

# Gene Regulation

Michael Schatz

October 9, 2024

Lecture 13. Applied Comparative Genomics



# Assignment 4

## Due: Monday Oct 14, 2024 by 11:59pm

The screenshot shows a GitHub repository page for 'appliedgenomics2024' under the 'assignments' folder. The 'assignment4' folder is selected. The README.md file is open, displaying the assignment details:

**Assignment 4: RNA-seq and Machine Learning**

Assignment Date: Monday, October 7, 2024  
Due Date: Monday, October 14 @ 11:59pm

**Assignment Overview**

In this assignment you will explore a couple of key aspects of RNA-seq and introduce the key concepts of machine learning. For this assignment, we will provide a Jupyter notebook with code for you to use and complete your assignment in.

As a reminder, any questions about the assignment should be posted to [Piazza](#).

See the notebook here: [Assignment4.ipynb](#)

**Packaging**

The solutions to the above questions should be submitted as a single PDF document that includes your name, email address, and all relevant code, text, and figures (as needed). If you use ChatGPT for any of the code, also record the prompts used. Submit your solutions by uploading the PDF to [GradeScope](#), and remember to select where in your submission each question/subquestion is. The Entry Code is: Z3J8YY.

If you submit after this time, you will use your late days. Remember, you are only allowed 4 late days for the entire semester!

**Resources**

- Jupyter notebooks: <https://jupyter.org/>
- scikit-learn: <https://scikit-learn.org/stable/>
- pytorch: <https://pytorch.org/>

The screenshot shows the 'Assignment4.ipynb' file content in a Jupyter notebook interface. The code cell contains the following imports:

```
import torch
import torch.nn as nn
import torch.optim as optim
import pandas as pd
import seaborn as sns
import numpy as np
import torch.nn.functional as F
import matplotlib.pyplot as plt

from sklearn.model_selection import train_test_split
from sklearn.datasets import make_classification
from sklearn.preprocessing import StandardScaler
from torch.utils.data import DataLoader, TensorDataset
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import classification_report
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression
from sklearn.manifold import TSNE
from sklearn.decomposition import PCA
from sklearn.model_selection import GridSearchCV
```

**Question 1: Differential Expression**

**Question 1a - Sample 5000 rows**

In the files `data1.txt` and `data2.txt`, we provide an abstraction of RNA-seq data, randomly sample 5000 rows from each file. Sample 3 times for each file (this emulates making experimental replicates) and conduct a paired t-test for differential expression of each of the 15 genes. Which genes are significantly differentially expressed at the 0.05 level and what is their mean fold change?

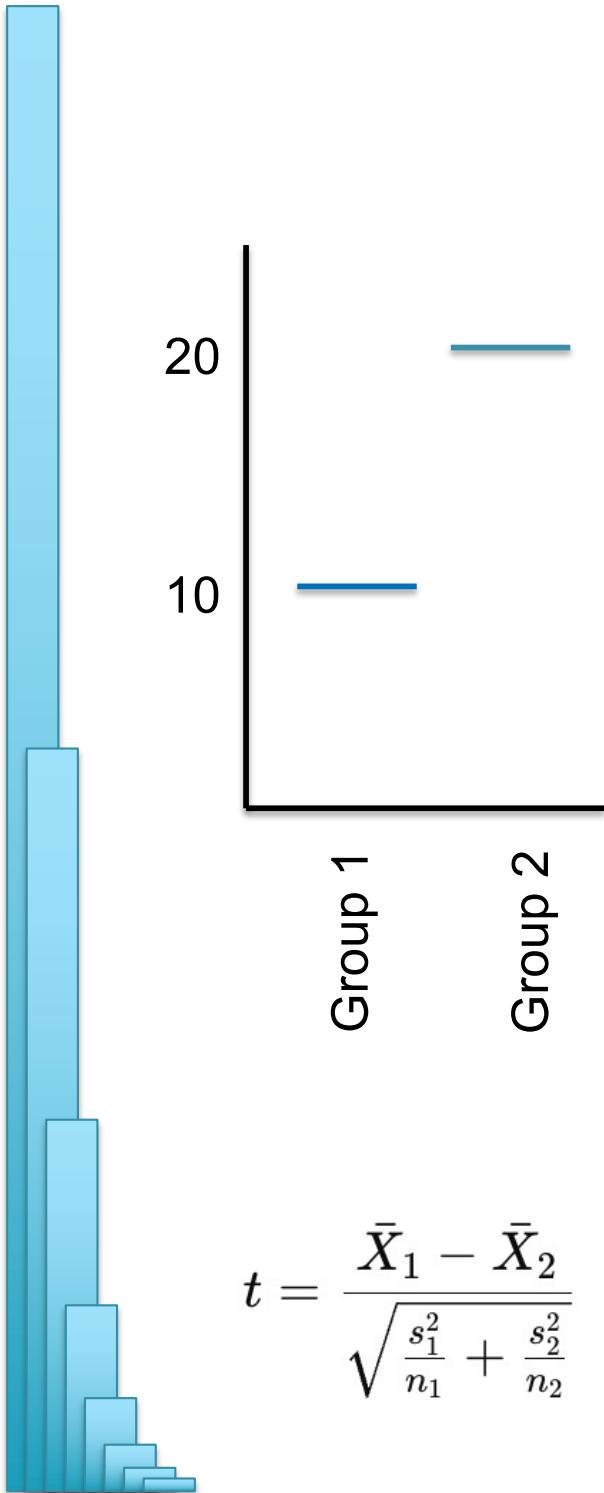
**Question 1b - Volcano plot**

Make a volcano plot of the data from part a: x-axis= $\log_2(\text{fold change of the mean expression of gene}_i)$ ; y-axis= $-\log_{10}(p\text{-value comparing the expression of gene}_i)$ . Label all of the genes that show a statistically significant change

<https://schatz-lab.org/appliedgenomics2024/assignments/assignment4/>

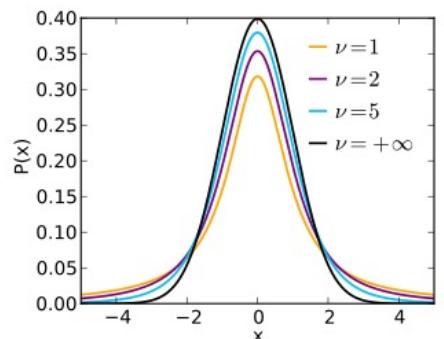
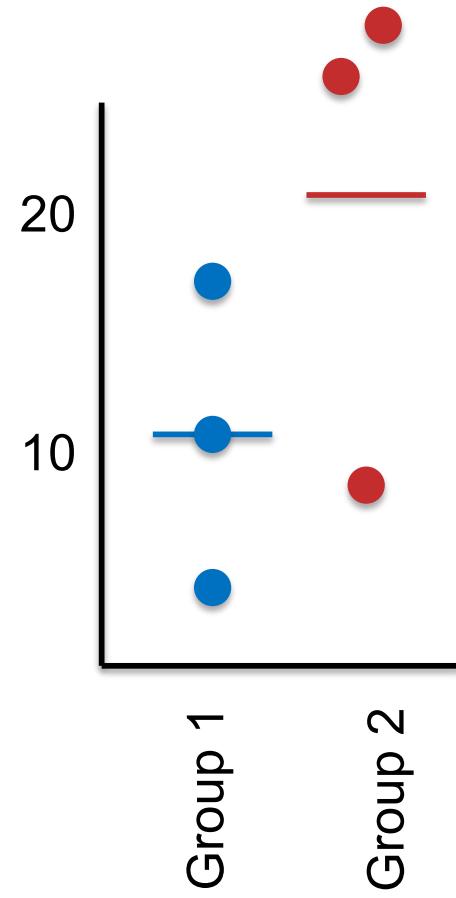
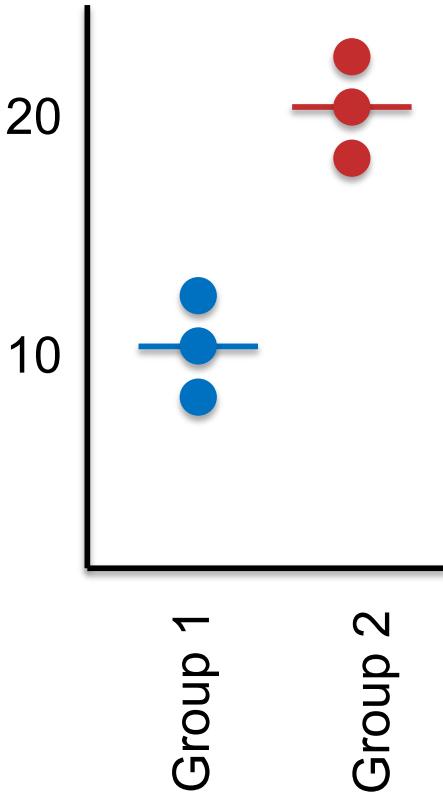
Check Piazza for questions!

# T-test



Where:

- $\bar{X}_1$  and  $\bar{X}_2$  are the means of the two samples.
- $s_1^2$  and  $s_2^2$  are the variances of the two samples.
- $n_1$  and  $n_2$  are the sizes of the two samples.



# “AlexNet”

---

## ImageNet Classification with Deep Convolutional Neural Networks

---

**Alex Krizhevsky**  
University of Toronto  
kriz@cs.utoronto.ca

**Ilya Sutskever**  
University of Toronto  
ilya@cs.utoronto.ca

**Geoffrey E. Hinton**  
University of Toronto  
hinton@cs.utoronto.ca

### Abstract

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called “dropout” that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.



**Alex Krizhevsky**  
U. Toronto/Google



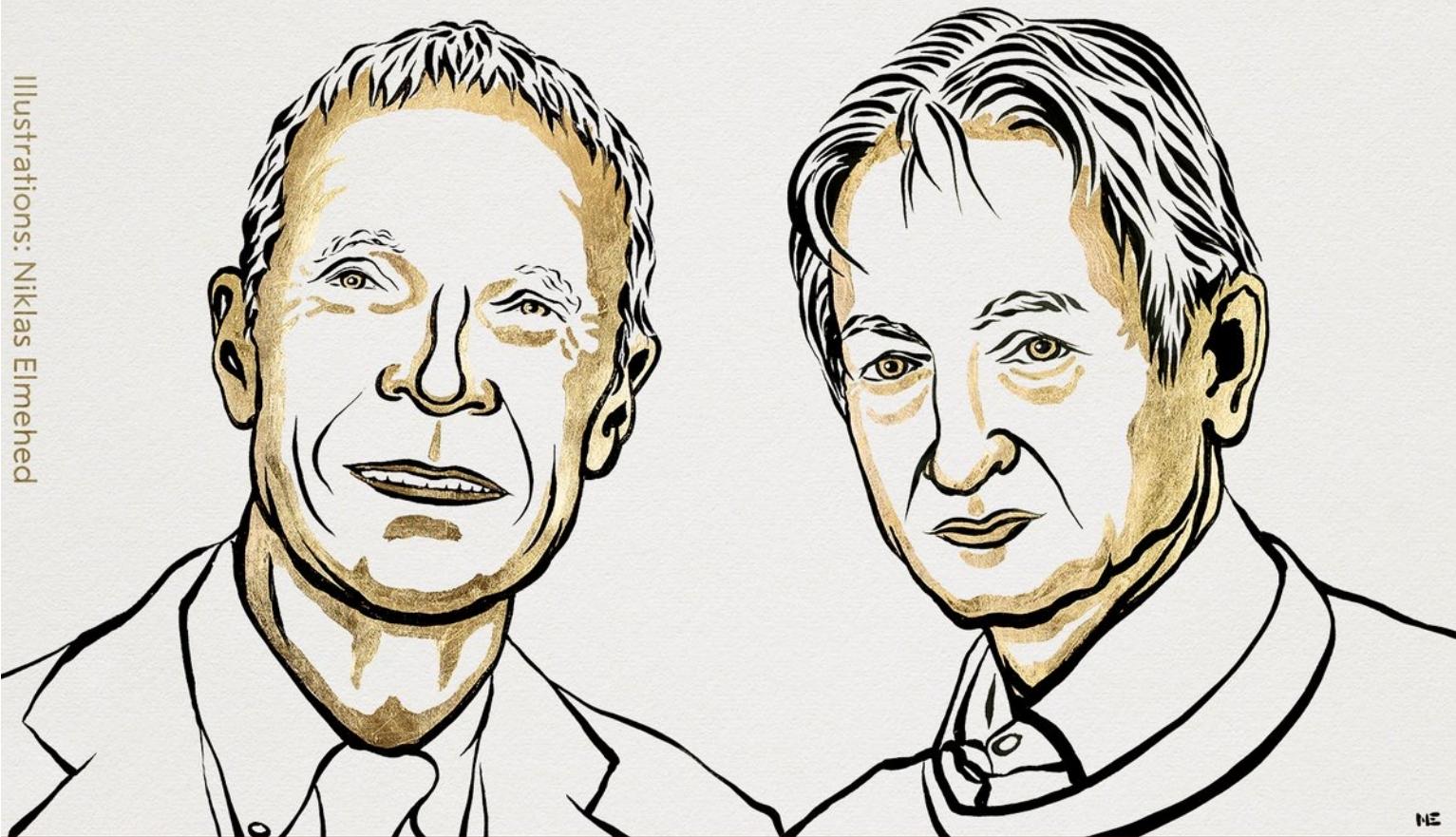
**Ilya Sutskever**  
U. Toronto/OpenAI/  
Safe Superintelligence Inc



**Geoffrey Hinton**  
U. Toronto/Vector Institute

# THE NOBEL PRIZE IN PHYSICS 2024

Illustrations: Niklas Elmehed



John J. Hopfield

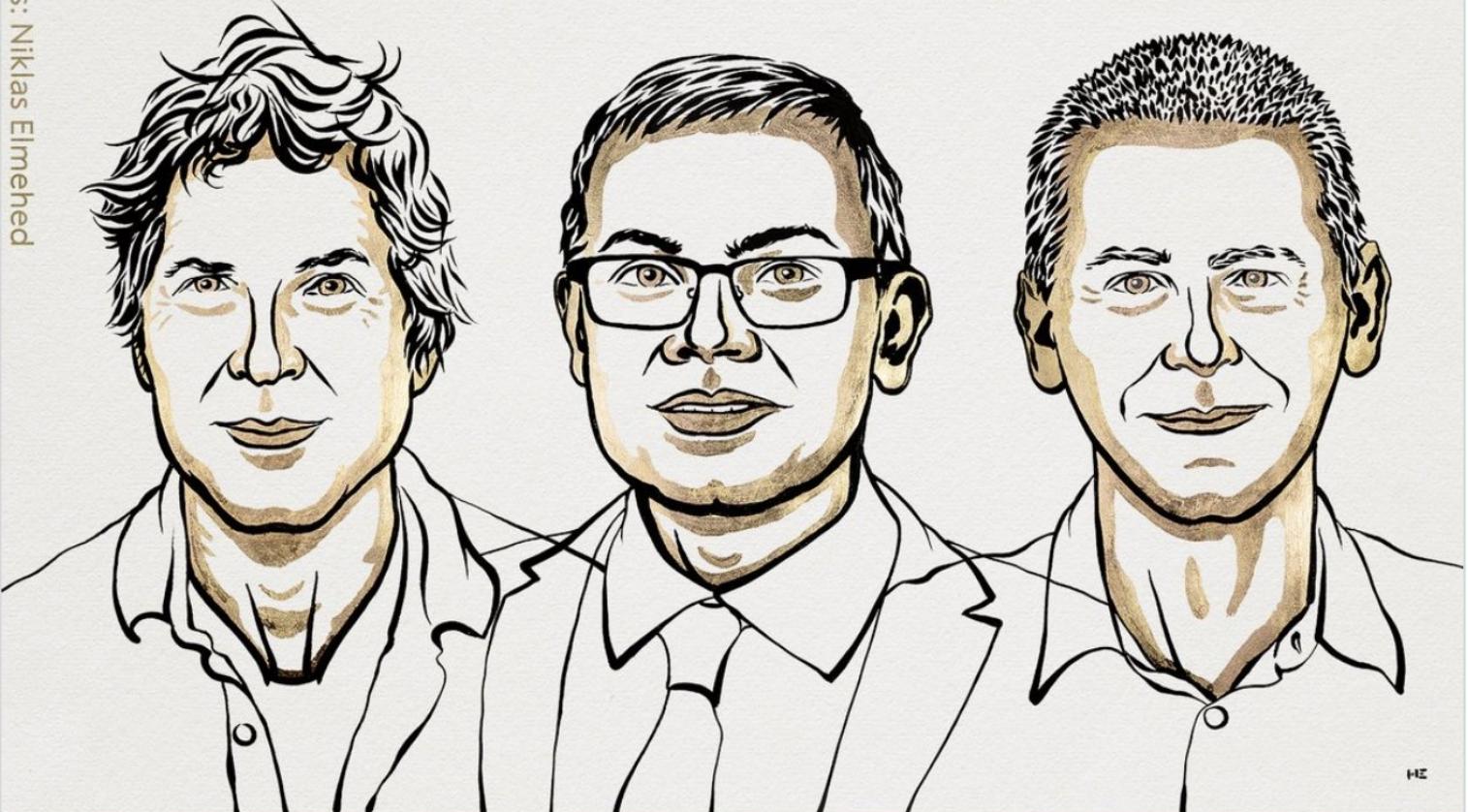
Geoffrey E. Hinton

"for foundational discoveries and inventions  
that enable machine learning  
with artificial neural networks"

