

AI LAB 2023

1.2 – THE DATA

FRANCESCA M. BUFFA



LAB STRUCTURE

- INTRO AND BACKGROUND
- THE DATA
- THE AI-LAB CHALLENGE – PART 1
- PART 1 - SHARING AND DISCUSSION

- UNSUPERVISED LEARNING EXAMPLES
- THE AI-LAB CHALLENGE – PART 2
- PART 2 - SHARING AND DISCUSSION

- SUPERVISED LEARNING EXAMPLES
- THE AI-LAB CHALLENGE – PART 3
- LARGE PROJECTS AND DATABASES

- THE AI-LAB CHALLENGE PARTS 1-3, SHARING AND DISCUSSION
- DATA INTERPRETATION
- DISCUSS AND PREPARE WORKSHOP PRESENTATIONS

THE DATA

The screenshot shows a web page from the MIT Sloan School of Management website. At the top, the MIT logo and the tagline "Smart. Open. Grounded. Inventive." are visible, along with a call-to-action button "Read our Ideas Made to Matter." On the right side of the header are navigation icons for a menu, search, and social media. Below the header, the text "Credit: Jennifer Tapias Derch" is displayed. The main content area features a large, bold title: "Why it's time for 'data-centric artificial intelligence'". Below the title, the author is listed as "by Sara Brown | Jun 7, 2022". A sidebar on the left contains a "Why It Matters" section with a red icon and the text: "Machine learning pioneer Andrew Ng argues that focusing on the quality of data fueling AI systems will help unlock its full power." At the bottom, there is a "Share" button with a red icon.

IDEAS MADE TO MATTER | ARTIFICIAL INTELLIGENCE

Why it's time for 'data-centric artificial intelligence'

by [Sara Brown](#) | Jun 7, 2022

Why It Matters

Machine learning pioneer Andrew Ng argues that focusing on the quality of data fueling AI systems will help unlock its full power.

Share

Andrew Ng: computer scientist focusing on artificial intelligence, among others cofounder and head of Google Brain.

WHAT IS BIOMEDICAL RAW DATA?

- SEQUENCING DATA
- MICROARRAY DATA
- PROTEIN STRUCTURE DATA
- MORPHOLOGICAL DATA
- IMAGES FROM MICROSCOPES
- CLINICAL IMAGING
- CLINICAL DATA
- ECOLOGICAL DATA
- BIOGEOGRAPHICAL DATA
- DEMOGRAPHIC DATA
-

Oversight:



Funding:



Participants



Biorepository



Sequencing



Fire wall

Data



Clinical Data

- Identifiable clinical data
- Longitudinal
- Linked to genomic data

Existing Clinical Data

Cancer & RD registries, HES, Mortality data, etc



Research Data

- Pseudonymised
- GeCIP and industry partners work within data centre

Data and Analysis

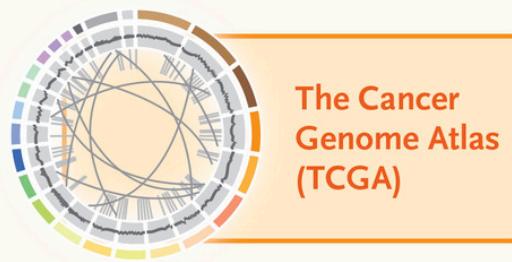
- Improvement
- Annotation & QC
 - Scientists/SMEs
 - Product comparison

Clinicians & Academics

Training
Health Education England

Industry

Genomic and Clinical Data Sources



**The Cancer
Genome Atlas
(TCGA)**



TARGET
(Therapeutically Applicable
Research to Generate
Effective Treatments)



**International
Cancer Genome
Consortium**



**NCI clinical
trials**



User-submitted studies

Genomic Data Commons (GDC)

1. Import and standardize genomic and clinical data from legacy and current NCI programs
2. Harmonize mapping of sequence data to the most current genome and transcriptome build
3. Implement state-of-the-art methods for derived data:
 - Mutation calls
 - Copy number
 - Structural variants
 - Digital gene expression
4. Maintain data security and manage authorized access
5. Provide data for browsing, download, or analysis on a colocalized computer cluster
6. Open GDC for upload of new cancer genomic data from researchers worldwide for comparison with existing data and sharing



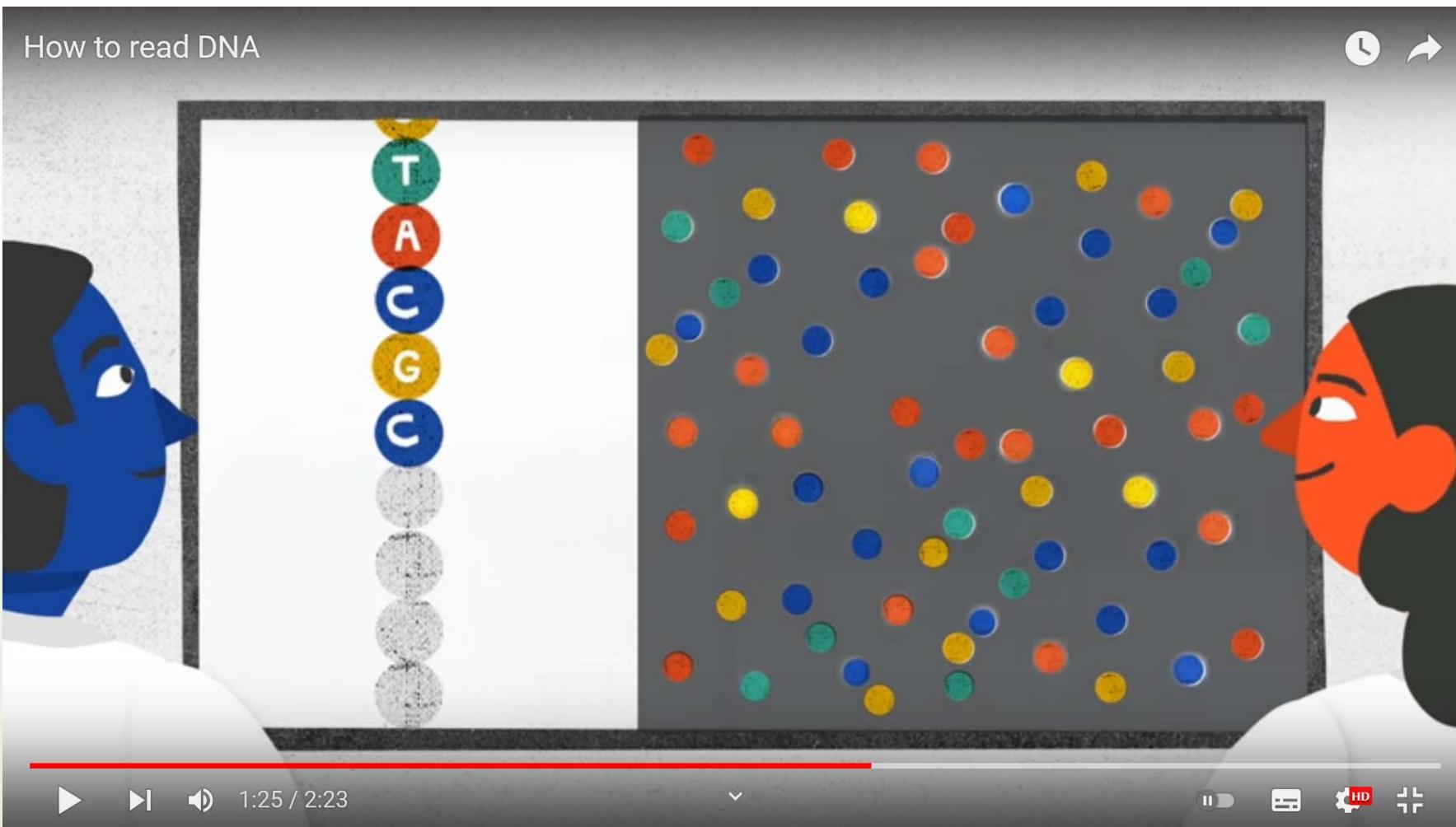
GDC Use Cases

**Identify low-frequency
cancer drivers**

**Define genomic
determinants of
treatment response**

**Compose clinical trial
cohorts sharing targetable
genetic lesions**

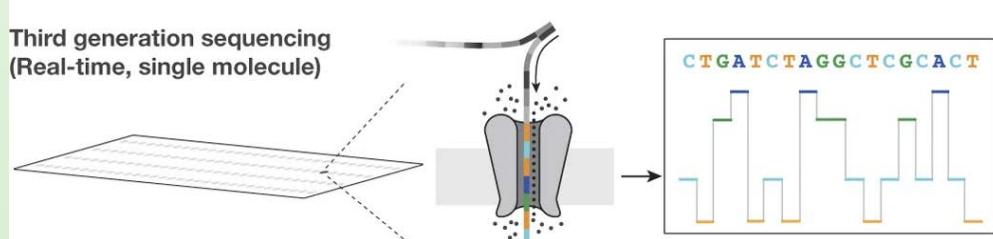
DNA: WE CAN READ IT AND WRITE IT



DNA SEQUENCING TECHNOLOGIES

[1977]

<https://www.youtube.com/watch?v=ONGdehkB8jU>



<https://www.nature.com/articles/nature24286/>

Second generation sequencing (massively parallel)

1 Genomic DNA



2 Fragmented DNA

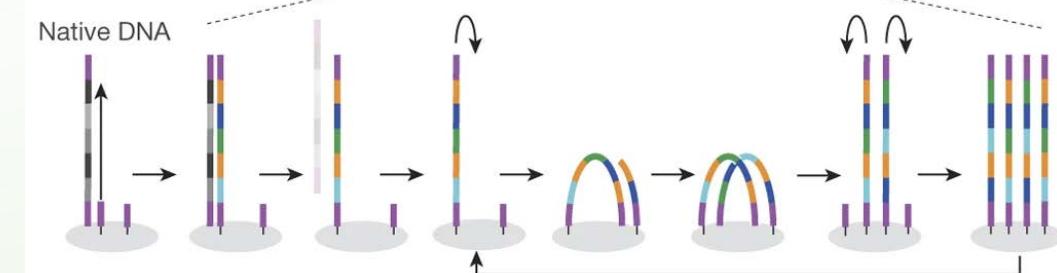


3 Adaptor ligation

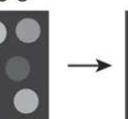
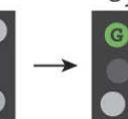
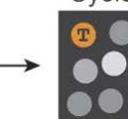


