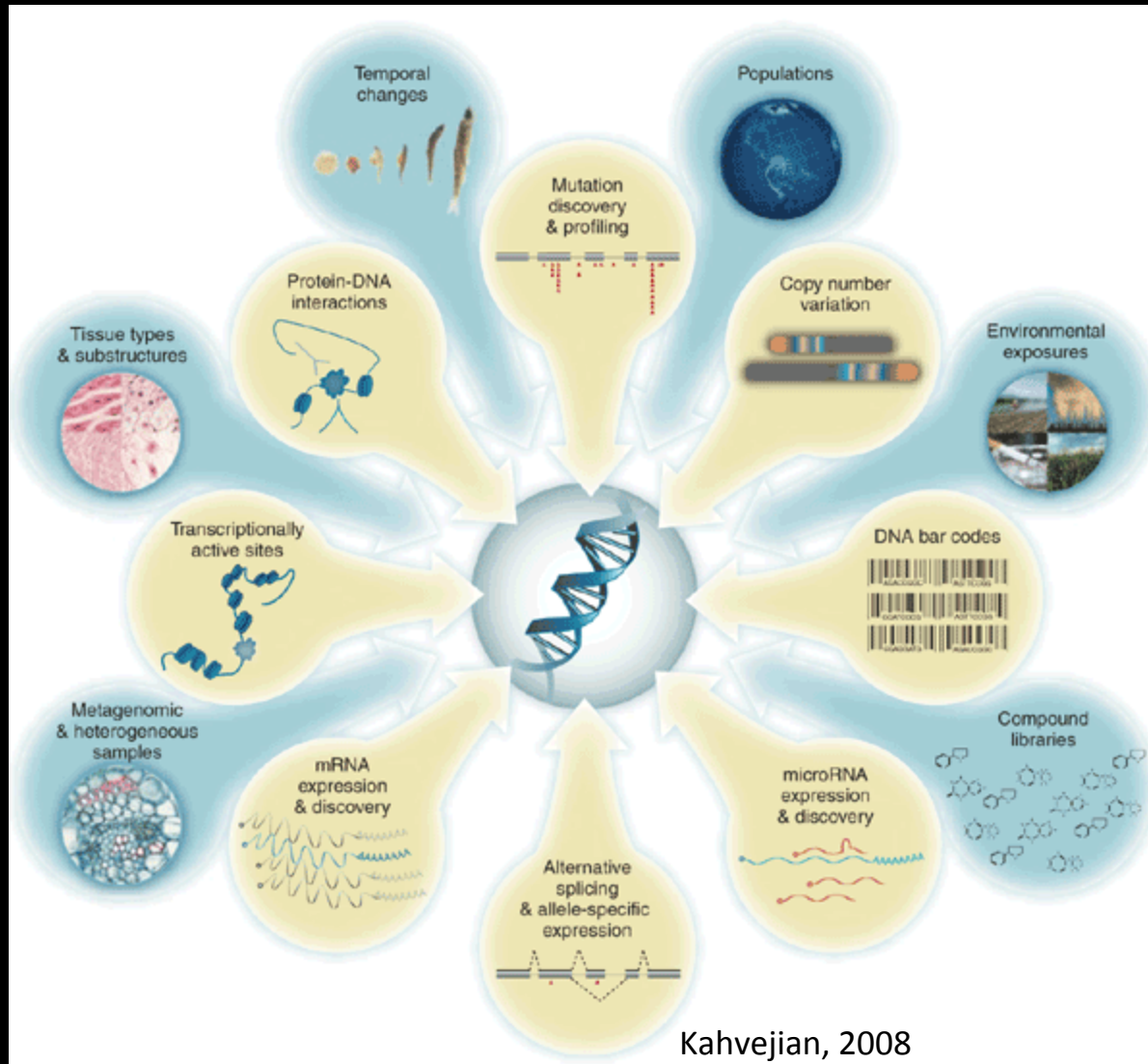




Since DNA defines the biochemical recipe for the genesis of organisms, sequencing allows us to create molecular portraits of development and disease at single-base resolution.



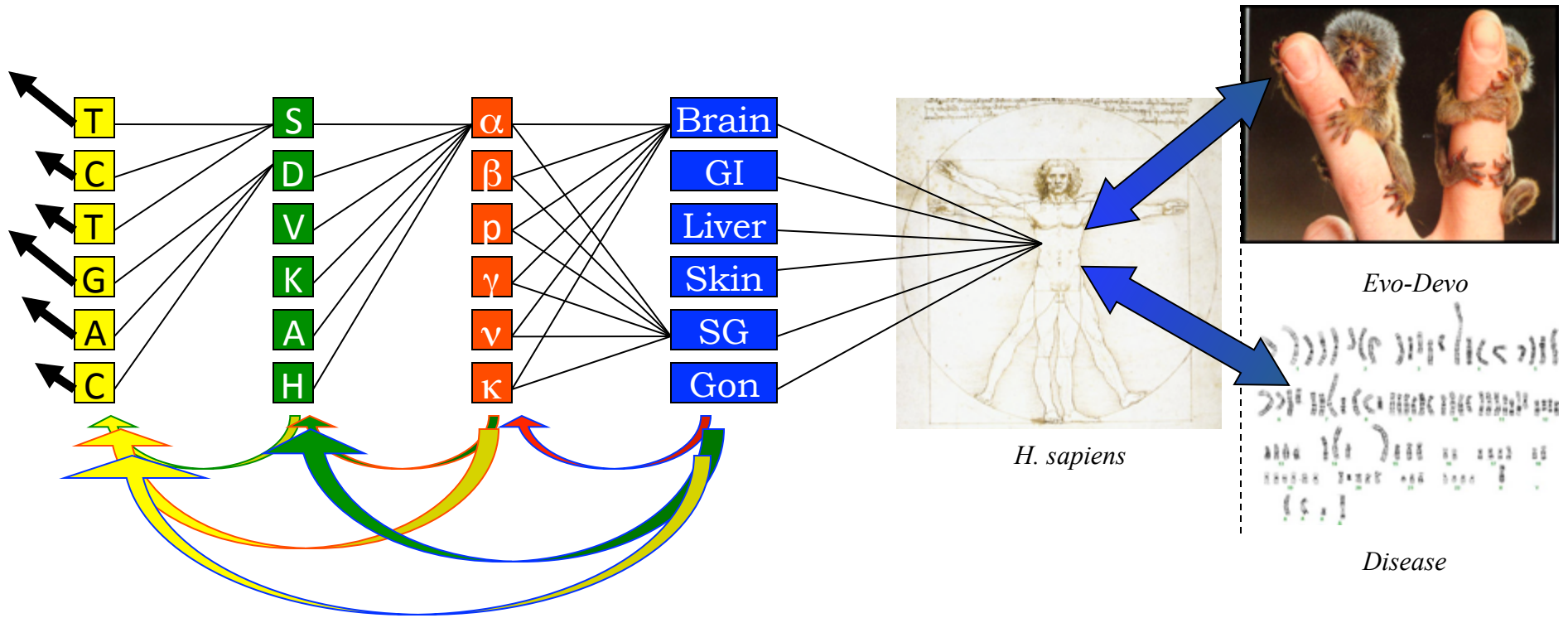
Kahvejian, 2008

PHASE TWO: INTERPRETATION

SEEBAN *the Ledger*

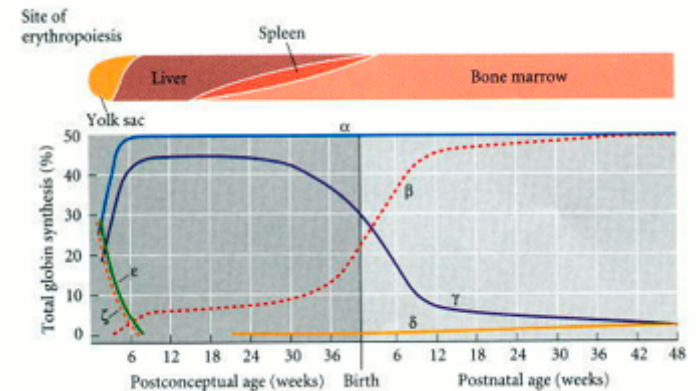
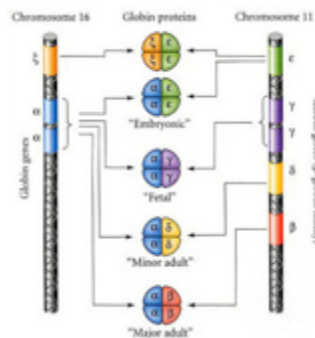


# Understanding the genome's mutation, selection, and/or drift

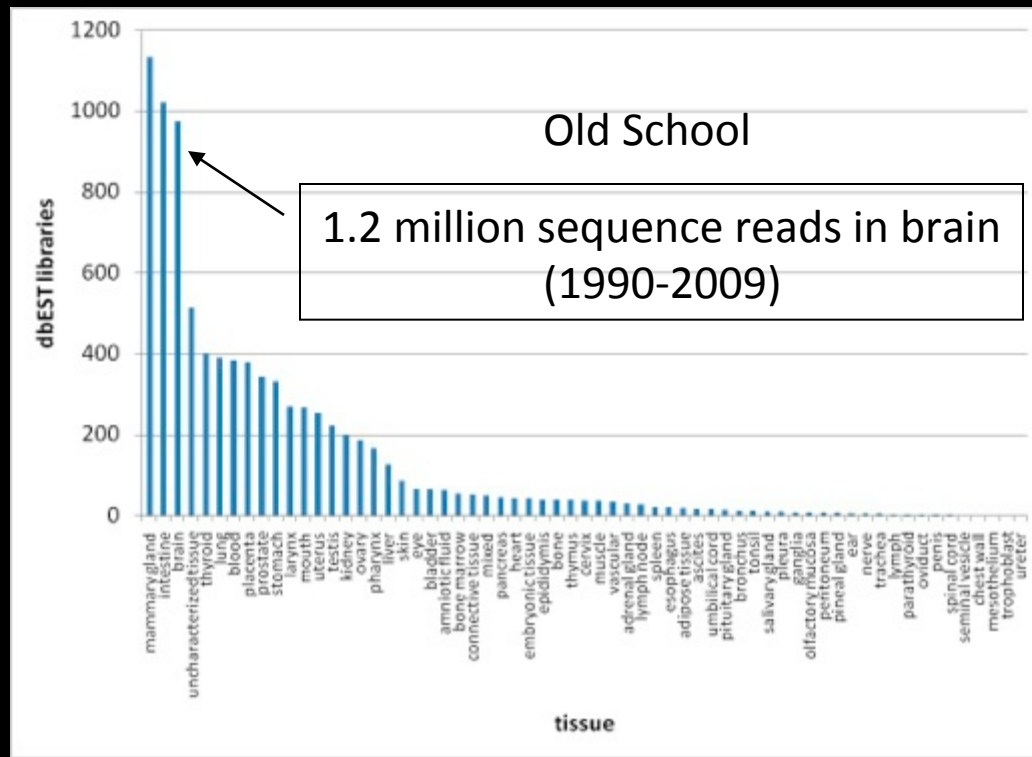


How can we understand these interactions?

Integrated, spatiotemporal molecular profiling.



# What erudition do we have now on the functional elements?



➤ Currently limited amount of EST info at NCBI

➤ EST data is expensive, time-consuming (cloning), and exhibits 3' bias.

➤ Much EST and cDNA data is for whole brains, and few libraries exist with region-specific data.

**New School:**

**One run of a NGS machine = billions of sequence reads in days**

# Description/Discussion of the Various Technologies

- The goal of the Archon X prize in Genomics is to enable a \$1,000 genome,
- Currently at \$3,000-\$50,000
- Certain platforms are better suited for certain tasks:
  - Counting applications (ChIP-Seq, RNA-Seq) need more reads
  - *De novo* assembly work needs longer reads
  - Whole genome re-sequencing requires lower errors rate and high processivity

# But, there are many options:

Platform	Library/ template preparation	NGS chemistry	Read length (bases)	Run time (days)	Gb per run	Machine cost (US\$)	Pros	Cons	Biological applications	Refs
Roche/454's GS FLX Titanium	Frag, MP/ emPCR	PS	330*	0.35	0.45	500,000	Longer reads improve mapping in repetitive regions; fast run times	High reagent cost; high error rates in homo- polymer repeats	Bacterial and insect genome <i>de novo</i> assemblies; medium scale (<3 Mb) exome capture; 16S in metagenomics	D. Muzny, pers. comm.
Illumina/ Solexa's GA <sub>II</sub>	Frag, MP/ solid-phase	RTs	75 or 100	4 <sup>†</sup> , 9 <sup>§</sup>	18 <sup>†</sup> , 35 <sup>§</sup>	540,000	Currently the most widely used platform in the field	Low multiplexing capability of samples	Variant discovery by whole-genome resequencing or whole-exome capture; gene discovery in metagenomics	D. Muzny, pers. comm.
Life/APG's SOLiD 3	Frag, MP/ emPCR	Cleavable probe SBL	50	7 <sup>†</sup> , 14 <sup>§</sup>	30 <sup>†</sup> , 50 <sup>§</sup>	595,000	Two-base encoding provides inherent error correction	Long run times	Variant discovery by whole-genome resequencing or whole-exome capture; gene discovery in metagenomics	D. Muzny, pers. comm.
Polonator G.007	MP only/ emPCR	Non- cleavable probe SBL	26	5 <sup>§</sup>	12 <sup>§</sup>	170,000	Least expensive platform; open source to adapt alternative NGS chemistries	Users are required to maintain and quality control reagents; shortest NGS read lengths	Bacterial genome resequencing for variant discovery	J. Edwards, pers. comm.
Helicos BioSciences HeliScope	Frag, MP/ single molecule	RTs	32*	8 <sup>†</sup>	37 <sup>†</sup>	999,000	Non-bias representation of templates for genome and seq-based applications	High error rates compared with other reversible terminator chemistries	Seq-based methods	91
Pacific Biosciences (target release: 2010)	Frag only/ single molecule	Real-time	964*	N/A	N/A	N/A	Has the greatest potential for reads exceeding 1 kb	Highest error rates compared with other NGS chemistries	Full-length transcriptome sequencing; complements other resequencing efforts in discovering large structural variants and haplotype blocks	S. Turner, pers. comm.

Michael  
Metzker,  
2010

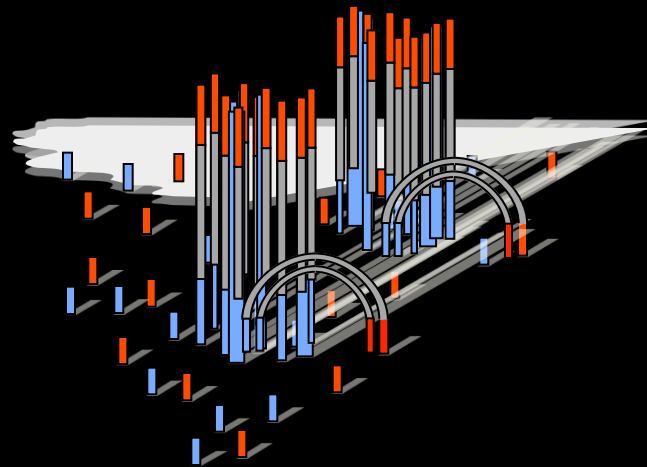
# Illumina SBS Technology

Reversible Terminator Chemistry Foundation

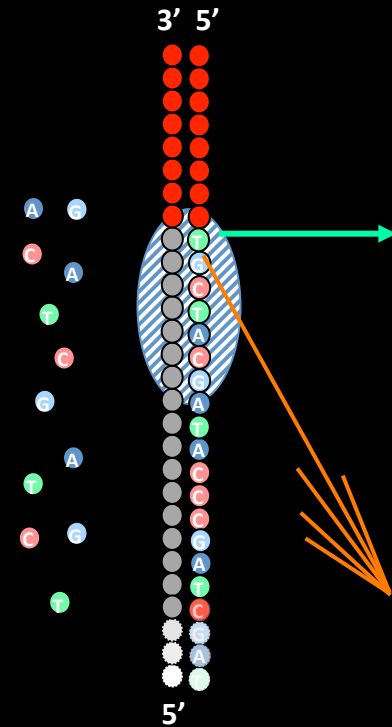
DNA  
(0.1-1.0 ug)



Sample  
preparation



Cluster growth



Sequencing

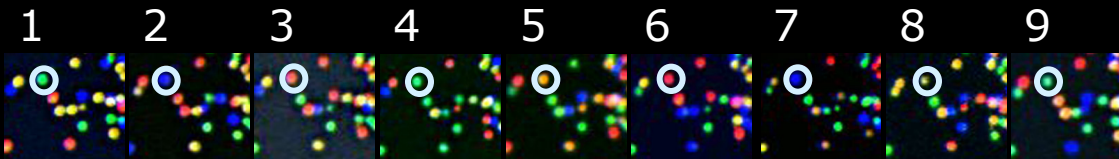


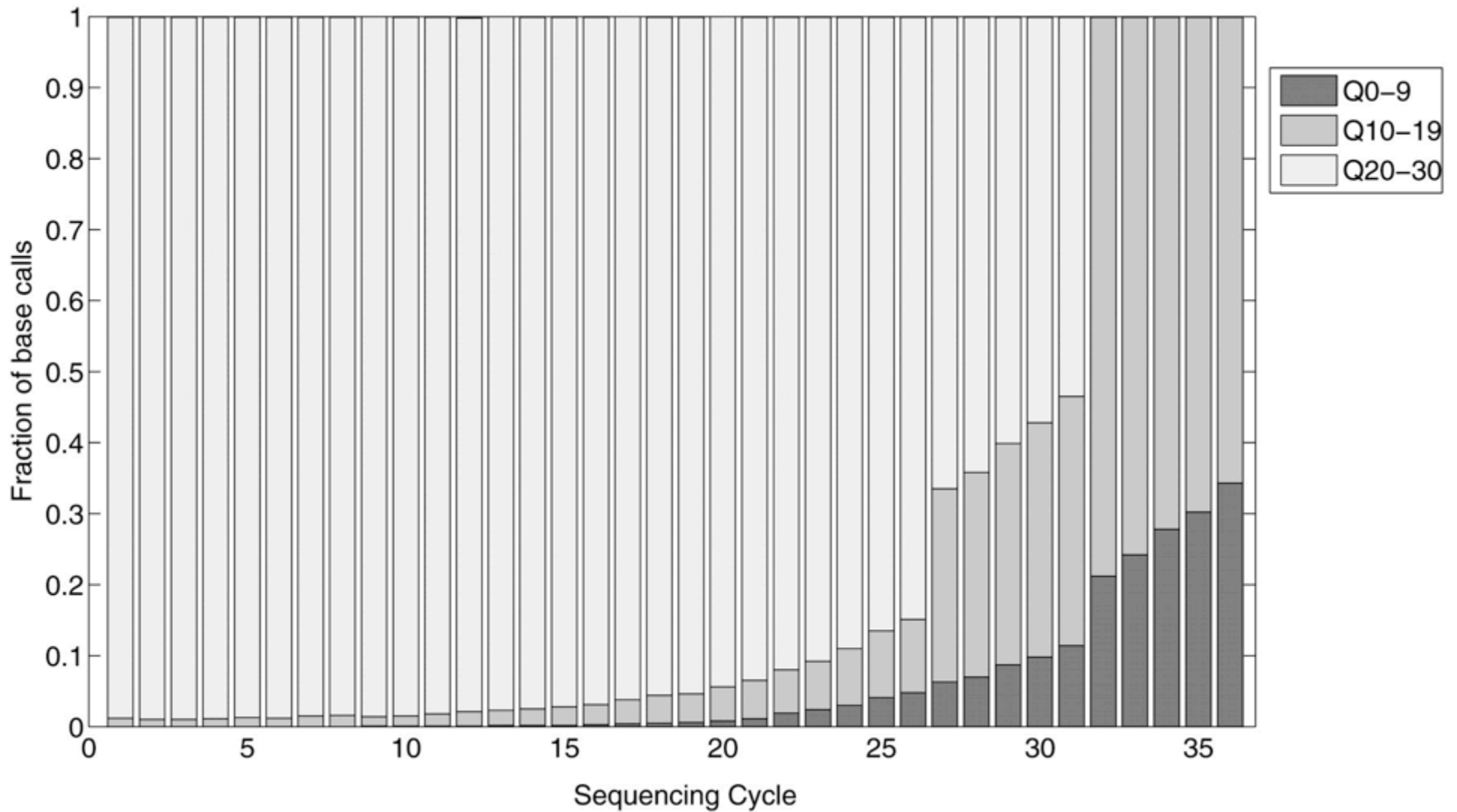
Image acquisition

T G C T A C G A T ...

