

# 生物芯片和高通量测序 技术及其应用

肖华胜博士

生物芯片上海国家工程研究中心  
上海伯豪生物技术有限公司

Email: [huasheng\\_xiao@shbio.com](mailto:huasheng_xiao@shbio.com)

# 主要内容

第一部分 概述

第二部分 基因芯片

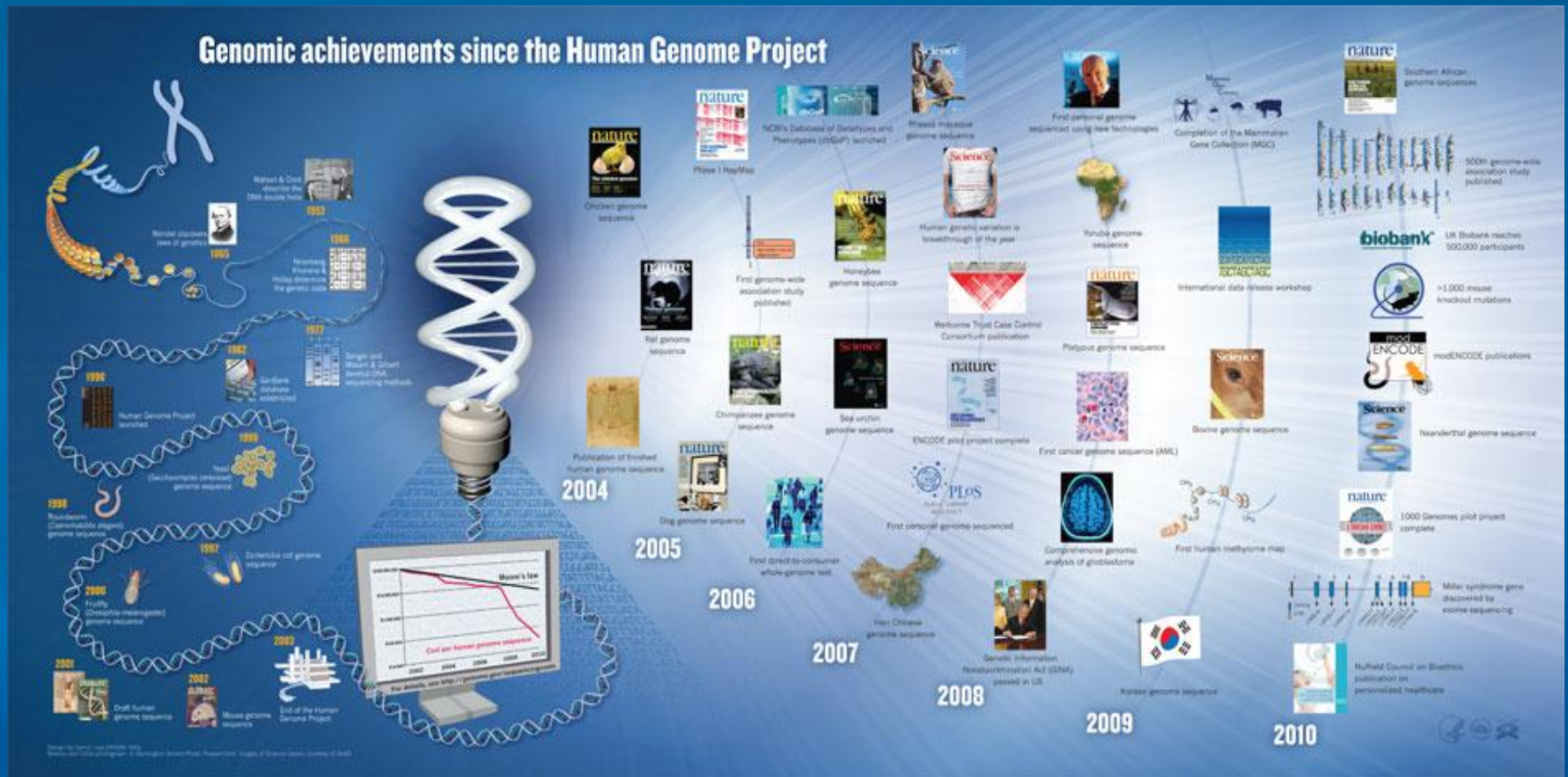
第三部分 蛋白芯片

第四部分 组织芯片

第五部分 高通量测序技术及其应用

# 第一部分 概述

“人类基因组计划”和“曼哈顿”原子弹计划，“阿波罗”登月计划，一起被誉为自然科学史上的“三大计划”。



图片来源: [Charting a course for genomic medicine from base pairs to bedside](#) Nature Volume: 470, Pages: 204 – 213 Date published: (10 February 2011)

人类基因组研究战略和实验技术源源不断地产生了日益庞大及复杂的基因组数据，这些数据已被载入公共数据库，并改变了对几乎所有生命过程的研究。



# 人类基因组计划(HGP)

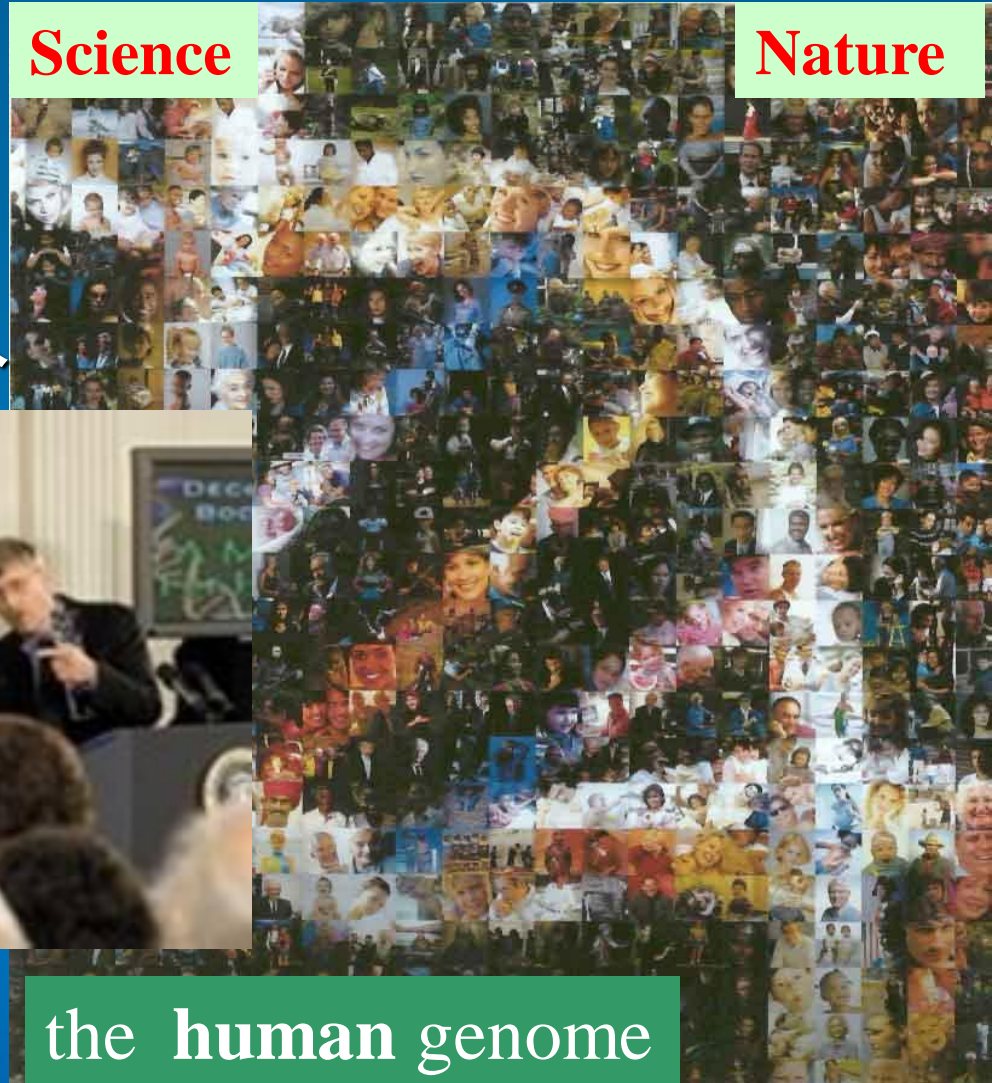
在1990年左右项目启动。

2000年6月26日，国际基因组组织  
宣布基因组草图完成。

2003年4月4日Collin Francis宣布人  
类基因组项目完成。



Xinhua



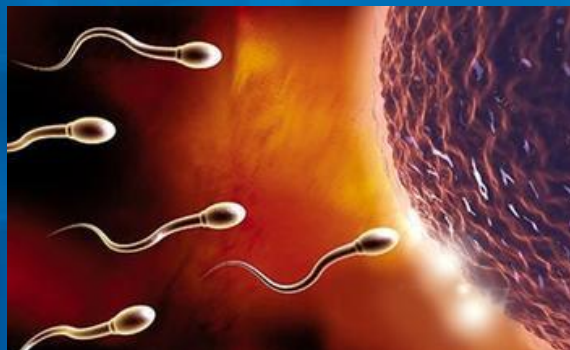
Science

Nature

the human genome



# 一些重要生物的基因组测序和重测序



# 生物芯片和测序技术是生命科学和医学研究的重要技术

基因组  
(DNA)

转录组  
(RNA)

蛋白质组  
(Protein)

功能基因组学



➤ Northern Blot

➤ 一代测序技术

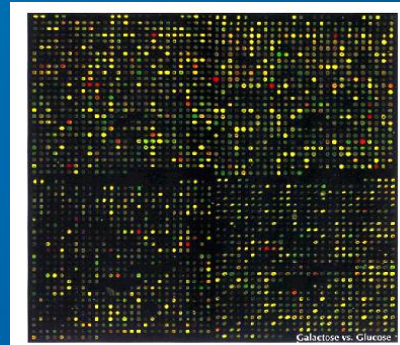
➤ 新一代测序技术

➤ 差异显示、RDA (代表性差异分析)

➤ 基因芯片

➤ SAGE

➤ 蛋白质芯片



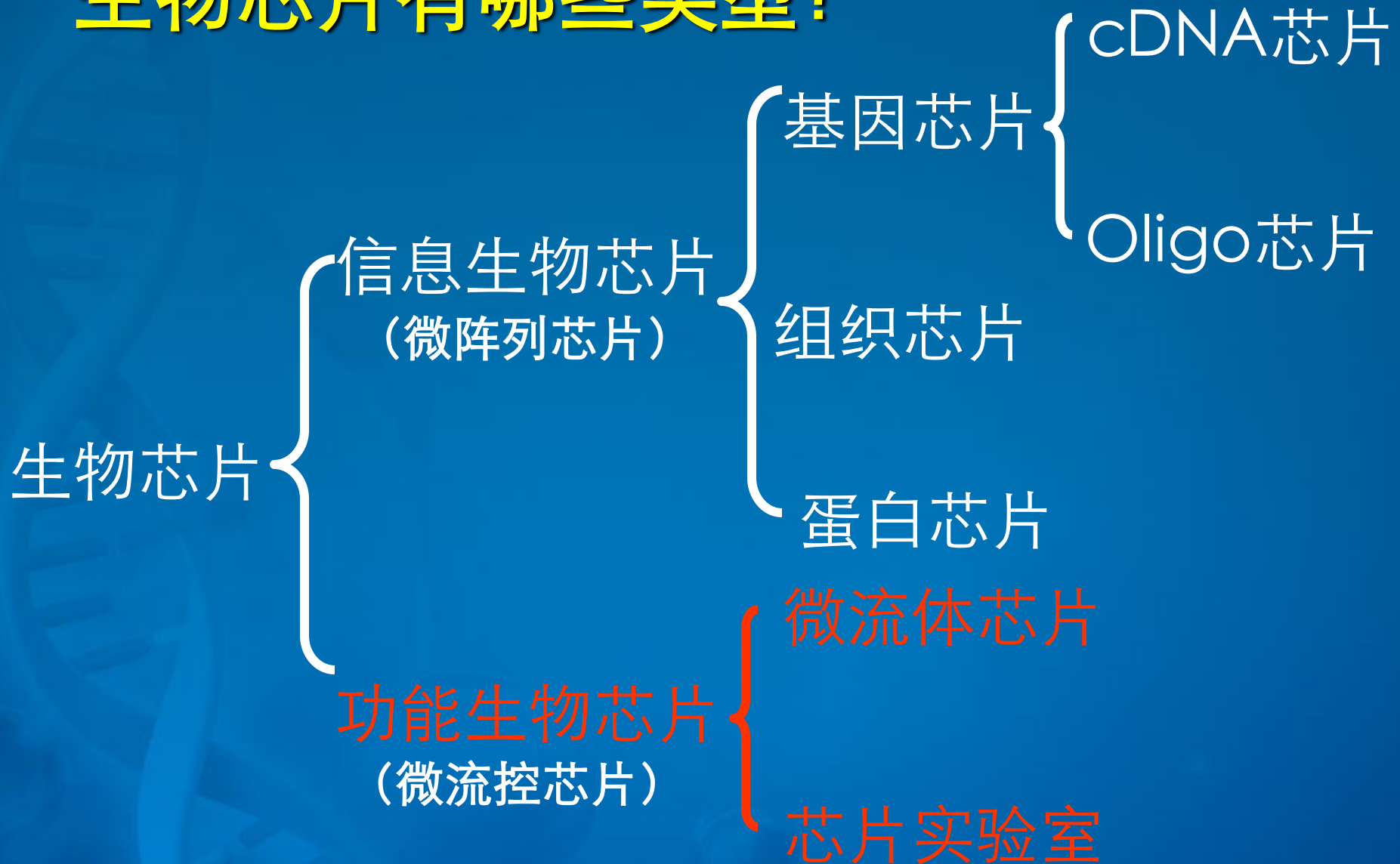


# 什么是生物芯片?

- 生物芯片的概念是Fodor等人于1991年提出(Fodor et al., 1991, Science)。
- 生物芯片的概念：是借用电子芯片的概念，是指能够快速并行处理多个样品并对其所包含的各种生物信息进行解剖的微型器件，它的加工运用了微电子工业和微机电系统加工中所采用的一些方法，只是由于其所处理和分析的对象是生物样品，所以叫生物芯片(Biochip)。
- 高通量，平行化，微量化的分析手段。



# 生物芯片有哪些类型？



# 生物芯片技术的应用领域



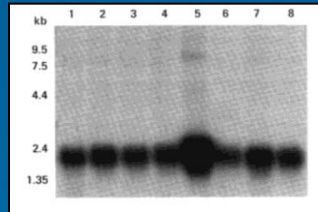
## 第二部分 基因芯片

基因芯片(Gene chip, DNA microarray)技术是指通过微阵列(Microarray)技术将高密度DNA片段通过高速机器人或原位合成方式以一定的顺序或排列方式使其附着在如膜、玻璃片等固相表面, 以同位素或荧光标记的DNA探针, 借助碱基互补杂交原理, 进行大量的基因表达及监测等方面研究的最新革命性技术。

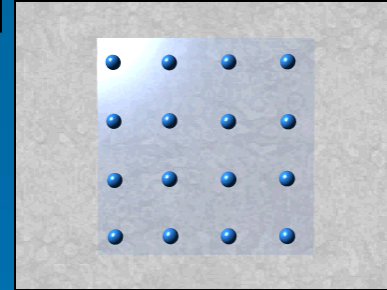


# 基因芯片发展过程

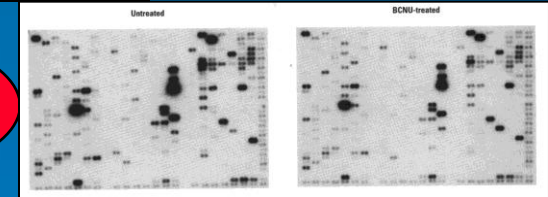
Southern & Northern Blot



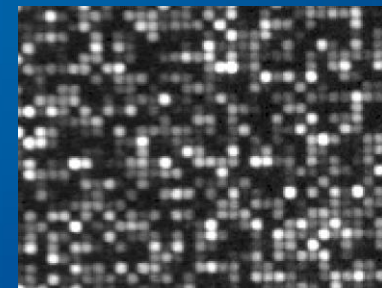
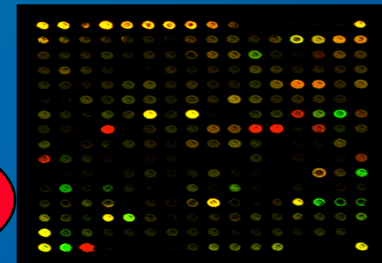
Dot Blot



