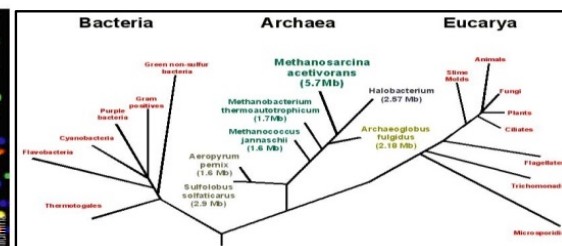
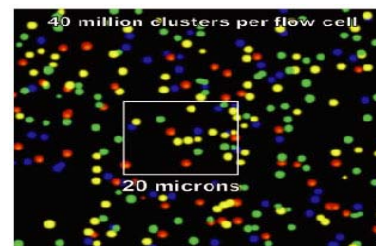




TAACCCTAACCCCTAACCCCTAACCCCTAACCCCTA  
 CCTAACCCCTAACCCCTAACCCCTAACCCCTAACCC  
 CCTAACCCCTAACCCCTAACCCCTAACCCCTAAC  
 AACCCCTAACCCCTAACCCCTAACCCCTAACCCCTA  
 ACCCTAACCCCTAACCCCTAACCCCTAACCCCTAAC  
 CTACCCTAACCCCTAACCCCTAACCCCTAACCCCTA  
 ACCCTAACCCCTAACCCCTAACCCCTAACCCCTAA

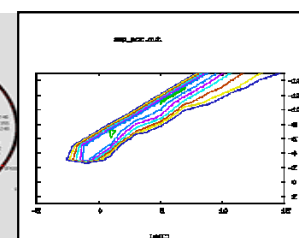
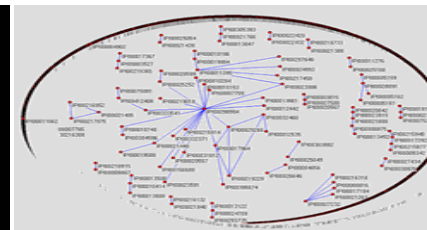
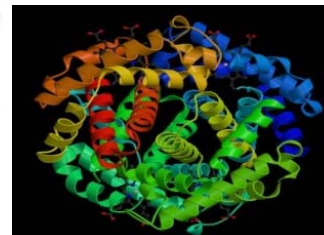
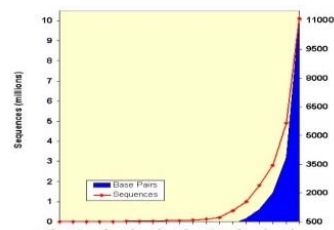
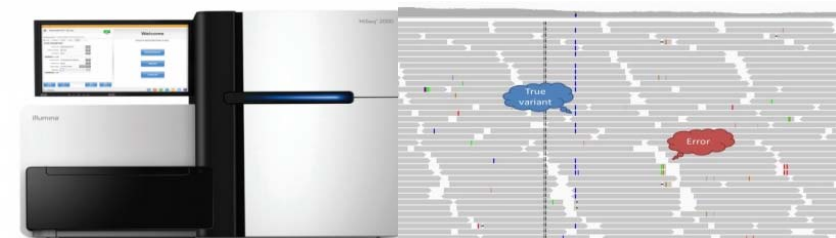


# Transcriptome Analysis with noncoding RNAs

北京大学生物信息学中心 高歌

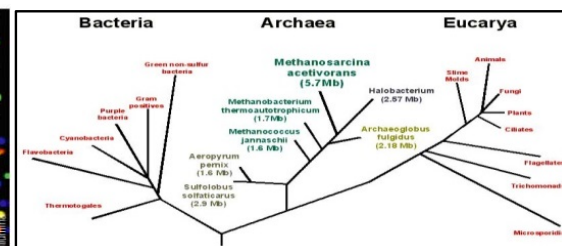
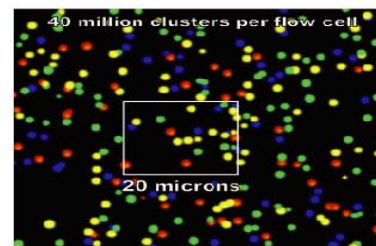
Ge Gao, Ph.D.

Center for Bioinformatics, Peking University





TAACCCTAACCCTAACCCTAACCCTAACCCTA  
CCTAACCCTAACCCTAACCCTAACCCTAACC  
CCCTAACCCTAACCCTAACCCTAACCCTAAC  
AACCCTAACCCTAACCCTAACCCTAACCCTA  
ACCCTAACCCTAACCCTAACCCTAACCCTAAC  
CTACCCTAACCCTAACCCTAACCCTAACCCTA  
ACCCTAACCCTAACCCTAACCCTAACCCTAA

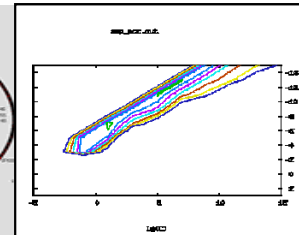
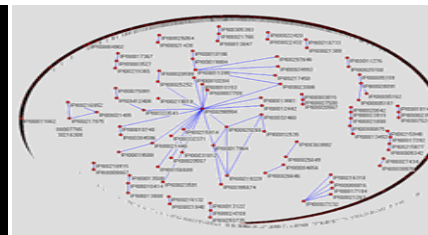
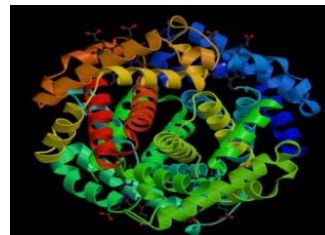
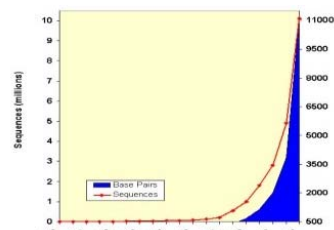
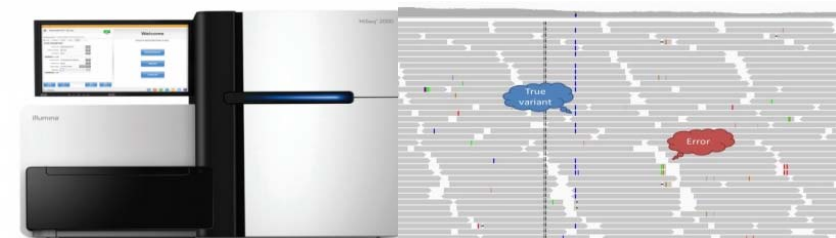


# Unit 3: Data Mining: Differential Expression and Clustering

北京大学生物信息学中心 高歌

Ge Gao, Ph.D.

Center for Bioinformatics, Peking University



How many non-coding transcripts?

What are the functional roles of those ncRNAs?