Base calling



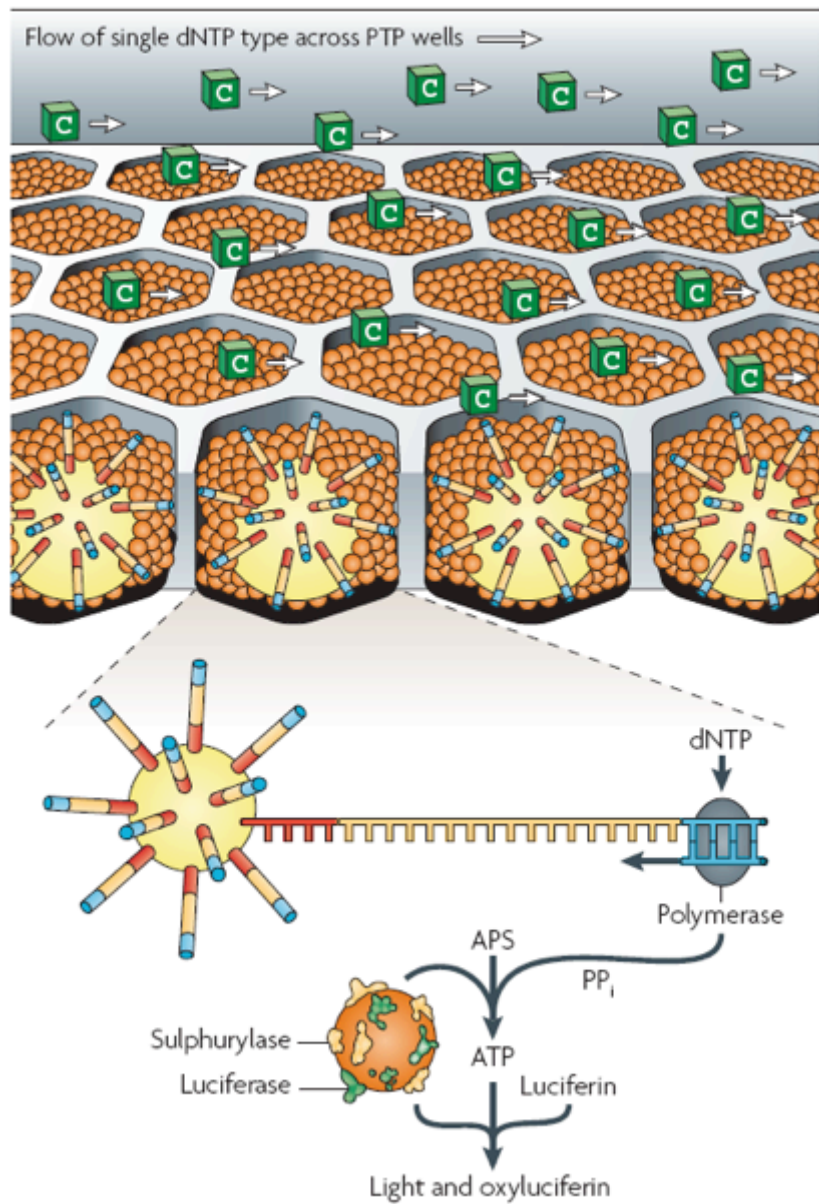
# Quality Scores vary



# Pyrosequencing

Roche/454 — Pyrosequencing

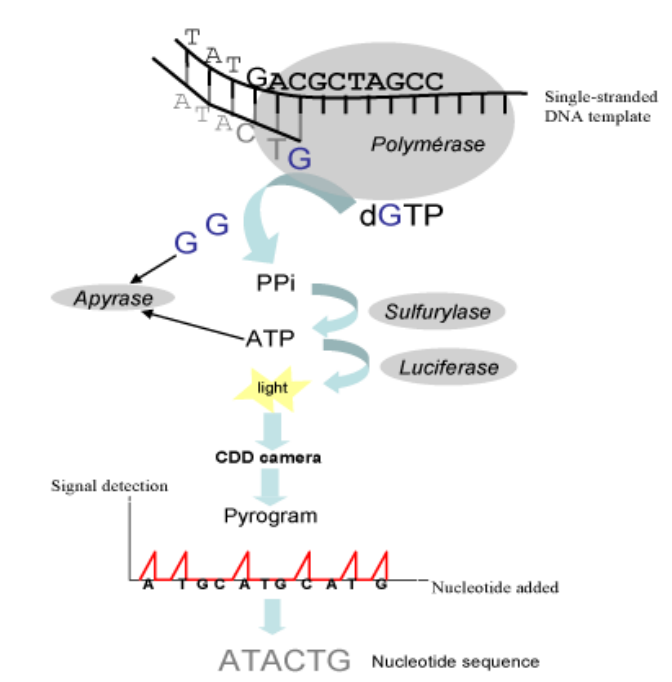
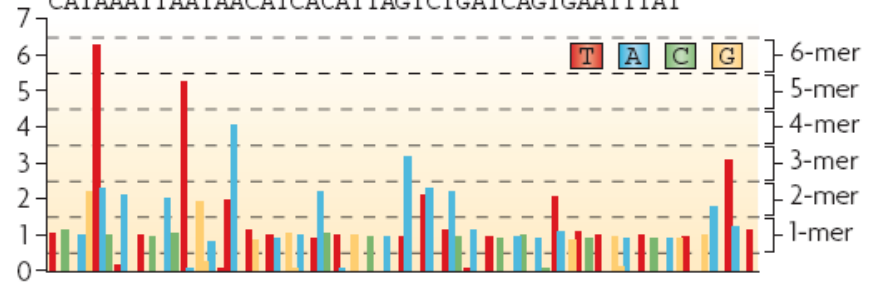
1–2 million template beads loaded into PTP wells



d

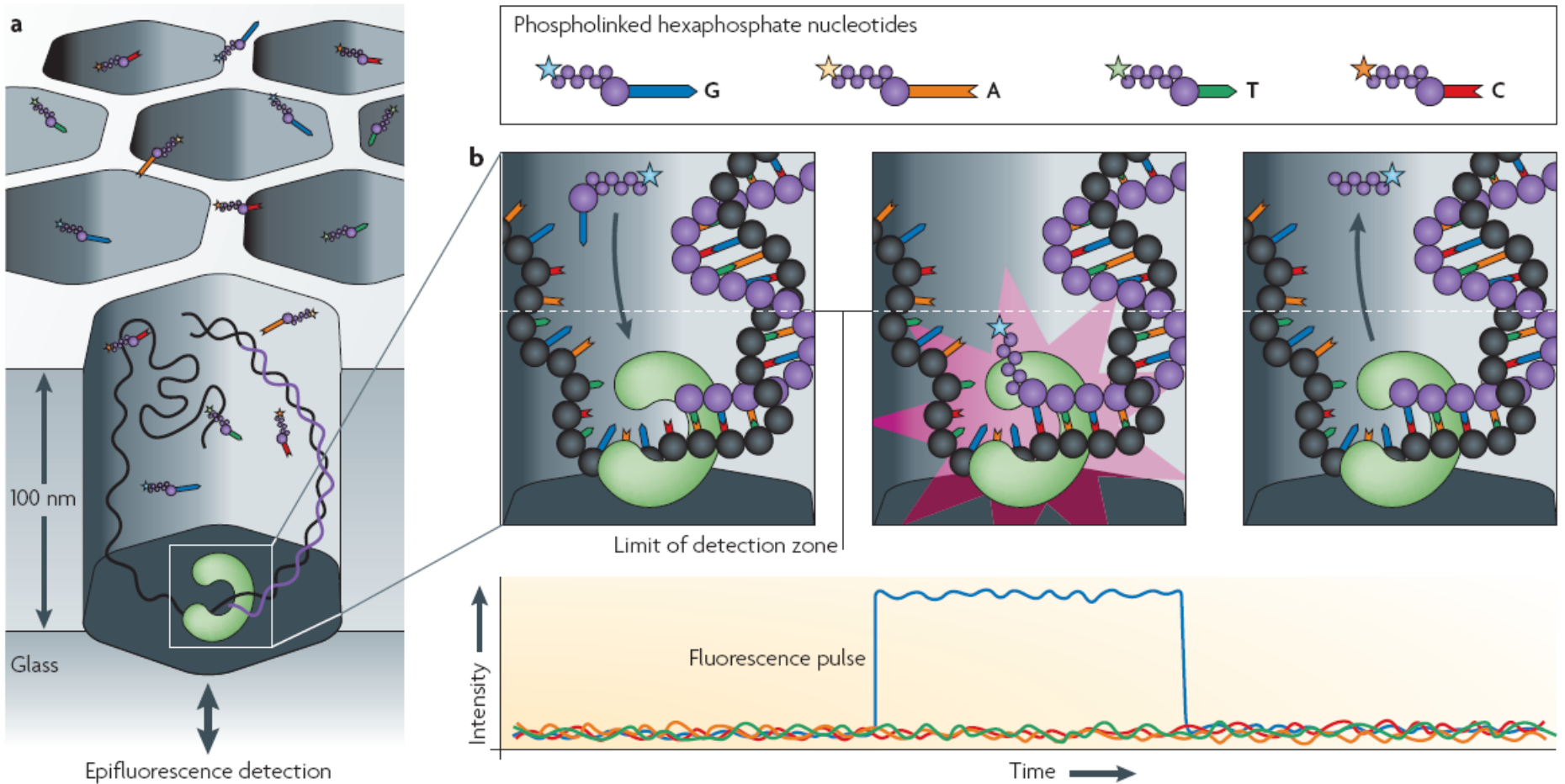
Flowgram

TCAGGTTTTTTAAACAATCAACTTTTTGGATTAAAATGTAGATAACTG  
CATAAATTAATAACATCACATTAGTCTGATCAGTGAATTTAT



# Single Molecule Real-Time

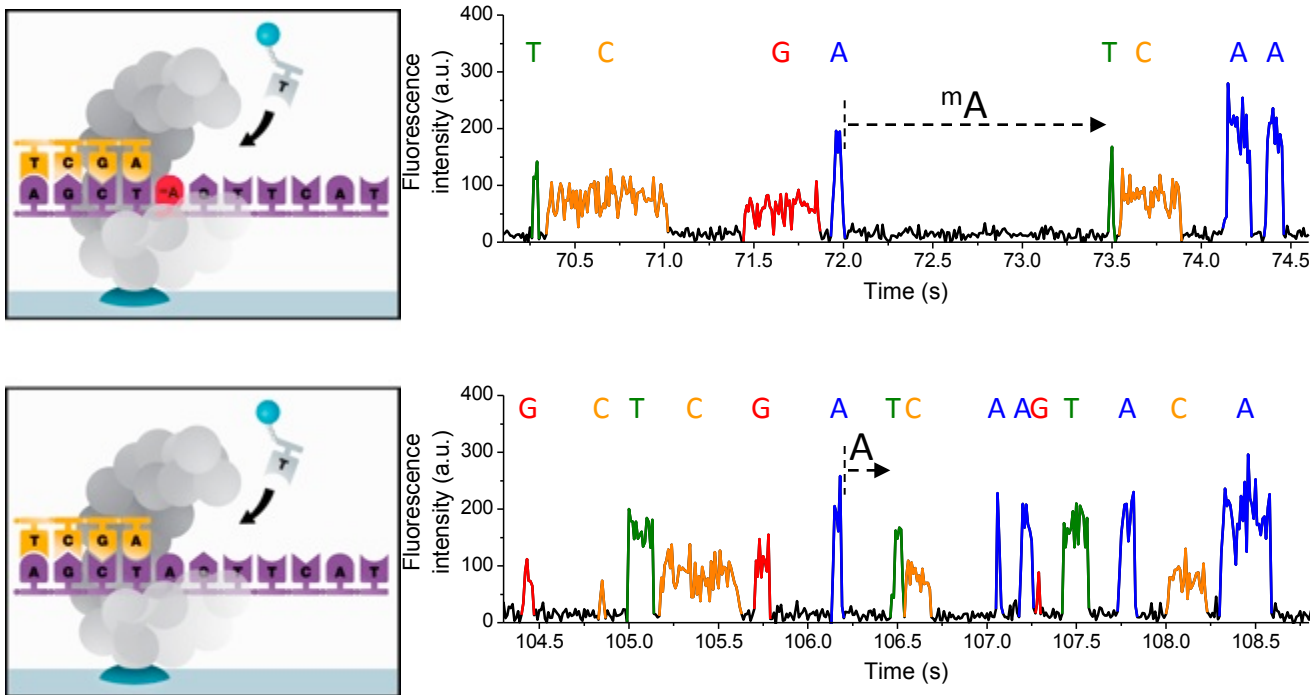
Pacific Biosciences — Real-time sequencing



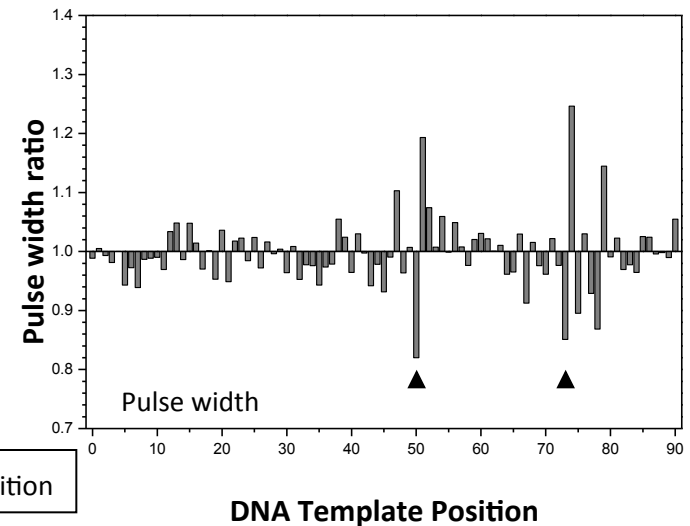
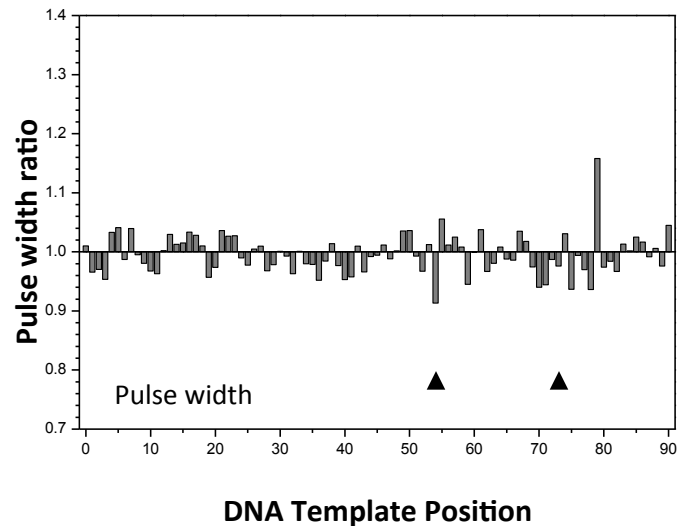
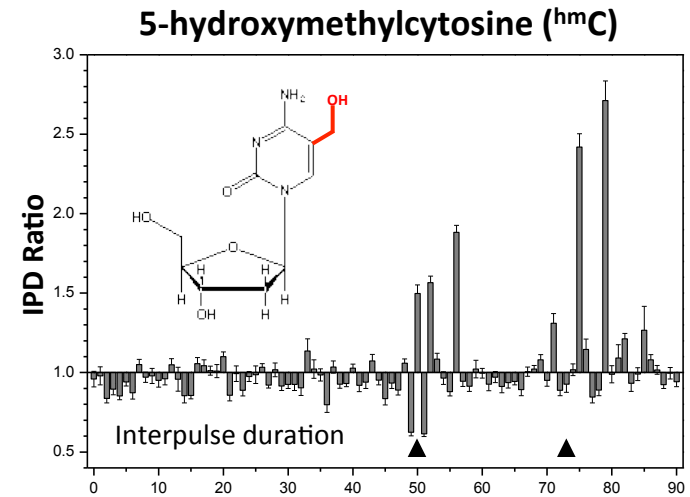
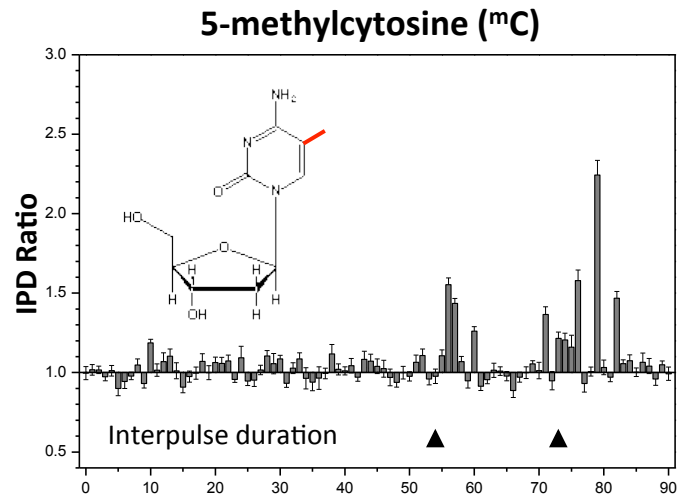
# Direct Detection of Methylation