4 Amplification

Native DNA



5 Detection



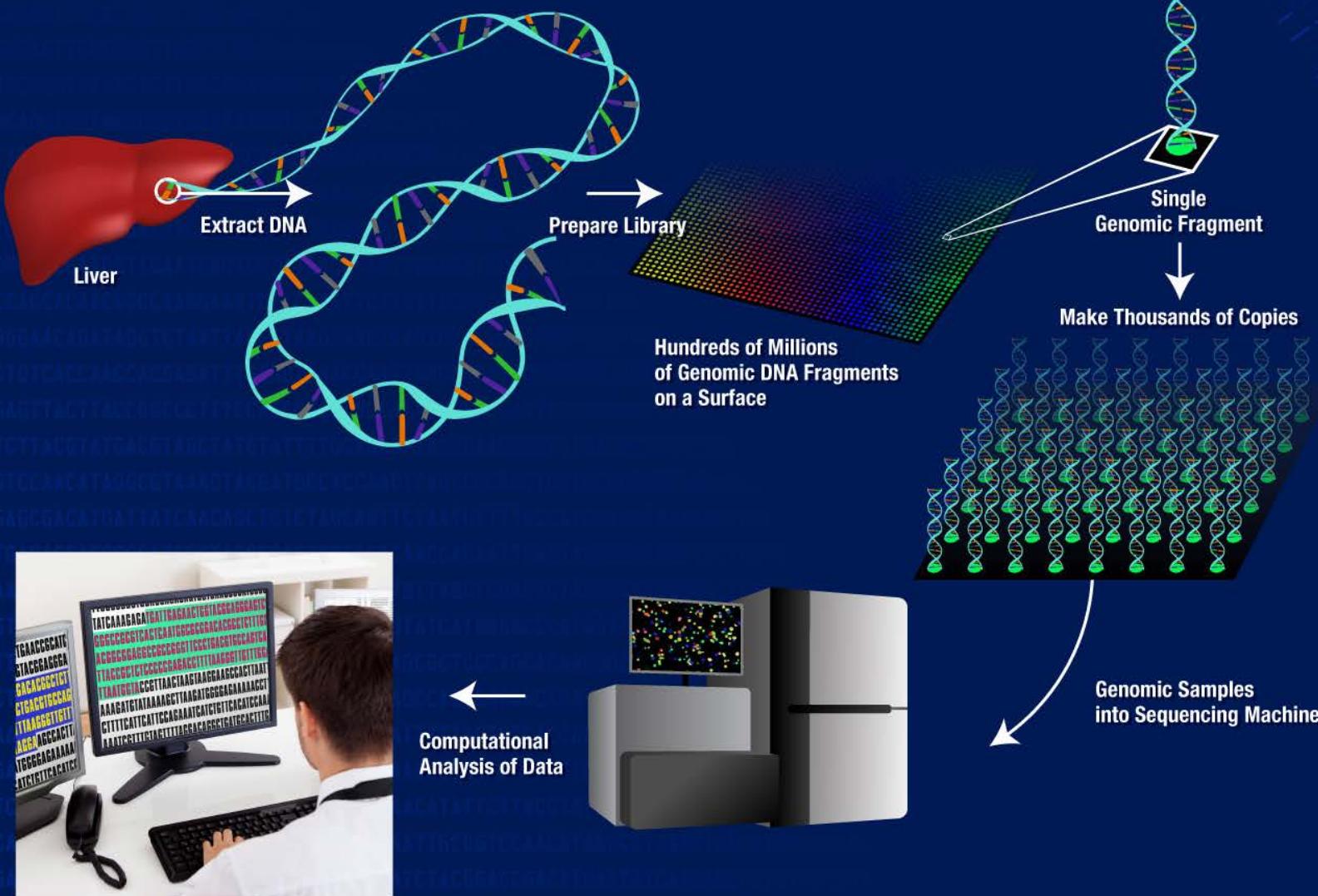
3' ... G A C T A G A T C C G A G C G T G A ... 5'

5' ... C T G A ...

<https://www.youtube.com/watch?v=v10bUR2aL5g>

Dna Sequencing

NHGRI FACT SHEETS
genome.gov



DNA/RNA STRING REPRESENTATION

Sequence: GTTCTCTTCTTCCCTAGCGGATAACAGAGTTGGCCCCACTGTCCCCTTCA

[1] + Sequence Length: 50

[2] + Nucleotide Frequency: {'A': 8, 'C': 18, 'G': 9, 'T': 15}

[3] + DNA/RNA Transcription: GUUCUCUUUCUUCCCUAGCGGAUAACAGAUUCGGCCCACUGUCCCCUUCAG

[4] + DNA String + Complement + Reverse Complement:

5' GTTCTCTTCTTCCCTAGCGGATAACAGAGTTGGCCCCACTGTCCCCTTCA 3'

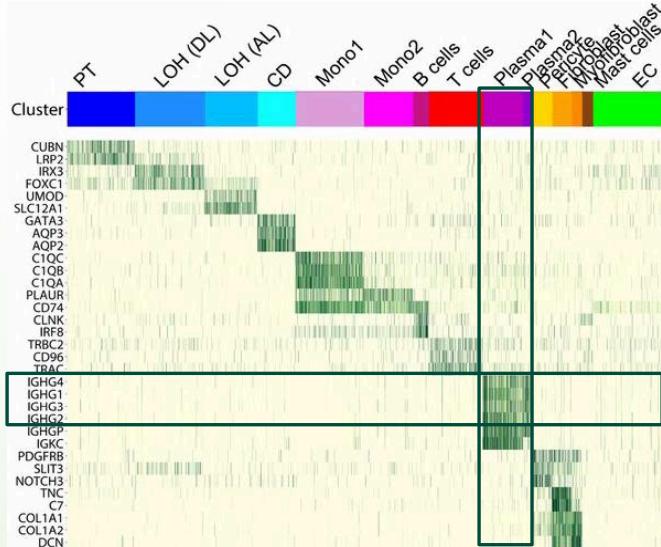
 |||||||||||||||||||||||||||||||||||||||||||

3' CAAGAGAAGAAGGGATCGCCTATTGTCTAAGCCGGTGACAGGGGAAGTC 5' [Complement]

5' CTGAAGGGGACAGTGGGCCGAATCTGTTATCCGCTAGGGAAAGAAGAGAAC 3' [Rev. Complement]

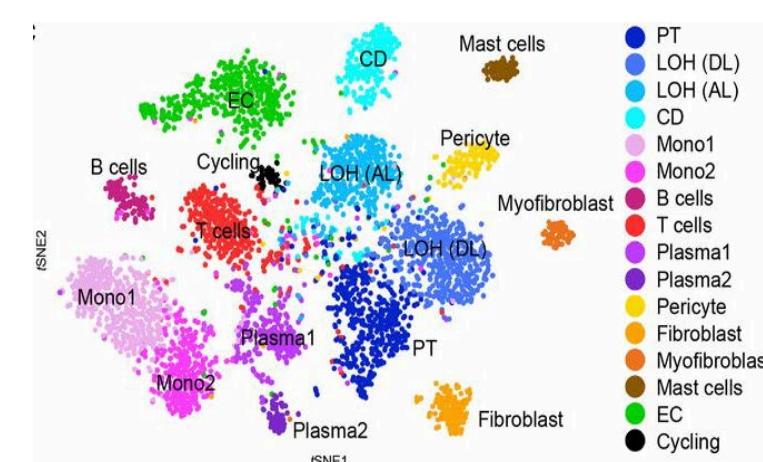
We can sequence the DNA and RNA of single cells

Identifying cell-type marker genes

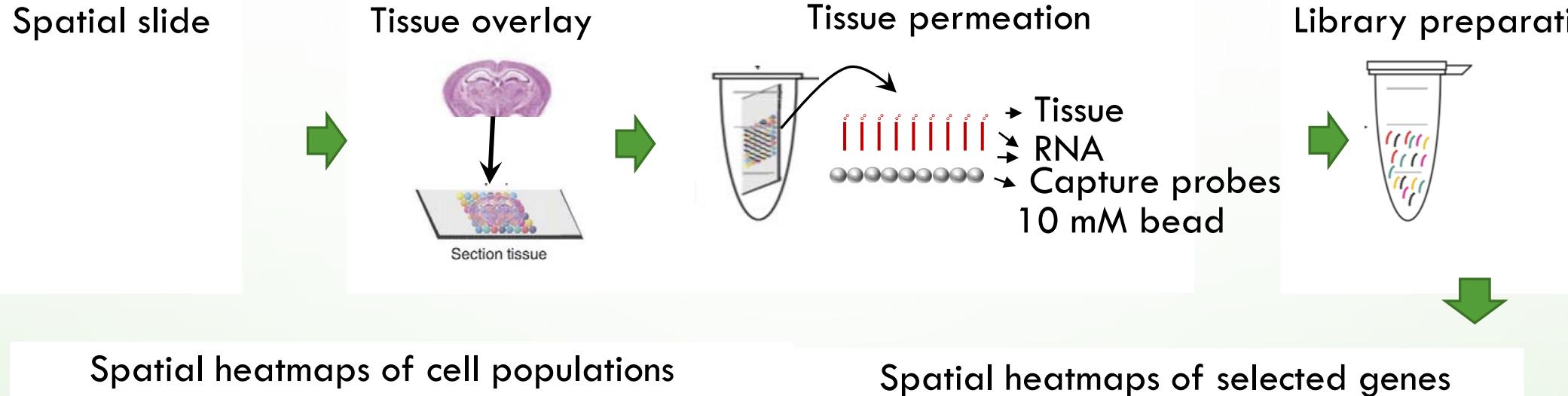


Uncovering tissue dynamics

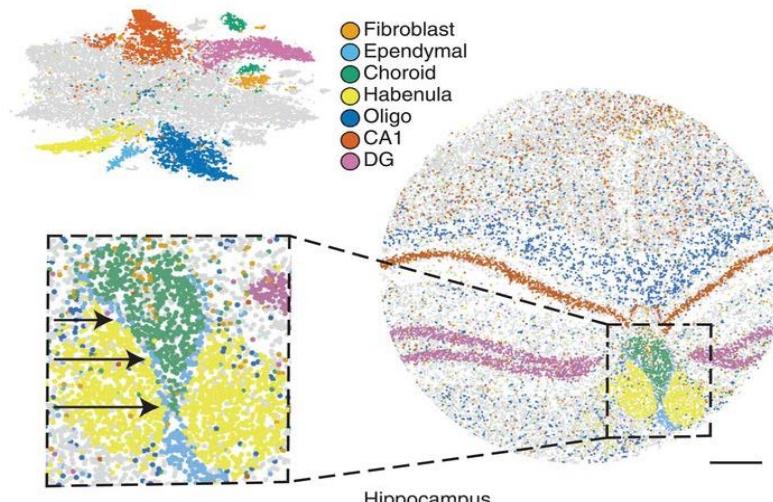
Discovering sample heterogeneity



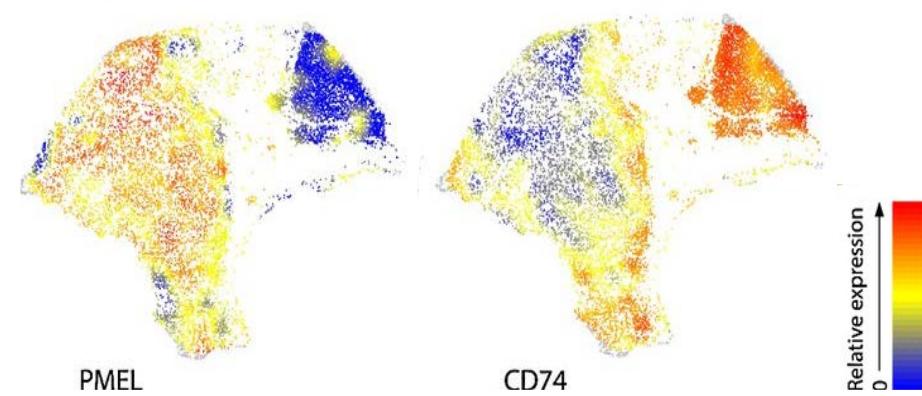
We can acquire a spatial image of single-cell expression