## microRNA (miRNA)

- single-stranded RNAs of 21-23 (or some say 20-25) bp RNAs with regulatory

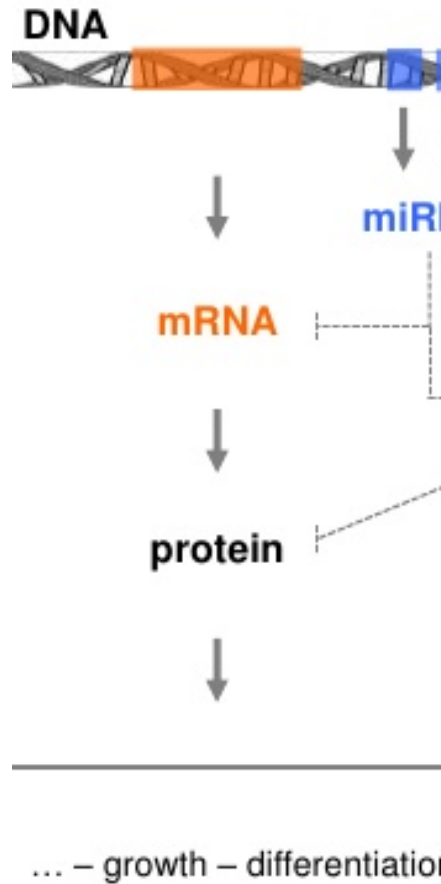
Table 4.2 Computational algorithms for microRNA target prediction

Name of the software	URL or availability	Supported organism(s)	Reference(s)
TargetScan, TargetScanS	<a href="http://genes.mit.edu/targetscan/">http://genes.mit.edu/targetscan/</a>	Vertebrates	Lewis <i>et al.</i> , 2003, 2005
miRanda	<a href="http://www.microrna.org/">http://www.microrna.org/</a>	Flies, vertebrates	Enright <i>et al.</i> , 2003, John <i>et al.</i> , 2004
DIANA-microT	<a href="http://diana.pcbi.upenn.edu/DIANA-microT/">http://diana.pcbi.upenn.edu/DIANA-microT/</a>	Vertebrates	Kiriakidou <i>et al.</i> , 2004
RNAhybrid	<a href="http://bibiserv.techfak.uni-bielefeld.de/rnahybrid/">http://bibiserv.techfak.uni-bielefeld.de/rnahybrid/</a>	Flies	Rehmsmeier <i>et al.</i> , 2004
GUUGle	<a href="http://bibiserv.techfak.uni-bielefeld.de/guugle/">http://bibiserv.techfak.uni-bielefeld.de/guugle/</a>	Flies	Gerlach <i>et al.</i> , 2006
PicTar	<a href="http://pictar.bio.nyu.edu/">http://pictar.bio.nyu.edu/</a>	Nematodes, flies, vertebrates	Grun <i>et al.</i> , 2005, Krek <i>et al.</i> , 2005, Lall <i>et al.</i> , 2006
MicroInspector	<a href="http://mirna.imbb.forth.gr/microinspector/">http://mirna.imbb.forth.gr/microinspector/</a>	Any	Rusinov <i>et al.</i> , 2005
MovingTargets	Available by request on DVD	Flies	Burgler <i>et al.</i> , 2005
FastCompare	<a href="http://tavazoielab.princeton.edu/mirnas/">http://tavazoielab.princeton.edu/mirnas/</a>	Nematodes, flies	Chan <i>et al.</i> , 2005
miRU	<a href="http://bioinfo3.noble.org/miRNA/miRU.htm">http://bioinfo3.noble.org/miRNA/miRU.htm</a>	Plants	Zhang 2005
TargetBoost	<a href="https://demo1.interagon.com/demo/">https://demo1.interagon.com/demo/</a>	Nematodes, flies	Saetrom <i>et al.</i> , 2006
rna22	<a href="http://cbcsrv.watson.ibm.com/rna22.html">http://cbcsrv.watson.ibm.com/rna22.html</a>	Nematodes, flies, vertebrates	Miranda <i>et al.</i> , 2006
miTarget	<a href="http://cbit.snu.ac.kr/~miTarget/">http://cbit.snu.ac.kr/~miTarget/</a>	Any	Kim <i>et al.</i> , 2006

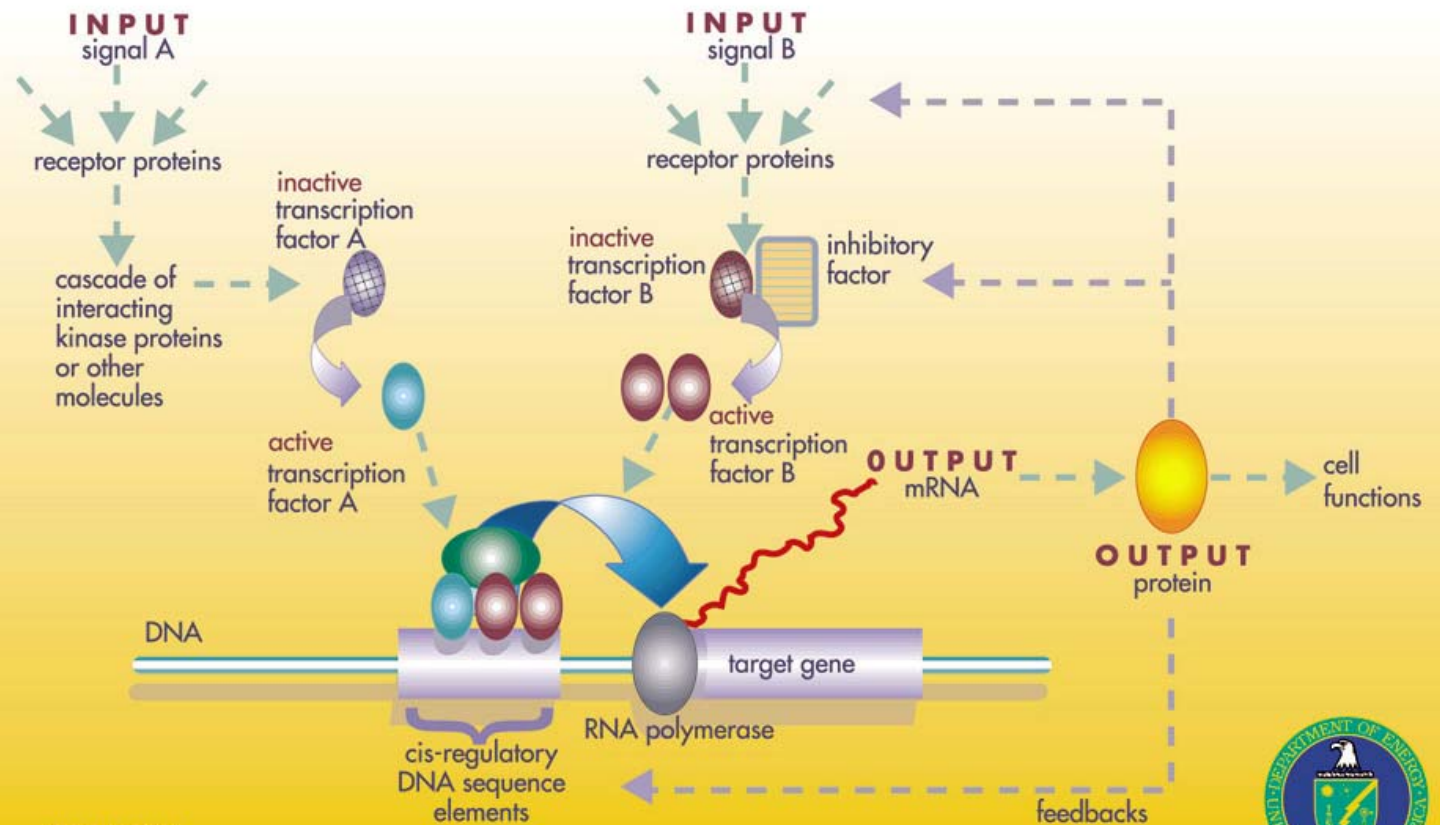
Target mRNAs from loci unrelated to miRNA gene

(source: *Cell* 116:281)

- the transcriptome



## A GENE REGULATORY NETWORK



YGG 01-0083

(Modified from [public.ornl.gov/site/gallery/highres/REGNET.jpg](http://public.ornl.gov/site/gallery/highres/REGNET.jpg))