Macroarray

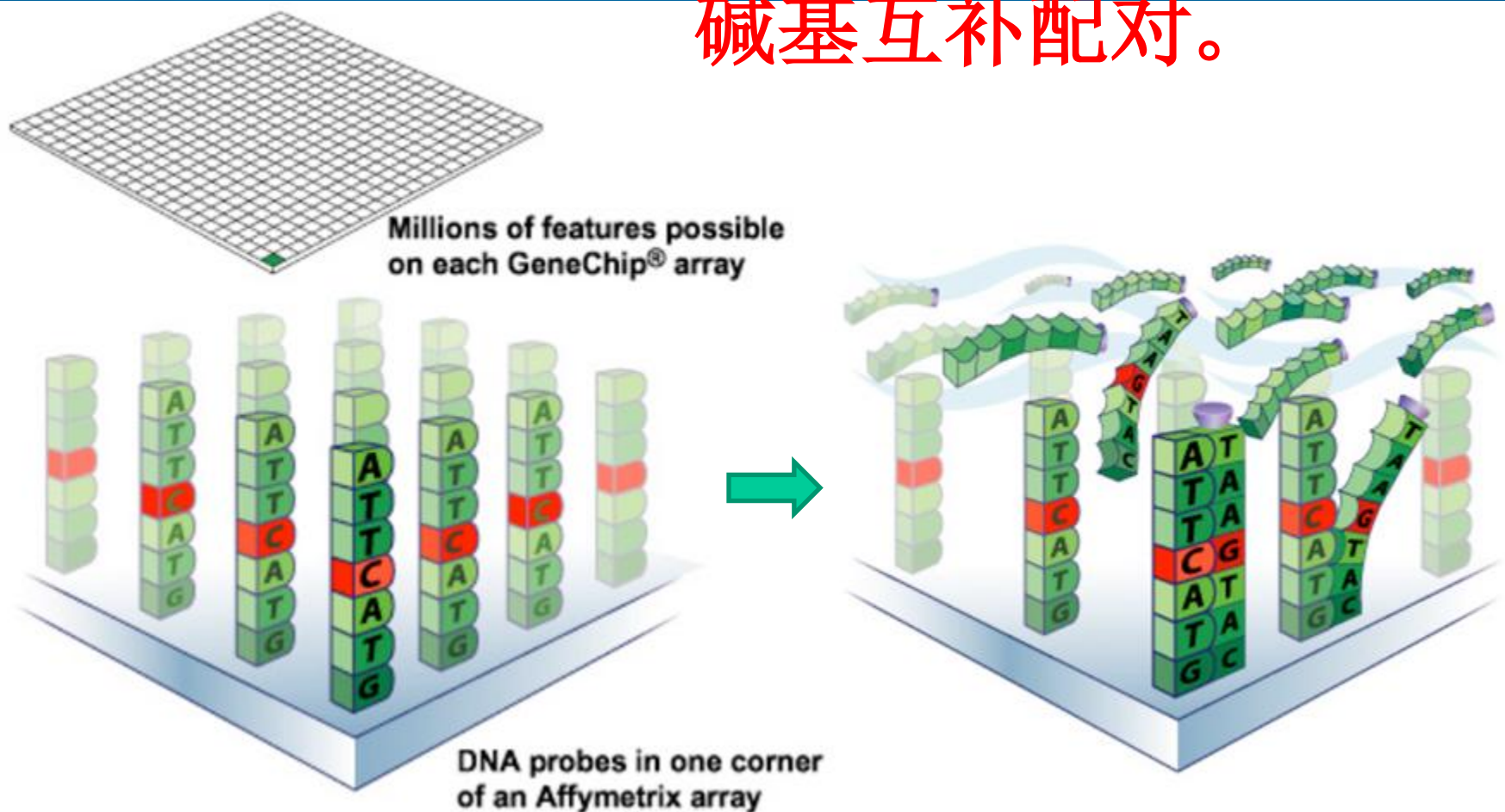


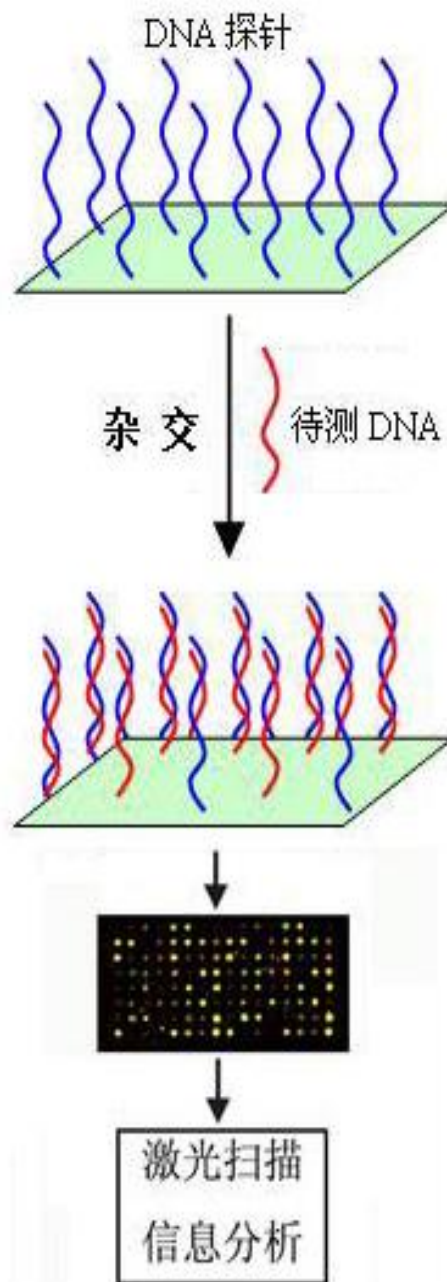
Microarray



# 基因芯片检测依赖核酸杂交原理

碱基互补配对。



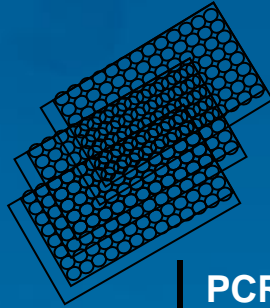


来自于待检测样本



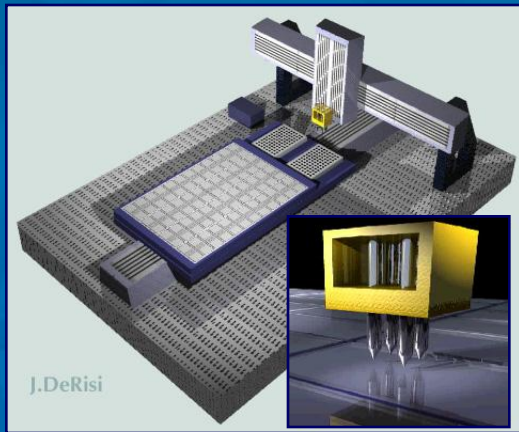
# 双色荧光基因芯片的实验流程

cDNA clones  
(target)



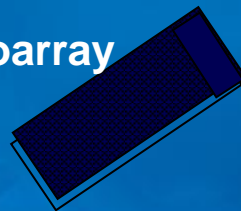
PCR product amplification  
purification

printing

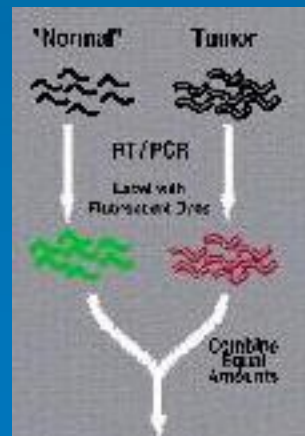


0.1nl/spot

microarray



mRNAprobe)



Hybridise target  
to microarray

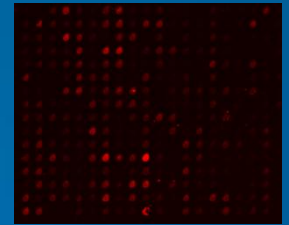
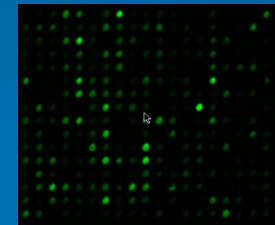
excitation

scanning

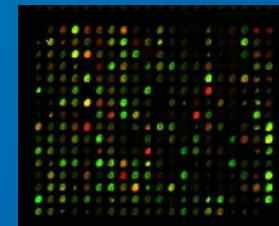
laser 2

laser 1

emission



overlay images and normalise



analysis

# 基因芯片的制作方法

- 原位合成
  - 光蚀保护 Affymetrix
  - Ink-jet Agilent
  - Digital Micromirror Device  
Nimbergene
  - 分子印章 东南大学
- 点样
  - cDNA PCR 产物
  - Oligo
  - 其它形式的基因产物

# 已有的微阵列制备技术

技术	探针类型	每个阵列的点数	实验室制备
接触点制	任意	~30,000	是
喷墨点制	任意	~20,000	是
无掩膜光刻技术	寡核苷酸	15,000	是
喷墨合成	寡核苷酸	240,000	否
基于掩膜的光刻技术	寡核苷酸	6,200,000	否
无掩膜光刻技术	寡核苷酸	2,100,000	否
微珠阵列	寡核苷酸	~50,000	否
电化学阵列	寡核苷酸	~15,000	否



# 主要微阵列供应商，以及他们是否提供定制芯片服务

公司	网站	定制芯片
Affymetrix	<a href="http://www.affymetrix.com/">http://www.affymetrix.com/</a>	是 *
Agilent Technologies	<a href="http://www.home.agilent.com/">http://www.home.agilent.com/</a>	是 *
Applied Biosystems	<a href="http://www.appliedbiosystems.com/">http://www.appliedbiosystems.com/</a>	否
Applied Microarrays	<a href="http://www.appliedmicroarrays.com/">http://www.appliedmicroarrays.com/</a>	是 *
CombiMatrix Corporation	<a href="http://www.combimatrix.com/">http://www.combimatrix.com/</a>	是
Illumina	<a href="http://www.illumina.com/">http://www.illumina.com/</a>	是
NimbleGen Systems*	<a href="http://www.nimblegen.com/">http://www.nimblegen.com/</a>	是 *
Oxford Gene Technologies	<a href="http://www.ogt.co.uk/">http://www.ogt.co.uk/</a>	是 *
Phalanx Biotech	<a href="http://www.phalanxbiotech.com/">http://www.phalanxbiotech.com/</a>	是
SuperArray Bioscience	<a href="http://www.superarray.com/">http://www.superarray.com/</a>	是

\* 已于2012年6月年关闭这个系统

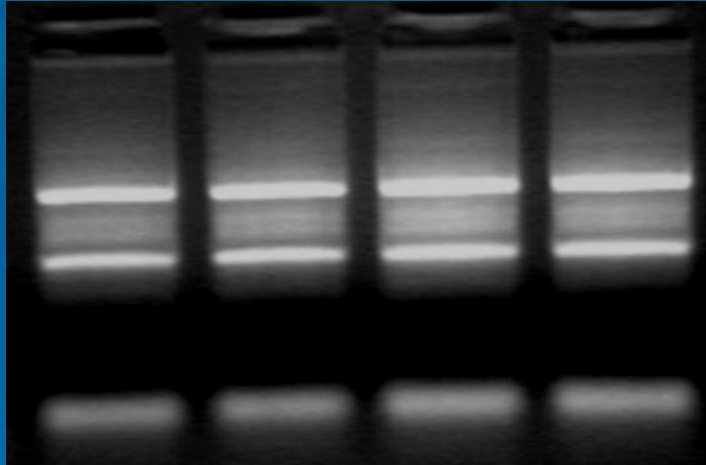
# 基因芯片使用样品的种类、处理和纯化

**DNA** 用常规的抽提方法的可以，要保证DNA的纯度，  
要求： $A_{260}/A_{280}$ 在1.6-1.8之间。

**RNA** { 总RNA：完整性和纯度，忠实性最好。  
mRNA：完整性和纯度，操作比较繁琐，容易降解。  
cRNA：需要通过线性放大，适用于样品量较少，如血液中分离的细胞，激光显微切割获得的细胞等。

新鲜的组织样本，细胞，FFPE组织，显微切割组织等

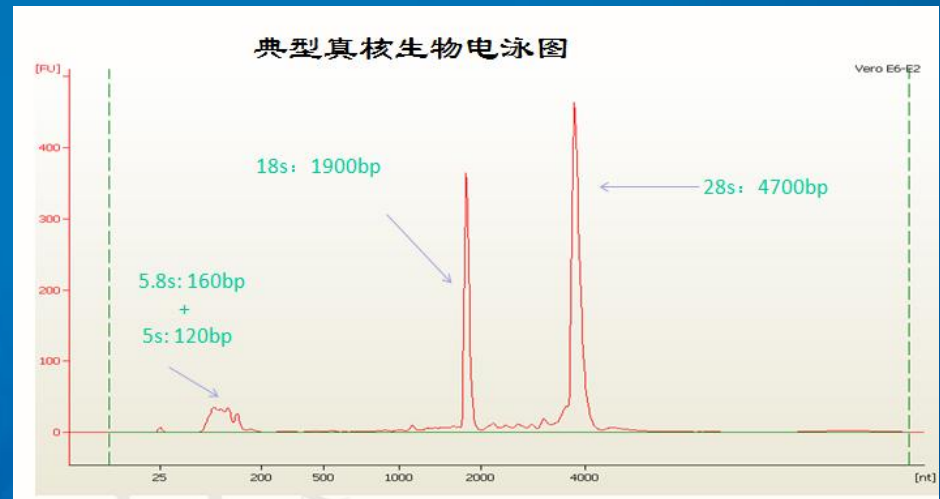
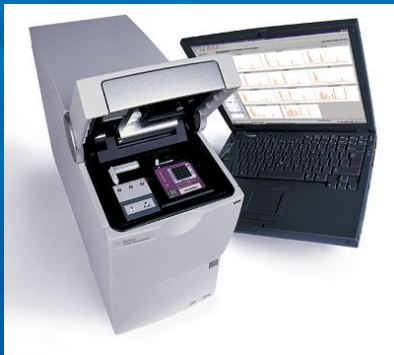
# 样本的质量检测



28S  
18S  
 $28S:18S > 2.0$



Total RNA



# 基因芯片的杂交

## ■ 样品的标记

标记方法

DNA: PCR, 随机引物, 缺口翻译。

RNA: 逆转录, 线性放大。

标记物

同位素

荧光染料(Cy3, Cy5)

化学发光

## ■ 芯片杂交

杂交体积(使核酸浓度增加

10万倍)

