

Introduction to genome biology and DNA microarray experiments

Sandrine Dudoit and Robert Gentleman

Statistics and Genomics - Lecture 1, Part I

Department of Biostatistics

Harvard School of Public Health

January 23-25, 2002

Outline of lecture 1

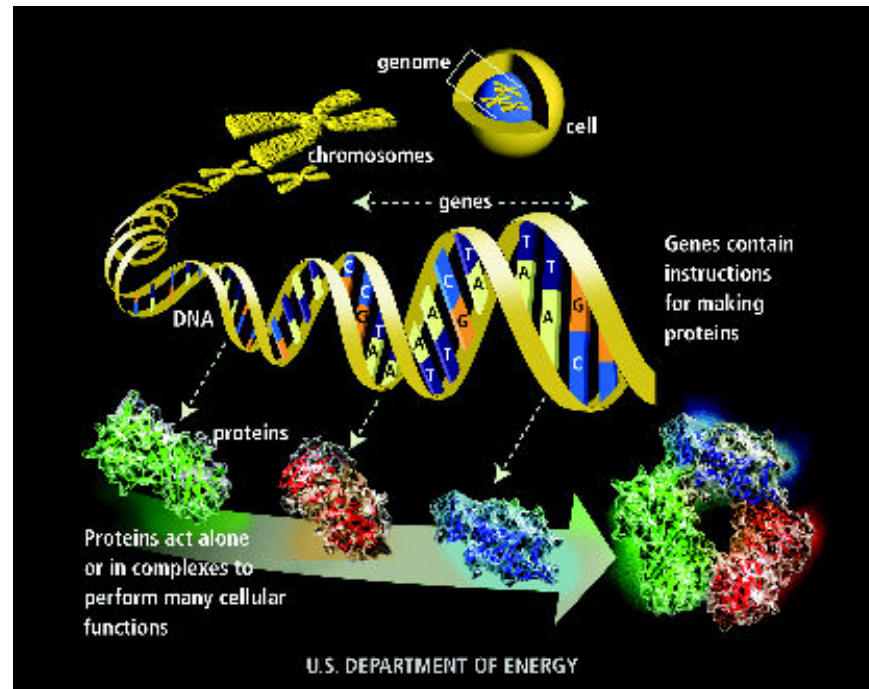
Part I:

- Introduction to genome biology;
- Introduction to microarray experiments.

Part II:

- Image analysis (cDNA microarrays);
- Normalization (cDNA microarrays);
- Experimental design.

Introduction to genome biology



The human genome

- The **cell** is the fundamental working unit of every living organism.
- Humans: trillions of cells (metazoa); other organisms like yeast: one cell (protozoa).
- Cells are of many different types (e.g. blood, skin, nerve cells), but all can be traced back to a single cell, the fertilized egg.

The human genome

- The **genome**, or blueprint for all cellular structures and activities in our body, is encoded in **DNA** molecules.
- Each cell contains a complete copy of the organism's **genome**.

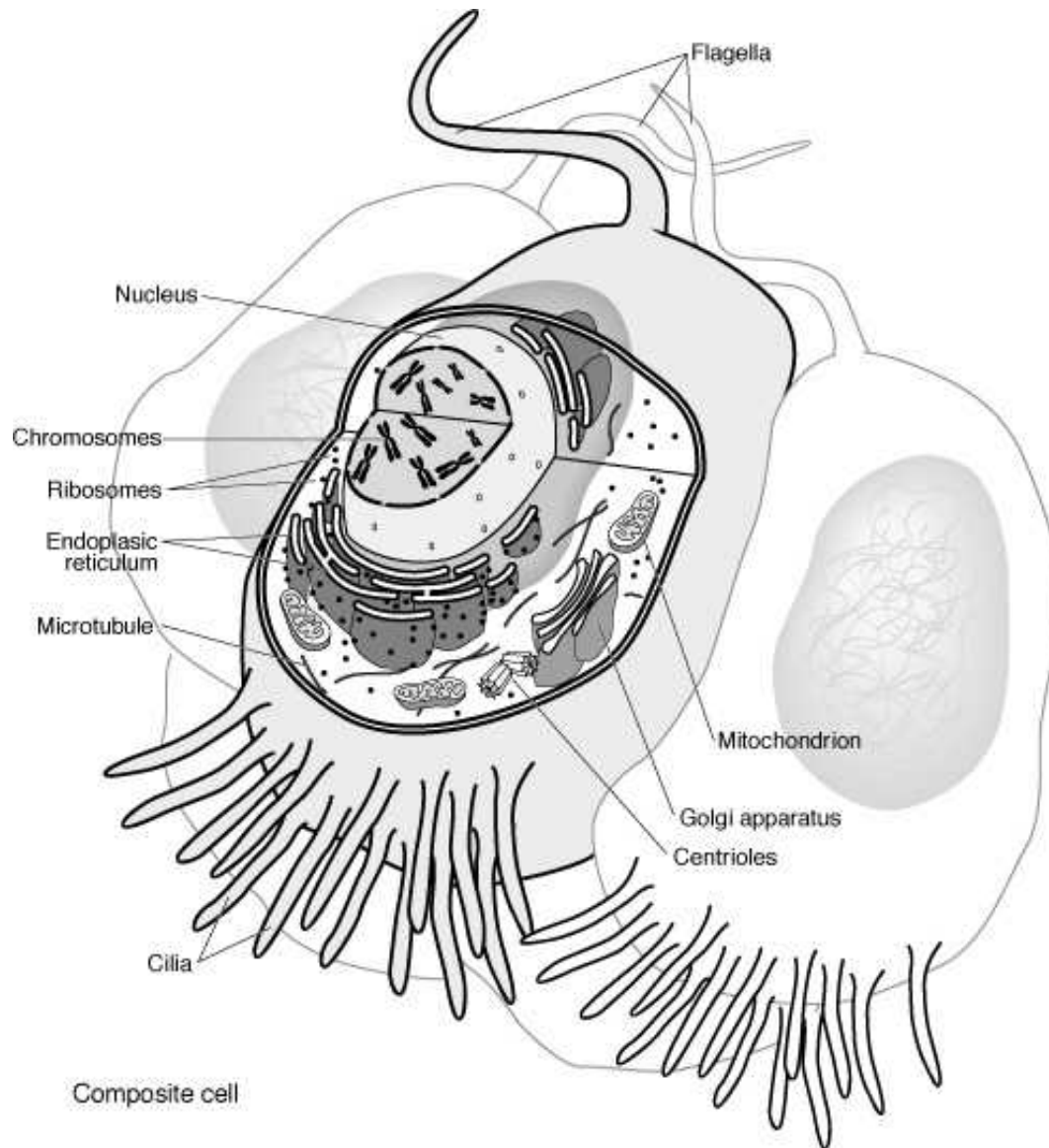
The human genome

- The human genome is distributed along 23 pairs of **chromosomes**
 - 22 autosomal pairs;
 - the sex chromosome pair, XX for females and XY for males.
- In each pair, one chromosome is paternally inherited, the other maternally inherited (cf. meiosis).

The human genome

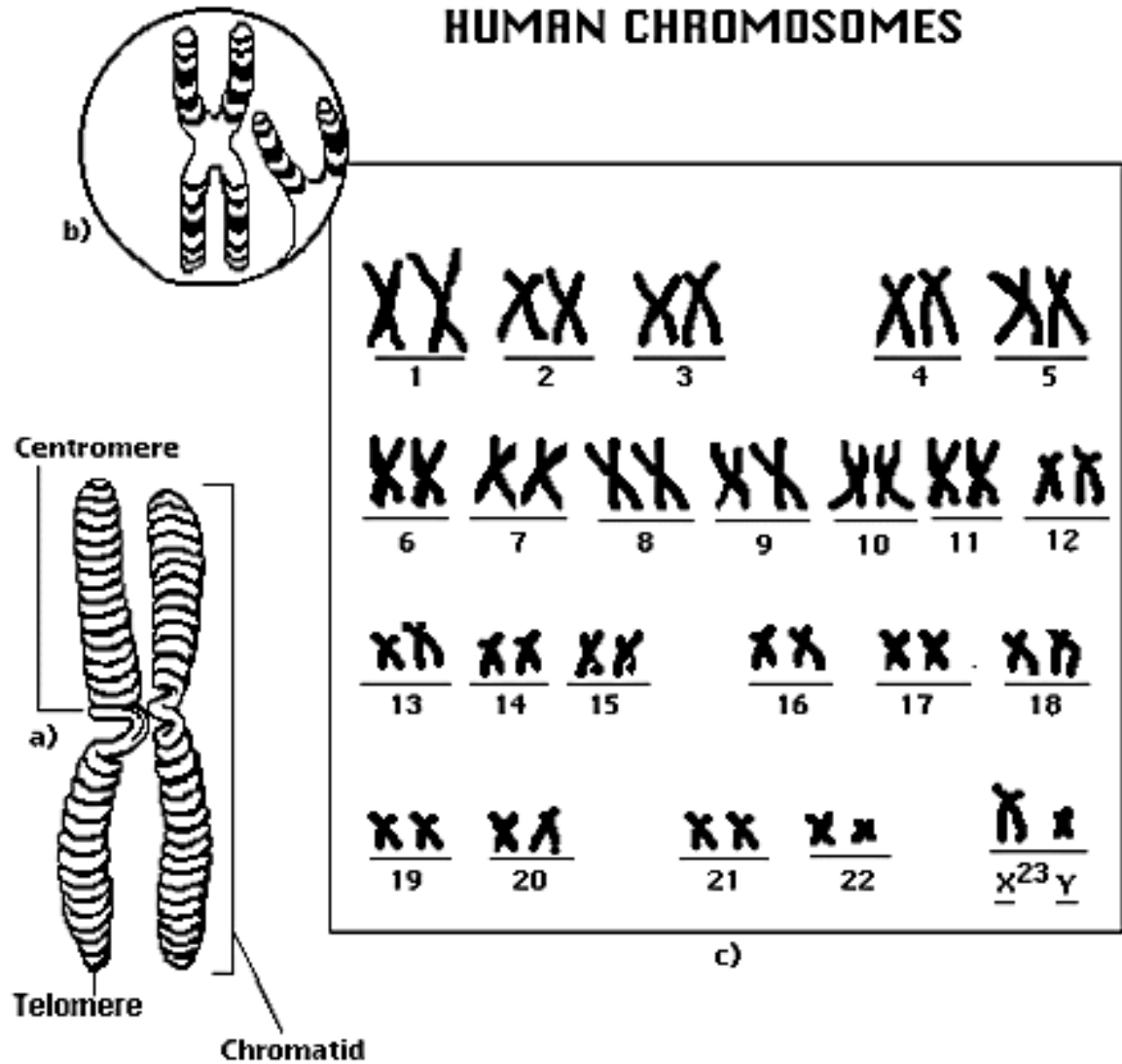
- Chromosomes are made of compressed and entwined **DNA**.
- A (protein-coding) **gene** is a segment of chromosomal **DNA** that directs the synthesis of a **protein**.