THE ROYAL SWEDISH ACADEMY OF SCIENCES

Illustrations: Niklas Elmehed

# THE NOBEL PRIZE IN CHEMISTRY 2024



David  
Baker

"for computational  
protein design"

Demis  
Hassabis

"for protein structure prediction"

John M.  
Jumper

THE ROYAL SWEDISH ACADEMY OF SCIENCES

# THE NOBEL PRIZE IN PHYSIOLOGY OR MEDICINE 2024

Illustrations: Niklas Elmehed



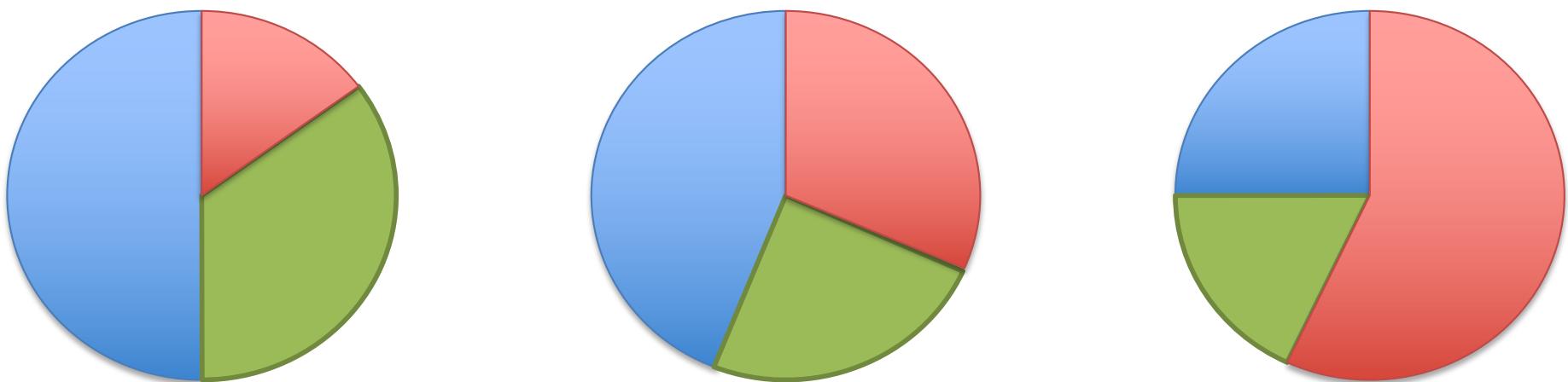
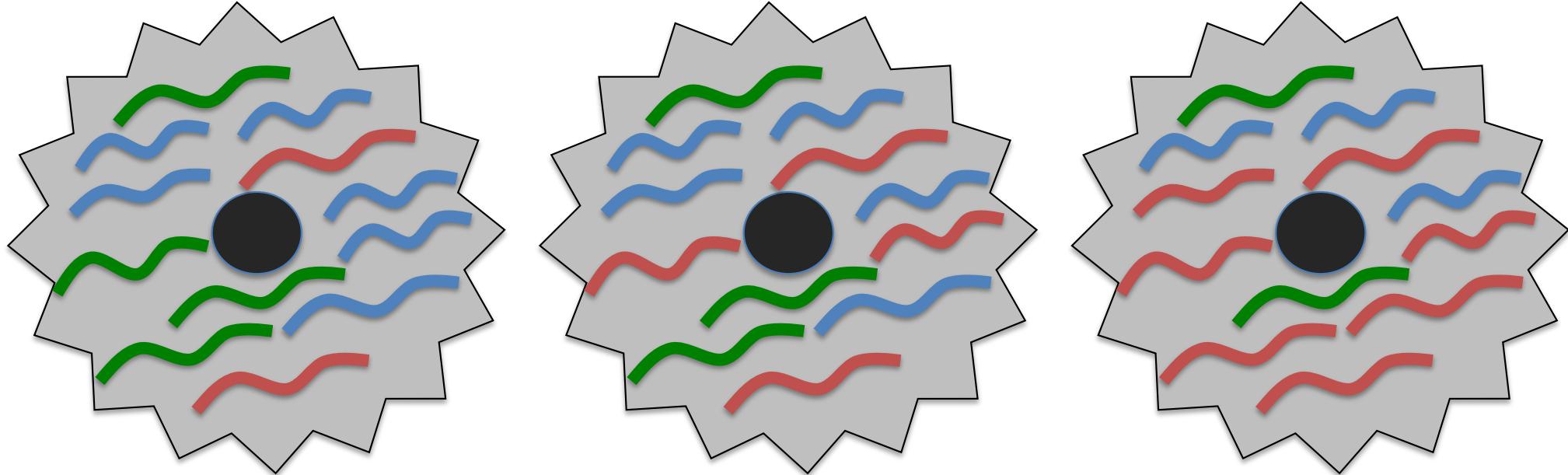
Victor Ambros

Gary Ruvkun

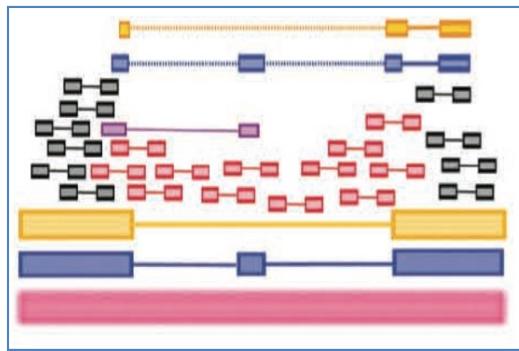
"for the discovery of microRNA and its role  
in post-transcriptional gene regulation"

THE NOBEL ASSEMBLY AT KAROLINSKA INSTITUTET

# RNA-seq Overview

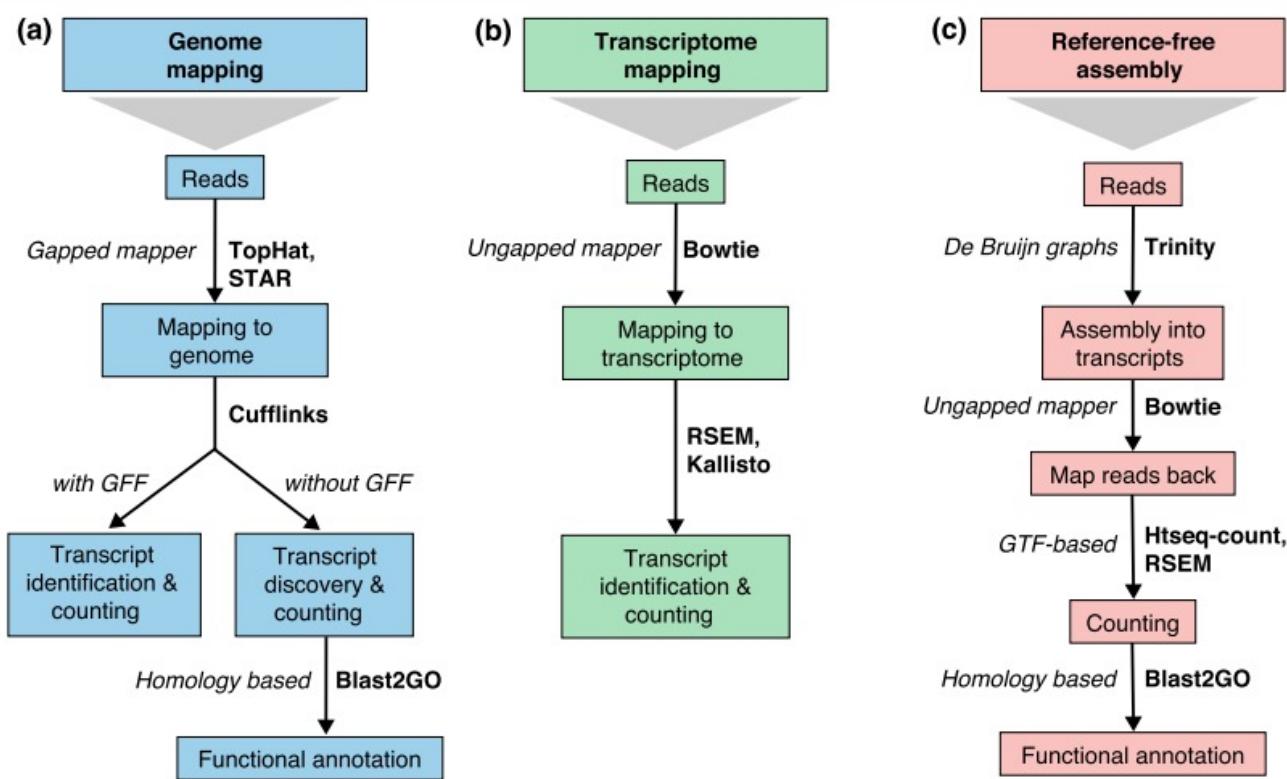


# RNA-seq Challenges



**Challenge I: Eukaryotic genes are spliced**

# RNA-Seq Approaches



**Fig. 2** Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the genome using a gapped aligner (TopHat, STAR). Transcript identification and quantification can proceed with or without an annotation file (GFF). **b** If no novel transcript discovery is needed, reads can be mapped to the reference transcriptome using an ungapped aligner (Bowtie). Transcript identification and quantification can occur simultaneously. **c** When no genome is available, reads need to be assembled first into contigs or transcripts. For quantification, reads are mapped back to the novel reference transcriptome and further analyzed using a transcriptome assembler (Trinity). This is followed by the functional annotation of the novel transcripts as in **(a)**. Representative software that can be used at each analysis step are indicated in bold text. Abbreviations: GFF General Feature Format, GTF gene transfer format, RSEM RNA-Seq by Expectation Maximization

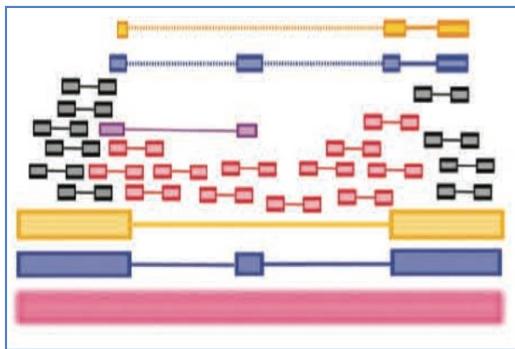
Which approach should we use?

It depends....

A survey of best practices for RNA-seq data analysis

Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8

# RNA-seq Challenges



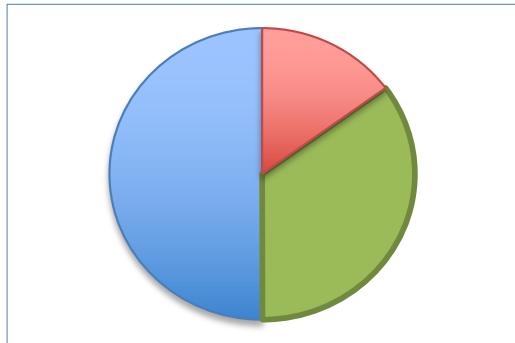
## Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

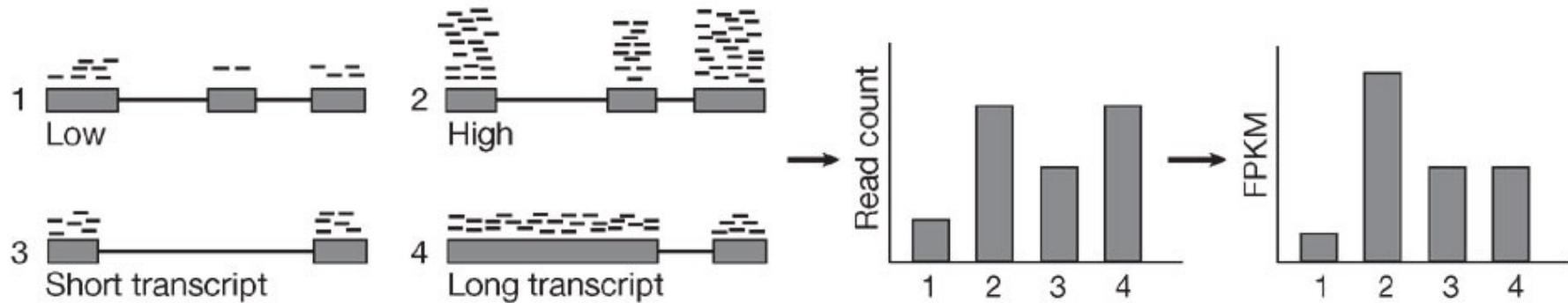
**TopHat: discovering spliced junctions with RNA-Seq.**

Trapnell et al (2009) *Bioinformatics*. 25:0 1105-1111

## Challenge 2: Read Count != Transcript abundance



# RPKM, FPKM, TPM



**Counting Reads that align to a gene DOESN'T work!**

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

**1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)**

=> Wait a second, reads in a pair aren't independent!

**2. FPKM: Fragments Per Kilobase of Exon Per Million Reads Mapped (Trapnell et al, 2010)**

=> Wait a second, FPKM depends on the average transcript length!