玻片: 2-200 $\mu$ l

滤膜: 5-50ml

杂交液和杂交液的组份(

杂交温度、时间

## ■ 芯片的洗涤

洗涤液的组成

洗涤的温度、时间



# 基因芯片的扫描

## ■ 激光共聚焦扫描

光源：特定波长的光

激发面积：<100平方微米

如：Agilent Scanner

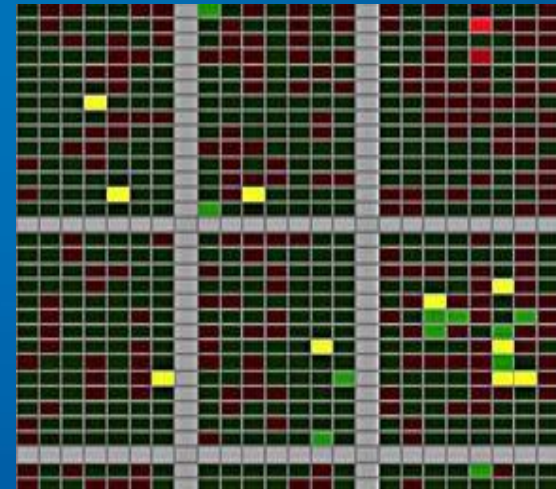
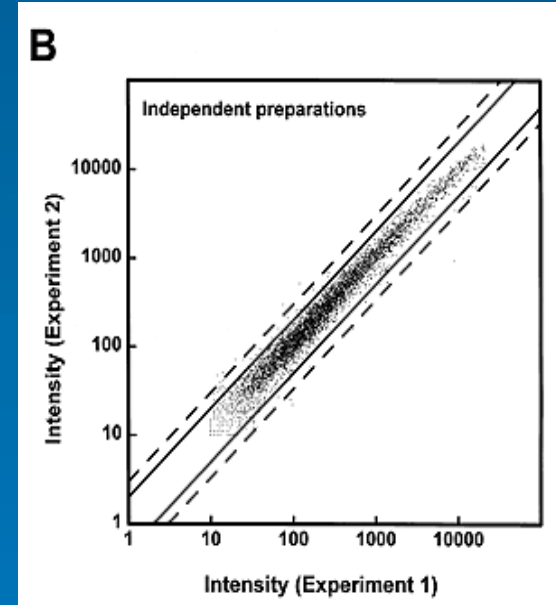
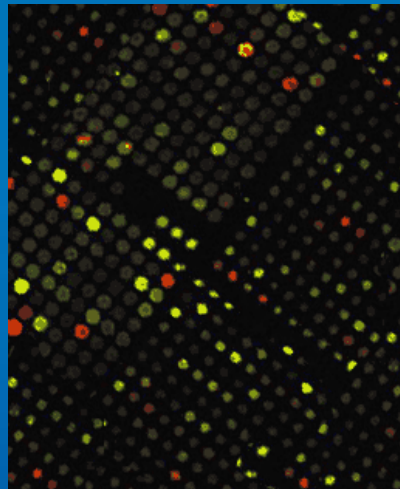
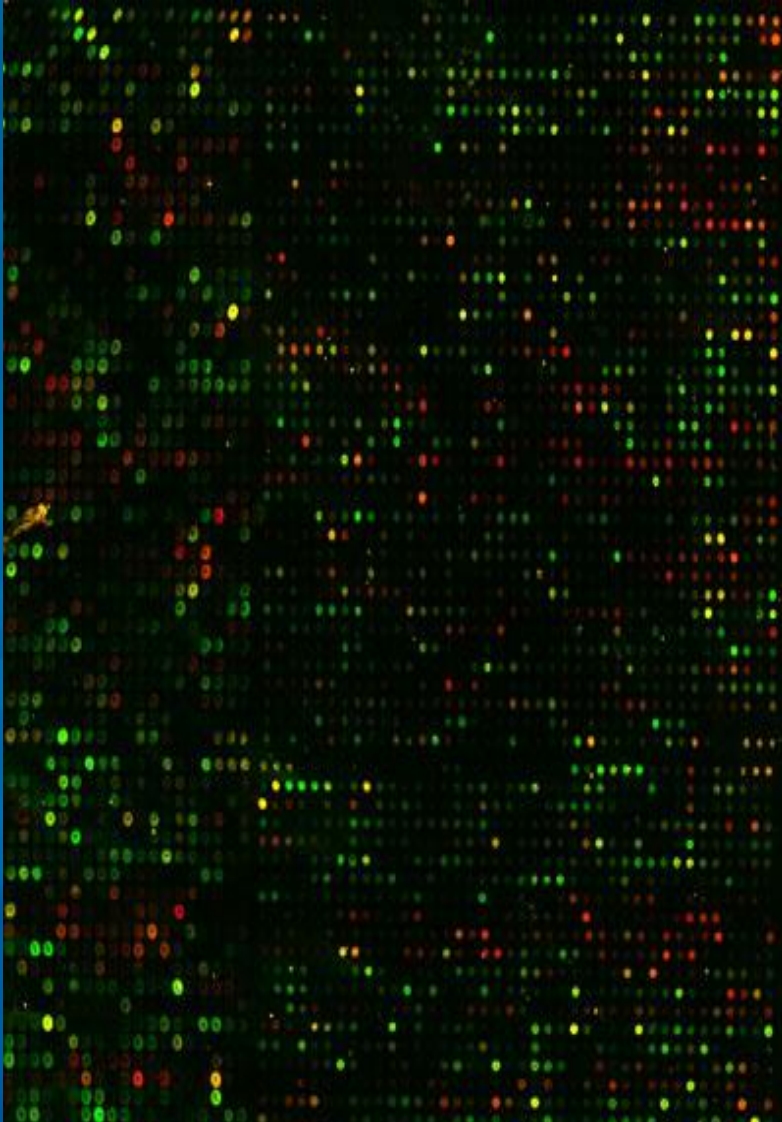


## ■ CCD（电荷耦合器件）成像

光源：连续波长的光（如弧光灯）

激发面积：同时激发多个1平方厘米的面积

# 基因芯片的数据的处理和分析



# 基因芯片的数据分析和解读

01

芯片图像转换

02

芯片数据的预  
处理

03

数据的质控  
(MAQC)

04

数据的归一化

05

差异基因

06

数据可视化

07

基因集分析

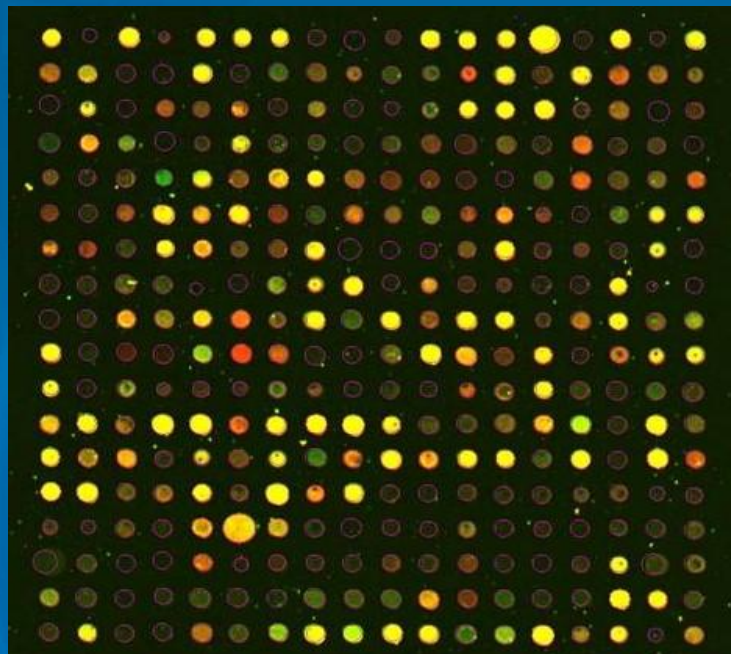
08

数据的存储和  
数据仓

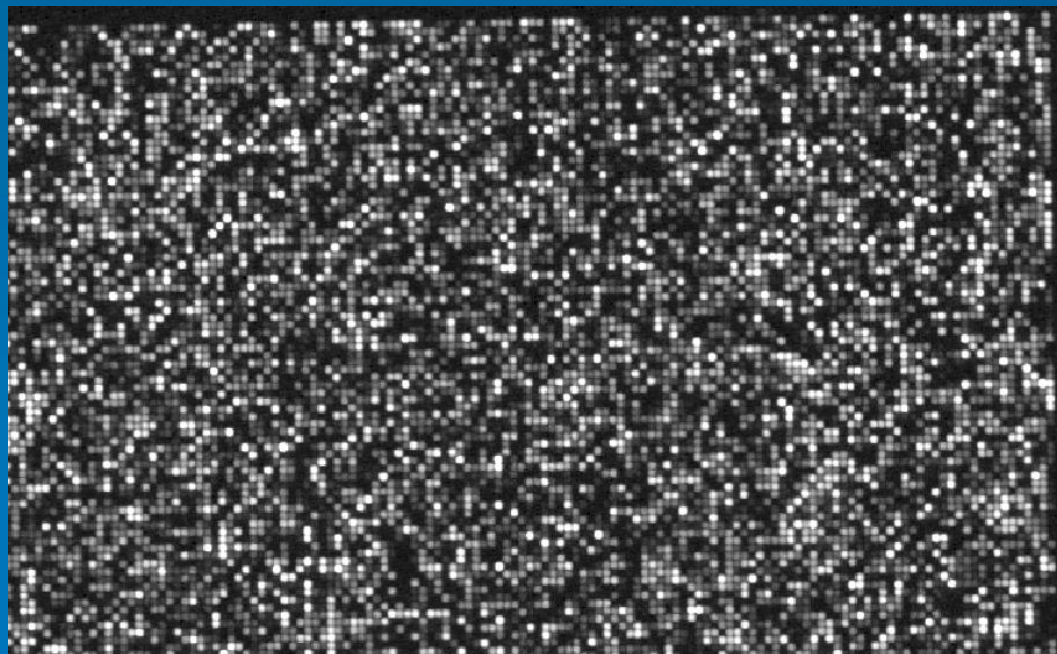


# 1. 芯片图像转换

双色荧光图



单色黑白图



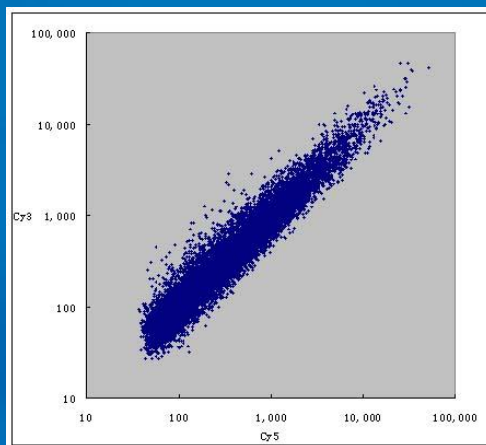
通过扫描仪带的图像软件转换成数字。



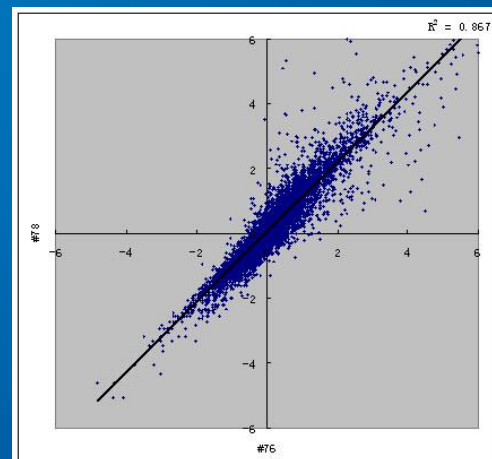
## 2. 数据的预处理

1. 阳性信号检出率：取阴性质控点如DMSO的信号值，分别计算其Cy3、Cy5信号平均值与2倍平均方差的和 ( $\text{mean} + 2 * \text{STDEV}$ )，Cy3、Cy5信号值分别大于这一值的基因点为检出点。
2. 将一些质量很低，数据可能不准确的点予以滤出。
3. 数据的分布：一般用散点图来表示。
4. 重复性一般用 $R^2$ 来表示。

片内



片间



### 3. 数据的质控

芯片数据的质量包括两个方面：整个芯片的质量和芯片上点的质量。

评价整张芯片质量的最简单也是最常用的方法是计算整个芯片的信噪比。信噪比低表示整张芯片背景高，芯片的质量差。

### 3. 数据的质控

信号点的质量影响体现在5个方面

- 信号点的大小和规则度
- 信噪比
- 信号点周围的背景强度
- 信号点背景的均一程度
- 信号的饱和程度

一般在分析时，会先确定质量低的点，并滤除这些点，这是数据处理的重要步骤。

### 3. 数据的质控

生物芯片质量控制协会（MAQC）于2006年在Nature Biotechnology杂志发表一系列的文章，提出了关于生物芯片数据可信度的重要发展问题。MAQC的主要任务包括：

- (1) 不同平台、不同芯片产生对照数据集
- (2) 建立对照RNA样本。
- (3) 芯片数据重复性的质量衡量标准。
- (4) 评估各种数据分析工具。

主要结论：经过规划慎密的实验设计，配合适当的数据转换，归一化及分析，基因芯片数据可以具有很好的重复性，以及不同设计，不同实验室及不同平台之间的数据具有可比性。



## 4. 数据的归一化

影响芯片原始数据的因素很多，在对芯片数据分析之前，必须进行数据的标准化（normalization）。数据的标准化是要减少芯片在处理过程中技术（系统）因素的影响，使检测的结果能真实地反映生物功能地差别，芯片的数据只有经过标准化处理后才具有可比性。

Normalization方法有按参与校正的基因分：

- 全部基因参与
- Housekeeping 基因参与
- Reference 基因参与

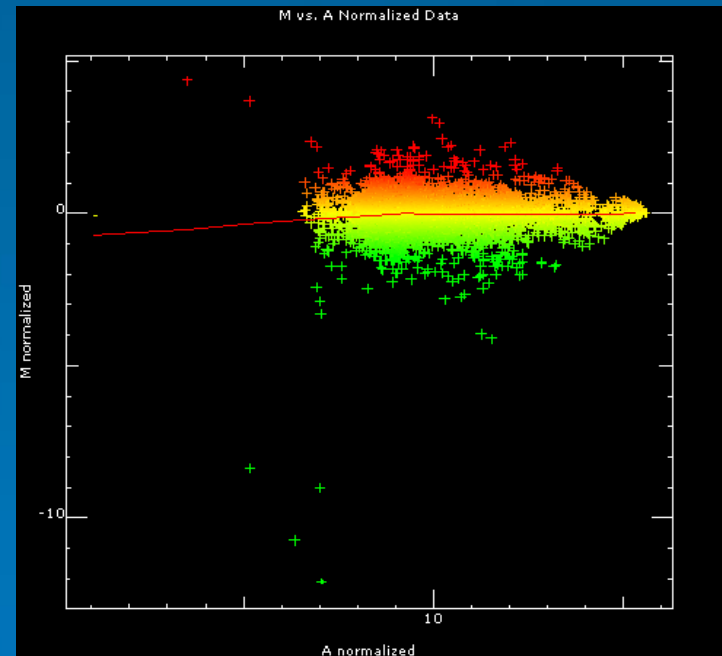
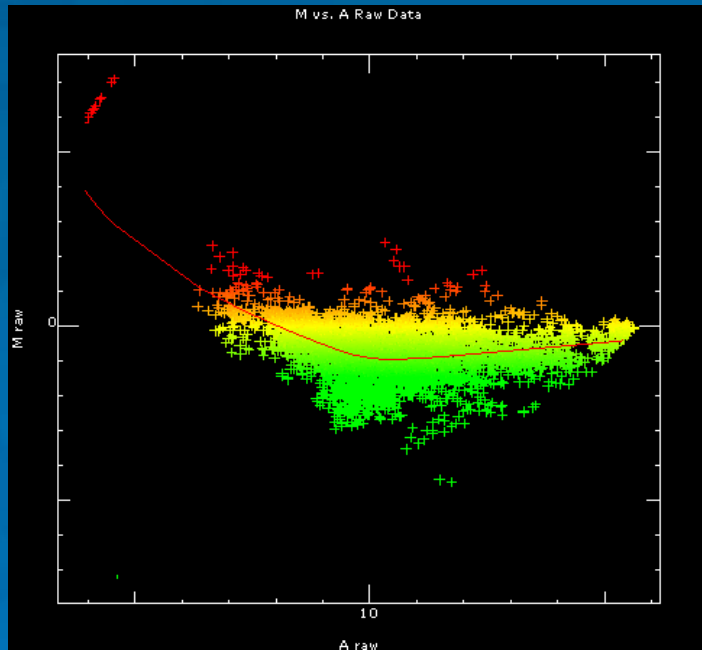
### 3. 数据的归一化