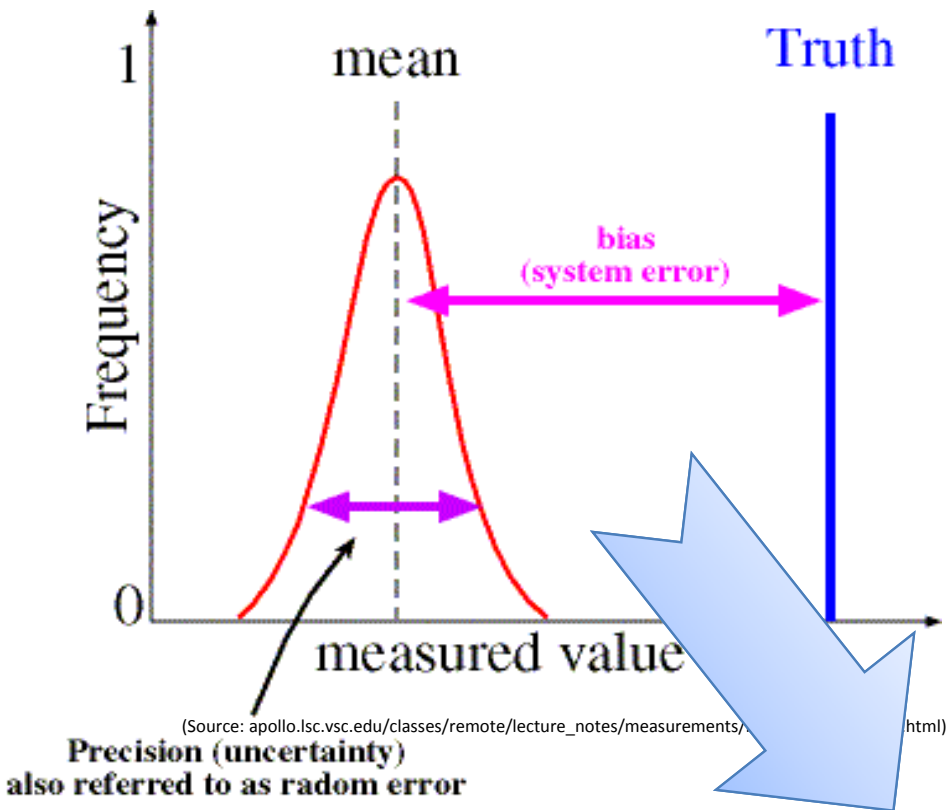


- Differentially expressed genes
- Co-expressed genes

## Data Mining: Differentially Expression Calling

- Identify the genes with **biological-significant difference** in expression levels across samples
- Differences in expression values can result from many non-biological sources (e.g. experiment error/bias)
  - The ‘real’ differences are the differences that can **NOT** be explained by the various errors introduced during the experimental phase

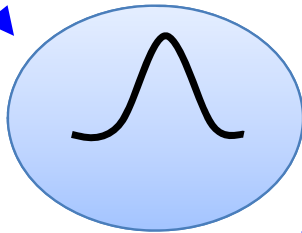
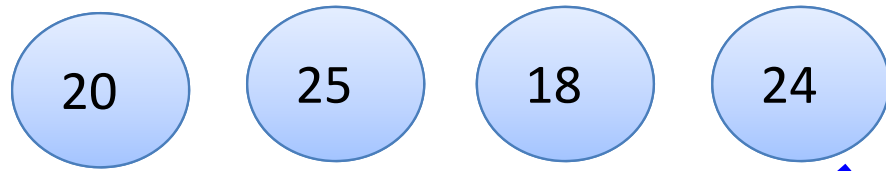




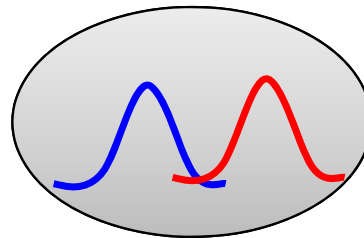
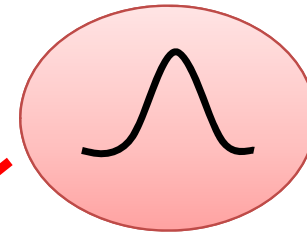
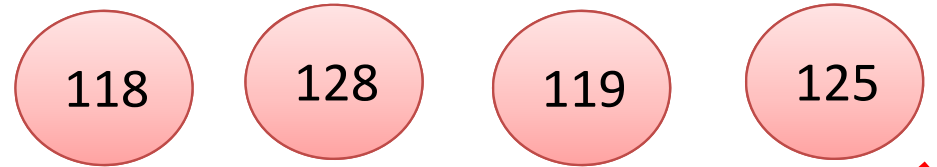
- **Random errors** arise from random fluctuations in the measurements
- It could be reduced by repeating experiment many times (and get a mean value)
- Random errors could be modeled statistically by **variance**.