Approach: Kinetic detection of methylated bases during SMRT DNA sequencing

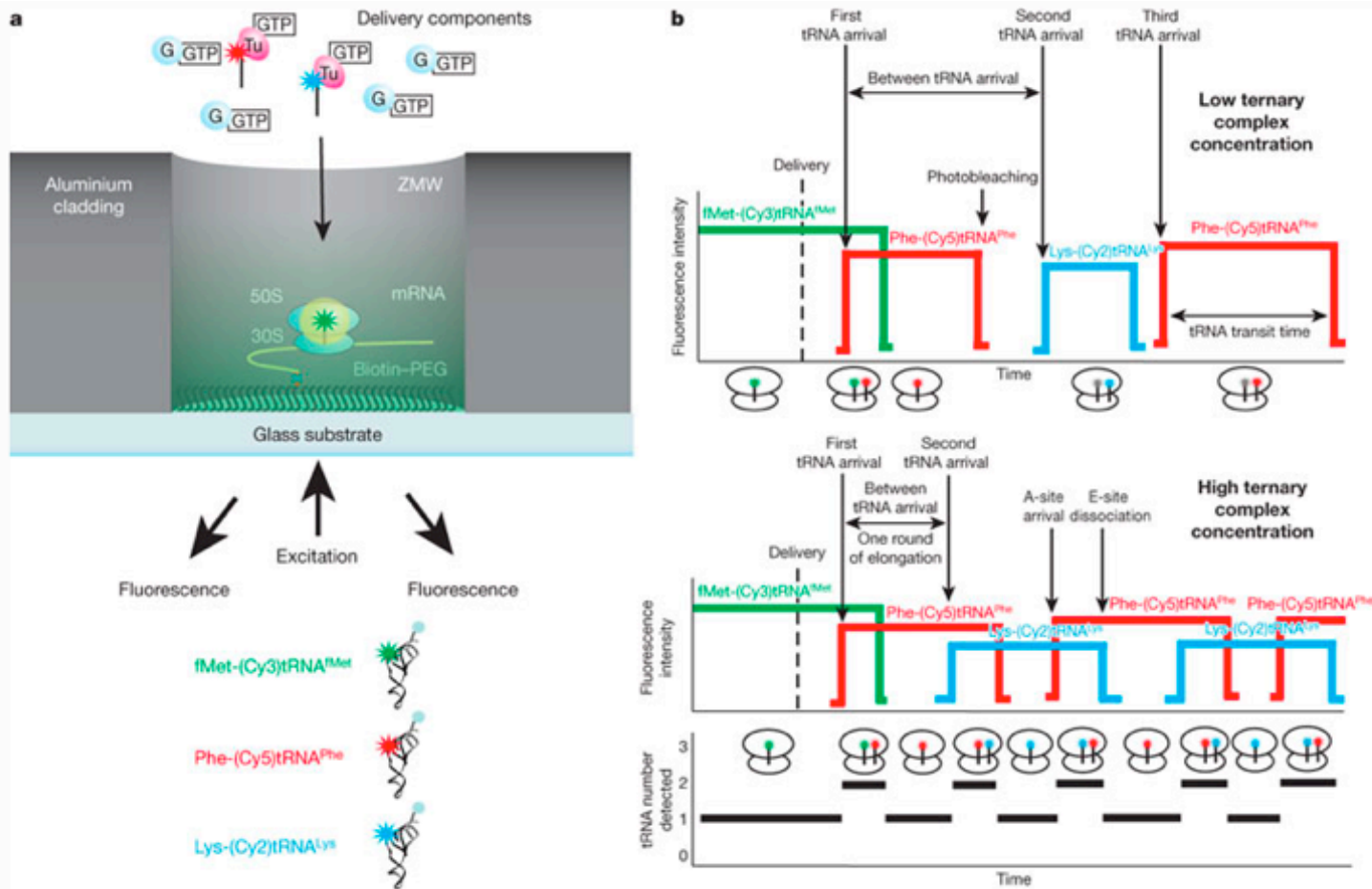
Example: N<sup>6</sup>-methyladenosine (m<sup>6</sup>A)



# detect other base modifications

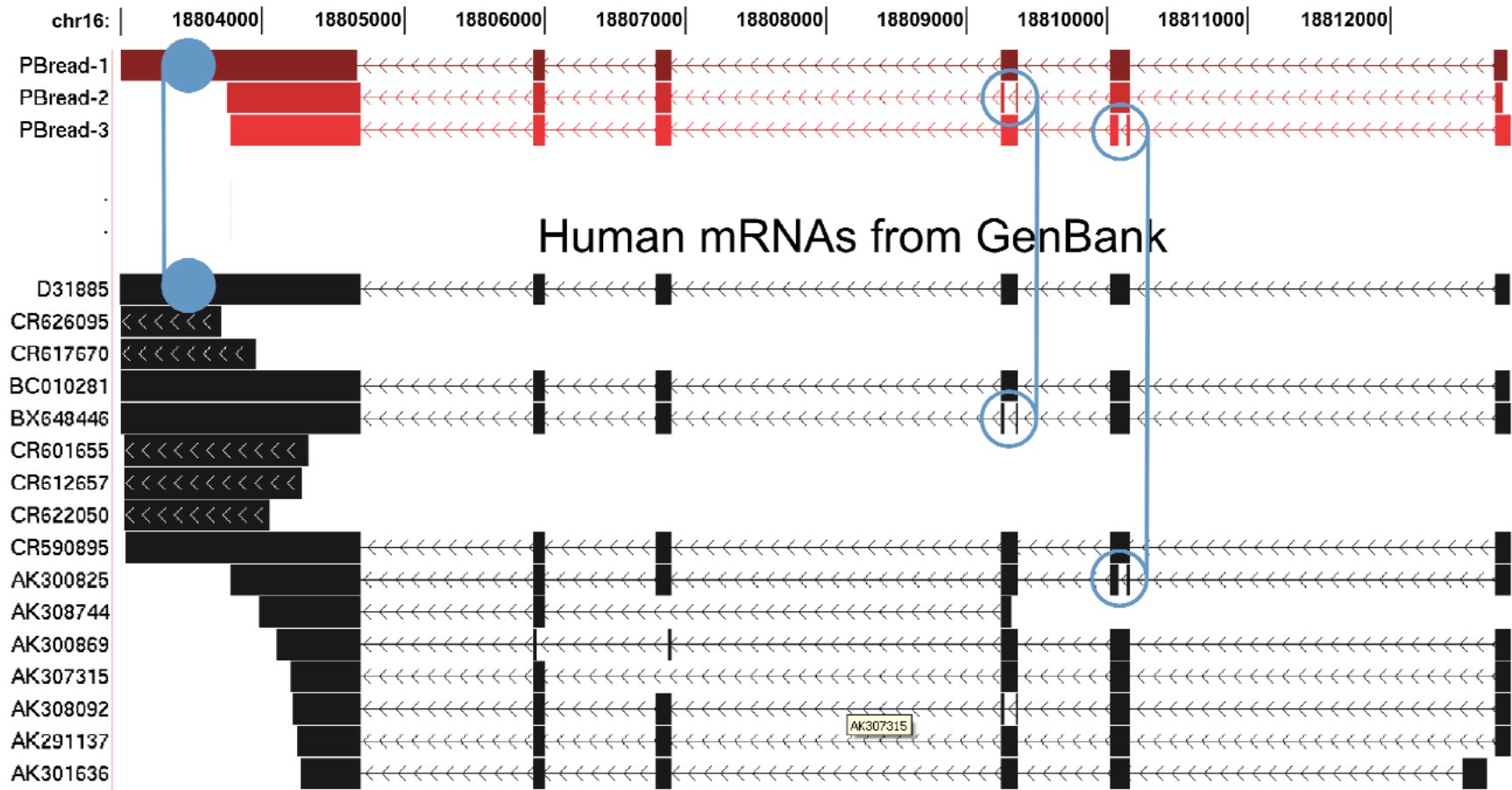


# Watch Translation in Real-Time

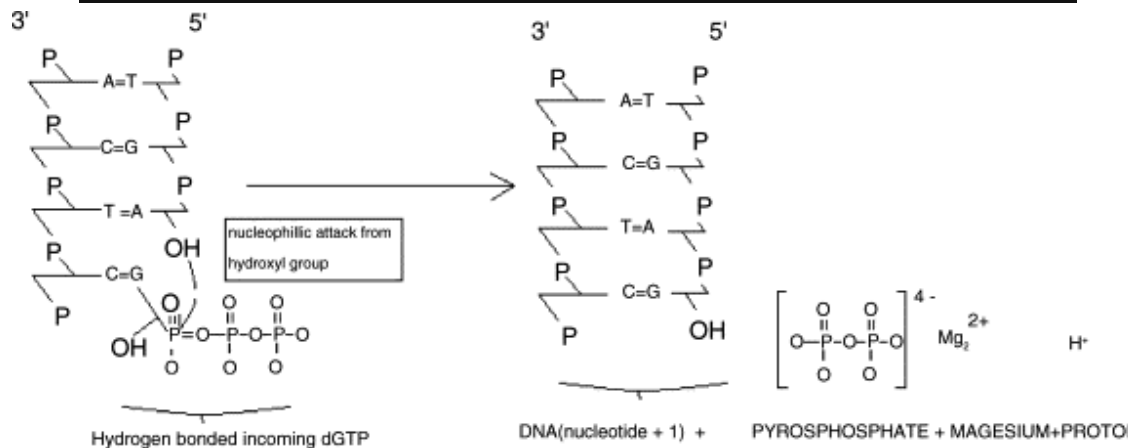
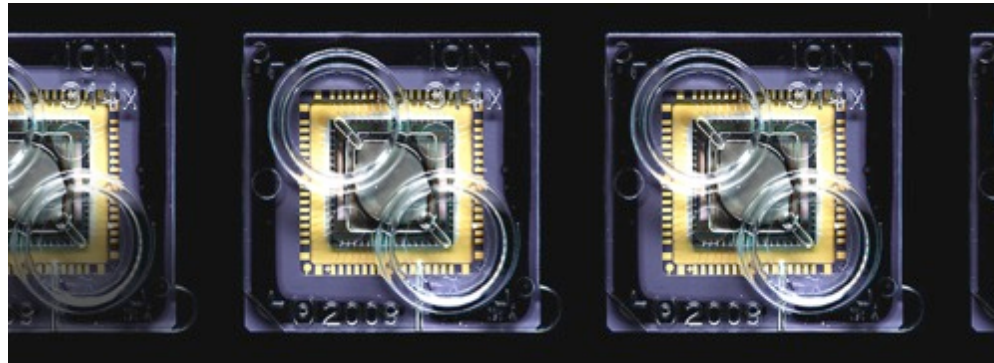
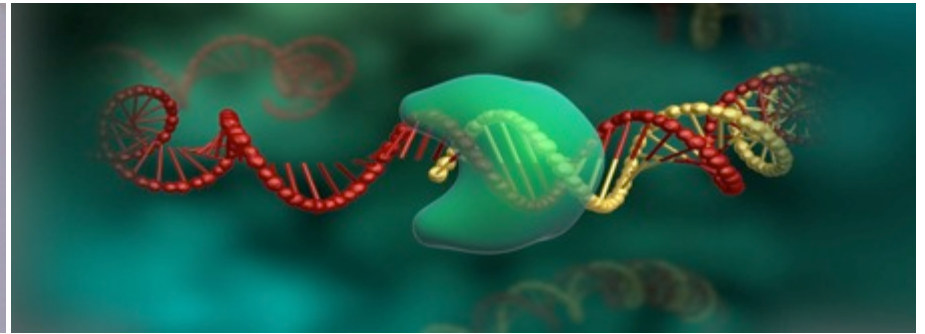
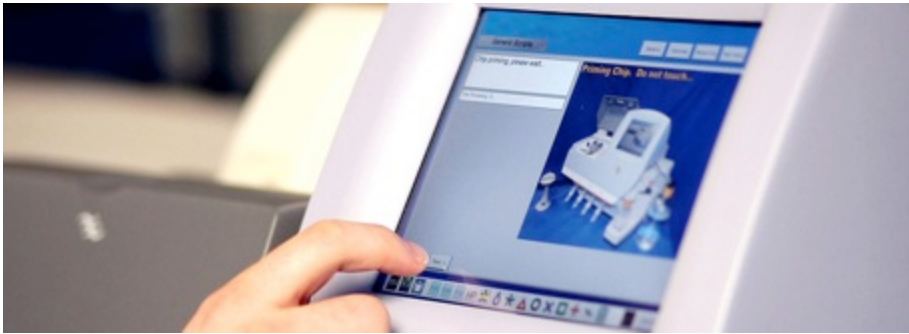


# Full-length cDNA sequencing

## PacBio Reads Mapping to ARL6IP1 mRNA Splice Variants



# “Post-Light,” Semi-Conductor Sequencing



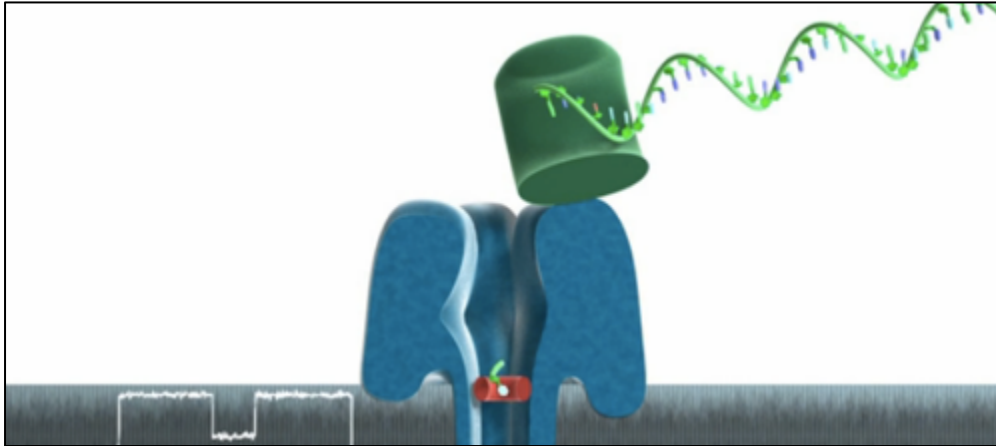
Essentially,  
11 million  
very small  
pH meters

Purushothaman *et al*, 2005  
IonTorrent, Inc.

