# Lecture 1
## Introduction to Micorarrays and Concepts of Molecular Biology

M. Saleet Jafri

Program in Bioinformatics and Computational Biology

George Mason University

# Lecture 1

## Overview of Molecular and Cellular Biology
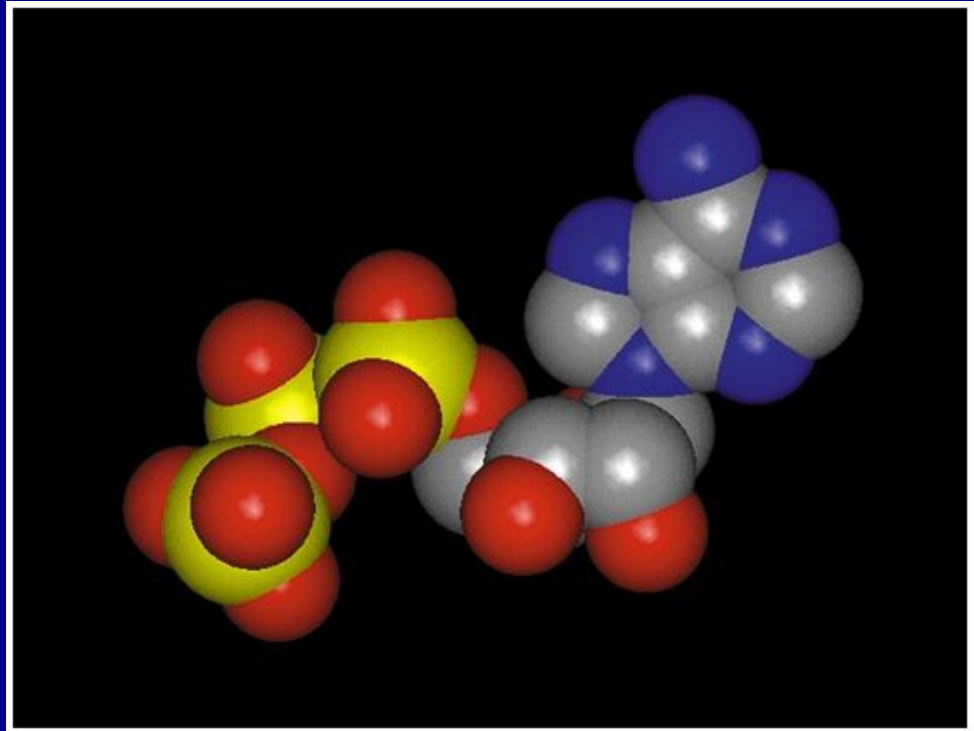
# Biological References

É*Molecular Biology of the Cell*
  **by Bruce Alberts**
  **(1994 or newer edition)**

É*Molecular Cell Biology*
  **by Darnell, Lodish, and Baltimore**
  **(1995 or newer edition)**
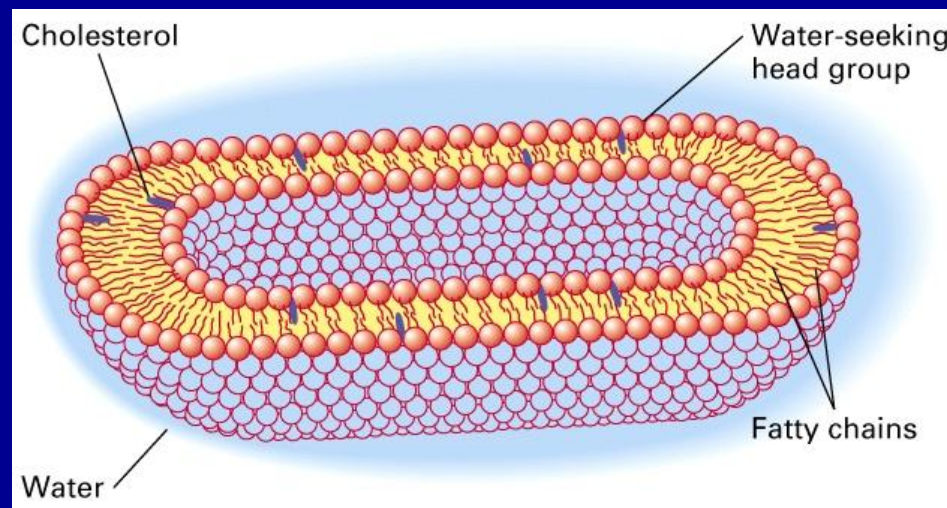
# Part I: Molecular Biology Review