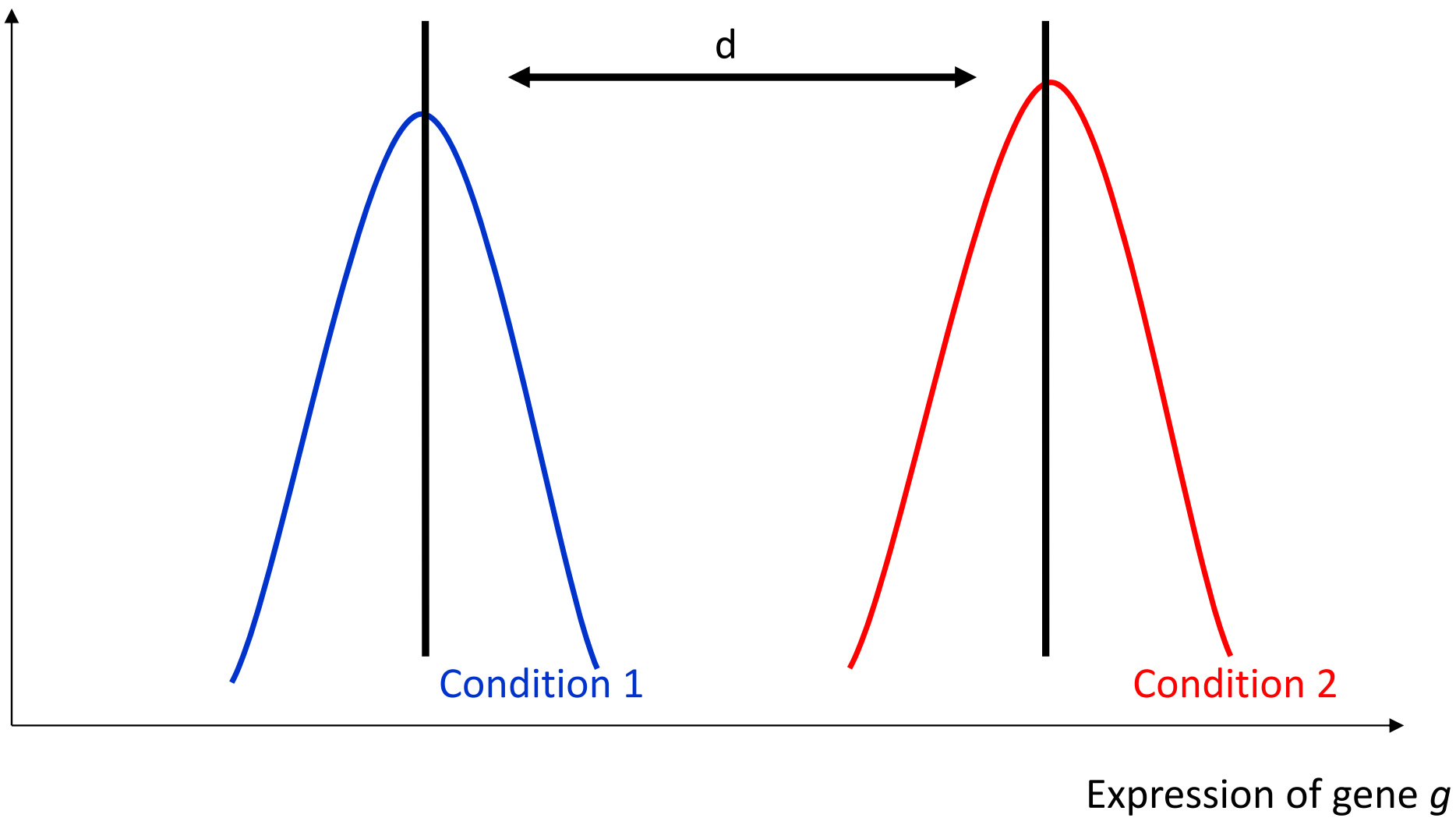
Condition 1

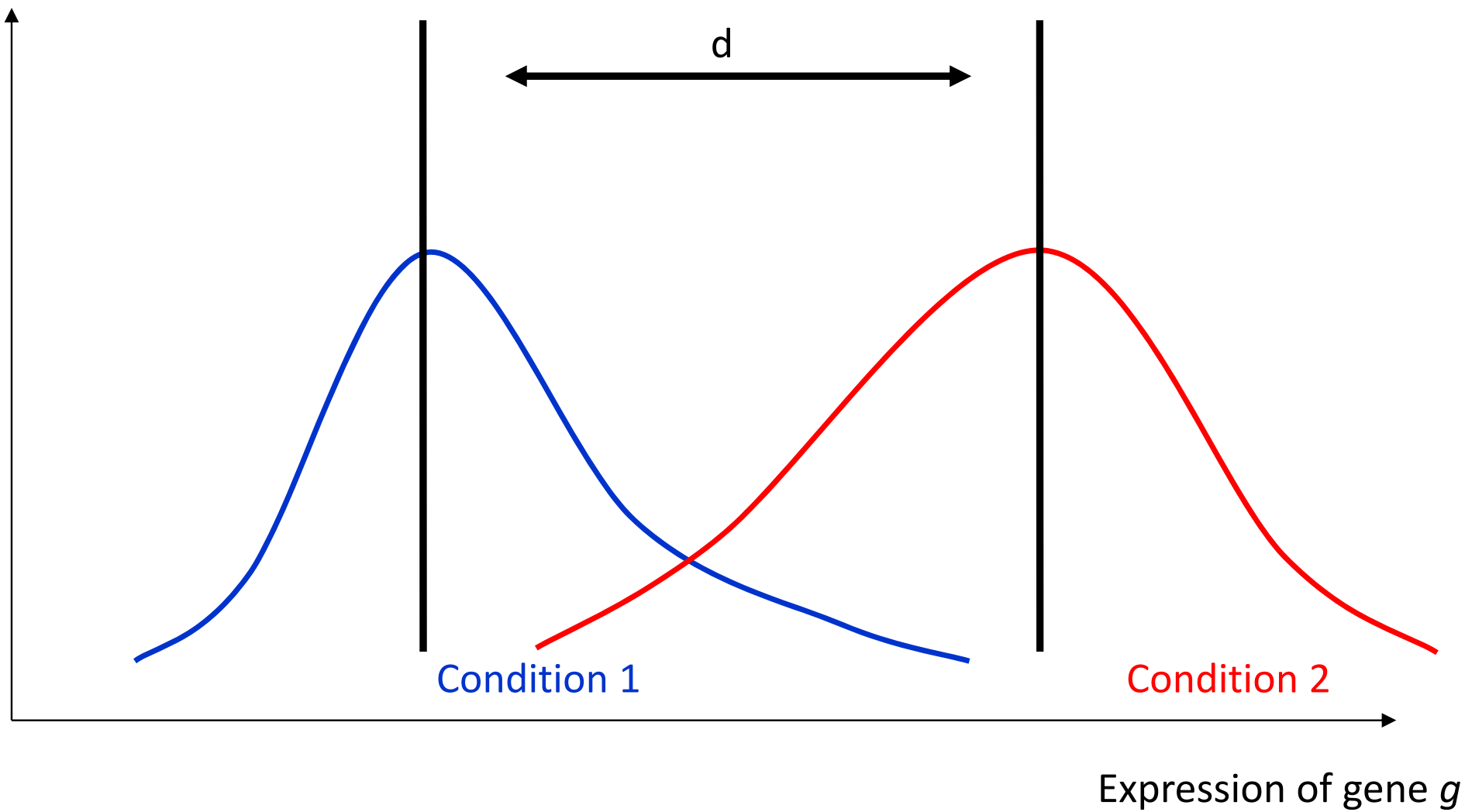


Condition 2



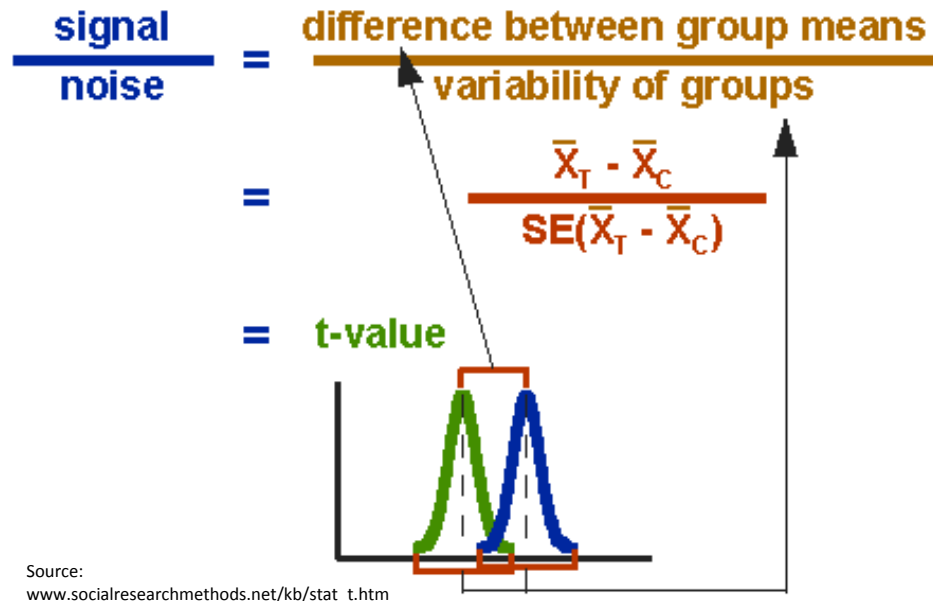
Distribution of differential expression statistic





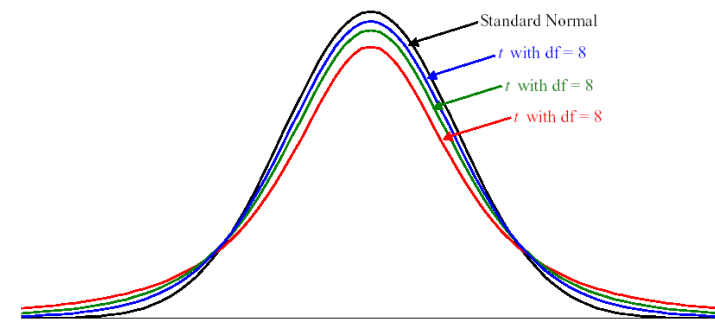
# Statistical calling

1. Select a **statistic** which takes the variance into account, and will rank the genes in order of supporting strength for “differential expression”.
2. Derive the p-value for each gene, based on the **NULL distribution** of the statistic.
3. Choose a **critical-value** for the gene with p-value less than which being called as “**being statistically significant**”.