Spatial heatmaps of cell populations

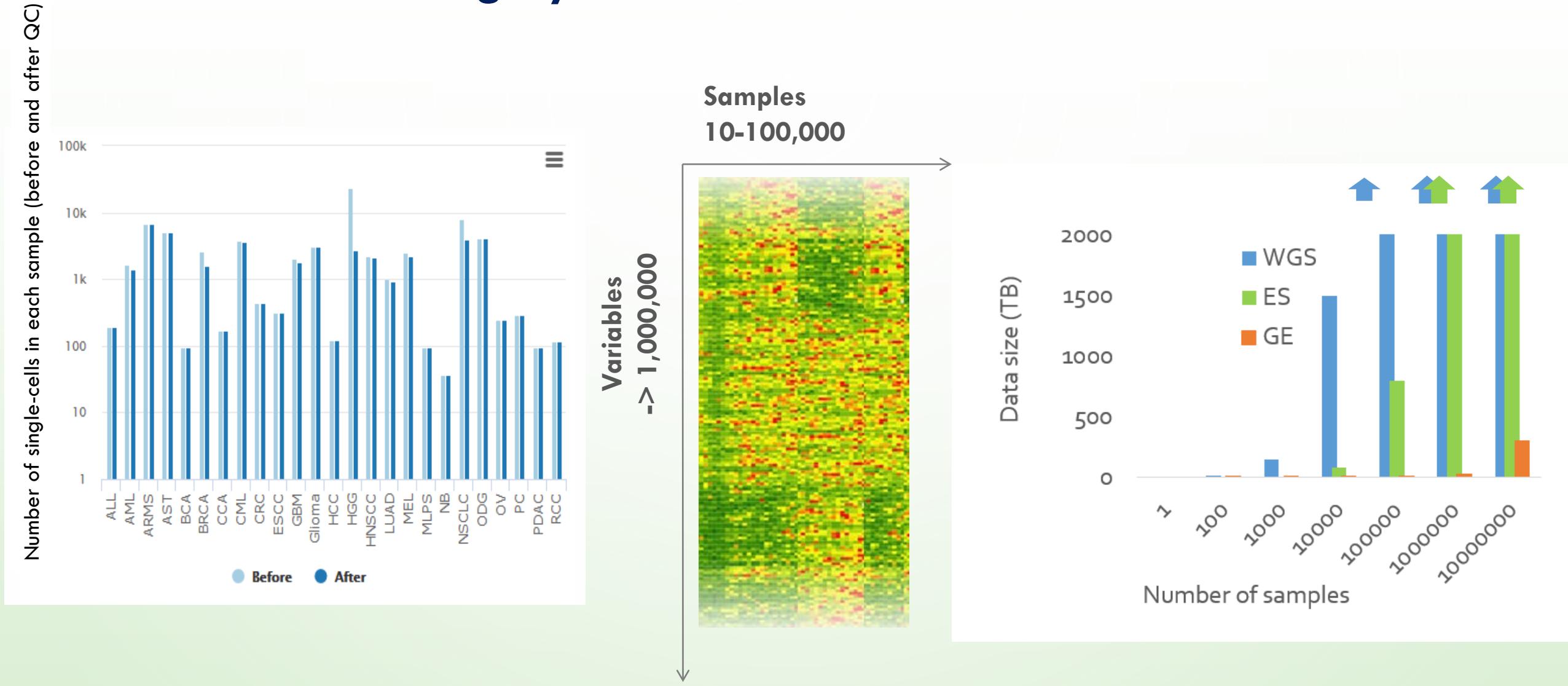


Spatial heatmaps of selected genes



Images from: Vickovic et al (2019) Nature Methods, Thrane (2018) Cancer Res and Rodrigues et al (2019) Science

Highly multidimensional datasets



DNA SEQUENCING DATA



Base Calling
(Vendor tool, non-text)

FASTQ files
Sequencing Reads

A green arrow points from the text "Base Calling (Vendor tool, non-text)" to the text "FASTQ files Sequencing Reads".

FASTQ

De facto standard for storing high-throughput sequencing data

- Text-based file containing raw + quality
- Can contain millions of entries (size ~GBs)

```
@ML-P2-14:9:000H003HG:1:11102:17290:1073 1:N:0:TCCTGAGC+GCGATCTA
TTTGGTAAACAGCATGAATTATTCTAGCCACTAAAACCTATGAACATCTTGTGAAGGTTTCAGATAGAGCCTGAAGTACACAGAGAACATTCTTAAAAAA
+
AAAAAAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE<EEEEEE
```

- 1) Begins with @, followed by sequence identifier and optional additional information/description
- 2) The sequence: the base calls A, C, T, G and N
- 3) A separator, which is simply a plus (+) sign.
- 4) The base call quality scores. It represents the probability of an error in base call. These are Phred +33 encoded, using ASCII characters to represent the numerical quality scores.

DNA SEQUENCING DATA



Base Calling
(Vendor tool, non-text)



FASTQ files
Sequencing Reads



Alignment or assembly

SAM/BAM files:
Aligned sequencing reads
.sam: uncompressed text file
.bam: compressed text file

ALIGNMENT: SEQUENCE MATCHING PROBLEM

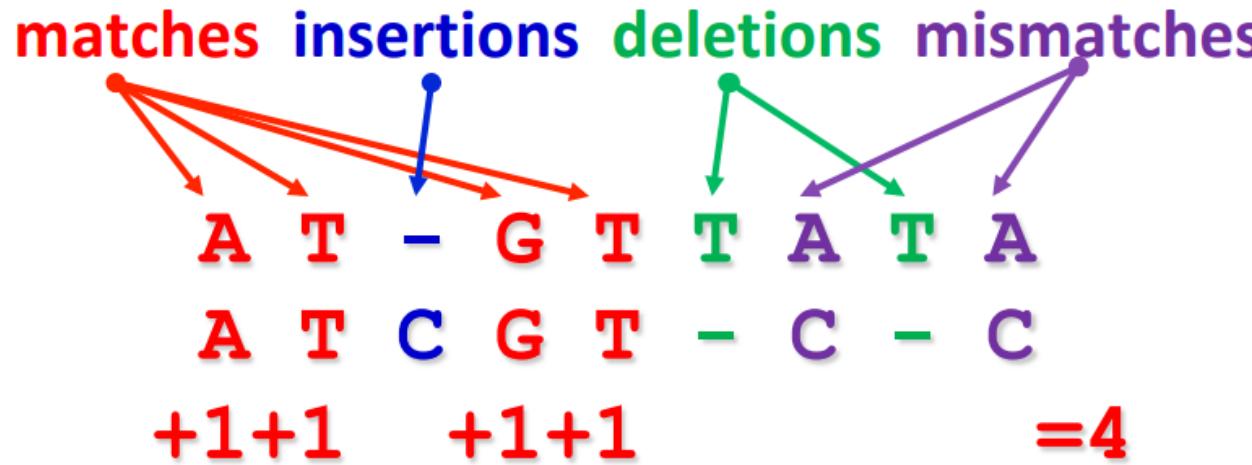
```
ATGTTATA
ATCGTCCC
```

Finding regions of similarities and dissimilarities between sequences, and infer a measure of relatedness, is vital for phylogenetic research:

Comparison of genetic material across organisms allows for example to:

- Infer functional relationships
- Infer evolutionary relationships

ALIGNMENT: SEQUENCE MATCHING PROBLEM



LONGEST COMMON SEQUENCE

Alignment of two sequences is a two-row matrix:

1st row: symbols of the 1st sequence (in order) interspersed by “-”

2nd row: symbols of the 2nd sequence (in order) interspersed by “-”

LONGEST COMMON SUBSEQUENCE

A	T	-	G	T	T	A	T	A
A	T	C	G	T	-	C	-	C

Matches in alignment of two sequences (ATGT) form their Common Subsequence

Longest Common Subsequence Problem:

Find a longest common subsequence of two strings.

- Input: Two strings.
- Output: A longest common subsequence of these strings.

For arbitrary number of input sequences, the problem is NP-hard.

When the number of sequences is constant, the problem is solvable in polynomial time with a number of approaches

SEQUENCE ALIGNMENT MAP - SAM FILE

TAB-delimited text format for storing alignment sequences against a reference

```
@HD VN:1.0 SO:coordinate
@SQ SN:chr20 LN:64444167
@PG ID:TopHat VN:2.0.14 CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5 --read-realign-edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat/index/chr20 /data/user446/mapping_tophat/L6_18_GTGAAA_L007_R1_001.fastq
HWI-ST1145:74:C101DACXX:7:1102:4284:73714 16 chr20 190930 3 100M * 0 0
    CCGTGTAAAGGTGGATCGGGTCACCTTCCCAGCTAGGCTTAGGGATTCTAGTTGCCCTAGGAAATCCAGCTAGTCCTGTCAGTCCCCCTCT
C BBDCDDCCDDDCDDDDCDC?DDDDDDDDDDDDCDCDDDDDDDDCCCCEDDC?DDDDDDDDDDDDDDDDDBDHFFFFDC@_
AS:i:-15 XM:i:3 X0:i:0 XG:i:0 MD:Z:55C20C13A9 NM:i:3 NH:i:2 CC:Z:= CP:i:55352714 HI:i:0
HWI-ST1145:74:C101DACXX:7:1114:2759:41961 16 chr20 193953 50 100M * 0 0
    TGCTGGATCATCTGGTTAGGGCTCTGACTCAGAGGACCTCGTCCCTGGGGCAGTGGACCTCCAGTGATTCCCTGACATAAGGGGATGGACGA
G DCDDDDDEDDDDDDCDCDDDDDDCCDDDCDDDEEC>DFFEJJJJJIGJJJIHGBHHGJIJJJJJJGJJJIJJJJJIHJJJJJJHHHHFFFFCCC
AS:i:-16 XM:i:3 X0:i:0 XG:i:0 MD:Z:60G16T18T3 NM:i:3 NH:i:1
HWI-ST1145:74:C101DACXX:7:1204:14760:4030 16 chr20 270877 50 100M * 0 0
    GGCTTATTGGTAAAAAAGGAATAGCAGATTAAATCAGAAATTCCACCTGGCCCAGCAGCACCAACCAGAAAGGAAGAGACAGGAAAAAACCA
C DDDDDDDDDCDCDDDDDDDEEEEEEFFFEFFEGHHHFGDJJIHJJJIJJIIIGGFJJJIHIIIIJJJJJIGHFAHGFHJHFGGHFFFDD@BB
AS:i:-11 XM:i:2 X0:i:0 XG:i:0 MD:Z:0A85G13 NM:i:2 NH:i:1
HWI-ST1145:74:C101DACXX:7:1210:11167:8699 0 chr20 271218 50 50M4700N50M * 0
    0 GTGGCTTCCACAGGAATGTTGAGGATGACATCCATGTCTGGGTGCACTTGGGTCTCCGAAGCAGAACATCCTCAAATATGACCTCTCG
accepted_hits.sam
```