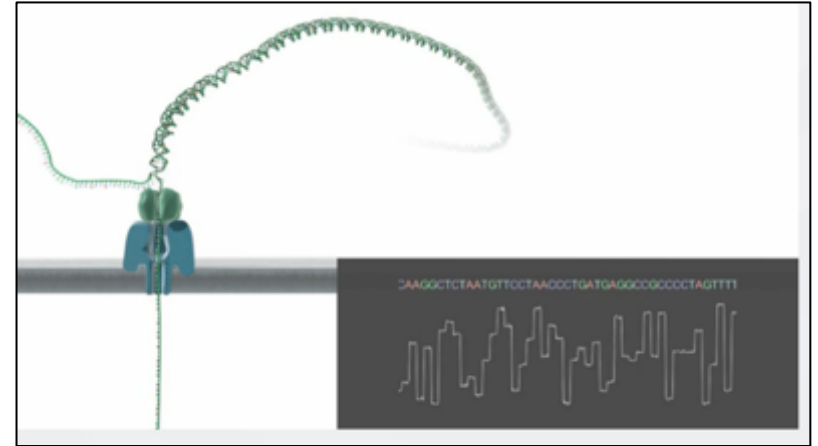




# DNA Sequencing



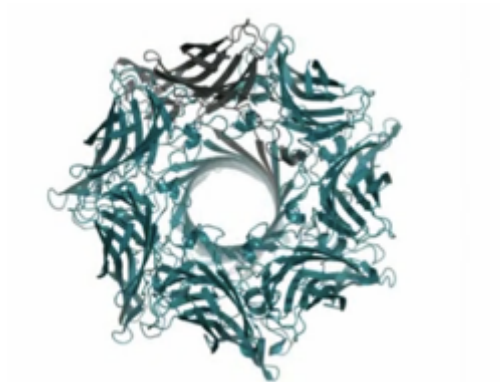
Exonuclease-Seq



Strand-Seq



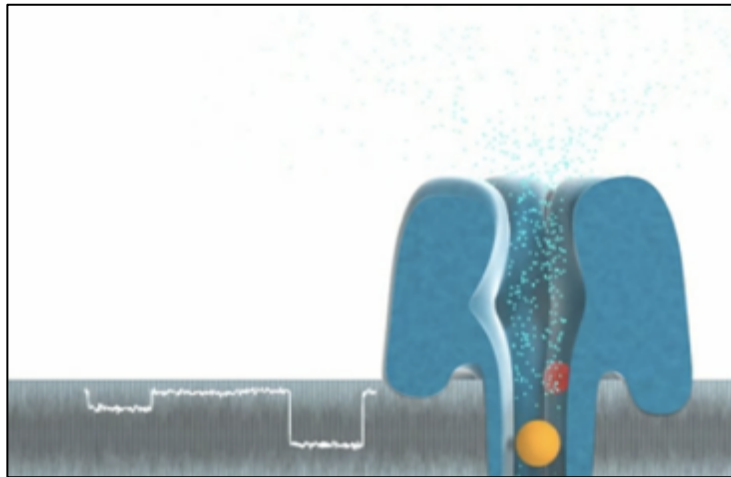
MinION



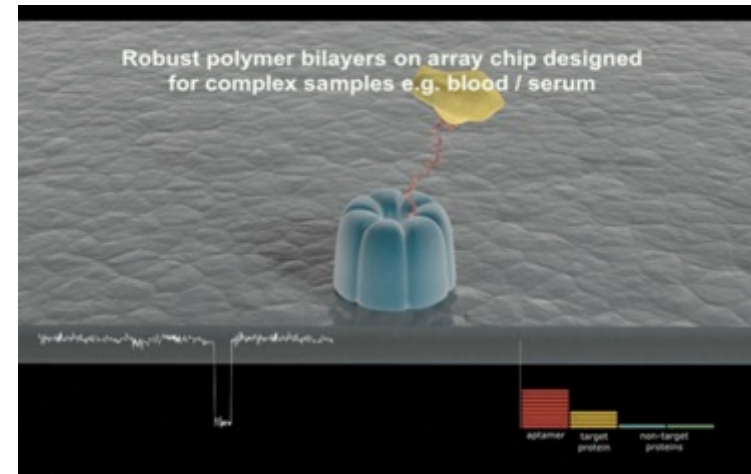
GridION



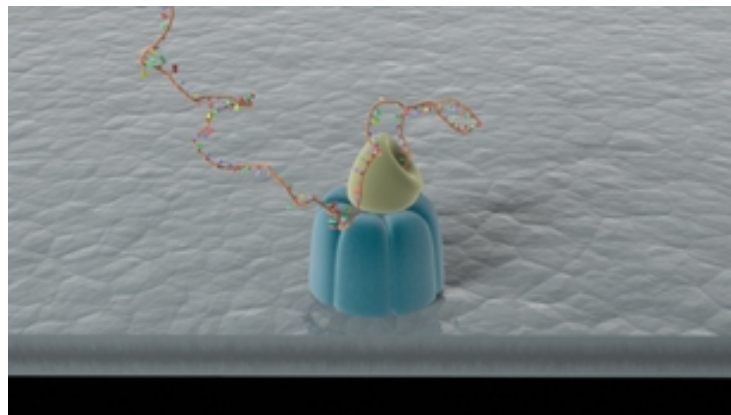
# Other (Maybe Killer) Apps



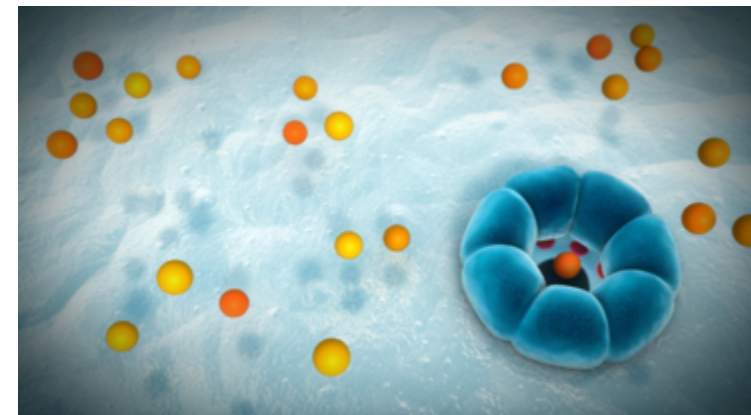
Analyte



Protein Aptamer



Direct RNA Sequencing

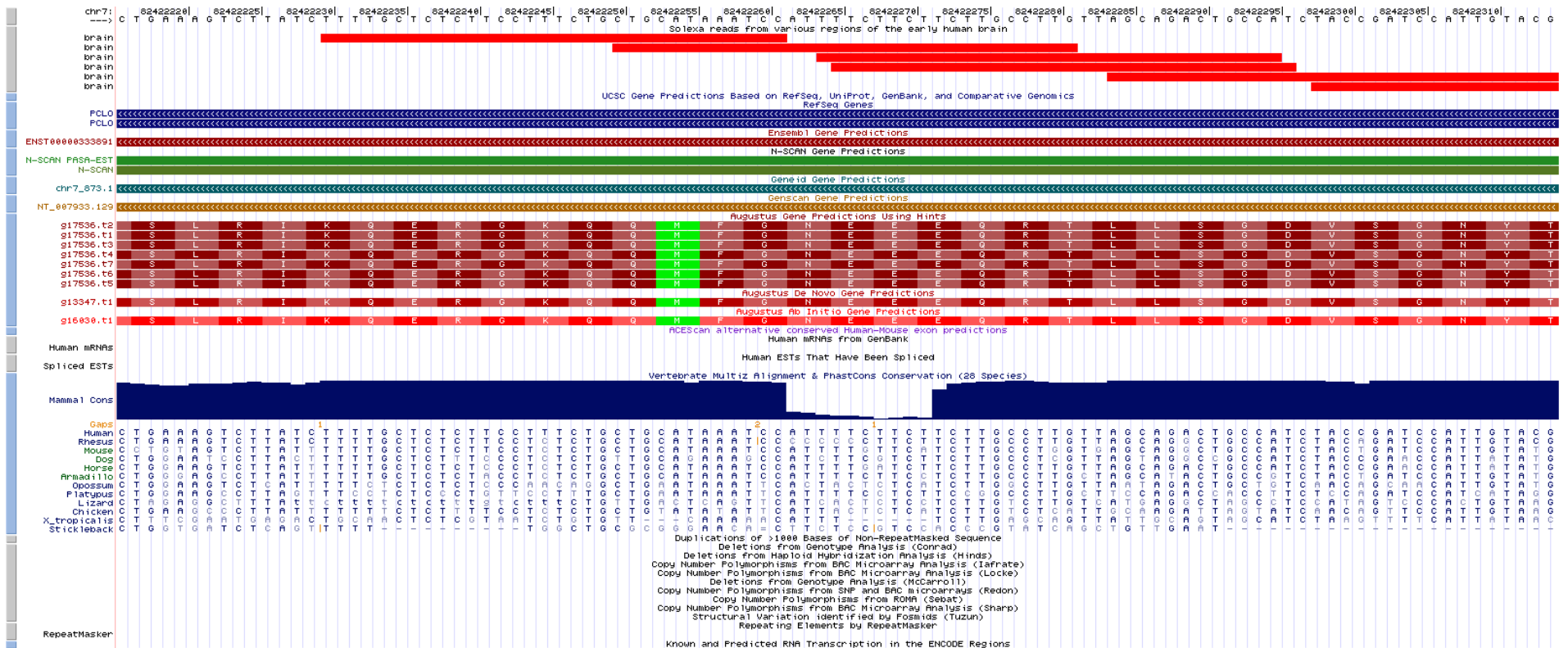


Small molecule

# Each Platform has various sources of noise, and thus Error

- De-Phasing
  - Lagging strand dephasing from incomplete extension
  - Leading strand dephasing from over-extension
- Dark Nucleotides
- Polymerase errors ( $10^{-5}$  to  $10^{-7}$ )
- Platform-specific errors
  - Illumina more likely to have error after 'G'
  - PCR-based methods miss GC- and AT-rich regions

# Alignment to the genome





# Analyzing High-Resolution Data

- **Bayesian Methods**
- Hidden Markov Models
- Permutation Testing
- Circular Binary Segmentation
- Seed-seeking
- Least Squares Regression
- Democratic Voting

# Prior

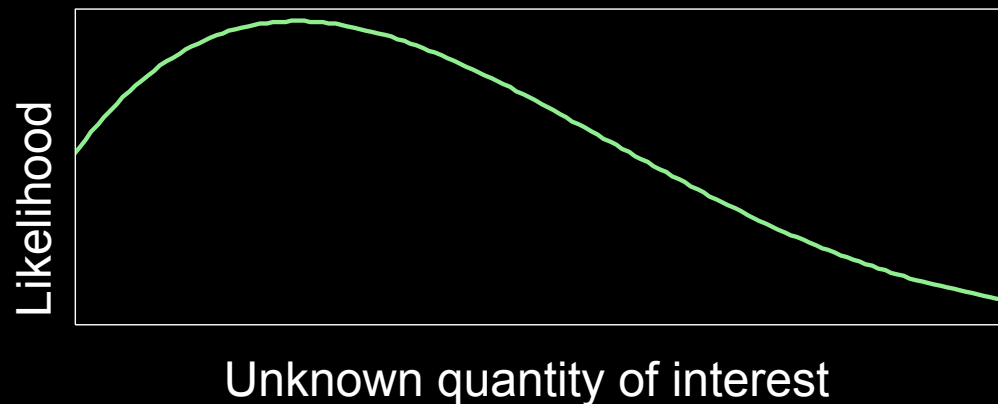
- The prior function  $\Pr(H)$  gives the probability of different possible values of the quantity of interest before the data are considered – that is, it represents the state of knowledge *prior* to the data.
- Prior may be broad or flat if we have few data (non-informative prior), or peaked if we have more information (informative prior).



(Population size, length at sexual maturity,  
haplotype frequency, model parameter)

# Likelihood

- The likelihood function  $\Pr(\text{data} | H)$  gives the probability of obtaining the data, given different possible values of the unknown quantity of interest (the “hypothesis”  $H$ ).
- The likelihood is calculated using a **statistical model**, which represents the process that produced the data. The likelihood function connects the parameters of the model to the data.

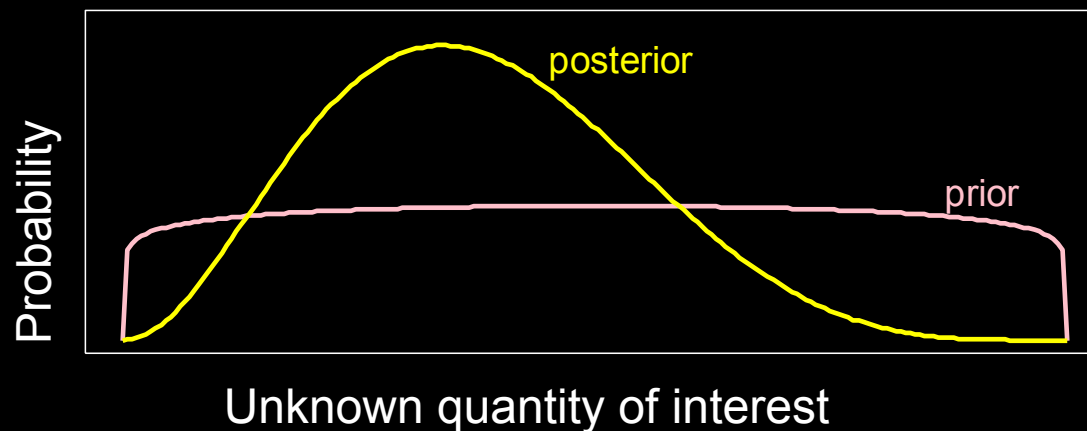


(Population size, length at sexual maturity,  
haplotype frequency, model parameter)



# Posterior

- The posterior function  $\Pr(H \mid \text{data})$  gives the probability of different possible values of the quantity of interest after the data are considered – that is, it represents the state of knowledge *posterior* to the data.
- The posterior is a combination of the prior (what we knew before) and the likelihood (what the data told us).
- The difference between the *prior* and the *posterior* indicates how much we learned from the data.



(Population size, length at sexual maturity,  
haplotype frequency, model parameter)

# Paradigm for Bayesian inference

posterior likelihood x prior

new state of knowledge information from new data x current state of knowledge

Thus, in Bayesian reasoning, new data update the current state of knowledge through Bayes' Theorem. The result is a new state of knowledge represented by the posterior.

# Bayes' Theorem

(H) The Hypothesis (the unknown) and  $x = \text{data}$

$$p(\mathbf{H} | \text{data}) = \frac{p(\text{data} | \mathbf{H}) p(\mathbf{H})}{p(\text{data})}$$

$$\text{posterior} = \text{likelihood} \times \text{prior}$$

## 2 fundamental Bayesian concepts:

1. Things that are unknown are represented by probability distributions.
2. Things that are known (data) are used to improve the knowledge of unknowns through Bayes' Theorem.

# Bayesian view of cancer

- 1% of women at 40 yrs have breast cancer (Prior)
- Mammography diagnoses 8/10 correctly (true positive rate, false negative rate)
- 10% of mammographies are false positives
- If you get a positive result, what are your odds of having breast cancer?

10,000 patients



100 patients

9,900 patients

8/10 true positives

10% false negatives



80 patients 990 patients

$$\frac{80}{1070} = 7.5\%$$

Bayes Theorem depends on the prior probability ( $\text{pr}(A)$ ):

$$\text{Pr}(A|B) = \frac{\text{Pr}(B|A) \text{Pr}(A)}{\text{Pr}(B)}$$

A = have cancer  
B = positive result

$$\frac{.8*.01}{.8*.01 + .1*.99} = 7.5\%$$

1,000,000 patients



healthy

999,999 patients

1 patient

8/10 true positives

10% false negatives



~.8 patients 9,999 patients

$$\frac{0.8}{10,000} = 0.0008\%$$



Q. What is the Bayesian Conspiracy?

A. The Bayesian Conspiracy is a multinational, interdisciplinary, and shadowy group of scientists that controls publication, grants, tenure, and the illicit traffic in grad students. The best way to be accepted into the Bayesian Conspiracy is to join the Campus Crusade for Bayes in high school or college, and gradually work your way up to the inner circles. It is rumored that at the upper levels of the Bayesian Conspiracy exist nine silent figures known only as the Bayes Council.

# Applications of Bayes to DNA-Seq

# GATK single sample genotype likelihoods

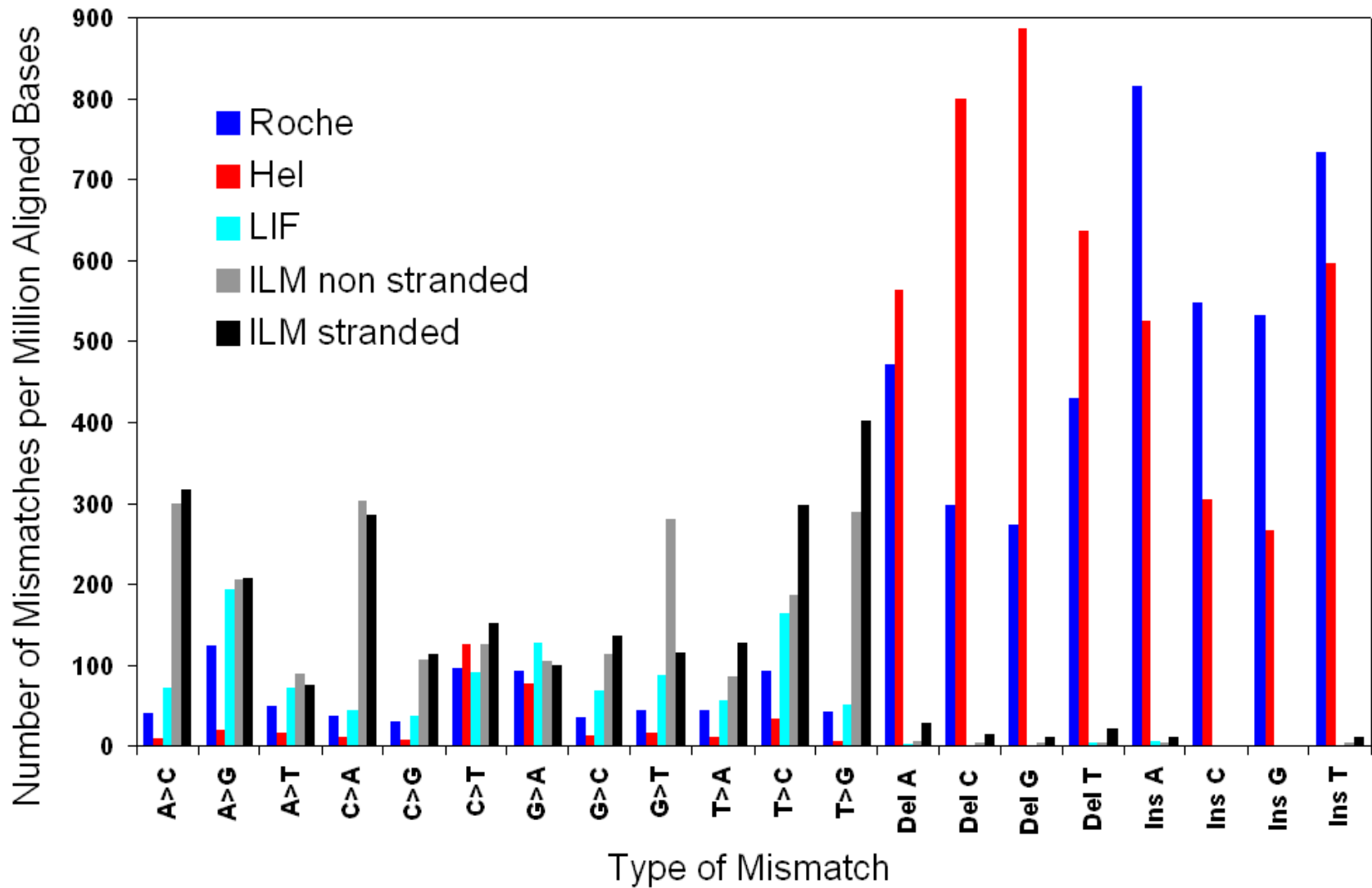
Bayesian model

Likelihood for the genotype    Prior for the genotype    Likelihood of the data given the genotype    Independent base model

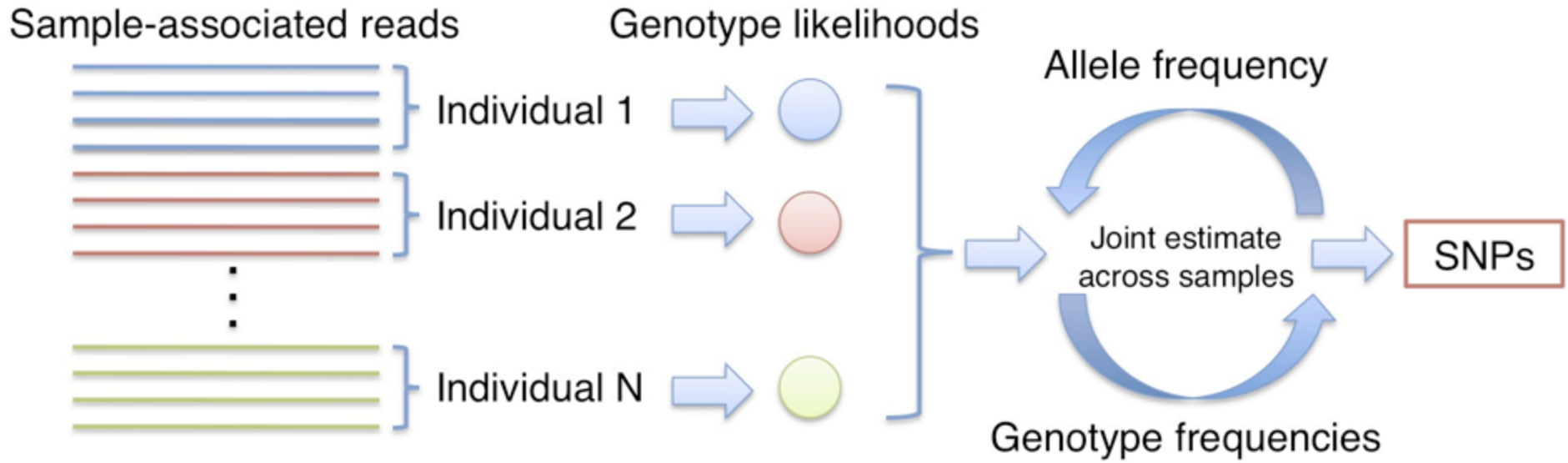
$$L(G | D) = P(G)P(D | G) = \prod_{b \in \{good\_bases\}} P(b | G)$$