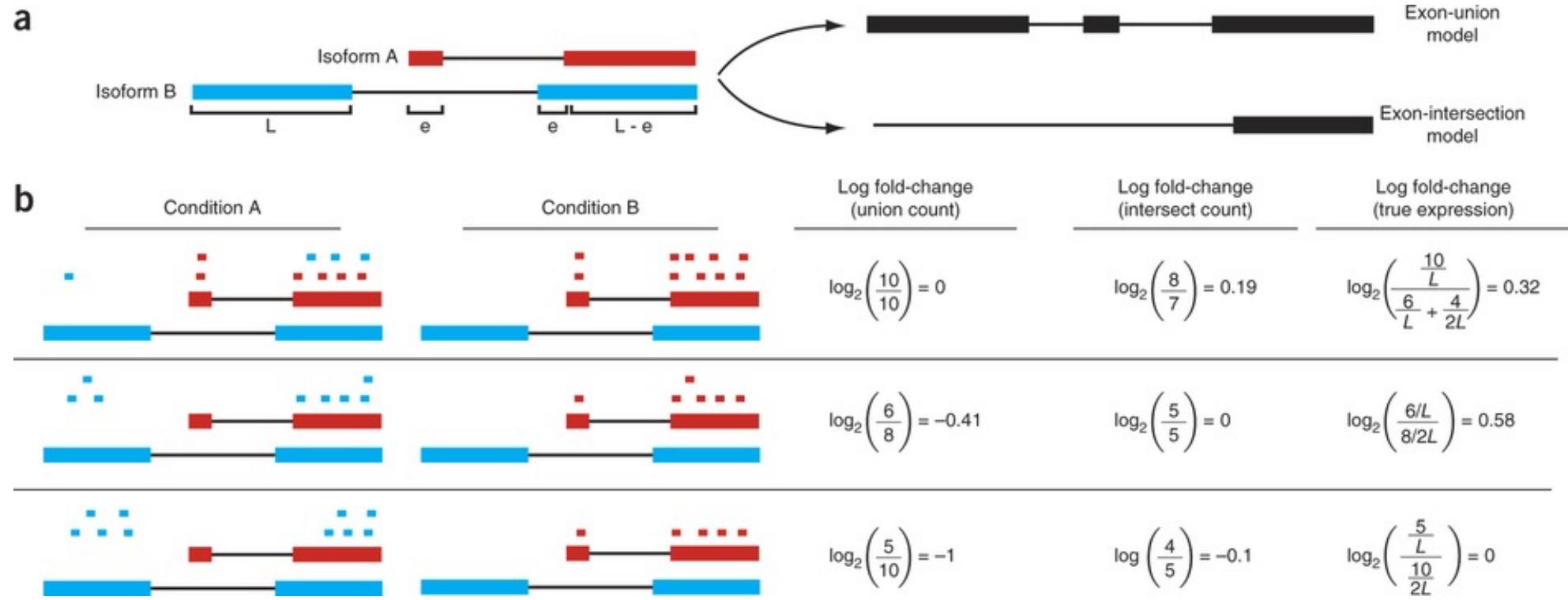
**3. TPM: Transcripts Per Million (Li et al, 2011)**

⇒ If you were to sequence one million full length transcripts, TPM is the number of transcripts you would have seen of type i, given the abundances of the other transcripts in your sample

=> Recommend you use TPM for all analysis, easy to compute given FPKM

$$TPM_i = \left( \frac{FPKM_i}{\sum_j FPKM_j} \right) \cdot 10^6$$

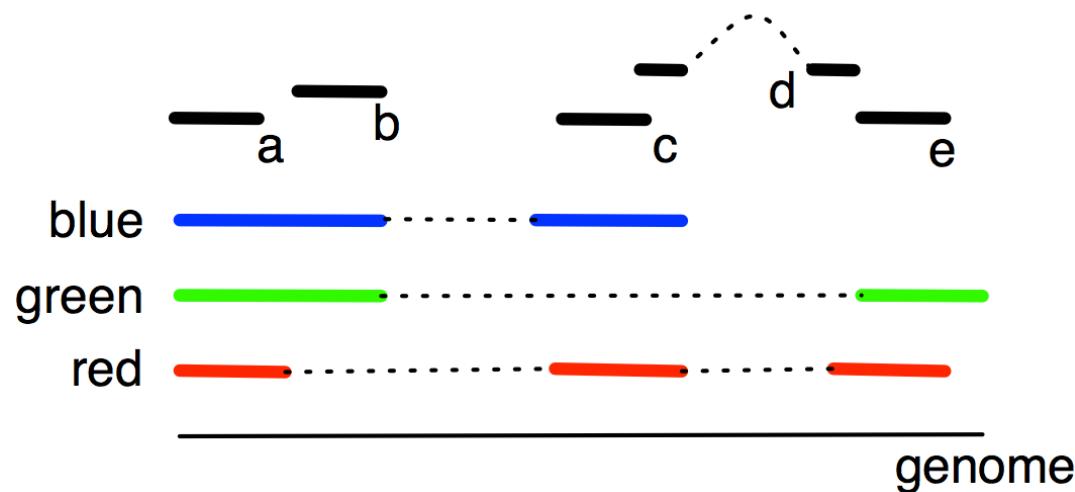
# Gene or Isoform Quantification?



**Key point : The length of the actual molecule from which the fragments derive is crucially important to obtaining accurate abundance estimates.**

**Differential analysis of gene regulation at transcript resolution with RNA-seq**  
Trapnell et al (2013) Nature Biotechnology 31, 46–53. doi:10.1038/nbt.2450

# Multi-mapping? Isoform ambiguity? Expectation Maximization to the Rescue



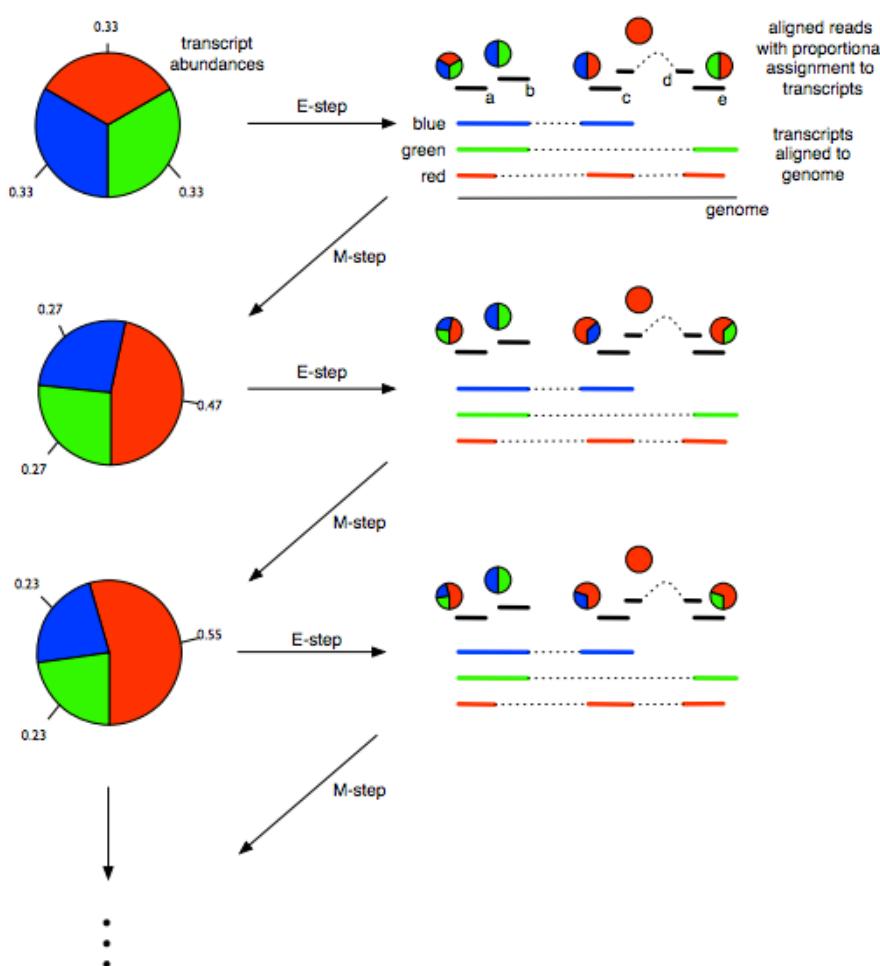
The gene has three isoforms (red, green, blue) of the same length.  
Our initial expectation is all 3 isoforms are equally expressed

There are five reads (a,b,c,d,e) mapping to the gene.

- Read a maps to all three isoforms
- Read d only to red
- Reads b,c,e map to each of the three pairs of isoforms.

What is the most likely expression level of each isoform?

# Multi-mapping? Isoform ambiguity? Expectation Maximization to the Rescue



The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

There are five reads (a,b,c,d,e) mapping to the gene. Read a maps to all three isoforms, read d only to red, and the other three (reads b,c,e) to each of the three pairs of isoforms.

During the expectation (E) step reads are proportionately assigned to transcripts according to the (current) isoform abundances (RGB): a=(.33,.33,.33), b=(0,.5,.5), c=(.5,.5), d=(1,0,0), e=(.5,.5,0)

Next, during the maximization (M) step isoform abundances are recalculated from the proportionately assigned read counts:

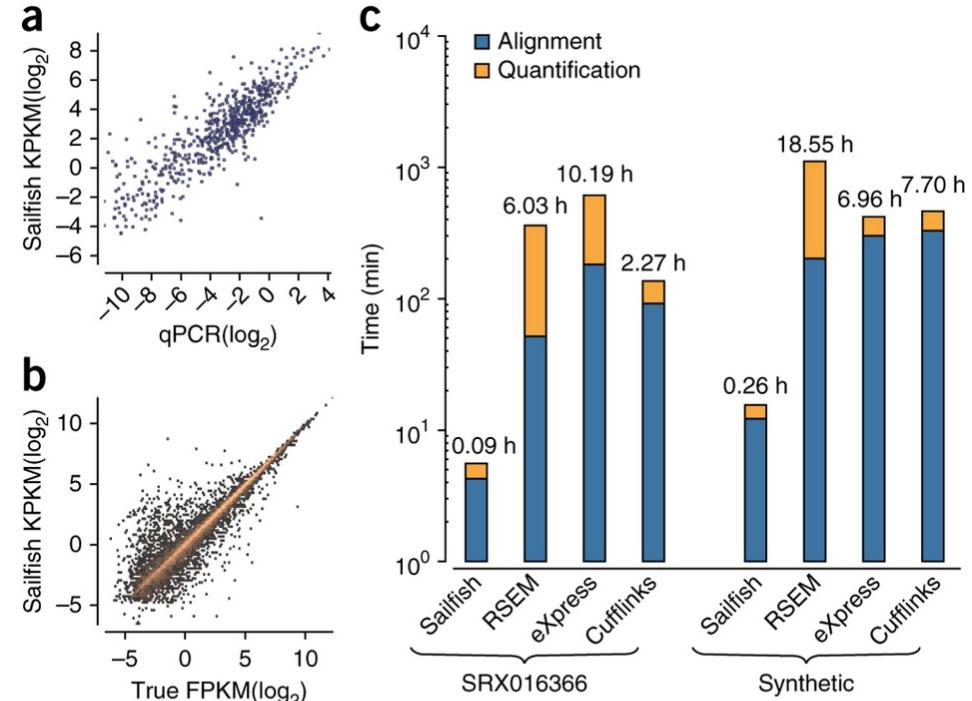
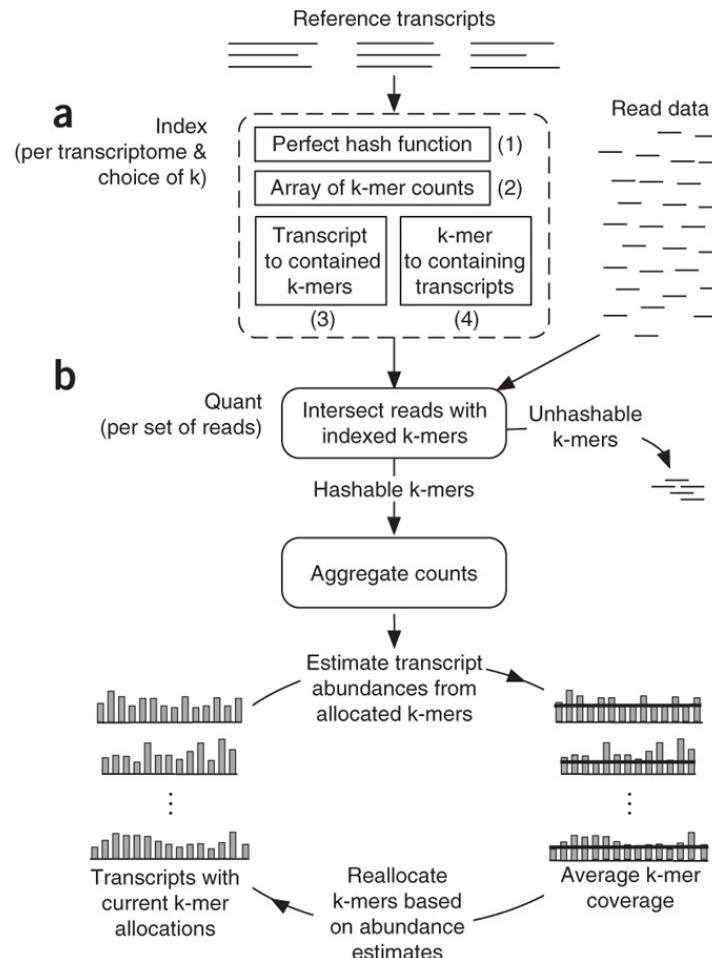
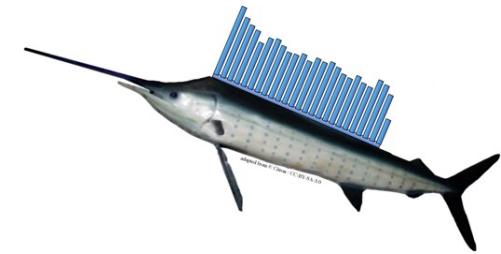
$$\text{red: } 0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)$$

$$\text{blue: } 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)$$

$$\text{green: } 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)$$

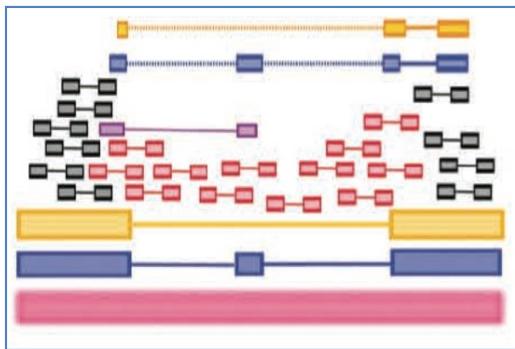
Repeat until convergence!