## Where do biological sequences come from?

- É  Life and evolution
- É  Proteins
- É  Nucleic Acids
- É  Central dogma
- É  Genetic code
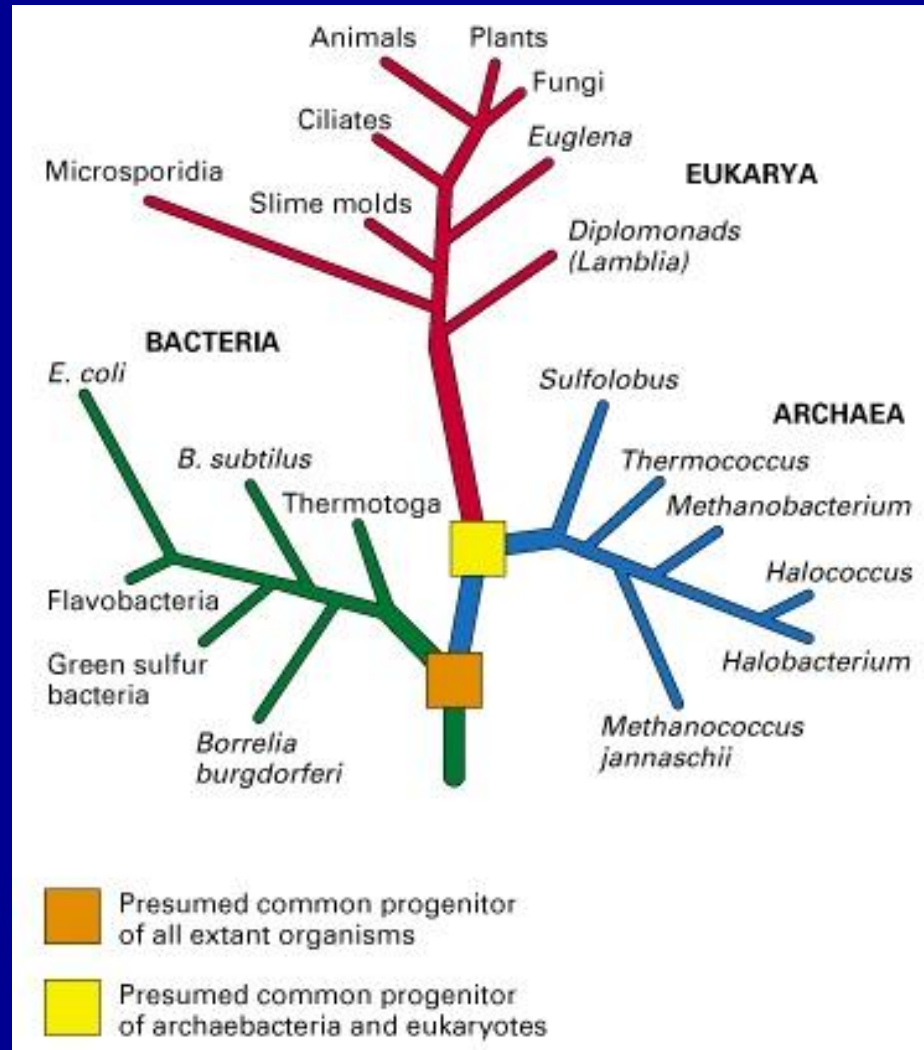- É  DNA structure
- É  Mitochondrial DNA

# Life

- É   Evolved from common origin
- É   ''3.5 billion years ago
- É   All life shares similar biochemistry
    - ó   Proteins:  active elements
    - ó   Nucleic acids: informational elements
- É    Molecular Biology: the study of structure and function of proteins and nucleic acids



Cholesterol — Water-seeking head group — Fatty chains — Water
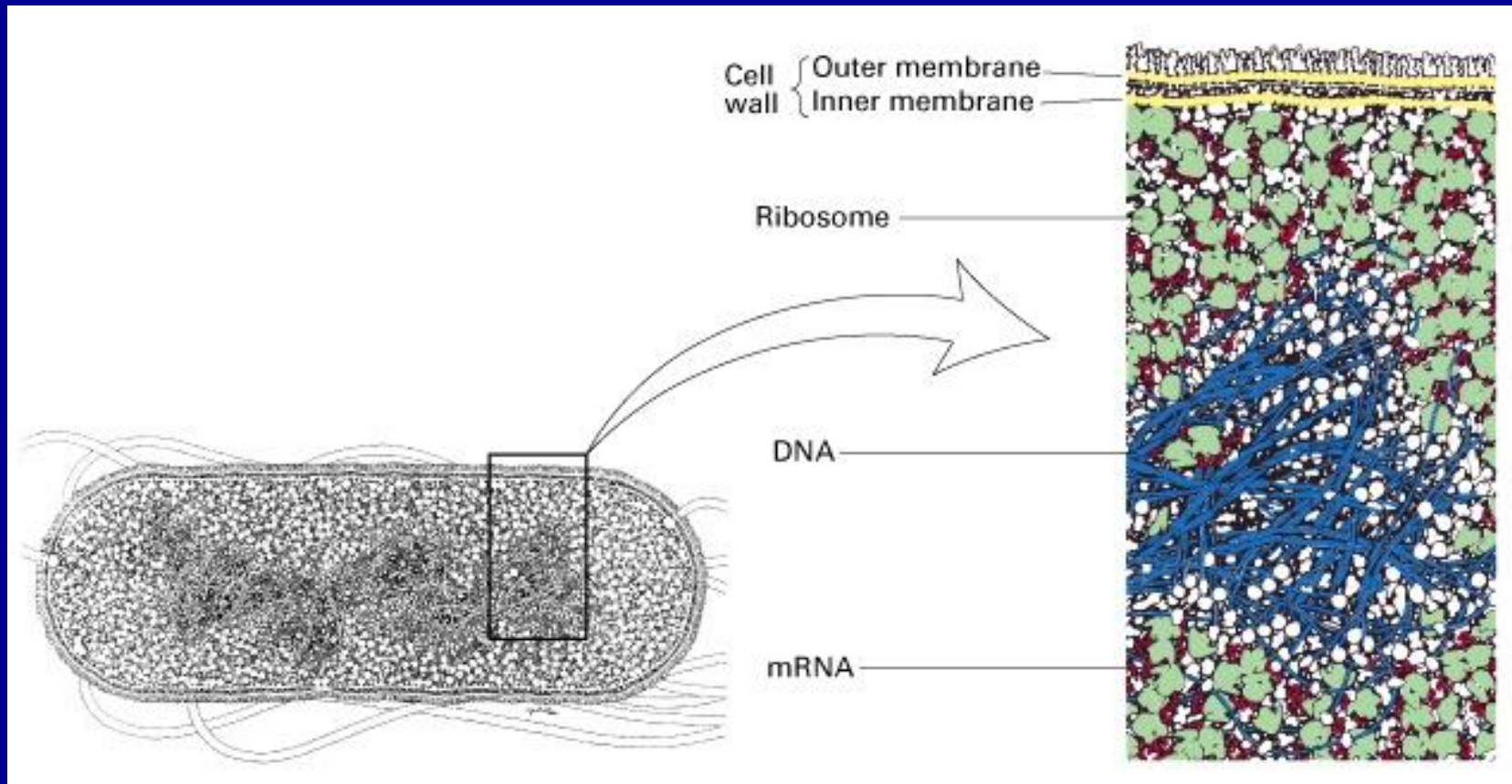
# Terrestrial Life