Binary Alignment/Map (BAM) file (*.bam) is a compressed binary SAM file (smaller size + faster access)

SEQUENCE ALIGNMENT MAP - SAM FILE

TAB-delimited text format for storing alignment sequences against a reference

```
@HD VN:1.0 SD:coordinate → File-level metadata. Optional.  
@SQ SN:chr20 LN:64444167  
@PG ID:TopHat VN:2.0.14 CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5 --read-realign-edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat/index/chr20 /data/user446/mapping_tophat/L6_18_GTGAAA_L007_R1_001.fastq  
HWI-ST1145:74:C101DACXX:7:1102:4284:73714 16 chr20 190930 3 100M * 0 0  
CCGTGTTAAAGGTGGATCGGGTCACCTTCCCAGCTAGGCTTAGGGATTCTTAGTTGCCCTAGGAAATCCAGCTAGTCCTGTCAGTCCCCCTCT  
C BBDCDDCCDDDCDDDDCDCC?DDDDDDDDDDDDCCDCDDDDDDDDCCCCEDDC?DDDDDDDDDDDDDDDDDBDHFFFFDC@  
AS:i:-15 XM:i:3 X0:i:0 XG:i:0 MD:Z:55C20C13A9 NM:i:3 NH:i:2 CC:Z:= CP:i:55352714 HI:i:0  
HWI-ST1145:74:C101DACXX:7:1114:2759:41961 16 chr20 193953 50 100M * 0 0  
TGCTGGATCATCTGGTTAGGGCTCTGACTCAGAGGACCTCGTCCCTGGGGCAGTGGACCTCCAGTGA  
G TCCCTGACATAAGGGGCATGGACGA  
G DCDDDDDEDDDDDDCDDDDDDCCDDDCDDDEEC>DFFEJJJJJIGJJJIHGBHHGJIJJJJGJJJIJJJJJIHJJJJJHHHHFFFFCCC  
AS:i:-16 XM:i:3 X0:i:0 XG:i:0 MD:Z:60G16T18T3 NM:i:3 NH:i:1  
HWI-ST1145:74:C101DACXX:7:1204:14760:4030 16 chr20 270877 50 100M * 0 0  
GGCTTATTGGTAAAAAAGGAATAGCAGATTAAATCAGAAATTCCACCTGGCCCAGCAGCACCAACCAGAAAGGAAGAGACAGGAAAAACCA  
C DDDDDDDDDCDDDDDDDDDEEEEEEEFFFEFFEGHHHFGDJJIHJJJIJJIIIGGFJJJIHIIIIJJJJJIGHFAHGFHJHFGGHFFFDD@BB  
AS:i:-11 XM:i:2 X0:i:0 XG:i:0 MD:Z:0A85G13 NM:i:2 NH:i:1  
HWI-ST1145:74:C101DACXX:7:1210:11167:8699 0 chr20 271218 50 50M4700N50M * 0  
0 GTGGCTTCCACAGGAATGTTGAGGATGACATCCATGTCTGGGTGCACTTGGGTCTCCGAAGCAGAACATCCTCAAATATGACCTCTCG  
accepted_hits.sam
```

Binary Alignment/Map (BAM) file (*.bam) is a compressed binary SAM file (smaller size + faster access)

SEQUENCE ALIGNMENT MAP - SAM FILE

TAB-delimited text format for storing alignment sequences against a reference

```
@HD VN:1.0 SO:coordinate → Reference sequence dictionary
@SQ SN:chr20 LN:64444167
@PG ID:TopHat VN:2.0.14 CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5 --read-realign-edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat/index/chr20 /data/user446/mapping_tophat/L6_18_GTGAAA_L007_R1_001.fastq
HWI-ST1145:74:C101DACXX:7:1102:4284:73714 16 chr20 190930 3 100M * 0 0
    CCGTGTAAAGGTGGATCGGGTCACCTTCCCAGCTAGGCTTAGGGATTCTAGTTGCCCTAGGAAATCCAGCTAGTCCTGTCAGTCCCCCTCT
C BBDCDDCCDDDCDDDDCDC?DDDDDDDDDDDDCDCDDDDDDDDCCCEDDC?DDDDDDDDDDDDDDDDDBDHFFFFDC@_
AS:i:-15 XM:i:3 X0:i:0 XG:i:0 MD:Z:55C20C13A9 NM:i:3 NH:i:2 CC:Z:= CP:i:55352714 HI:i:0
HWI-ST1145:74:C101DACXX:7:1114:2759:41961 16 chr20 193953 50 100M * 0 0
    TGCTGGATCATCTGGTTAGGGCTCTGACTCAGAGGACCTCGTCCCCTGGGGCAGTGGACCTCCAGTGATTCCCTGACATAAGGGGCATGGACGA
G DCDDDDDEDDDDDDCDCDDDDDDCCCDDDCDDDEEC>DFFEJJJJJIGJJJIHGBHHGJIJJJJJJGJJJIJJJJJIHJJJJJJHHHHFFFFCCC
AS:i:-16 XM:i:3 X0:i:0 XG:i:0 MD:Z:60G16T18T3 NM:i:3 NH:i:1
HWI-ST1145:74:C101DACXX:7:1204:14760:4030 16 chr20 270877 50 100M * 0 0
    GGCTTATTGGTAAAAAAGGAATAGCAGATTAAATCAGAAATTCCACCTGGCCCAGCAGCACCAACCAGAAAGAAGGGAAAGACAGGAAAAAACCA
C DDDDDDDDDCDCDDDDDDDEEEEEEEFFFEFGHHHFGDJJIHJJJIJJJJIIIGGFJJJIHIIIIJJJJJJIGHFAHGFHJHFGGHFFFDD@BB
AS:i:-11 XM:i:2 X0:i:0 XG:i:0 MD:Z:0A85G13 NM:i:2 NH:i:1
HWI-ST1145:74:C101DACXX:7:1210:11167:8699 0 chr20 271218 50 50M4700N50M * 0
    0 GTGGCTTCCACAGGAATGTTGAGGATGACATCCATGTCTGGGTGCACTTGGGTCTCCGAAGCAGAACATCCTCAAATATGACCTCTCG
accepted_hits.sam
```

Binary Alignment/Map (BAM) file (*.bam) is a compressed binary SAM file (smaller size + faster access)

SEQUENCE ALIGNMENT MAP - SAM FILE

TAB-delimited text format for storing alignment sequences against a reference

```
@HD VN:1.0 SO:coordinate
@SQ SN:chr20 LN:64444167 → Program
@PG ID:TopHat VN:2.0.14 CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5 --read-realign-edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat/index/chr20 /data/user446/mapping_tophat/L6_18_GTGAAA_L007_R1_001.fastq
HWI-ST1145:74:C101DACXX:7:1102:4284:73714 16 chr20 190930 3 100M * 0 0
    CCGTGTAAAGGTGGATCGGGTCACCTTCCCAGCTAGGCTTAGGGATTCTAGTTGCCCTAGGAAATCCAGCTAGTCCTGTCAGTCCCCCTCT
C BBDCDDCCDDDCDDDDCDC?DDDDDDDDDDDDCDCDDDDDDDDCCCEDDC?DDDDDDDDDDDDDDDDDBDHFFFFDC@_
AS:i:-15 XM:i:3 X0:i:0 XG:i:0 MD:Z:55C20C13A9 NM:i:3 NH:i:2 CC:Z:= CP:i:55352714 HI:i:0
HWI-ST1145:74:C101DACXX:7:1114:2759:41961 16 chr20 193953 50 100M * 0 0
    TGCTGGATCATCTGGTTAGGGCTCTGACTCAGAGGACCTCGTCCCCTGGGGCAGTGGACCTCCAGTGATTCCCTGACATAAGGGGATGGACGA
G DCDDDDDEDDDDDDCDCDDDDDDCCCDDDCDDDEEC>DFFEJJJJJIGJJJIHGBHHGJIJJJJJJGJJJIJJJJJIHJJJJJJHHHHFFFFCCC
AS:i:-16 XM:i:3 X0:i:0 XG:i:0 MD:Z:60G16T18T3 NM:i:3 NH:i:1
HWI-ST1145:74:C101DACXX:7:1204:14760:4030 16 chr20 270877 50 100M * 0 0
    GGCTTATTGGTAAAAAAGGAATAGCAGATTAAATCAGAAATTCCACCTGGCCCAGCAGCACCAACCAGAAAGGAAGAGACAGGAAAAAACCA
C DDDDDDDDDCDCDDDDDDDEEEEEEEFFFEFFEGHHHFGDJJIHJJJIJJJJIIIGFJJIHIIIIJJJJJIGHFAHGFHJHFGGHFFFDD@BB
AS:i:-11 XM:i:2 X0:i:0 XG:i:0 MD:Z:0A85G13 NM:i:2 NH:i:1
HWI-ST1145:74:C101DACXX:7:1210:11167:8699 0 chr20 271218 50 50M4700N50M * 0
    0 GTGGCTTCCACAGGAATGTTGAGGATGACATCCATGTCTGGGTGCACTTGGGTCTCCGAAGCAGAACATCCTCAAATATGACCTCTCG
accepted_hits.sam
```

Binary Alignment/Map (BAM) file (*.bam) is a compressed binary SAM file (smaller size + faster access)

