

二代真核转录组

王鹏

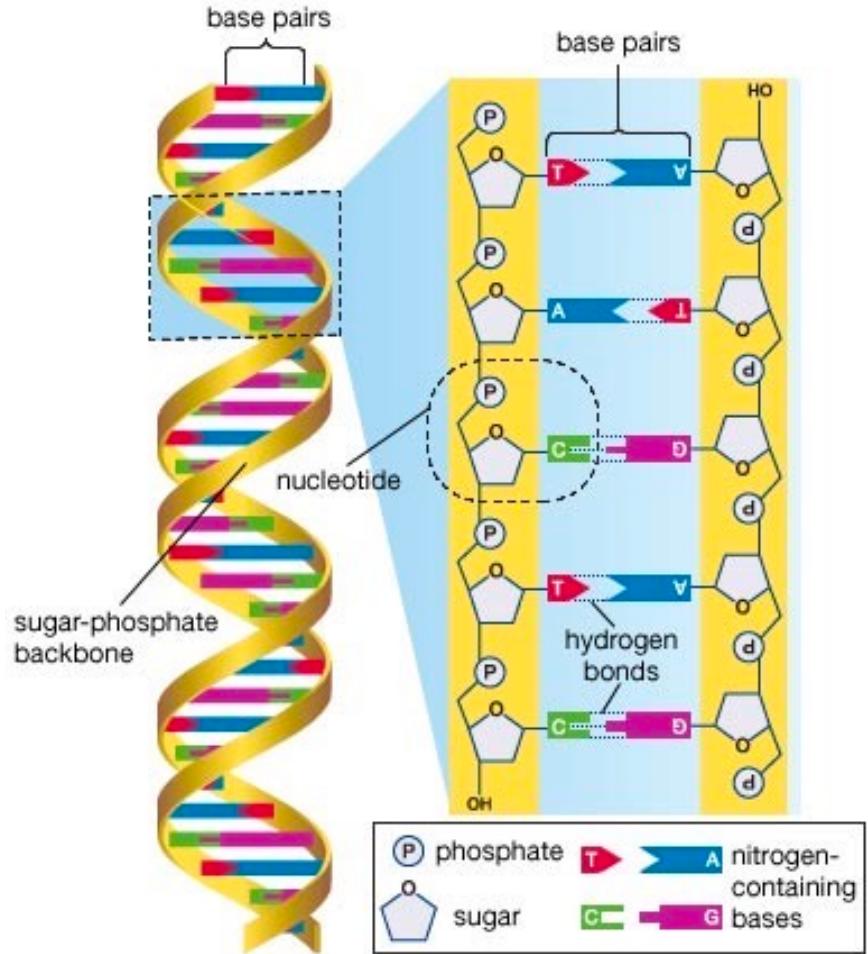
科技服务技术支持部

1



Illumina测序原理

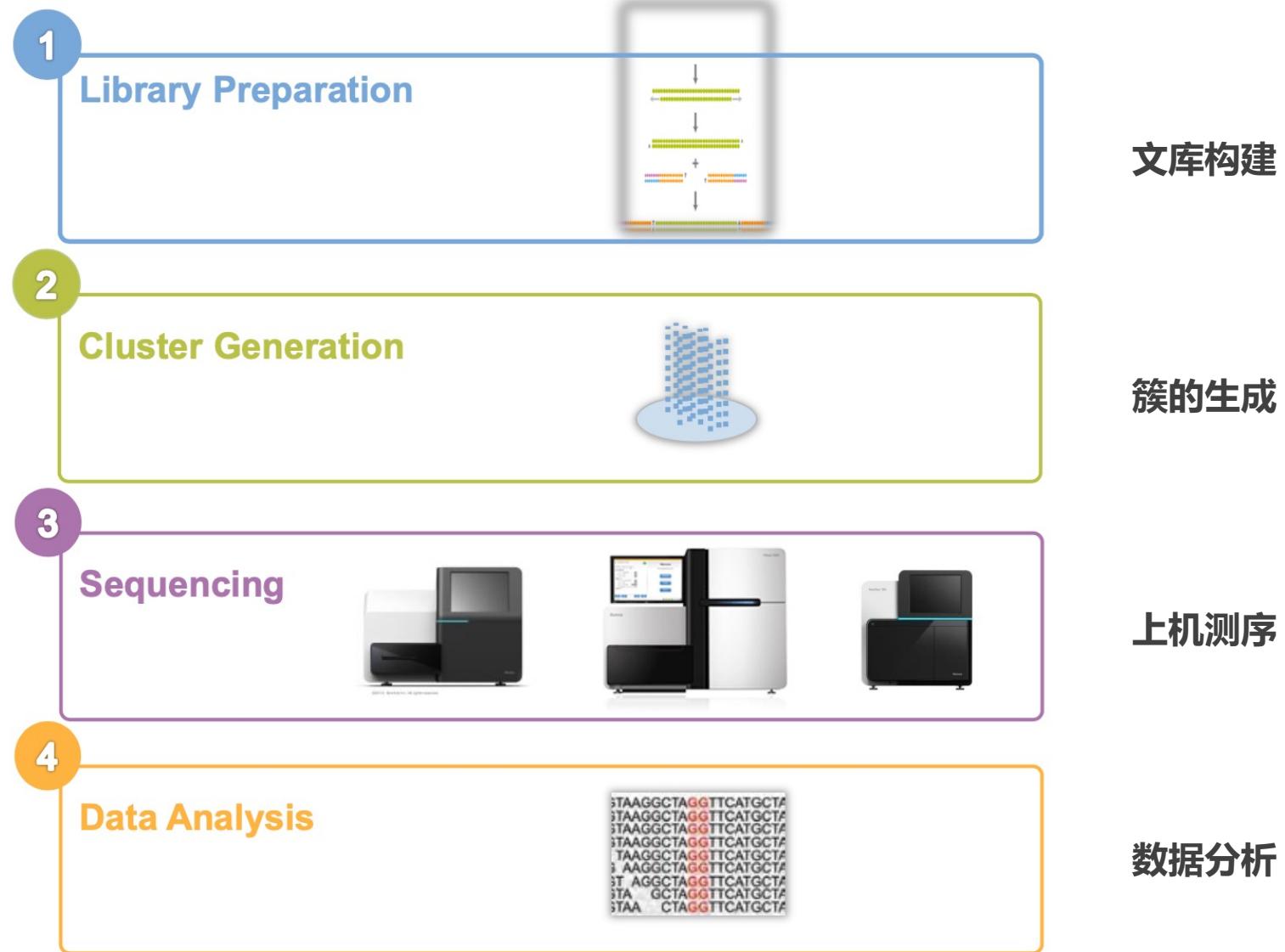
什么是测序



DNA测序：不可见的未知序列，转换成可见的ATGC
序列分析：可见的ATGC，转换成机器可读的0和1
统计绘图：机器识别的0和1，转换成人可读的图表

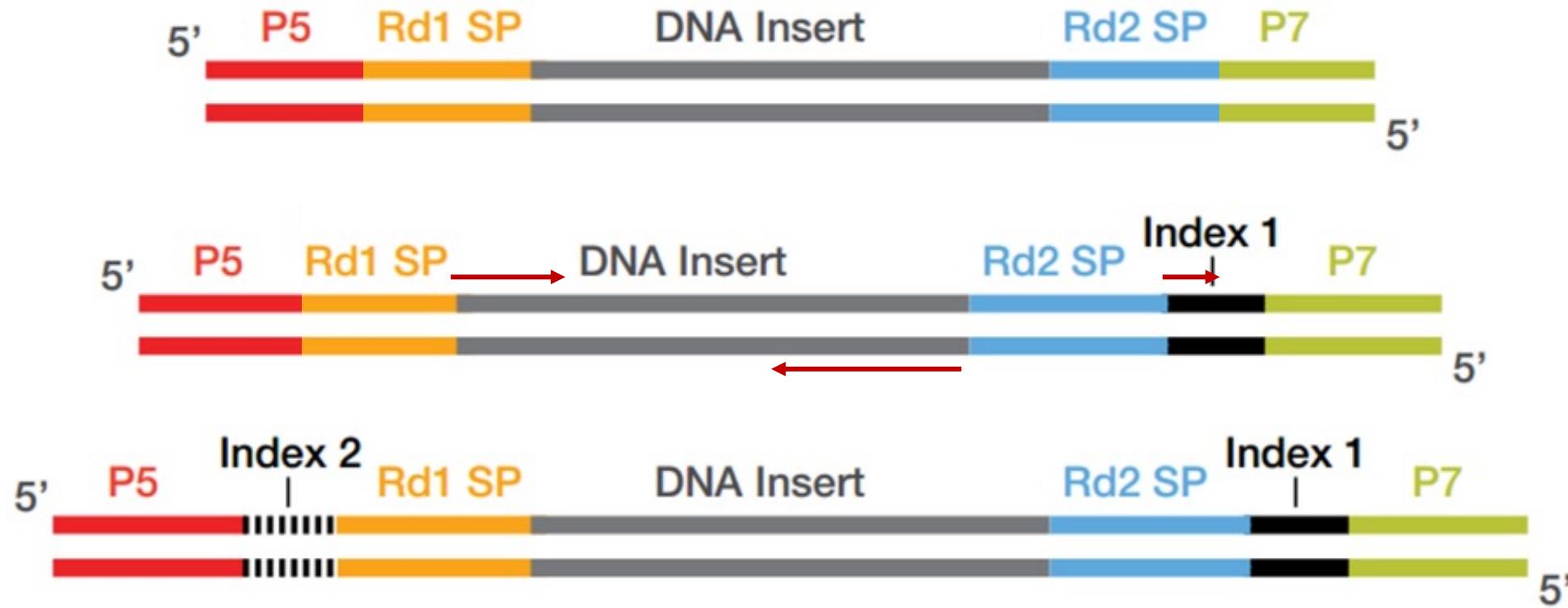
Illumina测序步骤

一代测序四步骤



Illumina文库结构

测序就像走路，首先要确定行程目标，其次要有个引路人



P5/P7接头：行程的起始和终止；Rd1SP：行程的引路人；Index：区分不同行人的标签

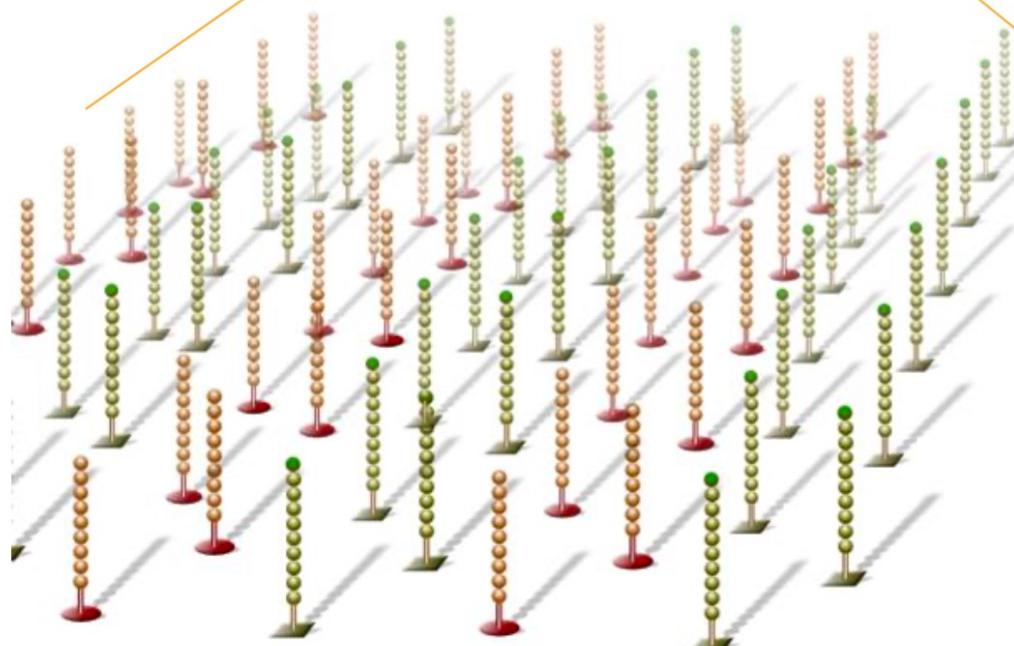
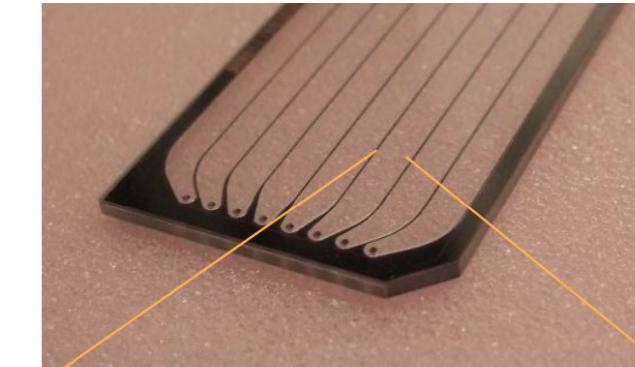
Illumina flowcell

走路要在一定区域行走，不能乱走，flowcell就类似行走的固定区域或载体

Cluster generation occurs on a flow cell

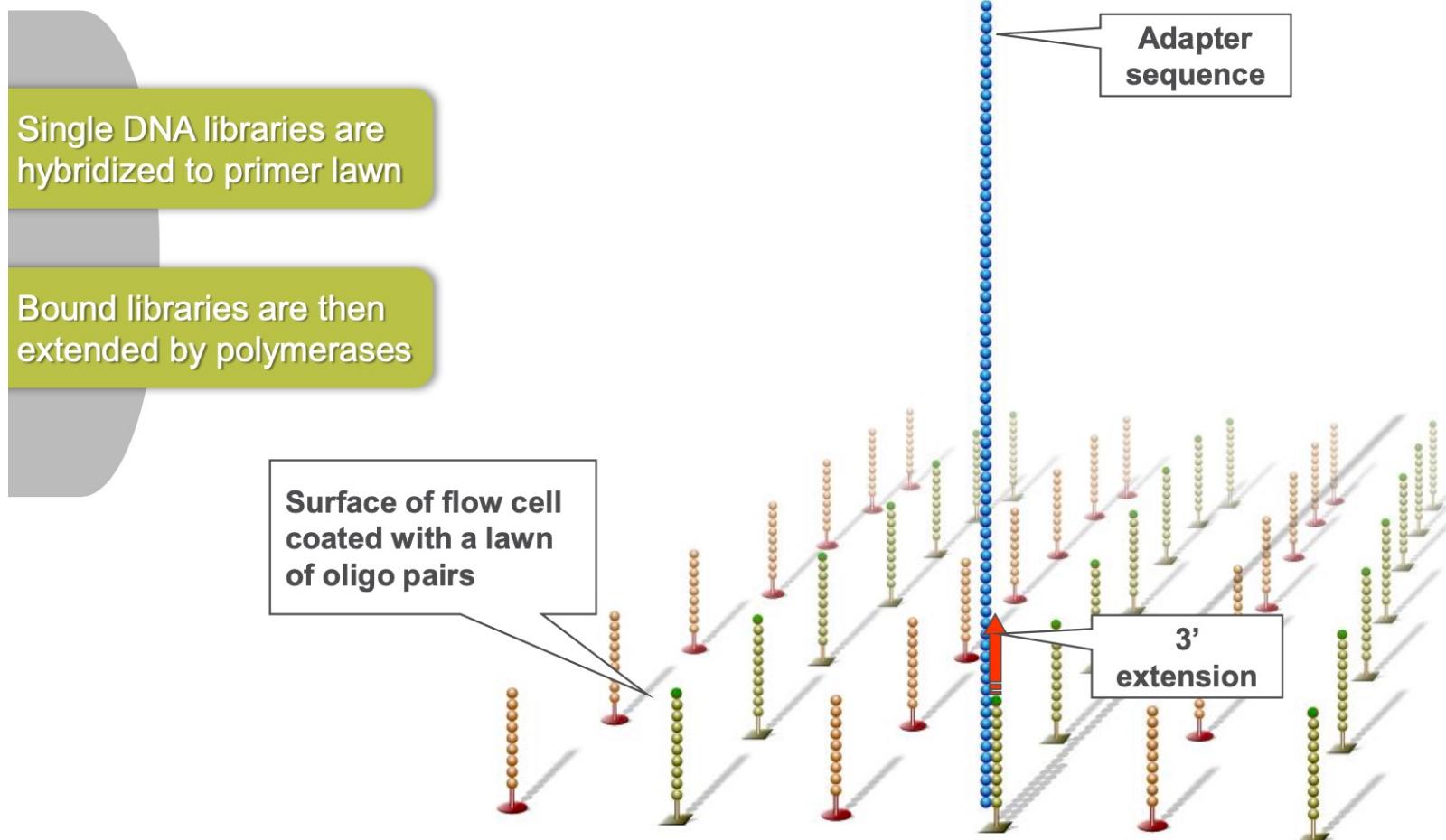
A flow cell is a thick glass slide with channels or lanes

Each lane is randomly coated with a lawn of oligos that are complementary to library adapters



片段杂交&延伸

Hybridize Fragment & Extend



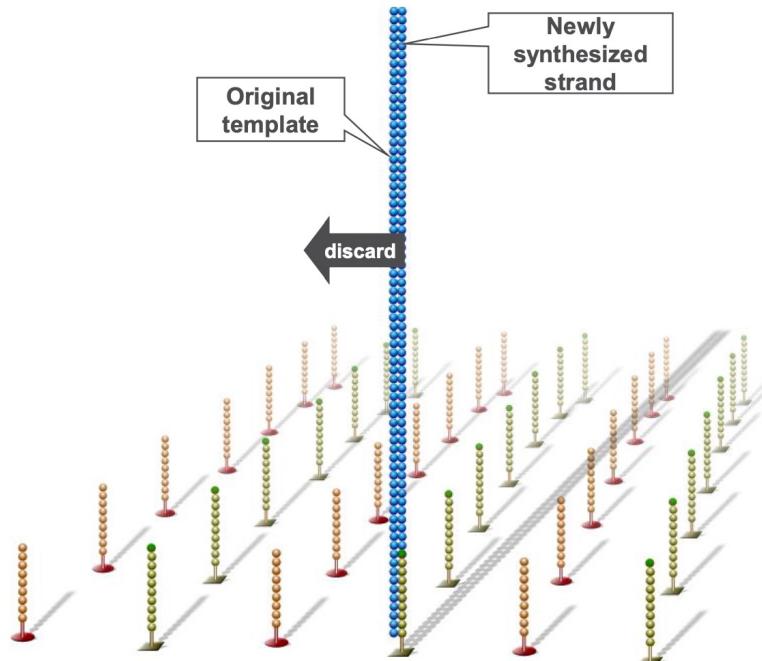
双链DNA变性丢掉原始DNA链

Denature Double-Stranded DNA

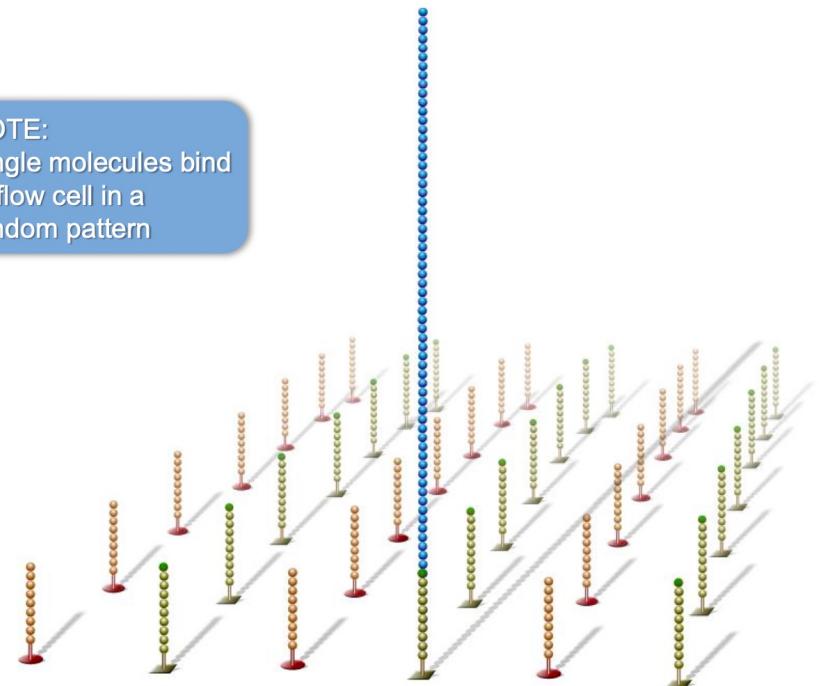
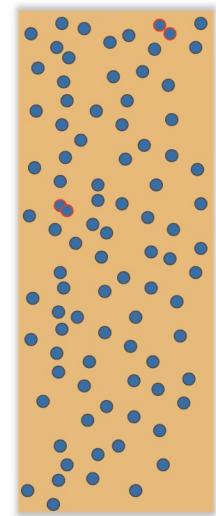
Double-stranded molecule is denatured