# Cell types

Prokaryotes ó no nuclear membrane, represented by cyanobacteria (blue-green algae) and common bacteria (*Escherichia coli*)

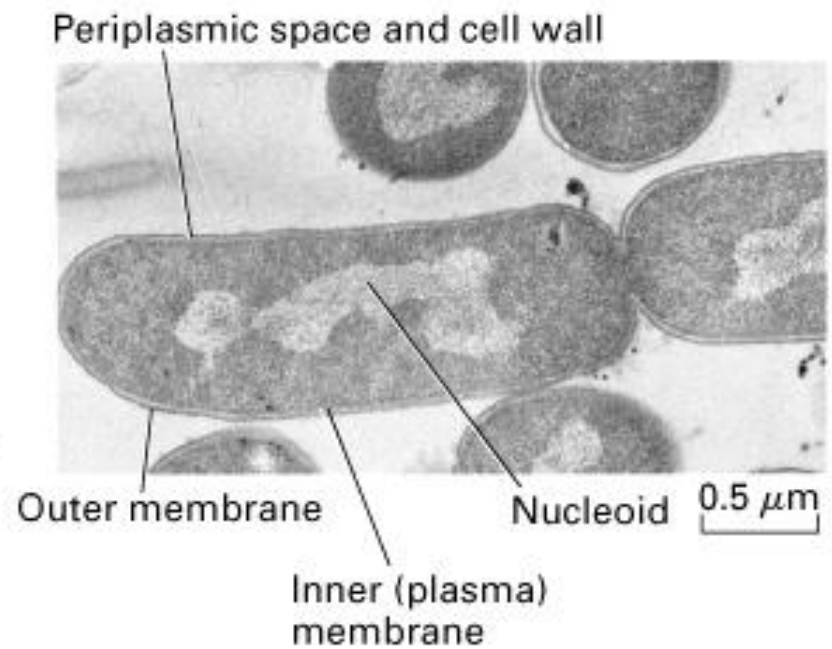Eukaryotes ó unicellular organisms such as yeast and multicellular organisms

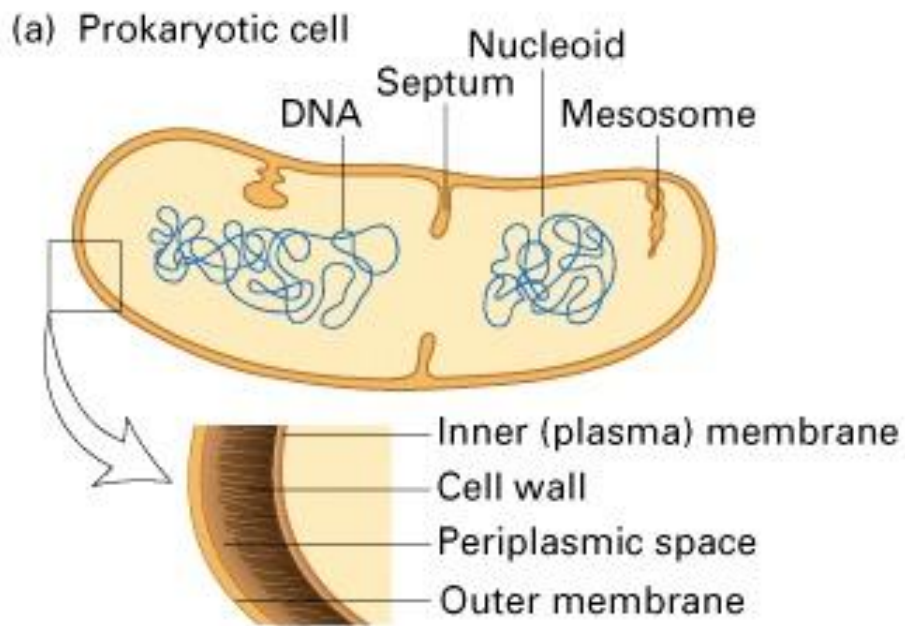Archaebacteria ó no nuclear membrane but similar to eukaryotes in transcription and translation mechanisms, discovered in deep sea thermal vents in 1982