- Priors applied during multi-sample calculation;  $P(G) = 1$
- Likelihood of data computed using pileup of bases and associated quality scores at given locus
- Only “good bases” are included: those satisfying minimum base quality, mapping read quality, pair mapping quality, NQS
- $P(b | G)$  uses a platform-specific confusion matrix
- $L(G | D)$  computed for all 10 genotypes

# Each Platform is slightly different, and so intrinsic errors are different



# The Broad Unified Genotyper SNP caller multiple-sample allele frequency and genotype estimates



- This approach allows us to combine weak single sample calls to discover variation among samples with high confidence

See [http://www.broadinstitute.org/gsa/wiki/index.php/Unified\\_genotyper](http://www.broadinstitute.org/gsa/wiki/index.php/Unified_genotyper) for more information

# SNIP-Seq SNP calling

For each potential variant site in the sequenced region:

Set the base quality value for each base call to the Illumina quality value

For  $k = 1, 2, \dots$

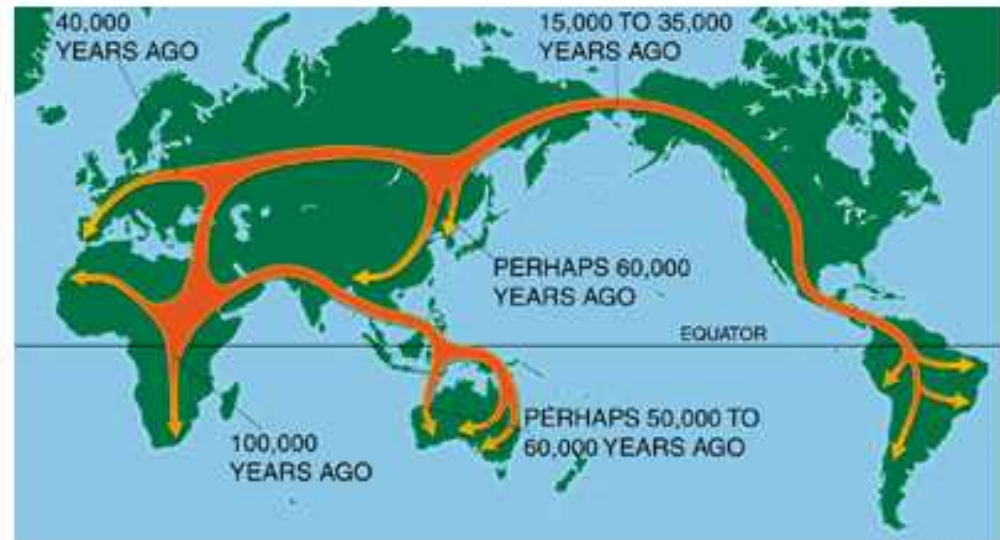
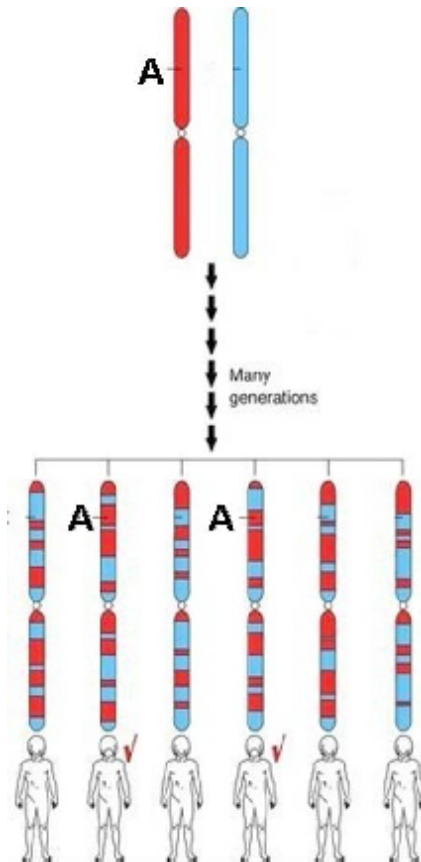
- a. Sample a genotype for each individual from the posterior distribution using a heterozygote prior of 0.001.
- b. Recalibrate the quality score for each base call using genotypes for all individuals.

If the genotype of any individual is different from the reference, identify position as a SNP.

Sample a genotype for each individual from the posterior distribution computed using a heterozygote prior of 0.2.

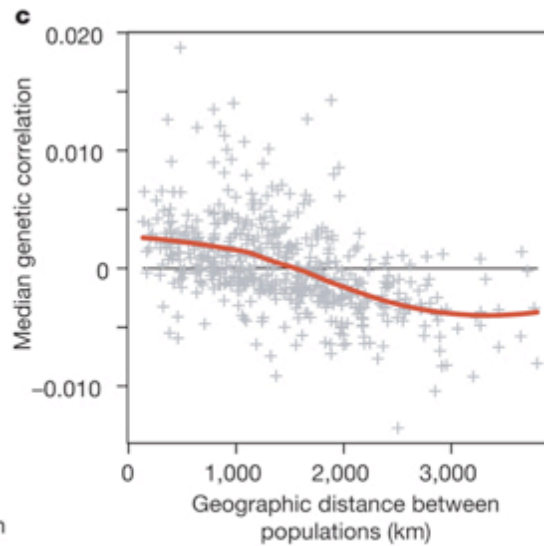
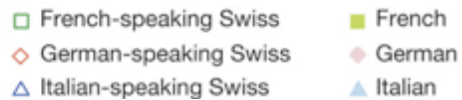
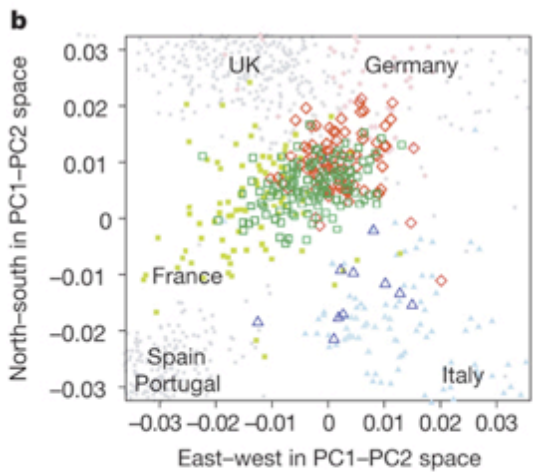
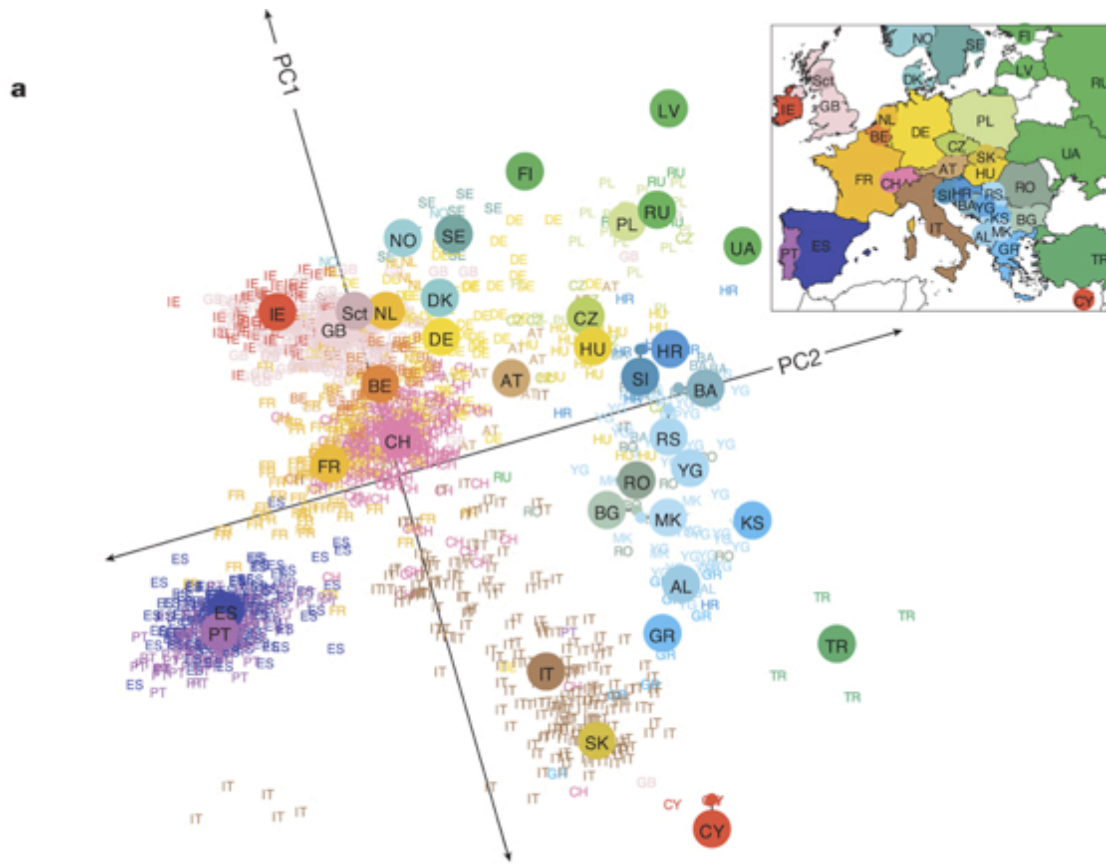
Population-based models can  
overcome systematic errors, but...

# Population Stratification is from the migration patterns of haplotypes throughout human history

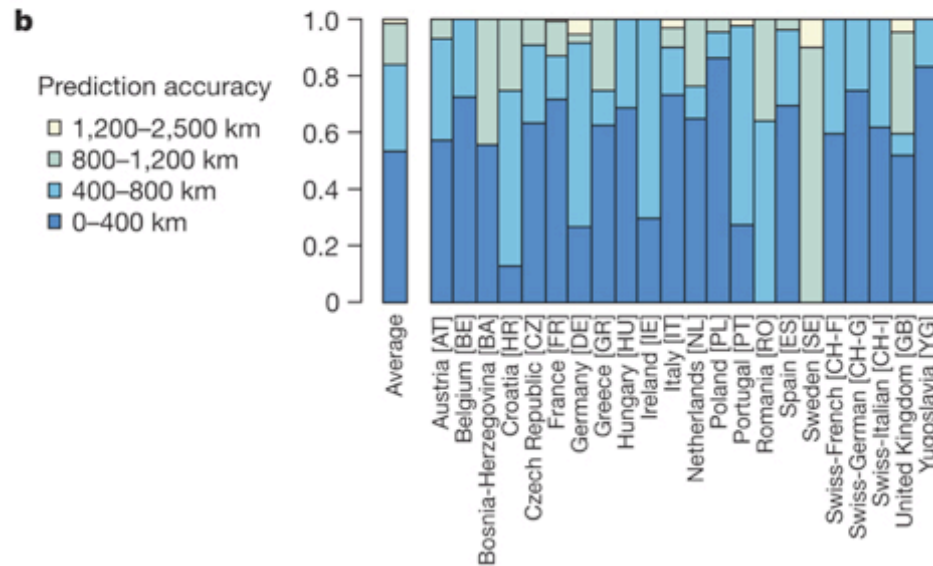
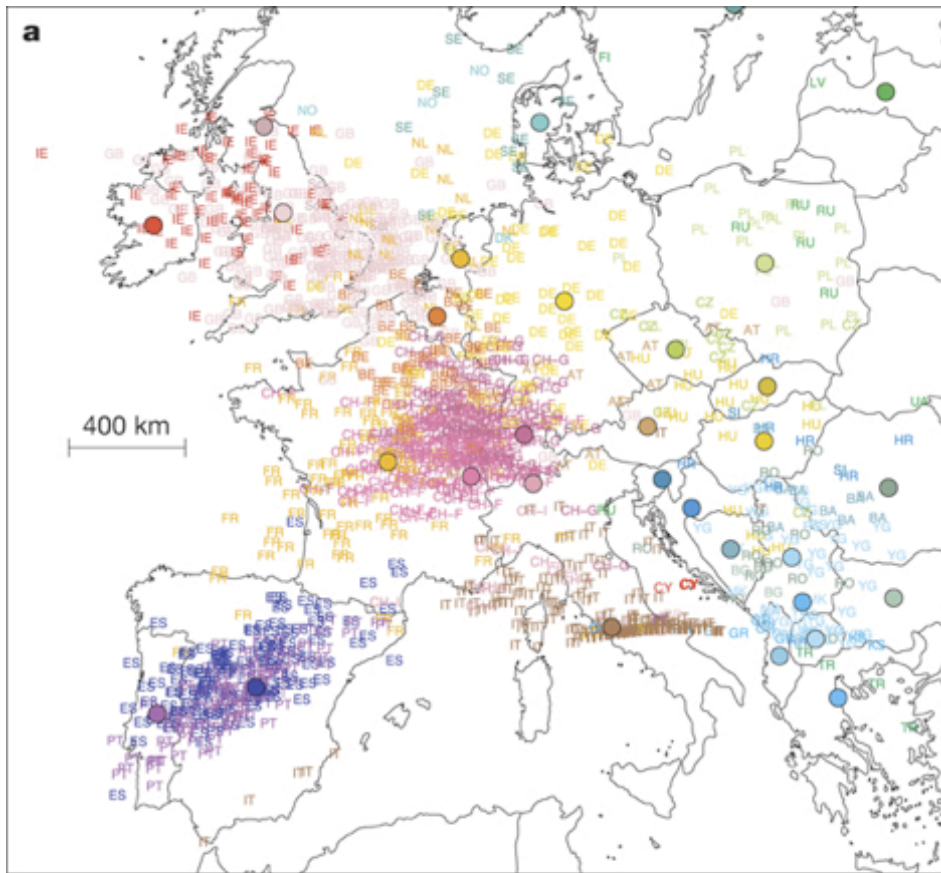


Tom Moore





Genes mirror geography within Europe  
 Novembre et al, 2008



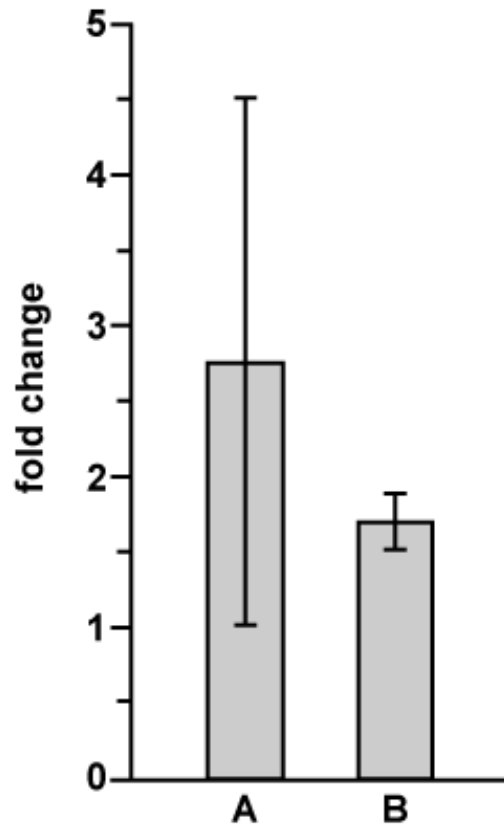
Genes mirror geography within Europe  
 Novembre et al, 2008

# Applications of Bayes to RNA-Seq



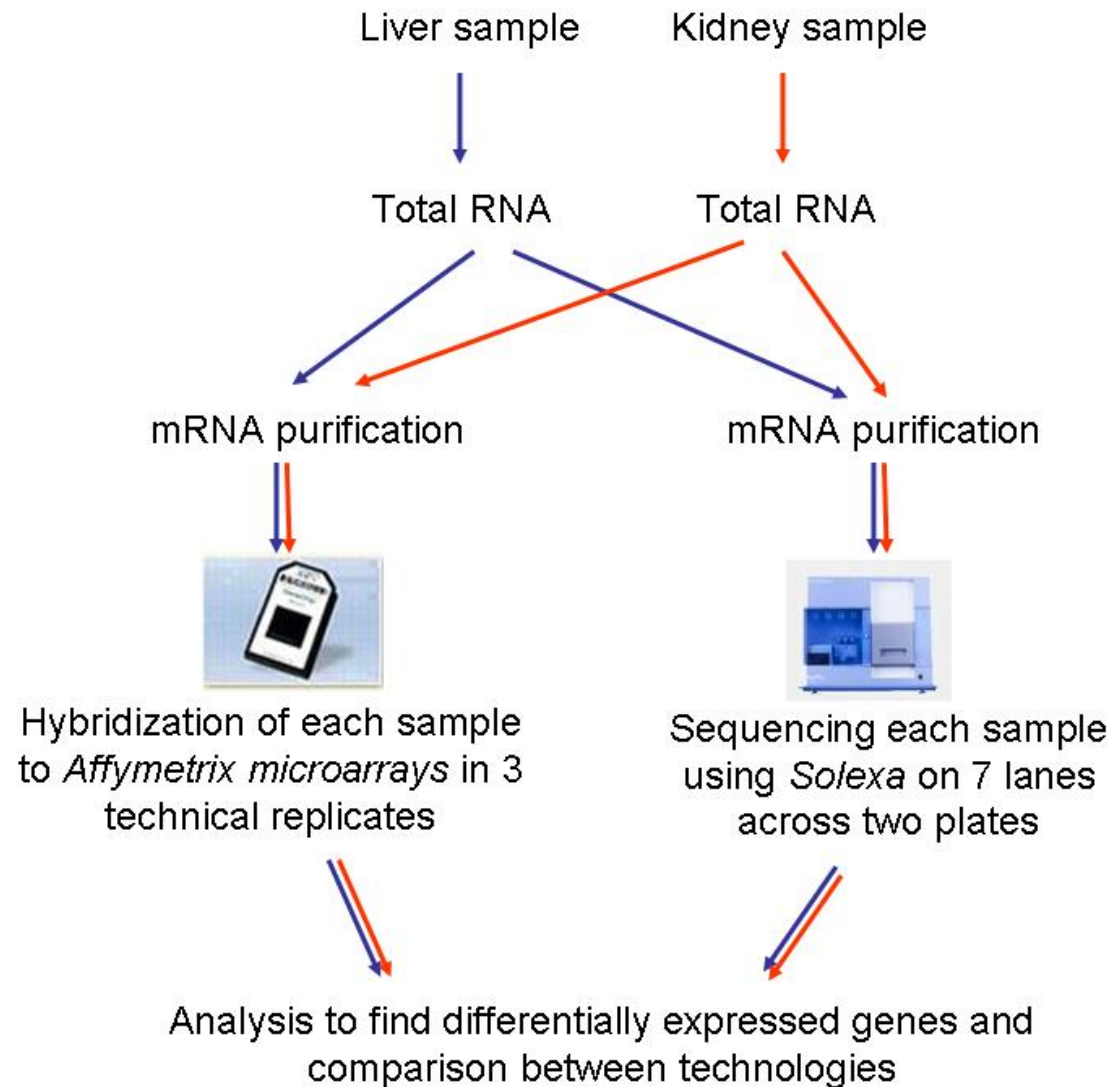
## Data Analysis: What genes are differentially expressed?