Original template washed away

Newly synthesized strand is covalently attached to flow cell surface



Single-Stranded DNA

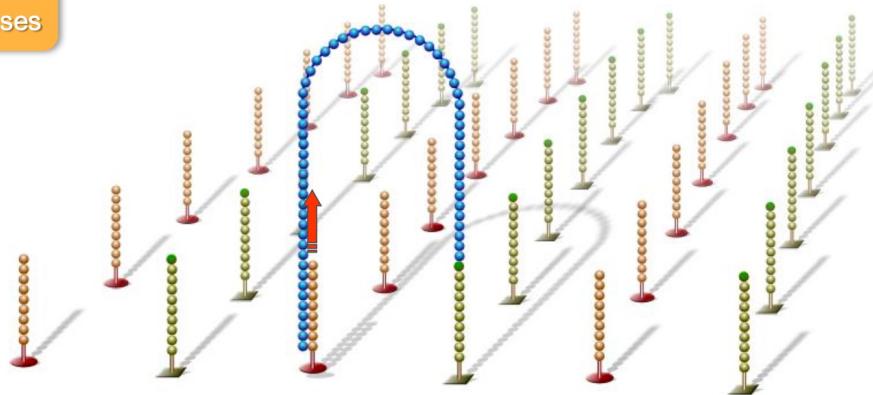


桥式PCR

Bridge Amplification

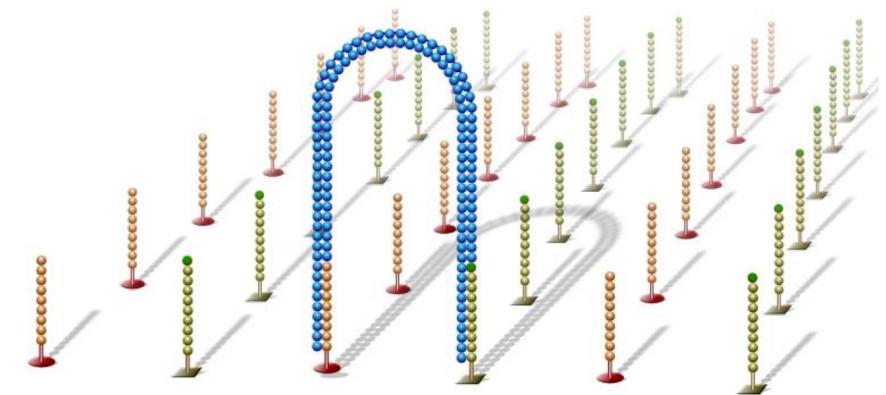
Single-stranded molecule flips over
and forms a bridge by hybridizing to
adjacent, complementary primer