按值的估计方法分有3种：

- 全局法 (global method)
- 线性回归 (linear regression)
- 局部加权最小二乘法 (LOWESS)

无论那种方法, 目的使调整芯片技术导致的误差, 但不能纠正由于mRNA样本和芯片上探针基因本身所带来的误差。

# 数据的预处理 – 标准化前后的散点图



- 左图是未标准化处理的散点图
- 右图是经LOWESS处理的散点图

任何芯片数据进行分析前，都必须进行数据的标准化。

## 5. 差异基因

- 一般来说， $\text{ratio} > 2$  或  $\text{ratio} < 0.5$  认为是在两种样本表达有差异。这方法没有考虑到差异表达的统计显著性。
- Z-score ( $(X - \mu) / \sigma$ ),  $Z > 2$  表示比率在平均比率加两倍方差之外，差异基因有了统计意义。
- T-检验 (t-test)，从重复芯片中识别差异的表达基因。
- SAM。
- ANOVA。

现在识别差异基因方法已经有很多方法，它仍是数据处理中的一个热点。

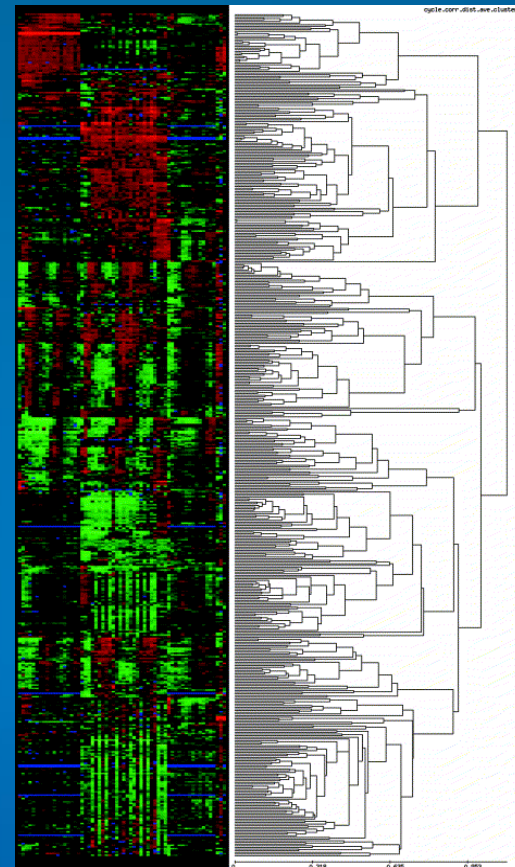


## 6. 数据的可视化(1) – 聚类分析

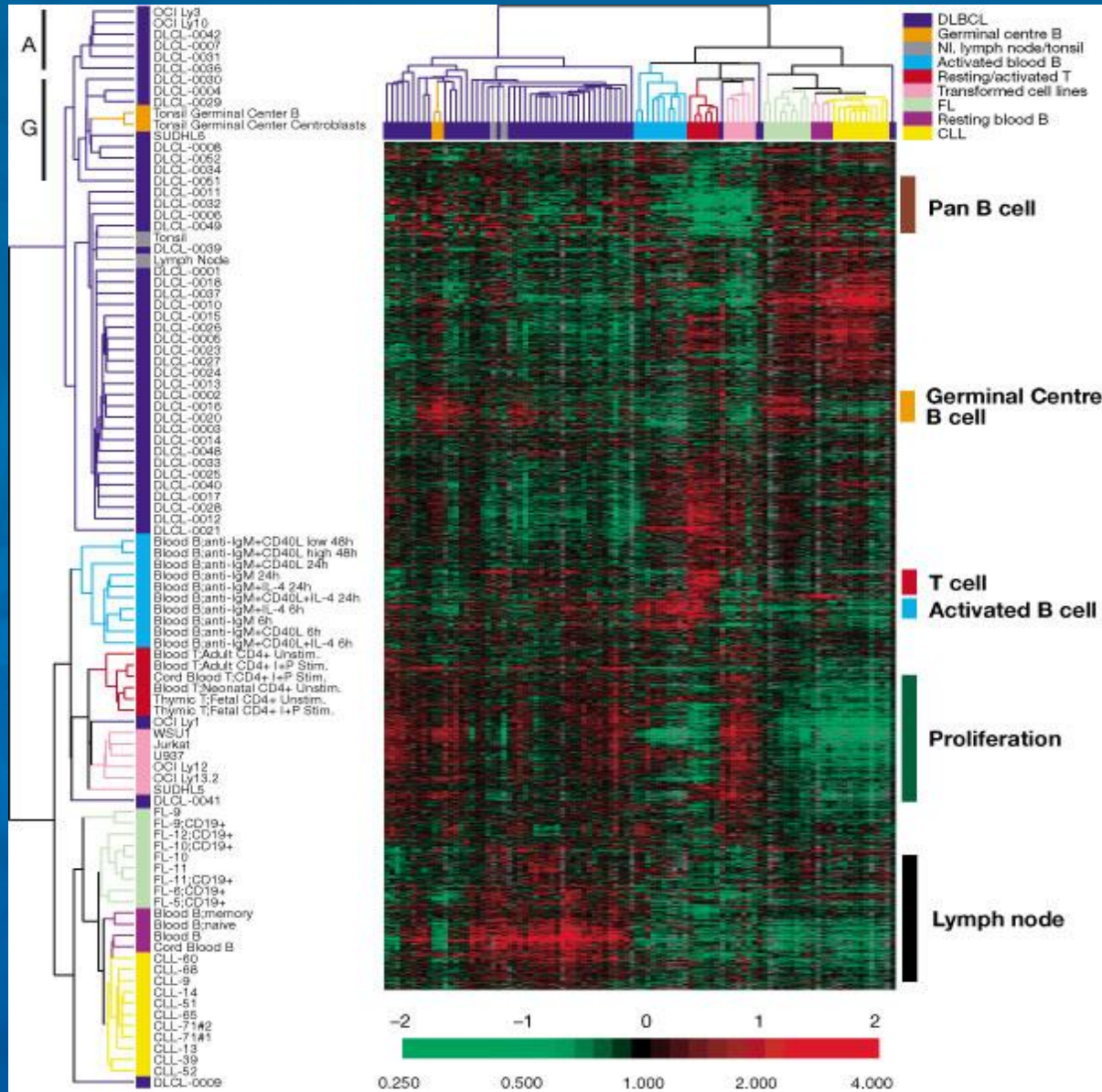
基因聚类分析的主要任务是确定相似表达模式的基因，相似的基因可能具有共同的特征。

探索完全未知的数据特征的方法

- 层次聚类(hierarchical cluster)
- K-means 基于向量的
- SOM
- 主成分分析(PCA)

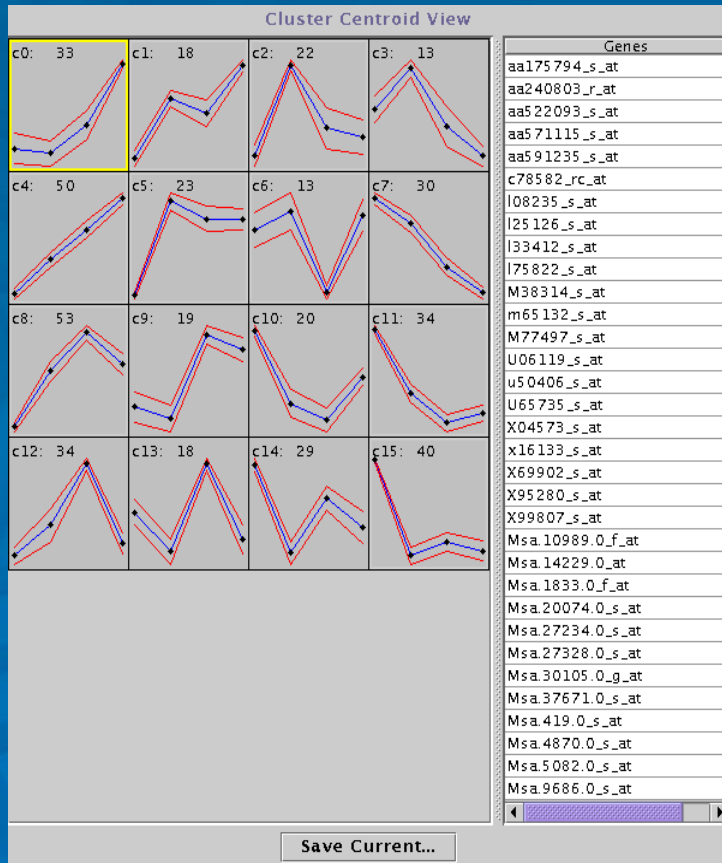


# 一个实例

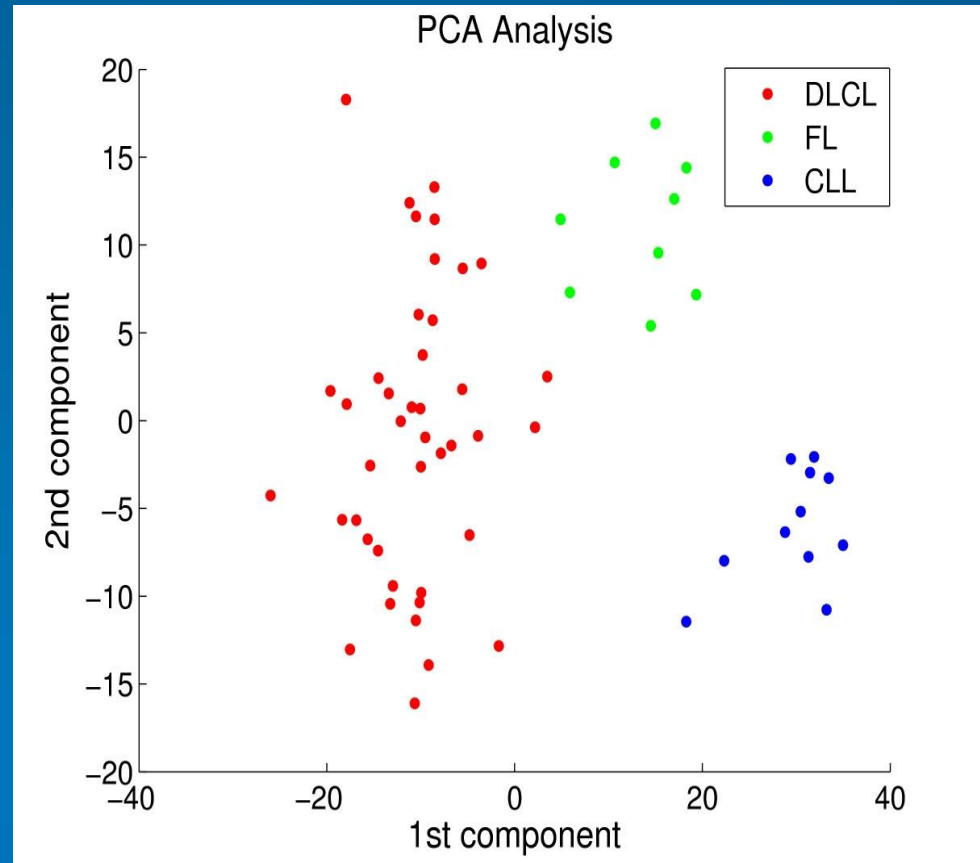


Nature Feb, 2000  
 Paper by  
 Allzadeh. A *et al*  
*Distinct types of  
 diffuse large  
 B-cell lymphoma  
 identified by  
 gene  
 expression  
 profiling*

# SOM



# PCA



## 7. 基因集分析

它是聚类分析的延伸，若有一些基因在聚类分析中被认为是同一群，我们可以判断 他们可能是同一个代谢途径

KEGG:

Kyoto Encyclopedia of Genes and Genomes

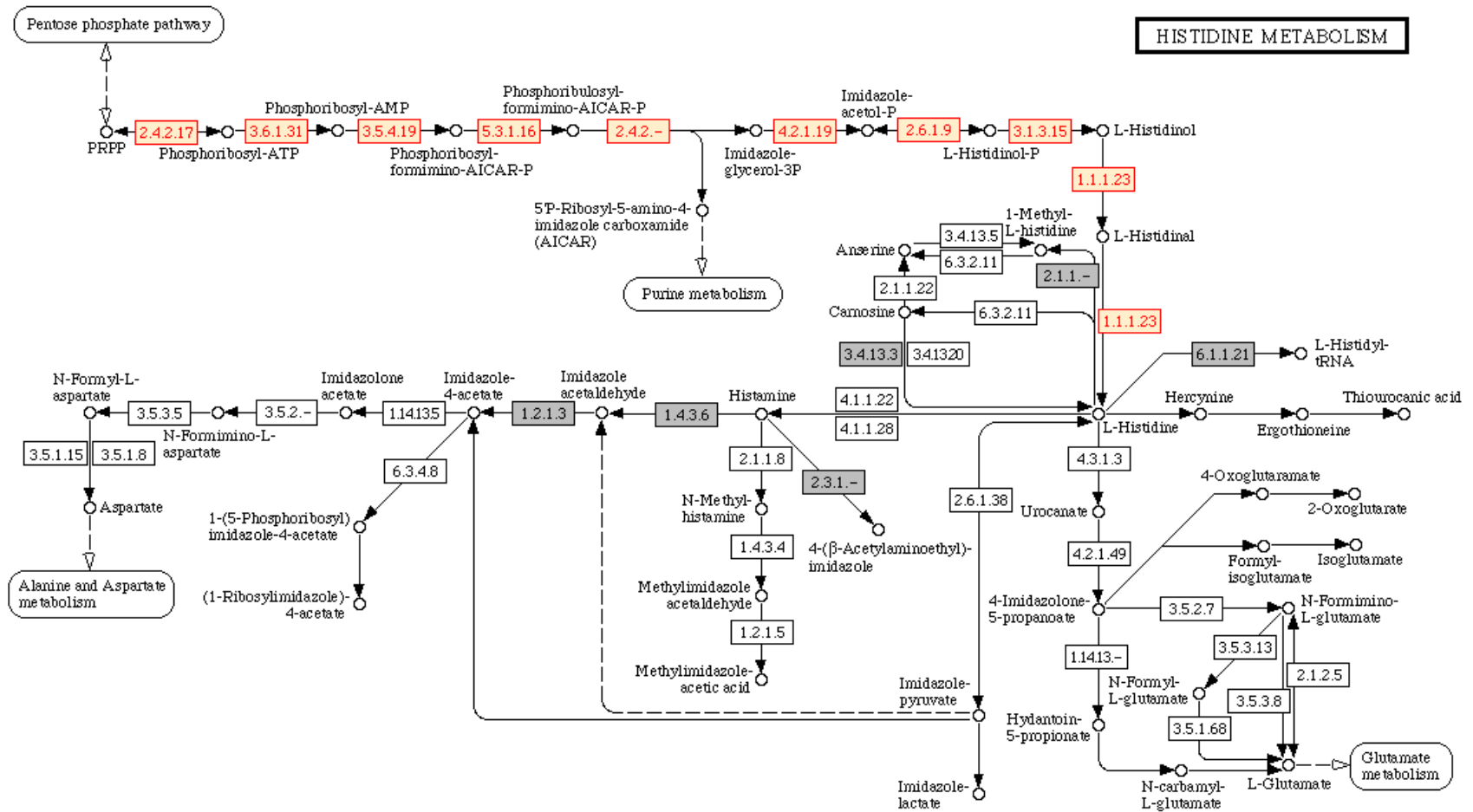
Software :

Pathway Assist; Pathway Editor; Pathway Finder

Network/pathway: 建议文献

Pilpel Y, etc. Identifying regulatory networks by combinatorial analysis of promoter elements. Nat Genet 2001 Oct;29(2):153-9.