The eukaryotic cell

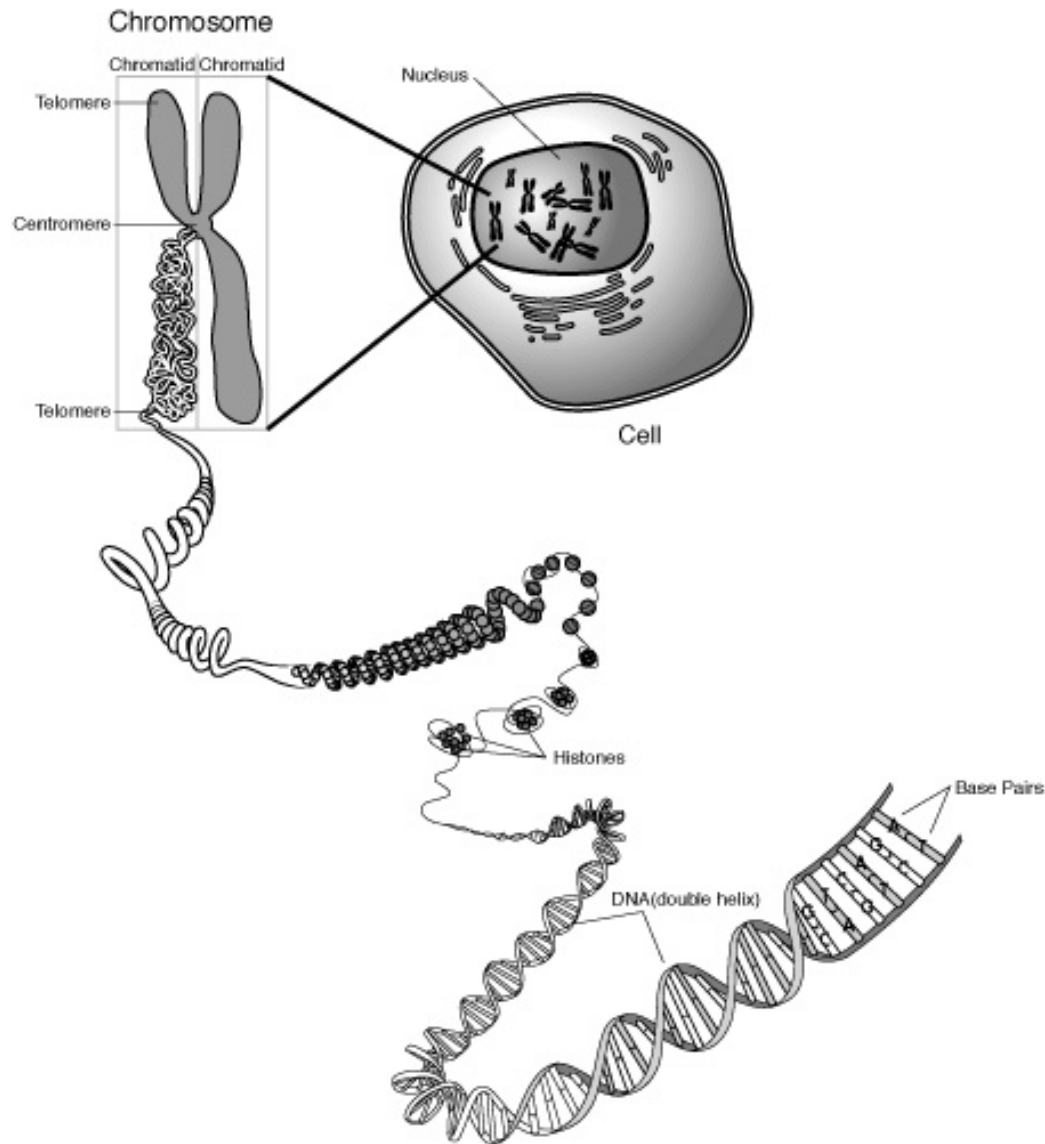


Chromosomes

HUMAN CHROMOSOMES



Chromosomes and DNA

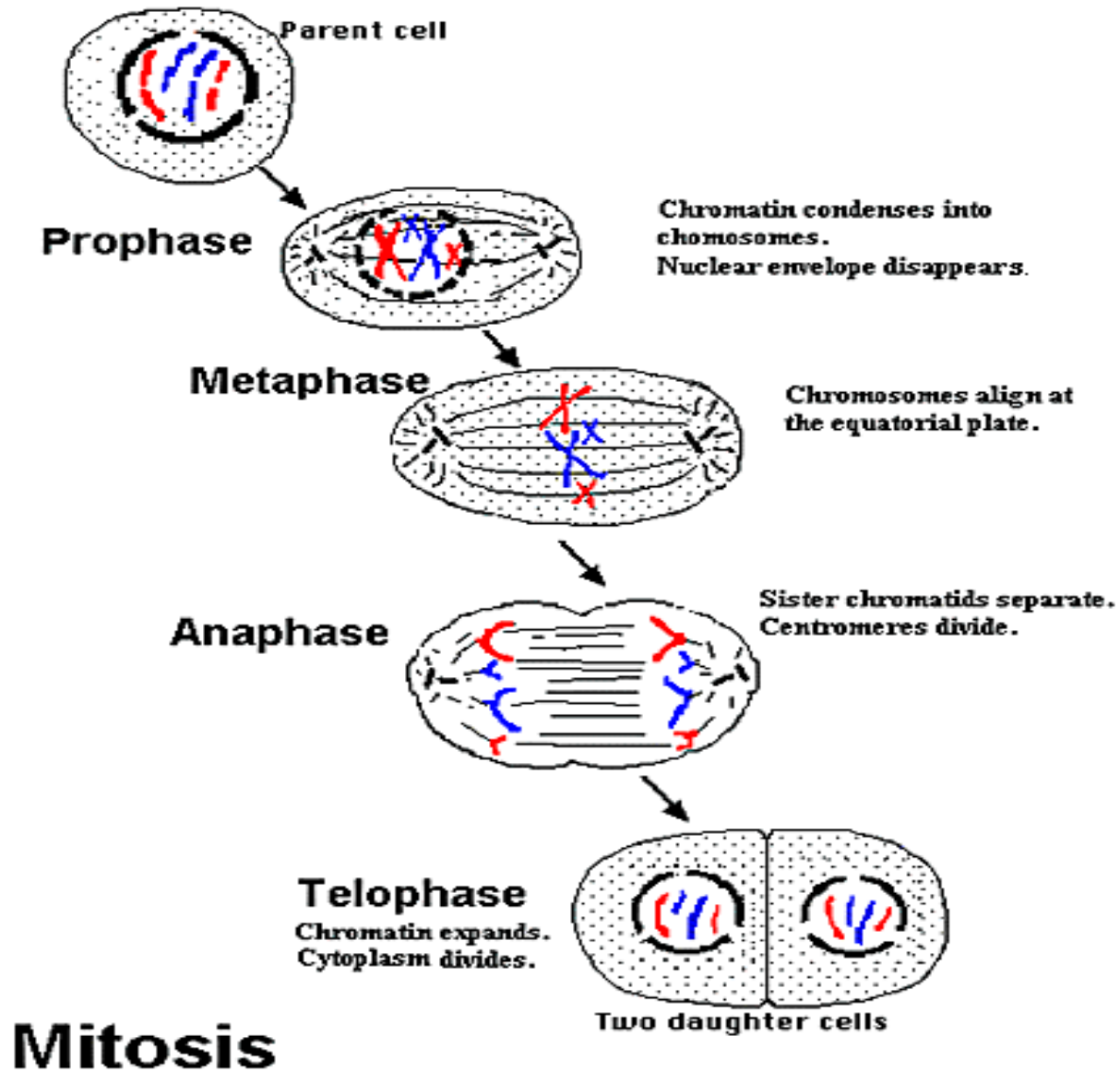


Cell divisions

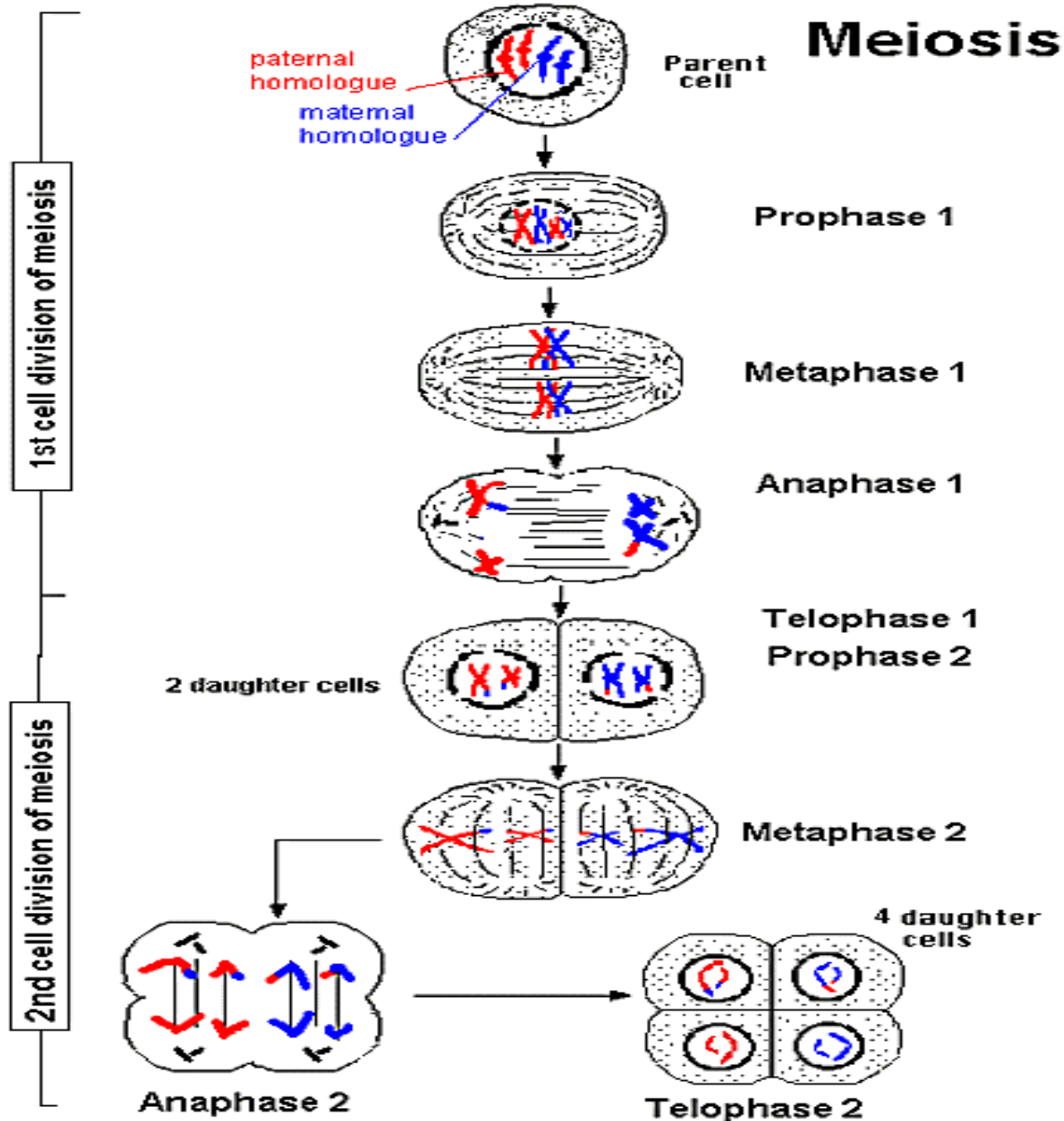
- **Mitosis.** One nuclear division produces two daughter **diploid** nuclei identical to the parent nucleus.
- **Meiosis.** Two successive nuclear divisions produces four daughter **haploid** nuclei, different from original cell.

Leads to the formation of gametes (egg/sperm).