- Early days—fold change cutoffs (e.g., 2x difference or better)
- not a very satisfying approach:
  - doesn't take into account variance
  - misses any small changes



Here, “A” has a fold change  $>2.5$ , but varies greatly between replicate experiments. “B” has a fold change of only 1.75, but changes reliably each time the experiment is performed.

# Experimental Design: Liver vs. Kidney



# Metric for RNA-Seq Expression

RPKM:

Reads per Kilobase per Million Reads

Normalizes for (1) gene size and (2) sequencing depth  
(~0.1-1 transcript/cell)

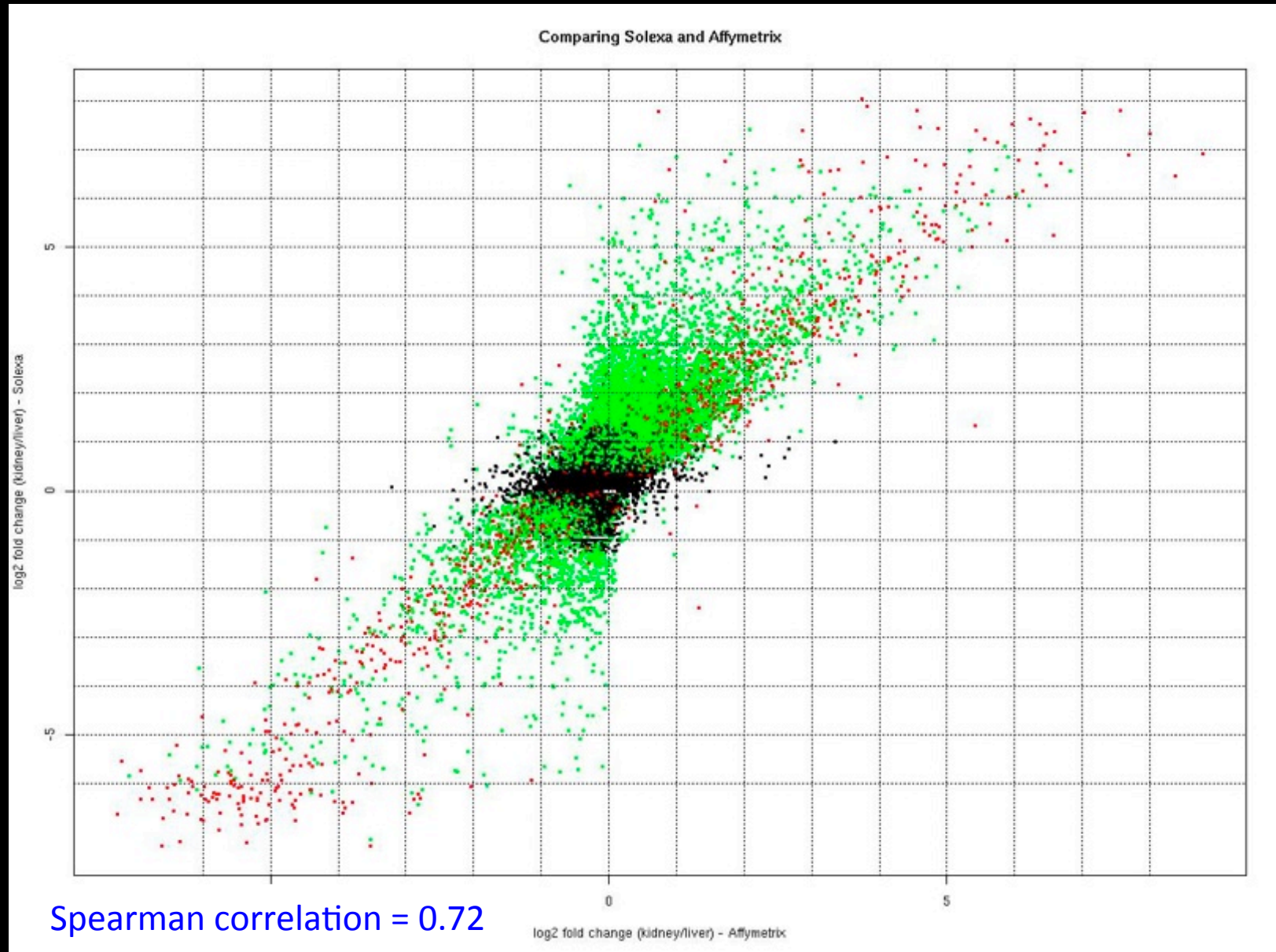
$$\text{RPKM} = \frac{N \text{ reads}}{1 \text{ gene}} \times \frac{1 \text{ gene}}{B \text{ bp}} \times \frac{1000 \text{ bp}}{1 \text{ Kb}} \times \frac{1 \text{ Million reads}}{Y \text{ total reads}}$$

Y = (exons, introns, intergenic reads)

FPKM=fragments-PKM  
is for paired-end data

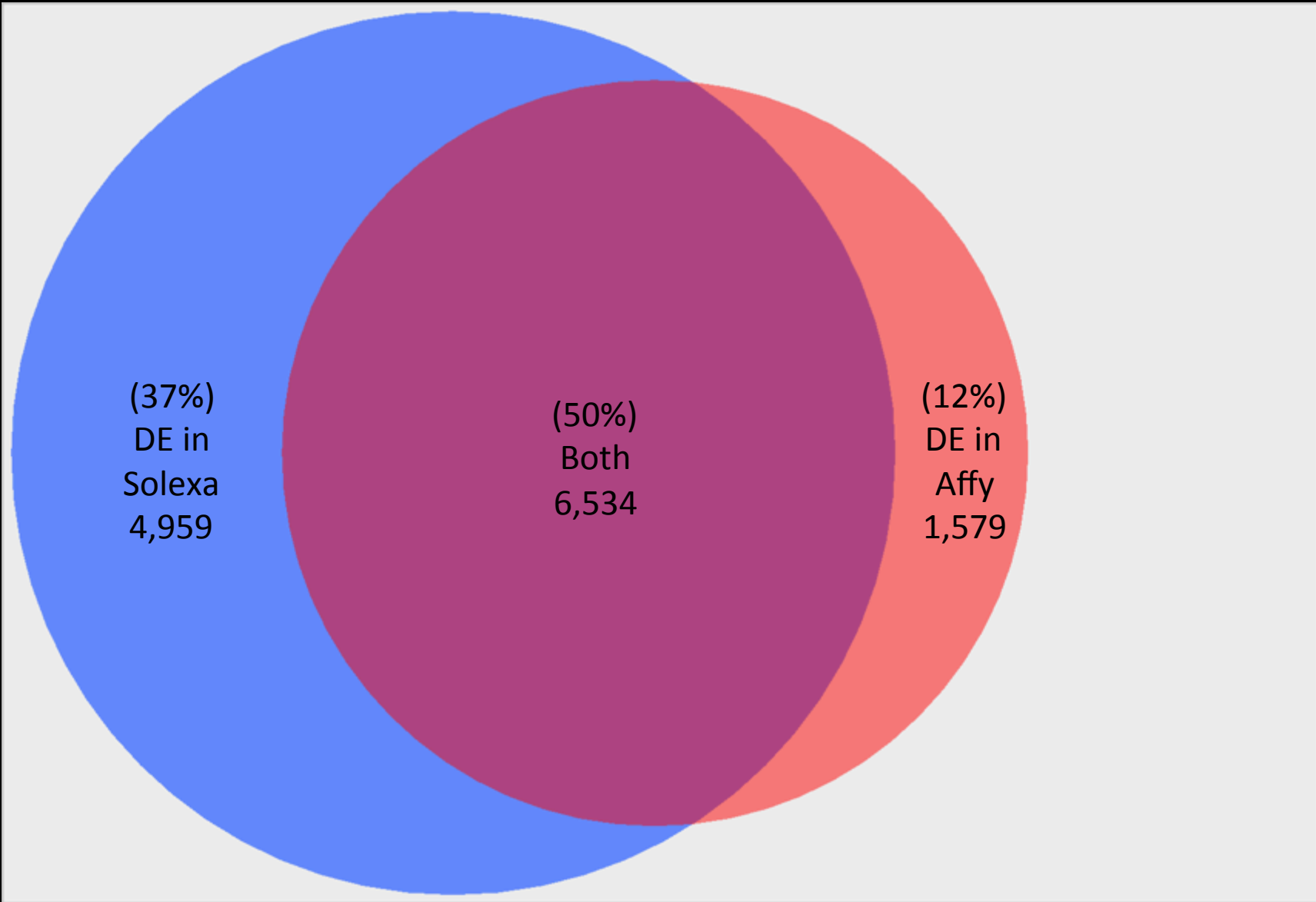
Mortazavi, Williams, *et al.*  
*Nature Methods*, 2008

# Comparing GA and Affy arrays





# 13,072 Differentially Expressed (DE) Genes



# Bias is introduced if these ratios are not kept:



**Good Run** RPKM = 128.4

**Bad Run** RPKM = 72.7

Mortazavi, Williams, *et al.*  
*Nature Methods*, 2008

# Normalization is needed

Define  $Y_{gk}$  as the observed count for gene  $g$  in library  $k$  summarized from the raw reads,  $\mu_{gk}$  as the true and unknown expression level (number of transcripts),  $L_g$  as the length of gene  $g$  and  $N_k$  as total number of reads for library  $k$ . We can model the expected value of  $Y_{gk}$  as:

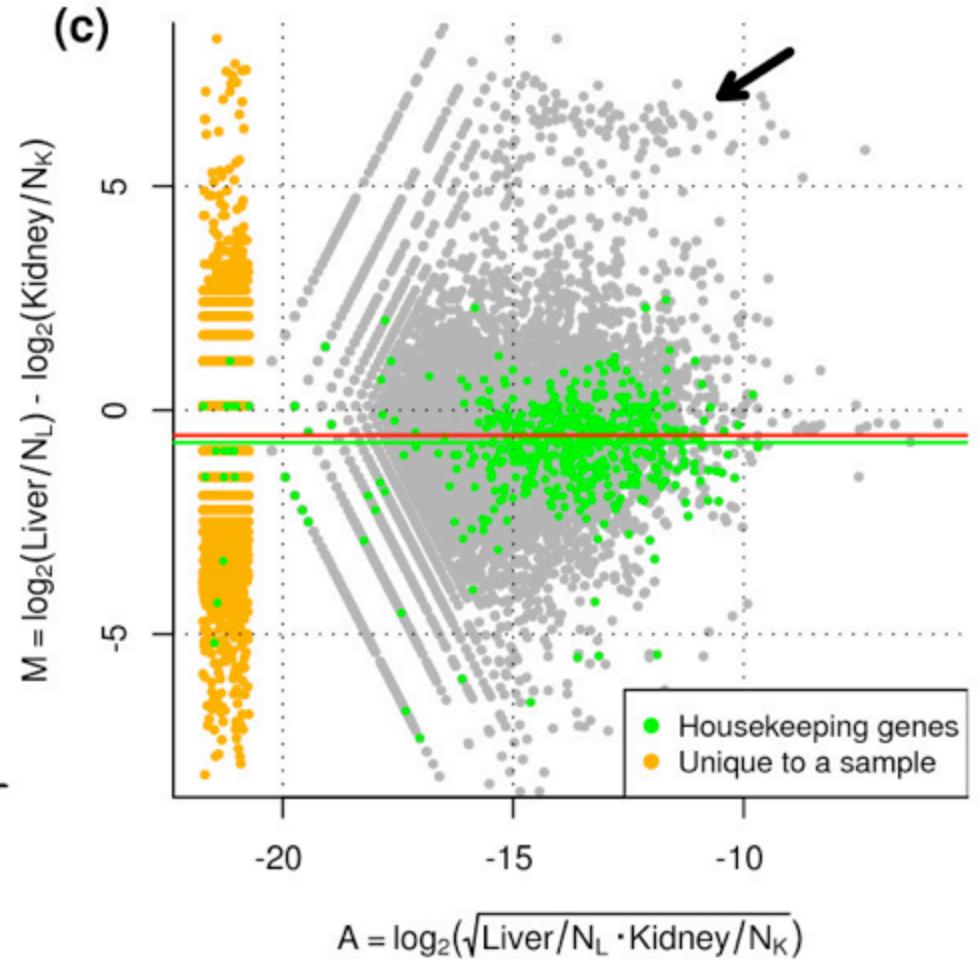
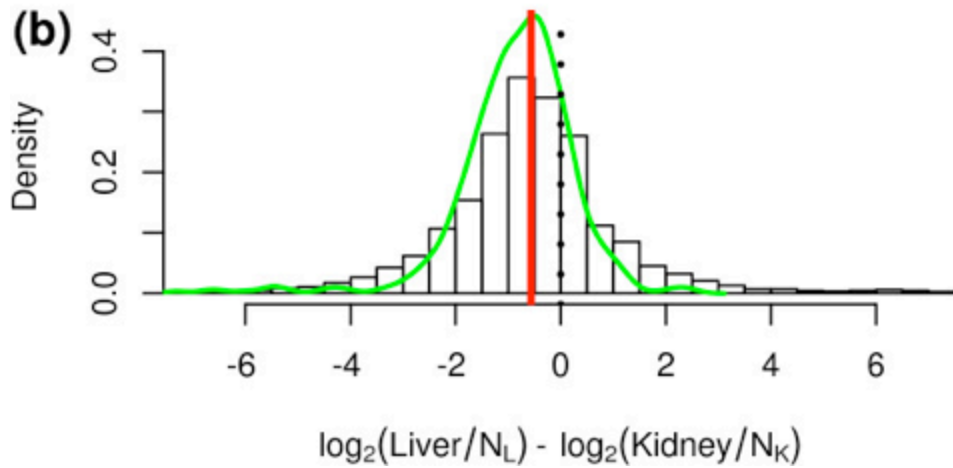
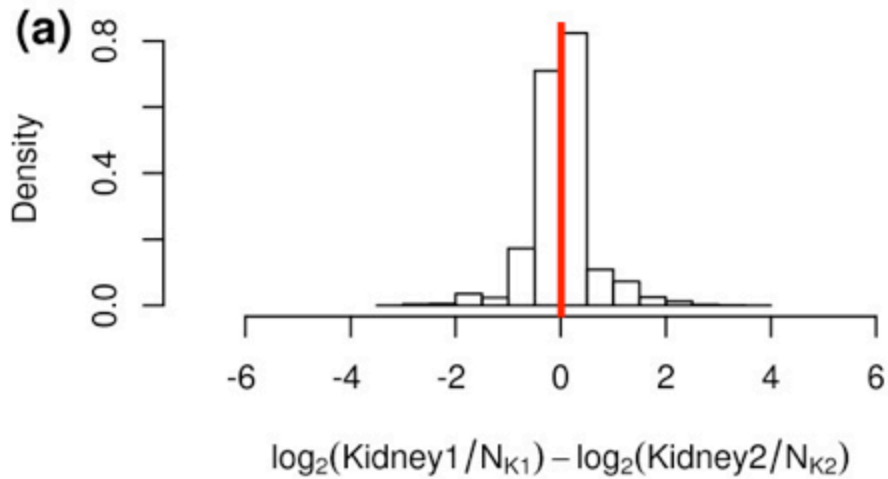
$$E[Y_{gk}] = \frac{\mu_{gk} L_g}{S_k} N_k$$

$$\text{where } S_k = \sum_{g=1}^G \mu_{gk} L_g$$

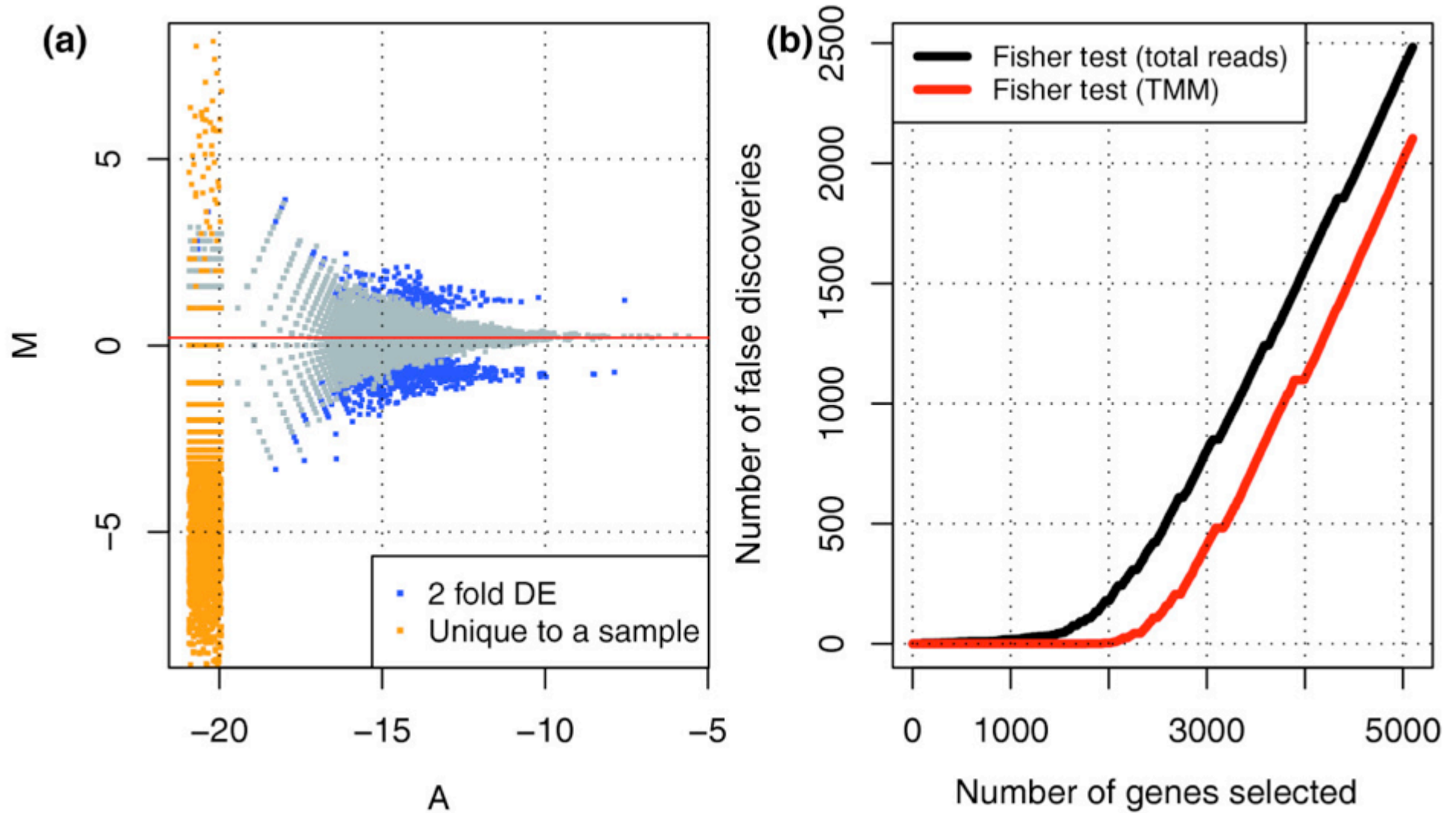
$S_k$  represents the total RNA output of a sample. The problem underlying the analysis of RNA-seq data is that while  $N_k$  is known,  $S_k$  is unknown and can vary drastically from sample to sample, depending on the RNA composition.