Hybridized primer is
extended by polymerases



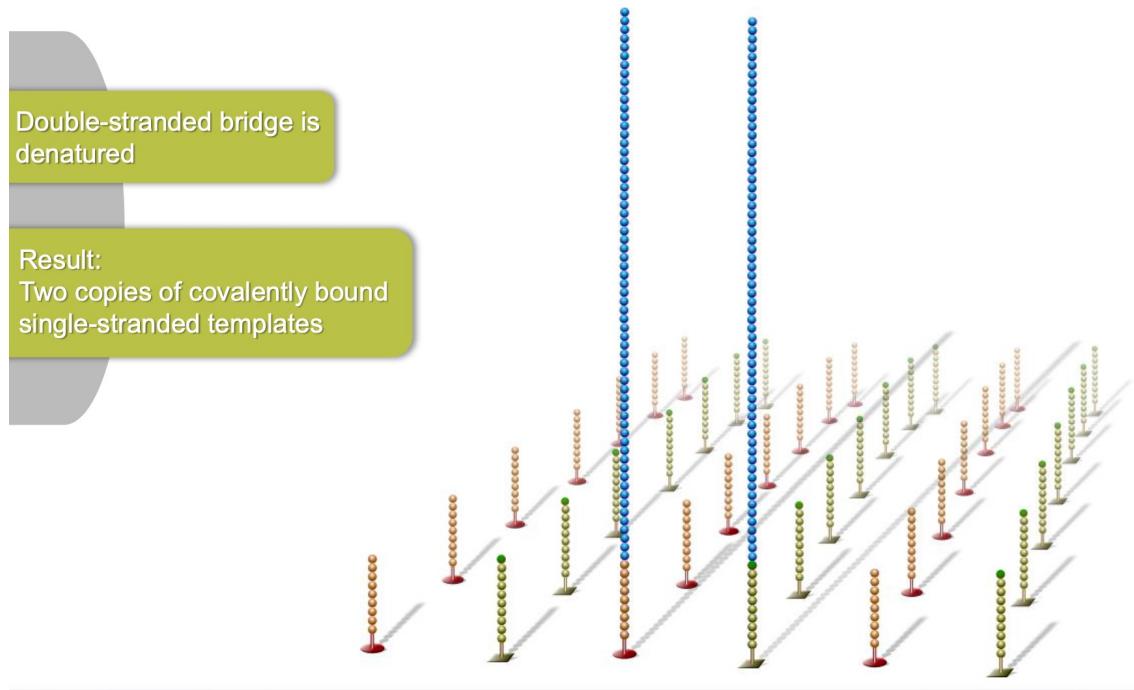
Bridge Amplification

Double-stranded bridge is formed

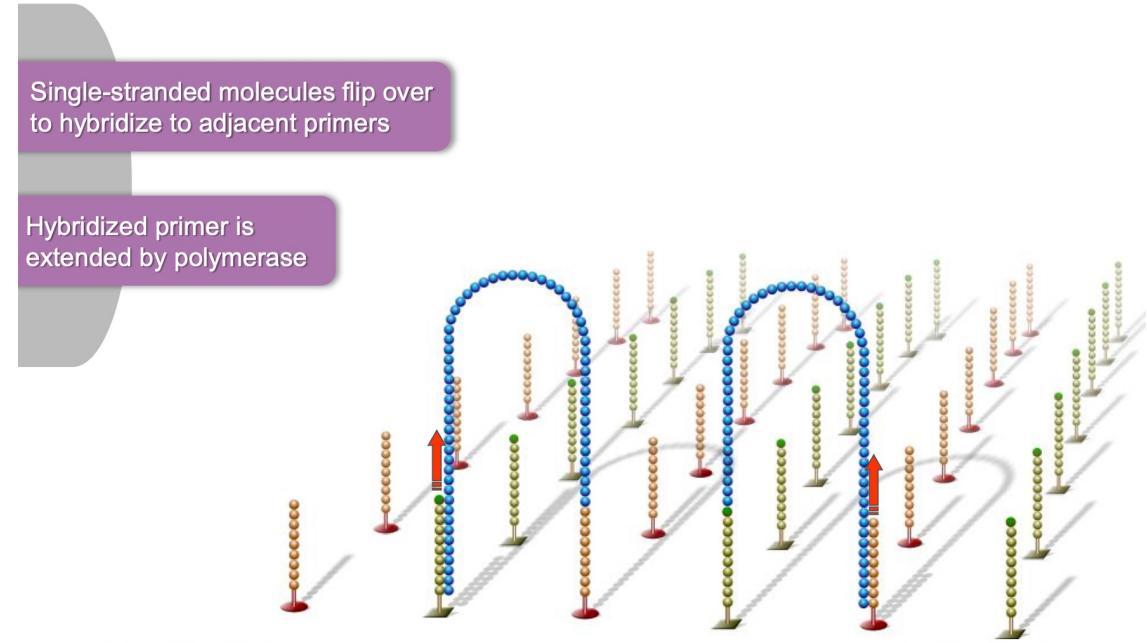


桥式PCR

Denature Double-Stranded Bridge

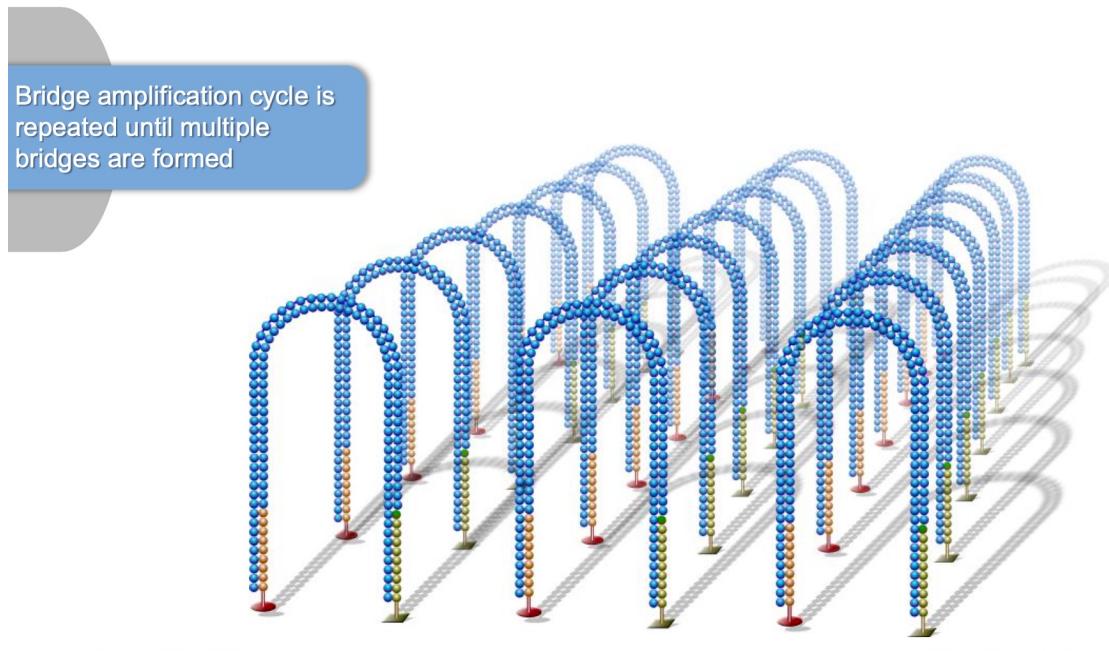


Bridge Amplification

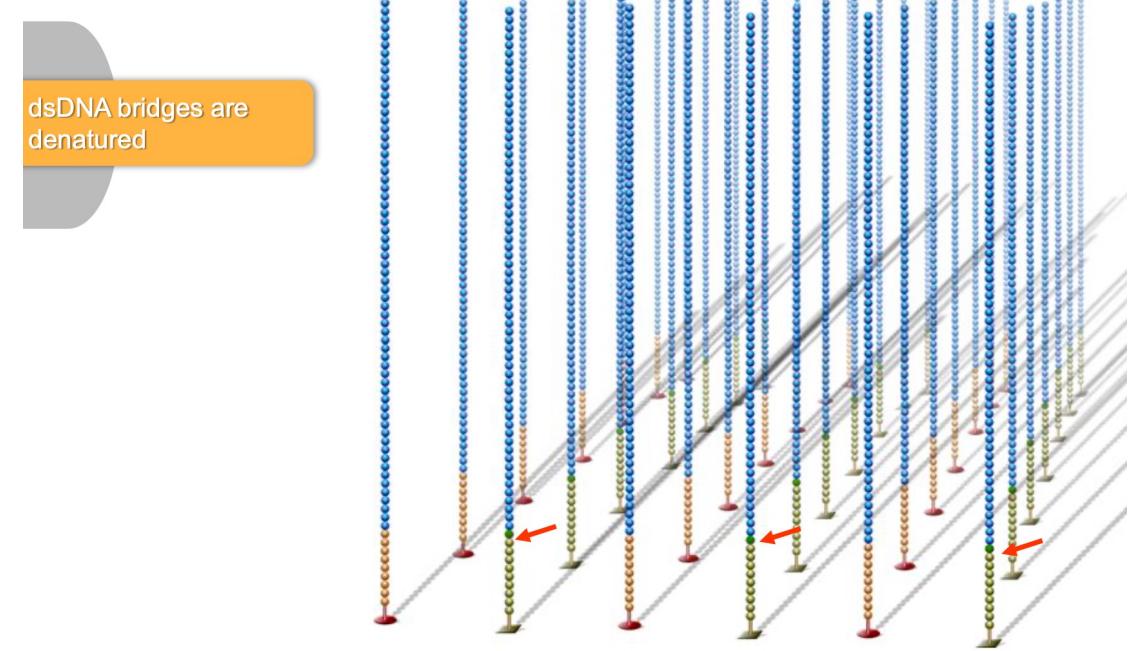


桥式PCR

Bridge Amplification

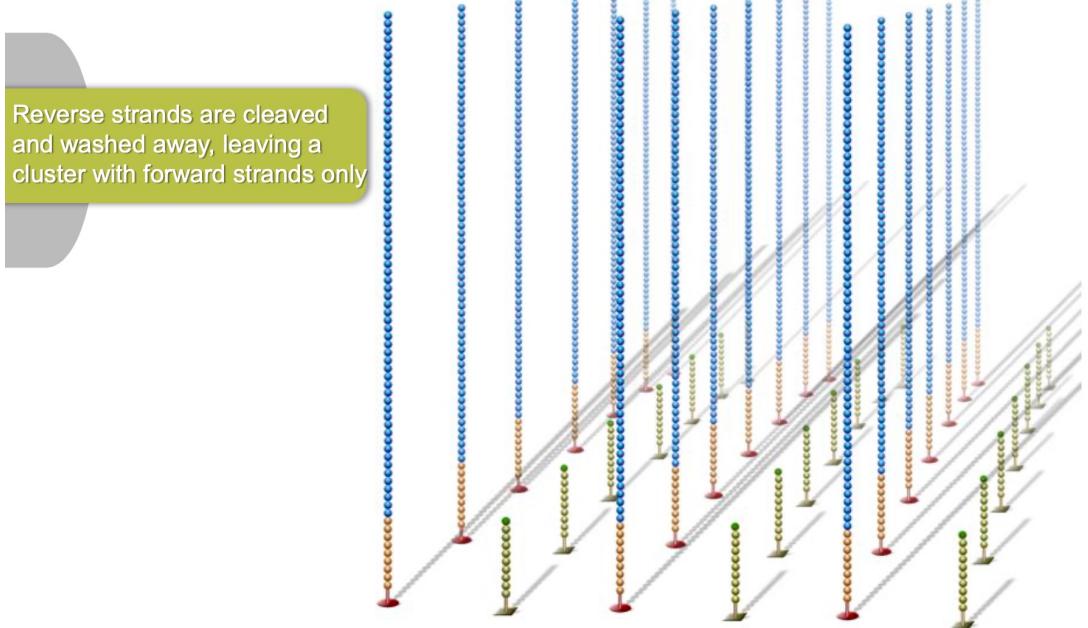


Linearization

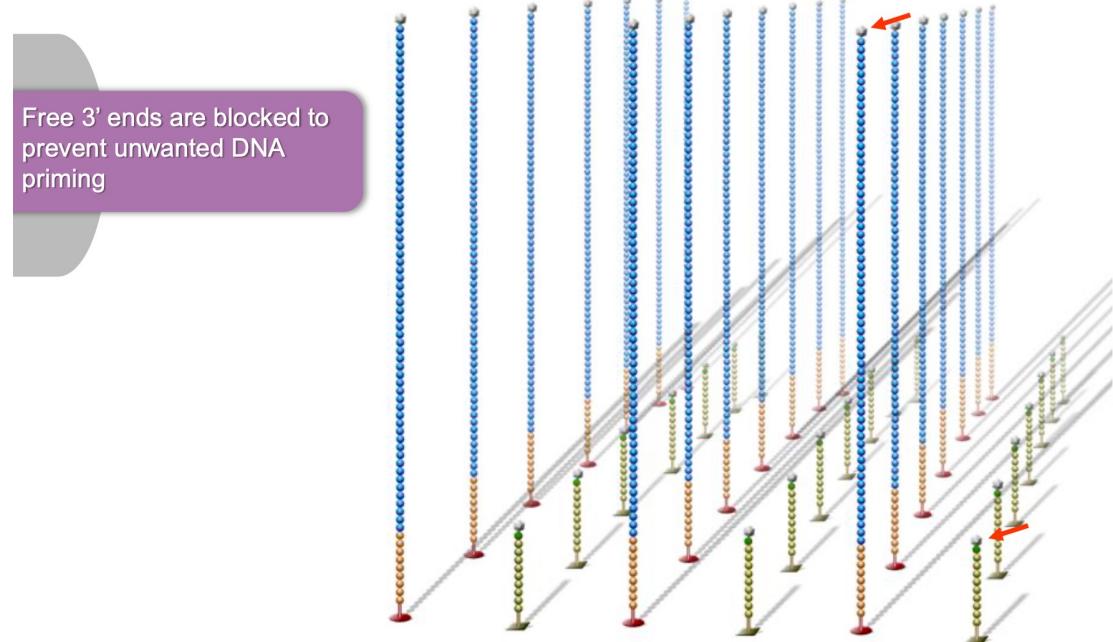


洗脱第二链，封闭3端

Reverse Strand Cleavage

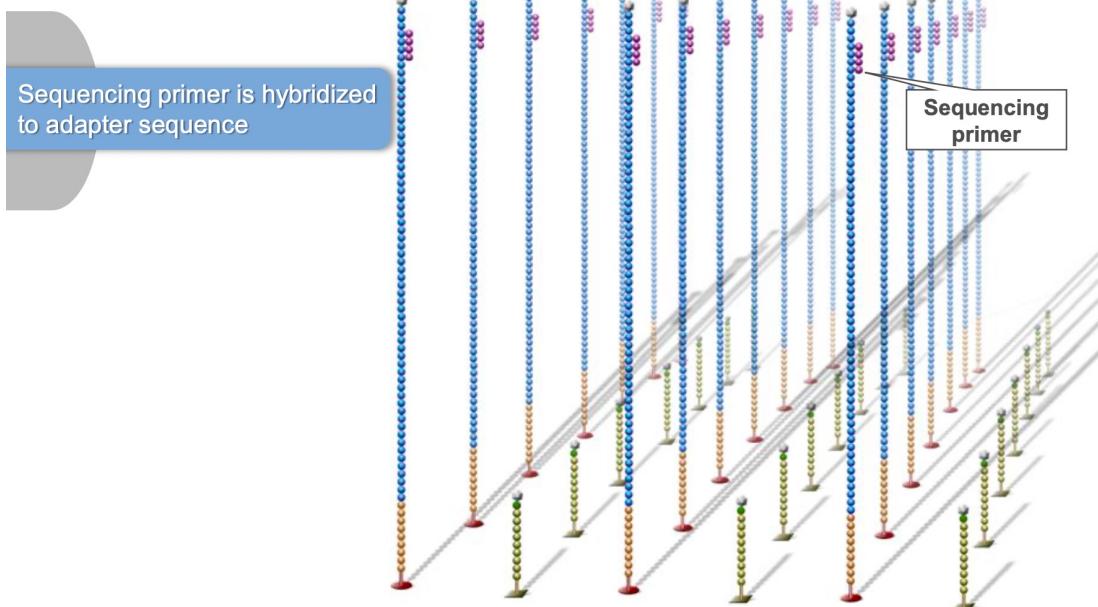


Blocking

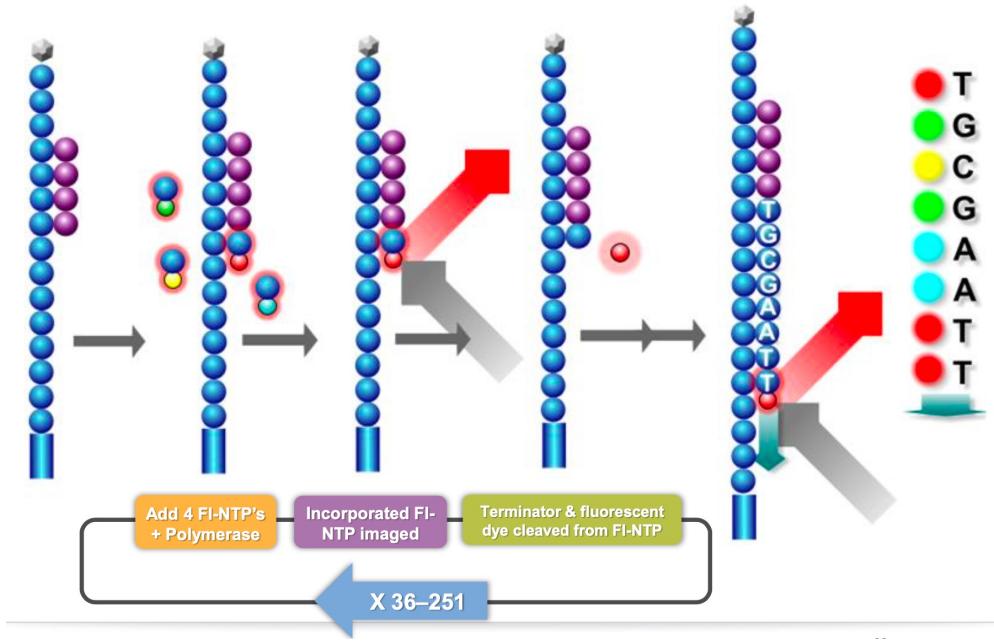


加入引物1，边合成边测序

Read 1 Primer Hybridization

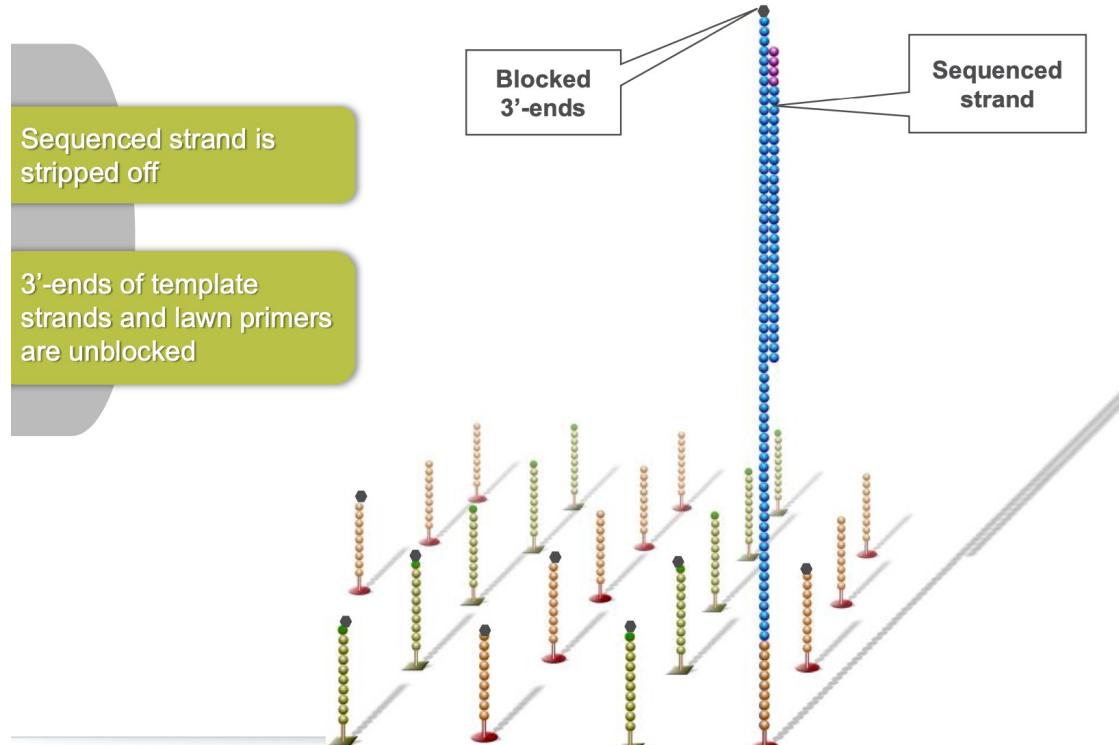


Sequencing By Synthesis (SBS)

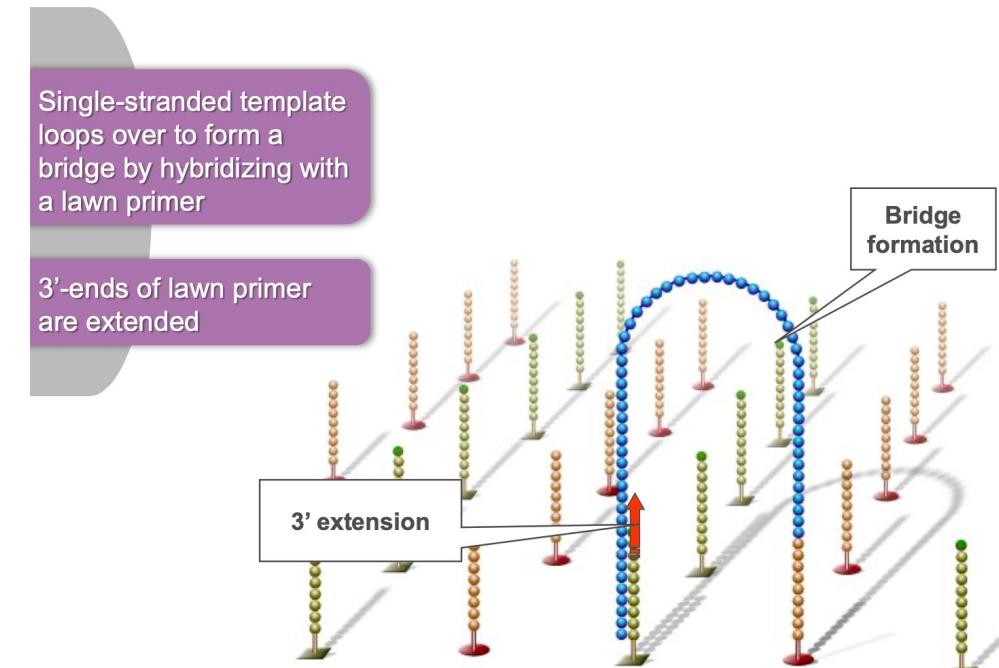


2端测序

Paired-End Sequencing

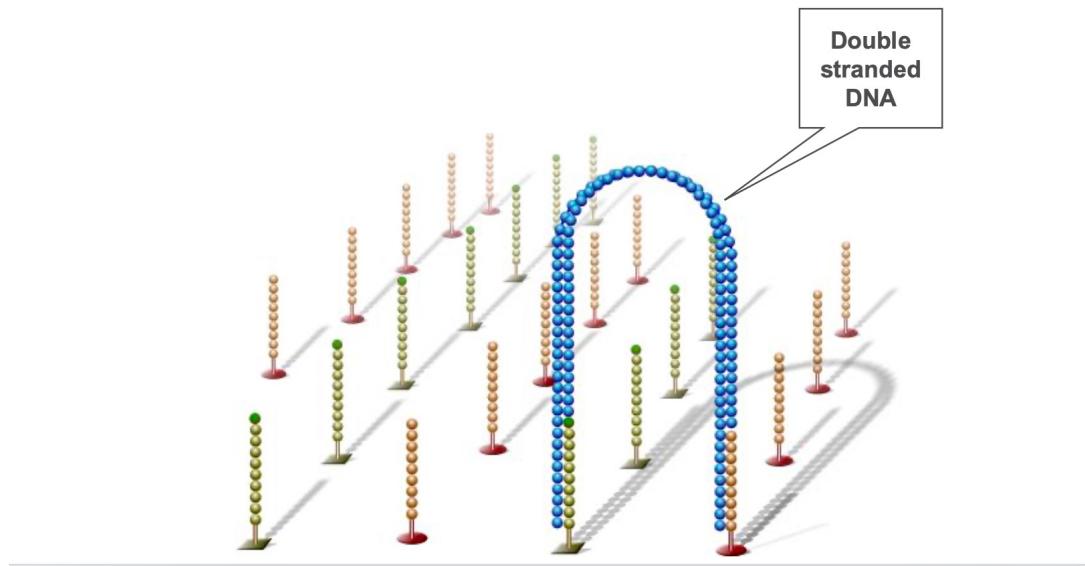


Paired-End Sequencing

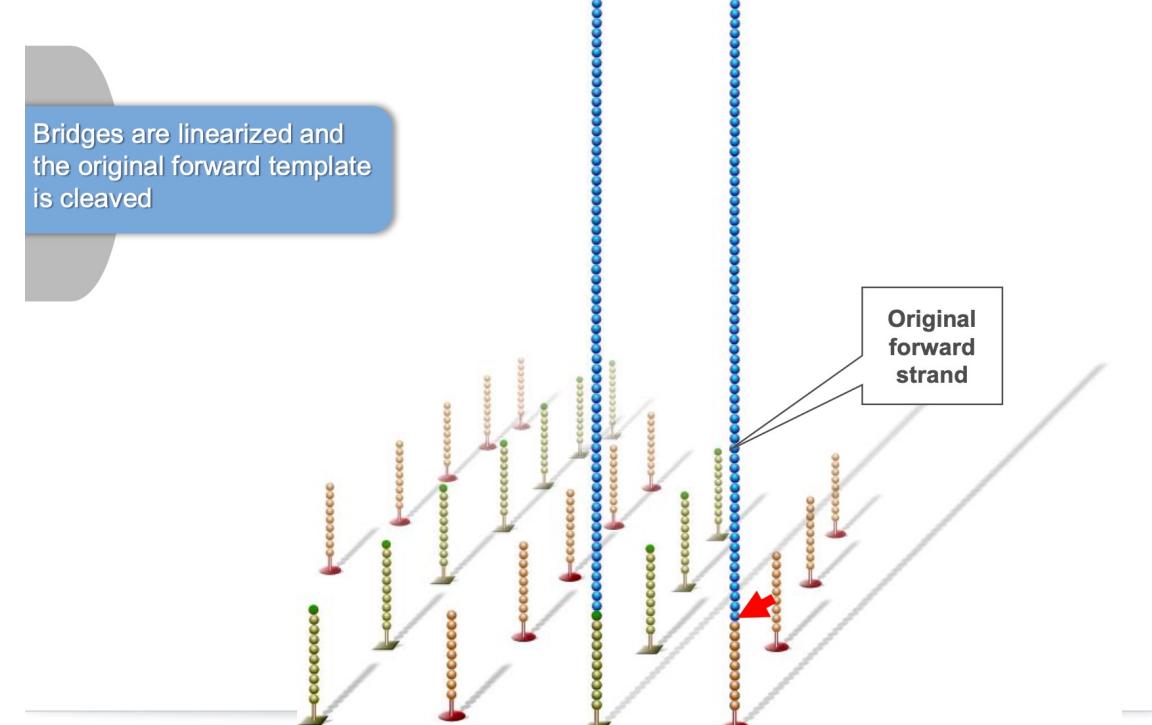


2端测序

Paired-End Sequencing



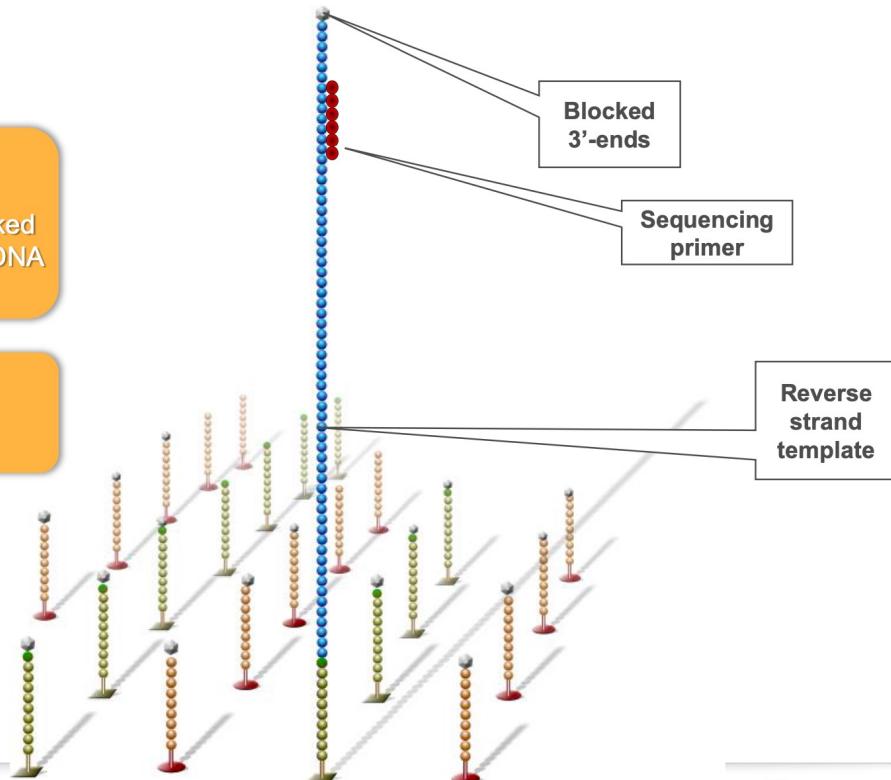
Paired-End Sequencing



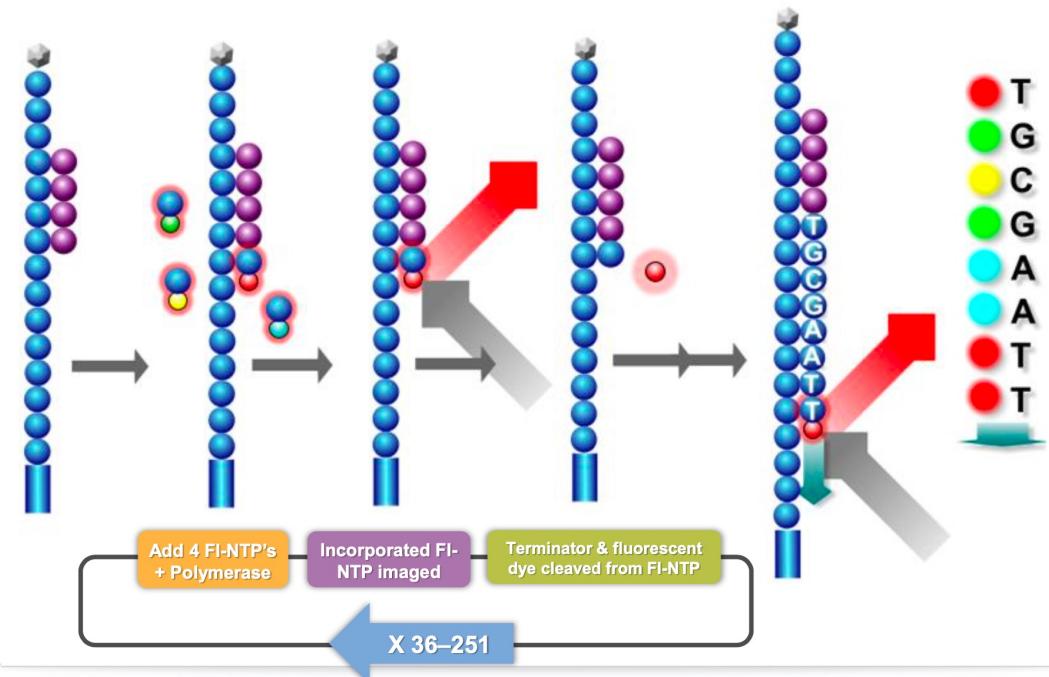
2端测序

Paired-End Sequencing

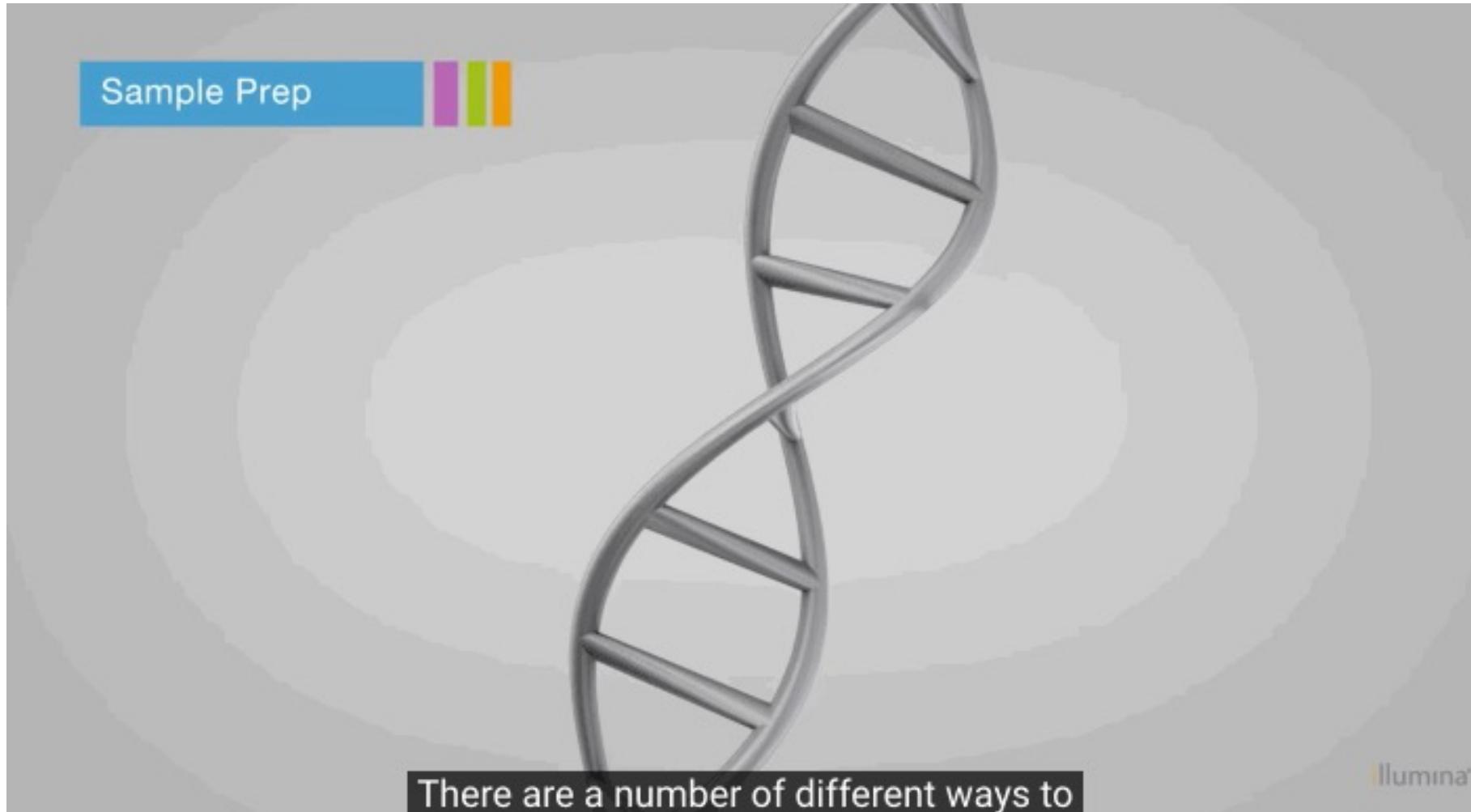
Free 3' ends of the reverse template and lawn primers are blocked to prevent unwanted DNA priming
Sequencing primer is hybridized to adapter sequence



