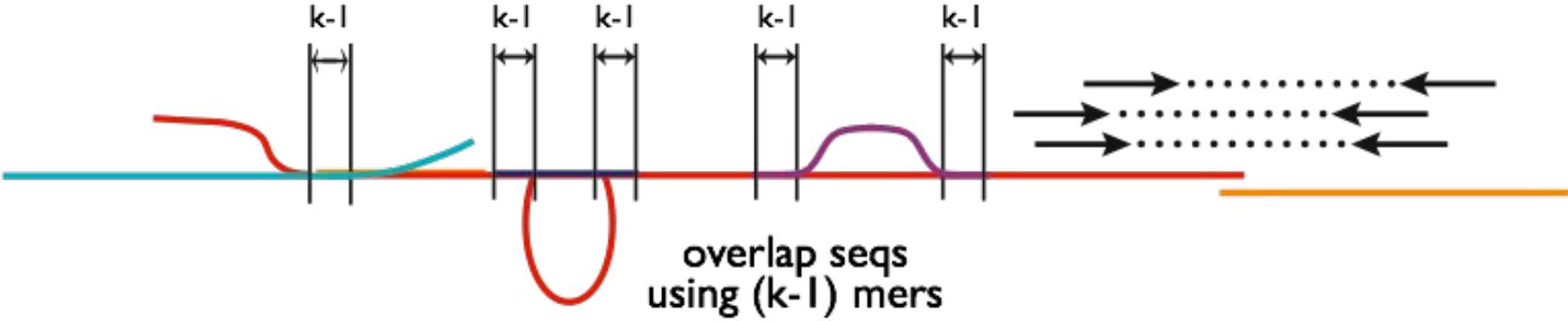
Chrysalis

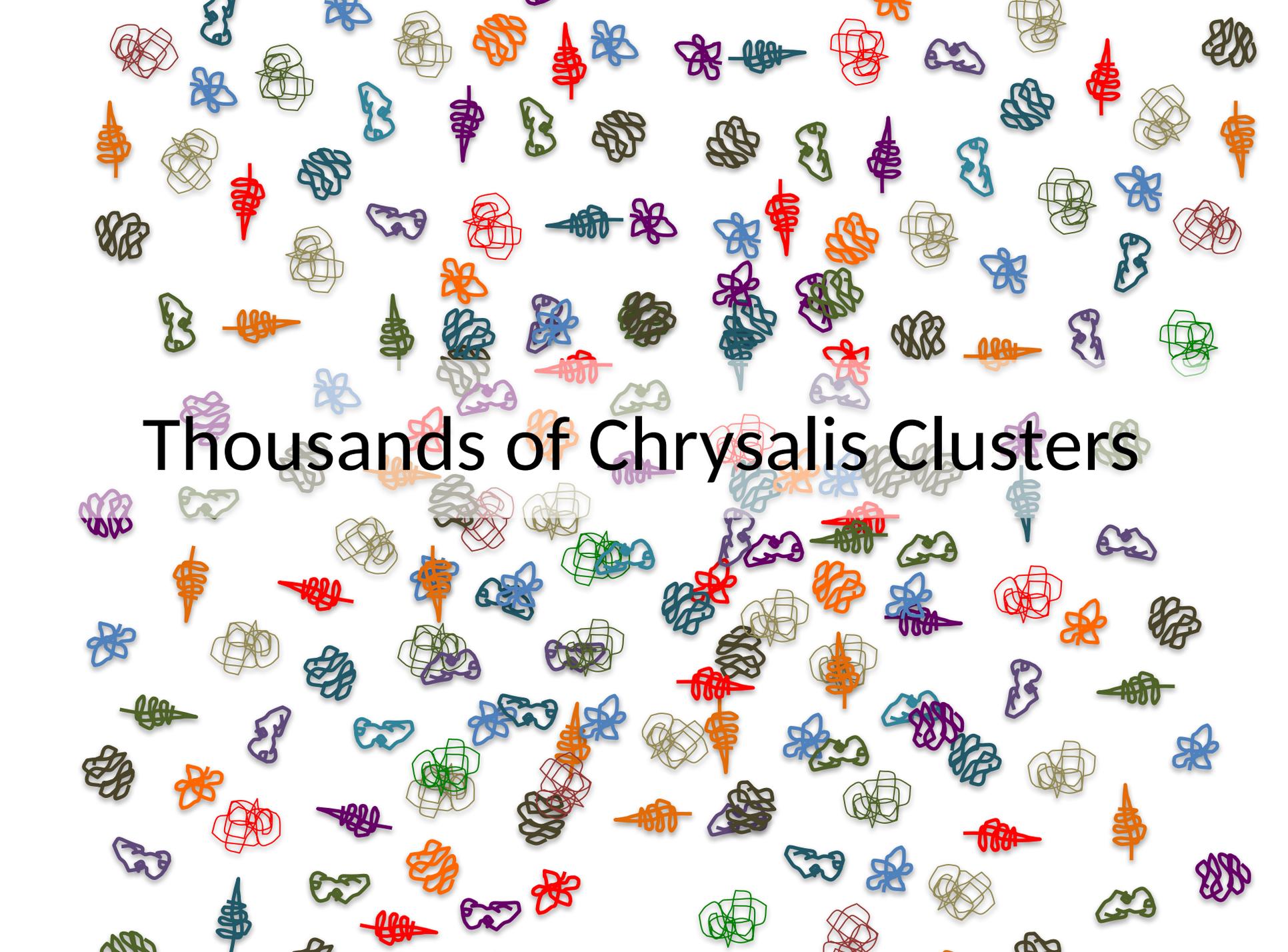
```
>a121:len=5845  
_____  
>a122:len=2560  
_____  
>a123:len=4443  
_____  
>a124:len=48  
_____  
>a125:len=8876  
_____  
>a126:len=68  
_____
```