# Prokaryotic Cell

# Prokaryotic Cell



(a) Prokaryotic cell

DNA — Septum — Nucleoid — Mesosome

Inner (plasma) membrane
Cell wall
Periplasmic space
Outer membrane

Periplasmic space and cell wall

Outer membrane — Nucleoid — 0.5 μm

Inner (plasma) membrane

# Eukaryotes

- In eukaryotes, transcription is complex:
  - Many genes contain alternating exons and introns
  - Introns are spliced out of mRNA
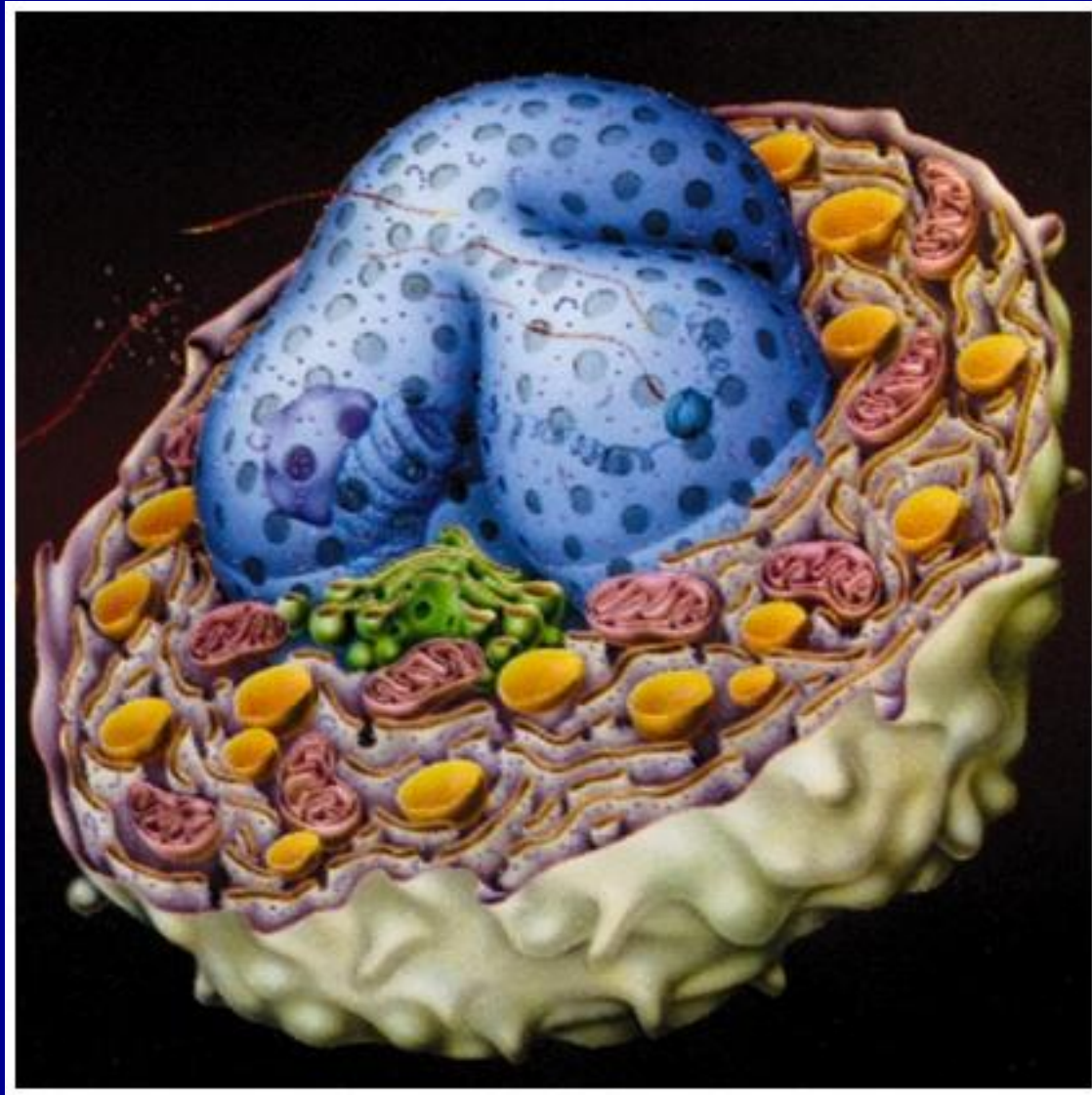  - mRNA then leaves the nucleus to be translated by ribosomes
- Genomic DNA: entire gene including exons and introns
  - The same genomic DNA can produce different proteins by alternative splicing of exons
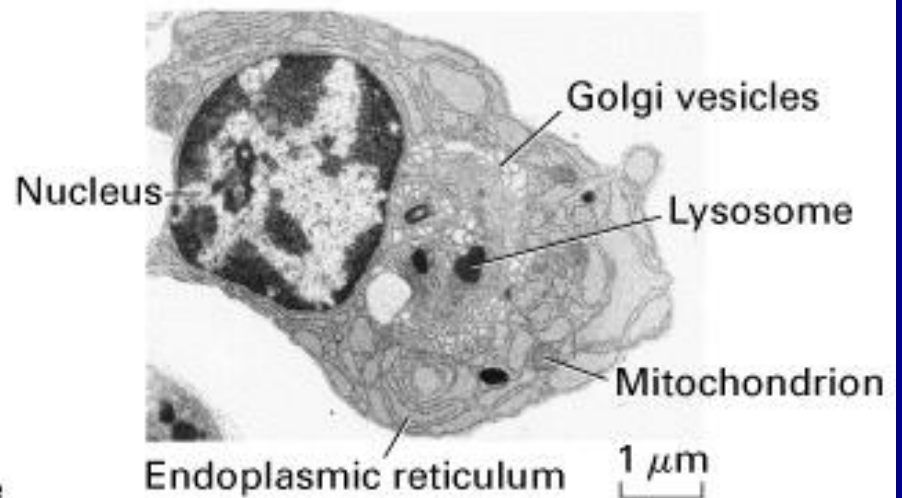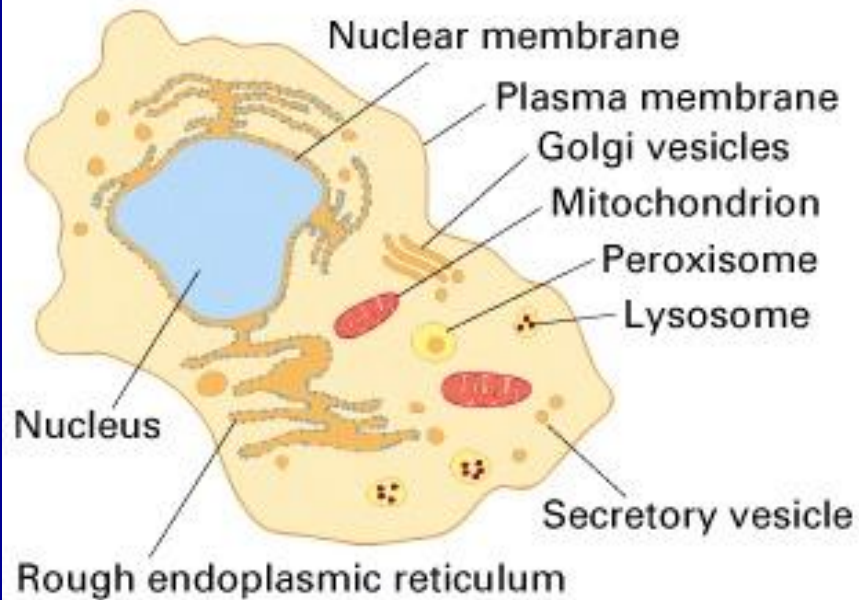- Complementary DNA (cDNA): spliced sequence containing only exons
  - cDNA can be manufactured by capturing mRNA and performing reverse transcription