# Sailfish: Fast & Accurate RNA-seq Quantification



**Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms**  
 Patro et al (2014) Nature Biotechnology 32, 462–464 doi:10.1038/nbt.2862

# RNA-seq Challenges

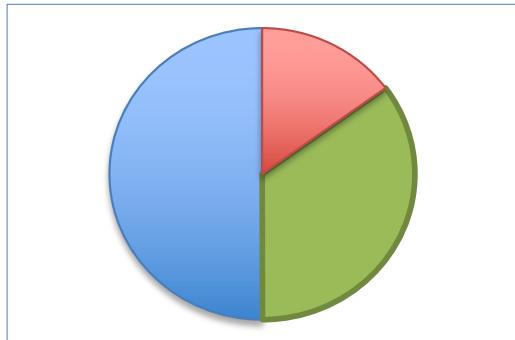


## Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

**TopHat: discovering spliced junctions with RNA-Seq.**

Trapnell et al (2009) *Bioinformatics*. 25:0 1105-1111

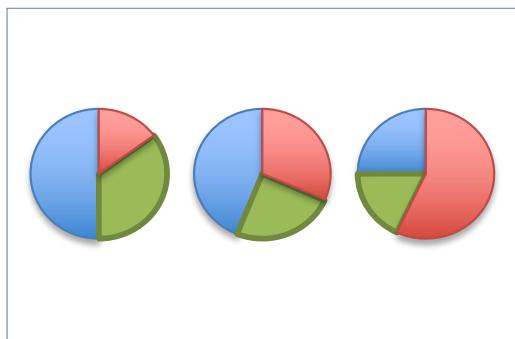


## Challenge 2: Read Count != Transcript abundance

Solution: Infer underlying abundances (e.g. TPM)

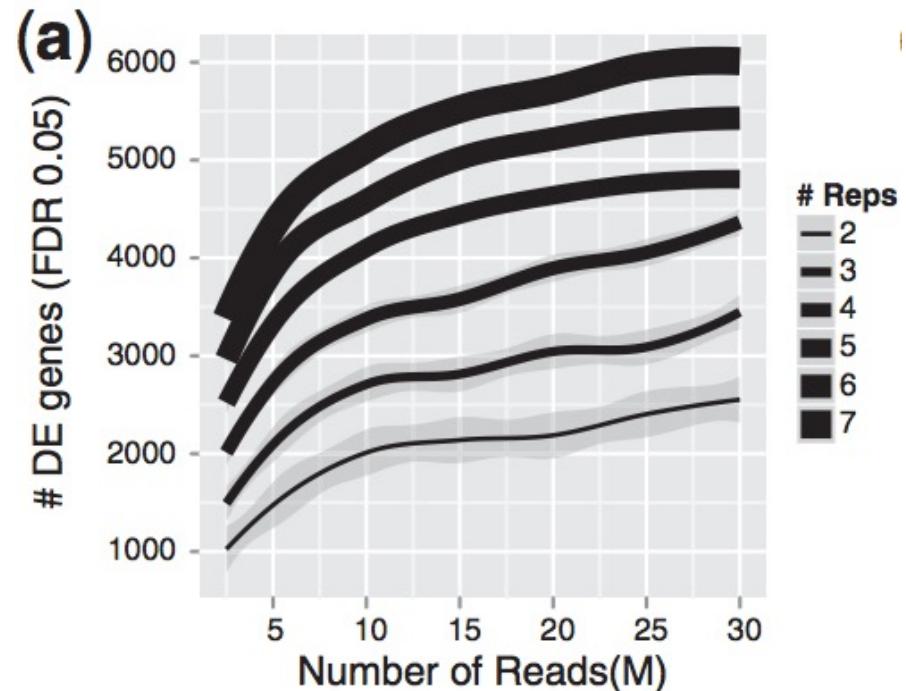
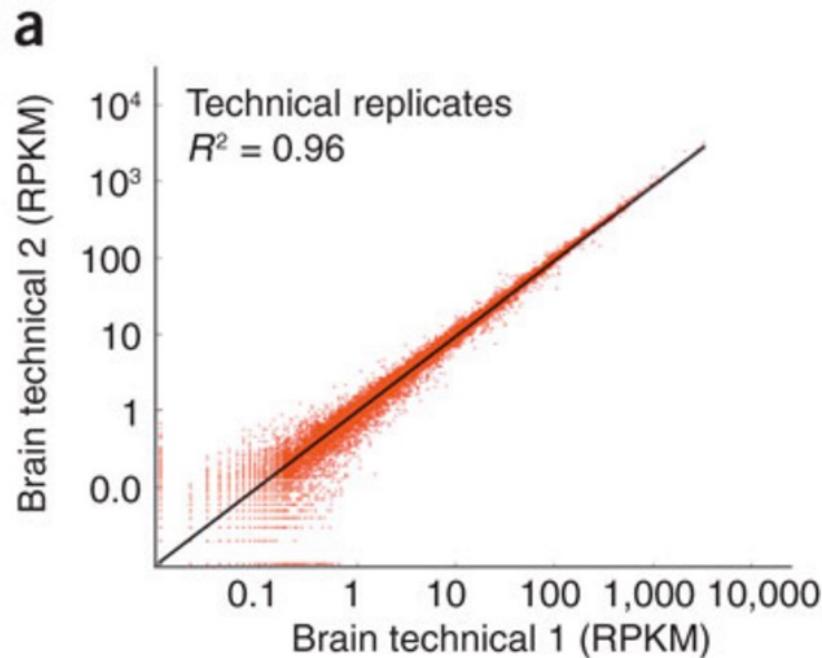
**Transcript assembly and quantification by RNA-seq**

Trapnell et al (2010) *Nat. Biotech.* 25(5): 511-515



## Challenge 3: Transcript abundances are stochastic

# How Many Replicates?

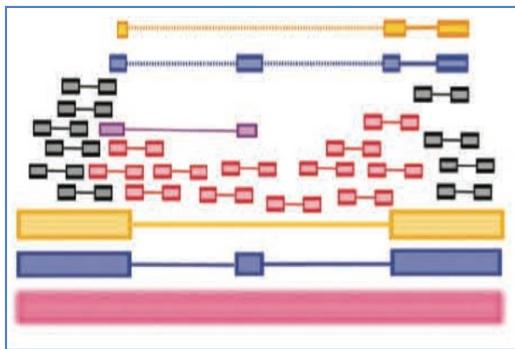


Why don't we have perfect replicates?

**Mapping and quantifying mammalian transcriptomes by RNA-Seq**  
Mortazavi et al (2008) Nature Methods. 5, 62-628

**RNA-seq differential expression studies: more sequence or more replication?**  
Liu et al (2013) Bioinformatics. doi:10.1093/bioinformatics/btt688

# RNA-seq Challenges

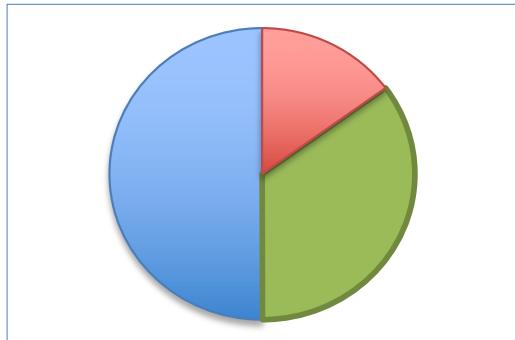


## Challenge 1: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

**TopHat: discovering spliced junctions with RNA-Seq.**

Trapnell et al (2009) *Bioinformatics*. 25:0 1105-1111

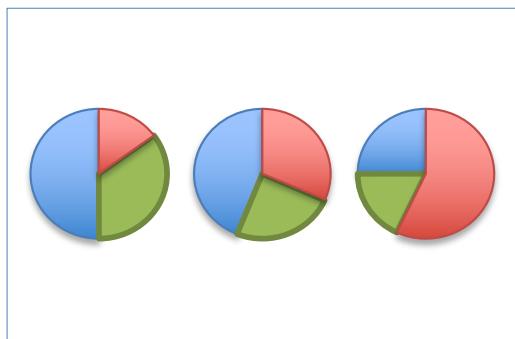


## Challenge 2: Read Count != Transcript abundance

Solution: Infer underlying abundances (e.g. TPM)

**Transcript assembly and quantification by RNA-seq**

Trapnell et al (2010) *Nat. Biotech.* 25(5): 511-515



## Challenge 3: Transcript abundances are stochastic

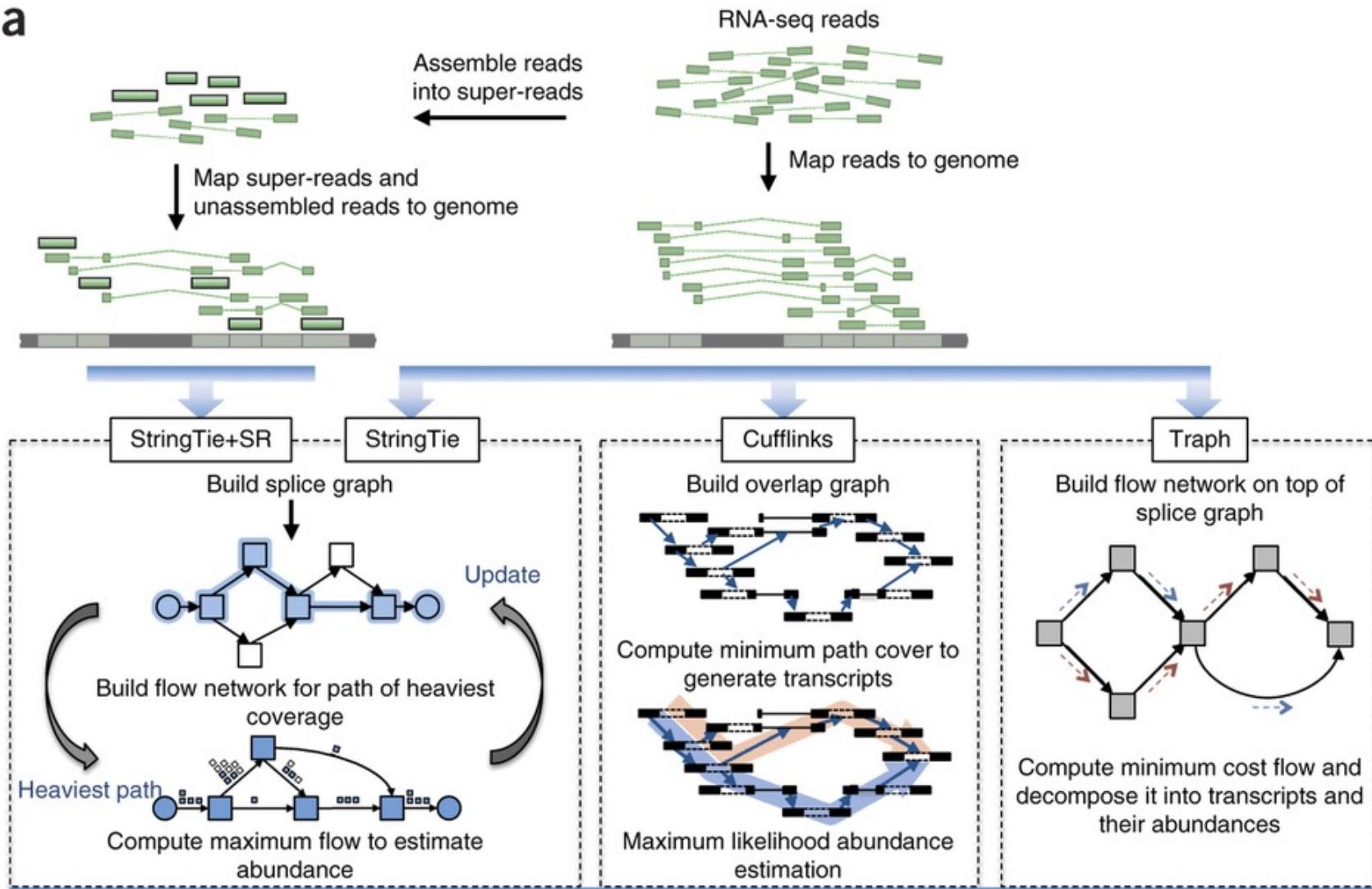
Solution: Replicates, replicates, and more replicates

**RNA-seq differential expression studies: more sequence or more replication?**

Liu et al (2013) *Bioinformatics*. doi:10.1093/bioinformatics/btt688

# Isoform Quantification Approaches

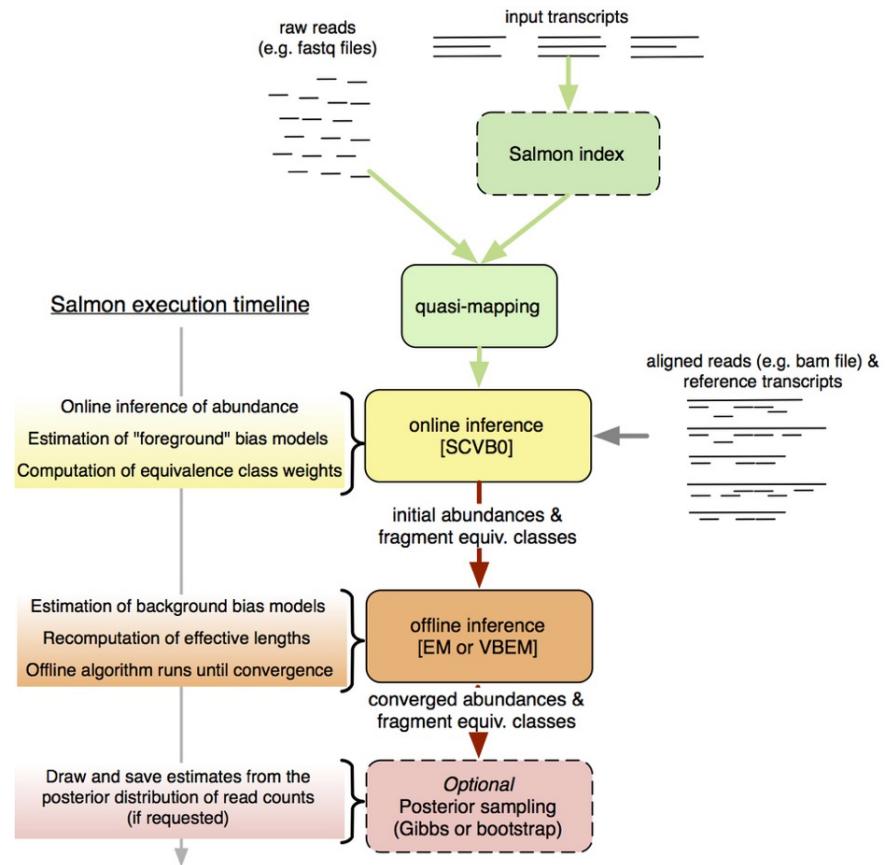
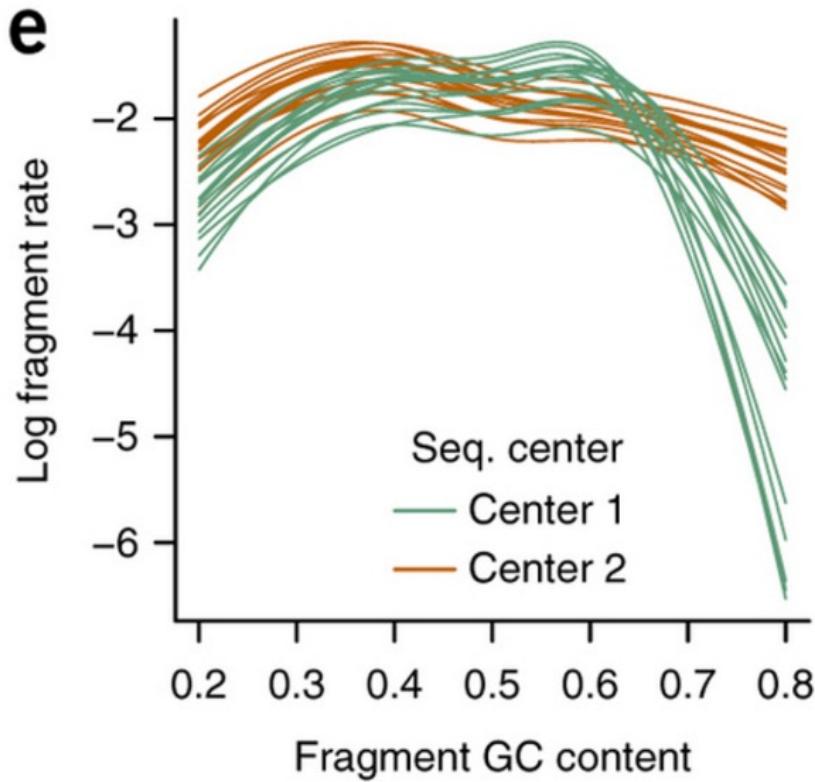
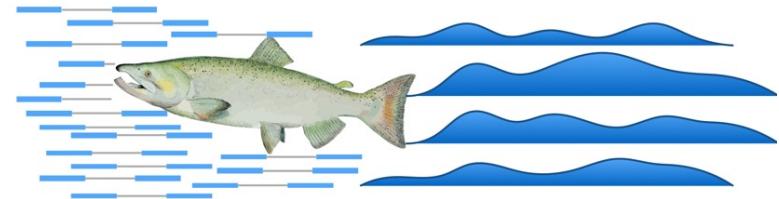
a



**StringTie enables improved reconstruction of a transcriptome from RNA-seq reads.**

Pertea M, et al. (2015) Nature Biotechnology. doi: 10.1038/nbt.3122.