Integrate isoforms
via $k-1$ overlaps

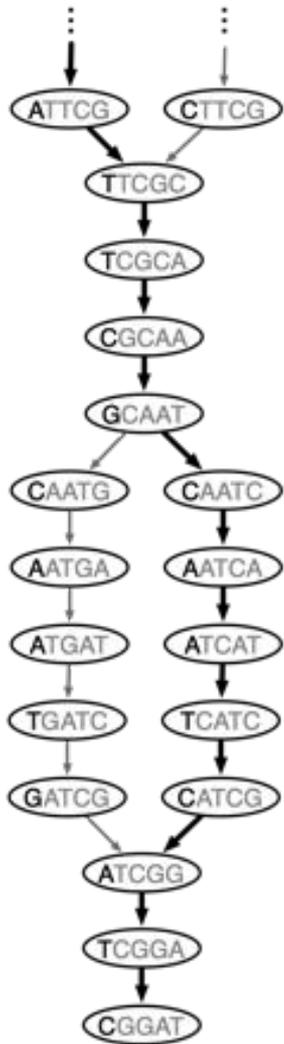
Build de Bruijn Graphs
(ideally, one per gene)



The image features a dense, scattered collection of thousands of small, hand-drawn, abstract shapes. These shapes are rendered in a variety of colors including red, blue, green, purple, orange, and black. Many of the shapes have a complex, multi-layered appearance, resembling tangled lines or intricate patterns that could be interpreted as stylized chrysalis clusters or abstract biological forms. The shapes are distributed across the entire page, with a central area where the text is overlaid.

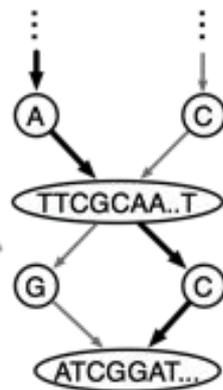
Thousands of Chrysalis Clusters

Butterfly



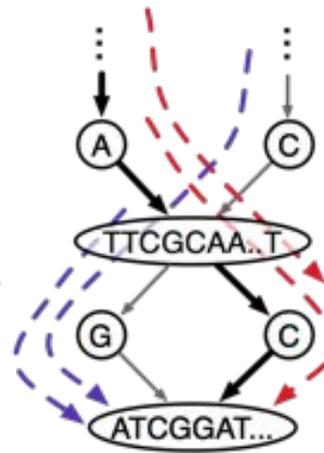
de Bruijn graph

compacting



compact graph

finding paths



compact graph with reads

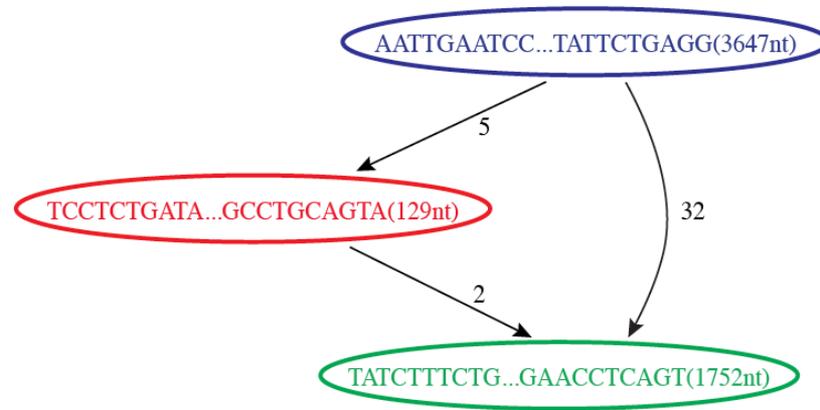
extracting sequences

..CTTCGCAA..TGATCGGAT..
..ATTGCAA..TCATCGGAT..