# Eukaryotic Cell

# Eukaryotic Cell



(b) Eukaryotic cell

Nuclear membrane
Plasma membrane
Golgi vesicles
Mitochondrion
Peroxisome
Lysosome
Nucleus
Secretory vesicle
Rough endoplasmic reticulum

Nucleus
Golgi vesicles
Lysosome
Mitochondrion
Endoplasmic reticulum
1 $\mu$m

# Eukaryotic Cell Organelles

É   Cell membrane

É   Nucleus

É   Cytoplasm

É   Endoplasmic Reticulum ó rough and smooth

É   Golgi Apparatus ó received newly formed proteins from the ER and modifies them and directs them to final destination
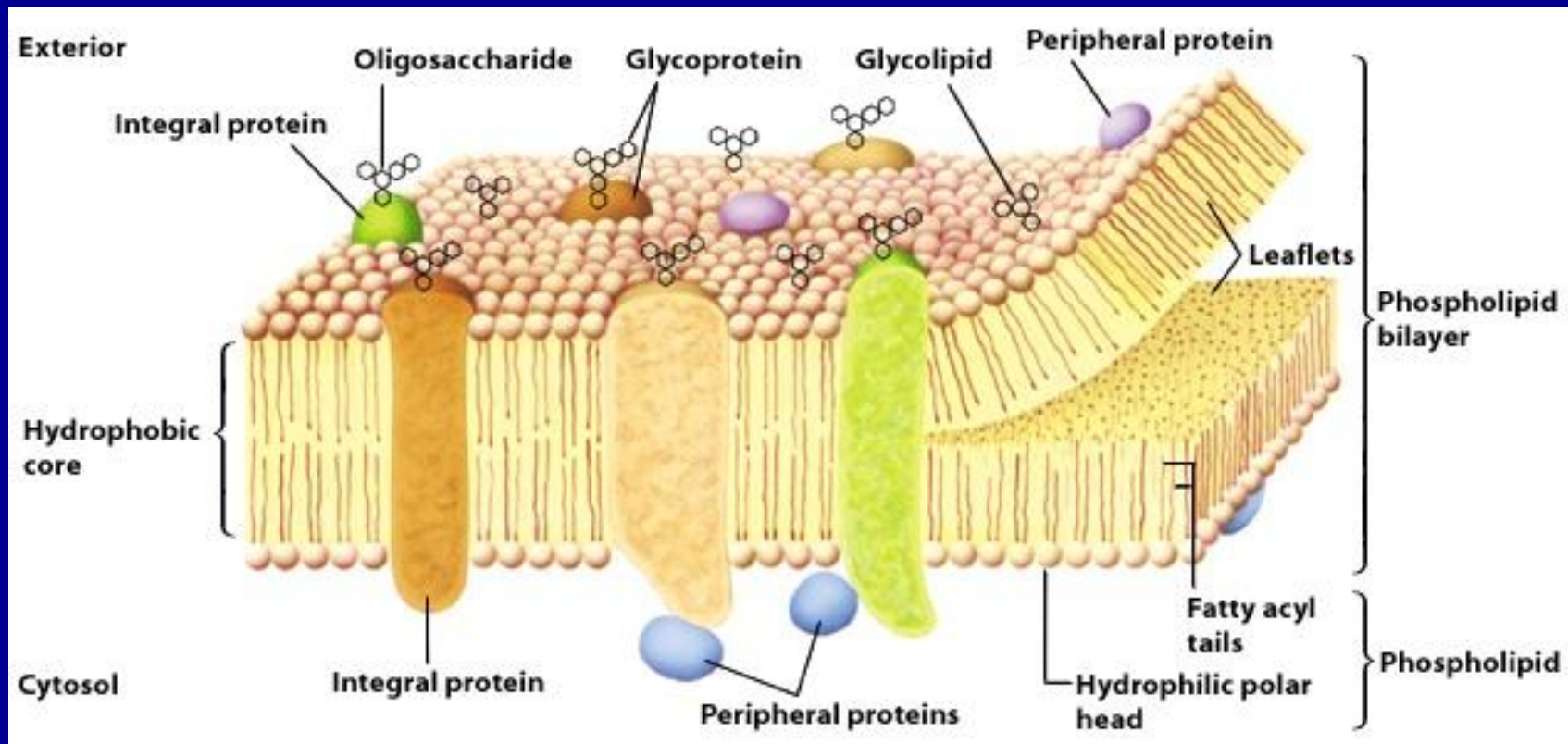
É   Mitochondria ó respiratory centers, have own circular DNA, bacterial origin