SEQUENCE ALIGNMENT MAP - SAM FILE

TAB-delimited text format for storing alignment sequences against a reference

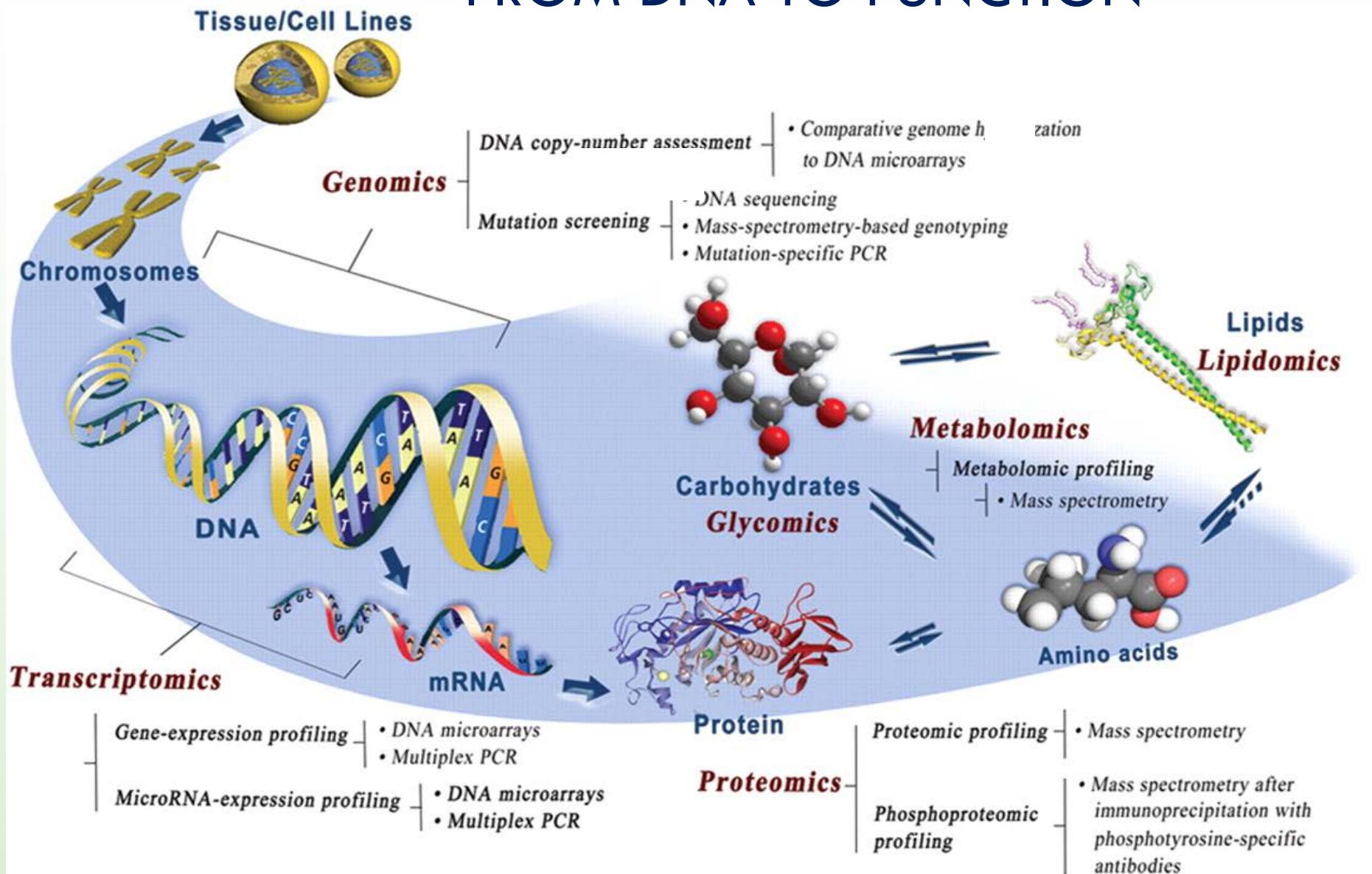
Alignment information										
@HD	VN:1.0	S0:coordinate								
@SQ	SN:chr20	LN:64444167								
@PG	ID:TopHat	VN:2.0.14	CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5 --read-realign-edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat/index/chr20 /data/user446/mapping_tophat/L6 18 GTGAAA L007 R1 001.fastq							
HWI-ST1145:74:C101DACXX:7:1102:4284:73714	16	chr20	190930	3	100M	*	0	0		
CCGTGTTAAAGGTGGATCGGGTCACCTTCCCAGCTAGGCTTAGGGATTCTAGTTGCCCTAGGAAATCCAGCTAGTCCTGTCAGTCCCCCTCTC										
C	BBDCDDCCDDDDCDDDDCDCCCDBC?DDDDDDDDDDDDDDCCDCDDDDDDDDCCCCEDDC?DDDDDDDDDDDDDDDDDDDBDHFFFFDC@									
AS:i:-15	XM:i:3	XO:i:0	XG:i:0	MD:Z:55C20C13A9	NM:i:3	NH:i:2	CC:Z:=	CP:i:55352714	HI:i:0	
HWI-ST1145:74:C101DACXX:7:1114:2759:41961	16	chr20	193953	50	100M	*	0	0		
TGCTGGATCATCTGGTTAGGGCTCTGACTCAGAGGACCTCGTCCCTGGGGCAGTGGACCTCCAGTGATTCCCTGACATAAGGGGCATGGACGA										
G	DCDDDDDEDDDDDDCDDDDDDCCCDDDCDDDDDEEC>DFFEJJJJJIGJJJIHGBHHGJIJJJJJJGJJJIJJJJJIHJJJJJJHHHHFFFFCCC									
AS:i:-16	XM:i:3	XO:i:0	XG:i:0	MD:Z:60G16T18T3	NM:i:3	NH:i:1				
HWI-ST1145:74:C101DACXX:7:1204:14760:4030	16	chr20	270877	50	100M	*	0	0		
GGCTTATTGGTAAAAAAGGAATAGCAGATTAAATCAGAAATTCCACCTGGCCCAGCAGCACCAACCAGAAAGAAGGGAAAGAACAGGAAAAACCA										
C	DDDDDDDDCCDDDDDDDDDEEEEEEEFFFEFFEGHHHFGDJJIHJJIIJJIIIGGFJJIHIIIIJJJJJJGHHFAHGFHJHFGGHFFFDD@BB									
AS:i:-11	XM:i:2	XO:i:0	XG:i:0	MD:Z:0A85G13	NM:i:2	NH:i:1				
HWI-ST1145:74:C101DACXX:7:1210:11167:8699	0	chr20	271218	50	50M4700N50M	*	0			
0	GTGGCTTCCACAGGAATGTTGAGGATGACATCCATGTCTGGGTGCACTTGGGTCTCCGAAGCAGAACATCCTCAAATATGACCTCTCG									

accepted_hits.sam

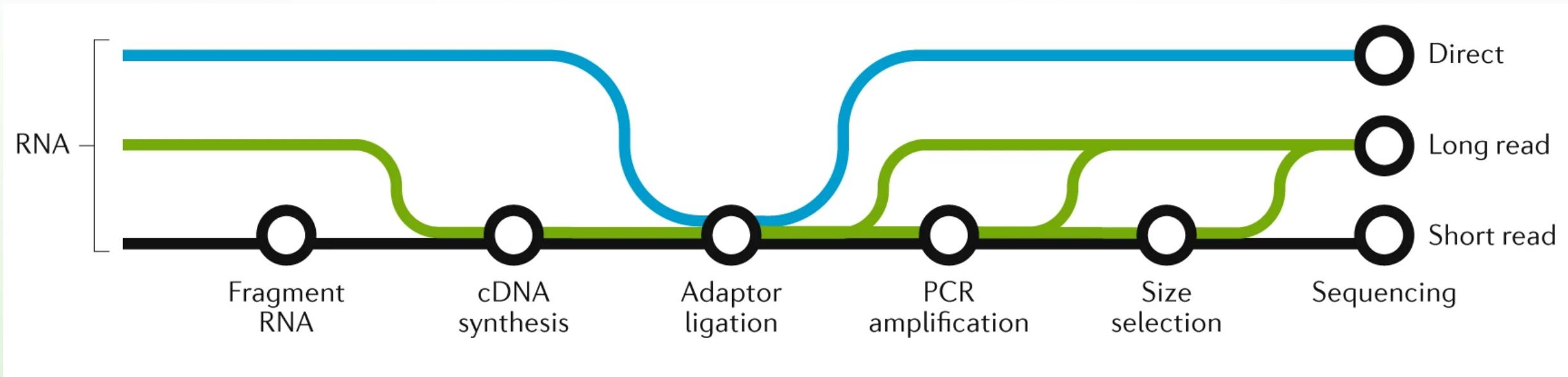
<http://samtools.github.io/hts-specs/SAMv1.pdf>

Binary Alignment/Map (BAM) file (*.bam) is a compressed binary SAM file (smaller size + faster access)