sequences
(isoforms and paralogs)

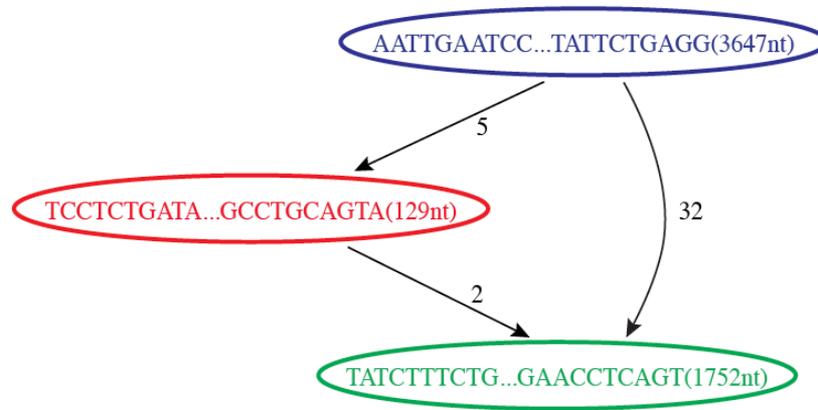
Butterfly Example 1: Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted
Sequence Graph



Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted
Sequence Graph

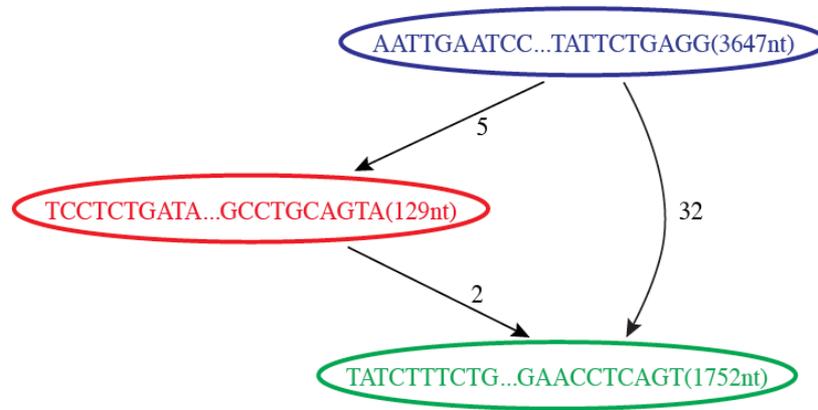


Reconstructed Transcripts



Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted
Sequence Graph

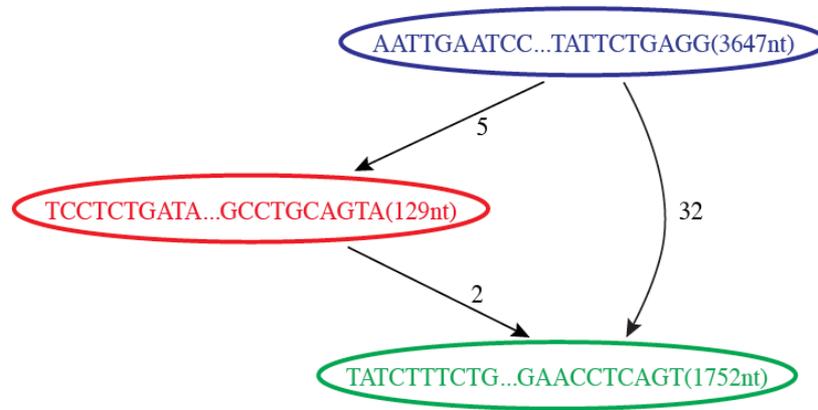


Reconstructed Transcripts



Reconstruction of Alternatively Spliced Transcripts

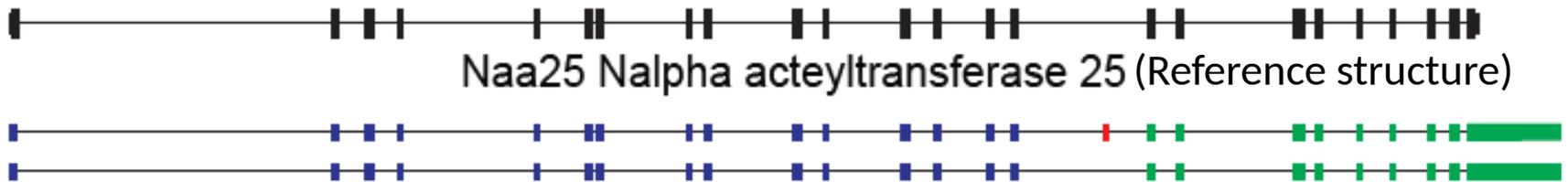
Butterfly's Compacted Sequence Graph



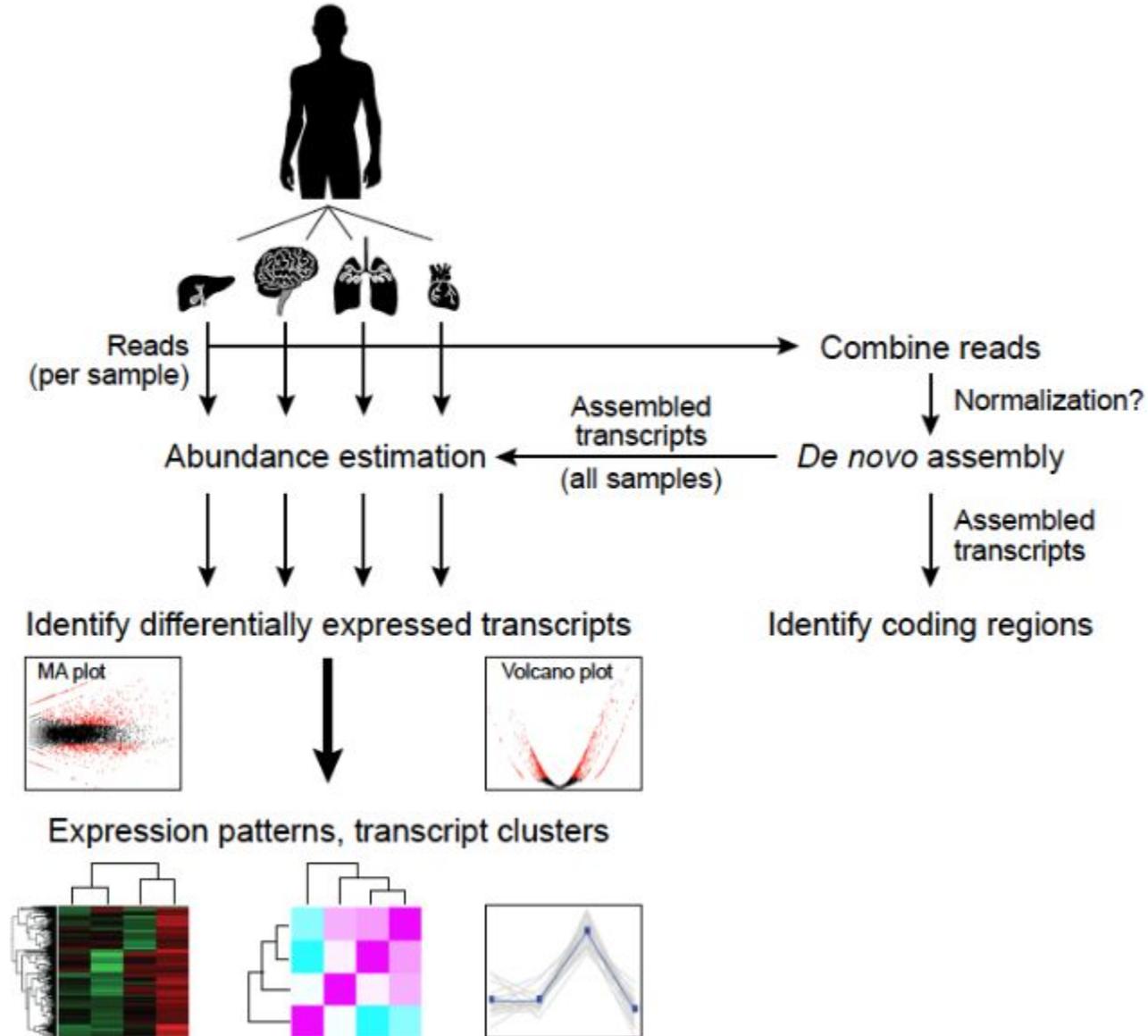