# Pathway 分析





# 基因的分类 (GO)

## 分子功能

the tasks performed by individual gene products; examples are *carbohydrate binding* and *ATPase activity*

## 生物过程

broad biological goals, such as *mitosis* or *purine metabolism*, that are accomplished by ordered assemblies of molecular functions

## 亚细胞定位

subcellular structures, locations, and macromolecular complexes; examples include *nucleus*, *telomere*, and *origin recognition complex*

## 8. 数据存储和数据仓

为了设计和利用计算机软件或者数据库处理和整合各种数据来源，也为了共享研究数据，需要有统一的标准，主要有：

(1) 数据记录：基因芯片数据有效地存储和检索需要一定的形式去获取芯片设计、探针序列，实验过程以及基因表达数据本身的详细信息，MIAME原则。向公共数据库发表和递交芯片数据时，只接手符合MIAME原则的数据。

(2) 数据描述：芯片数据要共享，需要用大家都能理解的术语进行描述。包括使用标准语，提供标准的词汇表和正式关系的描述，如：芯片实验，芯片特征，实验和分析步骤，基因，样本等。

(3) 编制软件：MAGE (MicroArray Gene Expression, MAGE) 开发了一个系统，统一存储和交换不同数据系统之间的芯片数据。

## 8. 数据存储和数据仓

1. ArrayExpress。 (2003)
2. Gene Expression Omnibus(GEO, 2002)。
3. Center fo Information Biology gene Expression database (CiBEX, 2004)。

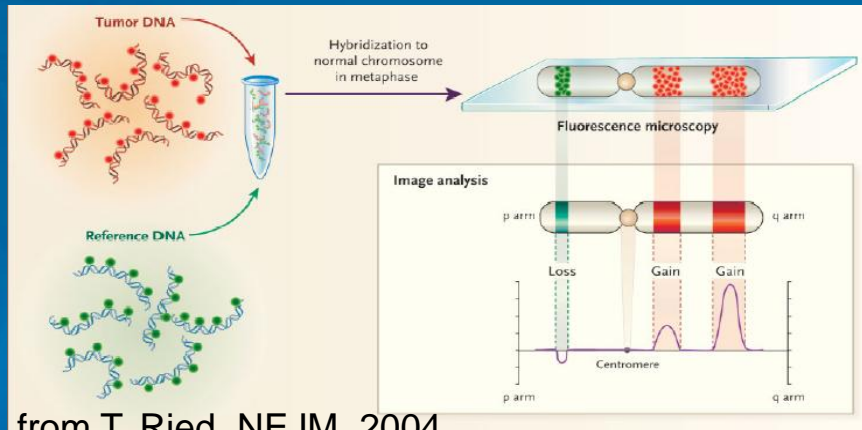
数据提交——数据管理——数据检索



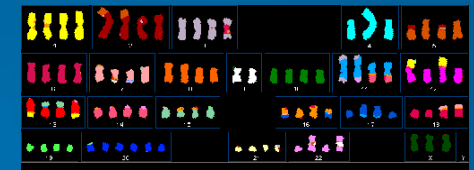
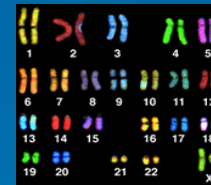
# 基因芯片在生命科学和医学 研究中的应用

# 在基因组水平的应用

aCGH : Array Based Comparative Genomic Hybridization

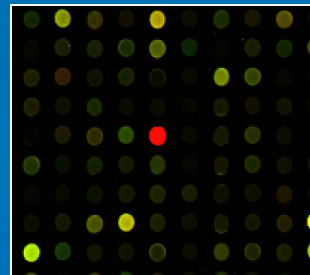
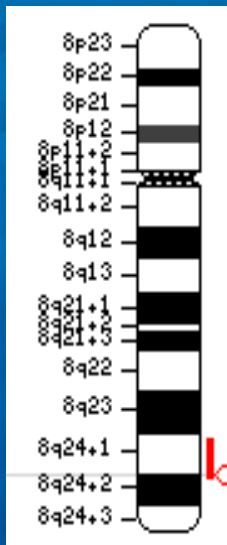


from T. Ried, NEJM, 2004



**Mechanisms of disease onset and progression,  
Identification of novel therapeutic targets,  
Drug resistance mechanisms and patient stratification  
Biomarkers for diagnostics and prognosis**

**Oligo aCGH**



- High-resolution
- High-throughput
- Quantitative
- Highly-flexible



# Correlation between genomic DNA copy number alterations and transcriptional expression in hepatitis B virus-associated hepatocellular carcinoma

## 1、HBV-induced HCC have unique pattern of DNA copy number alterations

(1) HCC (n=41) and HCC cell lines (n=12) occur significant DNA copy number alterations as compared with normal livers (n=2), and adjacent non-cancerous livers (n=5)

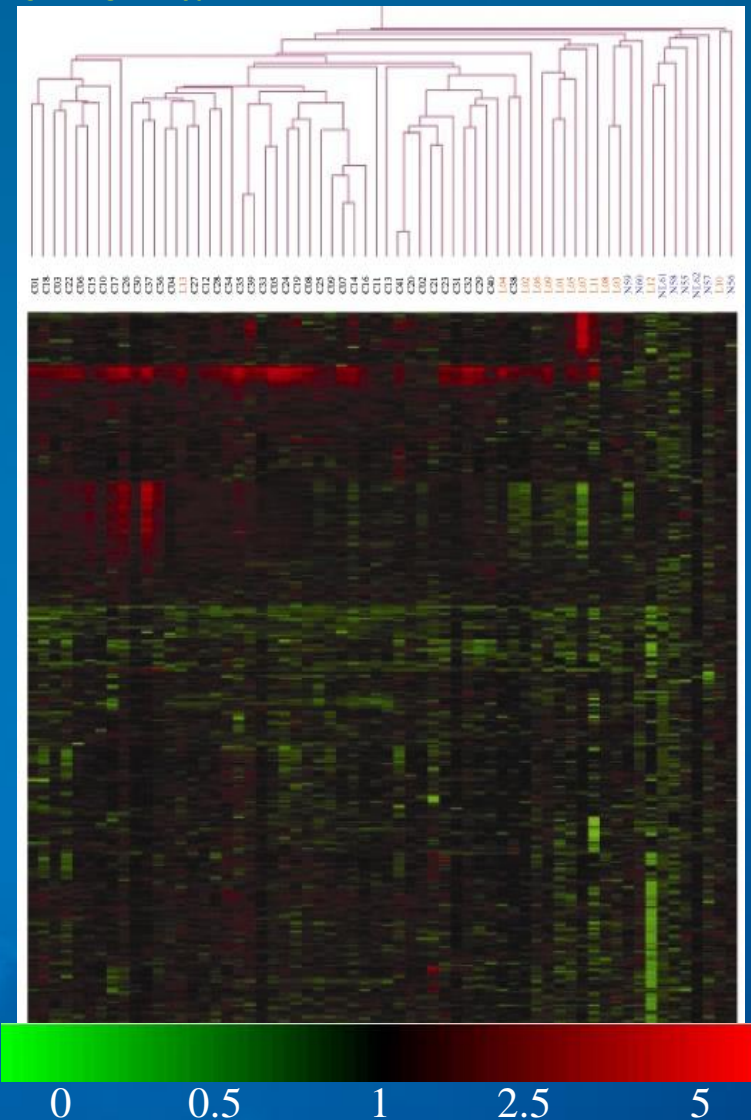
C=Cancer

N=non-cancerous liver

NL=normal liver

Red=amplification(gain)

Green=deletion(loss)



# High-resolution Survey of Human Lung Cancer

## High-resolution genomic profiles of human lung cancer

Giovanni Tonon<sup>\*,†</sup>, Kwok-Kin Wong<sup>\*,†,‡</sup>, Gautam Maulik<sup>\*,†</sup>, Cameron Brennan<sup>\*</sup>, Bin Feng<sup>\*</sup>, Yunyu Zhang<sup>\*</sup>, Deepak B. Khatry<sup>\*</sup>, Alexei Protopopov<sup>\*</sup>, Mingjian James You<sup>§</sup>, Andrew J. Aguirre<sup>\*</sup>, Eric S. Martin<sup>\*</sup>, Zhaoxui Yang<sup>\*</sup>, Hongbin Ji<sup>\*</sup>, Lynda Chin<sup>¶</sup>, and Ronald A. DePinho<sup>\*,§</sup>

<sup>\*</sup>Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA 02115; <sup>†</sup>Department of Pathology, Brigham and Women's Hospital, Boston, MA 02115; and Departments of <sup>‡</sup>Dermatology and <sup>§</sup>Genetics and Medicine, Harvard Medical School, Boston, MA 02115

Communicated by Webster K. Cavenee, University of California at San Diego, La Jolla, CA, May 18, 2005 (received for review April 13, 2005)

Lung cancer is the leading cause of cancer mortality worldwide, yet there exists a limited view of the genetic lesions driving this disease. In this study, an integrated high-resolution survey of regional amplifications and deletions, coupled with gene-expression profiling of non-small-cell lung cancer subtypes, adenocarcinoma and squamous-cell carcinoma (SCC), identified 93 focal copy-number alterations, of which 21 span <0.5 megabases and contain a median of five genes. Whereas all known lung cancer genes/loci are contained in the dataset, most of these recurrent copy-number alterations are previously uncharacterized and include high-amplitude amplifications and homozygous deletions. Notably, despite their distinct histopathological phenotypes, adenocarcinoma and SCC genomic profiles showed a nearly complete overlap, with only one clear SCC-specific amplicon. Among the few genes residing within this amplicon and showing consistent overexpression in SCC is *p63*, a known regulator of squamous-cell differentiation. Furthermore, intersection with the published pancreatic cancer comparative genomic hybridization dataset yielded, among others, two focal amplicons on 8p12 and 20q11 common to both cancer types. Integrated DNA-RNA analyses identified *WHSC1L1* and *TPX2* as two candidates likely targeted for amplification in both pancreatic ductal adenocarcinoma and non-small-cell lung cancer.

array comparative genomic hybridization | expression profiling | lung adenocarcinoma | squamous-cell lung carcinoma | TP73L

Lung cancer is the leading cause of cancer-related mortality in the United States, accounting for more than one-fourth of all cancer fatalities in 2004. Lung cancer is classified into two major subtypes, small-cell and non-small-cell lung cancer (NSCLC). NSCLC constitutes 75% of lung cancer cases and is subdivided further into three major histological subtypes: adenocarcinoma (AC), squamous-cell carcinoma (SCC), and large-cell carcinoma. The AC and SCC subtypes represent >85% of NSCLC cases. Although these NSCLC subtypes exhibit distinct pathological characteristics, the treatment approaches have remained generic and largely ineffective, despite advances in cytotoxic drugs, radiotherapy, and clinical management. For all stages of NSCLC, the 5-year survival rate has remained fixed at 15% for the last 15 years. The recent success of molecularly targeted therapies for a limited subset of cancer genotypes (1) has solidified the view that a more detailed knowledge of the spectrum of genetic lesions in lung cancer will, in turn, lead to meaningful therapeutic progress.

To date, the majority of lung cancer genetic studies have cataloged mutations or the promoter methylation status of known

particularly amplifications and deletions, suggests that only a small fraction of lung cancer genes has been identified. In particular, chromosomal CGH studies have revealed recurrent gains at 1q31, 3q25–27, 5p13–14, and 8q23–24 and deletions at 3p21, 8p22, 9p21–22, 13q22, and 17p12–13 (3–7). A recent array-CGH (aCGH) survey of known genes/loci using 348 BAC clones has confirmed recurrent chromosome-3p deletions and -3q gains and identified *PIK3CA* as a resident of the chromosome-3q amplicon (8).

Integrated CGH and expression profiling have emerged as effective entry points for cancer gene discovery, capable of providing a high-resolution view of the regional gains and losses throughout the cancer genome (9) and the associated copy-number-driven changes in gene expression (10, 11). In the microarray format, the resolution of CGH is dictated by the number and quality of mapped probes positioned along the genome (12). In this study, high-density gene-specific arrays were used to conduct high-resolution surveys of CNAs present in a collection of primary ACs and SCCs and of established NSCLC cell lines. Together with expression profiling, these datasets provide insights into the origins of, and genetic mechanisms driving, AC and SCC subtypes.

### Materials and Methods

**Cell Lines and Primary Tumors.** All of the primary tumors were acquired from the Cooperative Human Tissue Network (Philadelphia) and the Brigham and Women's Hospital tissue bank (Boston) under an approved institutional protocol. The tumor histology was confirmed by a pathologist (M.J.Y.) before inclusion in this study. All of the cell lines were obtained from the American Type Culture Collection. The characteristics of the primary tumors and cell lines are detailed in Tables 2 and 3, respectively, which are published as supporting information on the PNAS web site. Three independent, normal RNA references isolated from adjacent, histologically normal lung tissues were used as the normal control for the expression analysis.

**aCGH Profiling on Oligonucleotide and cDNA Microarrays.** Genomic DNAs from cell lines and primary tumors were extracted according to manufacturer's instructions (Gentra Systems). Genomic DNA was fragmented and random-prime labeled as described in ref. 11 and <http://genomic.dfci.harvard.edu/array.CGH.htm> and hybridized to either human cDNA or oligonucleotide microarrays. The cDNA microarray contains

Freely available online through the PNAS open access option.

Abbreviations: AC, adenocarcinoma; CGH, comparative genomic hybridization; aCGH,

## Application: Copy Number (aCGH)

- Integrating high-resolution gene amplification and deletion data with gene expression profiling
- 93 focal copy-number alterations
  - 74 amplifications
  - 19 deletions
  - median size = 1.53 MB
  - 21 within highly focal subset with median size < 0.5 MB (spanning ~ 5 genes)
  - 7 cross-tumor-type candidates
- Rational starting point for productive gene-discovery efforts



DANA-FARBER  
CANCER INSTITUTE

# Autism Association

Science**express**

Report

## Strong Association of De Novo Copy Number Mutations with Autism

Jonathan Sebat,<sup>1\*</sup> B. Lakshmi,<sup>1</sup> Dheeraj Malhotra,<sup>1†</sup> Jennifer Troge,<sup>1†</sup> Christa Lese-Martin,<sup>2</sup> Tom Walsh,<sup>3</sup> Boris Yanrom,<sup>1</sup> Seungtae Yoon,<sup>1</sup> Alex Krasnitz,<sup>1</sup> Jude Kendall,<sup>1</sup> Anthony Leotta,<sup>1</sup> Deepa Pai,<sup>1</sup> Ray Zhang,<sup>1</sup> Yoon-Ha Lee,<sup>1</sup> James Hicks,<sup>1</sup> Sarah J Spence,<sup>4</sup> Annette T. Lee,<sup>5</sup> Kaija Puura,<sup>6</sup> Terho Lehtimäki,<sup>7</sup> David Ledbetter,<sup>2</sup> Peter K. Gregersen,<sup>5</sup> Joel Bregman,<sup>8</sup> James S. Sutcliffe,<sup>9</sup> Vaidehi Jobanputra,<sup>10</sup> Wendy Chung,<sup>10</sup> Dorothy Warburton,<sup>10</sup> Mary-Claire King,<sup>3</sup> David Skuse,<sup>11</sup> Daniel H Geschwind,<sup>12</sup> T. Conrad Gilliam,<sup>13</sup> Kenny Ye,<sup>14</sup> Michael Wigler<sup>1\*</sup>

<sup>1</sup>Cold Spring Harbor Laboratory, 1 Bungtown Road, Cold Spring Harbor, NY 11724, USA. <sup>2</sup>Department of Human Genetics, Emory University School of Medicine, Atlanta, GA 30322, USA. <sup>3</sup>Department of Medicine and Genome Sciences, University of Washington, Seattle, WA 98195-7720, USA. <sup>4</sup>Pediatrics and Neurodevelopmental Psychiatry Branch, National Institute of Mental Health, National Institutes of Health, Bethesda, MD 20892-1255, USA. <sup>5</sup>Feinstein Institute for Medical Research, North Shore-Long Island Jewish Health System, Manhasset, NY 11030, USA. <sup>6</sup>Department of Child Psychiatry, University of Tampere, Medical School, Tampere, Finland. <sup>7</sup>Department of Clinical Chemistry, University Hospital of Tampere and University of Tampere, Medical School, Tampere, Finland. <sup>8</sup>Fay J. Lindner Center for Autism and Developmental Disorders, North Shore-Long Island Jewish Health System, 4300 Hempstead Turnpike, Bethpage, NY 11714, USA. <sup>9</sup>Center for Molecular Neuroscience, Vanderbilt University, Nashville, TN 37232-8548, USA. <sup>10</sup>Departments of Genetics and Development, and Pediatrics, Columbia University, New York, NY 10027, USA. <sup>11</sup>Behavioural and Brain Sciences Unit, Institute of Child Health, University College London, 30 Guilford Street, London, WC1N 1EH, UK. <sup>12</sup>Interdepartmental Program in the Neurosciences, Program in Neurogenetics, Neurology Department, David Geffen School of Medicine, University of California at Los Angeles, Los Angeles, CA 90095-1769, USA. <sup>13</sup>Department of Human Genetics, The University of Chicago, 920 East 58th Street, Chicago, IL 60637, USA. <sup>14</sup>Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx, NY 10461, USA.