$$M_g = \log_2 \frac{Y_{gk} / N_k}{Y_{gk'} / N_{k'}}$$

# Normalization is needed



# Normalization is needed



Trimmed Mean M-values (TMM)

Robinson and Oshlack Genome Biology 2010 11:R25



# MISO / Probabilistic analysis and design of RNA-Seq experiments for identifying mRNA isoform regulation

[Home](#) | [Paper](#) | [Software](#) | [Documentation](#) | [Datasets](#) | [Contact](#)

## About MISO

### About MISO

MISO (Mixture of Isoforms) is a probabilistic framework that quantitates the expression level of alternatively spliced genes from RNA-Seq data, and identifies differentially regulated isoforms or exons across samples. By modeling the generative process by which reads are produced from isoforms in RNA-Seq, the MISO model uses Bayesian inference to compute the probability that a read originated from a particular isoform.

MISO uses the inferred assignment of reads to isoforms to quantitate the abundances of the underlying set of alternative mRNA isoforms. Confidence intervals over estimates can be obtained, which quantify the reliability of the estimates.



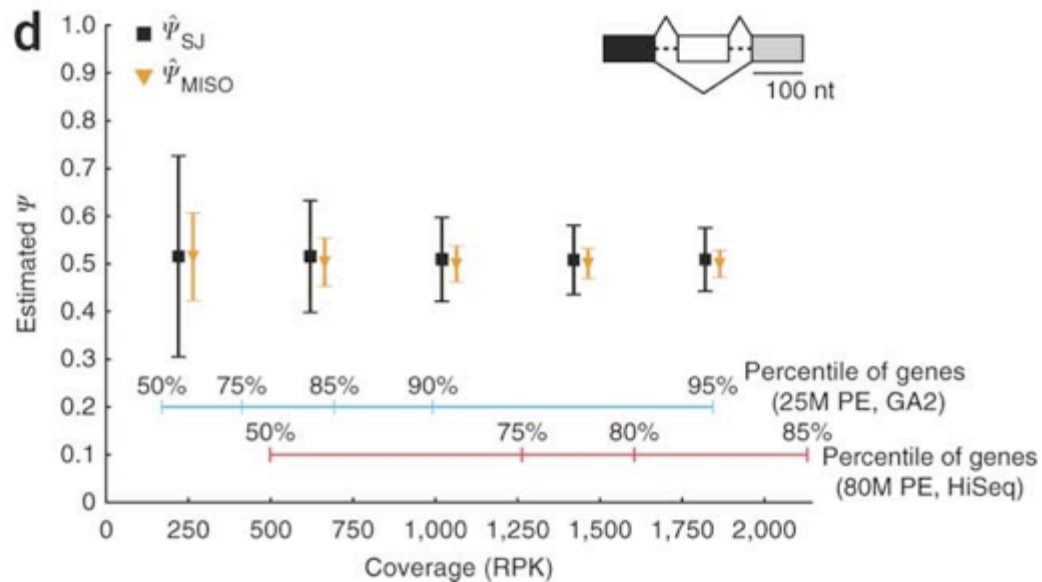
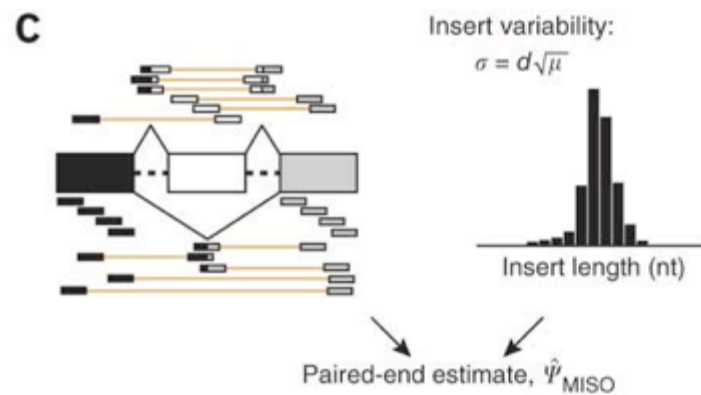
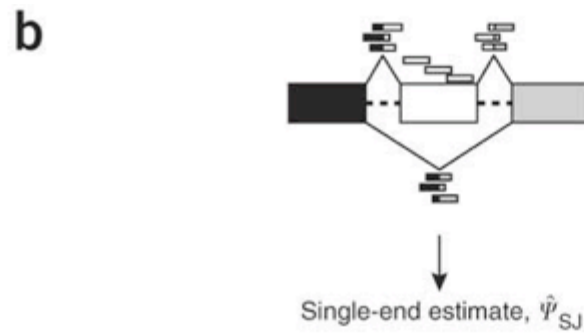
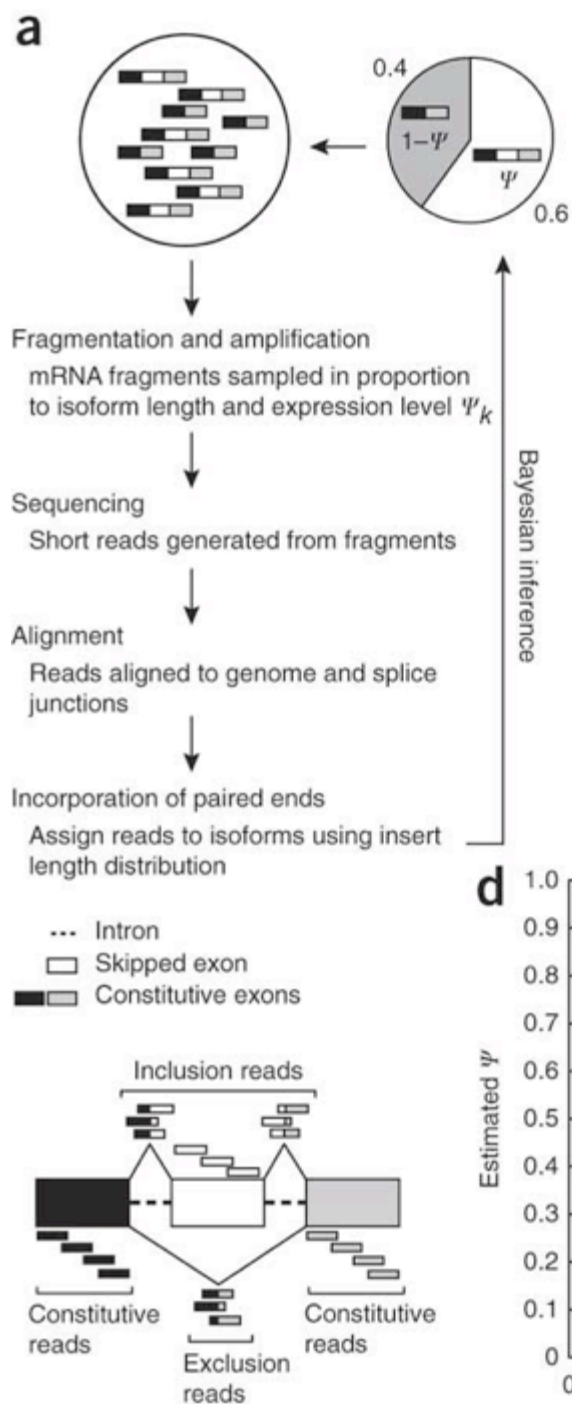
Depts. of [Biology](#) and [Biological Engineering](#)  
Dept. of [Brain and Cognitive Sciences](#)



Dept. of [Statistics](#)  
[FAS Center for Systems Biology](#)

### Contact

Department of Biology MIT  
31 Ames Street, 68-271A  
Cambridge, MA 02139-4307



# Coverage Requirements: How many lanes/plates/wells?

Depends on:

1. Read Length
2. Size of Transcriptome
3. Complexity of Transcriptome
4. Complexity of Tissue
5. Biological Variance
6. Errors (random and systematic)

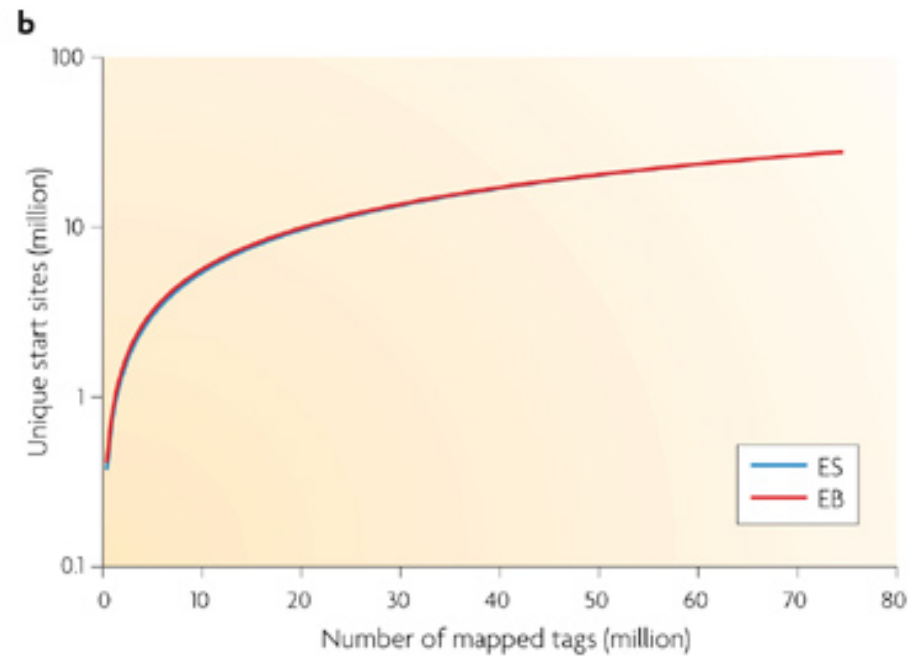
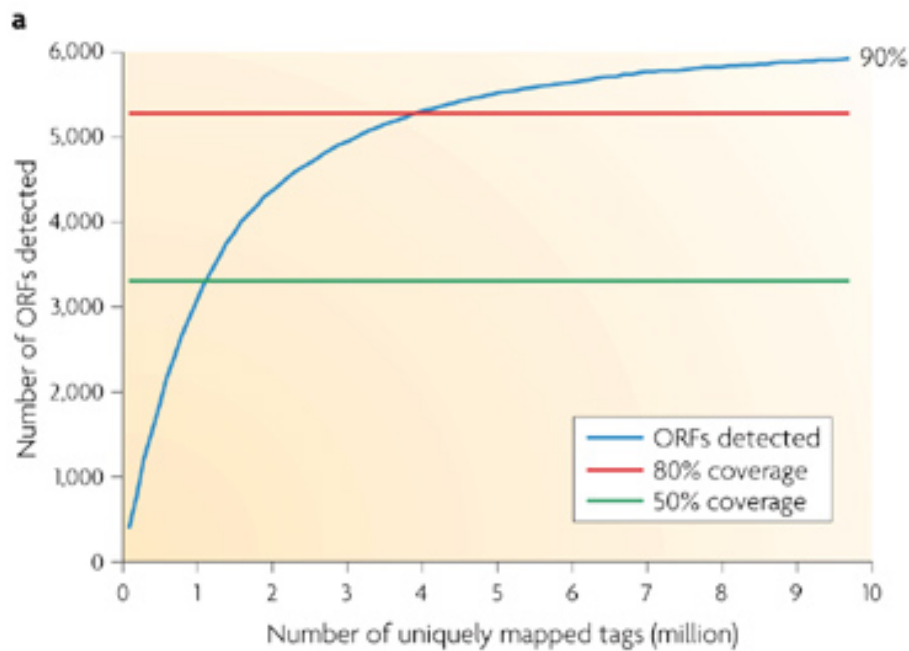


# Plateau of Information Starts @ ~500Mb

Number of lanes compared	Differentially expressed genes	Overlap with genes called from the array	Correlation of fold changes between Solexa and the array
One vs One	5670	4208	0.67
Two vs Two	7994	5340	0.70
Three vs Three	9482	5909	0.71
Four vs Four	10580	6278	0.72
Five vs Five	11493	6534	0.73

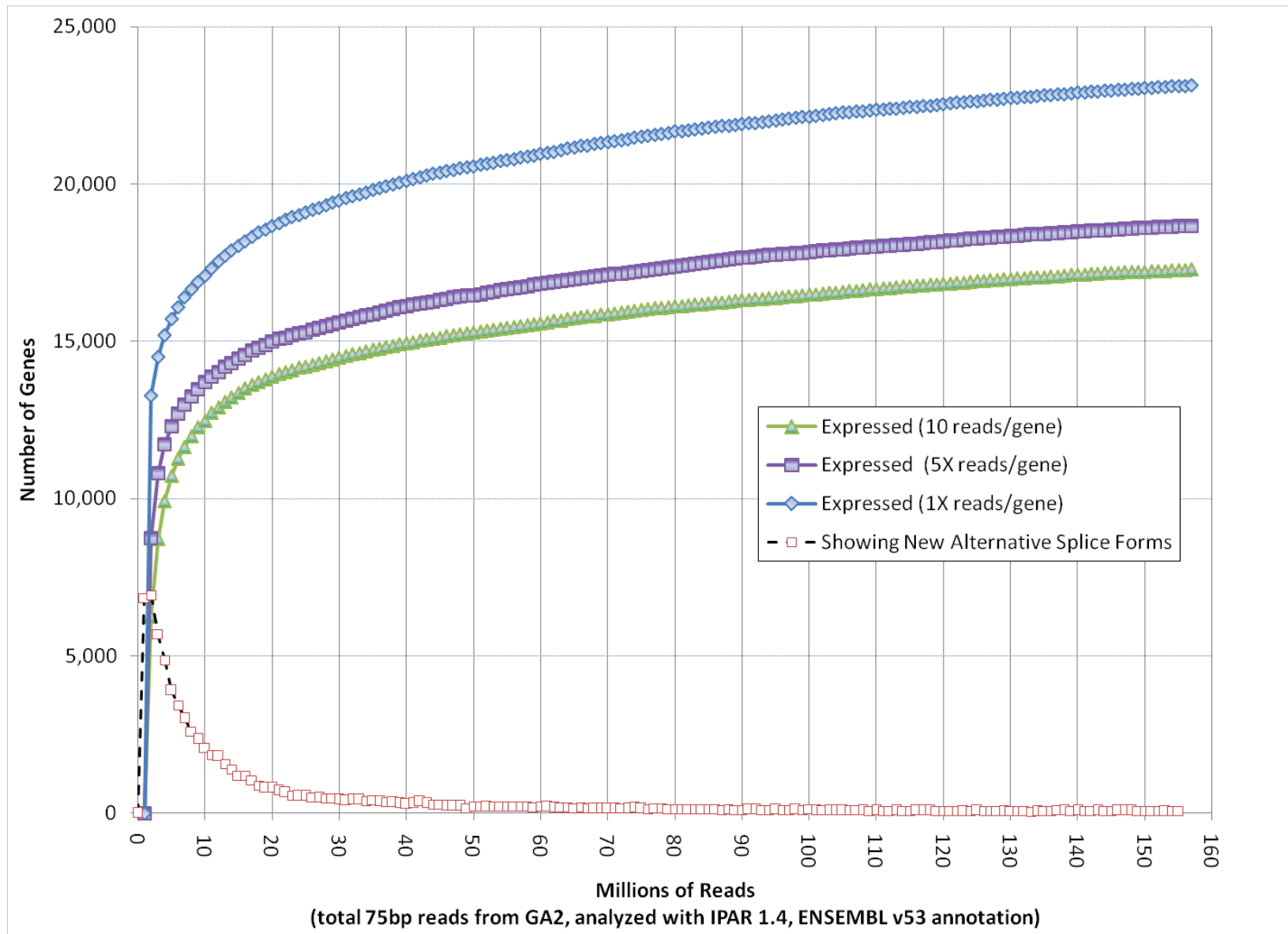
Liver			Kidney		
	No genes	Percentage		No genes	Percentage
Five Lanes	20080	100	Five Lanes	20921	100
Four Lanes	19695	97.9	Four Lanes	20552	98.2
Three Lanes	19170	95.5	Three Lanes	20064	96.0
Two Lanes	18390	91.6	Two Lanes	19355	92.5
One Lane	16973	84.5	One Lane	18080	86.4

# Coverage Requirements



Nature Reviews | Genetics

# No current visible end of gene discovery



# How many replicates do I need?

Calculation of the number of replicates depends on:

1. An estimate of  $\sigma^2$  obtained from previous experiments.
2. The size of the difference ( $\delta$ ) to be detected.
3. The assurance with which it is desired to detect the difference (i.e., Power of the test =  $1-\beta$ ).
4. The level of significance to be used in the actual experiment (i.e., Type I error).
5. The test required, whether a one-tail or two-tail test.