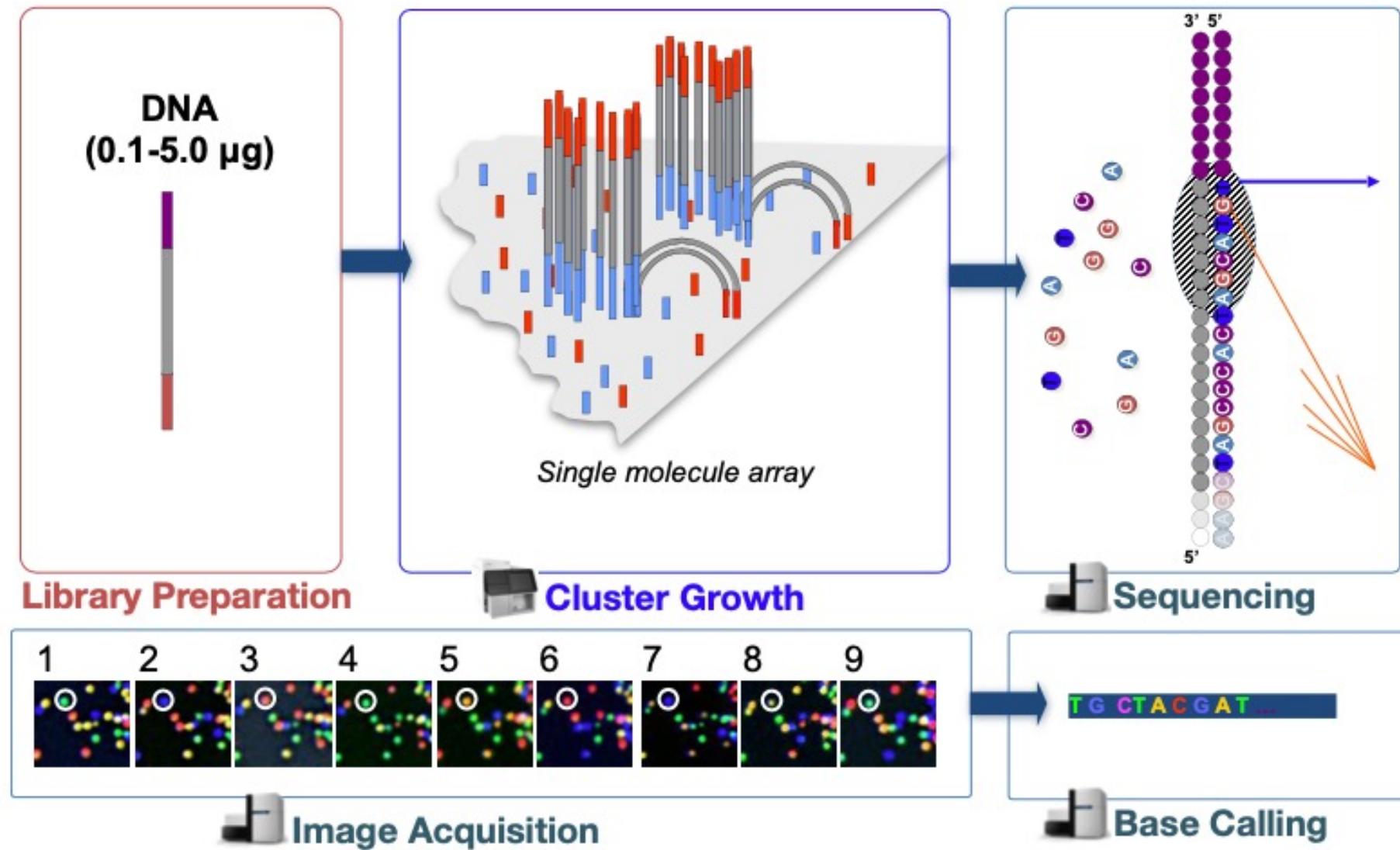
Sequencing By Synthesis 2nd Read



Illumina测序动态图



测序示意图

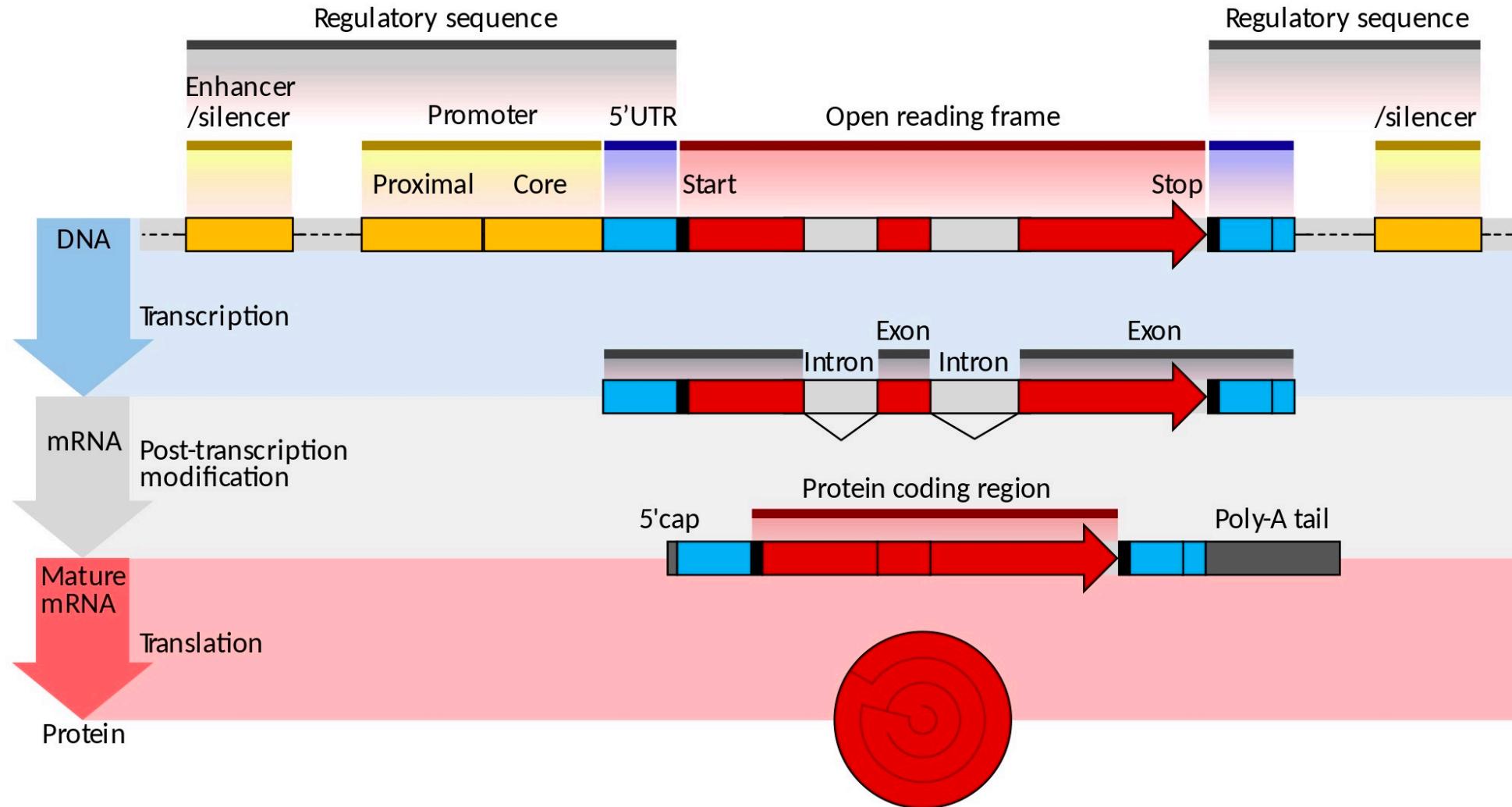


2

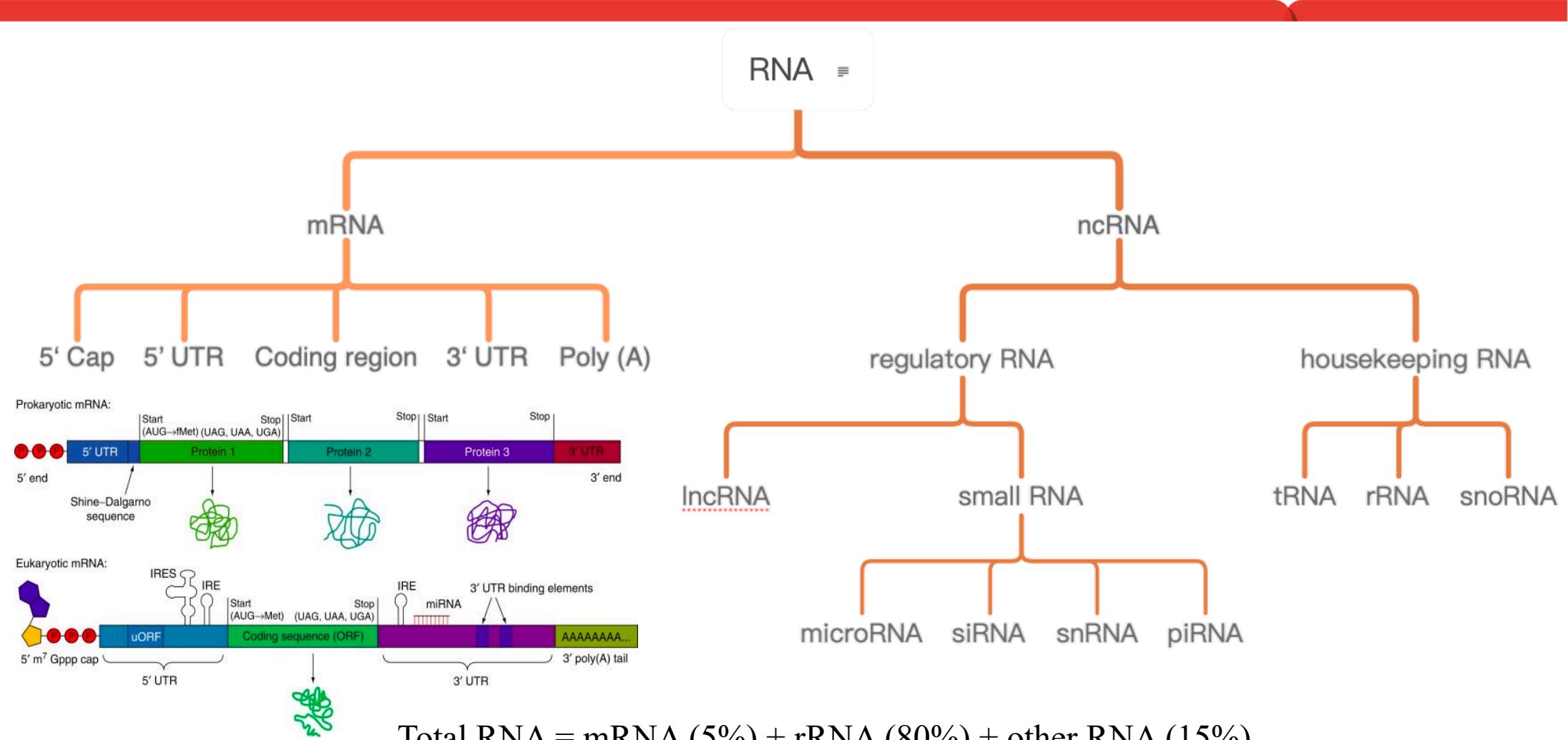


RNA常见的文库类型

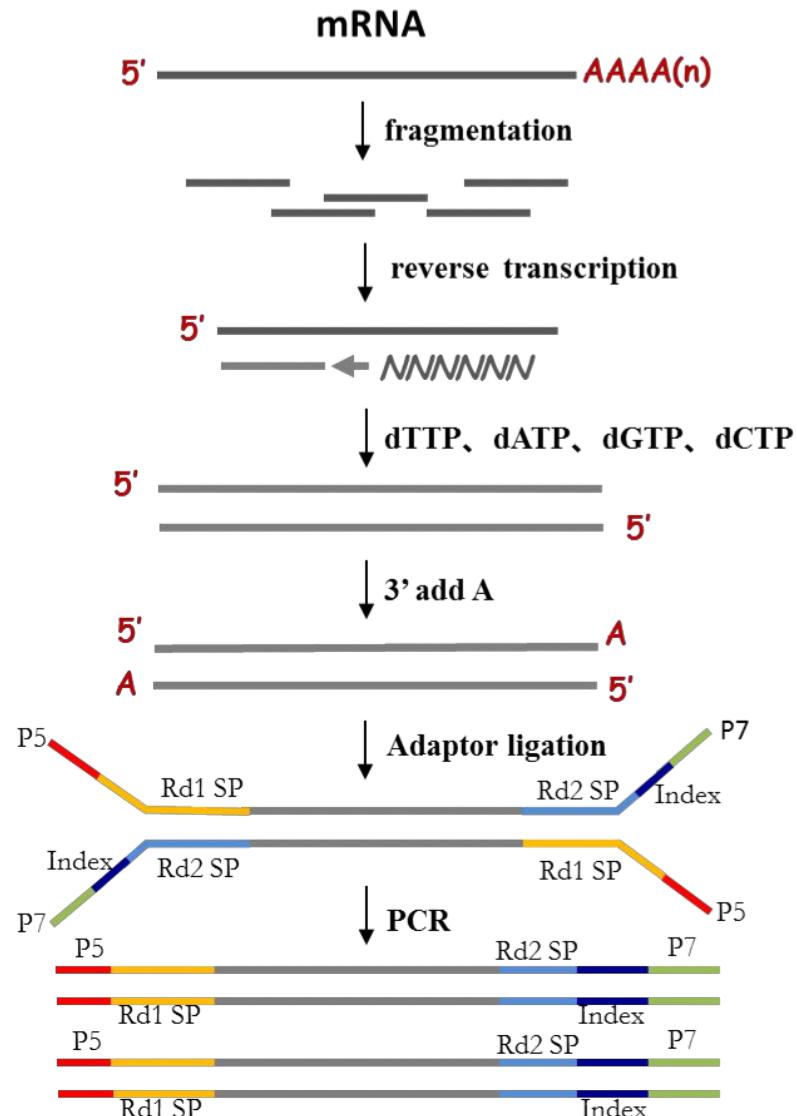
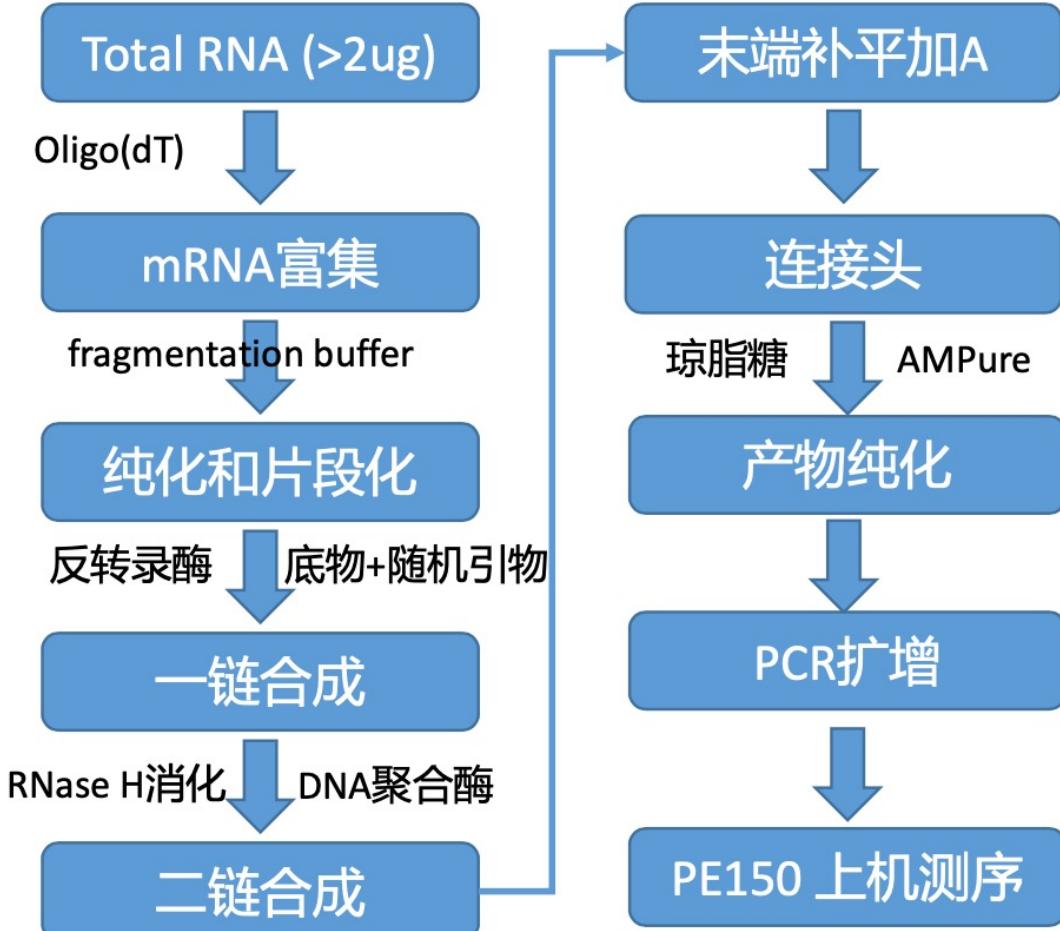
从DNA到mRNA