\*To whom correspondence should be addressed. E-mail: [sebat@cshl.edu](mailto:sebat@cshl.edu); [wigler@cshl.edu](mailto:wigler@cshl.edu)

†These authors contributed equally to this work

We tested the hypothesis that de novo copy number variation (CNV) is associated with autism spectrum disorders (ASDs). We performed comparative genomic hybridization (CGH) on the genomic DNA of patients and unaffected subjects to detect copy number variants not present in their respective parents. Candidate genomic regions were validated by higher-resolution CGH, paternity testing, cytogenetics, fluorescence in situ hybridization, and microsatellite genotyping. Confirmed de novo CNVs were significantly associated with autism ( $P = 0.0005$ ). Such CNVs were identified in 12 out of 118 (10%) of patients with sporadic autism, in 2 out of 77 (2%) of patients with an affected first-degree relative, and in 2 out of 196 (1.0%) of controls. Most de novo CNVs were smaller than microscopic resolution. Affected genomic regions were highly heterogeneous and included mutations of single genes. These findings establish de novo germline mutation as a more significant risk factor for ASD than previously recognized.

Autism spectrum disorders (ASDs) [MIM 209850] are characterized by language impairments, social deficits and repetitive behaviors. The onset of symptoms occurs by the

age of 3, and usually requires extensive support for the lifetime of the afflicted. The prevalence of ASD is estimated to be 1 in 166 (1), making it a major burden to society.

Genetics plays a major role in the etiology of autism. The concordance rates in monozygotic twins are 70% for autism and 90% for ASD, while the concordance rates in dizygotic twins are 5% and 10% respectively. Previous studies suggest autism displays a high degree of genetic heterogeneity. Efforts to map disease genes using linkage analysis have found evidence for autism loci on 20 different chromosomes. Regions implicated by multiple studies include 1p, 5q, 7q, 15q, 16p, 17q, 19p and Xq (2). Moreover, microscopy studies have identified cytogenetic abnormalities in >5% of affected children, involving many different loci on all chromosomes (3). In some rare syndromic forms of autism, such as Rett syndrome (4) and tuberous sclerosis (5), mutations in a single gene have been identified. Otherwise, neither linkage nor cytogenetics has unambiguously identified specific genes involved.

Genetic heterogeneity poses a considerable challenge to traditional approaches for gene mapping (6). Some of these limitations are overcome by methods which rely on the direct

## Utilized ROMA and Agilent 244k arrays for CGH analysis.

### Key conclusions:

- Each de novo CNV was rare in patient population
- Lesions at different loci can contribute to autism
- Clear evidence that the two classes of autism, familial and sporadic, are genetically distinct
- Methods for detecting CNVs genome wide are a powerful alternative to traditional gene mapping approaches



# CNV variation in copy number in the human genome

nature

Vol 444 | 23 November 2006 | doi:10.1038/nature05329

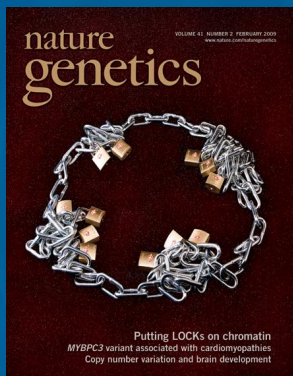
## ARTICLES

### Global variation in copy number in the human genome

Richard Redon<sup>1</sup>, Shumpei Ishikawa<sup>2,3</sup>, Karen R. Fitch<sup>4</sup>, Lars Feuk<sup>5,6</sup>, George H. Perry<sup>7</sup>, T. Daniel Andrews<sup>1</sup>, Heike Fiegler<sup>1</sup>, Michael H. Shapero<sup>4</sup>, Andrew R. Carson<sup>5,6</sup>, Wenwei Chen<sup>4</sup>, Eun Kyung Cho<sup>7</sup>, Stephanie Dallaire<sup>7</sup>, Jennifer L. Freeman<sup>7</sup>, Juan R. González<sup>8</sup>, Mònica Gratacòs<sup>8</sup>, Jing Huang<sup>4</sup>, Dimitrios Kalaitzopoulos<sup>1</sup>, Daisuke Komura<sup>3</sup>, Jeffrey R. MacDonald<sup>5</sup>, Christian R. Marshall<sup>5,6</sup>, Rui Mei<sup>4</sup>, Lyndal Montgomery<sup>1</sup>, Kunihiro Nishimura<sup>2</sup>, Kohji Okamura<sup>5,6</sup>, Fan Shen<sup>4</sup>, Martin J. Somerville<sup>9</sup>, Joelle Tchinda<sup>7</sup>, Armand Valsesia<sup>1</sup>, Cara Woodwark<sup>1</sup>, Fengtang Yang<sup>1</sup>, Junjun Zhang<sup>5</sup>, Tatiana Zerjal<sup>1</sup>, Jane Zhang<sup>4</sup>, Lluís Armengol<sup>8</sup>, Donald F. Conrad<sup>10</sup>, Xavier Estivill<sup>8,11</sup>, Chris Tyler-Smith<sup>1</sup>, Nigel P. Carter<sup>1</sup>, Hiroyuki Aburatani<sup>2,12</sup>, Charles Lee<sup>7,13</sup>, Keith W. Jones<sup>4</sup>, Stephen W. Scherer<sup>5,6</sup> & Matthew E. Hurles<sup>1</sup>

Copy number variation (CNV) of DNA sequences is functionally significant but has yet to be fully ascertained. We have constructed a first-generation CNV map of the human genome through the study of 270 individuals from four populations with ancestry in Europe, Africa or Asia (the HapMap collection). DNA from these individuals was screened for CNV using two complementary technologies: single-nucleotide polymorphism (SNP) genotyping arrays, and clone-based comparative genomic hybridization. A total of 1,447 copy number variable regions (CNVRs), which can encompass overlapping or adjacent gains or losses, covering 360 megabases (12% of the genome) were identified in these populations. These CNVRs contained hundreds of genes, disease loci, functional elements and segmental duplications. Notably, the CNVRs encompassed more nucleotide content per genome than SNPs, underscoring the importance of CNV in genetic diversity and evolution. The data obtained delineate linkage disequilibrium patterns for many CNVs, and reveal marked variation in copy number among populations. We also demonstrate the utility of this resource for genetic disease studies.

# GWAS应用案例：银屑病易感基因LCE研究



## 研究背景：

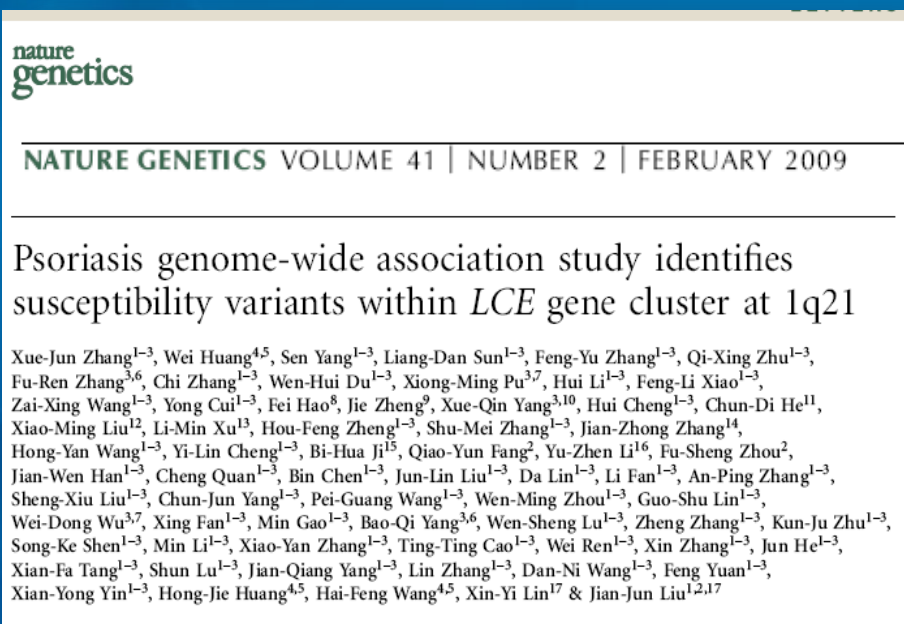
银屑病的发病率占世界人口的0.1%~3%，黄种人发病率为0.1%~0.3%。截至2007年，我国银屑病患者已经达到458万人，但发病的基因背景尚不清楚。SBC与安徽医科大学合作，采用全基因组关联分析手段，在大量样本中进行了中国汉族人银屑病易感基因的搜寻和鉴定工作，探讨这类疾病的发病机理和遗传易感性。

Nat Genet. 2009 Feb;41(2):205-10.

## 研究结果：

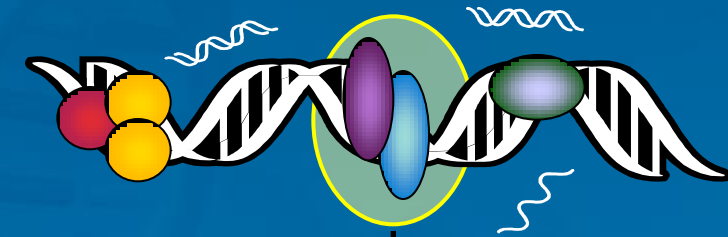
该研究发现了一个新的位于染色体1q21的LCE易感基因变异体，同时验证了已知的两个易感位点：MHC和IL12B。

LCE基因编码表皮终末分化角质外膜蛋白，该研究发现LEC基因变异与皮肤表皮细胞更新速度异常有关，其基因变异可以增加患病风险。

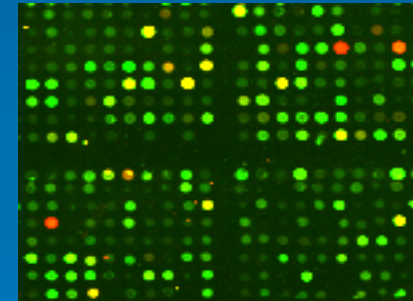




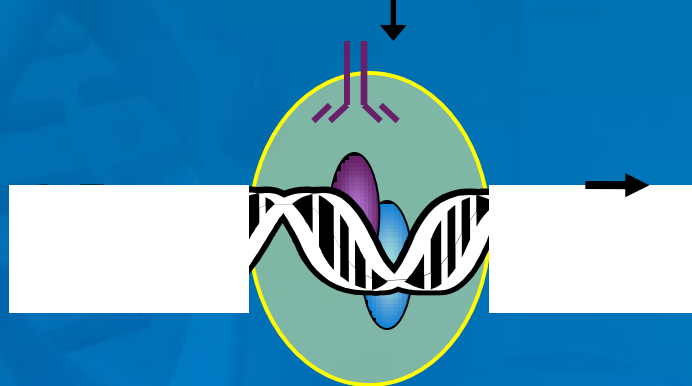
# 研究DNA和蛋白的相互作用



Hybridize to microarray  
for detection

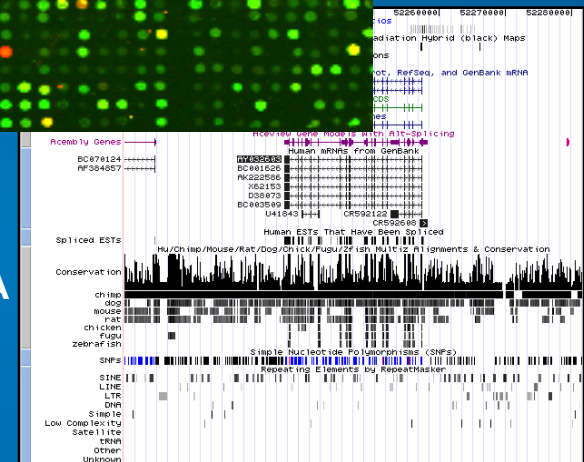


to promoter DNA regions



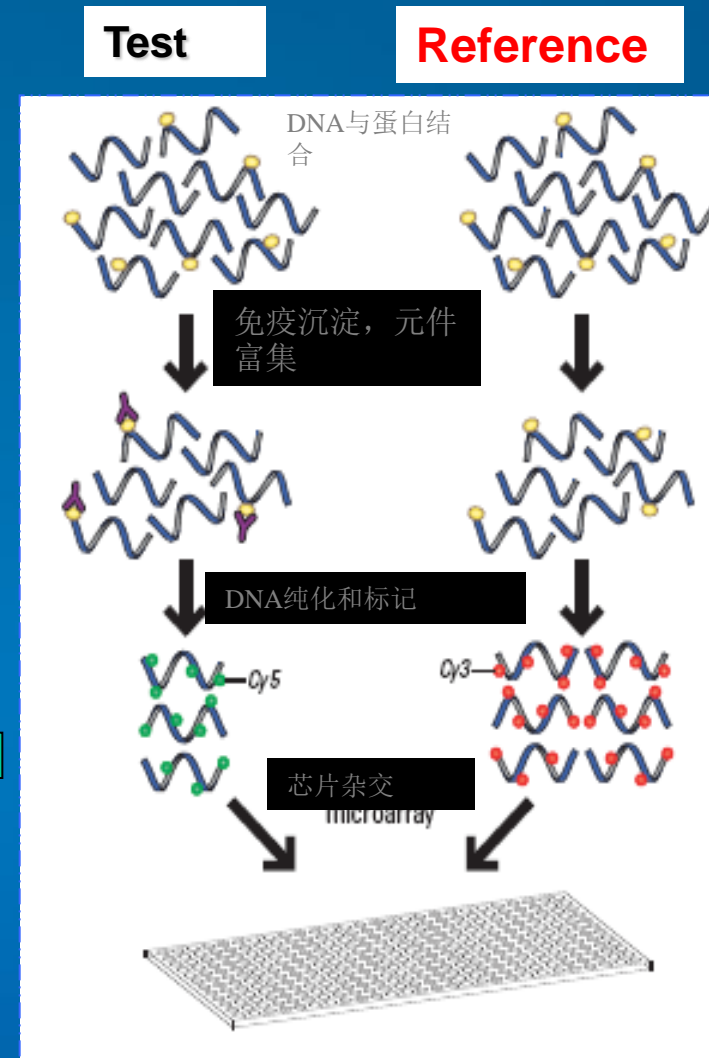
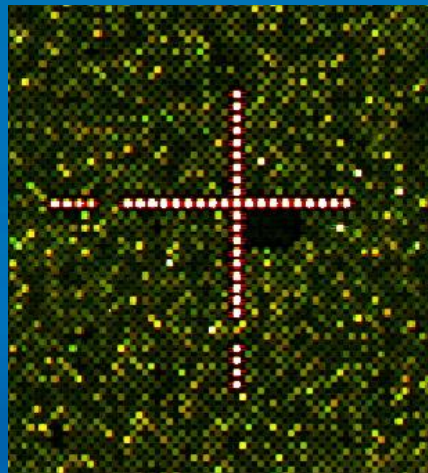
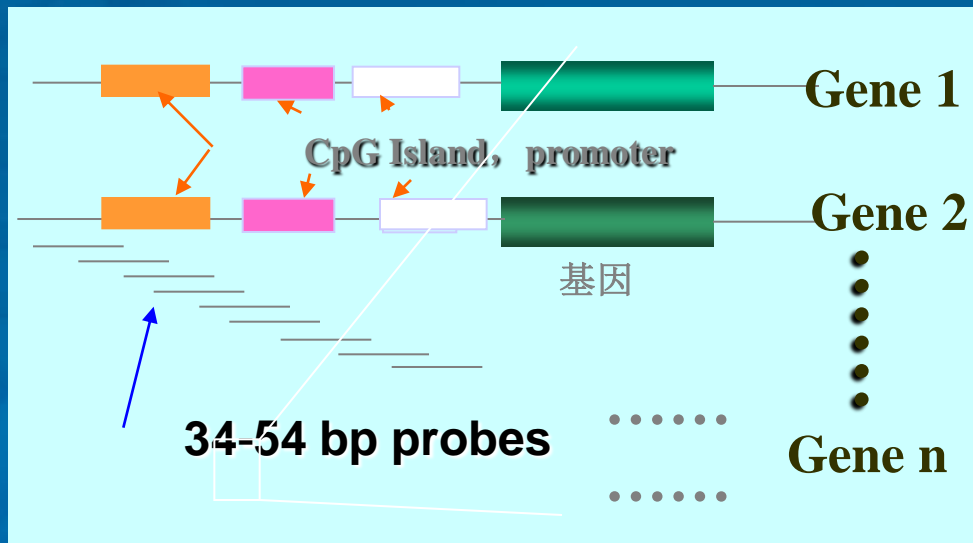
Enriched pool  
of protein-bound DNA