Reconstructed Transcripts



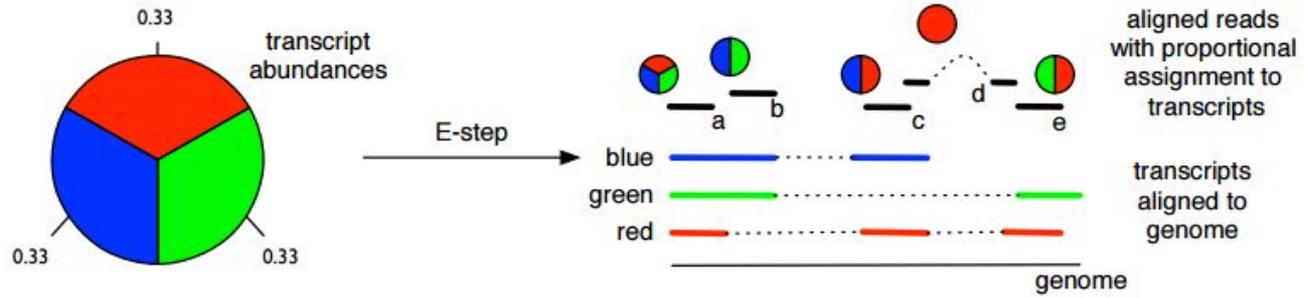
Aligned to Mouse Genome



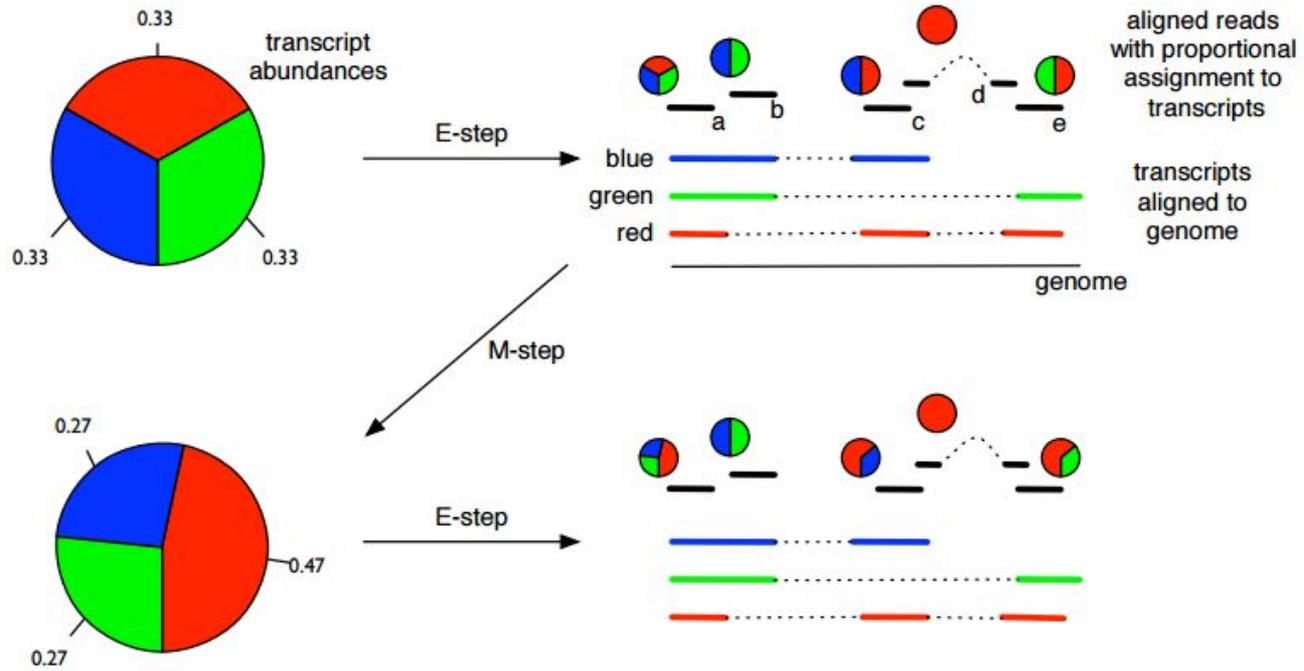
Trinity Demo



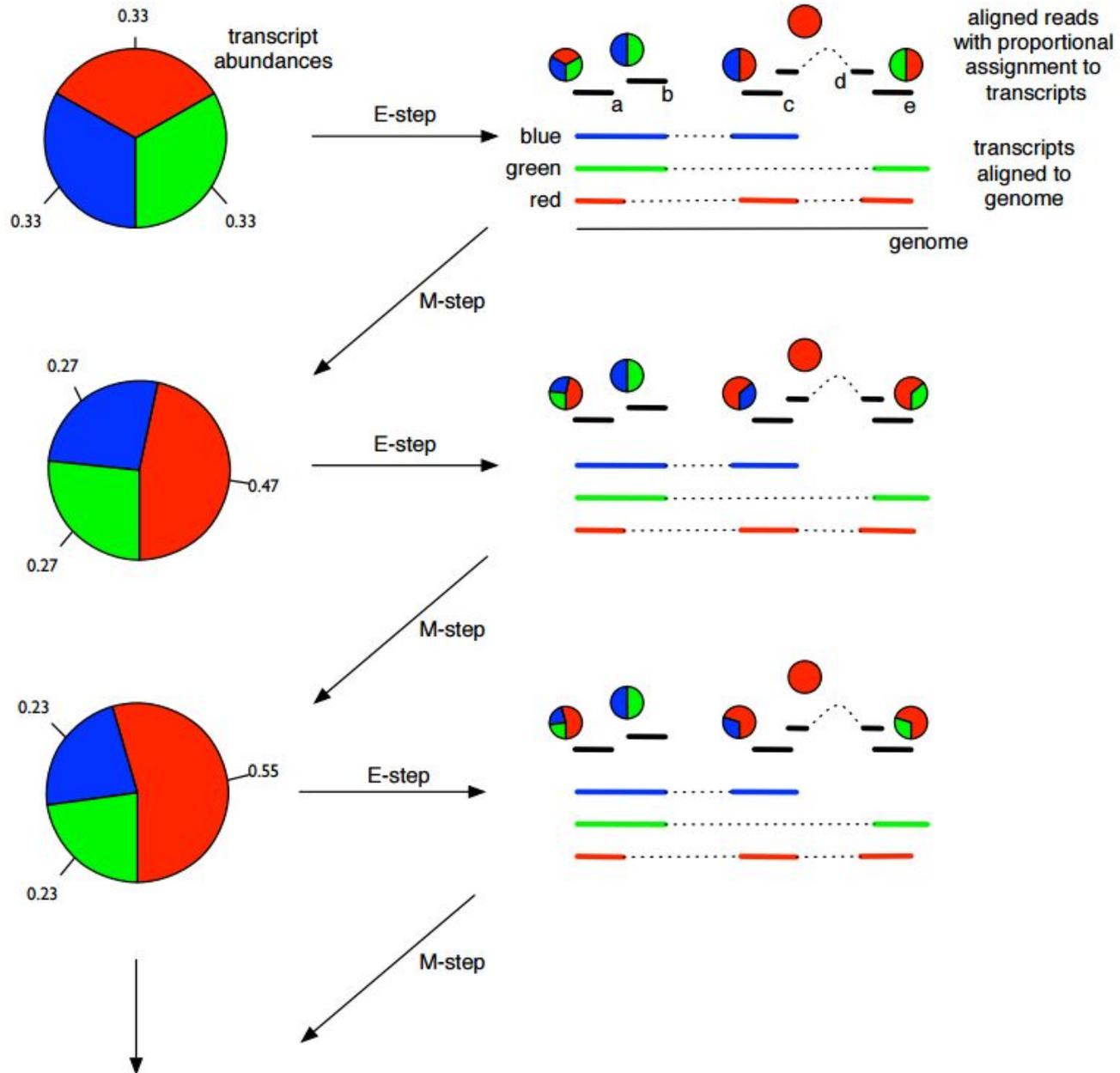
Expectation maximization used in rsem



Expectation maximization used in rsem