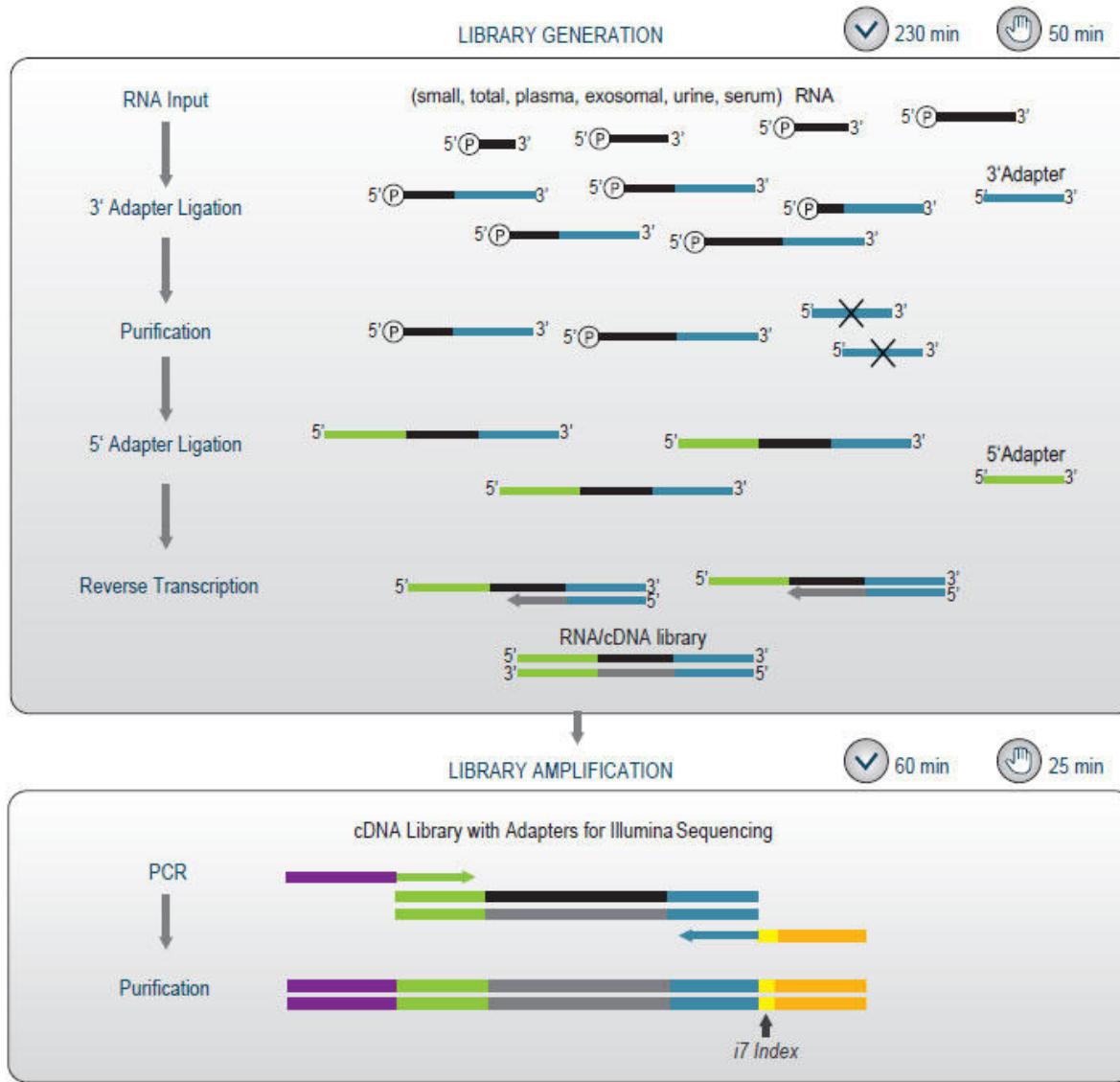
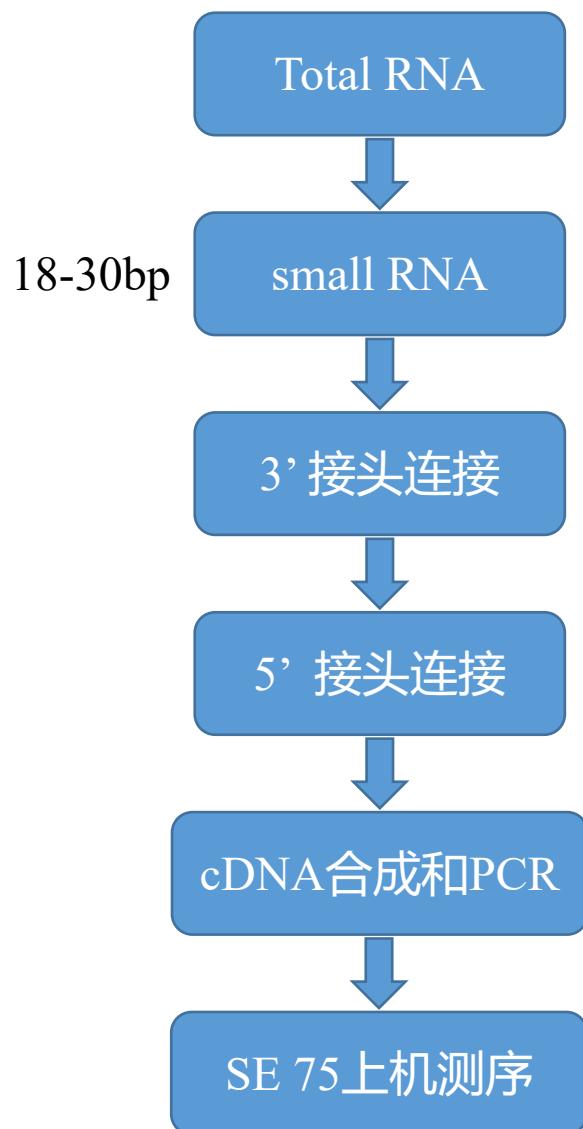
RNA分类



RNA-seq文库构建



Small RNA文库构建

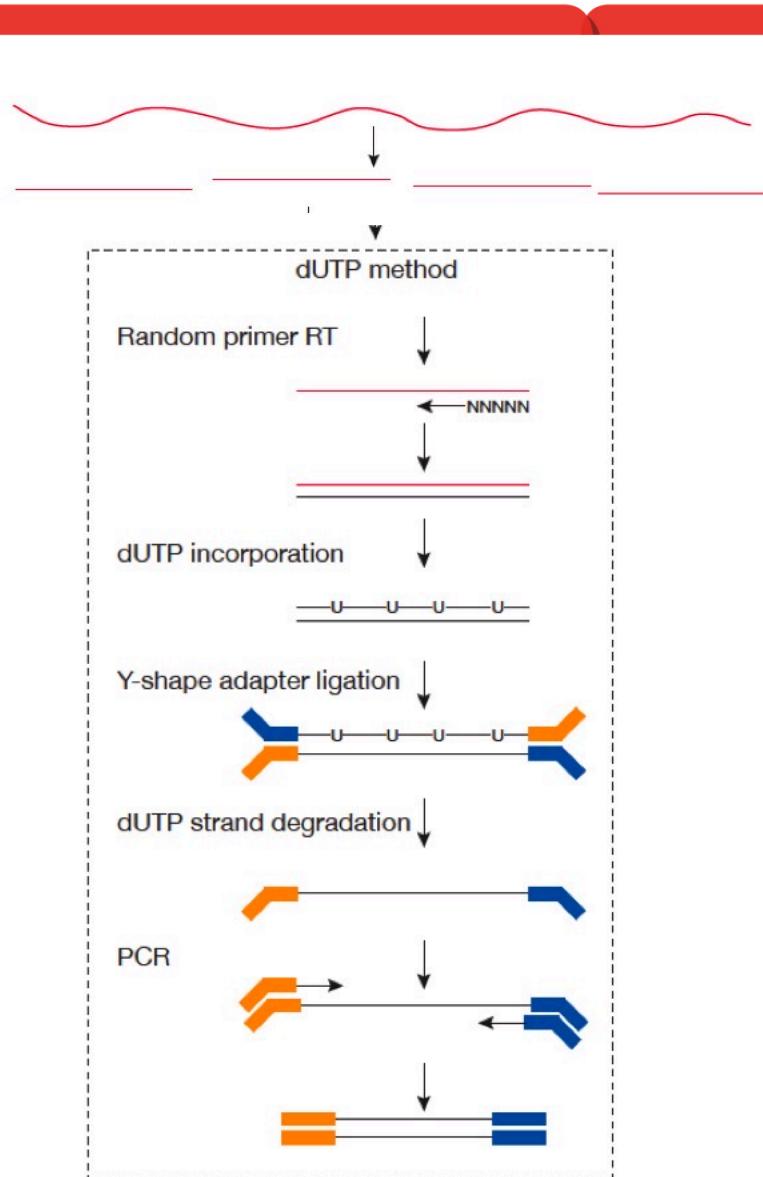
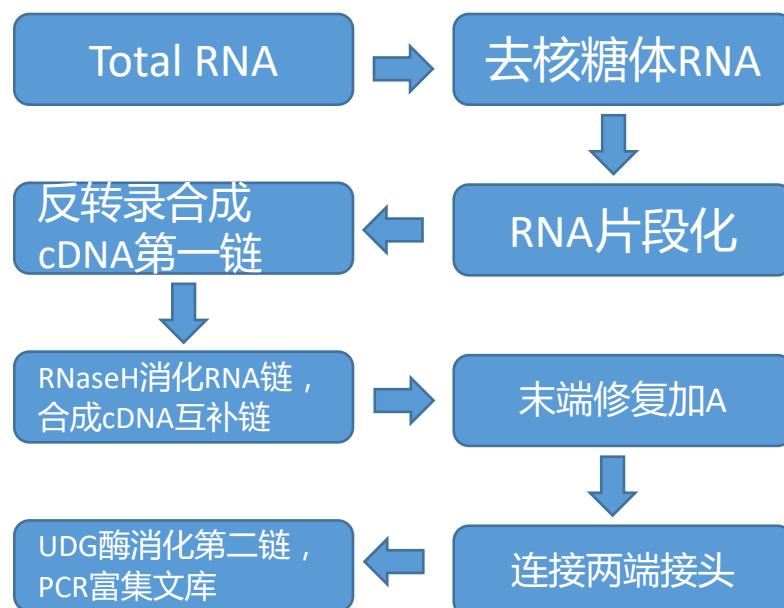


LncRNA文库构建

>200bp

链特建库可区分Reads的正负方向

- 定量更加准确
- 对可变剪切的检测更准确
- 无参转录组更加真实



3

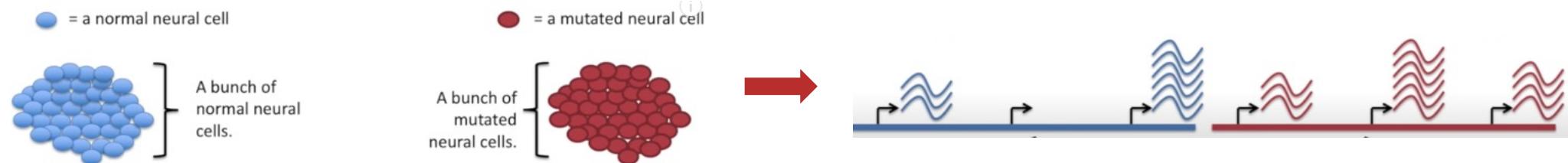
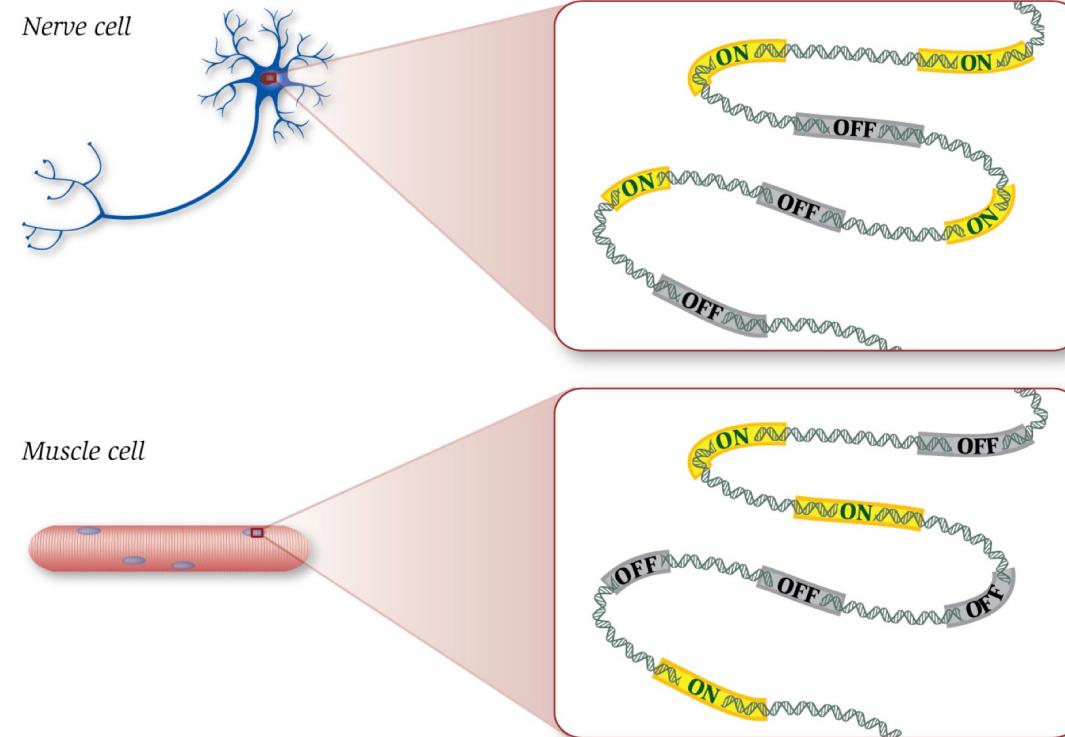


RNA-seq分析内容

RNA-seq解决的问题

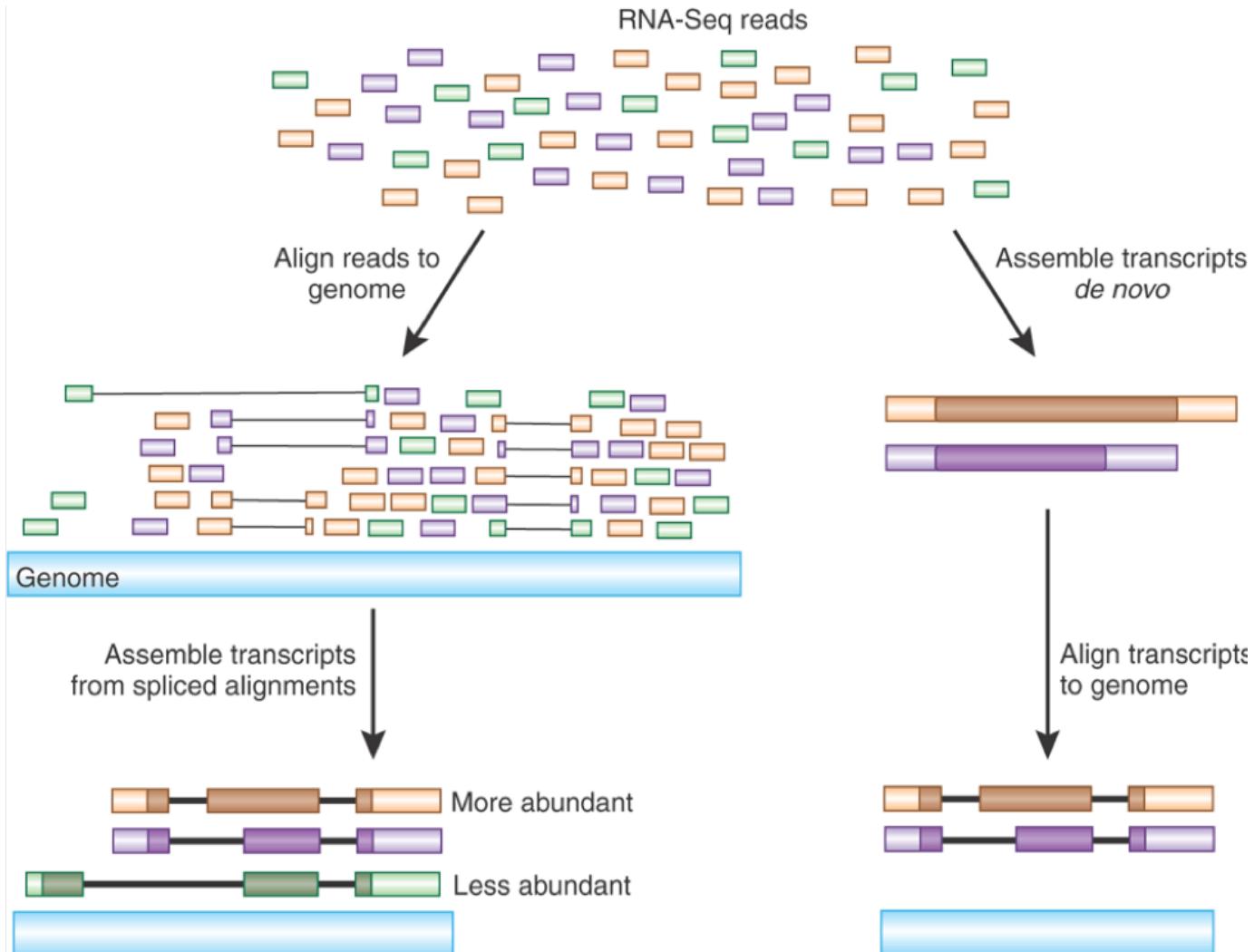
基因的表达具
有时空特异性

- 在某些样本中，哪些基因是活跃的（转录表达）？
- 不同样本中，差异表达的基因有哪些？

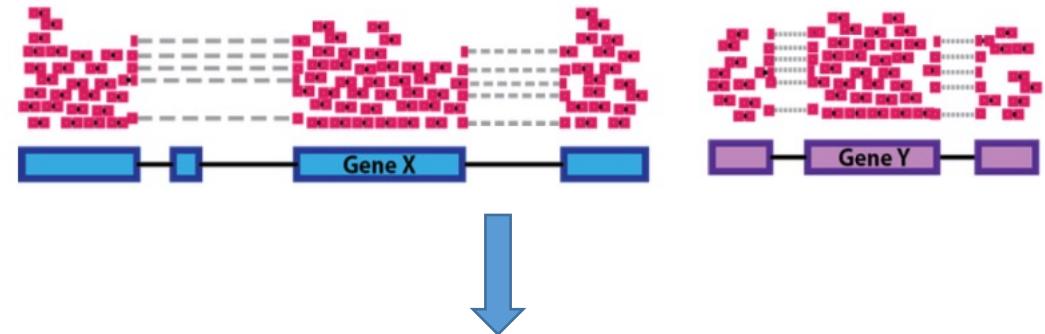


RNA-seq分析流程

有参转录组



基因表达矩阵



	ctrl_1	ctrl_2	exp_1	exp_1
geneA	10	11	56	45
geneB	0	0	128	54
geneC	42	41	59	41
geneD	103	122	1	23
geneE	10	23	14	56
geneF	0	1	2	0
...
...

Count matrix

获得基因表达矩阵后，

可以从哪些维度分析？

不同样本中活跃基因的数

目；样本之间的相似性...