Additional Information

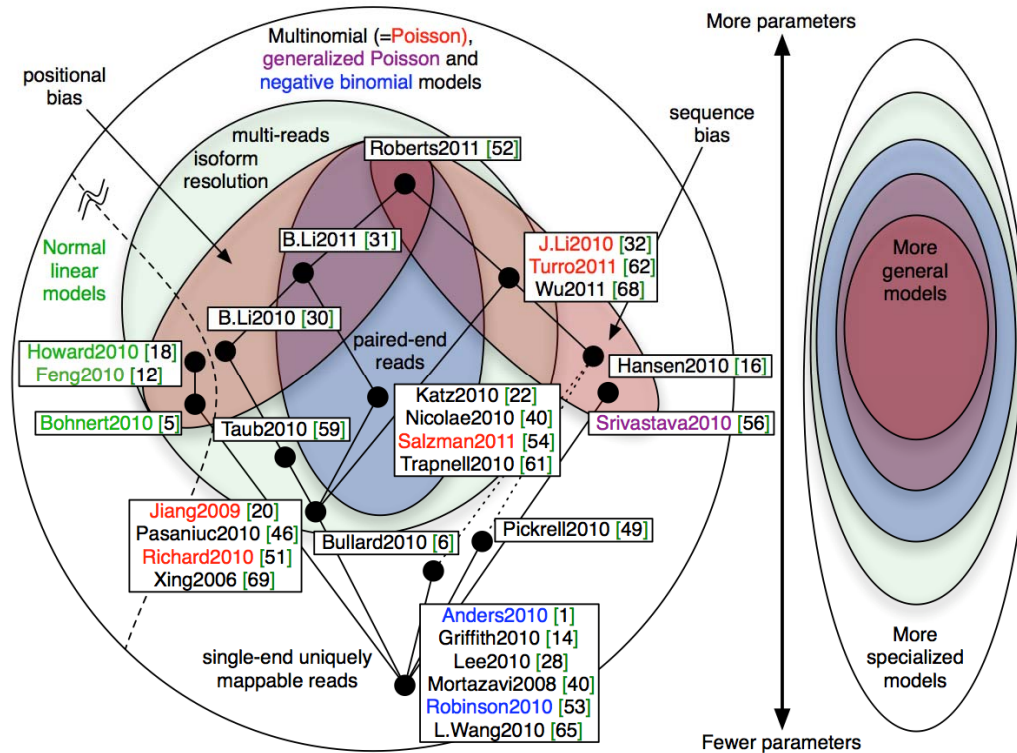
Student's *t*-distribution



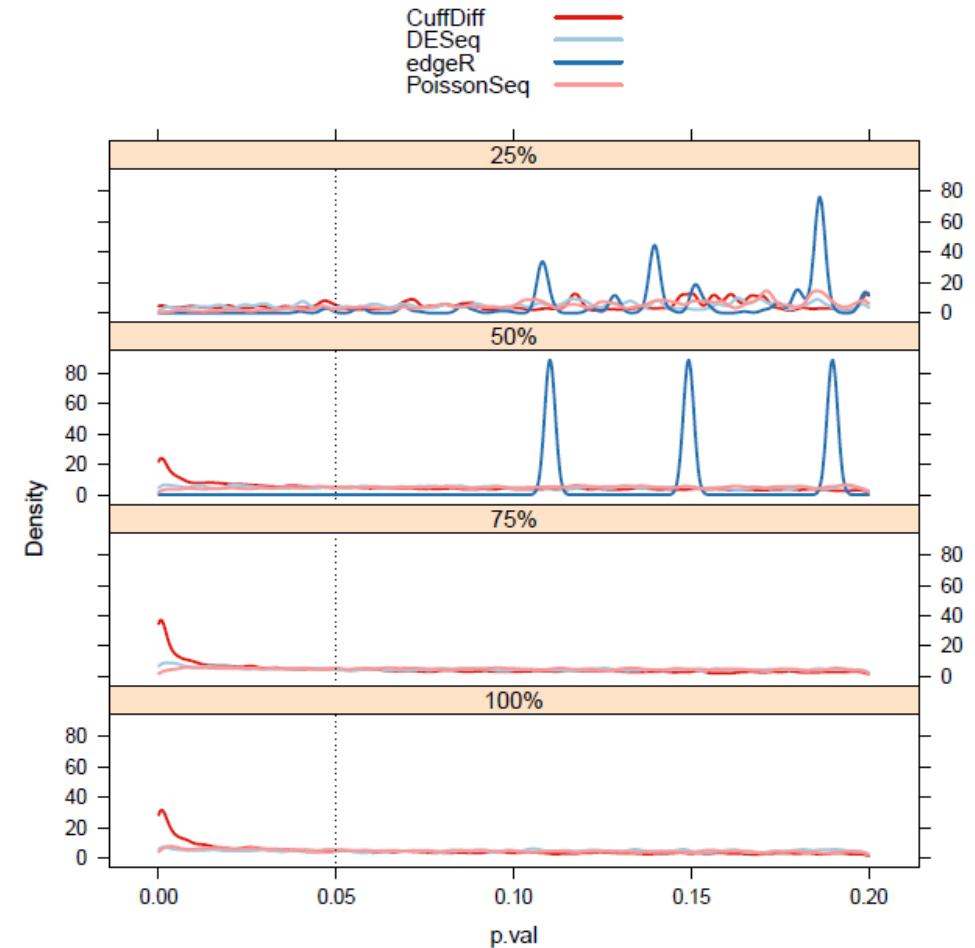
Source: [projectile.sv.cmu.edu/research/public/talks/t-test.htm](http://projectile.sv.cmu.edu/research/public/talks/t-test.htm)

- The t-test assesses whether the means of two groups are **statistically different** from each other
  - Take the **variance** into account through Standard Error (SE)
- Need to estimate the SE correctly
  - But the correct estimation depends on **prior distribution** (Normal) as well as **the number of replicates** (>10)

# Model the data in RNA-Seq



Patcher 2011, arXiv:1104.3889 [q-bio.GN]

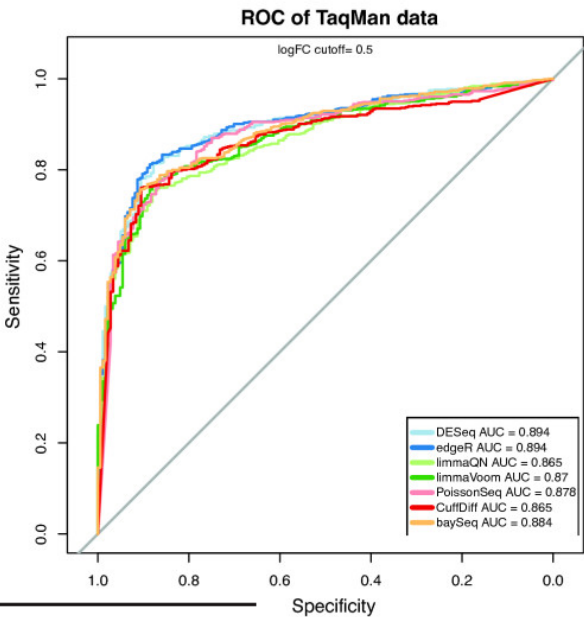


(Genome Biology 14:R95)

**METHOD** **Open Access**

# Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data

Franck Rapaport<sup>1</sup>, Raya Khanin<sup>1</sup>, Yupu Liang<sup>1</sup>, Mono Pirun<sup>1</sup>, Azra Krek<sup>1</sup>, Paul Zumbo<sup>2,3</sup>, Christopher E Mason<sup>2,3</sup>, Nicholas D Socci<sup>1</sup> and Doron Betel<sup>3,4\*</sup>