Expectation maximization used in rsem



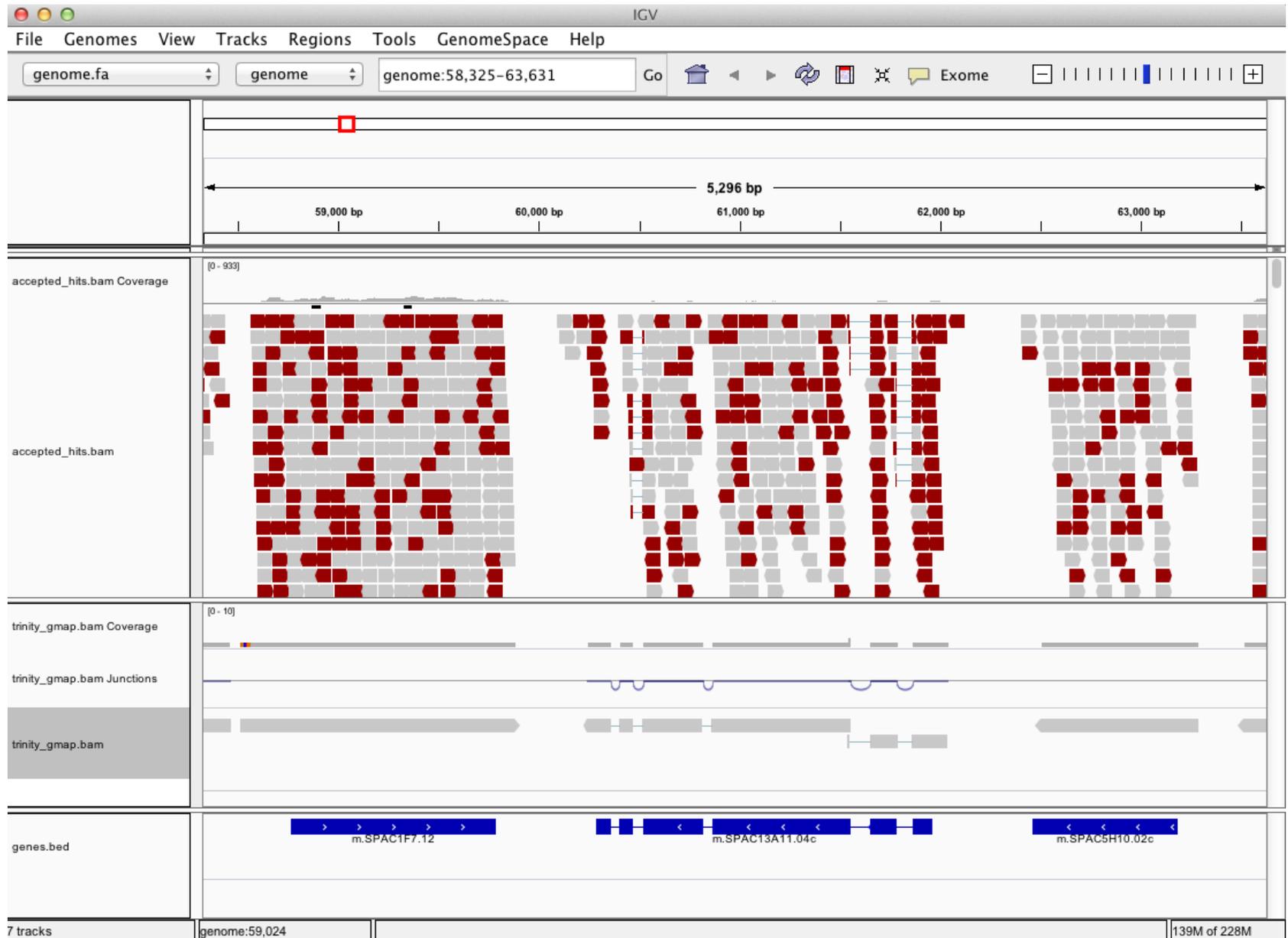
Trinity Demo

- Assemble RNA-Seq using Trinity
- Examine Trinity in context of a genome:
 - Align Trinity transcripts to the genome using GMAP
 - Align rna-seq reads to genome using Tophat
 - Visualize all alignments using IGV