# Eukaryotic Cell Organelles

É  Chromosomes ó chromatin, histones, centromeres and arms (2 pairs in eukaryotes)

É  Lysosomes ó contain acid hydrolases ó nucleases, proteases, glycodidases, lipases, phosphatases, sulfatases, phospholipases

É  Peroxisomes ó use oxygen to remove hydrogen from substrates forming $H_2O_2$, abundant in kidney and liver ódetoxification
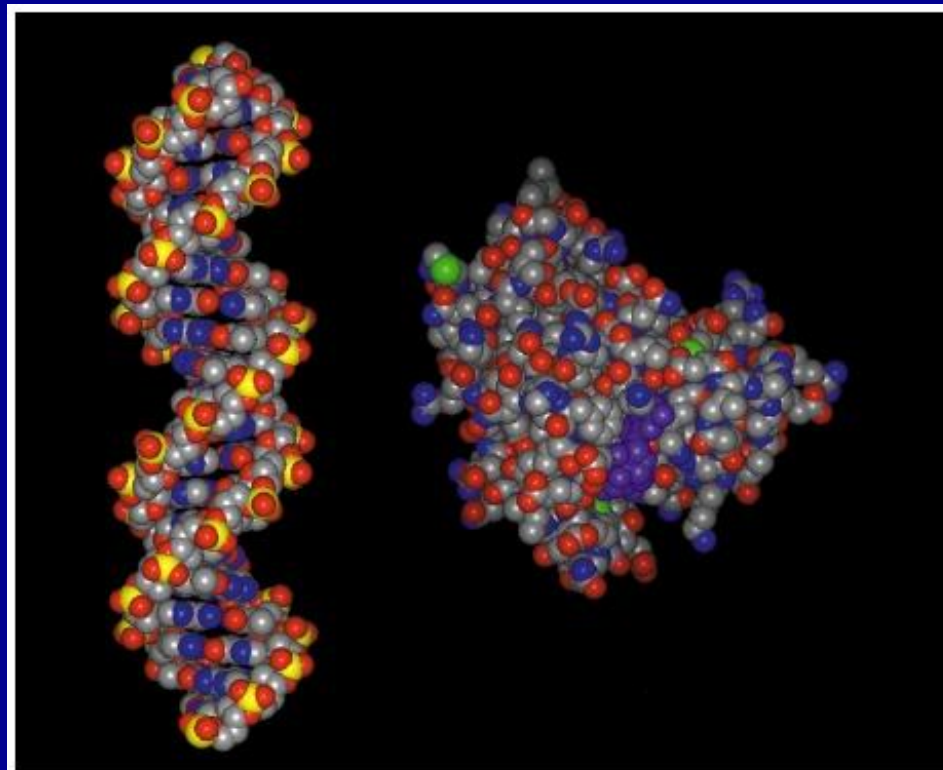
É  Cytoskeleton

# Eukaryotic Membrane

# Nucleic Acids

É   Two kinds:

    ó   RNA: ribonucleic acid

    ó   DNA: deoxyribonucleic acid



Nucleic acid
(DNA)

Protein
(Ras protein)

1 nm

# Nucleic Acid Structure

# DNA

- É Double stranded
- É Four bases: adenine (A), guanine (G), cytosine (C) and thymine (T)
- É A and G are purines
- É C and T are pyrimidines
- É A always paired with T (complementary)
- É C always paired with G (complementary)
- => Watson-Crick base pairs (bp)
- É DNA may consist of hundreds of millions bp
- É A short sequence (<100) is called an oligonucleotide

# RNA

- É Different sugar (ribose instead of 2ø-deoxyribose)
- É Uracil (U) instead of thymine (U binds with A)
- É RNA does not form a double helix
- É RNA may have a complex three-dimensional structure

# Central Dogma

DNA     RNA     Protein

DNA = Deoxyribonucleic Acid

RNA = Ribonucleic Acid

Protein = Functional and Structural units of cells

**Flow of Information is unidirectional**

The Central Dogma of Molecular Biology

# Gene Transcription or DNA Transcription

- RNA molecules synthesized by *RNA polymerase*
- RNA polymerase found in free and bound form
- RNA polymerase binds very tightly to *promoter* region on DNA
- Promoter region contains *start site*
- Transcription ends at *termination signal* site.
- Primary transcript ó direct coding of RNA from DNA
- RNA splicing ó introns removed to  make the mRNA
- mRNA ó contains the sequence of *codons* that code for a protein
- uracil replaces thymine
- splicing and alternative splicing

# Transcription of DNA to Messenger RNA

# Translation

- Ribosomes made of protein and ribosomal RNA (rRNA)
- Transfer RNA (tRNA) make connection between specific codons in mRNA and amino acids
  - As tRNA binds to the next codon in mRNA, its amino acid is bound to the last amino acid in the protein chain
- When a STOP codon is encountered, the ribosome releases the mRNA and synthesis ends

# Gene Translation

- *tRNA* ó links an amino acid to the codon on the mRNA via the *anit-codon*
- *rRNA* ó RNA found in ribosomes
- *ribosomes* ó large and small subunit, made of protein and rRNA
- *initiator tRNA* always carries methionine
- *initiation factors* ó proteins that catayze the start of transcription
- stop codon
- Endoplasmic Reticulum
- Posttranscriptional modification