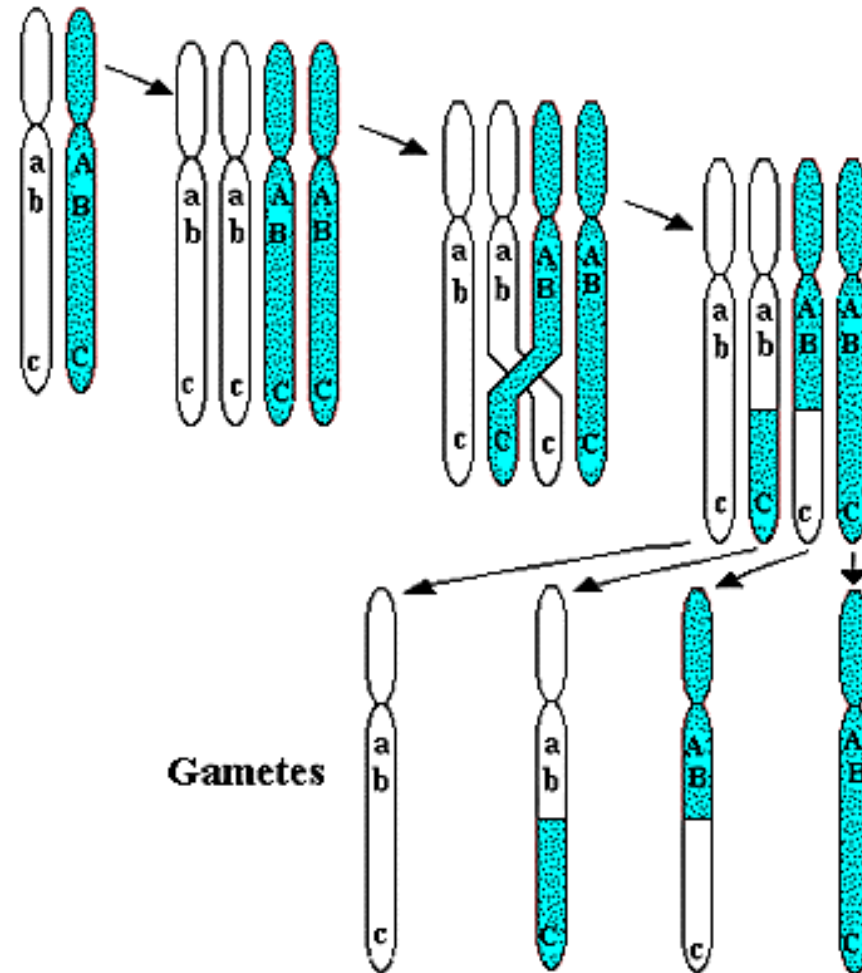
Mitosis



Meiosis



Recombination

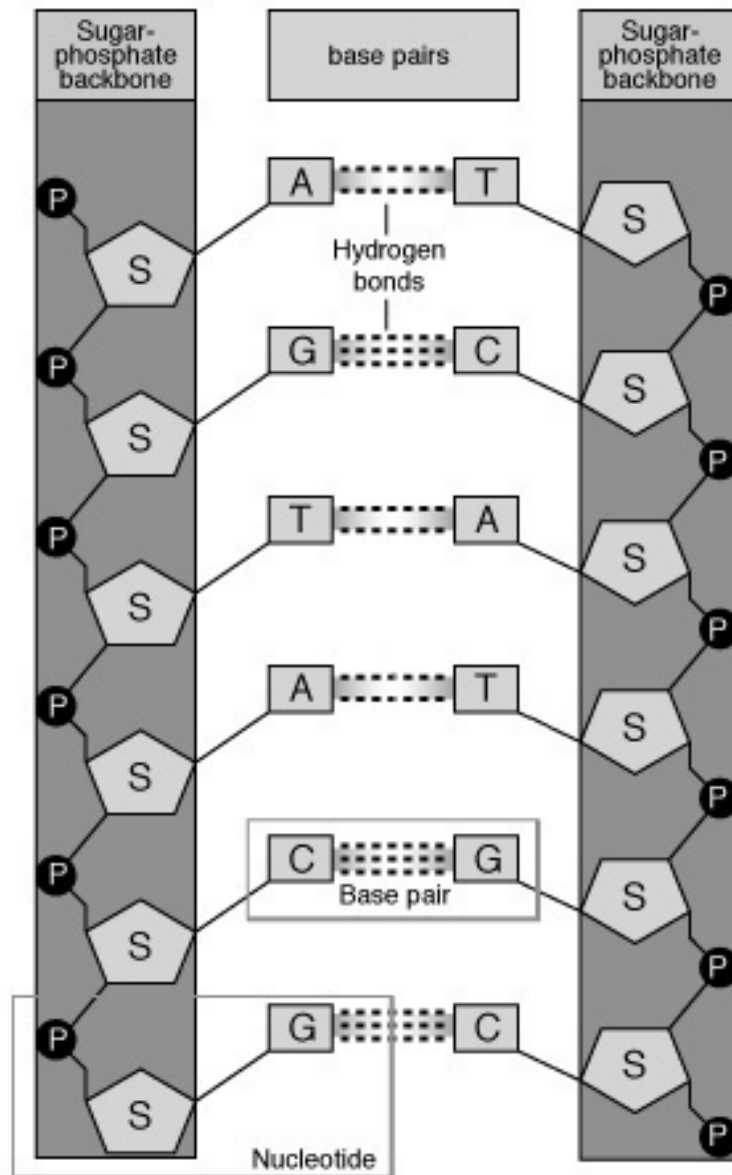


Crossing-over and recombination during meiosis

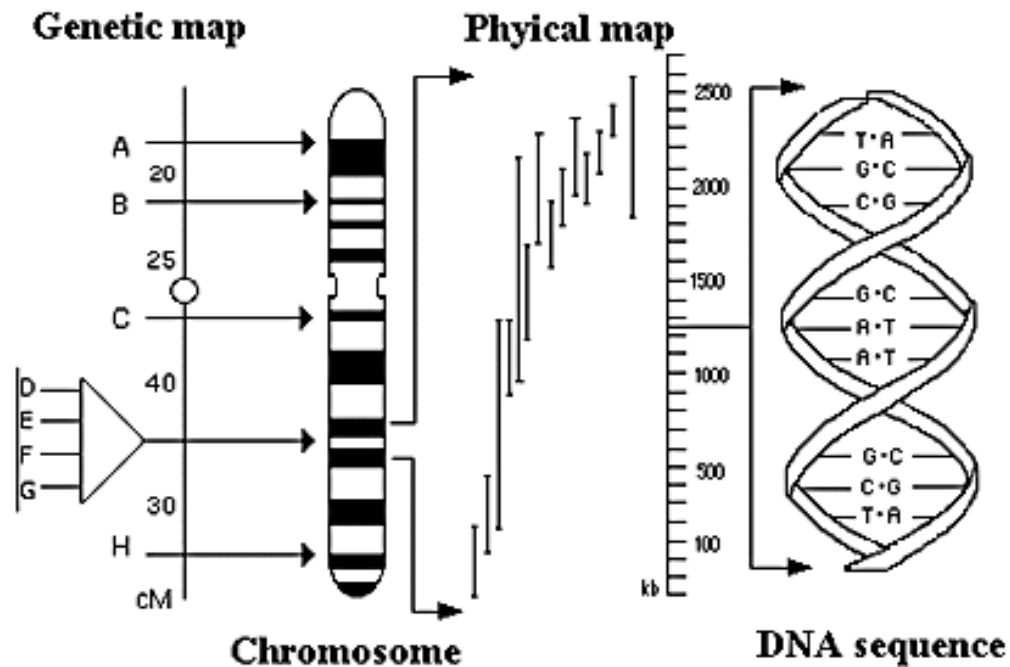
DNA

- A **deoxyribonucleic acid** or **DNA** molecule is a double-stranded polymer composed of four basic molecular units called **nucleotides**.
- Each nucleotide comprises a phosphate group, a deoxyribose sugar, and one of four nitrogen bases: **adenine (A)**, **guanine (G)**, **cytosine (C)**, and **thymine (T)**.
- The two chains are held together by hydrogen bonds between nitrogen bases.
- Base-pairing occurs according to the following rule: **G pairs with C**, and **A pairs with T**.

DNA



Genetic and physical maps



Sequences of base pairs mapping

Genetic and physical maps

- **Physical distance:** number of base pairs (bp).
- **Genetic distance:** expected number of crossovers between two loci, per chromatid, per meiosis.
Measured in Morgans (M) or centiMorgans (cM).
- 1cM \sim 1 million bp (1Mb).

The human genome in numbers

- 23 pairs of chromosomes;
- 2 meters of DNA;
- 3,000,000,000 bp;
- 35 M (males 27M, females 44M);
- 30,000-40,000 genes.

Proteins

- **Proteins:** large molecules composed of one or more chains of amino acids.
- **Amino acids:** class of 20 different organic compounds containing a basic amino group ($-\text{NH}_2$) and an acidic carboxyl group ($-\text{COOH}$).
- The order of the amino acids is determined by the **base sequence** of nucleotides in the **gene** coding for the protein.
- E.g. hormones, enzymes, antibodies.

Amino acids

FAMILIES OF AMINO ACIDS

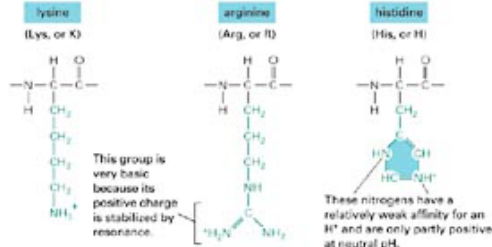
The common amino acids are grouped according to whether their side chains are:

- acidic
- basic
- uncharged polar
- nonpolar

These 20 amino acids are given both three-letter and one-letter abbreviations.

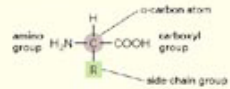
Thus: alanine = Ala = A

BASIC SIDE CHAINS

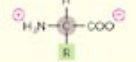


THE AMINO ACID

The general formula of an amino acid is

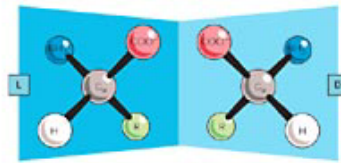


R is commonly one of 20 different side chains. At pH 7 both the amino and carboxyl groups are ionized.



OPTICAL ISOMERS

The α-carbon atom is asymmetric, which allows for two mirror image (or stereo) isomers, L and D.

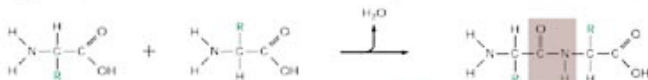


Proteins consist exclusively of L-amino acids.

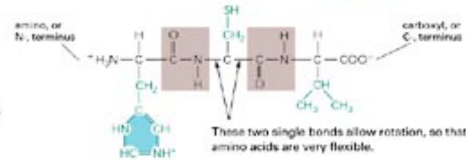
PEPTIDE BONDS

Amino acids are commonly joined together by an amide linkage, called a peptide bond.

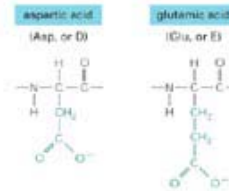
Peptide bond: The four atoms in each gray box form a rigid planar unit. There is no rotation around the C-N bond.



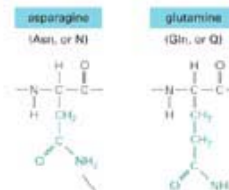
Proteins are long polymers of amino acids linked by peptide bonds, and they are always written with the N-terminus toward the left. The sequence of this tripeptide is histidine-cysteine-valine.



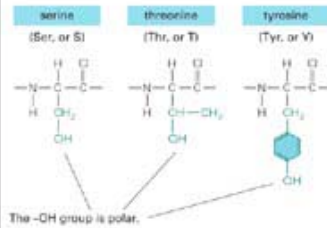
ACIDIC SIDE CHAINS



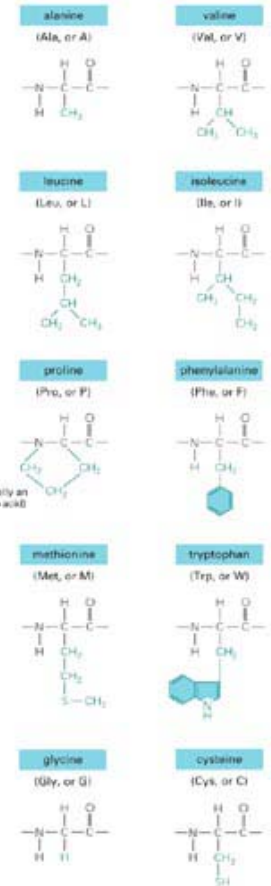
UNCHARGED POLAR SIDE CHAINS



Although the amide N is not charged at neutral pH, it is polar.



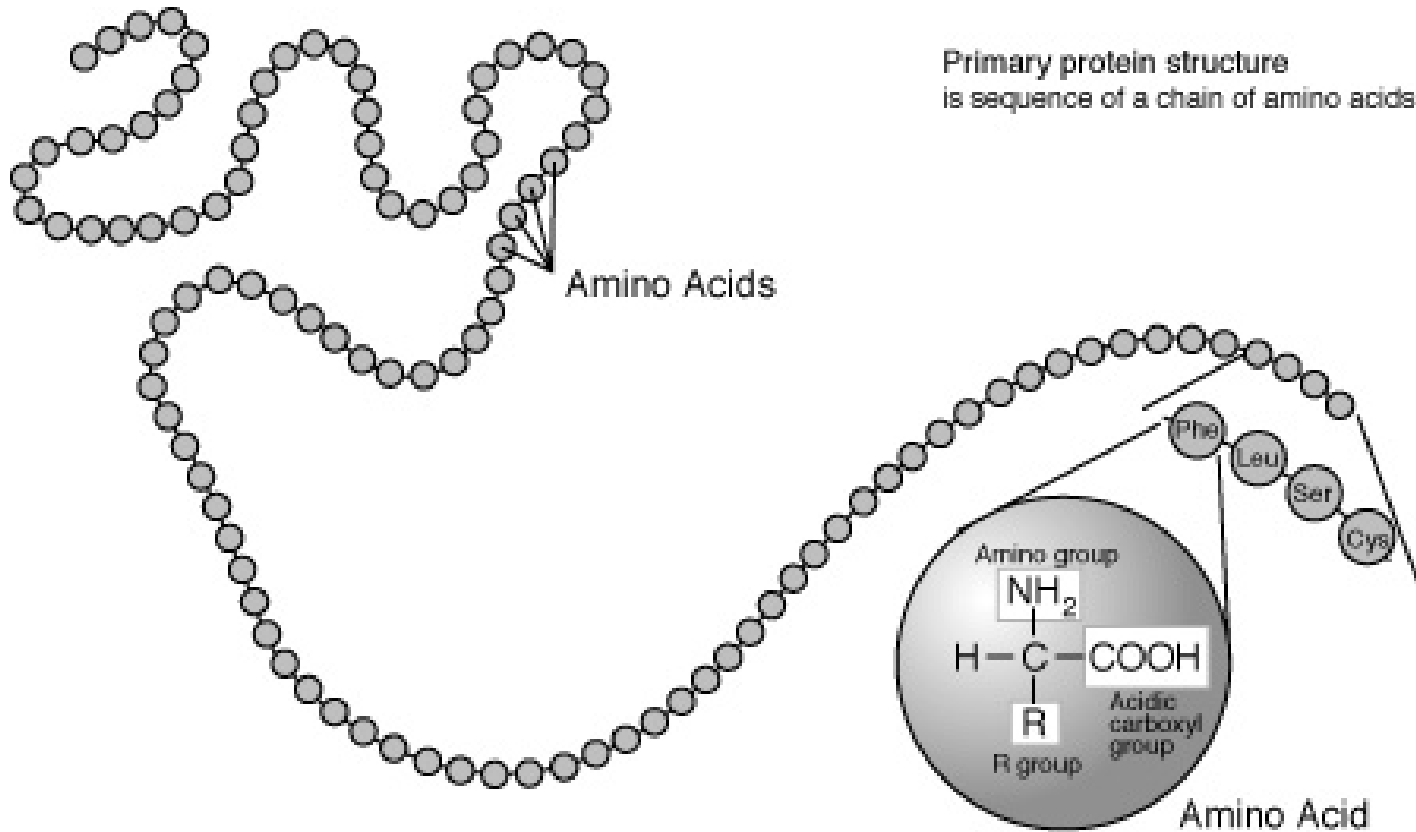
NONPOLAR SIDE CHAINS



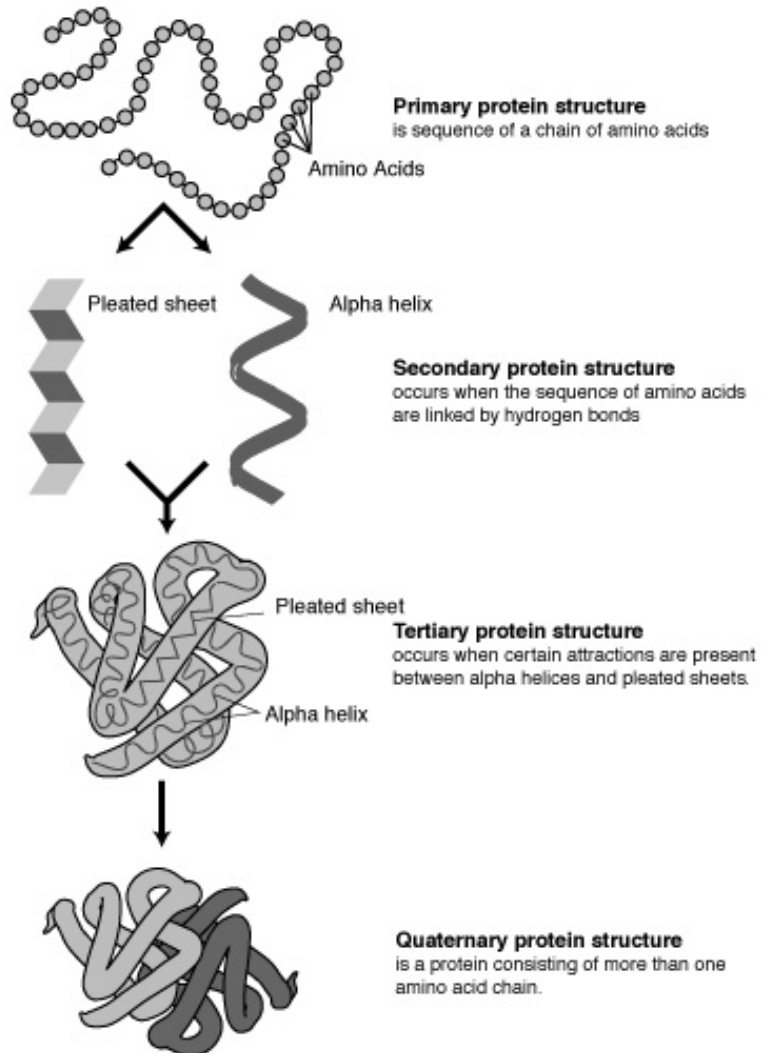
Disulfide bonds can form between two cysteine side chains in proteins.