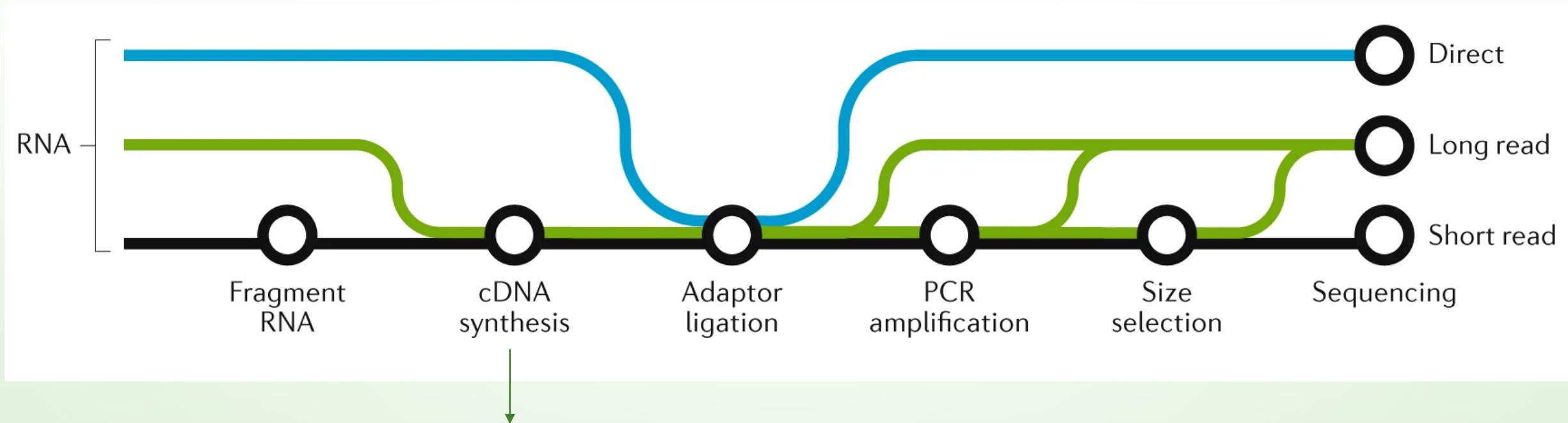
FROM DNA TO FUNCTION



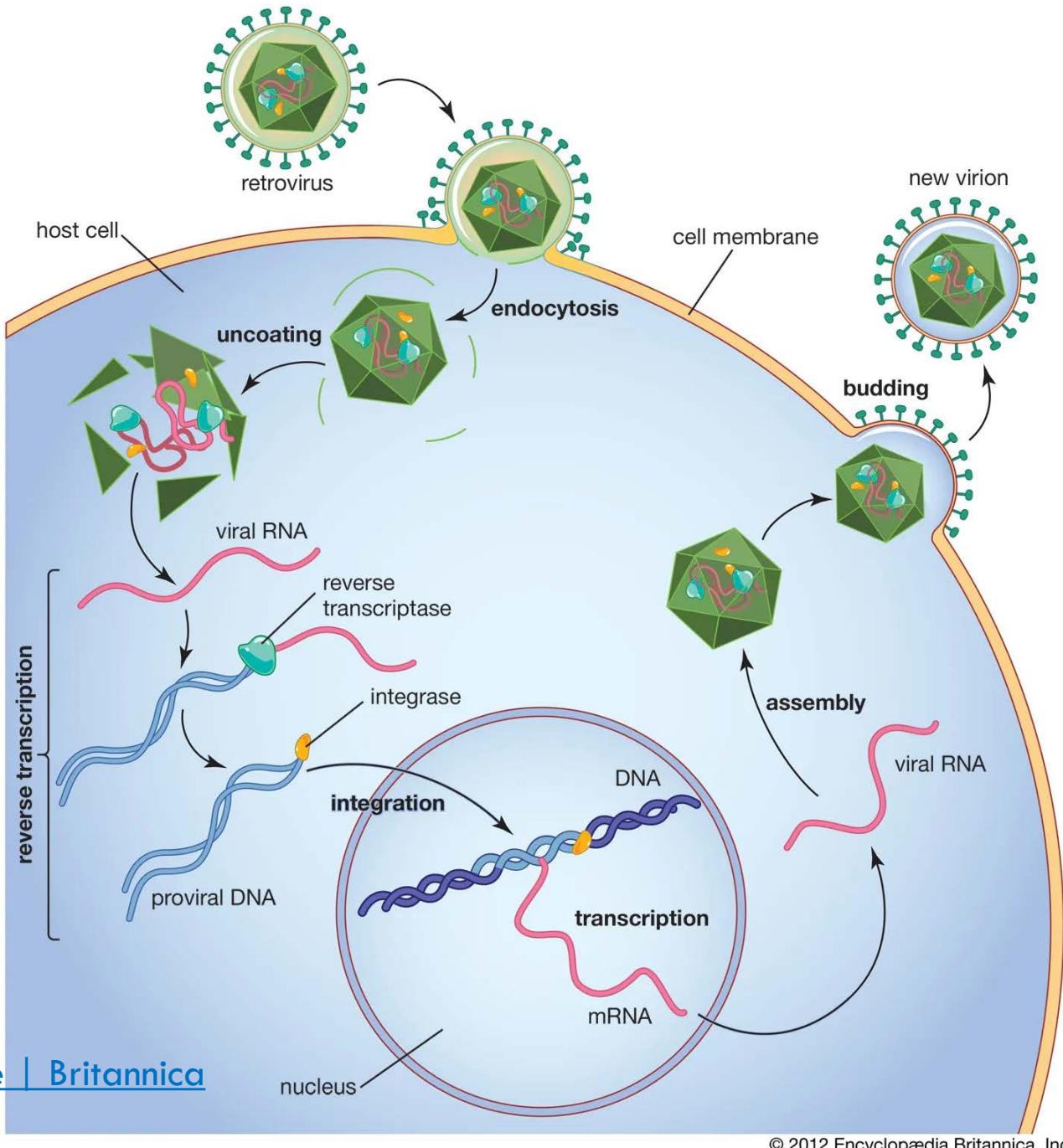
RNA SEQUENCING



RNA SEQUENCING



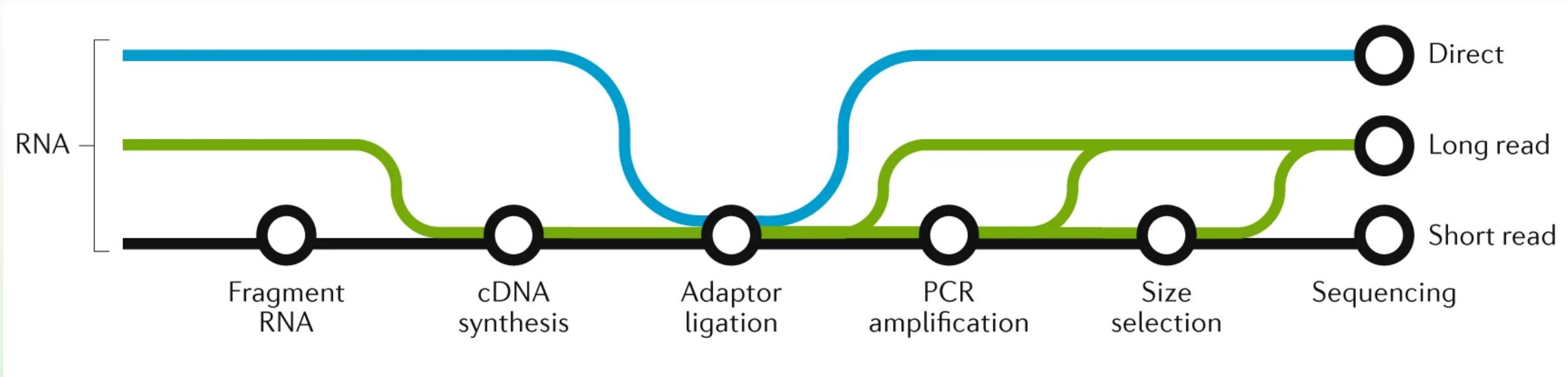
Retrovirus infection and reverse transcription



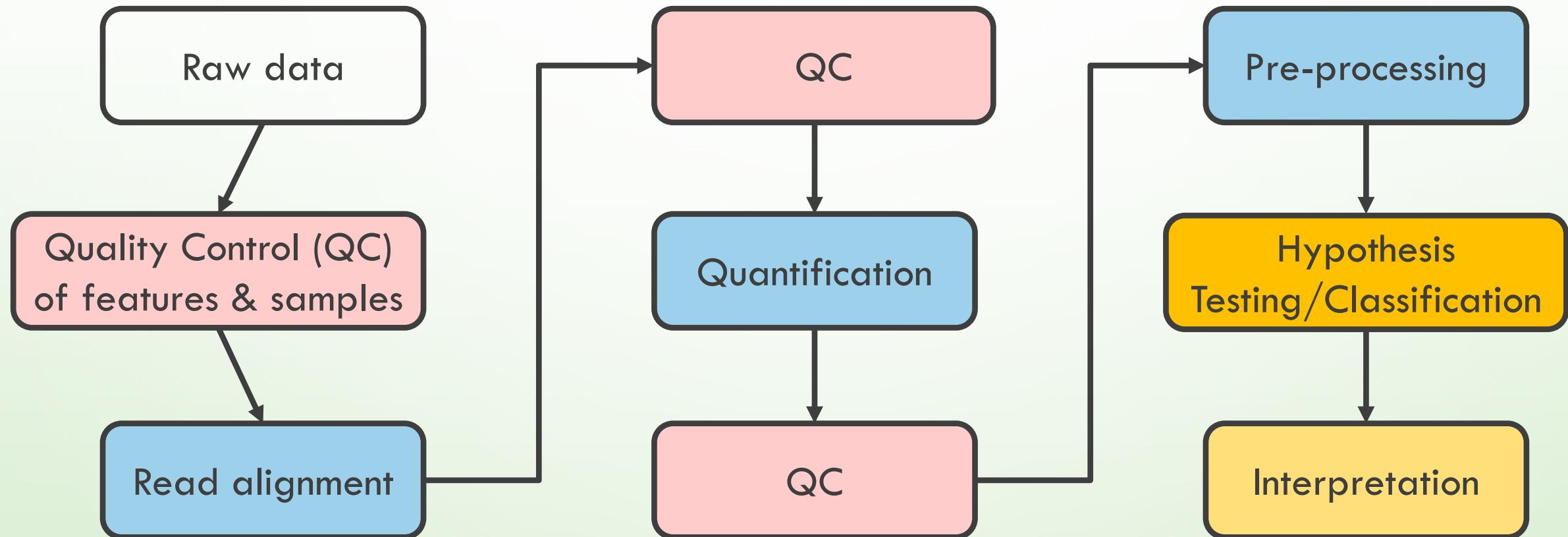
[Encyclopædia Britannica](#)