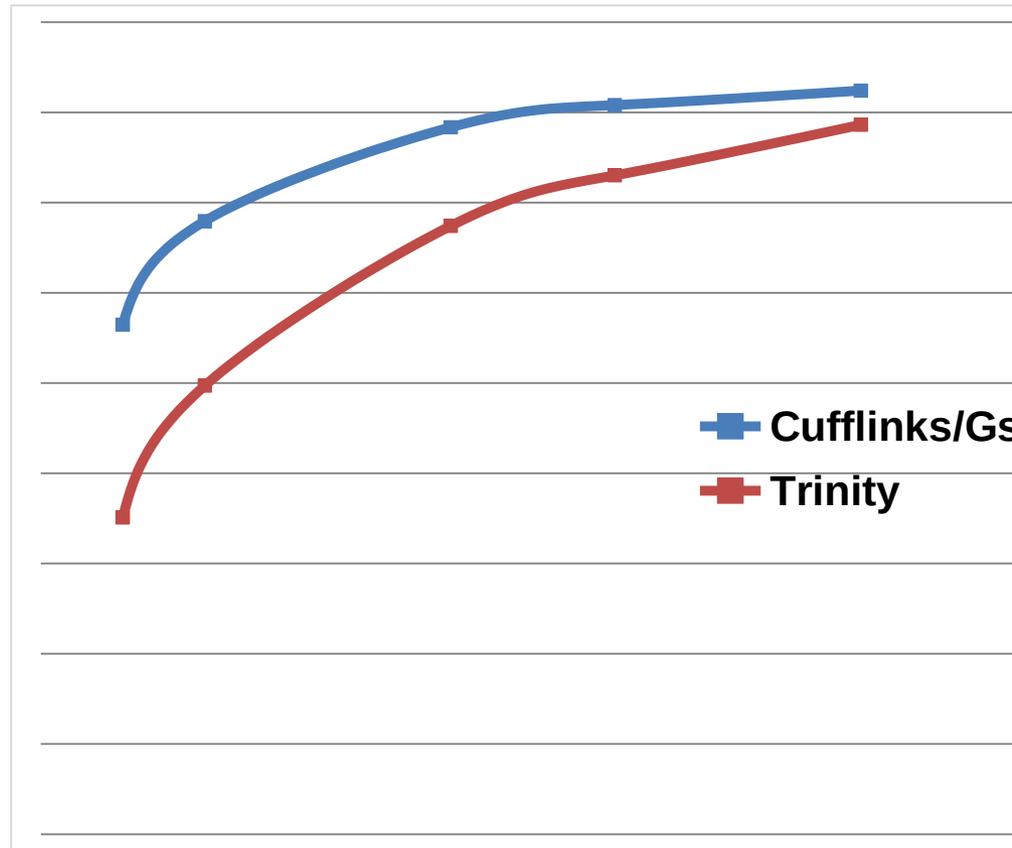
Try yourself:

```
echo 'export TRINITY_HOME=/home/reczko/tools/trinityrnaseq-v2.15.2' >> ~/.bashrc
source ~/.bashrc
export PATH=/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/snap/bin
cd ~/rnaseq_workshop
cp /home/reczko/tools/runTrinityDemo.pl .
./runTrinityDemo.pl
```