Proteins



Proteins



Cell types

CELL TYPES

There are over 200 types of cells in the human body. These are assembled into a variety of types of tissue such as

- epithelia
- connective tissue
- muscle
- nervous tissue

Most tissues contain a mixture of cell types.

EPITHELIA

Epithelial cells form coherent cell sheets called epithelia, which line the inner and outer surfaces of the body. There are many specialized types of epithelia.

Absorptive cells have numerous hairlike projections called microvilli on their free surface to increase the area for absorption.

Ciliated cells have cilia on their free surface that beat in synchrony to move substances such as mucus over the epithelial sheet.

Secretory cells are found in most epithelial layers. These specialized cells secrete substances onto the surface of the cell sheet.

Adjacent epithelial cells are bound together by cell junctions that give the sheet mechanical strength and also make it impermeable to small molecules. The sheet rests on a basal lamina.

CONNECTIVE TISSUE

The spaces between organs and tissues in the body are filled with connective tissue made principally of a network of tough protein fibers embedded in a polysaccharide gel. This **extracellular matrix** is secreted mainly by **fibroblasts**.

Two main types of extracellular protein fibers are **collagen** and **actin**.

fibroblasts in loose connective tissue

bone is made by cells called **osteoblasts**. These secrete an extracellular matrix in which crystals of calcium phosphate are later deposited.

Cells are held together by cell processes.

adipose cells among the largest cells in the body, are responsible for the production and storage of fat. The nucleus and cytoplasm are squeezed by a large lipid droplet.

NERVOUS TISSUE

Dendrites receive signals from other neurons or sensory receptors.

Cell body contains the nucleus and organelles.

Axon conducts electrical signals away from the cell body.

Myelin sheath is formed by specialized glial cells.

Synapse is where a neuron forms a specialized junction with another neuron or with a muscle cell.

SECRETORY EPITHELIAL CELLS

Secretory epithelial cells are often collected together to form a gland that specializes in the secretion of a particular substance. As illustrated, **exocrine glands** secrete their products (such as tears, mucus, and gastric juices) into ducts. **Endocrine glands** secrete hormones into the blood.

MUSCLE

Muscle cells produce mechanical force by their contraction. In vertebrates there are three main types:

skeletal muscle—this moves joints by its strong and rapid contraction. Each muscle is a bundle of muscle fibers, each of which is an enormous multinucleated cell.

smooth muscle—present in digestive tract, bladder, arteries, and veins. It is composed of thin elongated cells (not striated), each of which has one nucleus.

cardiac muscle—intermediate in character between skeletal and smooth muscle. It produces the heart beat. Adjacent cells are linked by electrically conducting junctions that cause the cells to contract in synchrony.

BLOOD

Erythrocytes (red blood cells) are very small cells, and in mammals have no nucleus or internal membranes. When mature they are stuffed full of the oxygen-binding protein hemoglobin.

Leucocytes (white blood cells) protect against infections. Blood contains about one leucocyte for every 100 red blood cells. Although leucocytes travel in the circulation, they can pass through the walls of blood vessels to do their work in the surrounding tissues. There are several different kinds, including:

- Lymphocytes**—responsible for immune responses such as the production of antibodies.
- Macrophages and neutrophils**—move to sites of infection, where they ingest bacteria and debris.

SENSORY CELLS

Among the most strikingly specialized cells in the vertebrate body are those that detect external stimuli. **Hair cells** of the inner ear are primary detectors of sound. They are modified epithelial cells that carry special microvilli (stereocilia) on their surface. The movement of these in response to sound vibrations causes an electrical signal to pass to the brain.

Rod cells in the retina of the eye are specialized to respond to light. The photosensitive region contains many membranous discs (ret) in whose membranes the light sensitive pigment rhodopsin is embedded. Light evokes an electrical signal (green arrow), which is transmitted to nerve cells in the eye, which relay the signal to the brain.

GERM CELLS

Both sperm and egg are haploid, that is, they carry only one set of chromosomes. A sperm from the male fuses with an egg from the female, which then forms a new diploid organism by successive cell divisions.

Differential expression

- Each cell contains a complete copy of the organism's genome.
- Cells are of many different types and states
E.g. blood, nerve, and skin cells, dividing cells, cancerous cells, etc.
- What makes the cells different?
- **Differential gene expression**, i.e., **when, where**, and in **what quantity** each gene is expressed.
- On average, 40% of our genes are expressed at any given time.

Central dogma

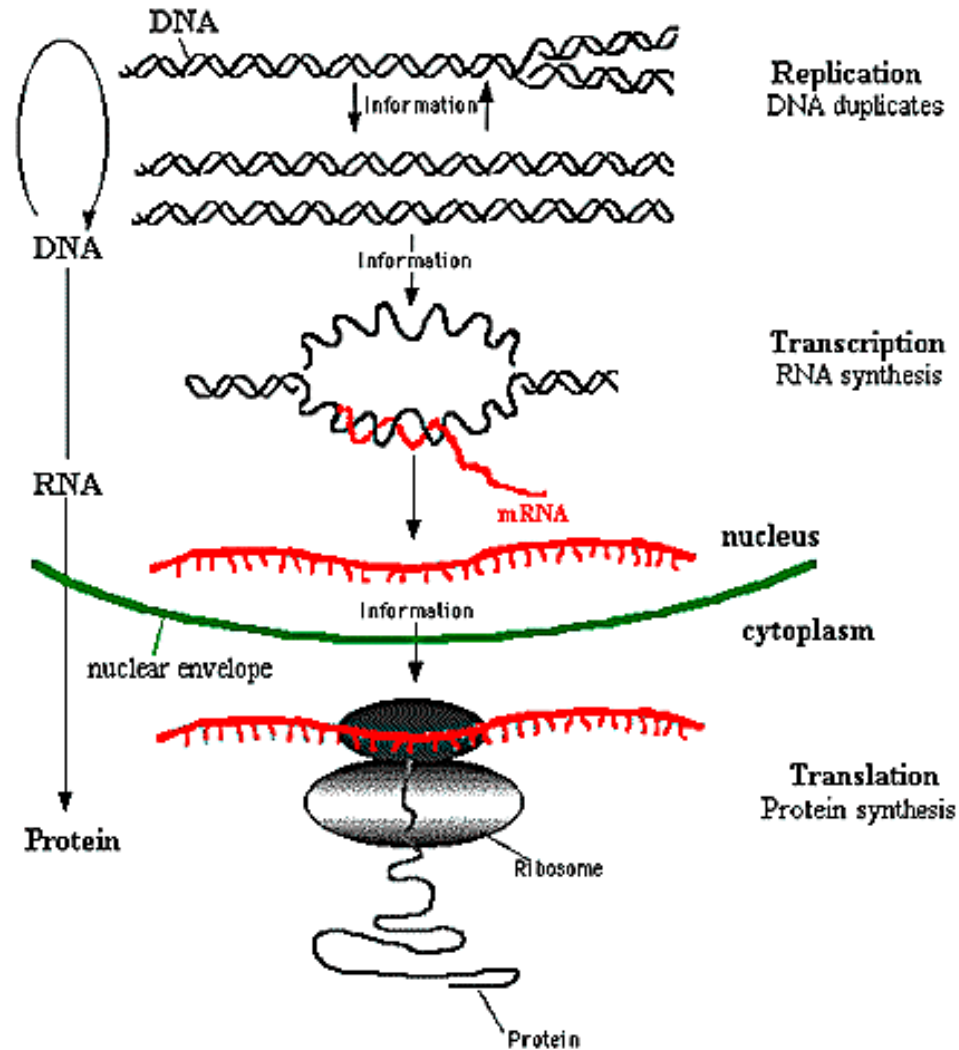
The **expression** of the genetic information stored in the DNA molecule occurs in two stages:

- (i) **transcription**, during which DNA is transcribed into mRNA;
- (ii) **translation**, during which mRNA is translated to produce a protein.

DNA → mRNA → protein

Other important aspects of regulation: methylation, alternative splicing, etc.

Central dogma



The Central Dogma of Molecular Biology

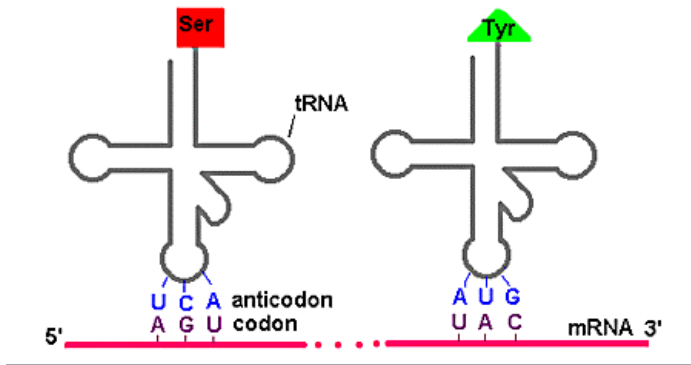
RNA

- A **ribonucleic acid** or **RNA** molecule is a nucleic acid similar to DNA, but
 - single-stranded;
 - ribose sugar rather than deoxyribose sugar;
 - **uracil (U)** replaces thymine (T) as one of the bases.
- RNA plays an important role in protein synthesis and other chemical activities of the cell.
- Several classes of RNA molecules, including **messenger RNA (mRNA)**, transfer RNA (tRNA), ribosomal RNA (rRNA), and other small RNAs.

The genetic code

- **DNA:** sequence of **four** different nucleotides.
- **Proteins:** sequence of **twenty** different amino acids.
- The correspondence between DNA's four-letter alphabet and a protein's twenty-letter alphabet is specified by the **genetic code**, which relates nucleotide triplets or **codons** to **amino acids**.

The genetic code



		2nd base in codon				
		U	C	A	G	
1st base in codon	U	Phe Phe Leu Leu	Ser Ser Ser Ser	Tyr Tyr STOP STOP	Cys Cys STOP Trp	U C A G
	C	Leu Leu Leu Leu	Pro Pro Pro Pro	His His Gln Gln	Arg Arg Arg Arg	U C A G
	A	Ile Ile Ile Met	Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	U C A G
	G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	U C A G
						3rd base in codon

The Genetic Code

Mapping between codons and amino acids is many-to-one: 64 codons but only 20 a.a..

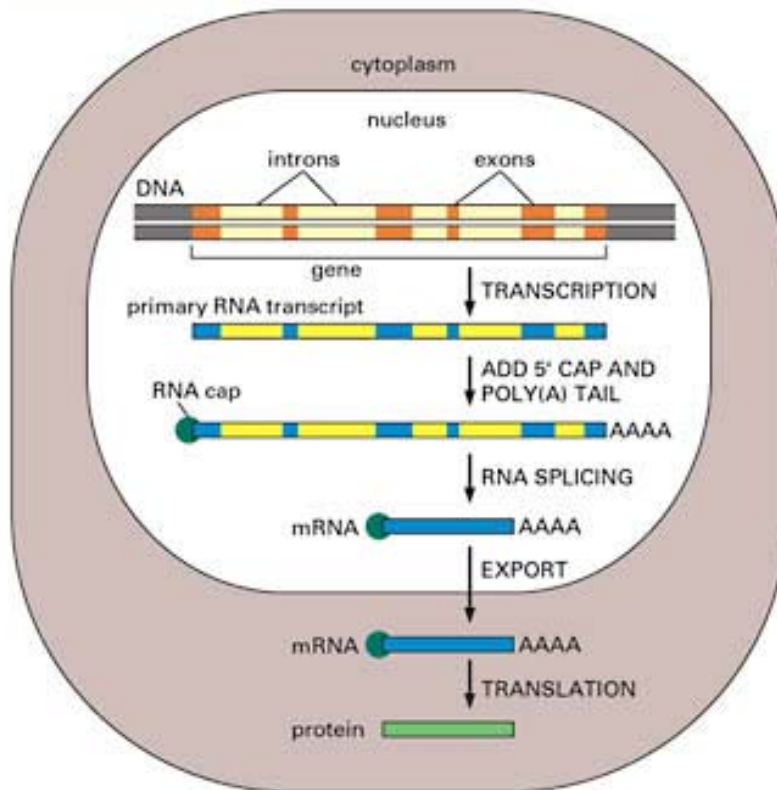
Third base in codon is often redundant,
e.g., stop codons.

Exons and introns

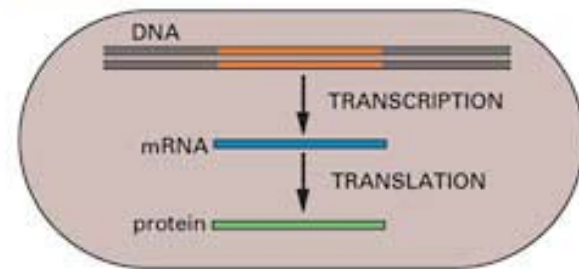
- Genes comprise only about 2% of the human genome; the rest consists of non-coding regions, whose functions may include providing chromosomal structural integrity and regulating when, where, and in what quantity proteins are made (regulatory regions).
- The terms **exon** and **intron** refer to coding (translated into a protein) and non-coding DNA, respectively.

Splicing

(A) EUCARYOTES



(B) PROCARYOTES



Alternative splicing

- There are more than 1,000,000 different human antibodies. How is this possible with only ~30,000 genes?
- **Alternative splicing** refers to the different ways of combining a gene's exons. This can produce different forms of a protein for the same gene,
- Alternative pre-mRNA splicing is an important mechanism for regulating gene expression in higher eukaryotes.
- E.g. in humans, it is estimated that approximately 30% genes are subject to alternative splicing.