[reverse transcriptase | enzyme](#) | Britannica

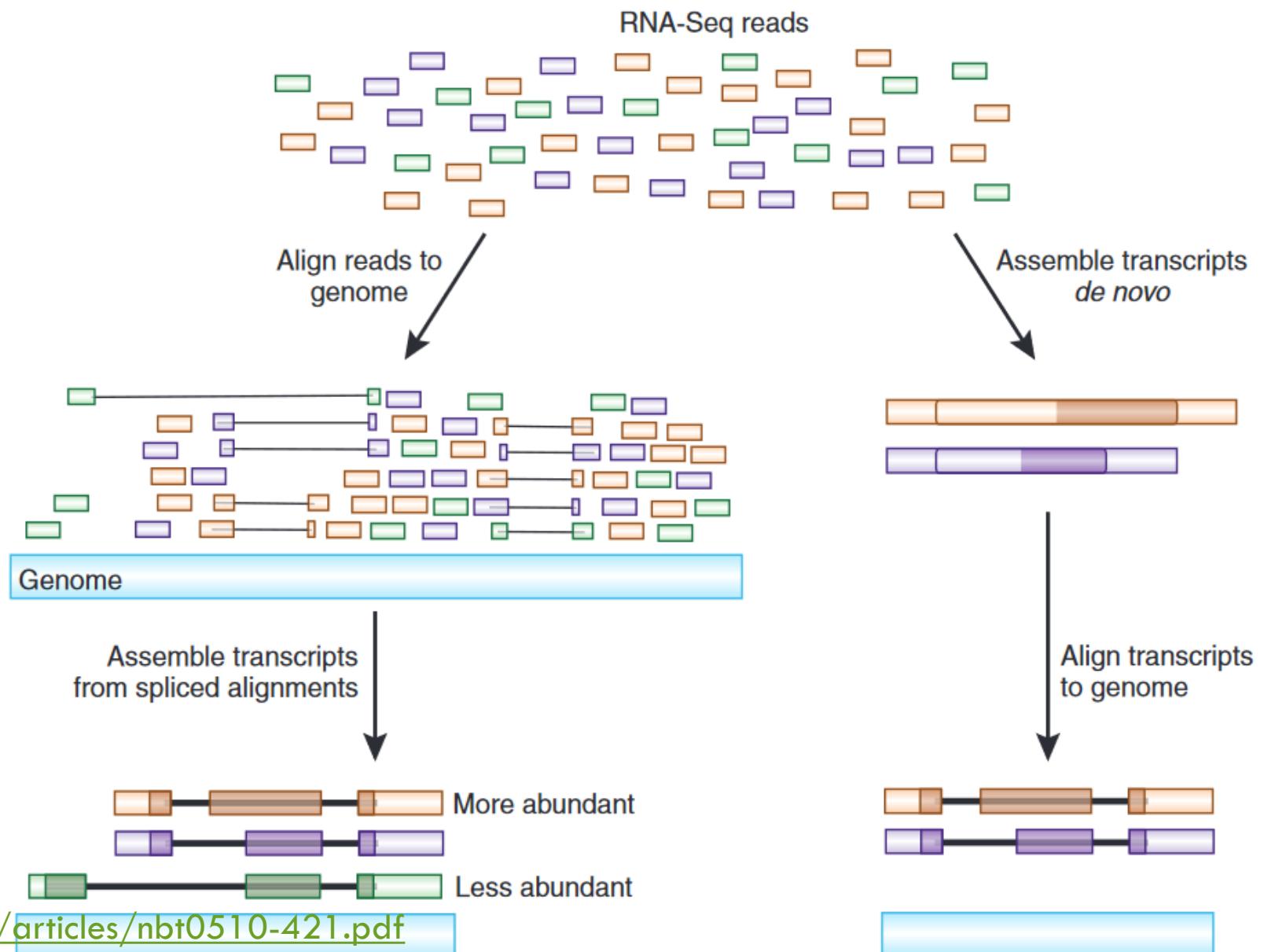
RNA SEQUENCING



RNA-SEQ: FROM LOOKING AT THE DATA TO ANALYSIS



RNA-SEQ ALIGNEMENT



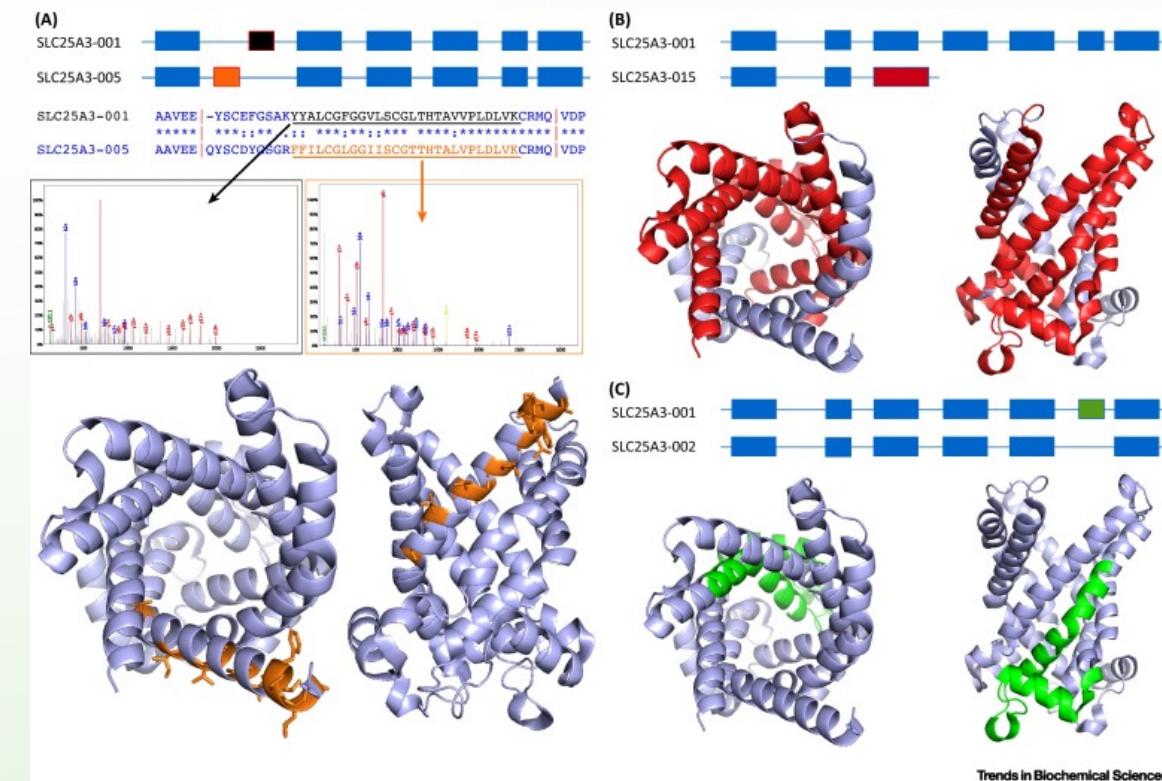
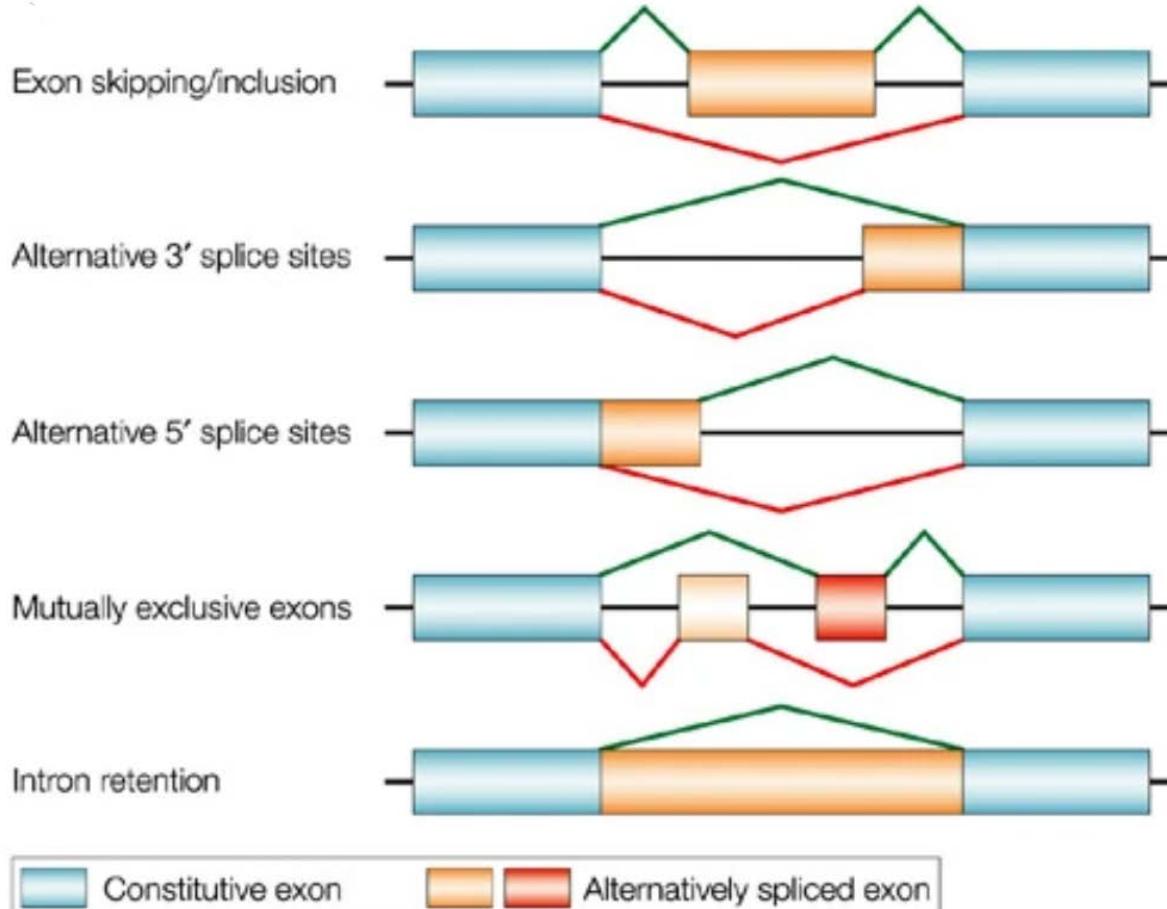
NUMBERS OF GENES IN GENOMES

- Simple assumption: the whole of the genome codes for genes of interest
- If we assume that the number of amino acids in a typical protein is roughly 300 (very simplistic!)
- Then the number of bases needed to code for our typical protein is ~ 1000 (3 base pairs per amino acid)
- Genes contained in a genome estimated as genome size/1000
- For bacterial genomes this works
- For eukaryotic genomes this completely fails!

Organism	# of protein-coding genes	# of genes naïve estimate: (genome size /1000)
HIV 1	9	10
<i>Influenza A virus</i>	10-11	14
Bacteriophage λ	66	49
<i>Epstein Barr virus</i>	80	170
<i>Buchnera sp.</i>	610	640
<i>T. maritima</i>	1,900	1,900
<i>S. aureus</i>	2,700	2,900
<i>V. cholerae</i>	3,900	4,000
<i>B. subtilis</i>	4,400	4,200
<i>E. coli</i>	4,300	4,600
<i>S. cerevisiae</i>	6,600	12,000
<i>C. elegans</i>	20,000	100,000
<i>A. thaliana</i>	27,000	140,000
<i>D. melanogaster</i>	14,000	140,000
<i>F. rubripes</i>	19,000	400,000
<i>Z. mays</i>	33,000	2,300,000
<i>M. musculus</i>	20,000	2,800,000
<i>H. sapiens</i>	21,000	3,200,000

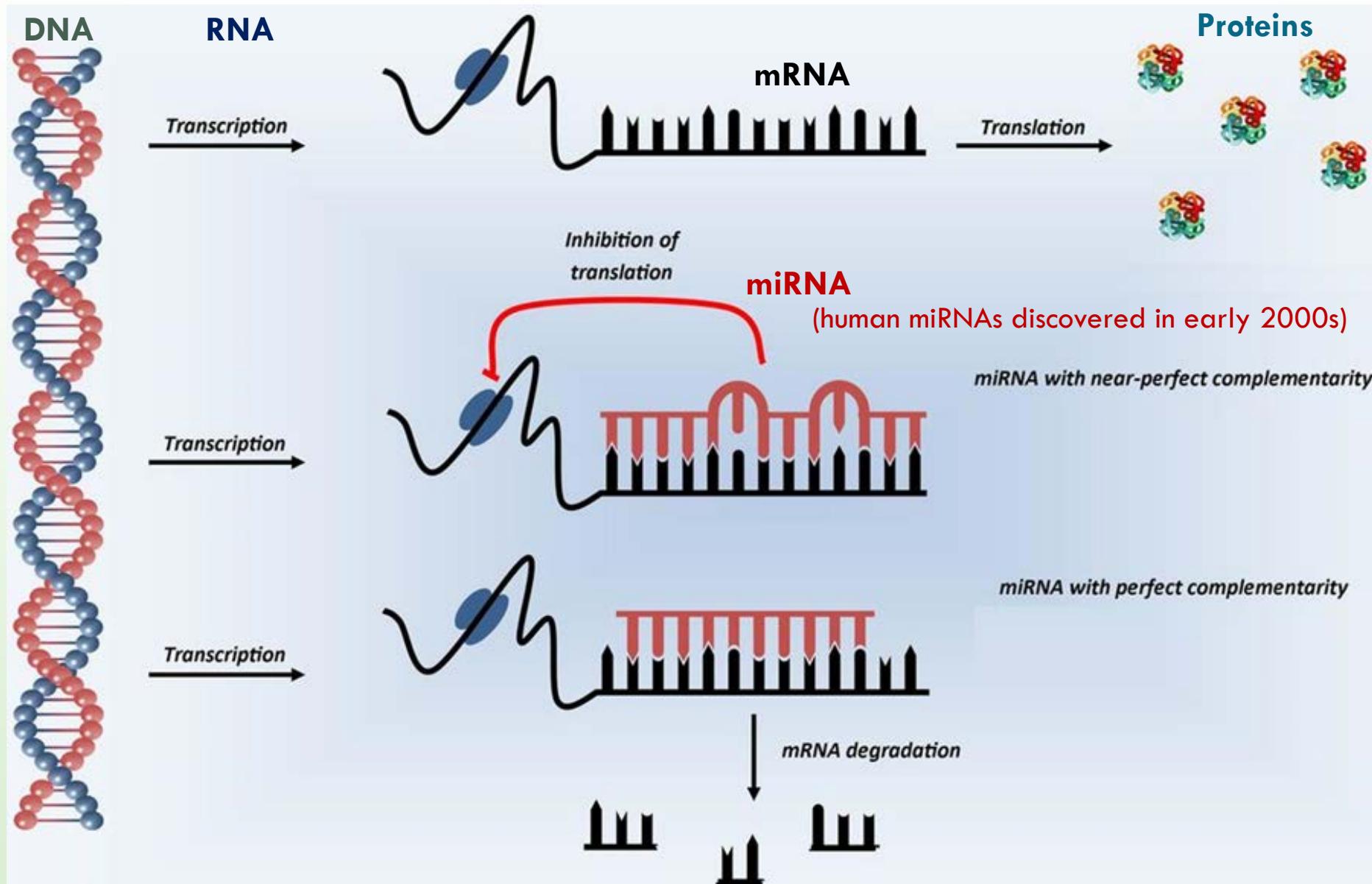
HOW CAN SO FEW GENES ENCODE FOR THE COMPLEXITY OF EUKARYOTES?