Crosslink protein-DNA complexes  
IP chromatin to enrich and capture bound  
DNA  
Purify and amplify DNA fragments



ChIP Analytics output

# DNA 甲基化芯片



# 基因芯片在基因表达水平的应用

## Identification of genes expression profile of dorsal root ganglion in the rat peripheral axotomy model of neuropathic pain

Normal rats      Rats 14 days peripheral nerve injury



Rat DRG cDNA libraries



7,523 genes and expressed sequence tags (ESTs)



cDNA array



NCBI rat DRG unigene library



RT-PCR

In situ hybridization

Northern blotting



Gene expression profile



Xiao et al., (2002) PNAS. 99: 8360-8365

Physiological and pharmacological  
analysis of important molecules

# Markedly regulated 122 genes and 51 novel ESTs

## Neuropeptides

## Receptors, membrane proteins

## Channels

## Signal transduction modulators and effectors

## Synaptic transmission

Accession NO.	Gene name	2d	7d	14d	28d	Accession NO.	Gene name	2d	7d	14d	28d
<b>A Neuropeptide</b>						<b>F Growth-associated protein</b>					
V01231	calcitonin gene-related peptide					U95001	developmentally-regulated carline factor				
NM_012659	calcitonin					AF271786	fibroblast growth factor 15				
M15191	tachykinin					M31897	insulin-like growth factor binding protein 5				
S70690	cholecystokinin					D45201	neurofibronin				
U95624	glucagon					D38629	adenomatous polyposis coli protein				
M20575	neuropeptide Y					M22427	basic fibroblast growth factor				
X00241	vasoactive intestinal polypeptide					U13253	DA11=152 kDa fatty acid binding protein				
<b>B Receptor and membrane protein</b>						X56551	fibroblast growth factor 7				
M58316	adenergic receptor alpha 2B					M16218	growth associated protein 43 (GAP43)				
D38450	G protein-coupled receptor, partial cds					L32391	growth arrest and DNA-damage-inducible protein 45				
A3061445	G protein-coupled receptor, LCR4					M69325	nerve growth factor-inducible protein				
M50518	metabotropic glutamate receptor 4					A306458	superior cervical ganglion 30				
X27121	metabotropic receptor type 2					<b>G Cytoskeleton and cell motility</b>					
Z11584	neuropeptide Y Y1 receptor					845817	high molecular weight neurofilament				
L20894	opioid receptor, non-type					Z11152	neBk molecular weight neurofilament				
L34237	prostaglandin F2 alpha receptor					A3031880	light molecular weight neurofilament				
X50132	serotonin receptor					V01217	cytoplasmic beta-actin				
NM_012789	adenergic receptor alpha 2A					U23369	ILM domain protein, CLP36				
J05122	benzodiazepine receptor, peripheral-type					X81198	ILM, muscle				
A304018	histamine receptor					A3040289	perlecan II				
M69418	cholecystokinin type B receptor					NM_006086	tubulin, beta 4				
NM_012959	CDNF receptor alpha					A3011639	class I beta-tubulin				
S77867	G protein-coupled receptor, UHR-1					<b>H Metabolism</b>					
L08494	GABA receptor alpha 5 subunit					U79118	ATPase				
X13722	low-density lipoprotein receptor precursor					D10041	long-chain fatty acid-CoA ligase				
A3004057	neuropeptide Y Y2 receptor					L05175	serine protease				
U86714	neuropeptide Y Y5 receptor					D34694	5G2 protease regulatory subunit 7				
L31612	nicotinic acetylcholine receptor alpha 7 subunit					NM_012777	acyl-protein D				
U22850	P2Y1 purinoceptor					S54526	brain keratinase				
<b>C Channel</b>						M34477	testis-specific farnesyl pyrophosphate synthetase				
T16002	potassium channel PCK4 subunit					D80215	NADH dehydrogenase oxidoreductase				
X83580	muscle potassium channel subunit 11					S52752	acyl-CoA oxidase				
X50184	sodium channel (SCN5A)					J04488	brain prostaglandin D synthetase				
M69601	L-type calcium channel alpha 2/delta 1 subunit					<b>I Protein modulation and protein synthesis</b>					
U157026	sodium channel beta 2 subunit					D29683	endothelin converting enzyme				
Y00766	sodium channel III					U67011	met cell protease 8				
<b>D Signal transduction modulator and effector</b>						M63247	alpha 1-antitrypsinase				
S55905	14-3-3 protein, gamma-subtype					M68870	endoplasmic reticulum stress protein				
L12380	ADP-ribosylation factor 1					M68589	heat shock 70 kDa protein				
BC066507	Abi-interactor 1					X83309	translation initiation factor eIF-4E				
BC066590	BAP31					D26307	Jun D, c-jun-related transcription factor				
X94351	Clc5 protein kinase					L12458	tyrosine				
U23651	6-phosphofructokinase muscle isozyme					X06148	ribosomal protein L5				
D10666	neural vesicular-like calcium-binding protein					M13422	ribosomal protein L7				
D34455	phospholipase C delta 4					M18655	large subunit ribosomal protein L36a				
X07286	protein kinase C alpha type					<b>J Others</b>					
X06889	Pal-S, Ras-related protein					Y17823	CEK109				
BC066438	mitotic activator protein regulator 1					A3131820	CtH-42 protein				
L13420	aktin-dependent protein kinase type II delta					L22191	glutamate-cysteine ligase in placental subunit				
BC066486	endothelial monocyte-activating polypeptide I					X54862	C-6-methyltransferase-DNA methyltransferase				
X52711	interferon-induced GTP-binding protein, rat					M21750	lipocortin V				
X58931	protein tyrosine kinase					D50093	protein protein				
L27843	nuclear tyrosine phosphatase PPL-1					X50267	clathrin A				
U67309	neuronal nitric oxide synthase					L34067	glycogen				
S49400	protein tyrosine phosphatase					M69056	monocyte/erythrophil chemokine inhibitor				
BC0664549	mb1, ras-related GTPase					Y00169	Th1-4 gene for fibroblast tropomyosin 4				
M67679	mb15, ras-related GTPase					NM_019004	beta-galactonide-binding lectin				
D38222	tyrosine phosphatase-like protein					M17085	major alpha-globin				
<b>E Proteins related to synaptic transmission</b>						M19667	lipocortin I				
X06832	chromogranin A					M20559	lipocortin III				
D32249	neurite generation associated protein 1					Y00480	class II MHC alpha chain, F11.D				
L11962	synaptic vesicle protein 2E					A3088934	class II MHC RT1.D3 beta chain				
M24104	muscle associated neurofilament protein 1 (TAMF-1)					X58899	microvascular endothelial differentiation gene 1				
A3003991	synaptosomal-associated protein 25 kDa (SNAP-25)					NM_006054	perlecan 3				
L38247	synaptophysin IV					U88282	histamine protein component 1				
<b>cDNA array ratio of nVol of axotomized DRG/control DRG</b>											
		10-5		5-2		2-0.5		0.5-0.2		<0.2	

## Growth-associated proteins

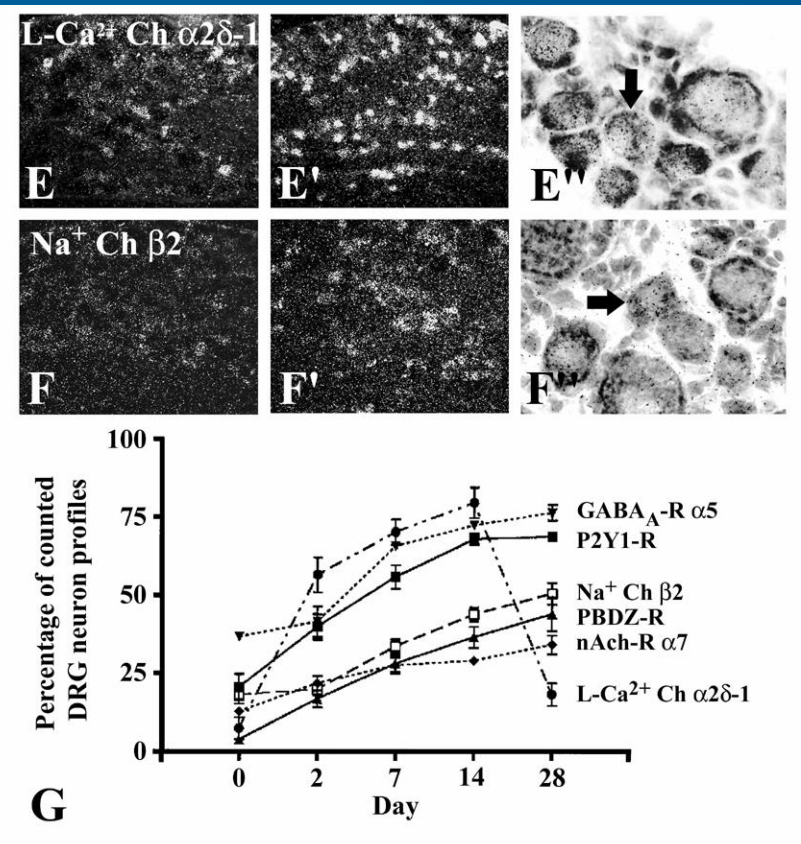
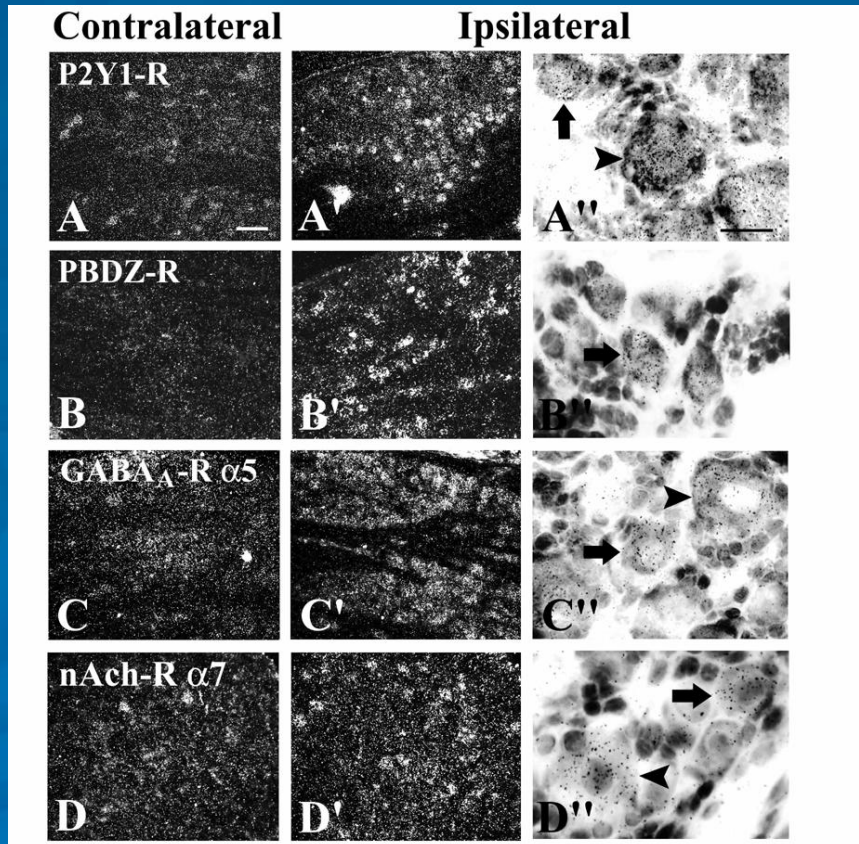
## Cytoskeleton and cell motility

## Metabolism

## Protein modulation and protein synthesis

## Others





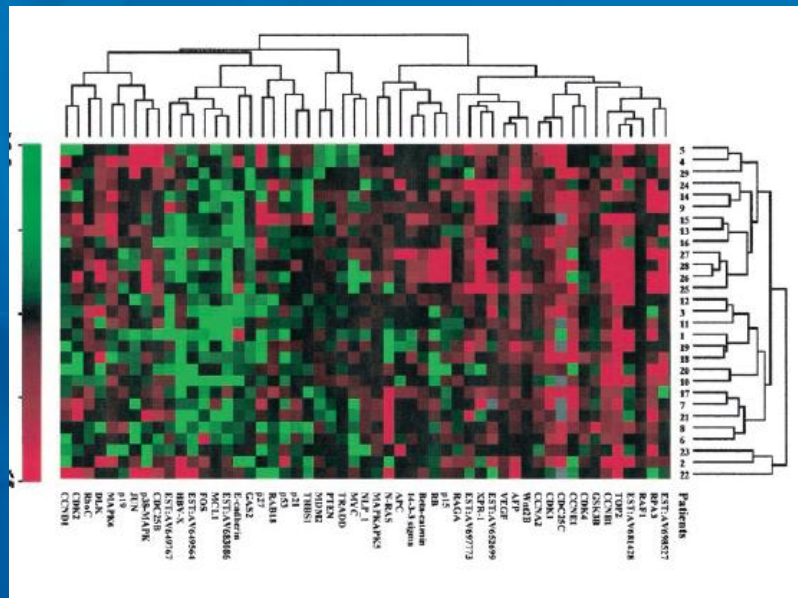
**Unexpectedly several antidepressant, antianxiety and anticonvulsant drug target molecules were up regulated in sensory neuron after axotomy. These provided the candidate genes for developing new drugs and therapeutic methods to treat neuropathic pain with minimal side effects.**



# Insight into hepatocellular carcinogenesis at transcriptome level by comparing gene expression profiles of hepatocellular carcinoma with those of corresponding noncancerous liver

Xiang-Ru Xu<sup>\*†</sup>, Jian Huang<sup>\*</sup>, Zhi-Gang Xu<sup>\*</sup>, Bin-Zhi Qian<sup>\*</sup>, Zhi-Dong Zhu<sup>\*</sup>, Qing Yan<sup>\*</sup>, Ting Cai<sup>\*</sup>, Xin Zhang<sup>\*</sup>, Hua-Sheng Xiao<sup>\*</sup>, Jian Qu<sup>\*</sup>, Feng Liu<sup>\*</sup>, Qiu-Hua Huang<sup>\*</sup>, Zhi-Hong Cheng<sup>\*</sup>, Neng-Gan Li<sup>\*</sup>, Jian-Jun Du<sup>\*</sup>, Wei Hu<sup>\*</sup>, Kun-Tang Shen<sup>\*</sup>, Gang Lu<sup>\*</sup>, Gang Fu<sup>\*</sup>, Ming Zhong<sup>\*</sup>, Shu-Hua Xu<sup>\*</sup>, Wen-Yi Gu<sup>\*</sup>, Wei Huang<sup>\*</sup>, Xin-Tai Zhao<sup>‡</sup>, Geng-Xi Hu<sup>§</sup>, Jian-Ren Gu<sup>‡</sup>, Zhu Chen<sup>\*†¶</sup>, and Ze-Guang Han<sup>\*¶</sup>