Alternative splicing



Primary isoform



Cryptic exon



Exon extension
(5' or 3')



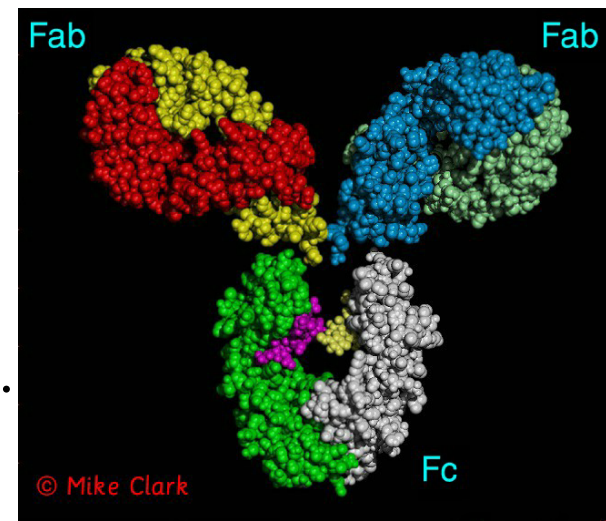
Exon skipping



Exon truncation

Immunoglobulin

- B cells produce antibody molecules called immunoglobulins (Ig) which fall in five broad classes.
- Diversity of Ig molecules
 - DNA sequence: recombination, mutation.
 - mRNA sequence: alternative splicing.
 - Protein structure: post-translational proteolysis, glycosylation.



IgG1

Functional genomics

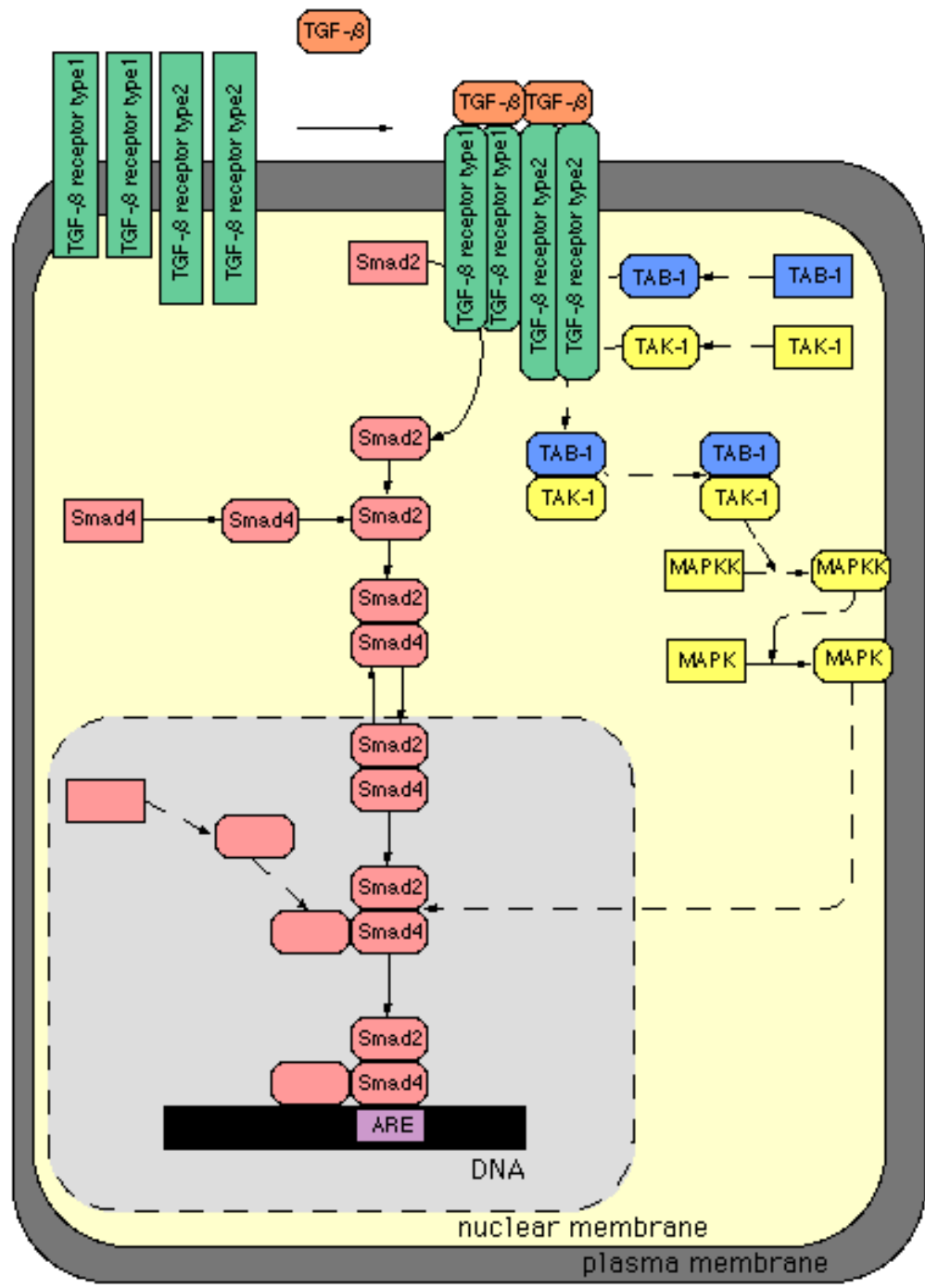
- The various **genome projects** have yielded the complete DNA sequences of many organisms.
 - E.g. human, mouse, yeast, fruitfly, etc.
 - Human: 3 billion base-pairs, 30-40 thousand genes.
- Challenge: **go from sequence to function**, i.e., define the role of each gene and understand how the genome functions as a whole.

Pathways

- The complete genome sequence doesn't tell us much about how the organism functions as a biological system.
- We need to study how different gene products function to produce various components.
- Most important activities are not the result of a single molecule but depend on the coordinated effects of multiple molecules.

TGF- β pathway

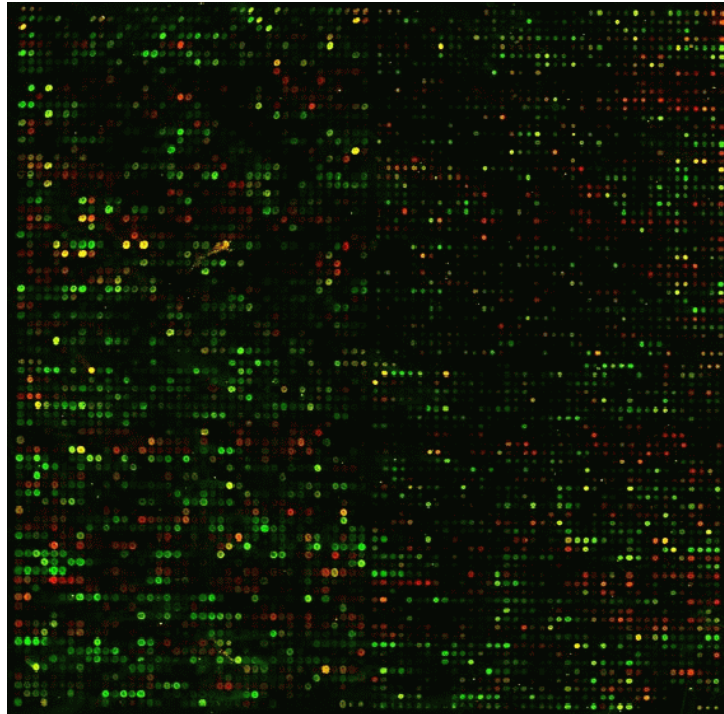
- TGF- β (transforming growth factor beta) plays an essential role in the control of development and morphogenesis in multicellular organisms.
- This is done through SMADS, a family of signal transducers and transcriptional activators.



Pathways

- <http://www.grt.kyushu-u.ac.jp/spad/>
- There are many open questions regarding the relationship between expression level and pathways.
- It is not clear whether expression level data will be informative.

DNA microarrays

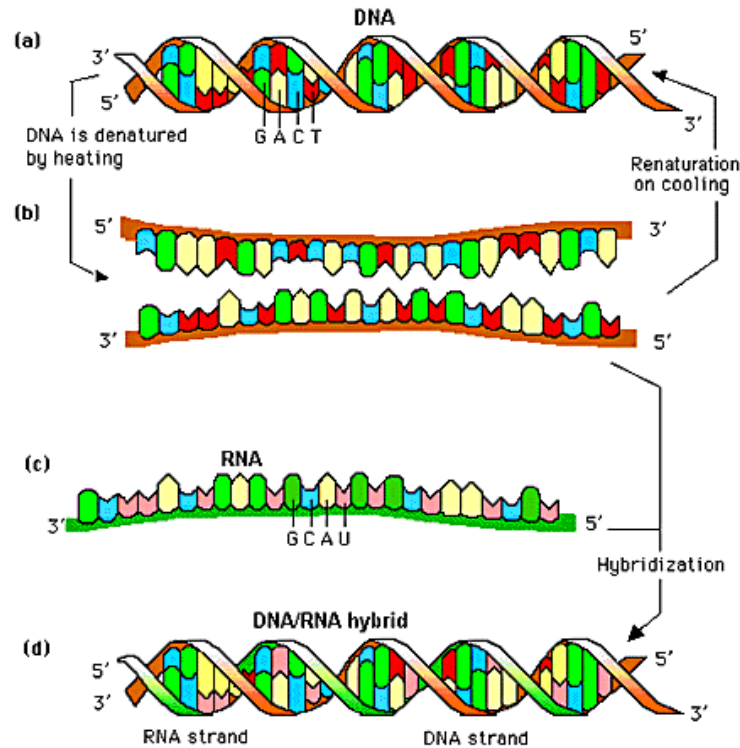


DNA microarrays

DNA microarrays rely on the **hybridization** properties of nucleic acids to monitor DNA or RNA abundance on a genomic scale in different types of cells.

The ancestor of microarrays: the Northern blot.

Nucleic acid hybridization



Nucleic Acid Hybridization

Gene expression assays

The main types of gene expression assays:

- Serial analysis of gene expression (SAGE);
- **Short oligonucleotide arrays (Affymetrix);**
- Long oligonucleotide arrays (Agilent Inkjet);
- Fibre optic arrays (Illumina);
- **cDNA arrays (Brown/Botstein).**