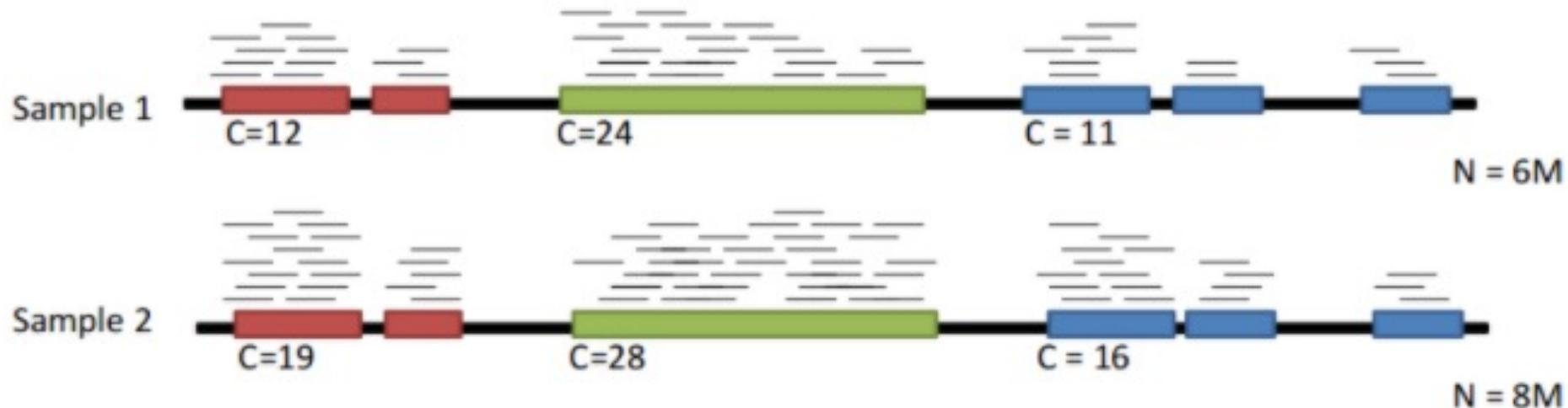
样本维度

差异表达基因有多少；基

因的具体功能是什么？

基因维度

数据的标准化

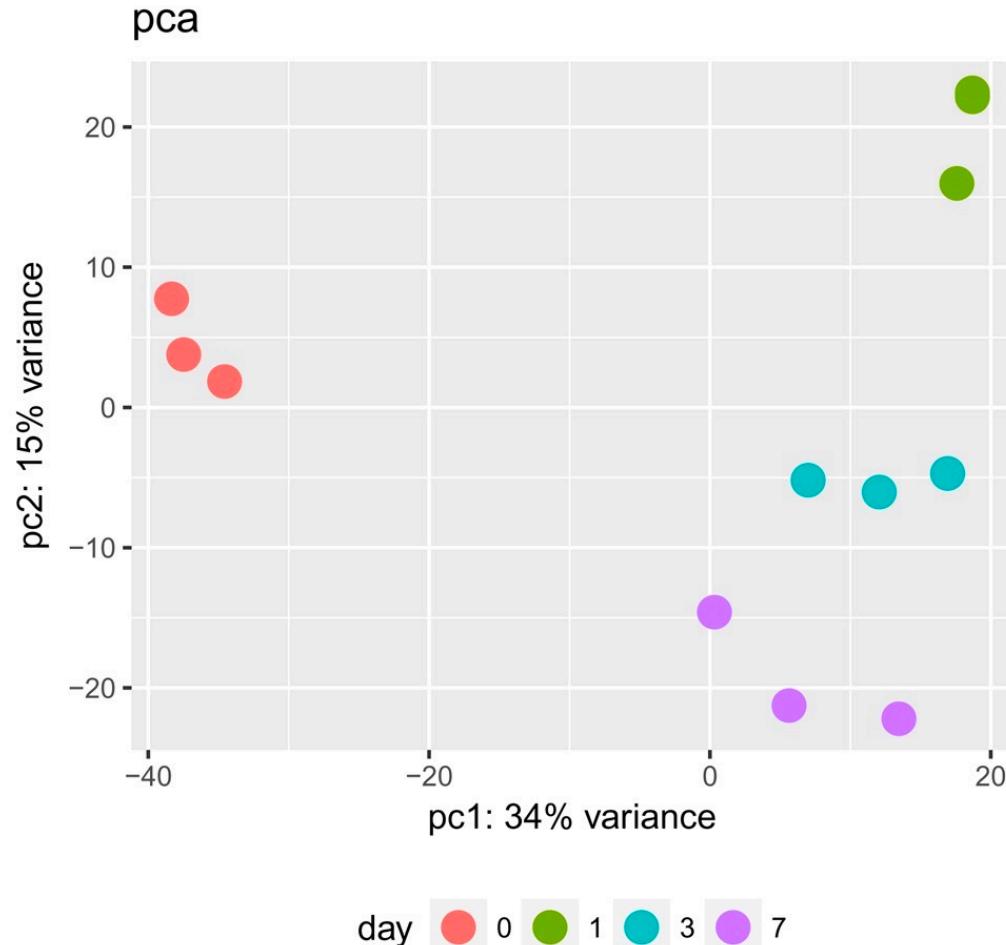


测序深度的影响
基因长度的影响

$$RPKM = \frac{10^6 * n_r}{L * N}$$
$$FPKM = \frac{10^6 * n_f}{L * N}$$
$$TPM = \frac{\frac{n_i}{L_i} * 10^6}{\sum_{i=1}^N \frac{n_i}{L_i}}$$

注：n : reads比对上的count数；N : 总的测序的reads数；L : 基因片段的长度

样本维度-PCA



不同类型的样本，是否有显著的分离？

分离：样本间的基因表达模式差异较大

不分离：样本间的基因表达较为相似

横坐标：主成分1及其方差解释率

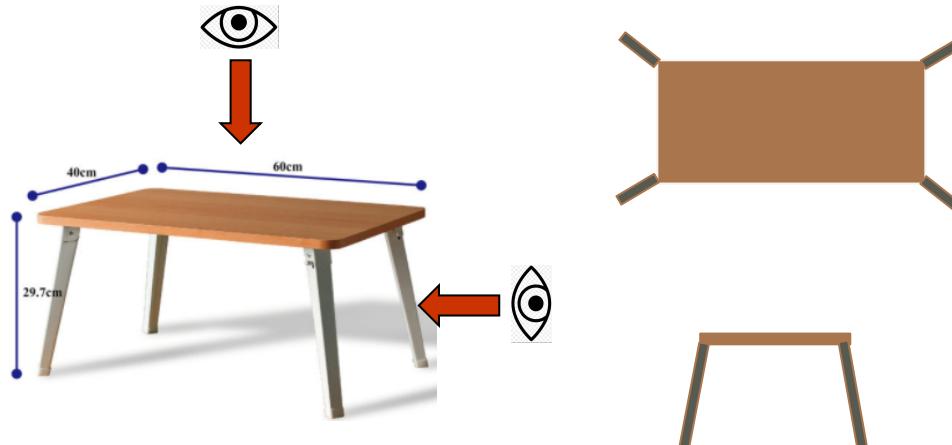
纵坐标：主成分2及其方差解释率

图中的点代表样本

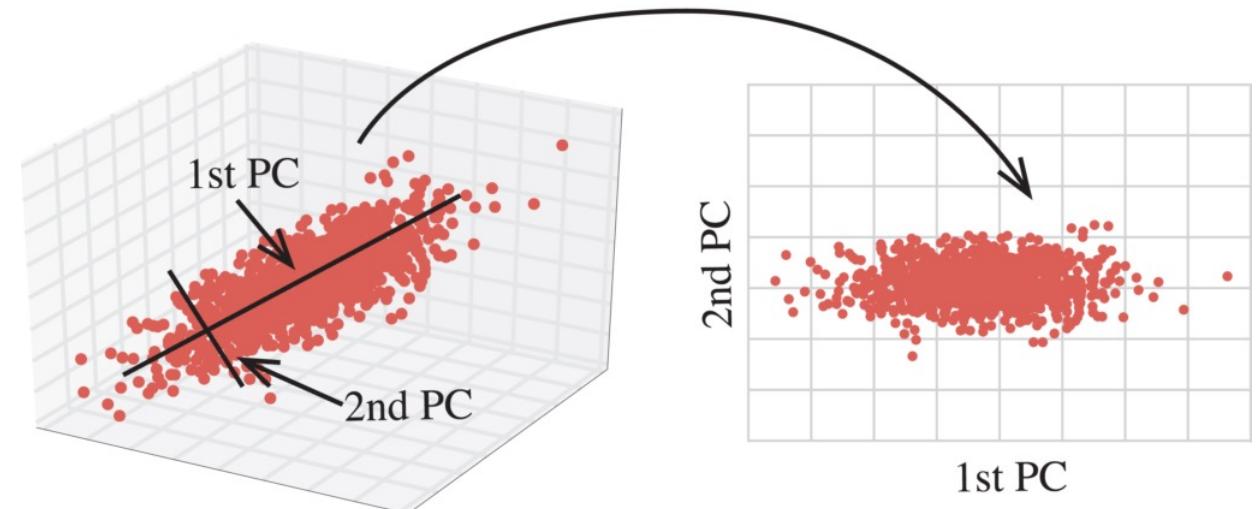
颜色表示不同分组

PCA是什么

Principal component analysis : 就是从最佳的角度看数据

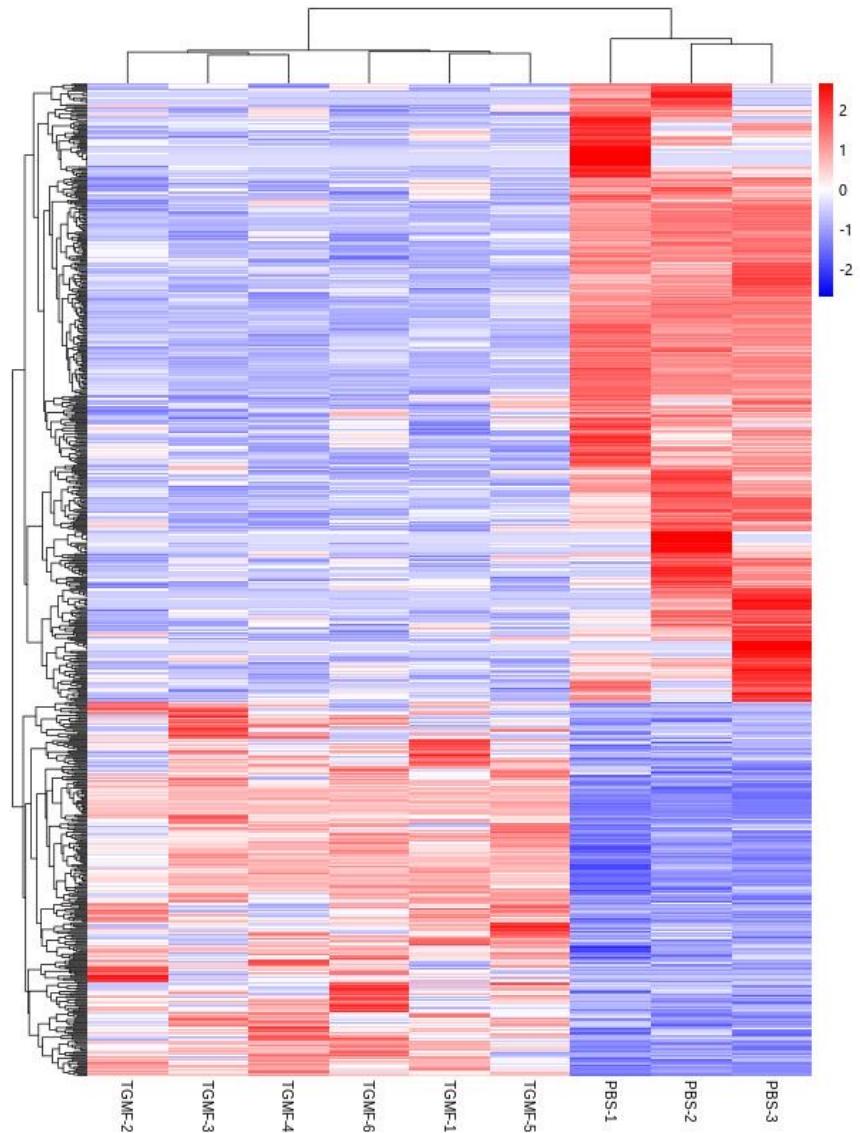


High dimension
Complexity
Simplicity
Low dimension



对数据的“降维打击”，以达到视觉可及的程度
也是对数据的简化，方便理解数据的组成

样本维度-表达热图



不同基因和样本的表达量差异

行：基因名；列：样本名

红色：高表达；绿色：低表达

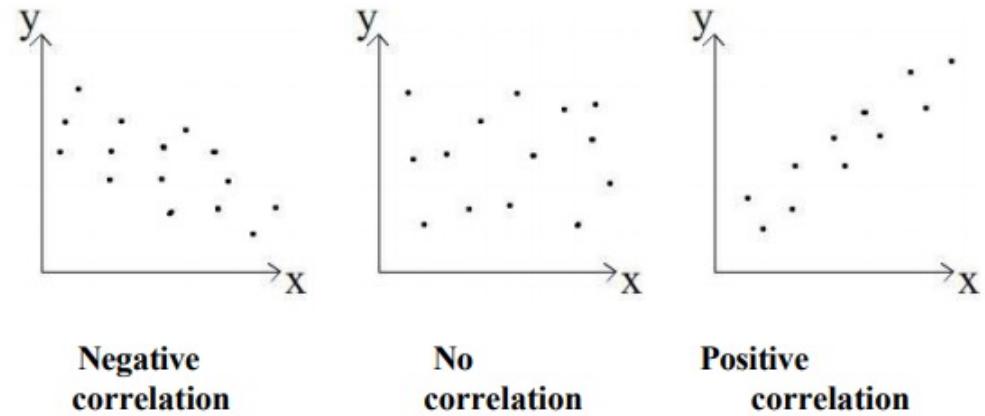
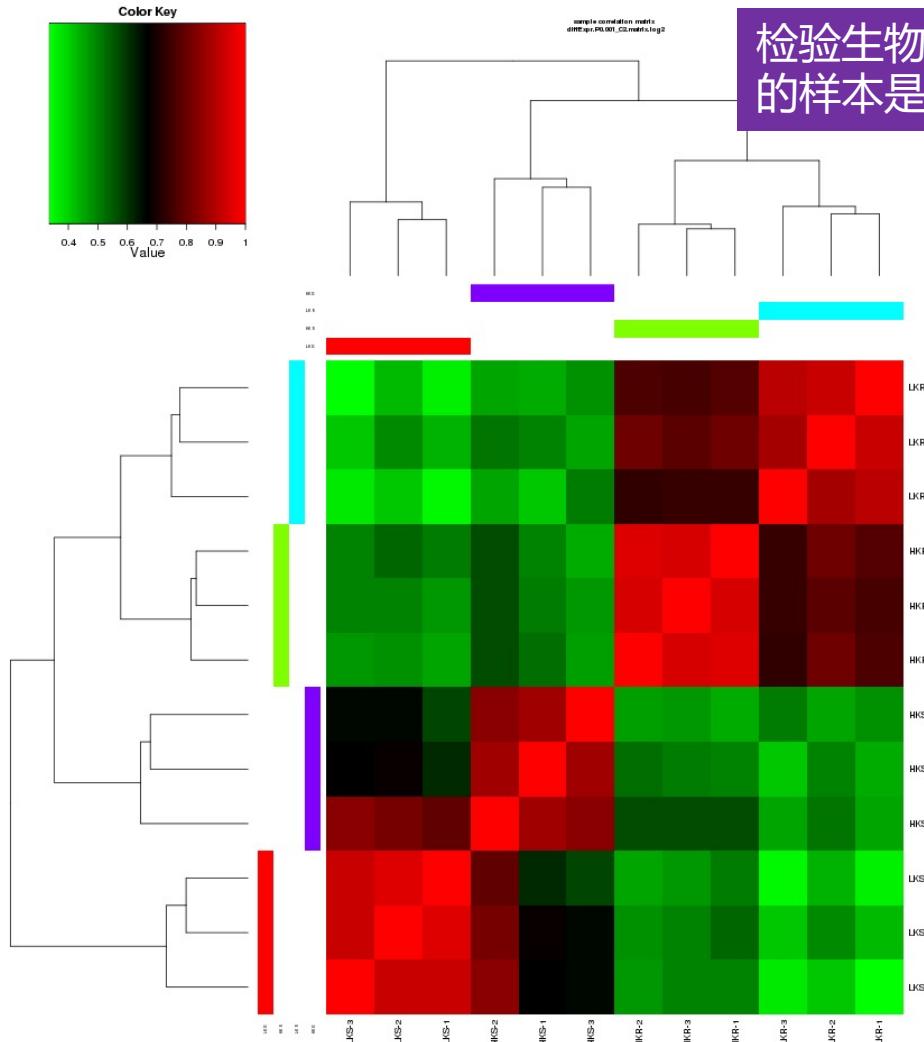
行聚类：基因在不同样本中的表达相似

列聚类：不同样本的基因表达比较相似

标准化 Z-score

$$z = \frac{x - \mu}{\sigma}$$

样本维度-样本相关性



1、只表相关，而不表因果 2、相关亦有强弱之分

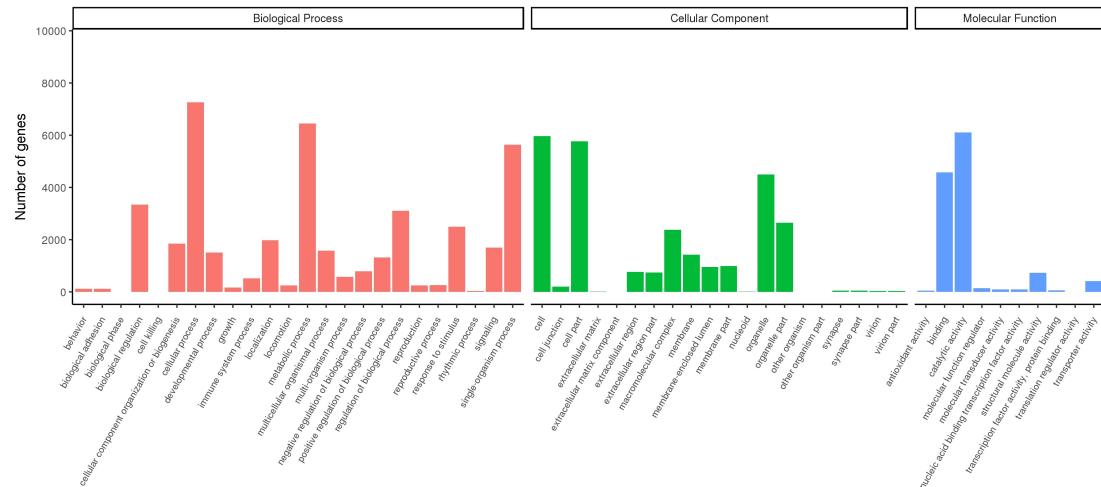


基因维度-基因功能注释

GO

Gene Ontology : 相似的基因，保守的功能

- 分子功能 (Molecular Function) 【催化或转运活性】
- 细胞组分 (Cellular Component) 【线粒体、核糖体】
- 生物学过程 (Biological Process) 【DNA修复或信号转导】