# Salmon: The ultimate RNA-seq Pipeline?

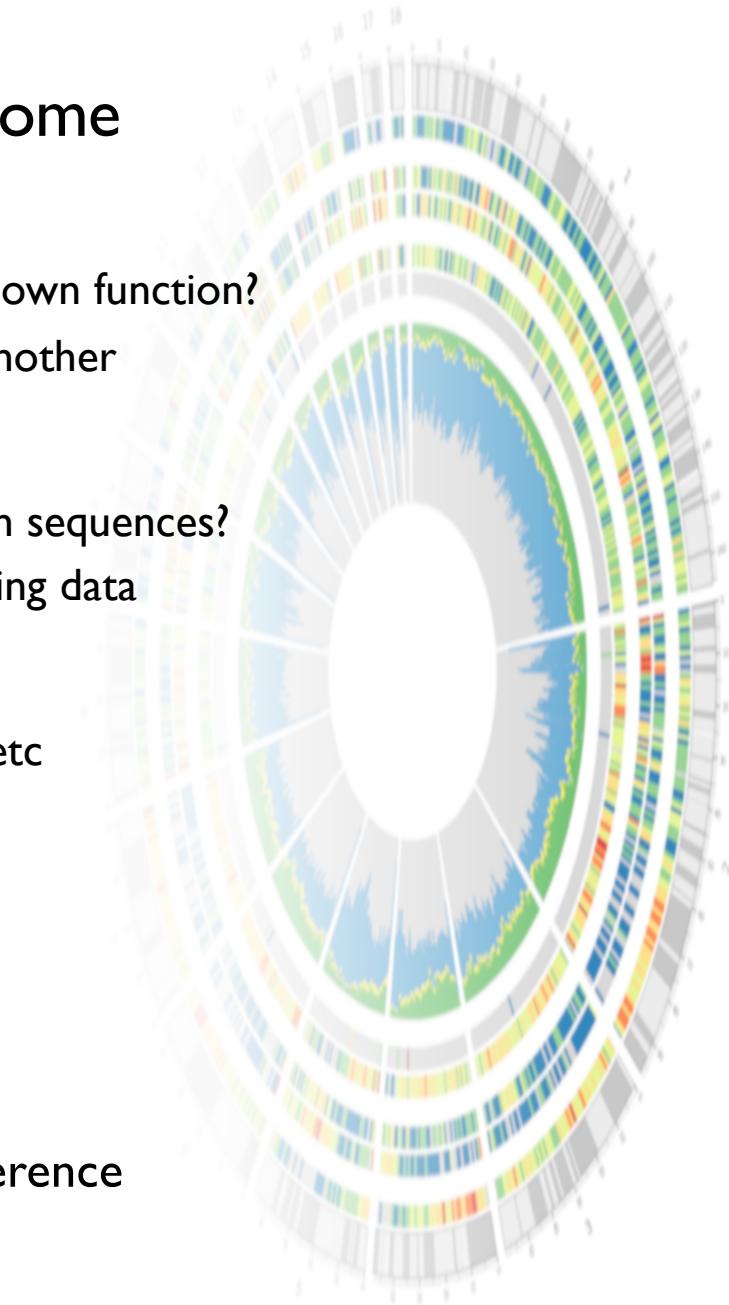


**Modeling of RNA-seq fragment sequence bias reduces systematic errors in transcript abundance estimation**  
Love et al (2016) Nature Biotechnology 34, 1287–1291 (2016) doi:10.1038/nbt.3682

**Salmon provides fast and bias-aware quantification of transcript expression**  
Patro et al (2017) Nature Methods (2017) doi:10.1038/nmeth.4197

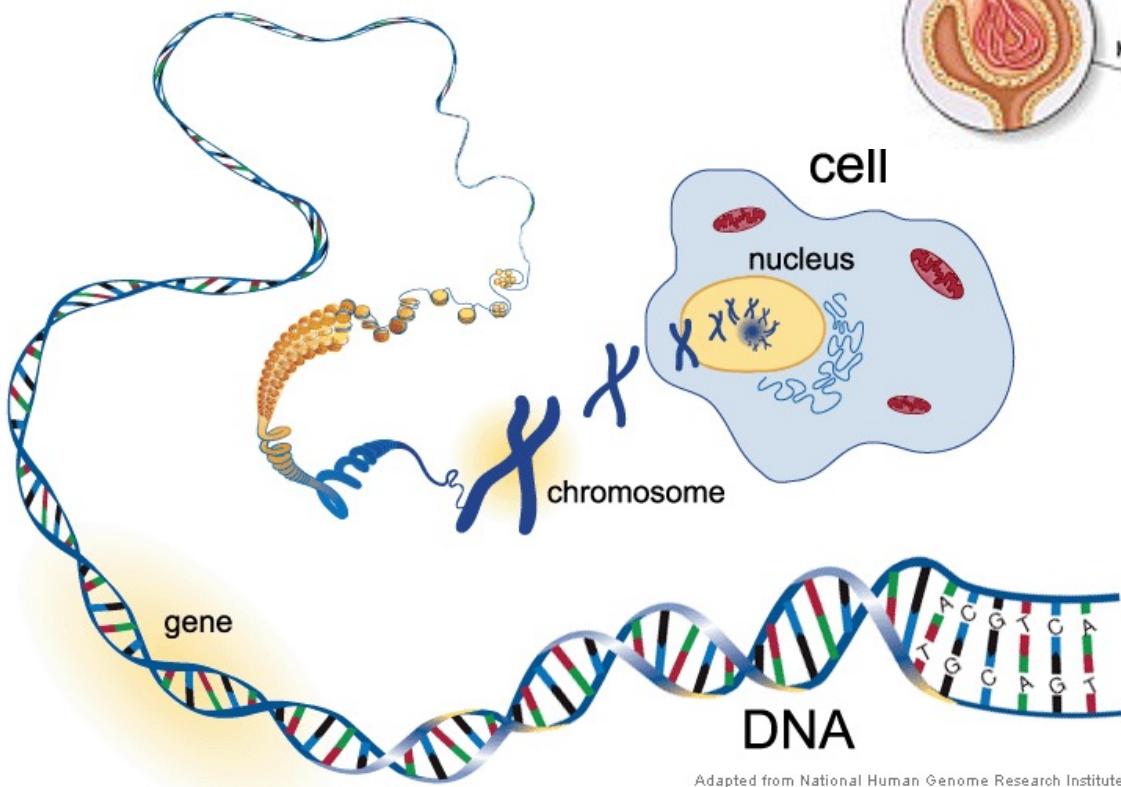
# Annotation Summary

- Three major approaches to annotate a genome
  - 1. Alignment:
    - Does this sequence align to any other sequences of known function?
    - Great for projecting knowledge from one species to another
  - 2. Prediction:
    - Does this sequence statistically resemble other known sequences?
    - Potentially most flexible but dependent on good training data
  - 3. Experimental:
    - Lets test to see if it is transcribed/methylated/bound/etc
    - Strongest but expensive and context dependent
- Many great resources available
  - Learn to love the literature and the databases
  - Standard formats let you rapidly query and cross reference
  - Google is your number one resource ☺

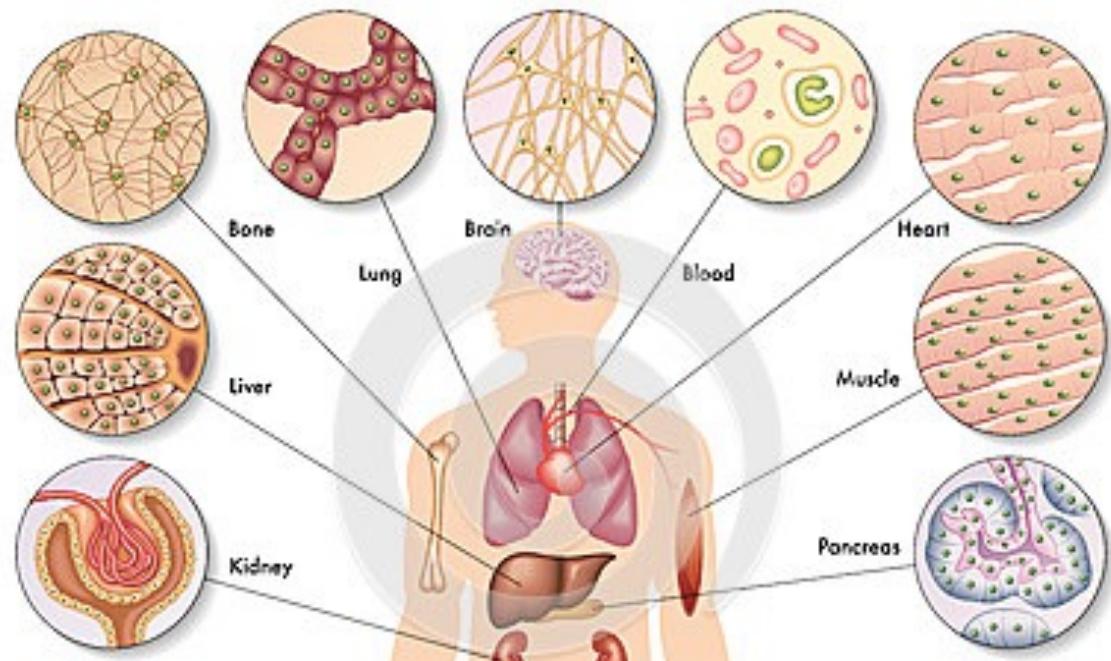


# Why Genes?

Each cell of your body contains an exact copy of your 3 billion base pair genome.



Adapted from National Human Genome Research Institute



Your body has a few hundred (thousands?) major cell types, largely defined by the gene expression patterns

# Human Evolution



~5 Mya

- Humans and chimpanzees shared a common ancestor ~5-7 million years ago (Mya)
- Single-nucleotide substitutions occur at a mean rate of 1.23% but ~4% overall rate of mutation: comprising ~35 million single nucleotide differences and ~90 Mb of insertions and deletions
- Orthologous proteins in human and chimpanzee are extremely similar, with ~29% being identical and the typical orthologue differing by only two amino acids, one per lineage

***Initial sequence of the chimpanzee genome and comparison with the human genome***  
(2005) *Nature* 437, 69-87 doi:10.1038/nature04072

# Human Evolution



~5 Mya

~75 Mya

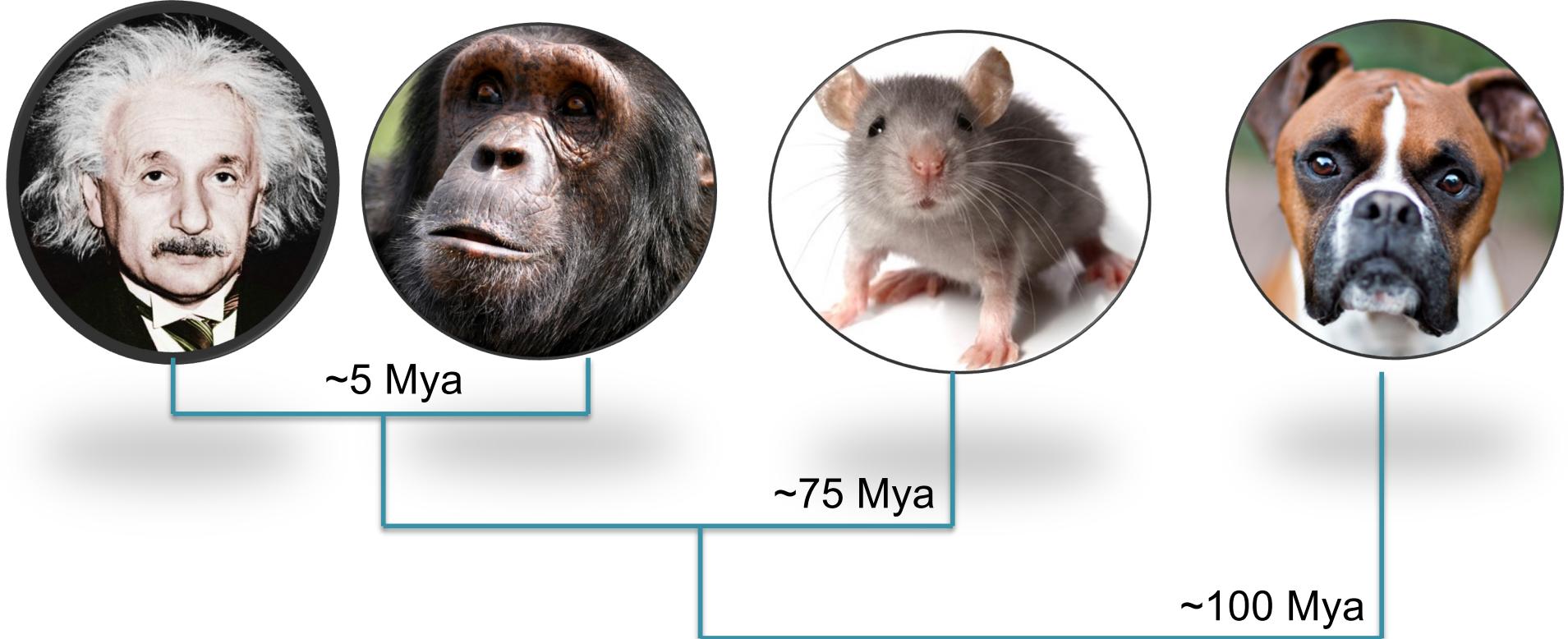
"In the roughly 75 million years since the divergence of the human and mouse lineages, the process of evolution has altered their genome sequences and caused them to diverge by ***nearly one substitution for every two nucleotides***"

***The mouse and human genomes each seem to contain about 30,000 protein-coding genes.*** These refined estimates have been derived from both new evidence-based analyses that produce larger and more complete sets of gene predictions, and new de novo gene predictions that do not rely on previous evidence of transcription or homology. The proportion of mouse genes with a single identifiable orthologue in the human genome seems to be approximately 80%. ***The proportion of mouse genes without any homologue currently detectable in the human genome (and vice versa) seems to be less than 1%.***"

***Initial sequencing and comparative analysis of the mouse genome***

Chinwalla et al (2002) Nature. 420, 520-562 doi:10.1038/nature01262

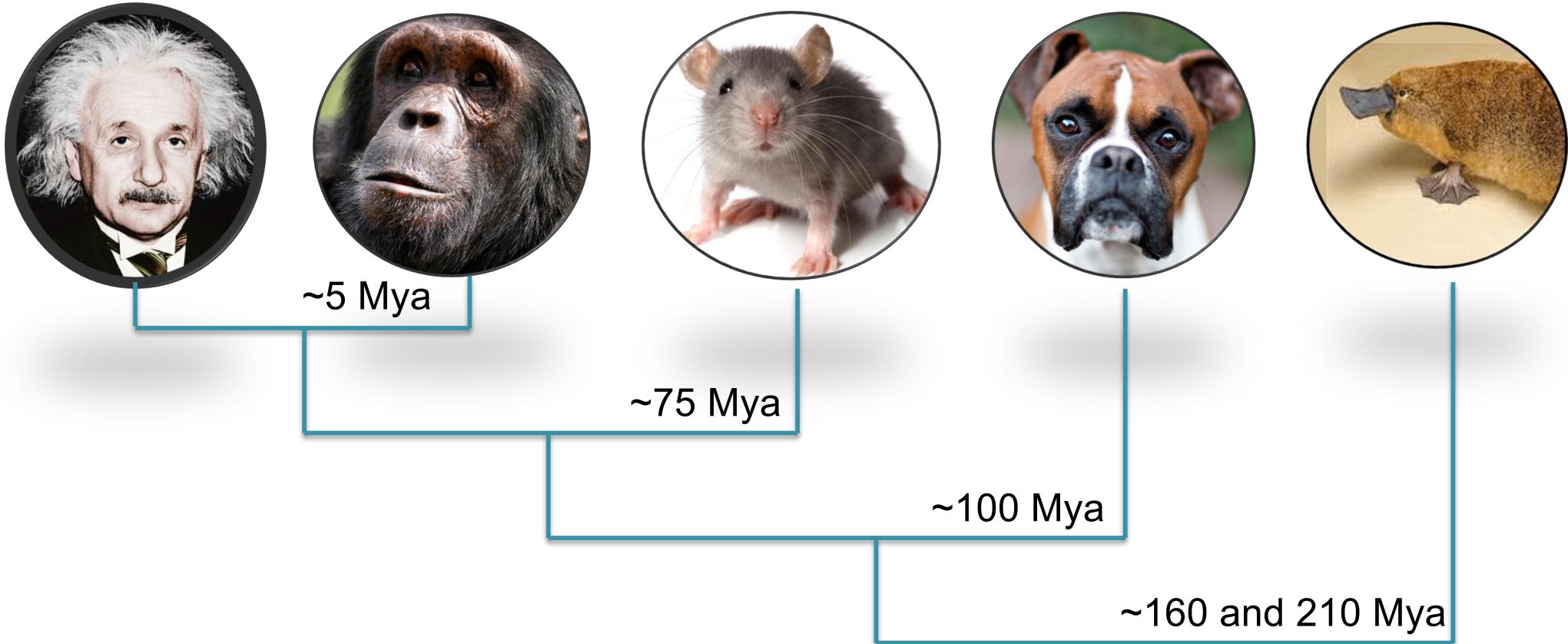
# Human Evolution



"We generated gene predictions for the dog genome using an evidence-based method (see Supplementary Information). The resulting collection contains **19,300 dog gene predictions, with nearly all being clear homologues of known human genes**. The dog gene count is substantially lower than the ~22,000-gene models in the current human gene catalogue (Ensembl build 26). For many predicted human genes, we find no convincing evidence of a corresponding dog gene. Much of the excess in the human gene count is attributable to **spurious gene predictions in the human genome**"