Evaluation	Cuffdiff	DESeq	edgeR	limmaVoom	PoissonSeq	baySeq
Normalization and clustering	All methods performed equally well					
DE detection accuracy measured by AUC at increasing qRT-PCR cutoff	Decreasing	Consistent	Consistent	Decreasing	Increases up to log expression change $\leq 2.0$	Consistent
Null model type I error	High number of FPs	Low number of FPs	Low number of FPs	Low Number of FPs	Low number of FPs	Low number of FPs
Signal-to-noise vs <i>P</i> value correlation for genes detected in one condition	Poor	Poor	Poor	Good	Moderate	Good
Support for multi-factored experiments	No	Yes	Yes	Yes	No	No
Support DE detection without replicated samples	Yes	Yes	Yes	No	Yes	No
Detection of differential isoforms	Yes	No	No	No	No	No
Runtime for experiments with three to five replicates on a 12 dual-core 3.33 GHz, 100 G RAM server	Hours	Minutes	Minutes	Minutes	Seconds	Hours

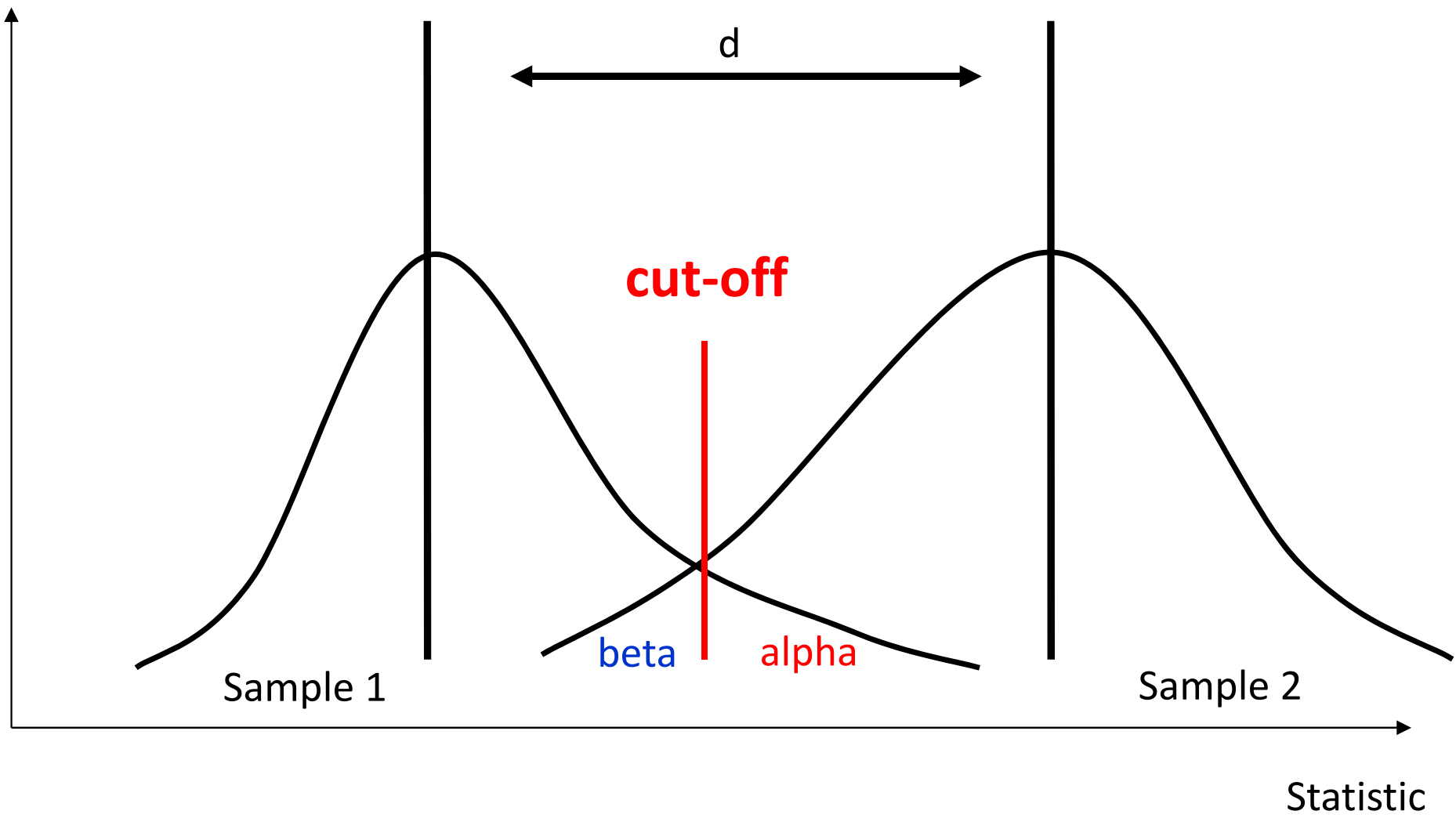
AUC, area under curve; DE, differential expression; FP, false positive.



		Hypothesis truth?	
		$H_1$ (active)	$H_0$ (inactive)
Output of statistical test	Reject $H_0$ (active)	Hit	Type I error
	Accept $H_0$ (inactive)	Type II error	Correct rejection

FUNCTIONAL MAGNETIC RESONANCE IMAGING, Figure 12.2 © 2004 Sinauer Associates, Inc.

- Type I Error (False Positive): **rejecting** the null hypothesis when **it is true**
- Type II Error (False Negative): **accepting** the null hypothesis when **it is false**



# Multiple Testing Issue

- If more than one test is made, then the collective FP value is **greater** than in the single-test
  - That is, **overall Type I error** increases
- E.g: you checked your RNA-Seq data and found 20 significantly different genes with a 0.05 threshold on each gene, then what is the chance that you making at least one error in overall?

- $\text{Pr}(\text{making a mistake}) = 0.05$
- $\text{Pr}(\text{not making a mistake}) = 1 - 0.05 = 0.95$
- $\text{Pr}(\text{not making any mistake}) = 0.95^{20} = 0.358$
- $\text{Pr}(\text{making at least one mistake}) = 1 - 0.358 = 0.642$

➔ There is a 64.2% chance of making at least one mistake

Multiple Testing Issue

# Bonferroni Correction

- Most straightforward and plain
- For  $n$  hypothesis tests, only call p-values less than  $\alpha/n$  as “being significant”.
  - Or, adjust the raw p-value as  $\min(n \cdot p, 1)$
- For example, if we want to have an experiment wide Type I error rate of 0.05 when we comparing 30000 genes, we’d need p-values less than  $0.05/30000 = 1.67 \times 10^{-6}$  so that the gene(s) could be called as “being significant”

# Type I (false positive) error rates

- Family-wise Error Rate

$$\text{FWER} = p(V \geq 1)$$

- Per-family Error Rate

$$\text{PFER} = E(V)$$

- Per-comparison Error Rate

$$\text{PCER} = E(V)/m$$

- False Discovery Rate

$$\text{FDR} = E(V/R)$$

- False Positive Rate

**Proportion** of false positives among the genes that are flagged as differentially expressed.

	# not rejected	# rejected	totals
# true H	U	V (False Positive)	$m_0$
# non-true H	T (False Negative)	S	$m_1$
totals	$m - R$	R	$m$