# Translation

˝ **Involves ribosomes, and RNA**

˝ **Ribosomes made of protein and RNA**

˝ **Messenger RNA (mRNA) is the sequence transcribed from the DNA**

˝ **The mRNA is 'threaded' through the ribosomes.**

˝ **Transfer RNA (tRNA) brings the different amino a cids to the ribosome complex so that the amino acids can be attached to the growing amino acid chain.**

tRNA

Protein – amino acid chain

mRNA

Amino acid

Ribosome

# Ribosome

# Eukaryotic and Prokaryotic

**The Genetic Code**

RNA

Ribonucleic acid

# DNA Structure

- É DNA contains:
  - ó Promoters
  - ó Genes
  - ó Junk DNA
- É Reading frames
- É An open reading frames (ORF):  a contiguous sequence of DNA starting at a start codon and ending at a  STOP codon

# Chromosomes

- A genome: a complete set of chromosomes within a cell
- Different species have different numbers of chromosomes in their genomes
- Prokaryotes usually have a single chromosome, often a circular DNA molecule
- Eukaryotic chromosomes appear in pairs (diploid), each inherited from one parent
  - Homologous chromosomes carry the same genes
  - Some genes are same in both parents
  - Some genes appear in different forms called alleles
    - E.g. human blood has three alleles: A, B, and O
- All genes are present in all cells, but a given cell types only expresses a small portion of the genes

# Chromosomes

# Chromatin Structure

# Gene Coding and Replication

- É  Double helix
- É  Nitrogenous bases A,T,G,C
- É  Sugar-Phosphate backbone
- É  Nucleotide ó sugar + base + phosphate group
- É  Nucleoside ó sugar + base
- É  Purines ó adenine, guanine
- É  Pyrimidines ó cysteine, thymine
- É  A-T ó 2 H bonds, G-C ó 3 H bonds

# Gene Coding and Replication

- 5ʹ end contains a phosphate group
- 3ʹ end is free
- DNA extended from 5ʹ to 3ʹ
- Gene is a segment of DNA that codes for a specific protein
- Exons are coding regions of the DNA
- Introns are "in between" regions, found in eukaryotes
- Codons
- Reading frame
- *Consensus sequences* are conserved regions found in a particular type of regulatory region

# Mitosis



Parental cell (2n) in G₁

DNA replication

Parental cell (4n) in G₂

Chromatid

Mitotic apparatus — Metaphase cell

Anaphase cell

Cytokinesis

Daughter cells (2n)

QuickTime™ and a
Sorenson Video decompressor
are needed to see this picture.

# Proteins

- Functions:
  - Structural proteins
  - Enzymes
  - Transport
  - Antibody defense
- Chains of amino acids
- Typical size ~300 residues

# Protein Folding

É    Primary structure ó amino acid sequence

É    Secondary structure ó local structure such as $\alpha$ helix and $\beta$ sheets

É    Tertiary structure ó 3-dimensional structure of a protein monomer

É    Quarternary structure ó 3-dimensional structure of a fully functional protein (protein complexes).

Primary structure:  residue
   sequence

Secondary structure:
   local structures
      (Helices, sheets, loops)

Tertiary structure:
   position of each atom

Quaternary structure:
   how groups of proteins pack
   together

**Primary protein structure**
is sequence of a chain of amino acids

Amino Acids

Pleated sheet          Alpha helix

**Secondary protein structure**
occurs when the sequence of amino acids
are linked by hydrogen bonds

Pleated sheet

**Tertiary protein structure**
occurs when certain attractions are present
between alpha helices and pleated sheets.

Alpha helix

**Quaternary protein structure**
is a protein consisting of more than one
amino acid chain.

Mechanism of enzyme activity

The protein's 3-d shape determines what molecules it can bind to

Unsolved problem: predicting a protein's folding pattern based on its sequence

# Cell Signaling and Biochemical Pathways

É Surface receptors

É G-proteins, kinases, etc

É Transcription factors

É Other biochemical reactions ó glycolysis, citric acid cycle, etc.

# Molecular Biology Summary

- É  Life and evolution
- É  Proteins
- É  Nucleic Acids
- É  Eukaryotes versus Prokaryotes
- É  Ribosome
- É  Translation
- É  Transcription
- É  Central dogma
- É  Genetic code
- É  DNA structure
- É  Chromosome
- É  Mitosis

# Polymerase Chain Reaction

In order to make enough RNA for microarray experiments, PCR is used to amplify the amount of RNA

Most PCR methods typically amplify nulceic fragments of up to 10 kilo base pairs (kb), although some techniques allow for amplification of fragments up to 40 kb in size.

# Polymerase Chain Reaction

A basic PCR set up requires several components and reagents. These components include:

É *DNA template* that contains the DNA region (target) to be amplified.

É Two *primers*, which are complementary to the DNA regions at the 5' (five prime) or 3' (three prime) ends of the DNA region.

É A DNA polymerase such as *Taq polymerase* or another DNA polymerase with a temperature optimum at around 70°C.

É *Deoxynucleoside triphosphates* (dNTPs; also very commonly and erroneously called deoxynucleotide triphosphates), the building blocks from which the DNA polymerases synthesizes a new DNA strand.

É *Buffer solution*, providing a suitable chemical environment for optimum activity and stability of the DNA polymerase.