Trinity transcripts aligned to genome scaffolds to examine intron/exon structures (Trinity transcripts aligned using GMAP)



Improved reconstruction with deeper sequencing depth and Genome-based reconstruction is more sensitive than de novo methods



Genes w/ fully reconstructed transcripts

■ Cufflinks/Gsnap
■ Trinity

Million PE reads

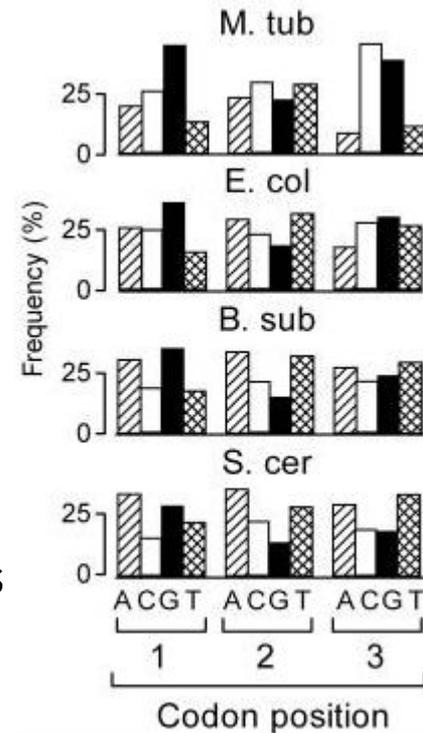


Mouse data

Prediction of coding potential

- Periodicity detection

- Coding sequences have an inherent periodicity of three
- Especially good on long coding sequences
- Auto-correlation
 - Seeking the strongest response when shifted sequence is compared with original
 - Michel (1986), *J. Theor. Biol.* **120**, 223-236.
- Fourier transformation: Spectral analysis
 - Detection of peak at position corresponding to 1/3 of the frequency
 - Silverman and Linsker (1986), *J. Theor. Biol.* **118**, 295-300.



Summary of Key Points

- RNA-Seq is a versatile method for transcriptome analysis enabling quantification and novel transcript discovery.
- Genome-based and genome-free methods exist for transcript reconstruction
- Expression quantification is based on sampling and counting reads derived from transcripts
- Fold changes based on few read counts lack statistical significance.
- Multiple analysis frameworks are available – alternative and often complementary approaches to support biological investigations.

Software Links

- Tuxedo
 - Bowtie: <http://bowtie-bio.sourceforge.net/index.shtml>
 - Tophat: <http://tophat.cbcb.umd.edu/>
 - Cufflinks: <http://cufflinks.cbcb.umd.edu/>
- Trinity
<http://trinityrnaseq.sourceforge.net/>
- IGV for Visualization
<http://www.broadinstitute.org/igv/>
- GMAP
<http://research-pub.gene.com/gmap/>
- Samtools
<http://samtools.sourceforge.net/>

Papers of Interest

- Next generation transcriptome assembly
 - <http://www.nature.com/nrg/journal/v12/n10/full/nrg3068.html>
- Tuxedo protocol
 - <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3334321/>
- Trinity
 - <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3571712/>
 - <http://www.nature.com/nprot/journal/v8/n8/full/nprot.2013.084.html>