# q-value

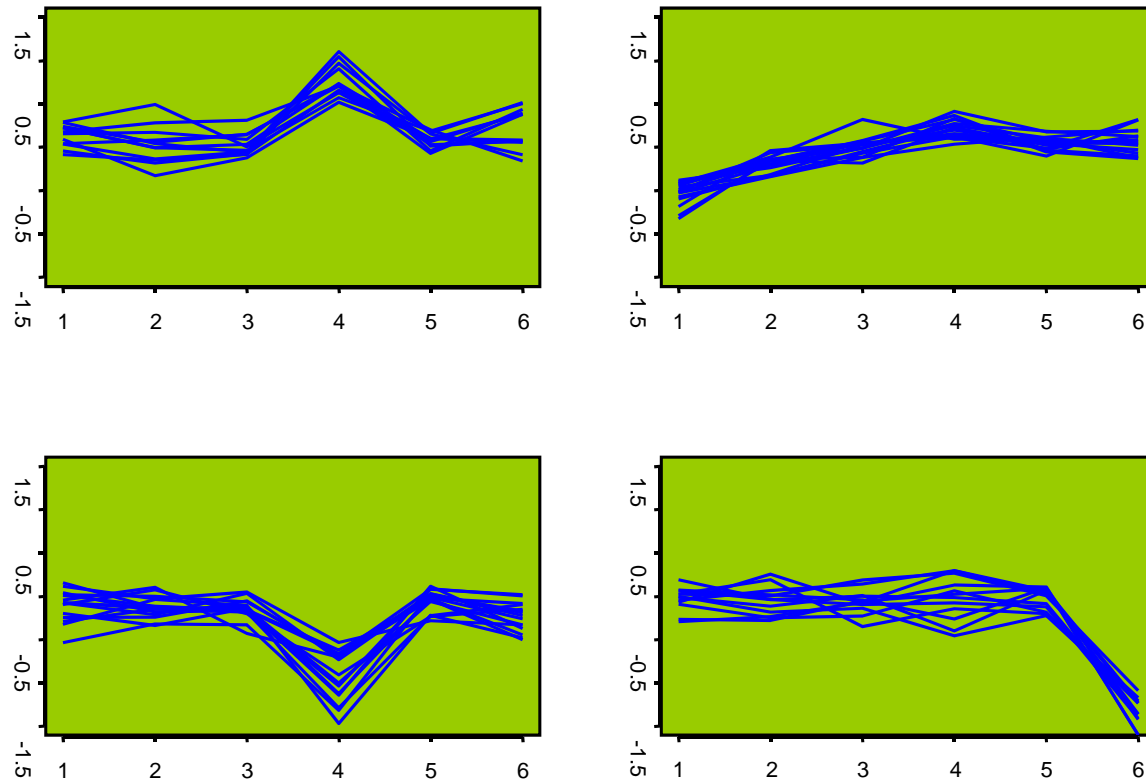
- **q-value** is an measure of False Discovery Rate (FDR)
  - Proposed by Storey *et al.* in 2002 and tuned for microarray analysis
- The **q-value** for a particular gene  $g$  is the **expected proportion** of false positives incurred when calling that gene  $g$  “significant”.
- In contrast, the **p-value** for a particular gene  $g$  is the **probability that a randomly generated expression profile** would be as or more extremely differentially expressed.

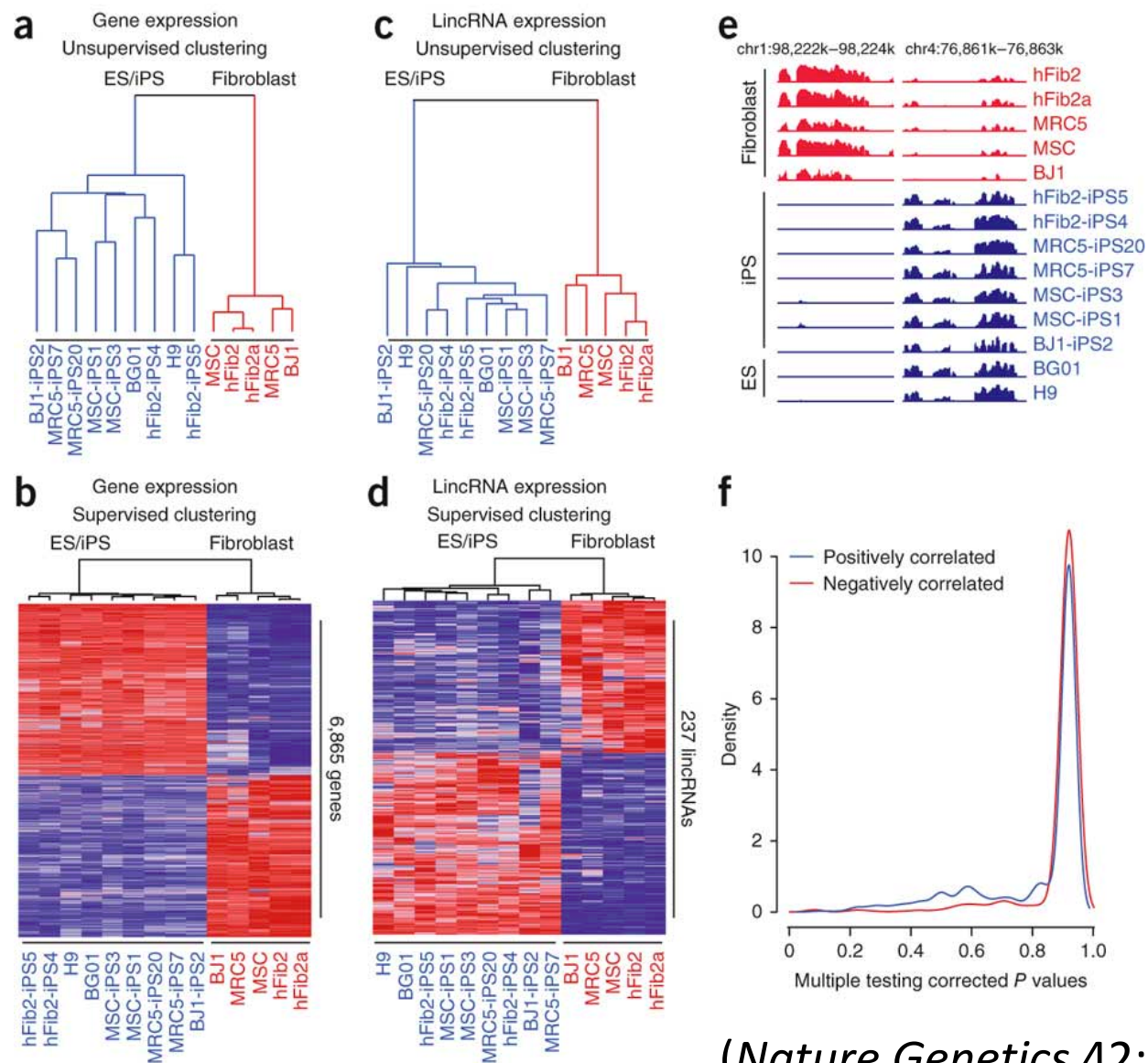


- Differentially expressed genes
- Co-expressed genes

**Clustering**: Group cases (genes/samples) with similar expression pattern/levels (**Unsupervised learning**)

- Hierarchical Cluster, k-mean Cluster, Self-Organizing Maps (SOM), etc





(*Nature Genetics* 42:1113)

Copyright © Peking University

**Distance measurement**: how “similar” between two genes’ profile

Euclidean distance  
(Absolution distance)