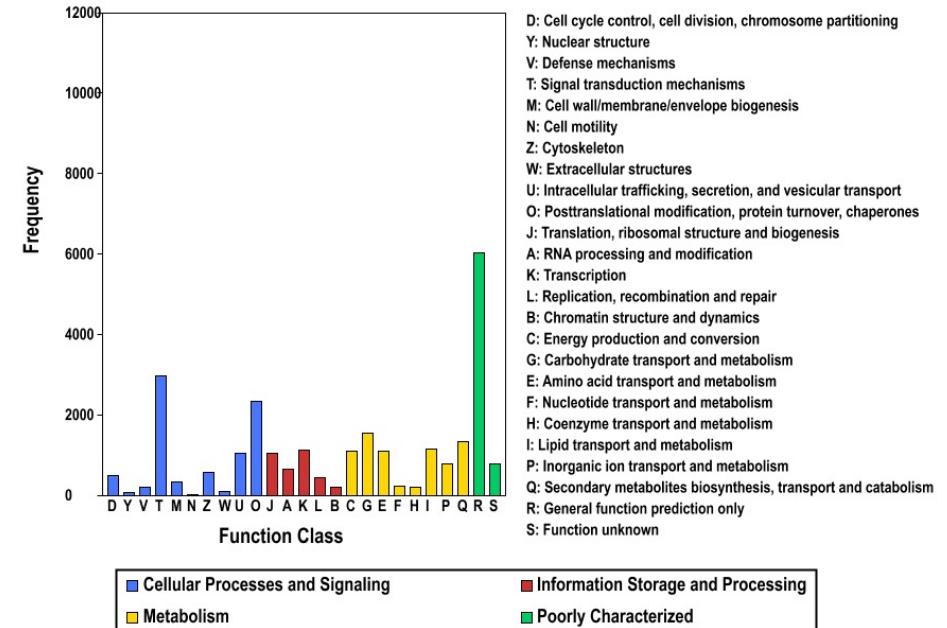


KOG

Eukaryotic Orthologous Groups of proteins

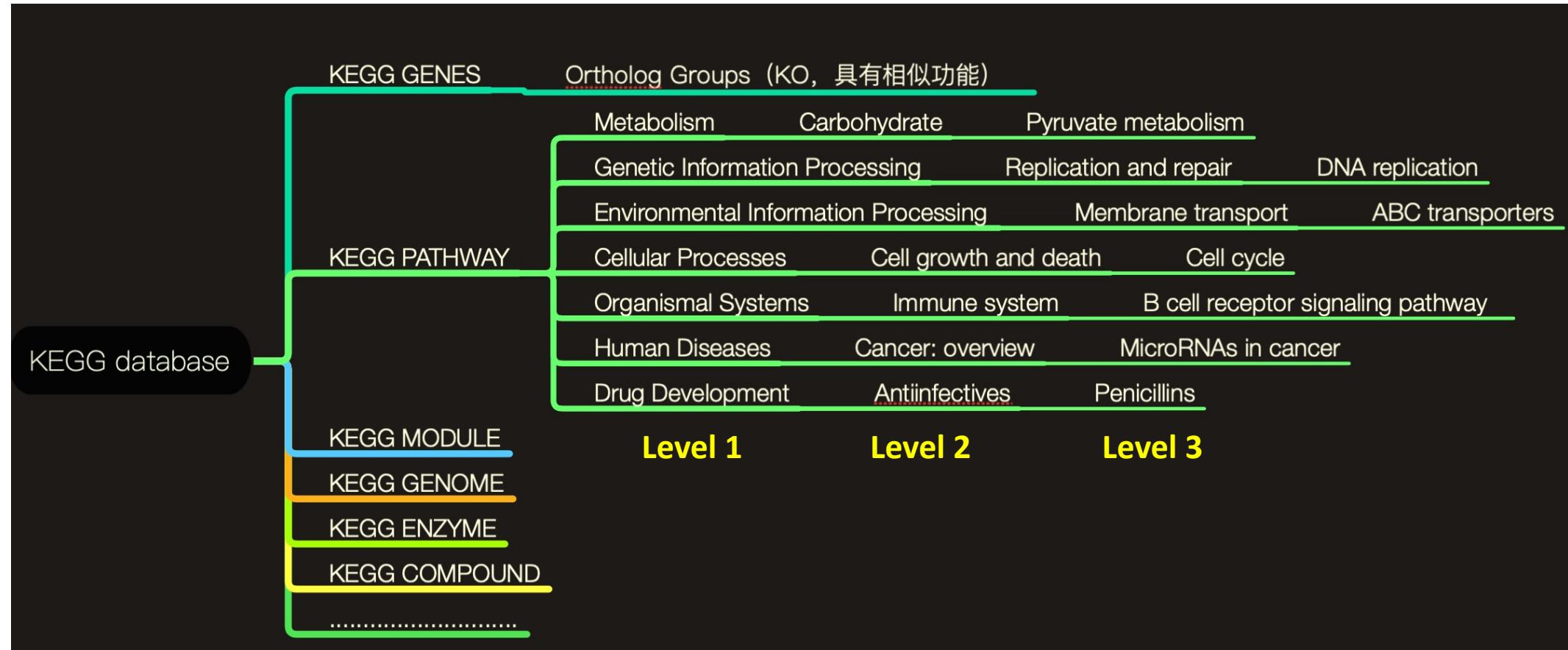
该数据库来源七个真核生物的全基因组，酿酒酵母、裂殖酵母、脑原虫、拟南芥、线虫、果蝇和人的完整基因组。共有4852个KOGs（来自于共同祖先的同源基因），包含60579个蛋白，分成25个功能单元。

KOG Function Classification of Consensus Sequence



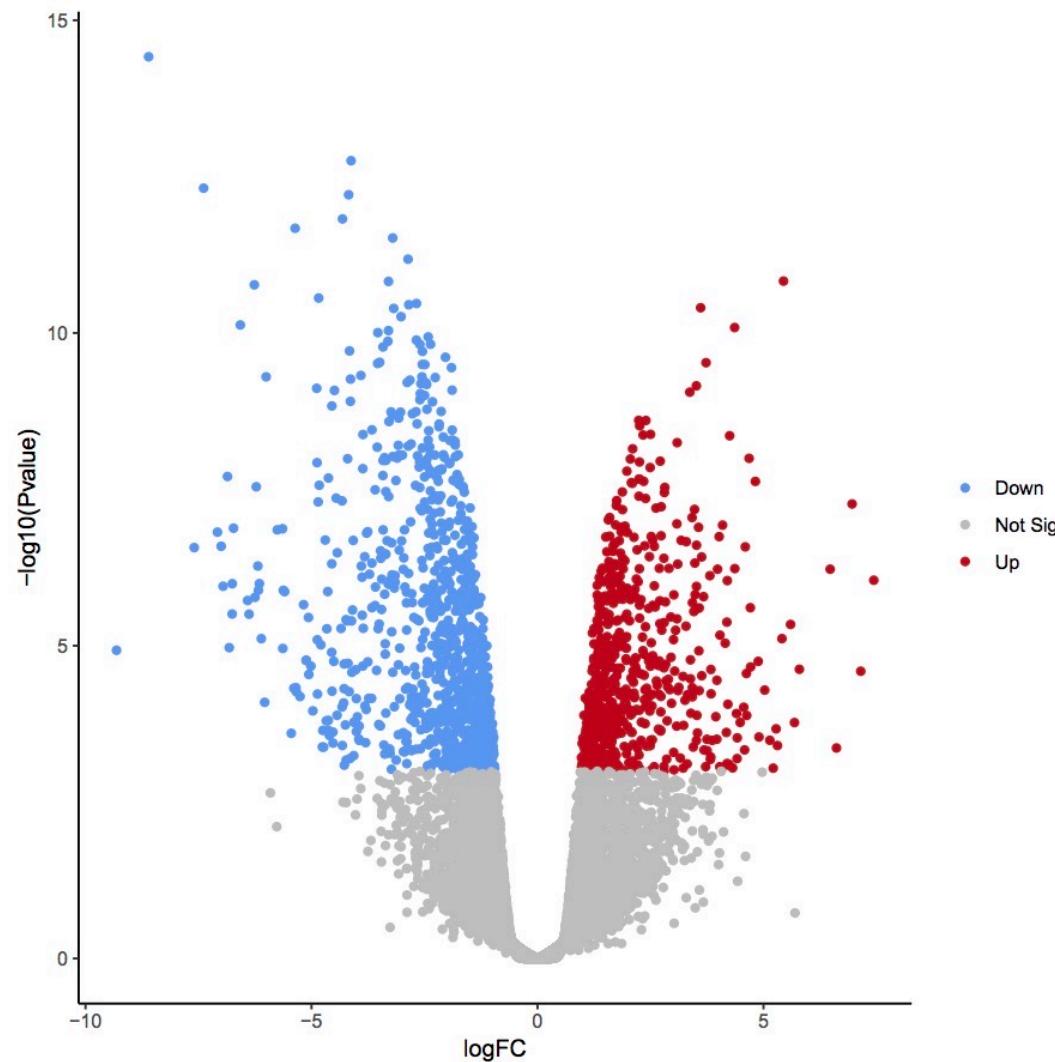
基因维度-基因功能注释

Kyoto Encyclopedia of Genes and Genomes



提供全面的功能注释信息，包括基因、酶、化学产物、所属物种等，并系统分类了代谢通路

基因维度-差异表达基因



对照-处理组的差异基因有多少、有哪些

圆点：某个基因；颜色：蓝色下调；红色：上调；灰色：差异不显著

横坐标：**log2 (fold change)**

纵坐标：纵坐标： $-\log_{10}(\text{FDR})$

着色： $\text{FDR} < 0.05 \text{ } \& \text{ } |\log_2(\text{FC})| > 1$

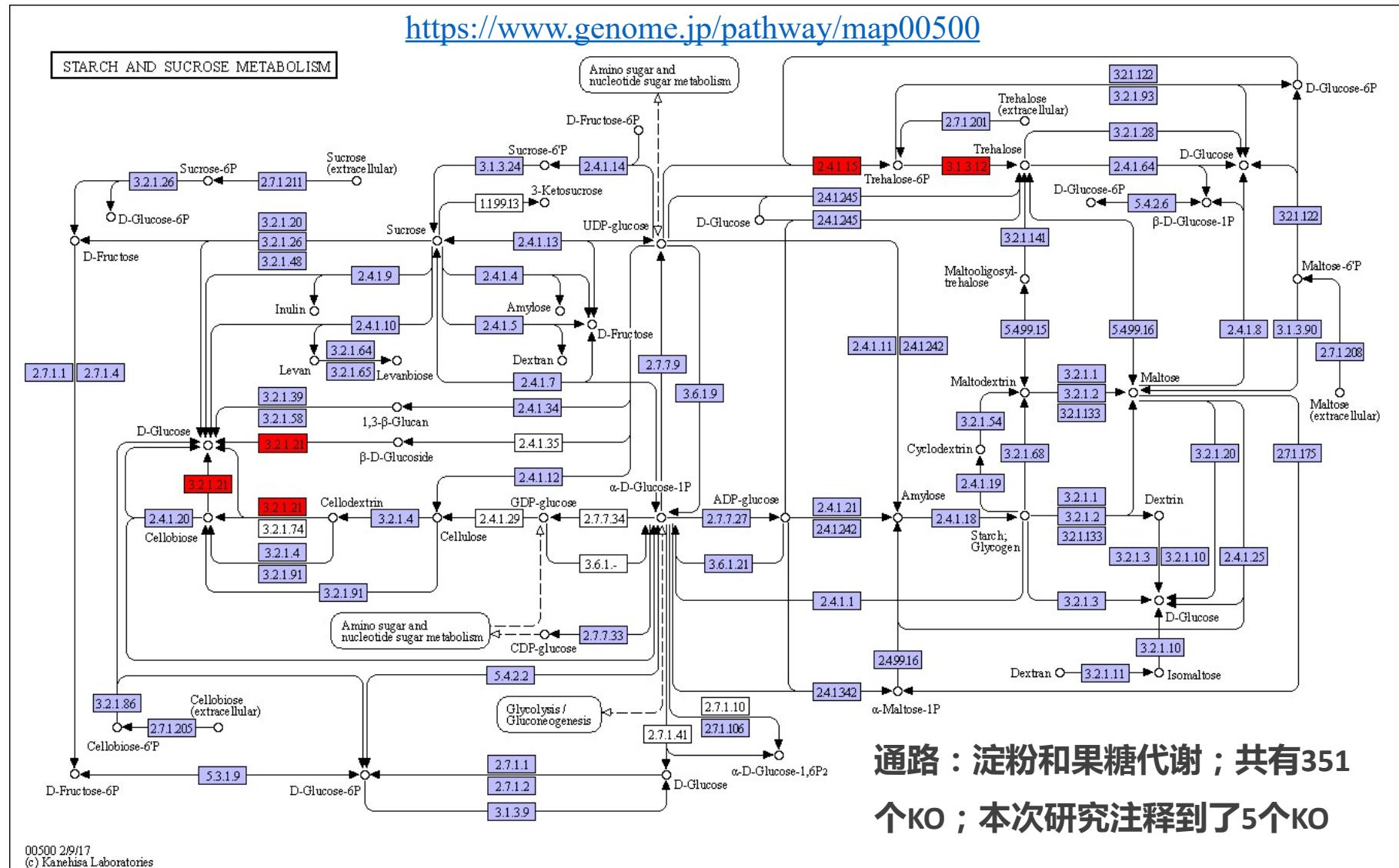
Deseq2 and edgeR

? Fold change为什么用常用log转化？

(0, 1)	(1, +∞)
(-∞, 0)	(0, +∞)

基因维度-KEGG代谢通路

<https://www.genome.jp/pathway/map00500>



圆圈：化学物质

矩形：基因产物 (酶或蛋白)