<sup>\*</sup>Chinese National Human Genome Center at Shanghai, 351 Guo Shou-Jing Road, Shanghai 201203, China; <sup>†</sup>Shanghai Institute of Hematology, Rui Jin Hospital, 197 Rui Jin Road II, Shanghai 200025, China; <sup>‡</sup>National Laboratory for Oncogenes and Related Genes, Shanghai Cancer Institute, 25 Lane 2200, Xietu Road, Shanghai 200032, China; <sup>§</sup>Shanghai Institute of Biochemistry and Cell Biology, Shanghai Institutes for Biological Sciences, Chinese Academy of Sciences, 320 Yueyang Road, Shanghai 200031, China; and <sup>¶</sup>Morgan-Tan International Center for Life Science, Institute of Genetics, Fudan University, 220 Han Dan Road, Shanghai 200433, China



acute myeloid leukemia

acute lymphoblastic leukemia

## Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring

T. R. Golub,<sup>1,2\*</sup> D. K. Slonim,<sup>1†</sup> P. Tamayo,<sup>1</sup> C. Huard,<sup>1</sup>  
M. Gaasenbeek,<sup>1</sup> J. P. Mesirov,<sup>1</sup> H. Coller,<sup>1</sup> M. L. Loh,<sup>2</sup>  
J. R. Downing,<sup>3</sup> M. A. Caligiuri,<sup>4</sup> C. D. Bloomfield,<sup>4</sup>  
E. S. Lander<sup>1,5\*</sup>

Science, 1999, 286:531-537

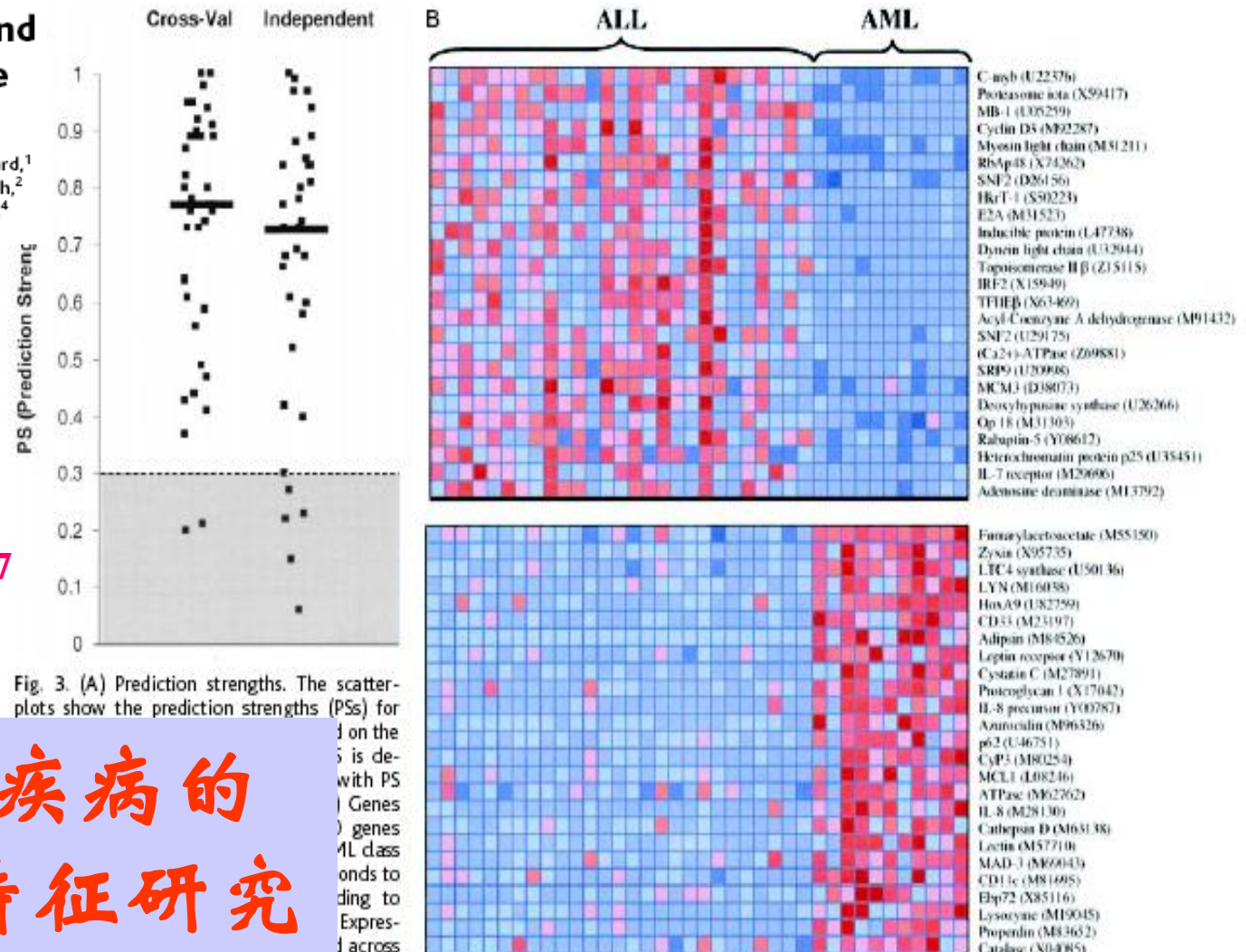


Fig. 3. (A) Prediction strengths. The scatter-plots show the prediction strengths (PSs) for genes on the y-axis. The top panel shows the PSs for genes in the ALL class, and the bottom panel shows the PSs for genes in the AML class. (B) Heatmaps of normalized expression levels for ALL and AML samples. The color scale indicates normalized expression from -3 (low) to 3 (high).

不同类型疾病的  
基因表达特征研究

the samples such that the mean is 0 and the SD is 1. Expression levels greater than the mean are shaded in red, and those below the mean are shaded in blue. The scale indicates SDs above or below the mean. The top panel shows



The New England  
Journal of Medicine

NUMBER 25



From the Divisions of Diagnostic Oncology (M.I.V., L.I.V., DWV., LL.P., D.v., A.M.V., A.C.D.), Radiotherapy (S.K.), Biometrics (L.V.), Surgical Oncology (R.B.), Netherlands Cancer Biomedical Genomics Amsterdam (I. land, Wash. Y.D.H., H.D., G.J.S., C.F. print requests to Dr. Bernards at the I Netherlands Cancer Institute, Plesman Netherlands, or at r.bernards@nki.nl.

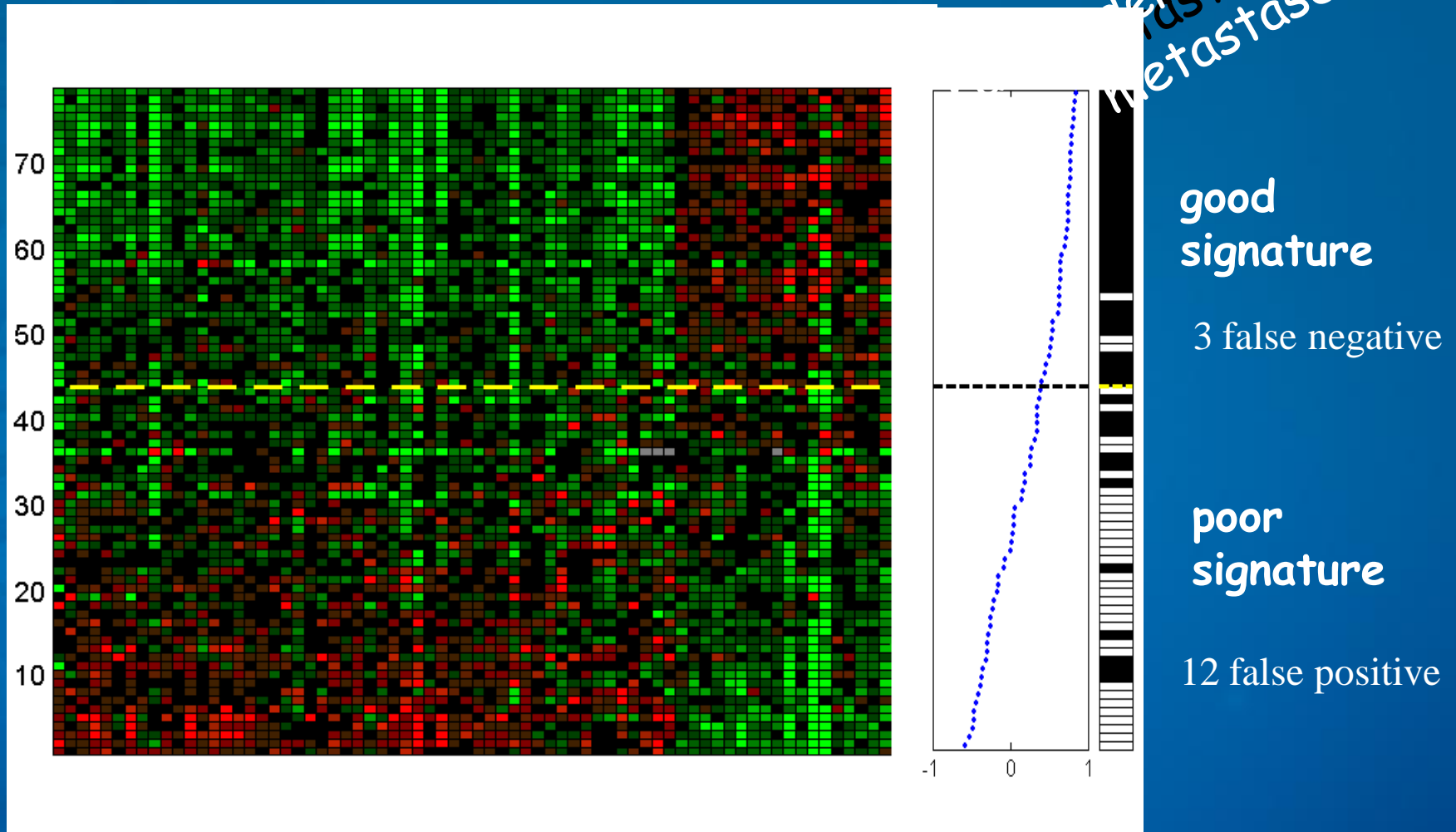


**Laura J. Van't Veer, PhD**  
COO and Co-founder, Agendia  
Head, Family Cancer Center  
Netherlands Cancer Institute



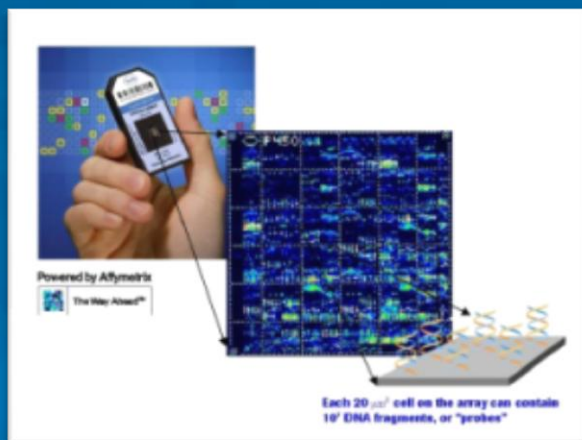
- Identifying gene expression signature that reflects biological behavior of a tumor
  - *70-gene prognostic profile tested*
  - *295 breast cancer samples*
  - *outperformed all clinical variables for predicting patient survival*
- Microarray classifiers to direct customized therapy
- Starting point for targeted and rational drug development

# Supervised Classification Prognosis

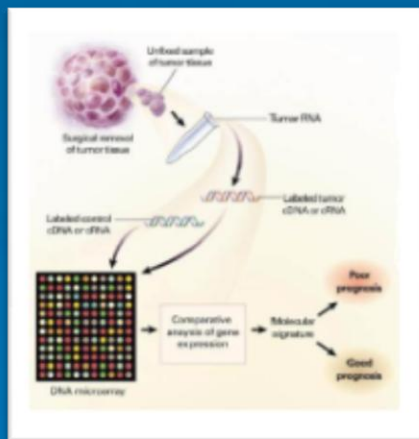


# 走向临床诊断的基因芯片

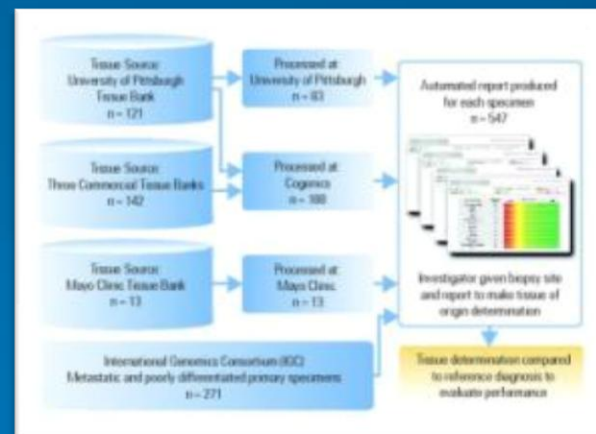
## 通过FDA或CE审核的诊断型芯片



2004年12月：罗氏公司 AmpliChip CYP450 Test 芯片诊断产品通过FDA审核。



2007年2月：荷兰Agendia公司乳腺癌预后诊断芯片MammaPrint通过FDA审核。



2010年：美国Pathwork公司的TOO芯片（Tissue of Origin Test）通过FDA审核。



2004年，Affymetrix专门用于体外诊断芯片的仪器系统——GeneChip System 3000Dx (GCS 3000Dx) 先后通过CE和FDA的审核。



# 国内已经批准的临床基因检测芯片

## 国家药品监督管理局注册产品

TORCH核酸检测试剂盒(基因芯片法)	葡萄糖6磷酸脱氢酶基因突变检测试剂盒(基因芯片法)
B-raf基因突变检测试剂盒（基因芯片法）	苯丙氨酸羟化酶基因突变检测试剂盒（基因芯片法）
IL28B基因多态性检测试剂盒（基因芯片法）	ALDH2 (Glu504Lys)基因检测试剂盒(DNA微阵列芯片法)
K-ras基因突变检测试剂盒（基因芯片法）	CYP2C19基因检测试剂盒（DNA微阵列芯片法）
CYP2C19基因多态性核酸检测试剂盒(基因芯片法)	分枝杆菌菌种鉴定试剂盒（DNA微阵列芯片法）
乙型肝炎病毒基因分型检测试剂盒(基因芯片法)	结核分枝杆菌耐药基因检测试剂盒（DNA微阵列芯片法）
载脂蛋白E（ApoE）基因型检测试剂盒（基因芯片法）	地中海贫血基因检测试剂盒（微阵列芯片法）
VKORC1和CYP2C9基因检测试剂盒（PCR-电化学基因芯片法）	九项遗传性耳聋基因检测试剂盒（微阵列芯片法）
人乳头瘤病毒分型检测试剂盒（基因芯片法）	人乳头瘤病毒（HPV）分型检测试剂盒（微阵列芯片法）
人乳头状瘤病毒基因分型检测试剂盒(基因芯片法)	十五项遗传性耳聋相关基因检测试剂盒（微阵列芯片法）
乙型肝炎病毒耐药突变位点检测试剂盒(基因芯片法)	抗栓治疗用药相关6个基因位点多态性检测试剂盒（微阵列芯片法）
CYP2D6*10、CYP2C9*3、ADRB1(1165G>C)、AGTR1(1166A>C)、ACE(I/D)检测试剂盒（基因芯片法）	乙型肝炎病毒（HBV）基因分型和耐药突变位点检测试剂盒（微阵列芯片法）
六项呼吸道病毒核酸检测试剂盒（恒温扩增芯片法）	呼吸道病原菌核酸检测试剂盒（恒温扩增芯片法）