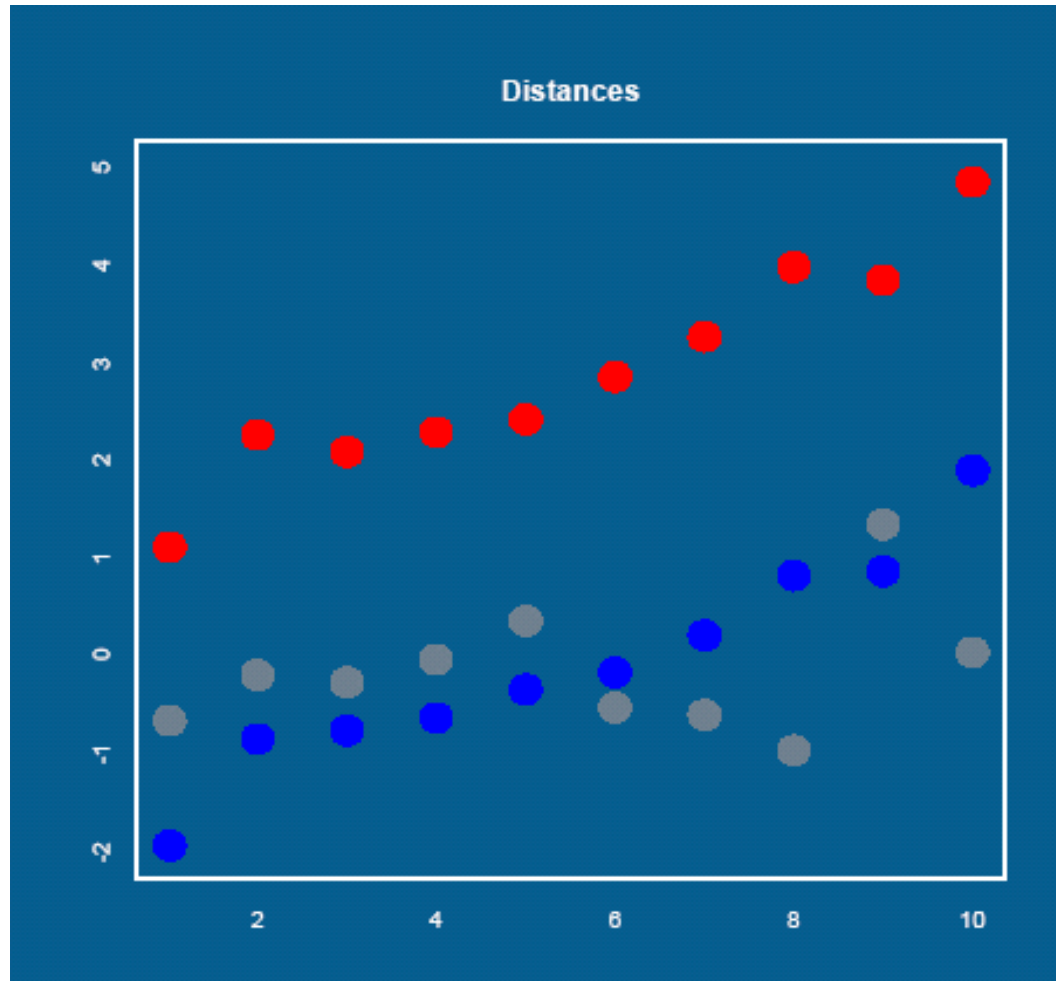
$$s(x_1, x_2) = \sqrt{\sum (x_{1k}^2 - x_{2k}^2)}$$

Pearson distance  
(Correlation distance)

$$s(x_1, x_2) = \frac{\sum_{k=1}^K (x_{1k} - \bar{x}_1)(x_{2k} - \bar{x}_2)}{\sqrt{\sum_{k=1}^K (x_{1k} - \bar{x}_1)^2 \sum_{k=1}^K (x_{2k} - \bar{x}_2)^2}}$$

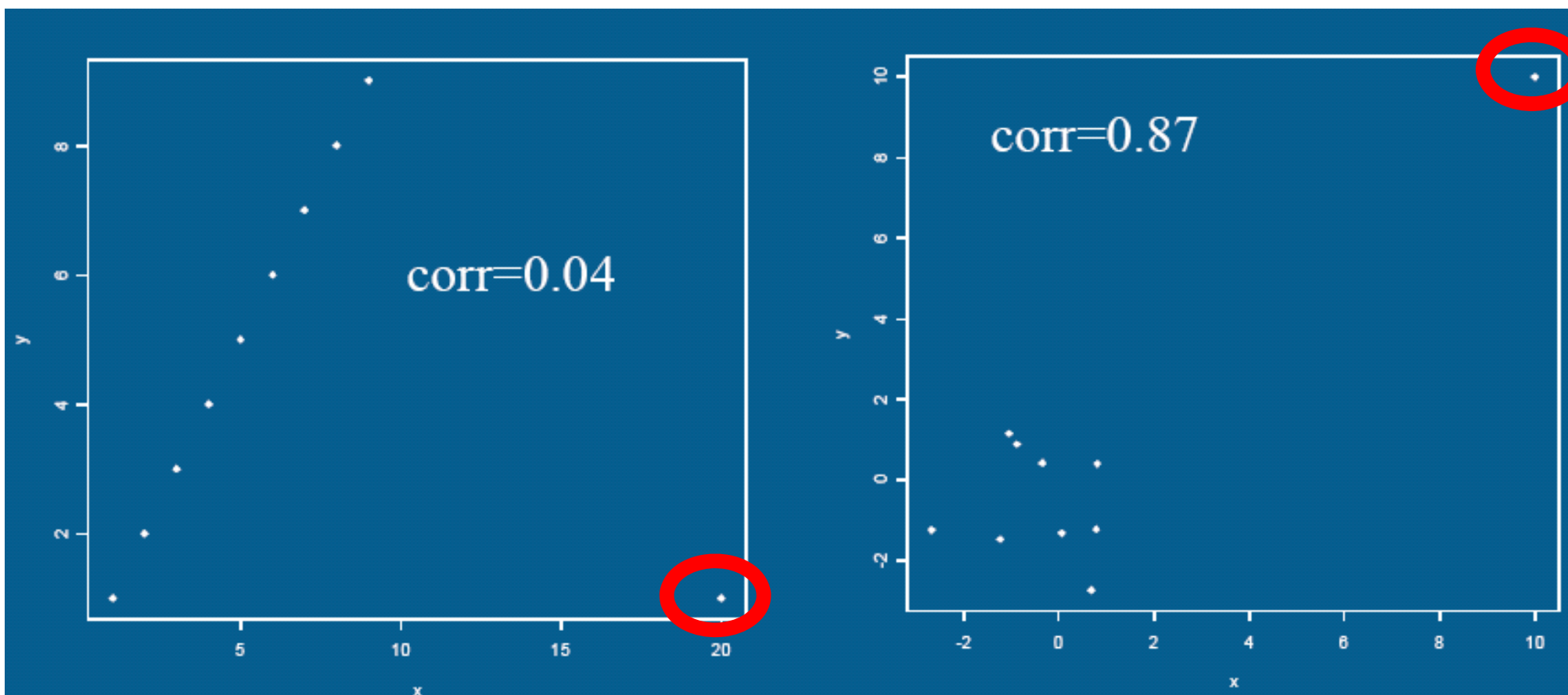
Pearson Distance:

- red-blue: .006
- red-gray: .768
- blue-gray: .7101



Eucl. Distance:

- red-blue: 9.45
- red-gray: 10.26
- blue-gray: 3.29



# Summary Question

- Do you think the classical t-test could be used in differential expression calling? Explain.
- Spearman coefficient is a more robust correlation measurement than Pearson coefficient. Could you write the distance formula with Spearman coefficient?



# 生物信息学：导论与方法

## Bioinformatics: Introduction and Methods



<https://www.coursera.org/course/pkubioinfo>