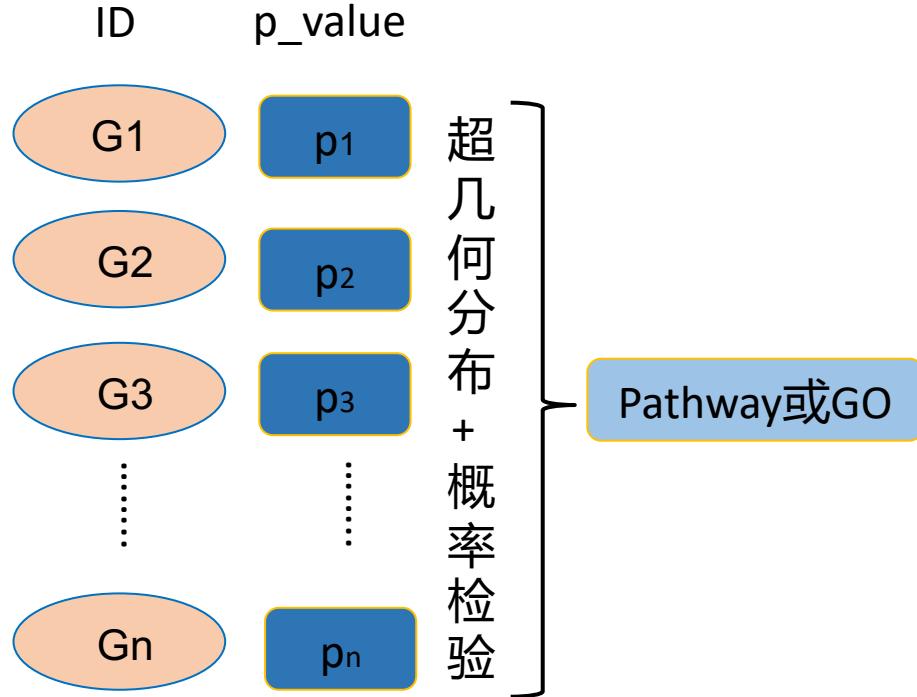
蓝：可超链接；

红：注释到该KO

白：无法链接

通路：淀粉和果糖代谢；共有351个KO；本次研究注释到了5个KO

基因维度-富集分析

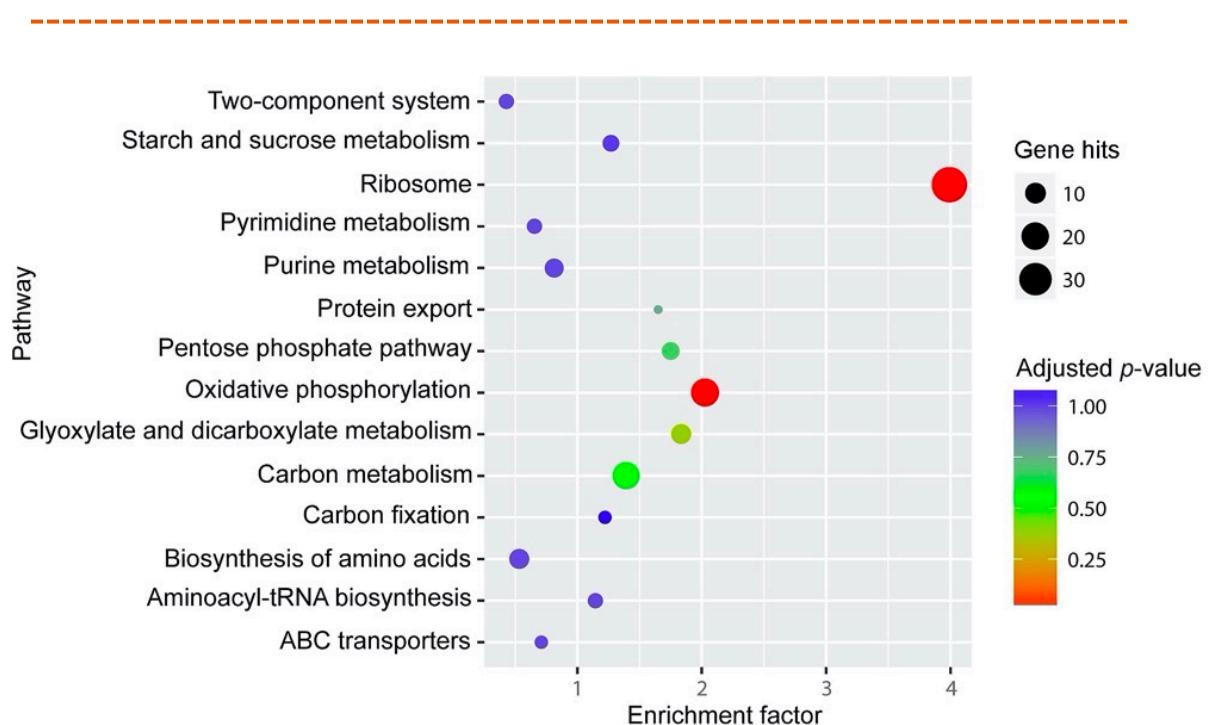


信号：在某个Pathway或GO term中的差异基因比例
背景：所有差异基因数占整个注释基因数目的比例

通俗举例，一个省区是否显著富裕？

信号：富裕人数占该省总人口的比例

背景：国家富裕人数占总人口的比例

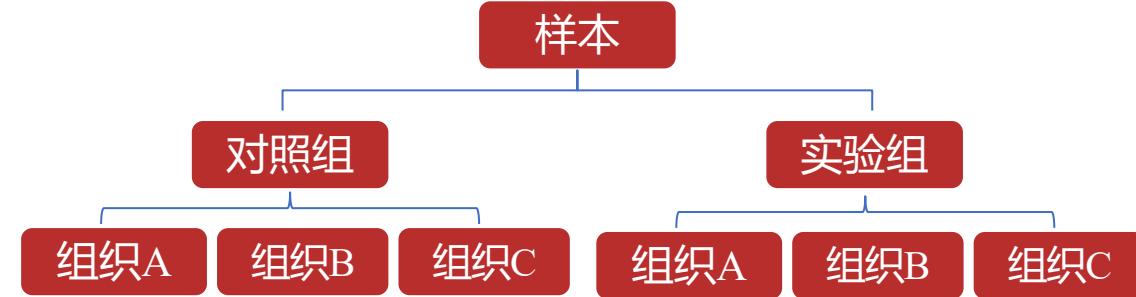


4



RNA-seq研究思路

差异转录组

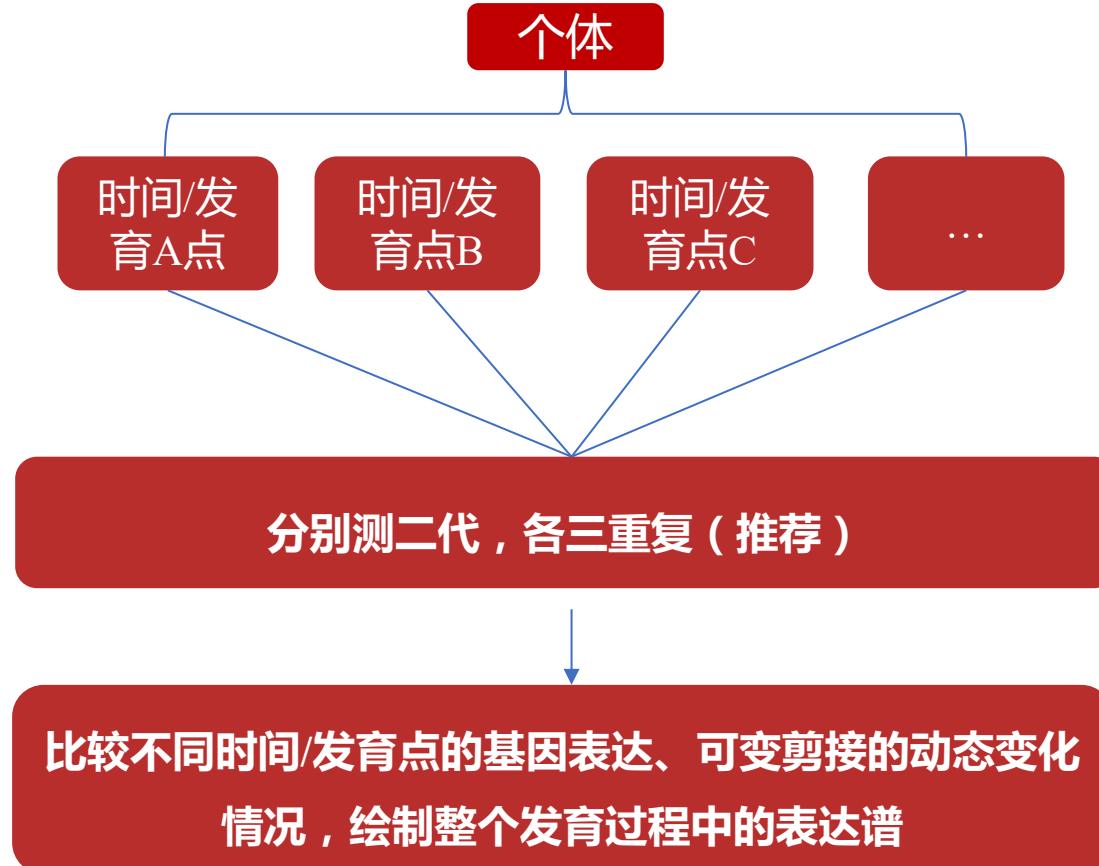


分别测二代，各三个重复

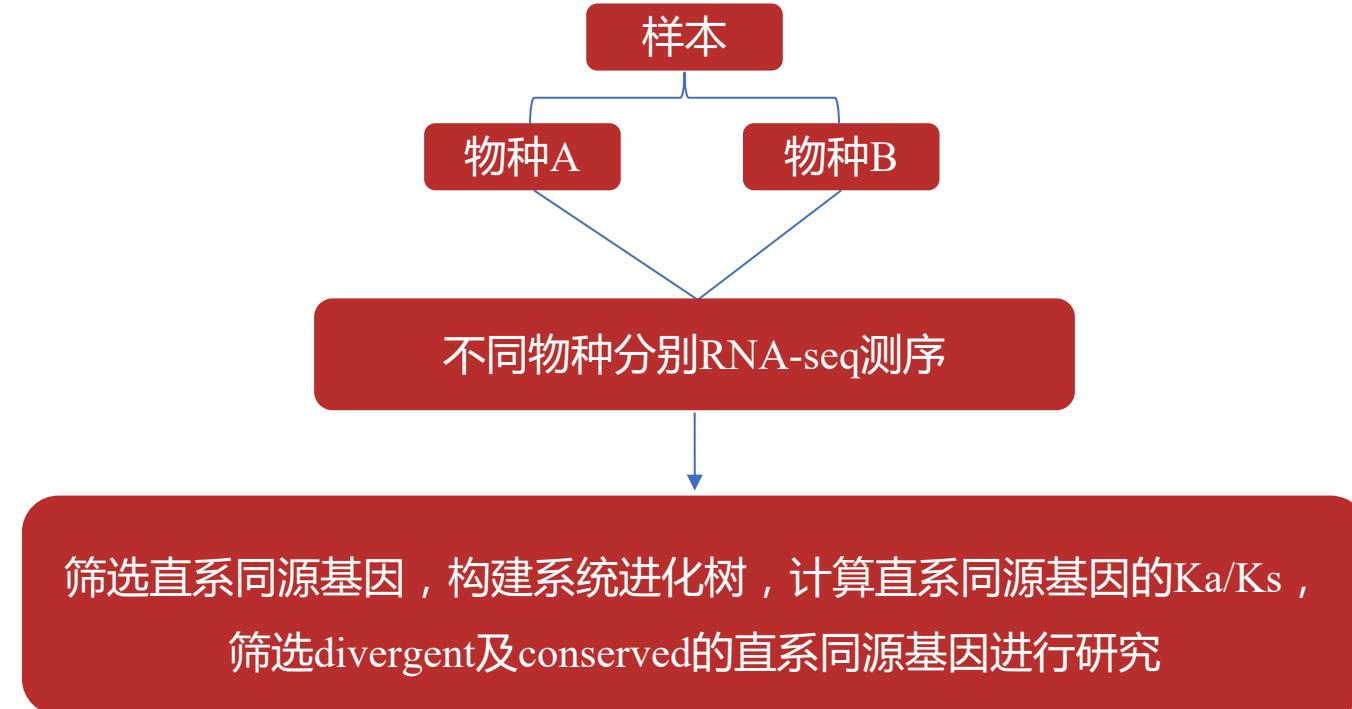
差异分析，既可分析组织间差异，也可分析胁迫和正常条件的差异，找出差异关键基因或转录本，功能性研究，可结合多组学进行研究，并联合实际处理条件，研究用药的作用靶点和动植物的胁迫响应机制

注：样本选择为同一物种的正常组和实验组。

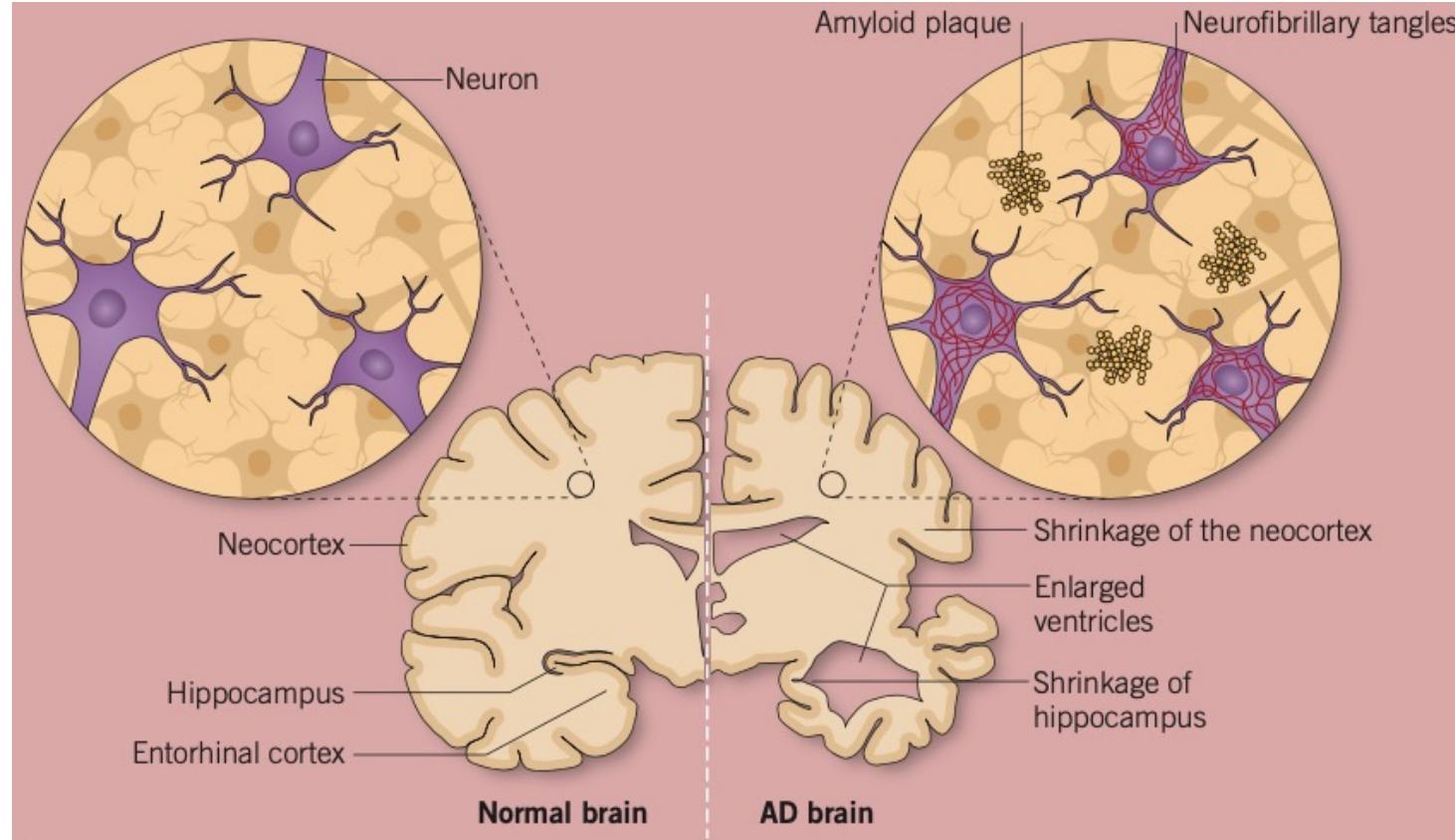
动态转录组



比较转录组



案例解读



阿尔兹海默症（Alzheimer's disease，AD）病例假说： β 淀粉样蛋白（A β ）和Tau（微管蛋白）；A β 形成斑块和Tau过度磷酸化形成神经纤维缠结，进而影响神经元健康，导致神经退行性症状。

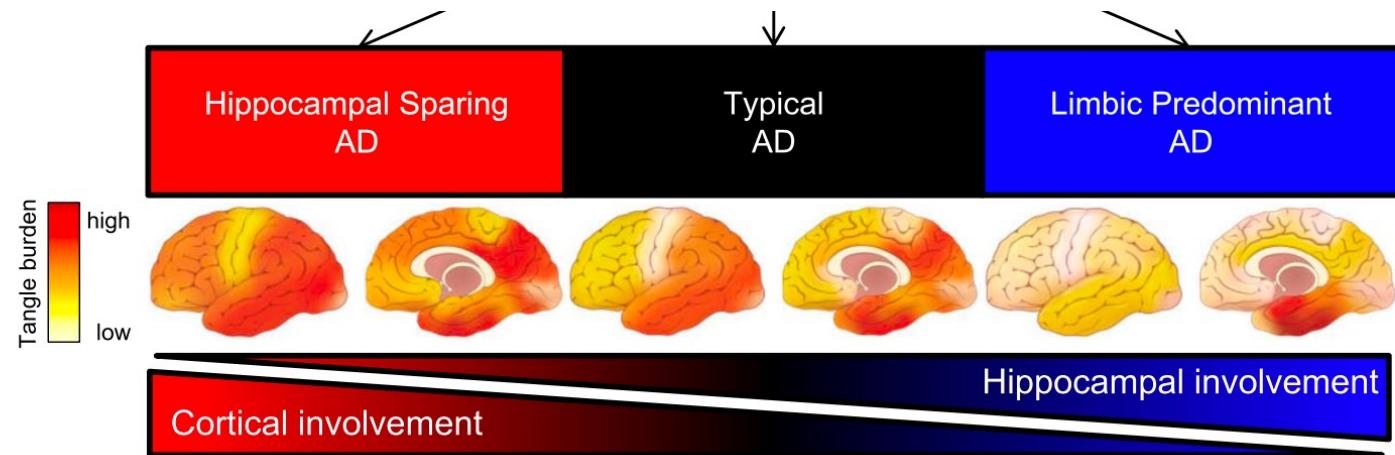
A β 沉积：皮层到边缘系统

Tau：皮层下核团/内嗅皮层到边缘系统



鉴定不同AD亚型中的基因表达谱

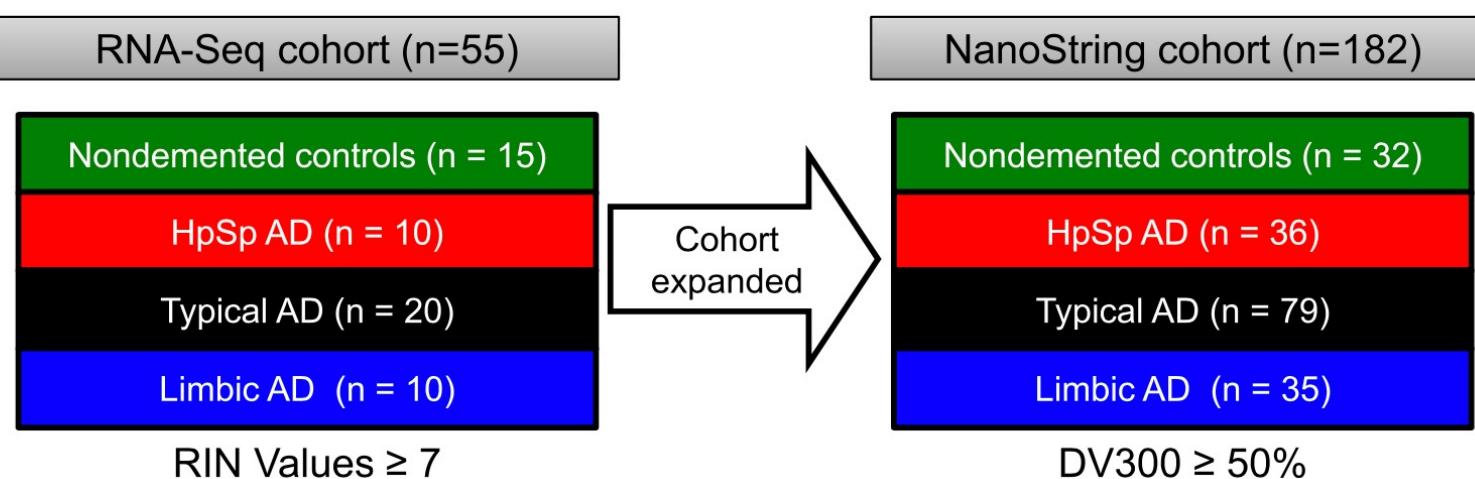
实验设计



三种AD的亚型：

海马保护型AD（皮层严重，海马较轻）；典型的AD（海马和皮层均有病变）；边缘主导型AD（海马严重，皮层较轻）

b

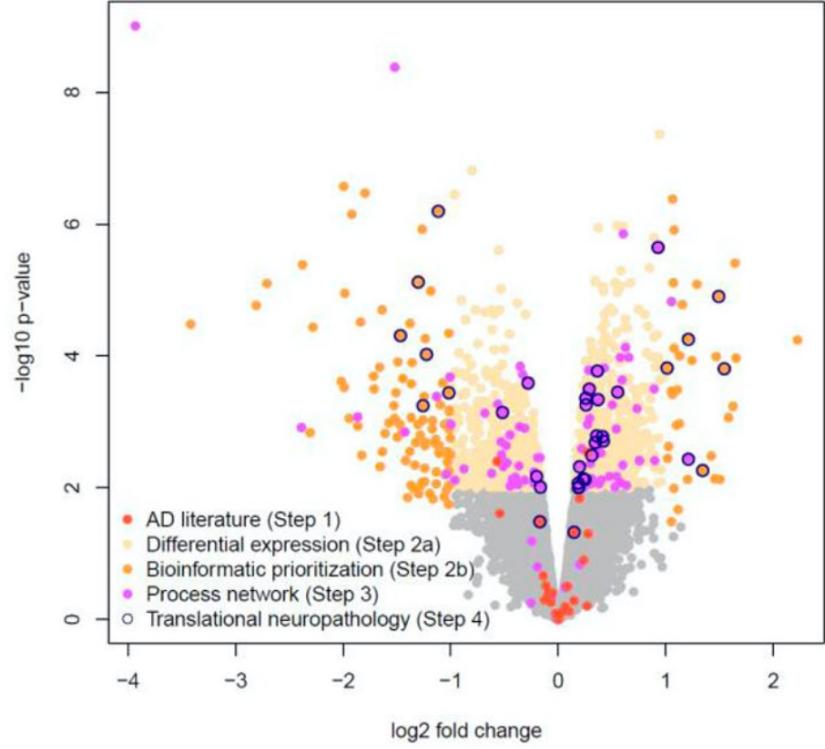


海马区域易感基因

主要结果

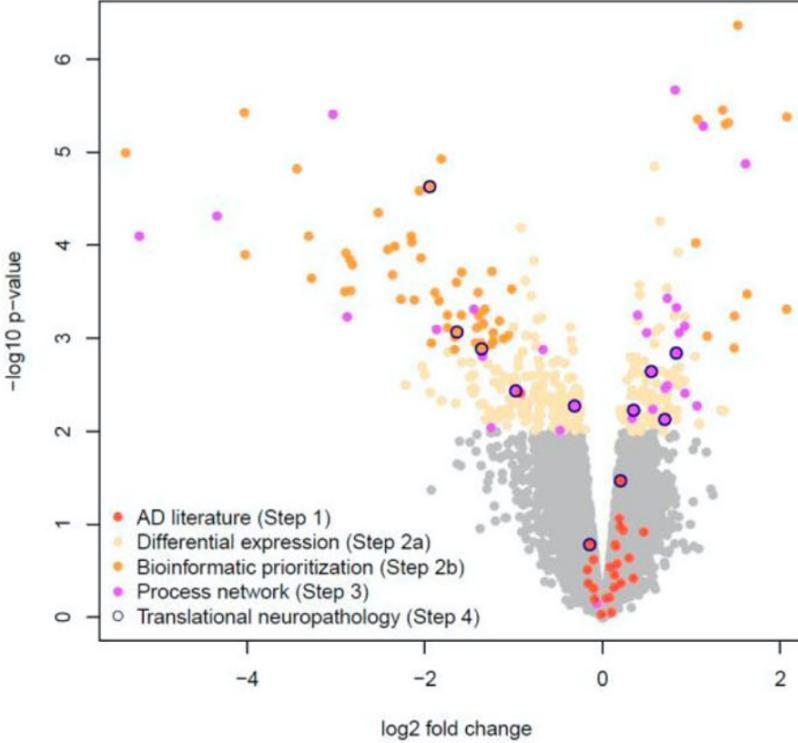
k

Representative Phenotype



l

Extreme Phenotype



AD中海马易感性的基因；发现了 SERPINA5, RYBP, SLC38A2, FEM1B和 PYDC1可以很好的区分AD和对照组

典型AD（与对照组相比）和边缘主导型AD（与海马保护型AD相比）中，大多数基因都是下调的。

在典型AD中的差异基因多与细胞生长和增殖相关，而在边缘主导型AD中大多数基因与炎症和免疫反应相关

今天讲了什么

Illumina测序原理

mRNA文库构建

RNA-seq解决的问题

RNA-seq的分析流程

RNA-seq的分析内容

RNA-seq的研究思路

一篇AD的案例



BerryGenomics
贝瑞基因

Thank You!



官方网站



官方微信

TCGATCGA GATCGATCGATCGATCGATCGATCG

TAGATCGATCGATCGATCGATCGATCGATCGATCG

CGATGATCGATCGATCGATCGATCGATCGATCGATCG

www.berrygenomics.com