**Genome sequence, comparative analysis and haplotype structure of the domestic dog**  
Lindblad-Toh et al (2005) Nature. 438, 803-819 doi:10.1038/nature04338

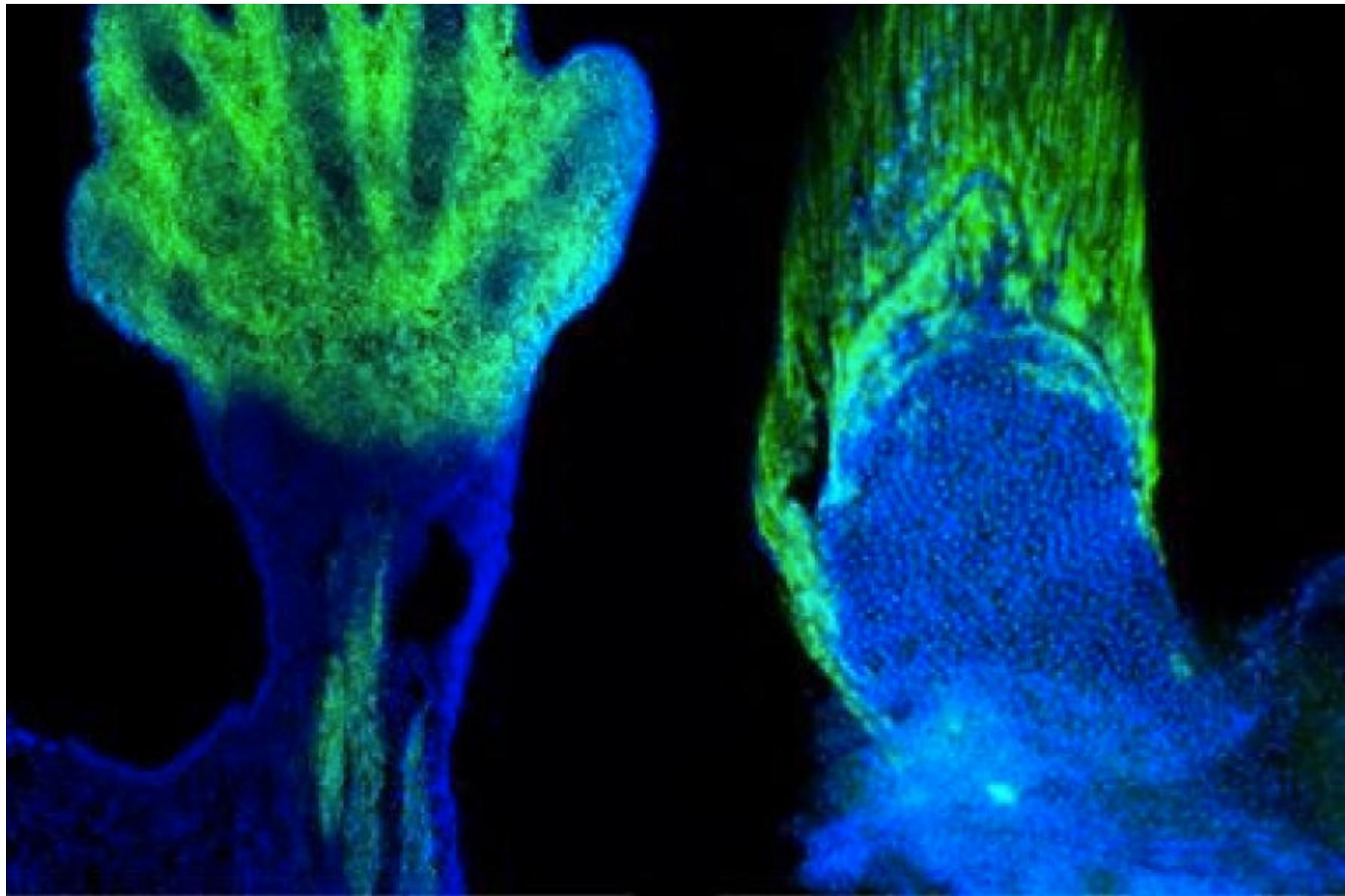
# Human Evolution



**As expected, the majority of platypus genes (82%; 15,312 out of 18,596) have orthologues in these five other amniotes** (Supplementary Table 5). The remaining 'orphan' genes are expected to primarily reflect rapidly evolving genes, for which no other homologues are discernible, erroneous predictions, and true lineage-specific genes that have been lost in each of the other five species under consideration.

**Genome analysis of the platypus reveals unique signatures of evolution**  
(2008) *Nature*. 453, 175-183 doi:10.1038/nature06936

# Animal Evolution



***Digits and fin rays share common developmental histories***

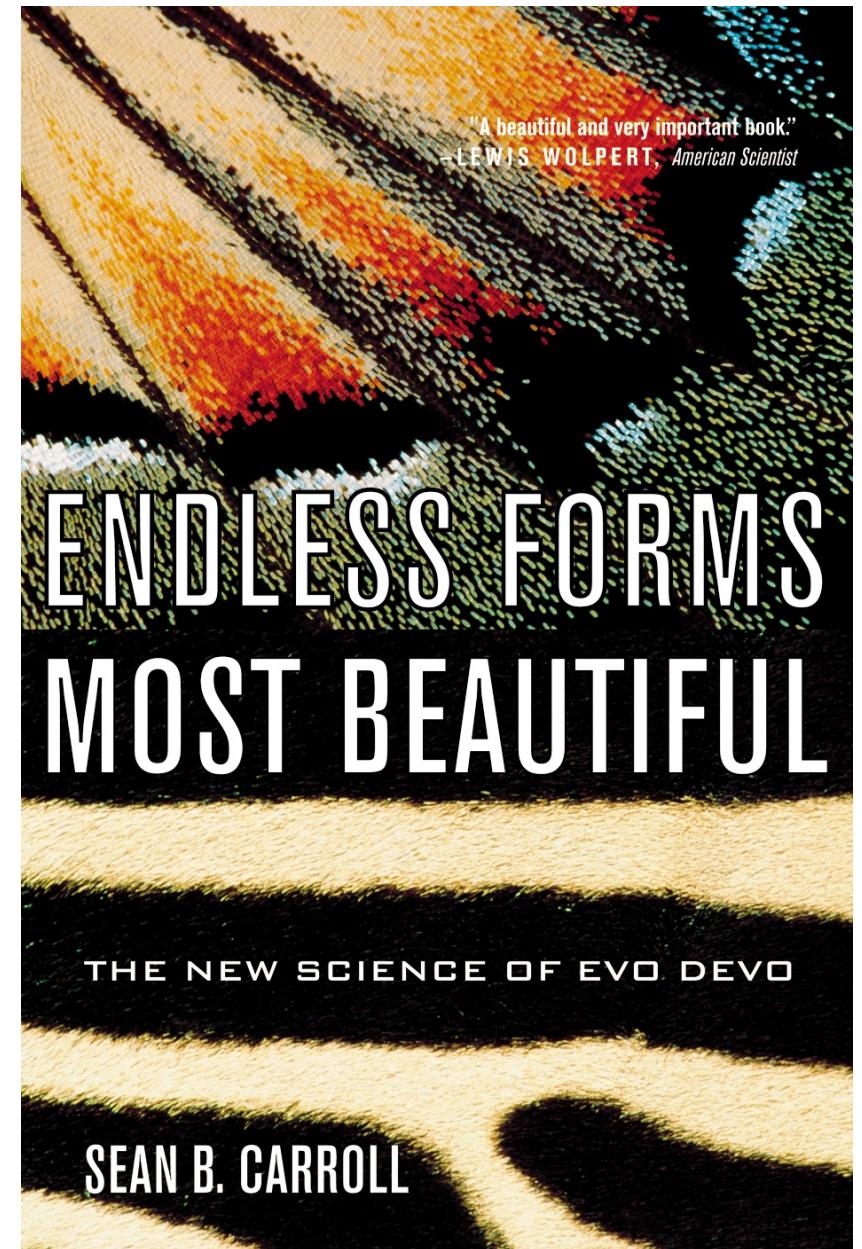
Nakamura et al (2016) *Nature*. 537, 225–228. doi:10.1038/nature19322

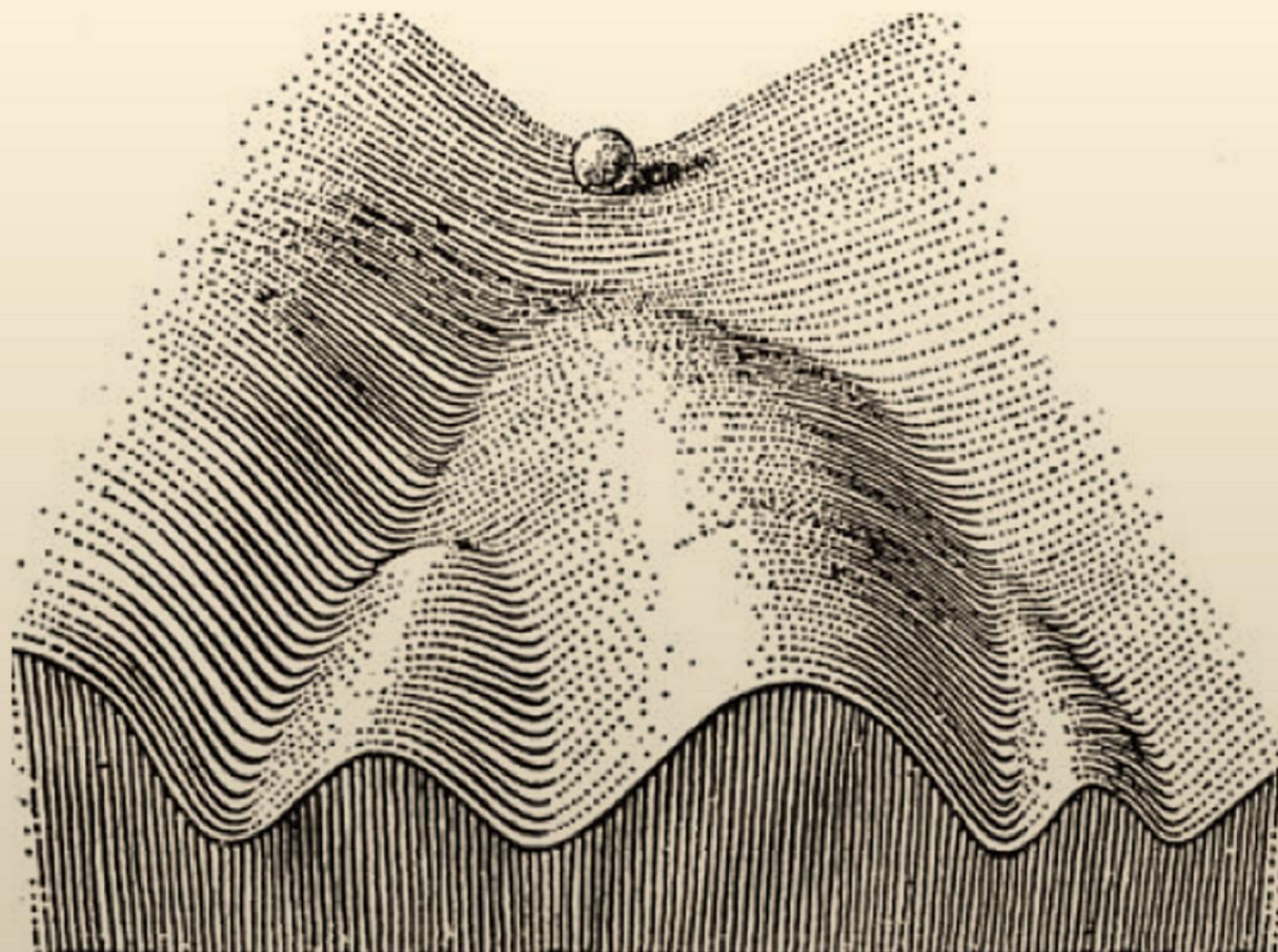
# More Information



*“Anything found to be true of  
*E. coli* must also be true of  
elephants”*

-Jacques Monod

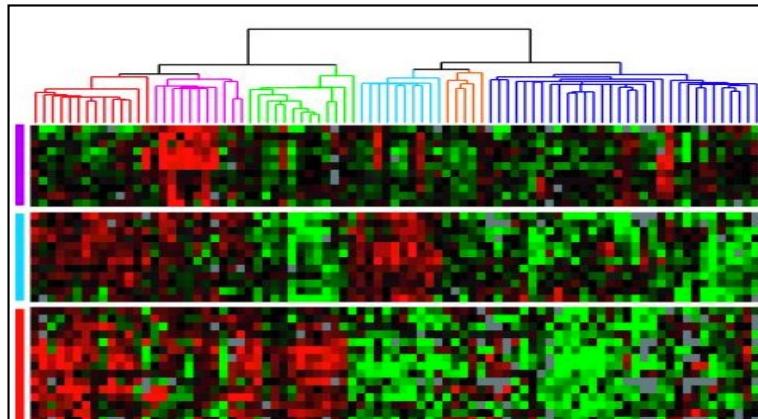




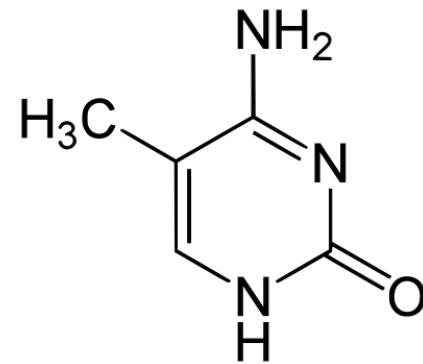
**The Strategy of the Genes**  
CH Waddington (1957)