To determine the number of replicates use the following formula :

$$\#reps = 2 \left( Z_{\alpha/2} + Z_{\beta} \right) \left( \frac{\sigma}{\delta} \right)^2$$

where:  $Z_{\alpha/2}$  is associated with the Type I error (two-tailed)

$Z_{\beta}$  is associated with the Type II error

$\delta$  is the true difference to be detected, and

$\sigma$  is the known variance obtained from previous experiments

# Bayes in Chip-seq too!

Research article

Highly accessed

Open Access

## BayesPeak: Bayesian analysis of ChIP-seq data

**Christiana Spyrou**<sup>1,3</sup> ✉, **Rory Stark**<sup>3</sup> ✉, **Andy G Lynch**<sup>4</sup> ✉ and **Simon Tavaré**<sup>2,4</sup> ✉

1 Statistical Laboratory, Centre for Mathematical Sciences, Wilberforce Road, Cambridge, UK

2 DAMTP, Centre for Mathematical Sciences, Wilberforce Road, Cambridge, UK

3 Cancer Research UK, Cambridge Research Institute, Li Ka Shing Centre, Robinson Way, Cambridge UK

4 Department of Oncology, University of Cambridge, Li Ka Shing Centre, Robinson Way, Cambridge, UK

✉ author email   ✉ corresponding author email

*BMC Bioinformatics* 2009, **10**:299   doi:10.1186/1471-2105-10-299

The electronic version of this article is the complete one and can be found online at: <http://www.biomedcentral.com/1471-2105/10/299>

Received: 8 May 2009

Accepted: 21 September 2009

Published: 21 September 2009

© 2009 Spyrou et al; licensee BioMed Central Ltd.

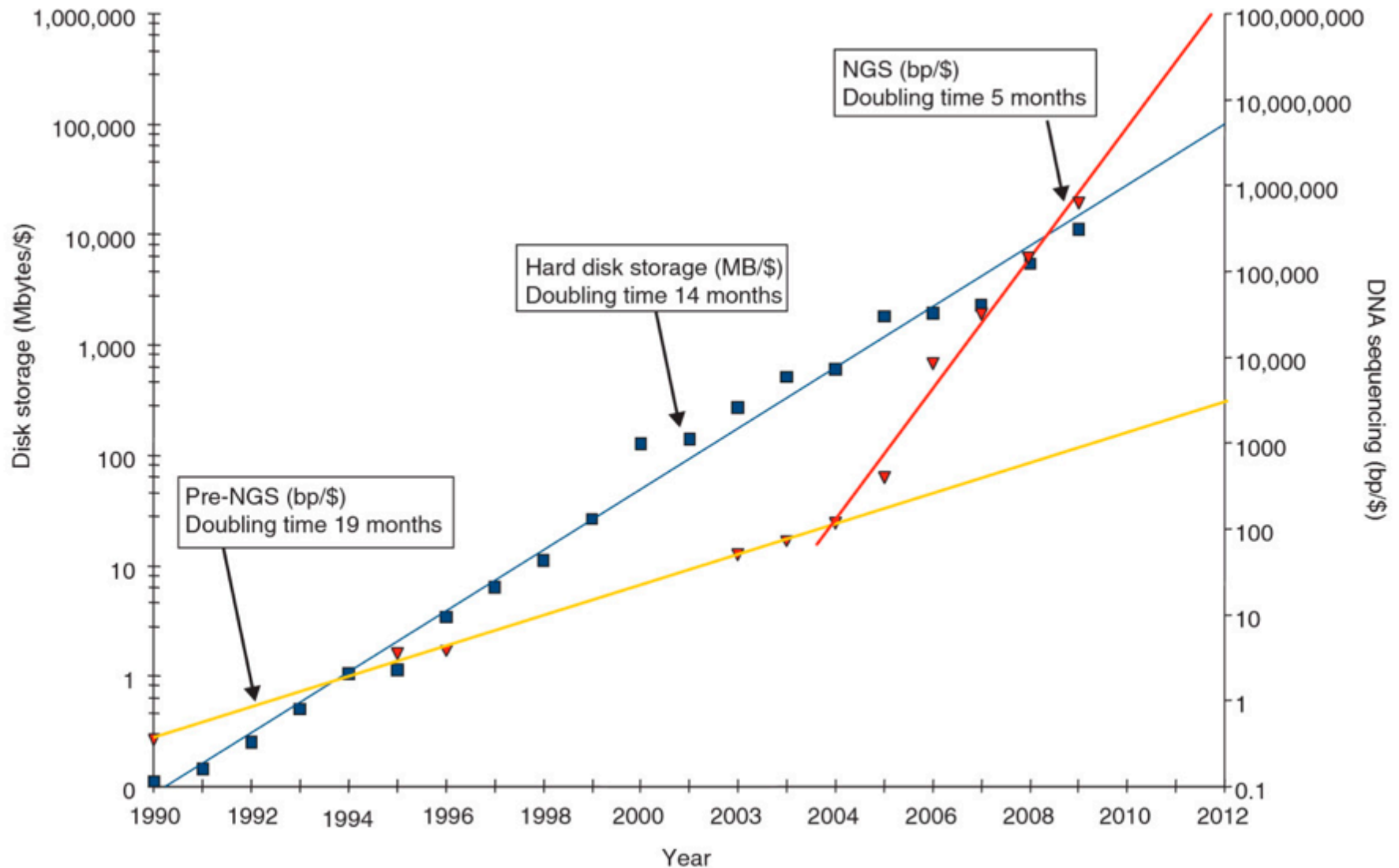
This is an Open Access article distributed under the terms of the Creative Commons Attribution License

(<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

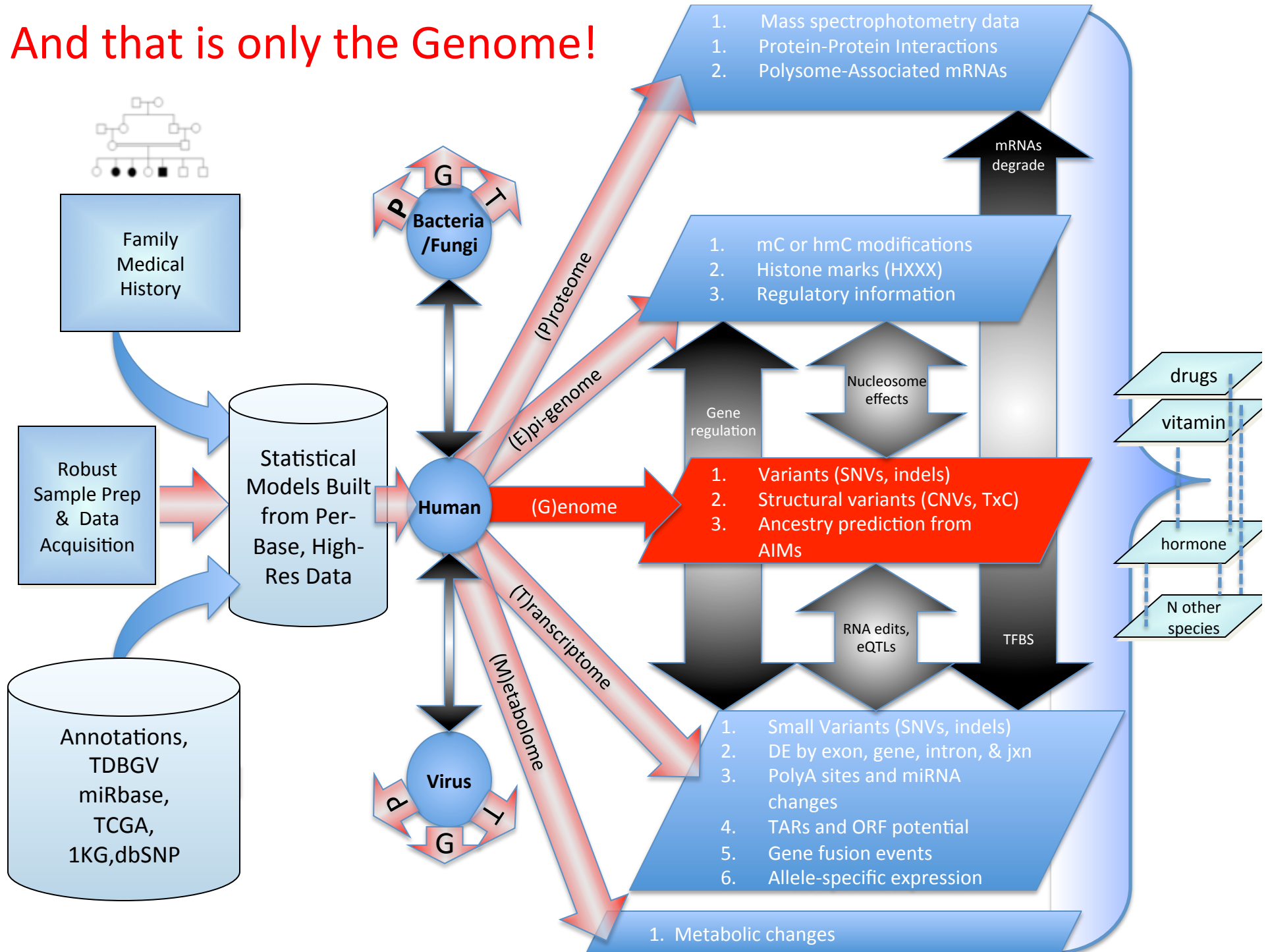
# What's the problem with Bayesian statistics? (according to non-Bayesians)

1. Priors can introduce subjective judgment into data analysis.
2. Priors affect the result. Different people can get different answers from the same data.
3. It's too hard. There are no simple point-and-click programs.

# This requires a lot of space

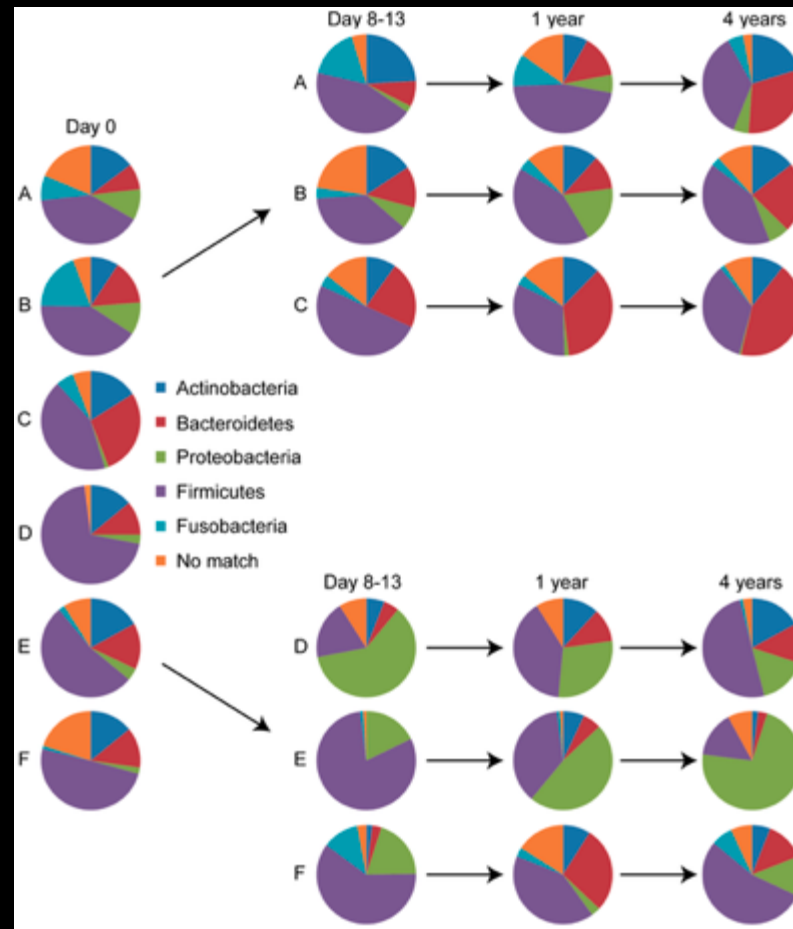


# And that is only the Genome!





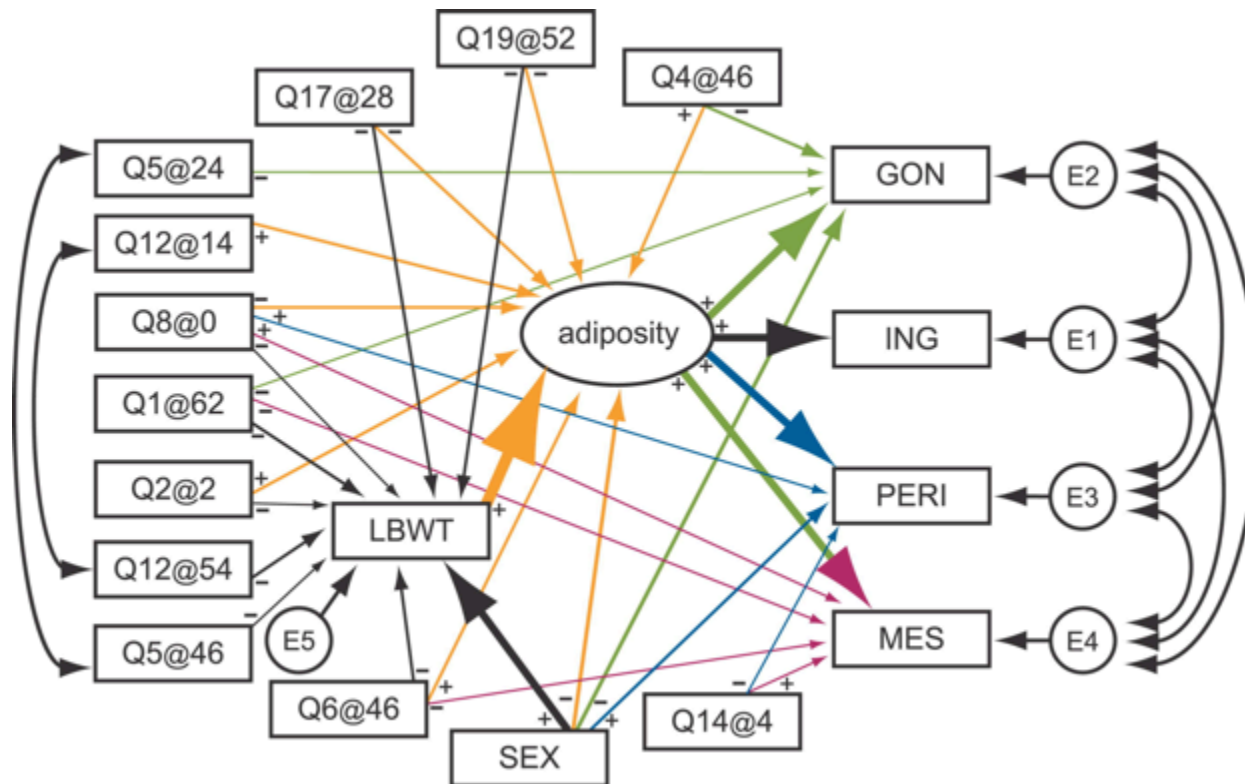
Meta-genomic phenotypes can persist for years, and “passenger genomes” can be a phenotype, as well as their distributions.



Normal throat

Throat + antibiotics

# Risk factors for diseases usually involve many genes and pathways



# There are other factors than these!



The image shows a screenshot of a Nature journal article page. At the top, the word "nature" is written in a large, white, serif font on a dark red background. To its right, the text "International weekly journal of science" is written in a smaller, white, sans-serif font. Below this, there is a search bar with the text "Search This journal". The main content area is white and contains the following text: "nature.com > journal home > current issue > letter > full text", "NATURE | LETTER", the title "Genome, epigenome and RNA sequences of monozygotic twins discordant for multiple sclerosis SHOW NO DIFFERENCE", a list of authors including Sergio E. Baranzini, Joann Mudge, Jennifer C. van Velkinburgh, Pouya Khankhanian, Irina Khrebtukova, Neil A. Miller, Lu Zhang, Andrew D. Farmer, Callum J. Bell, Ryan W. Kim, Gregory D. May, Jimmy E. Woodward, Stacy J. Caillier, Joseph P. McElroy, Refujia Gomez, Marcelo J. Pando, Leonda E. Clendenen, Elena E. Ganusova, Faye D. Schilkey, Thiruvarangan Ramaraj, Omar A. Khan, Jim J. Huntley, Shujun Luo, Pui-yan Kwok, Thomas D. Wu, and "et al.", and links for "Affiliations", "Contributions", and "Corresponding authors". At the bottom, it says "Nature 464, 1351–1356 (29 April 2010) | doi:10.1038/nature08990" and "Received 25 July 2009 | Accepted 11 March 2010".

nature.com > journal home > current issue > letter > full text

NATURE | LETTER

Genome, epigenome and RNA sequences of monozygotic twins discordant for multiple sclerosis **SHOW NO DIFFERENCE**

Sergio E. Baranzini, Joann Mudge, Jennifer C. van Velkinburgh, Pouya Khankhanian, Irina Khrebtukova, Neil A. Miller, Lu Zhang, Andrew D. Farmer, Callum J. Bell, Ryan W. Kim, Gregory D. May, Jimmy E. Woodward, Stacy J. Caillier, Joseph P. McElroy, Refujia Gomez, Marcelo J. Pando, Leonda E. Clendenen, Elena E. Ganusova, Faye D. Schilkey, Thiruvarangan Ramaraj, Omar A. Khan, Jim J. Huntley, Shujun Luo, Pui-yan Kwok, Thomas D. Wu  *et al.*

[Affiliations](#) | [Contributions](#) | [Corresponding authors](#)

Nature 464, 1351–1356 (29 April 2010) | doi:10.1038/nature08990  
Received 25 July 2009 | Accepted 11 March 2010

Systems biology requires spatiotemporal monitoring of the genome, epigenome, transcriptome, proteome, metabolome, and the environment, to see the interactome