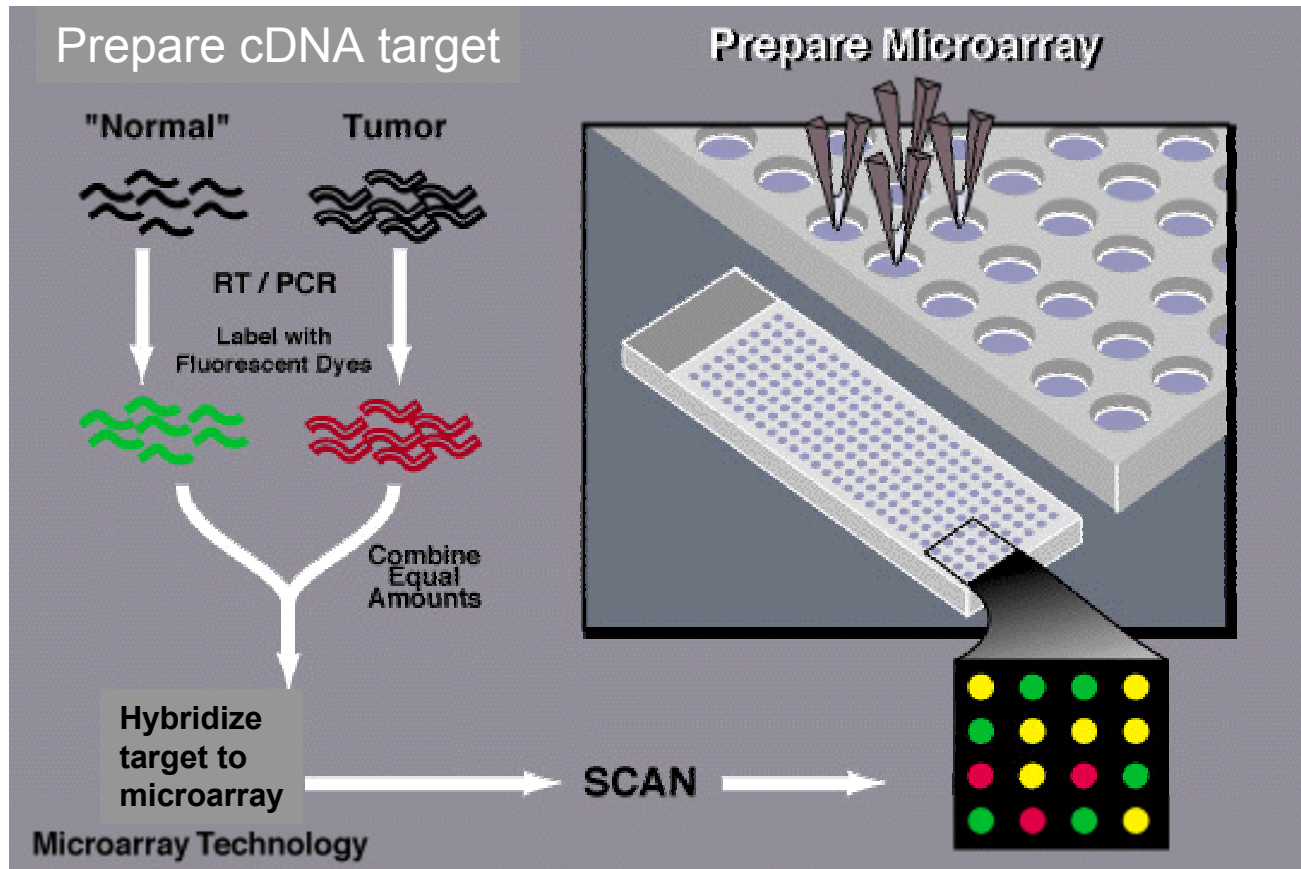
Applications of microarrays

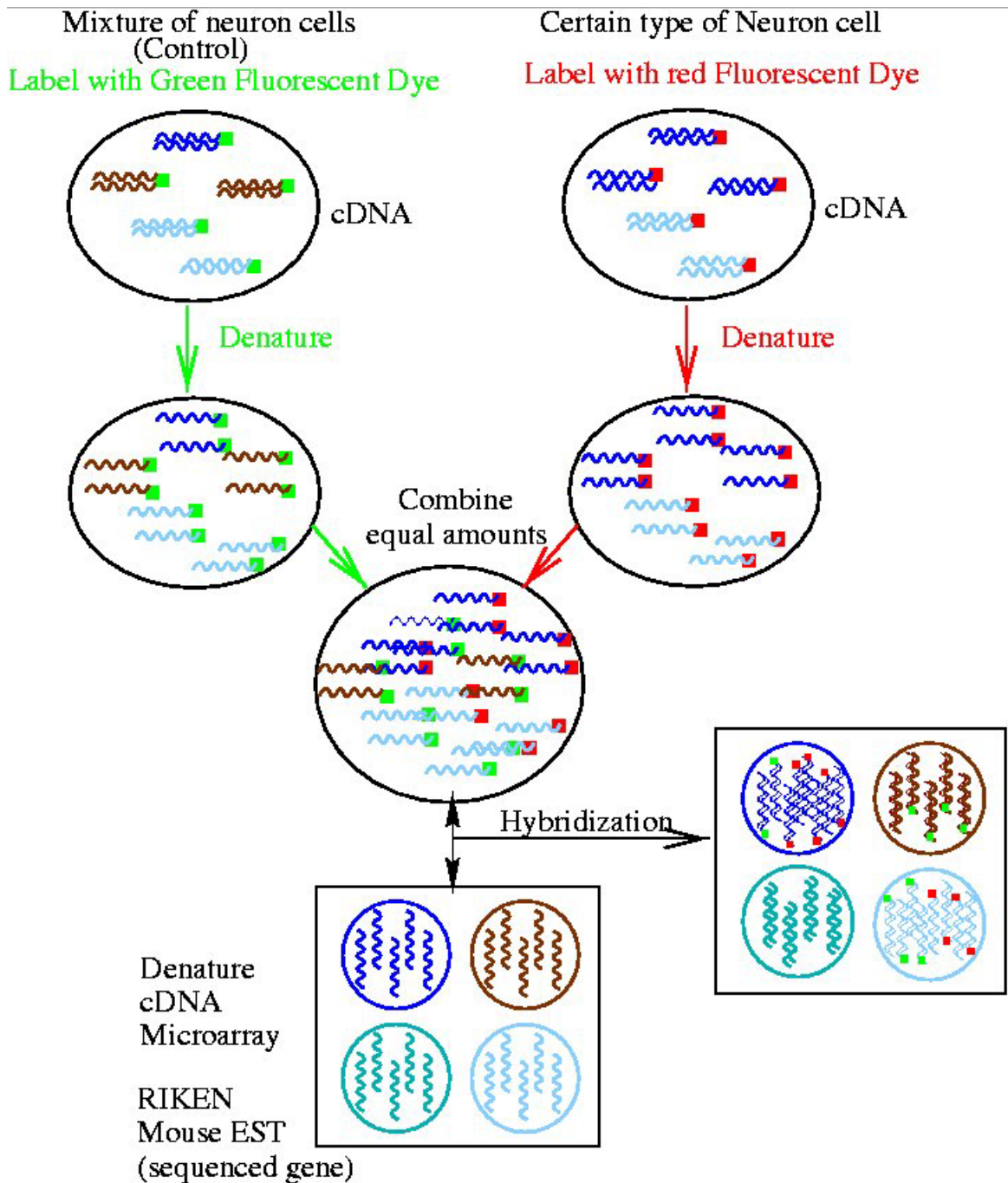
- **Measuring transcript abundance (cDNA arrays);**
- Genotyping;
- Estimating DNA copy number (CGH);
- Determining identity by descent (GMS);
- Measuring mRNA decay rates;
- Identifying protein binding sites;
- Determining sub-cellular localization of gene products;
- ...

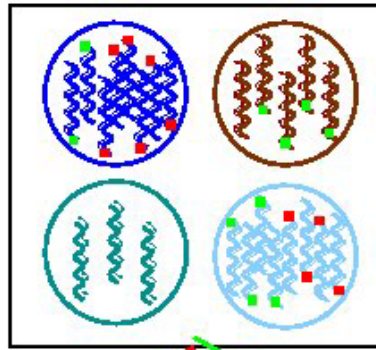
Applications of microarrays

- **Cancer research:** Molecular characterization of tumors on a genomic scale
→ more reliable diagnosis and effective treatment of cancer.
- **Immunology:** Study of host genomic responses to bacterial infections; reversing immunity.
- ...

cDNA microarray experiment







Scan for Red
Wavelength

Scan for Green
Wavelength

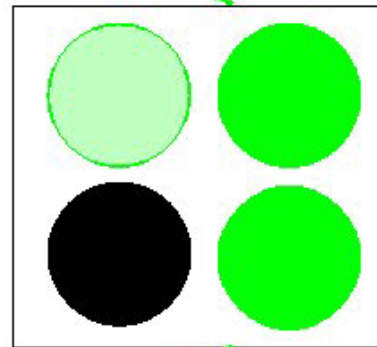
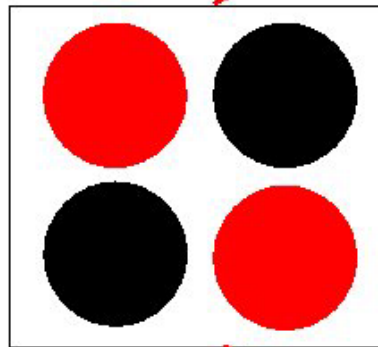
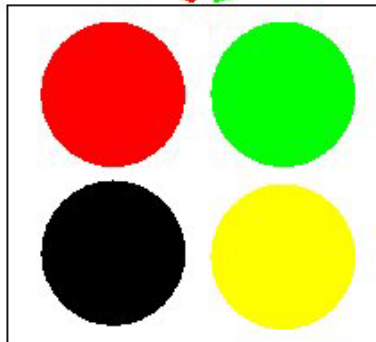


Image Programs
ScanAlyze



The process

Building the chip:

MASSIVE PCR



PCR PURIFICATION
AND PREPARATION



PREPARING
SLIDES



PRINTING



RNA preparation:

CELL CULTURE
AND HARVEST



RNA ISOLATION



cDNA PRODUCTION



Hybing the chip:

POST PROCESSING



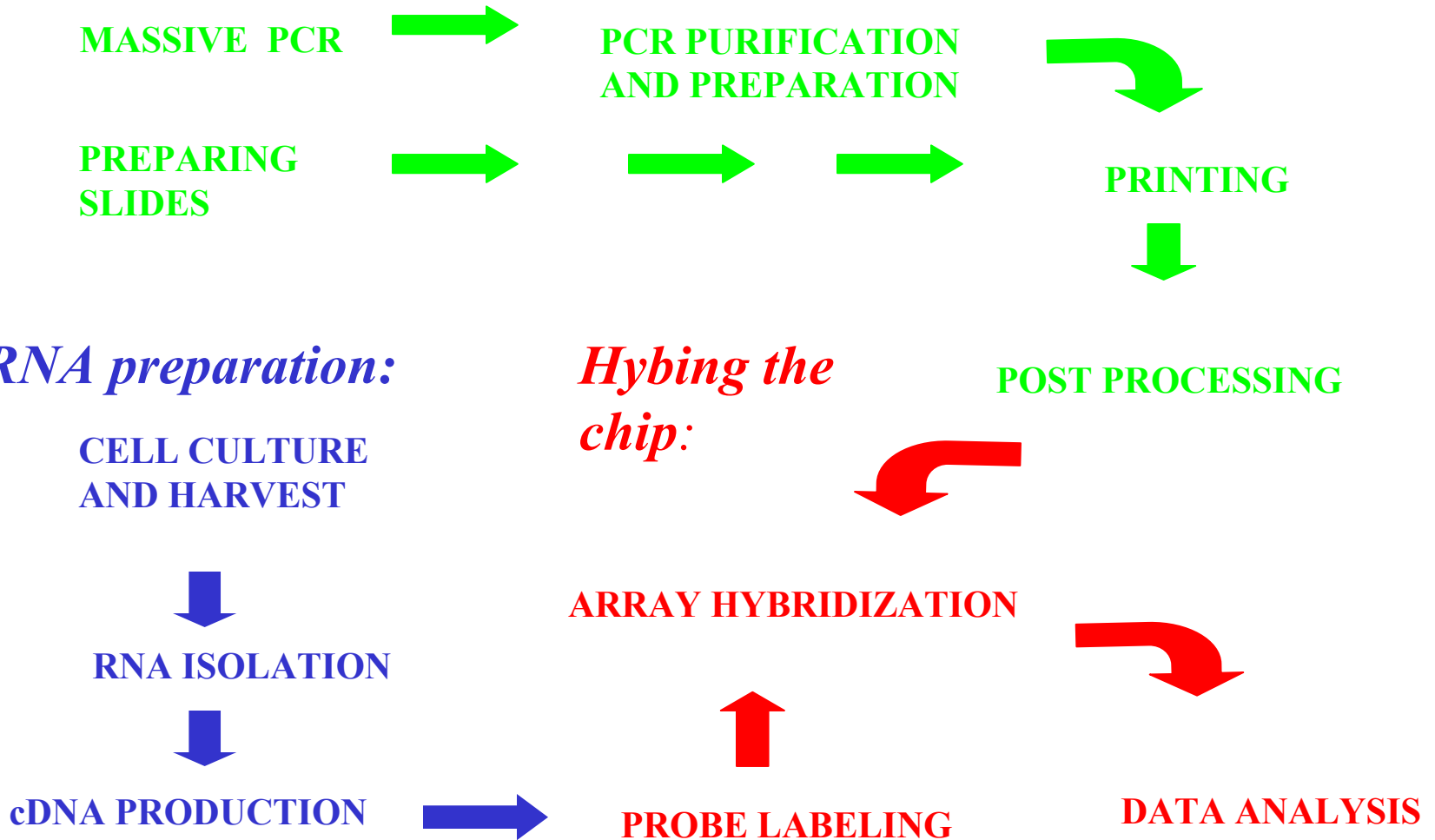
ARRAY HYBRIDIZATION



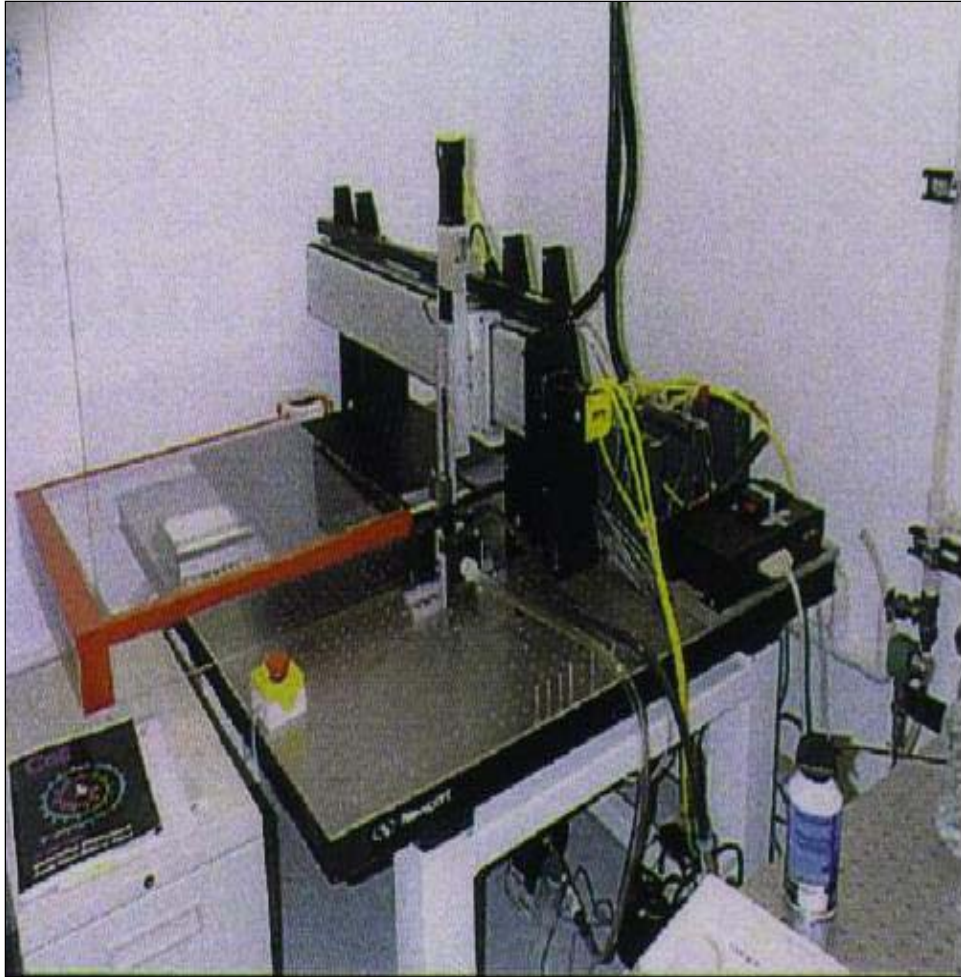
PROBE LABELING



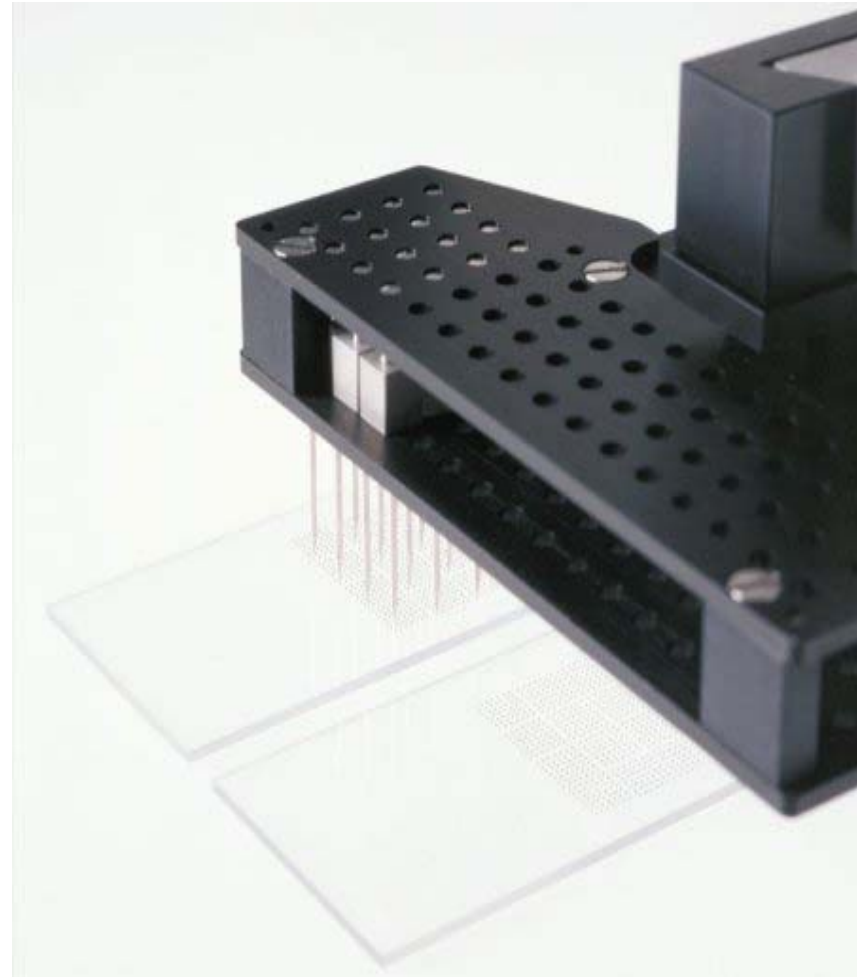
DATA ANALYSIS



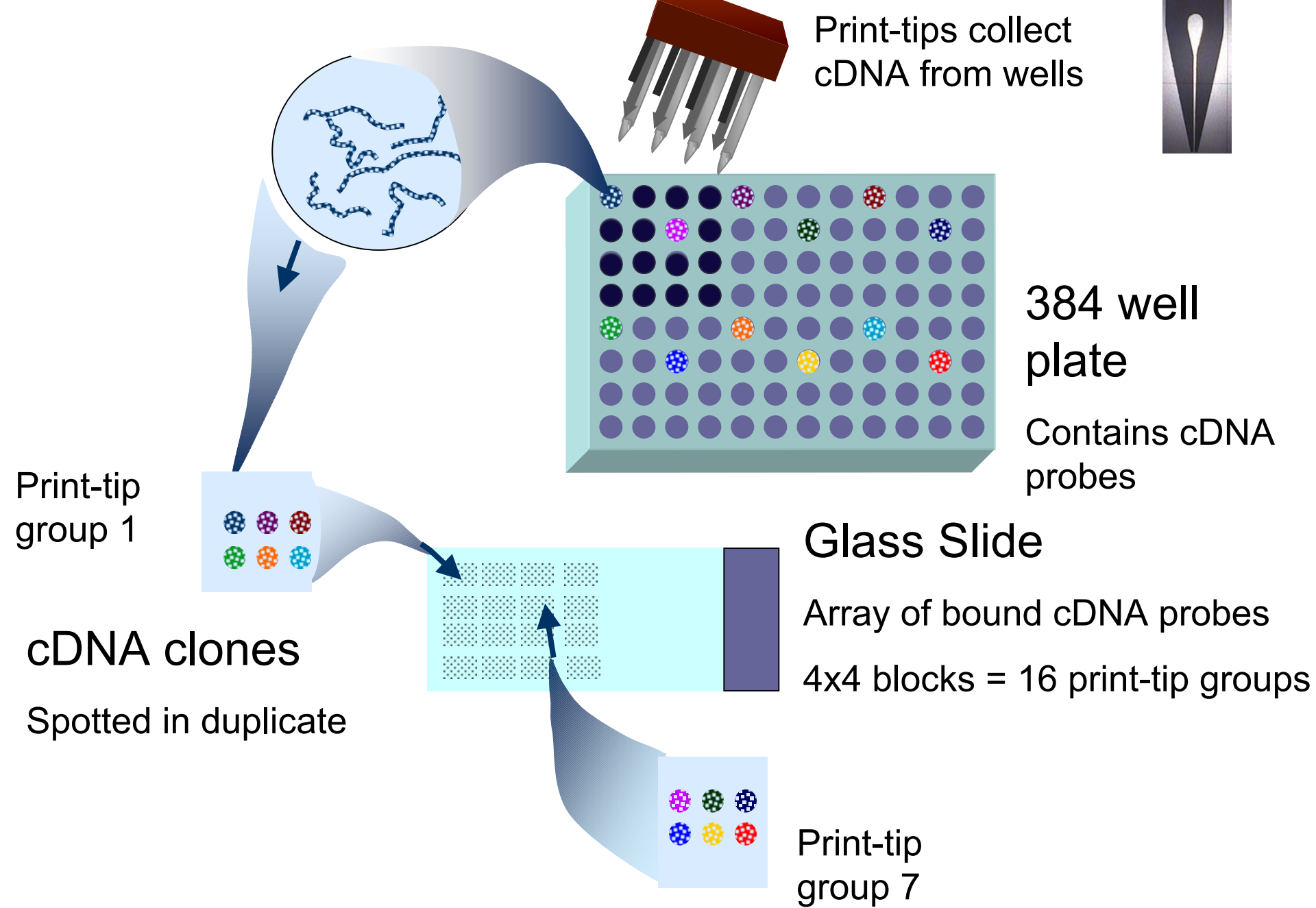
The arrayer



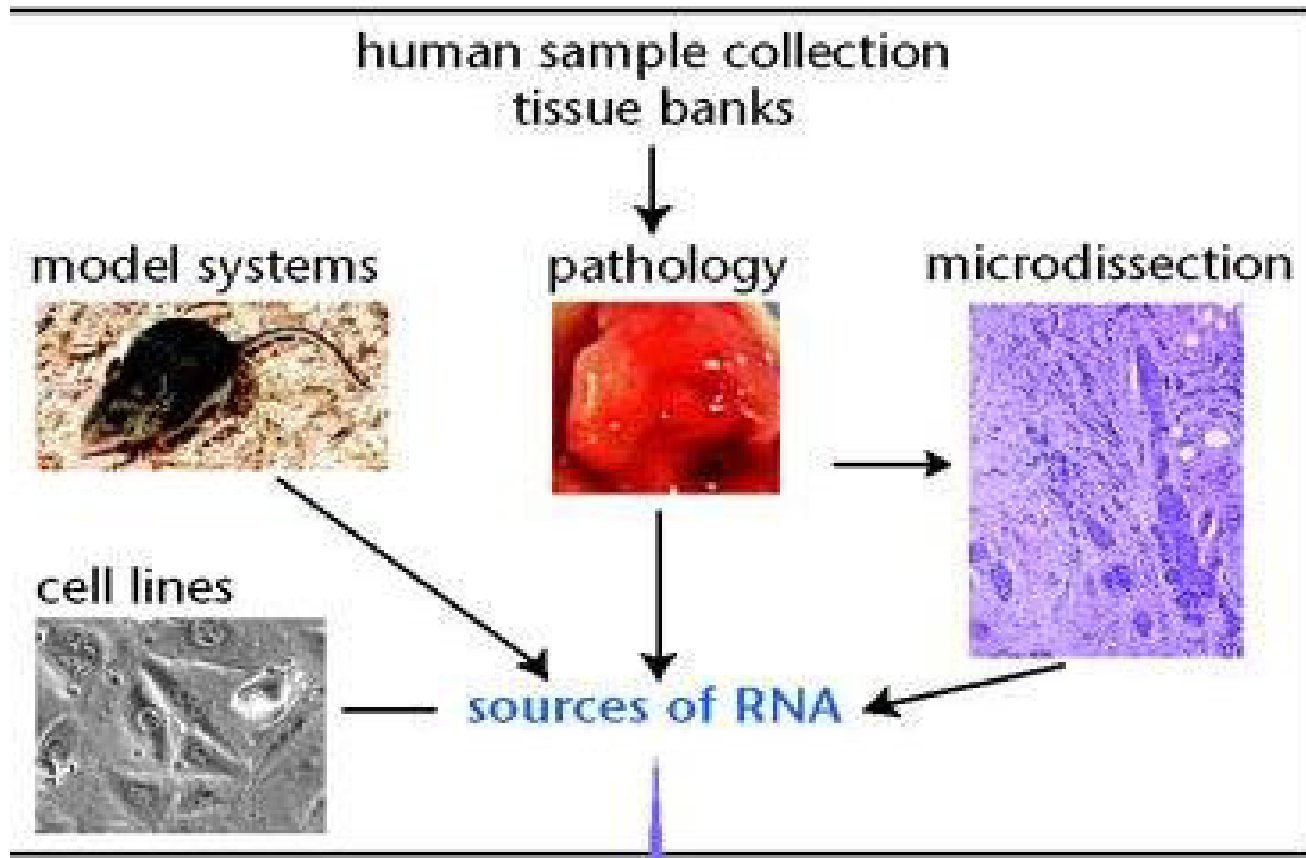
Ngai Lab arrayer , UC Berkeley



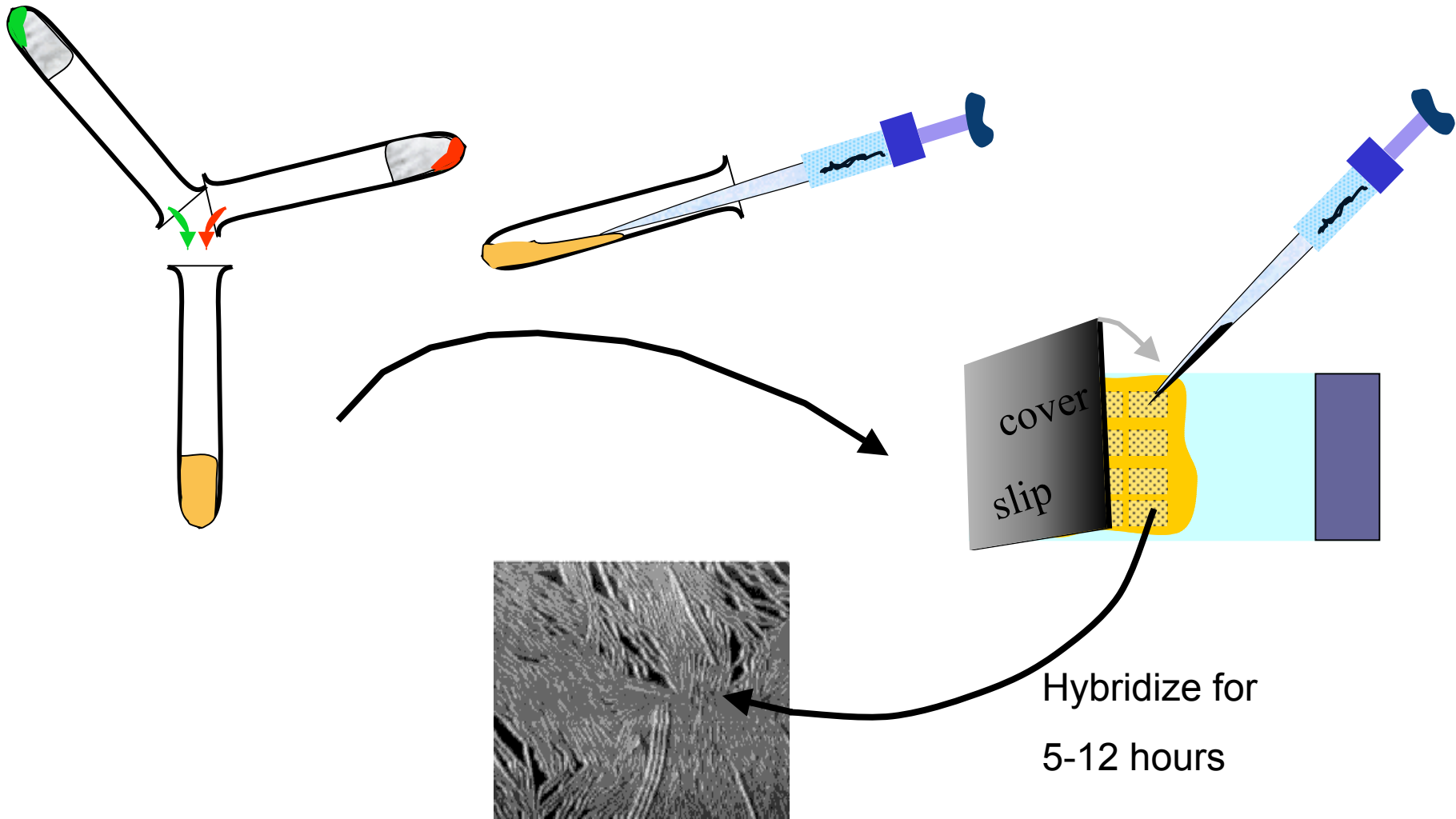
Print-tip head



Sample preparation

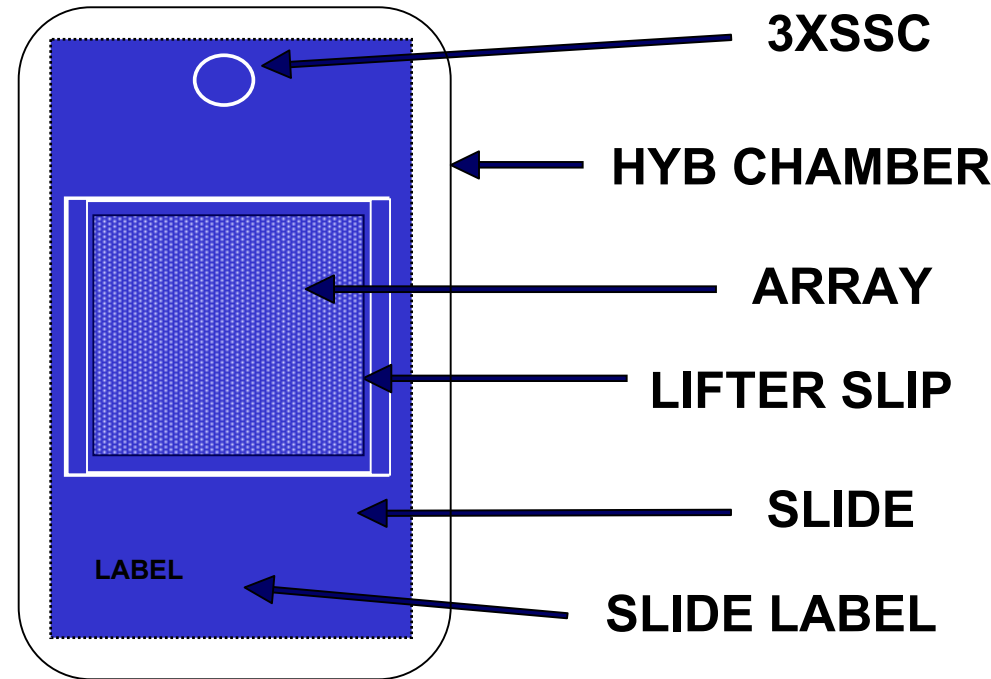


Hybridization



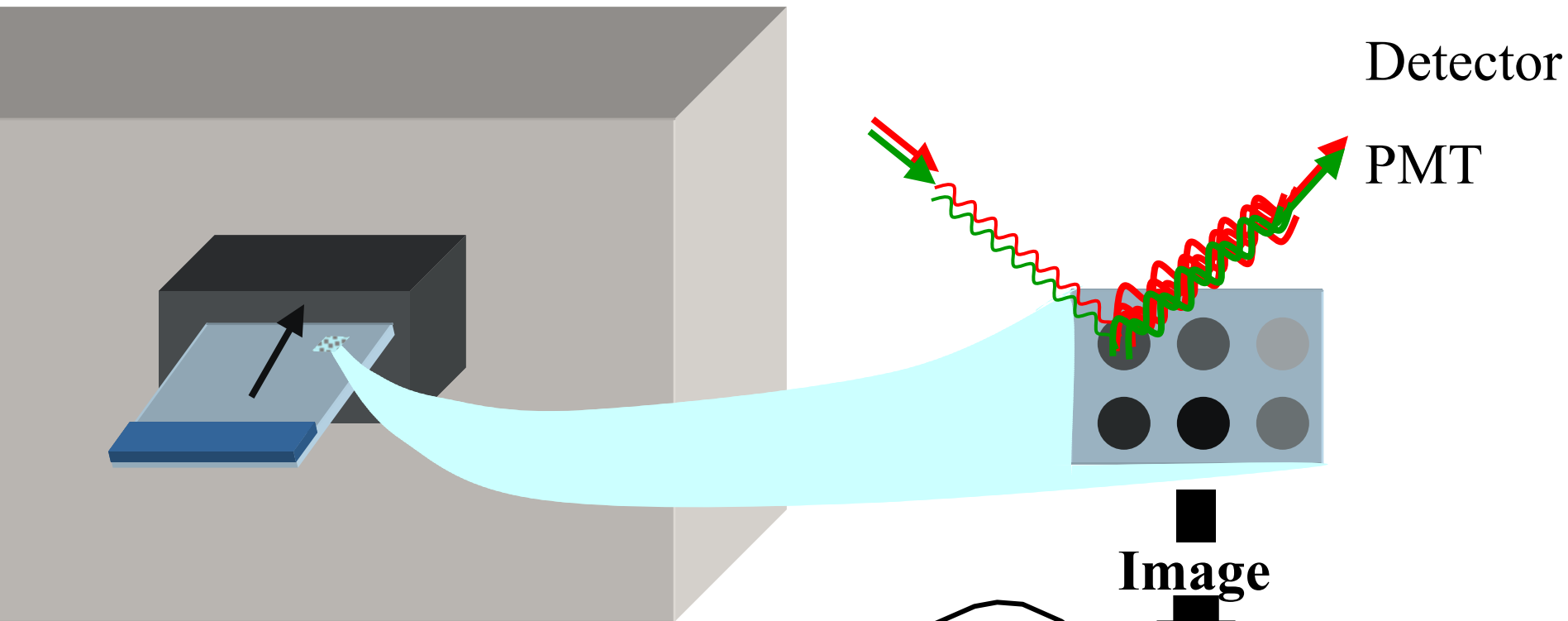
Binding of cDNA target samples to cDNA probes on the slide

Hybridization chamber



- Humidity
- Temperature
- Formamide
(Lowers the Tmp)

Scanning

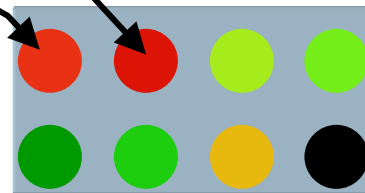


Image

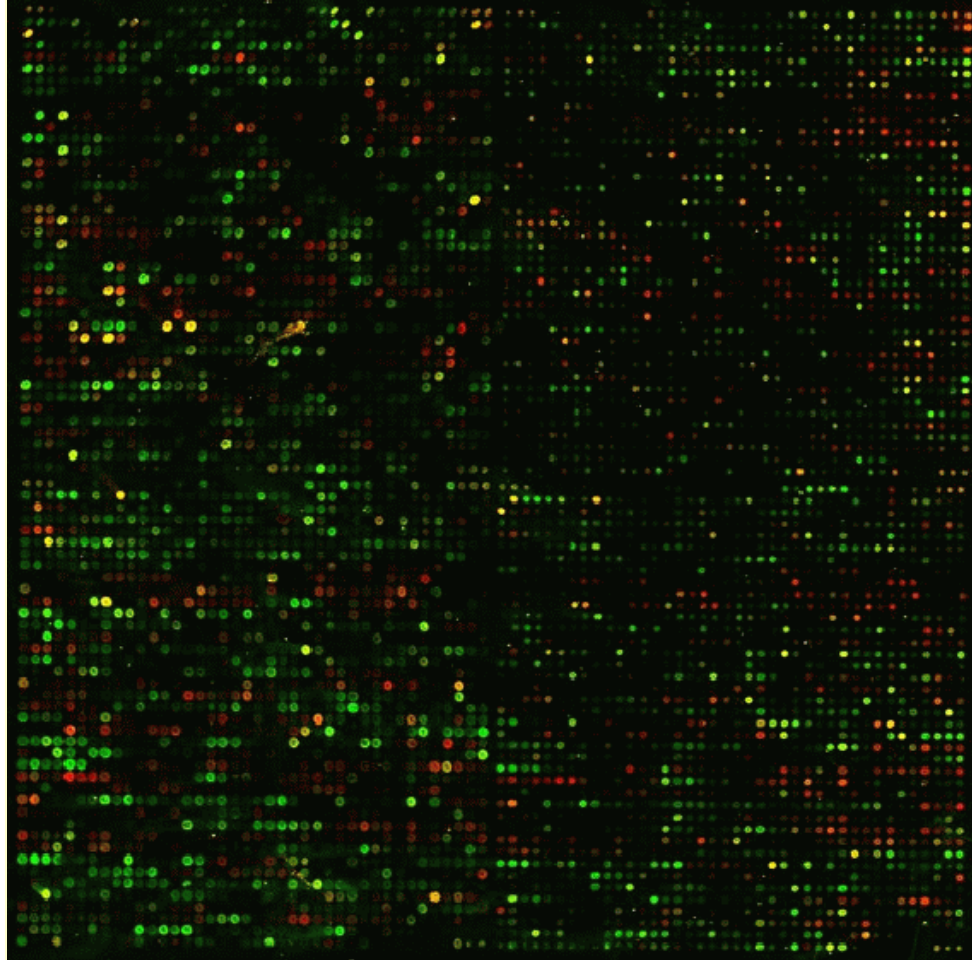
Duplicate spots

Cy5: 635nm

Cy3: 532nm



RGB overlay of Cy3 and Cy5 images



Raw data

- Human cDNA arrays
 - ~43K spots;
 - 16-bit TIFFs: ~ 20Mb per channel;
 - ~ 2,000 x 5,500 pixels per image;
 - Spot separation: ~ 136um;
 - For a “typical” array:
Mean = 43, med = 32, SD = 26 pixels per spots

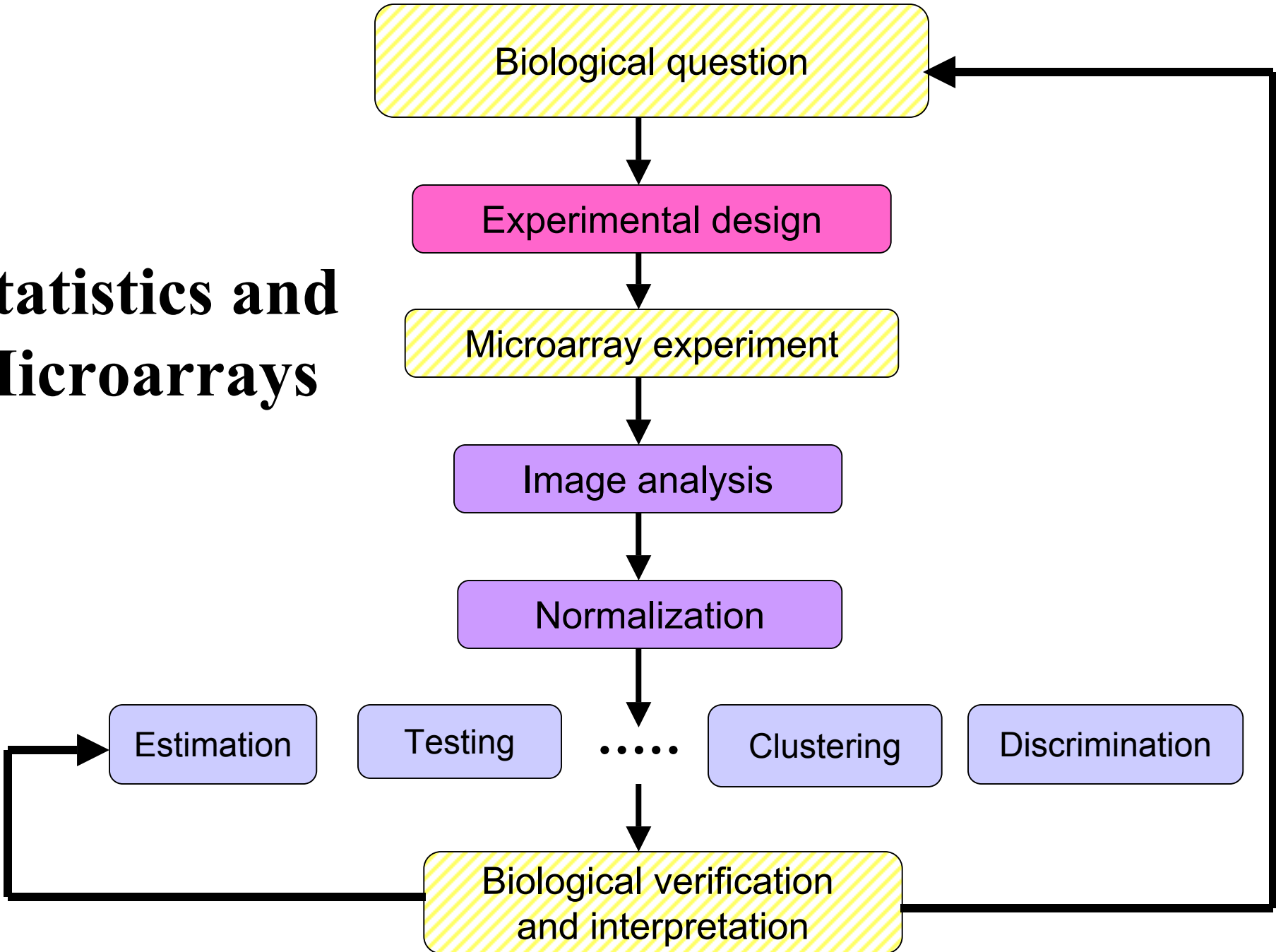
WWW resources

- Complete guide to “microarraying”
<http://cmgm.stanford.edu/pbrown/mguide/>
<http://www.microarrays.org>
 - Parts and assembly instructions for printer and scanner;
 - Protocols for sample prep;
 - Software;
 - Forum, etc.
- Animation:
<http://www.bio.davidson.edu/courses/genomics/chip/chip.html>

Integration of biological data

- Expression, sequence, structure, annotation.
- Integration will depend on our using a common language and will rely on database methodology as well as statistical analyses.
- This area is largely unexplored.

Statistics and Microarrays



Statistical computing

Everywhere ...

- for statistical design and analysis:
pre-processing, estimation, pattern discovery and recognition, etc.
- for integration with biological information resources
(in-house and external databases).

Road map

- Lecture 1, Part II: cDNA arrays
 - Pre-processing: Image analysis;
 - Pre-processing: Normalization;
 - Experimental design.

Road map

- Lecture 2: Differential expression.
- Lecture 3: Applications of HMMs to sequence analysis.
- Lecture 4: Affymetrix chips.
- Lecture 5: Classification.