# \*-seq in 4 short vignettes

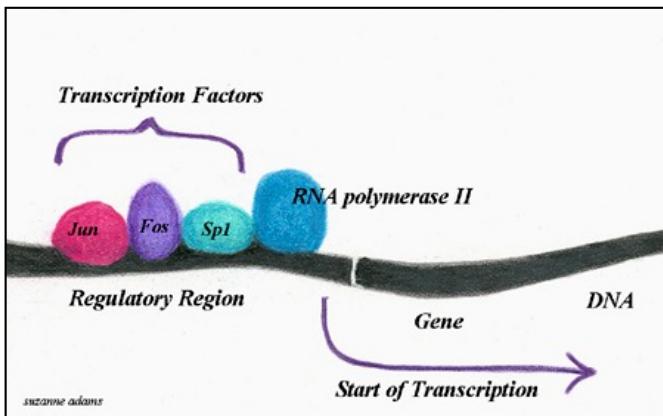
## RNA-seq



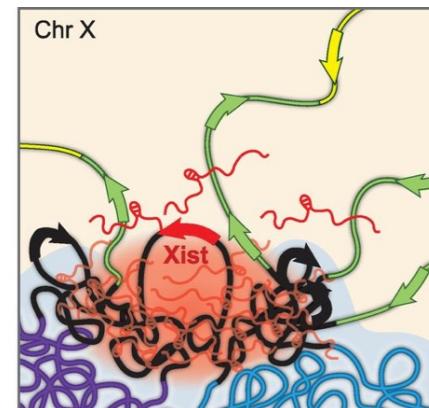
## Methyl-seq



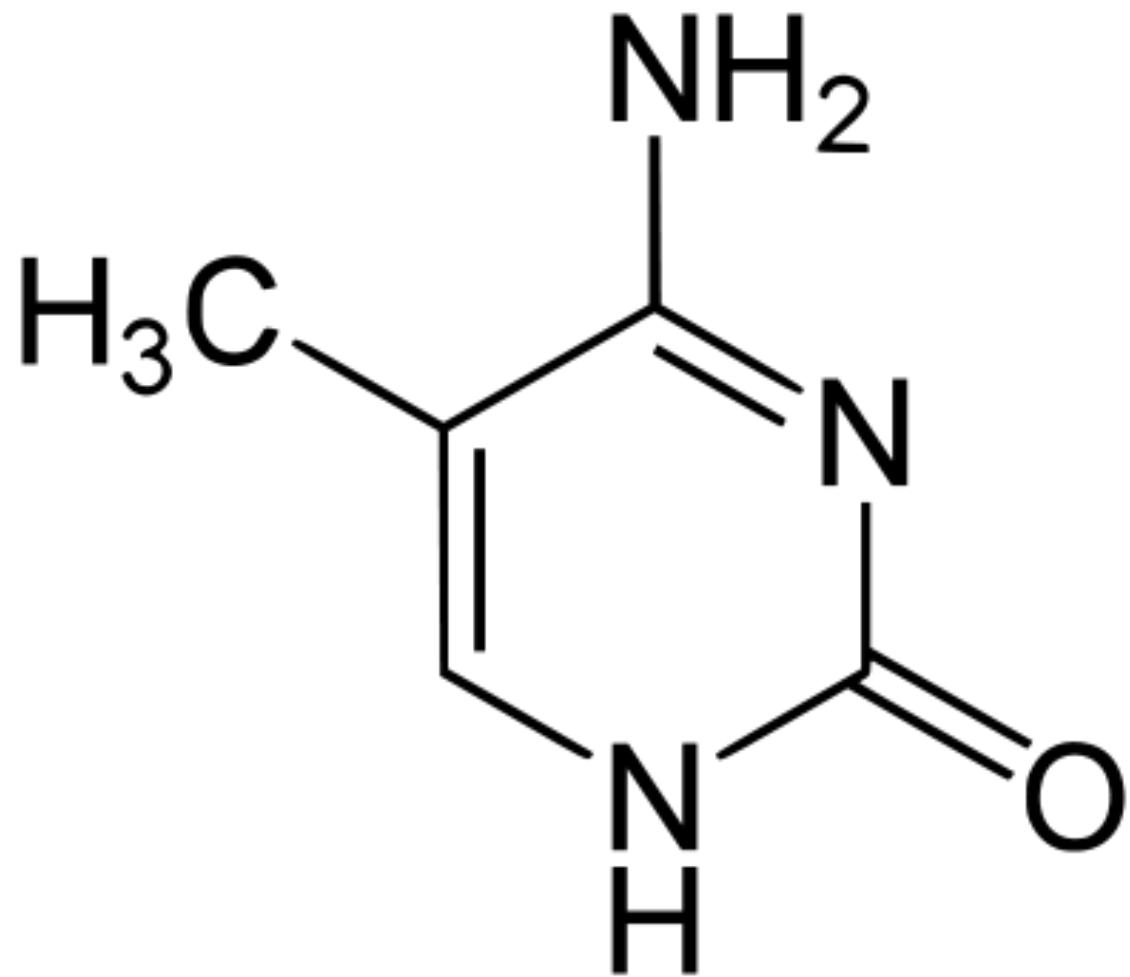
## ChIP-seq



## Hi-C

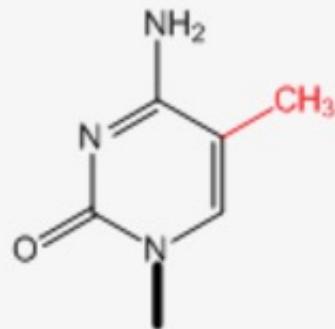


# Methyl-seq

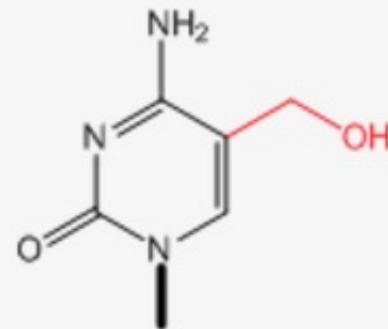


**Finding the fifth base: Genome-wide sequencing of cytosine methylation**  
Lister and Ecker (2009) *Genome Research.* 19: 959-966

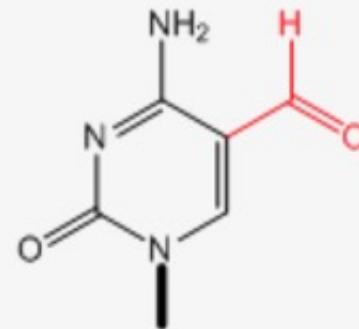
# Epigenetic Modifications to DNA



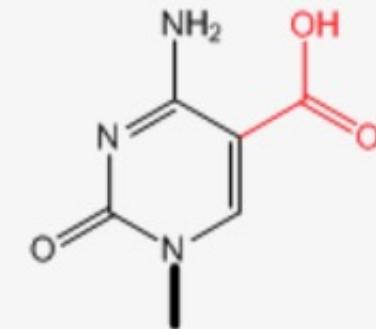
5-mC



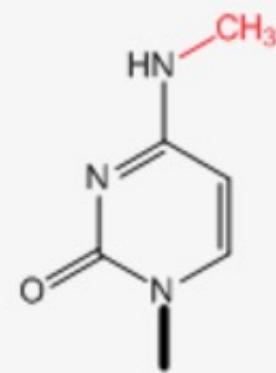
5-hmC



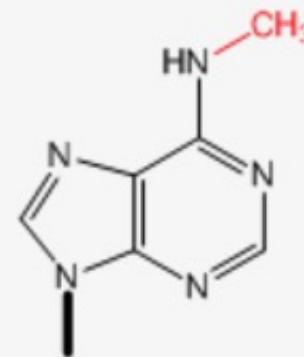
5-fC



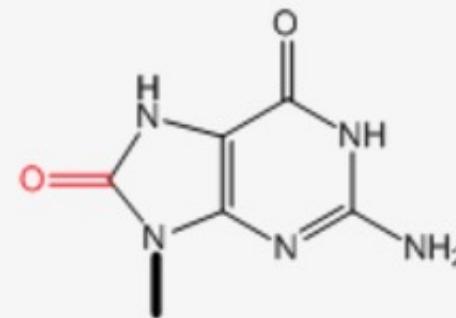
5-caC



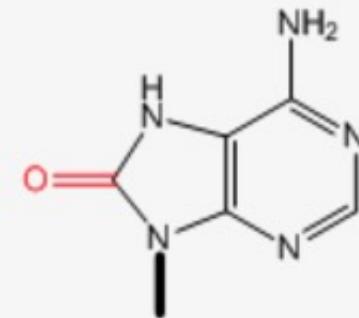
4-mC



6-mA



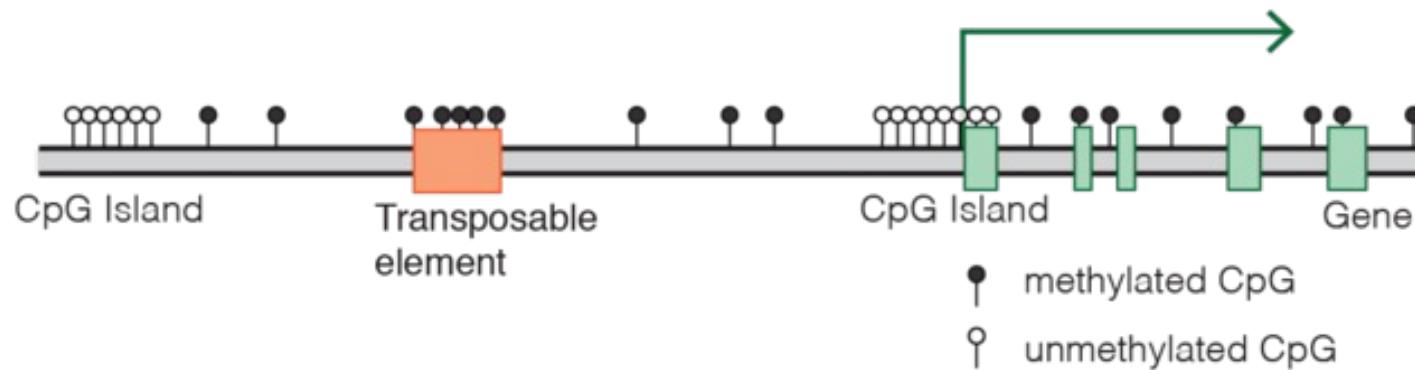
8-oxoG



8-oxoA

# Methylation of CpG Islands

Typical mammalian DNA methylation landscape



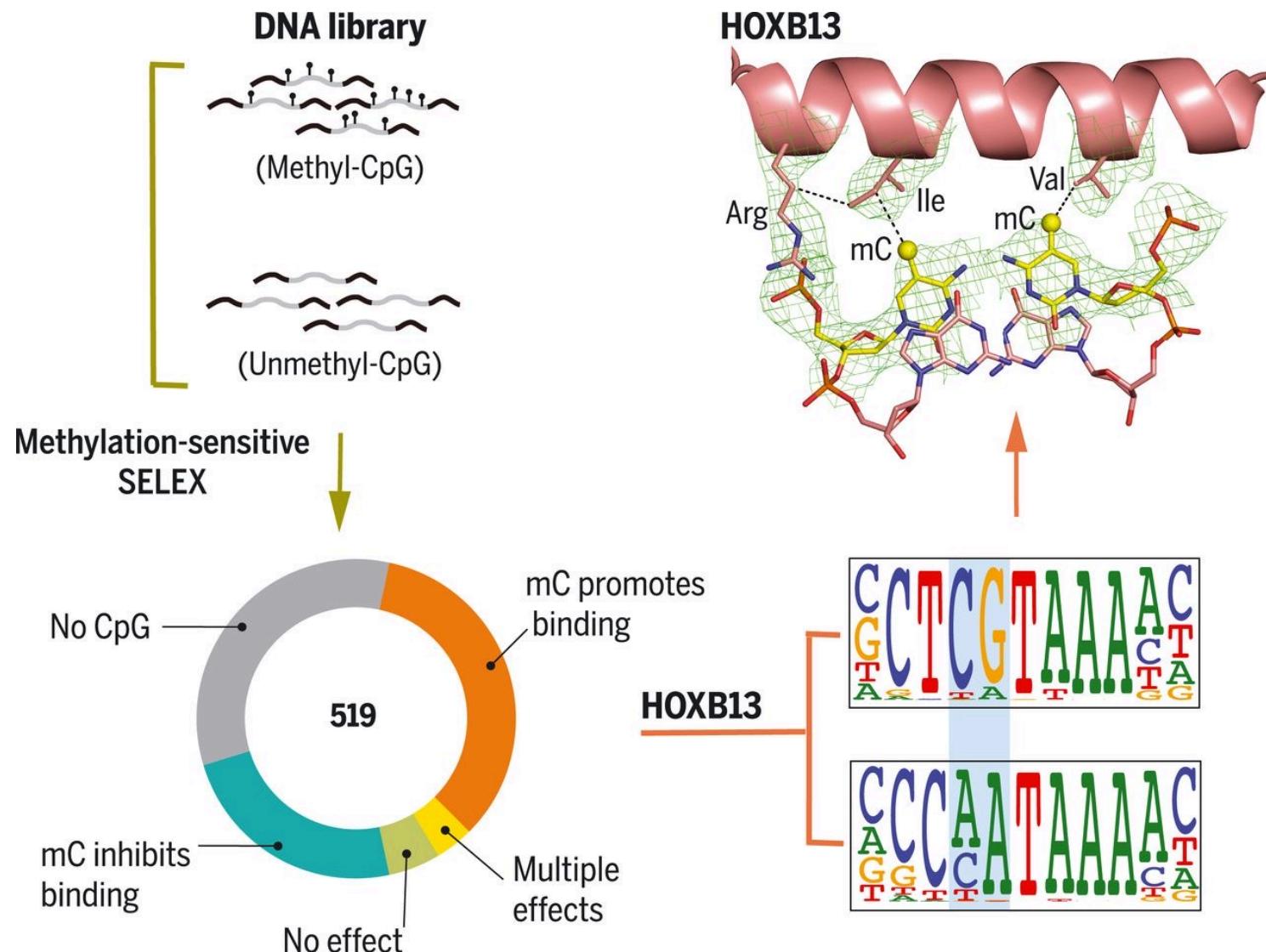
**CpG islands are (usually) defined as regions with**

- 1) a length greater than 200bp,
- 2) a G+C content greater than 50%,
- 3) a ratio of observed to expected CpG greater than 0.6

**Methylation in promoter regions correlates negatively with gene expression.**

- CpG-dense promoters of actively transcribed genes are never methylated
- In mouse and human, around 60-70% of genes have a CpG island in their promoter region and most of these CpG islands remain unmethylated independently of the transcriptional activity of the gene
- Methylation of DNA itself may physically impede the binding of transcriptional proteins to the gene
- Methylated DNA may be bound by proteins known as methyl-CpG-binding domain proteins (MBDs) that can modify histones, thereby forming compact, inactive chromatin, termed heterochromatin.

# Methylation of TF binding domains



**Impact of cytosine methylation on DNA binding specificities of human transcription factors**  
Yin et al (2017) Science. doi: 10.1126/science.aaj2239

# The Honey Bee Epigenomes: Differential Methylation of Brain DNA in Queens and Workers

Frank Lyko<sup>1</sup>\*, Sylvain Foret<sup>2</sup>\*, Robert Kucharski<sup>3</sup>, Stephan Wolf<sup>4</sup>, Cassandra Falckenhayn<sup>1</sup>, Ryszard Maleszka<sup>3</sup>\*

**1** Division of Epigenetics, DKFZ-ZMBH Alliance, German Cancer Research Center, Heidelberg, Germany, **2** ARC Centre of Excellence for Coral Reef Studies, James Cook University, Townsville, Australia, **3** Research School of Biology, the Australian National University, Canberra, Australia, **4** Genomics and Proteomics Core Facility, German Cancer Research Center, Heidelberg, Germany





**Loss of Karma transposon methylation underlies the mantled somaclonal variant of oil palm**  
Ong-Abdullah, et al (2015) *Nature*. doi:10.1038/nature15365



**Loss of Karma transposon methylation underlies the mantled somaclonal variant of oil palm**  
Ong-Abdullah, et al (2015) *Nature*. doi:10.1038/nature15365



Somaclonal variation arises in plants and animals when differentiated somatic cells are induced into a pluripotent state, but the resulting clones differ from each other and from their parents. In agriculture, somaclonal variation has hindered the micropropagation of elite hybrids and genetically modified crops, but the mechanism responsible remains unknown. The oil palm fruit 'mantled' abnormality is a somaclonal variant arising from tissue culture that drastically reduces yield, and has largely halted efforts to clone elite hybrids for oil production. Widely regarded as an epigenetic phenomenon, 'mantling' has defied explanation, but here we identify the MANTLED locus using epigenome-wide association studies of the African oil palm *Elaeis guineensis*. DNA hypomethylation of a LINE retrotransposon related to rice Karma, in the intron of the homeotic gene DEFICIENS, is common to all mantled clones and is associated with alternative splicing and premature termination. **Dense methylation near the Karma splice site (termed the Good Karma epiallele) predicts normal fruit set, whereas hypomethylation (the Bad Karma epiallele) predicts homeotic transformation, parthenocarpy and marked loss of yield.** Loss of Karma methylation and of small RNA in tissue culture contributes to the origin of mantled, while restoration in spontaneous revertants accounts for non-Mendelian inheritance. The ability to predict and cull mantling at the plantlet stage will facilitate the introduction of higher performing clones and optimize environmentally sensitive land resources.