From one RNA many possibilities for proteins



HOW CAN SO FEW GENES ENCODE FOR THE COMPLEXITY OF EUKARYOTES?

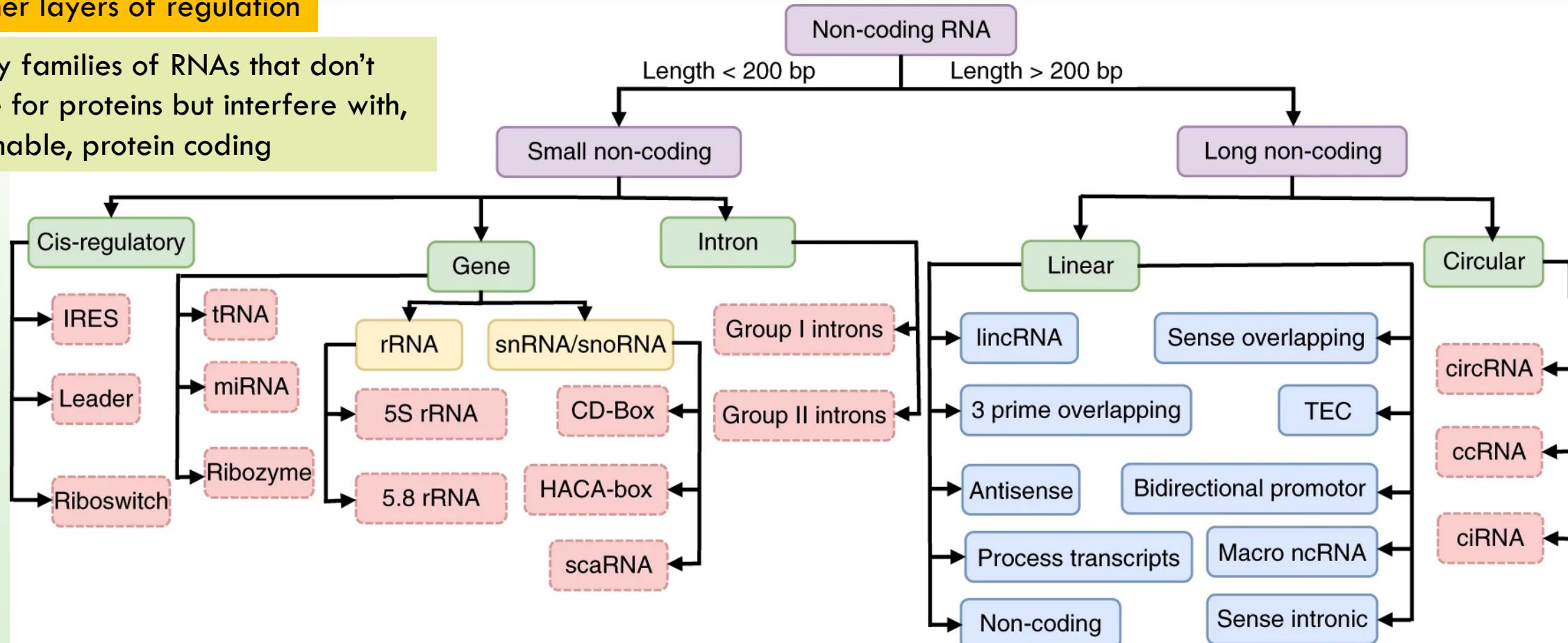
Further layers of regulation



HOW CAN SO FEW GENES ENCODE FOR THE COMPLEXITY OF EUKARYOTES?

Further layers of regulation

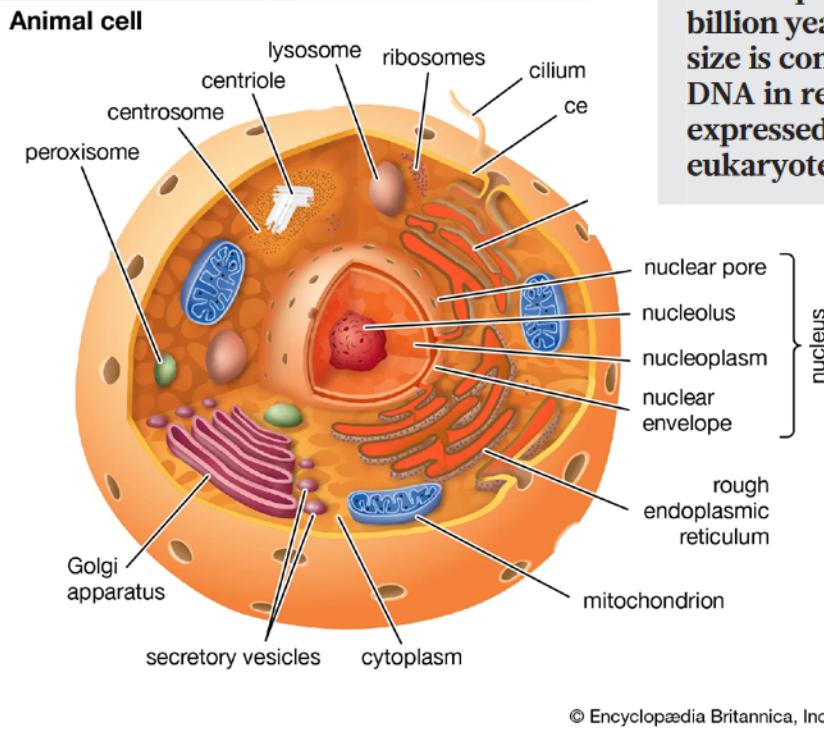
Many families of RNAs that don't code for proteins but interfere with, or enable, protein coding



HOW CAN SO FEW GENES ENCODE FOR THE COMPLEXITY OF EUKARYOTES?

The energetics of genome complexity

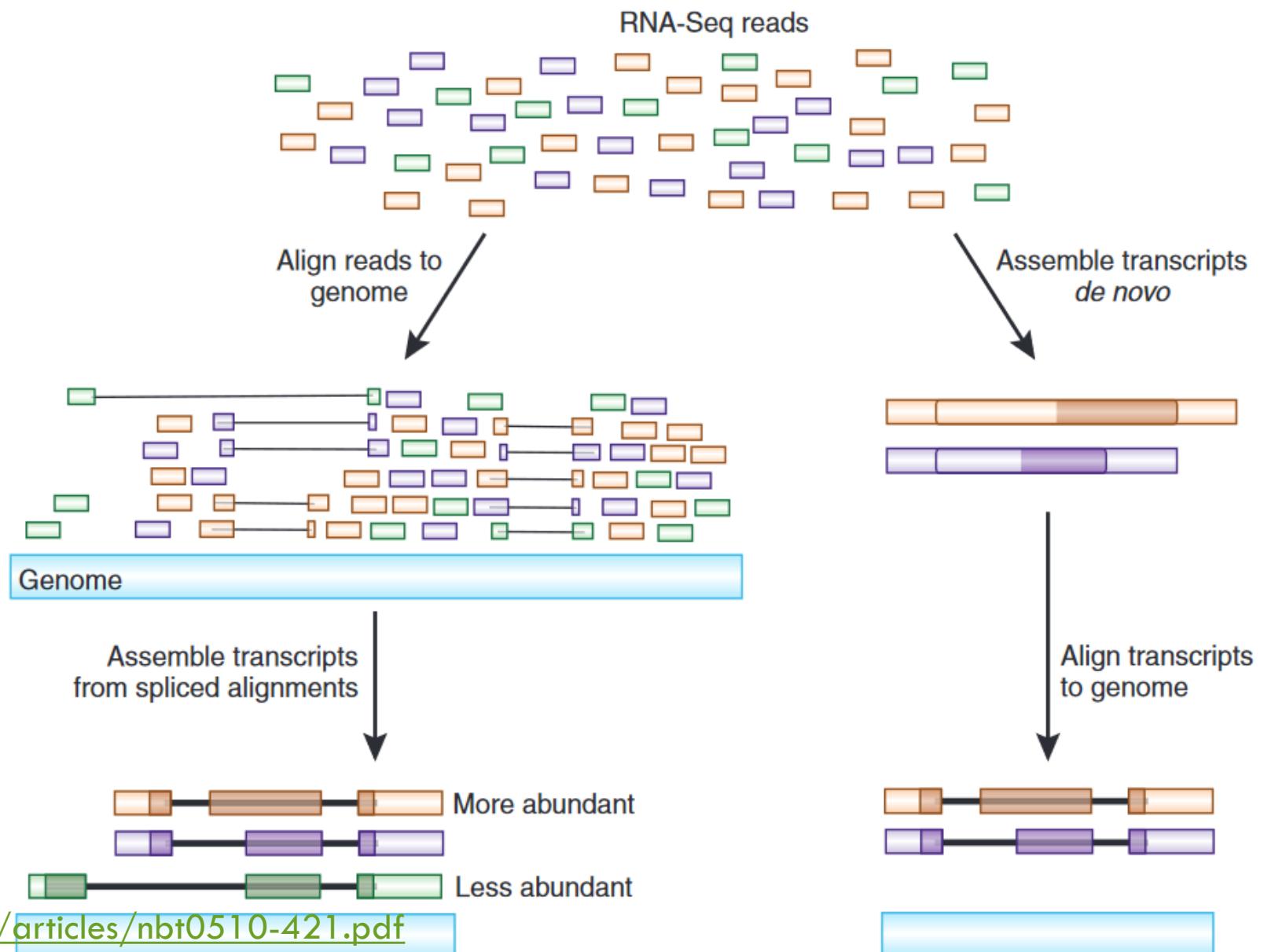
Nick Lane¹ & William Martin²



All complex life is composed of eukaryotic (nucleated) cells. The eukaryotic cell arose from prokaryotes just once in four billion years, and otherwise prokaryotes show no tendency to evolve greater complexity. Why not? Prokaryotic genome size is constrained by bioenergetics. The endosymbiosis that gave rise to mitochondria restructured the distribution of DNA in relation to bioenergetic membranes, permitting a remarkable 200,000-fold expansion in the number of genes expressed. This vast leap in genomic capacity was strictly dependent on mitochondrial power, and prerequisite to eukaryote complexity: the key innovation en route to multicellular life.

<https://www.nature.com/articles/nature09486>

RNA-SEQ ALIGNEMENT



QUANTIFYING GENE EXPRESSION

Sequencing reads can be counted on any feature (e.g. exons, introns, genes)

- Gene expression
- Expression of different isoforms (challenging)

Read
—

— - - —

Read across splice junctions



COUNTS/GENE EXPRESSION MATRIX



	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7
Gene A	345	0	23	56	76	4	3
Gene B	60	45	56	32	24	58	54
Gene C	0	0	0	0	0	0	0
Gene D	453	569	764	897	564	432	865

RNA-SEQ DATA ANALYSIS