**Loss of Karma transposon methylation underlies the mantled somaclonal variant of oil palm**  
Ong-Abdullah, et al (2015) *Nature*. doi:10.1038/nature15365

# Hypomethylation distinguishes genes of some human cancers from their normal counterparts

Andrew P. Feinberg & Bert Vogelstein

Cell Structure and Function Laboratory, The Oncology Center,  
Johns Hopkins University School of Medicine, Baltimore,  
Maryland 21205, USA

It has been suggested that cancer represents an alteration in DNA, heritable by progeny cells, that leads to abnormally regulated expression of normal cellular genes; DNA alterations such as mutations<sup>1,2</sup>, rearrangements<sup>3-5</sup> and changes in methylation<sup>6-8</sup> have been proposed to have such a role. Because of increasing evidence that DNA methylation is important in gene expression (for review see refs 7, 9-11), several investigators have studied DNA methylation in animal tumours, transformed cells and leukaemia cells in culture<sup>8,12-30</sup>. The results of these studies have varied; depending on the techniques and systems used, an increase<sup>12-19</sup>, decrease<sup>20-24</sup>, or no change<sup>25-29</sup> in the degree of methylation has been reported. To our knowledge, however, primary human tumour tissues have not been used in such studies. We have now examined DNA methylation in human cancer with three considerations in mind: (1) the methylation pattern of specific genes, rather than total levels of methylation, was determined; (2) human cancers and adjacent analogous normal tissues, unconditioned by culture media, were analysed; and (3) the cancers were taken from patients who had received neither radiation nor chemotherapy. In four of five patients studied, representing two histological types of cancer, substantial hypomethylation was found in genes of cancer cells compared with their normal counterparts. This hypomethylation was progressive in a metastasis from one of the patients.

and (3) *Hpa*II and *Hha*I cleavage sites should be present in the regions of the genes.

The first cancer studied was a grade D (ref. 43), moderately well differentiated adenocarcinoma of the colon from a 67-yr-old male. Tissue was obtained from the cancer itself and also from colonic mucosa stripped from the colon at a site just outside the histologically proven tumour margin. Figure 1 shows the pattern of methylation of the studied genes. Before digestion with restriction enzymes, all DNA samples used in the study had a size >25,000 base pairs (bp). After *Hpa*II cleavage, hybridization with a probe made from a cDNA clone of human growth hormone (HGH) showed that significantly more of the DNA was digested to low-molecular weight fragments in DNA from the cancer (labelled C in Fig. 1) than in DNA from the normal colonic mucosa (labelled N). In the hybridization conditions used, the HGH probe detected the human growth hormone genes as well as the related chorionic somatotropin

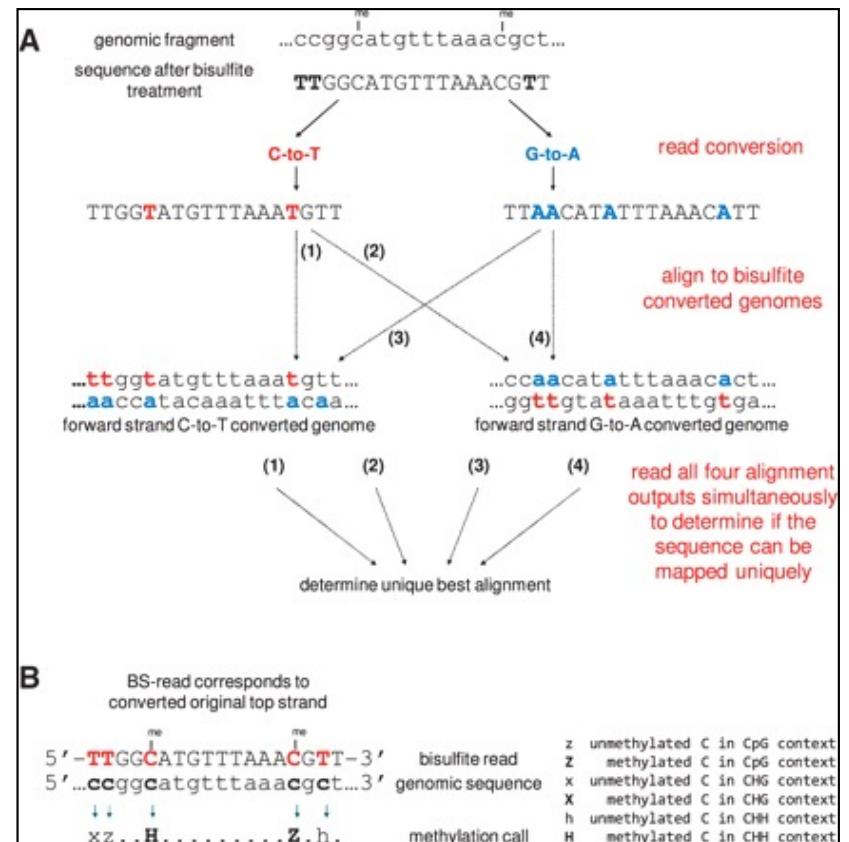
Table 1 Quantitation of methylation of specific genes in human cancers and adjacent analogous normal tissues

Patient	Carcinoma	Probe	Enzyme	% Hypomethylated fragments		
				N	C	M
1	Colon	HGH	{ <i>Hpa</i> II	<10	35	—
			{ <i>Hha</i> I	<10	39	—
			{ <i>Hpa</i> II	<10	52	—
		$\gamma$ -Globin	{ <i>Hha</i> I	<10	39	—
			{ <i>Hpa</i> II	<10	<10	—
			{ <i>Hha</i> I	<10	<10	—
2	Colon	HGH	{ <i>Hpa</i> II	<10	76	—
			{ <i>Hha</i> I	<10	85	—
		$\gamma$ -Globin	{ <i>Hpa</i> II	<10	58	—
			{ <i>Hha</i> I	<10	23	—
			{ <i>Hpa</i> II	<10	<10	—
			{ <i>Hha</i> I	<10	<10	—
3	Colon	HGH	{ <i>Hpa</i> II	<10	41	—
			{ <i>Hha</i> I	<10	38	—
		$\gamma$ -Globin	{ <i>Hpa</i> II	<10	50	—

# Bisulfite Conversion

**Treating DNA with sodium bisulfite will convert unmethylated C to T**

- 5-MethylC will be protected and not change, so can look for differences when mapping
- Requires great care when analyzing reads, since the complementary strand will also be converted (G to A)
- Typically analyzed by mapping to a “reduced alphabet” where we assume all Cs are converted to Ts once on the forward strand and once on the reverse

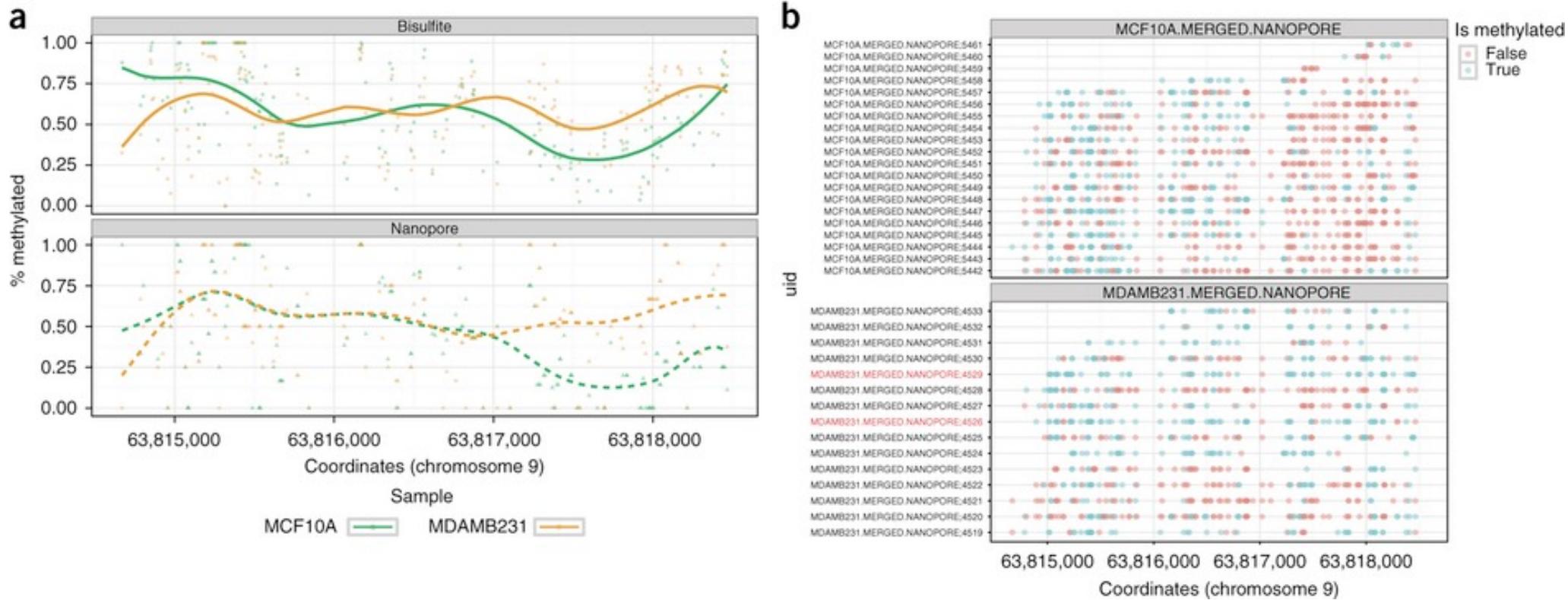


# Bisulfite Conversion



**Bismark: a flexible aligner and methylation caller for Bisulfite-Seq applications**  
Krueger and Andrews (2010) *Bioinformatics*. 27 (11): 1571-1572.

# Methylation changes in cancer detected by Nanopore Sequencing

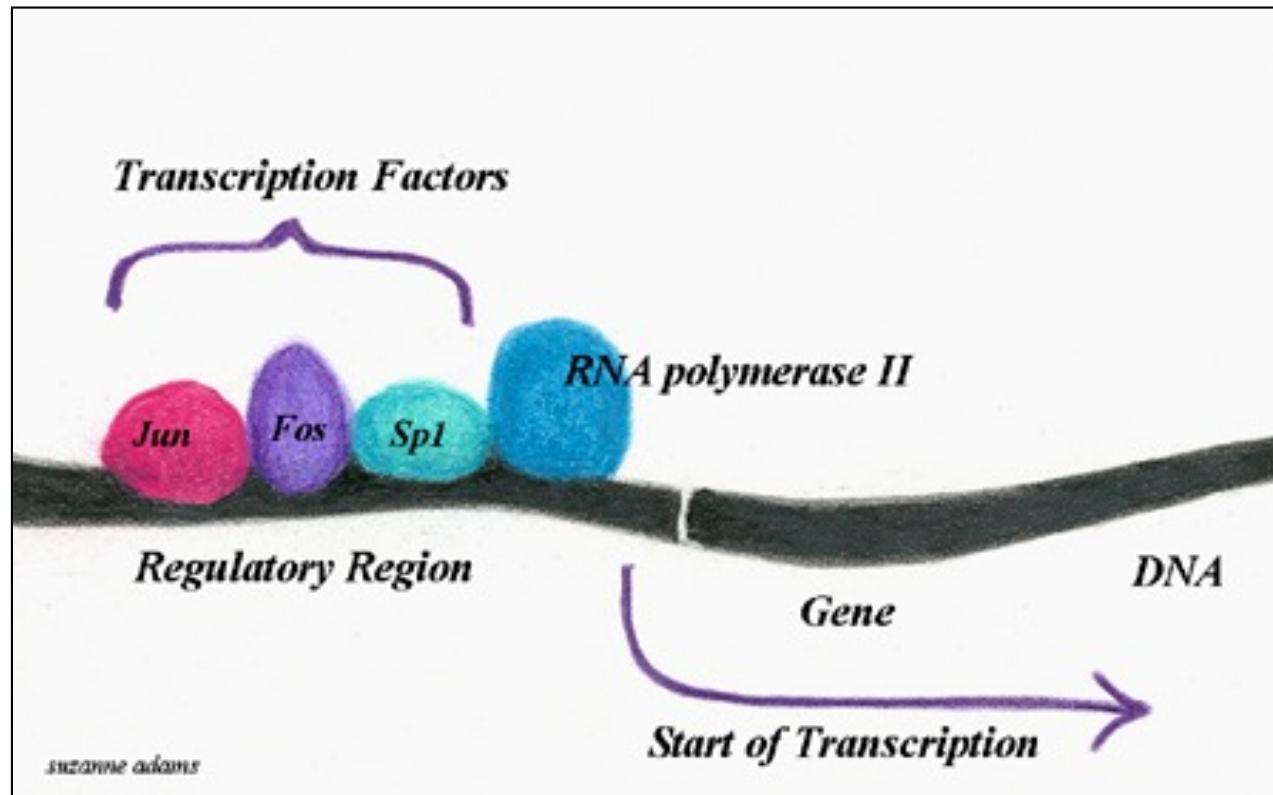


Comparison of bisulfite sequencing and nanopore-based R7.3 data in reduced representation data sets from cancer and normal cells. (a) Raw data (points) and smoothed data (lines) for methylation, as determined by bisulfite sequencing (top) and nanopore-based sequencing using an R7.3 pore (bottom), in a genomic region from the human mammary epithelial cell line MCF10A (green) and metastatic mammary epithelial cell line MDA-MB-231 (orange). (b) Same region as in a but with individual nanopore reads plotted separately. Each CpG that can be called is a point. Blue indicates methylated; red indicates unmethylated.

## Detecting DNA cytosine methylation using nanopore sequencing

Simpson, Workman, Zuzarte, David, Dursi, Timp (2017) Nature Methods. doi:10.1038/nmeth.4184

# ChIP-seq



**Genome-wide mapping of in vivo protein-DNA interactions.**

Johnson et al (2007) Science. 316(5830):1497-502

# Transcription

Transcription - YouTube

Secure https://www.youtube.com/watch?v=WsofH466lqk

Michael

Search

AUTOPLAY

Up next

Transcription and Translation: From DNA to Protein  
Professor Dave Explains 151K views

RNA polymerase reads 6:27

DNA - transcription and translation Wisam Kabaha 40K views

7:18

Transcription and mRNA processing | Biomolecules | Khan Academy 106K views

10:25

DNA transcription and translation Animation Haider abd 45K views

7:18

Translation ndsuvirtualcell 2.1M views

3:33

Transcription and Translation Overview Armando Hasudungan 611K views

13:18

DNA, Hot Pockets, & The Longest Word Ever: Crash CrashCourse 2.2M views

14:08

Transcription 1 khanacademymedicine 263K views

12:06

TRANSCRIPTION 1 KHAN ACADEMY 1:28

TRANSCRIPTION congthanhang 795K views

Moana - Best Scenes (FHD)

!行人

Transcription

2,018,430 views

4K 294

SUBSCRIBE 45K

ndsuvirtualcell

Uploaded on Jan 30, 2008

NDSU Virtual Cell Animations Project animation 'Transcription'. For more information please see <http://vcell.ndsu.edu/animations>

<https://www.youtube.com/watch?v=bKlpDtJdK8Q>

<https://www.youtube.com/watch?v=WsofH466lqk>