É *Divalent cations*, magnesium or manganese ions; generally Mg2+ is used, but Mn2+ can be utilized for PCR-mediated DNA mutagenesis, as higher Mn2+ concentration increases the error rate during DNA synthesis

É *Monovalent cation* potassium ions.  (From Wikipedia)

# Polymerase Chain Reaction

The PCR is commonly carried out in a reaction volume of 10-200 l in small reaction tubes (0.2-0.5 ml volumes) in a thermal cycler. The thermal cycler heats and cools the reaction tubes to achieve the temperatures required at each step of the reaction (see below). Many modern thermal cyclers make use of the Peltier effect which permits both heating and cooling of the block holding the PCR tubes simply by reversing the electric current. Thin-walled reaction tubes permit favorable thermal conductivity to allow for rapid thermal equilibration. Most thermal cyclers have heated lids to prevent condensation at the top of the reaction tube. Older thermocyclers lacking a heated lid require a layer of oil on top of the reaction mixture or a ball of wax inside the tube.

From Wikipedia

# Polymerase Chain Reaction

The PCR usually consists of a series of 20 to 40 repeated temperature changes called cycles; each cycle typically consists of 2-3 discrete temperature steps. Most commonly PCR is carried out with cycles that have three temperature steps (Fig. 2). The cycling is often preceded by a single temperature step (called *hold*) at a high temperature (>90°C), and followed by one hold at the end for final product extension or brief storage. The temperatures used and the length of time they are applied in each cycle depend on a variety of parameters. These include the enzyme used for DNA synthesis, the concentration of divalent ions and dNTPs in the reaction, and the melting temperature ($T_m$) of the primers

From Wikipedia

# Polymerase Chain Reaction

*Denaturation step*: This step is the first regular cycling event and consists of heating the reaction to 94-98°C for 20-30 seconds. It causes melting of DNA template and primers by disrupting the hydrogen bonds between complementary bases of the DNA strands, yielding single strands of DNA.

*Annealing step*: The reaction temperature is lowered to 50-65°C for 20-40 seconds allowing annealing of the primers to the single-stranded DNA template. Typically the annealing temperature is about 3-5 degrees Celsius below the Tm of the primers used. Stable DNA-DNA hydrogen bonds are only formed when the primer sequence very closely matches the template sequence. The polymerase binds to the primer-template hybrid and begins DNA synthesis.

From Wikipedia

# Polymerase Chain Reaction

*Extension/elongation step*: The temperature at this step depends on the DNA polymerase used; Taq polymerase has its optimum activity temperature at 75-80°C,[10][11] and commonly a temperature of 72°C is used with this enzyme. At this step the DNA polymerase synthesizes a new DNA strand complementary to the DNA template strand by adding dNTPs that are complementary to the template in 5' to 3' direction, condensing the 5'-phosphate group of the dNTPs with the 3'-hydroxyl group at the end of the nascent (extending) DNA strand. The extension time depends both on the DNA polymerase used and on the length of the DNA fragment to be amplified. As a rule-of-thumb, at its optimum temperature, the DNA polymerase will polymerize a thousand bases per minute. Under optimum conditions, i.e., if there are no limitations due to limiting substrates or reagents, at each extension step, the amount of DNA target is doubled, leading to exponential (geometric) amplification of the specific DNA fragment.

# Polymerase Chain Reaction

*Final hold*: This step at 4-15°C for an indefinite time may be employed for short-term storage of the reaction.

# Polymerase Chain Reaction

Ethidium bromide-stained PCR products after gel electrophoresis. Two sets of primers were used to amplify a target sequence from three different tissue samples. No amplification is present in sample #1; DNA bands in sample #2 and #3 indicate successful amplification of the target sequence. The gel also shows a positive control, and a DNA ladder containing DNA fragments of defined length for sizing the bands in the experimental PCRs.

To check whether the PCR generated the anticipated DNA fragment (also sometimes referred to as the amplimer or amplicon), agarose gel electrophoresis is employed for size separation of the PCR products. The size(s) of PCR products is determined by comparison with a DNA ladder (a molecular weight marker), which contains DNA fragments of known size, run on the gel alongside the PCR products.

# Lecture 1

## Introduction to Microarrays

# Microarray Data Analysis

Gene chips allow the simultaneous monitoring of the expression level of thousands of genes.  Many statistical and computational methods are used to analyze this data.  These include:

- statistical hypothesis tests for differential expression analysis
- principal component analysis and other methods for visualizing high-dimensional microarray data
- cluster analysis for grouping together genes or samples with similar expression patterns
- hidden Markov models, neural networks and other classifiers for predictively classifying sample expression patters as one of several types (diseased, ie. cancerous, vs. normal)

# What is Microarray Data?

In spite of the ability to allow us to simultaneously monitor the expression of thousands of genes, there are some liabilities with micorarray data. Each micorarray is very expensive, the statistical reproducibility of the data is relatively poor, and there are a lot of genes and complex interactions in the genome.

Microarray data is often arranged in an $n$ x $m$ matrix $\mathbf{M}$ with rows for the n genes and columns for the m biological samples in which gene expression has been monitored. Hence, $m_{ij}$ is the expression level of gene $i$ in sample $j$. A row $\mathbf{e}_i$ is the *gene expression pattern* of gene $i$ over all the samples. A column $\mathbf{s}_j$ is the expression level of all genes in a sample $j$ and is called the *sample expression pattern*.

# Types of Microarrays

É  cDNA microarray

É  Nylon membrane and plastic arrays (by Clontech)

É  Oligonucleotide silicon chips (by Affymetrix)

É  Note:  Each new version of a microarray chip is at least slightly different from the previous version.  This means that the measures are likely to change.  This has to be taken into account when analyzing data.

# cDNA Microarray

É  The expression level $e_{ij}$ of a gene $i$ in sample $j$ is expressed as a log ratio, $\log(r_{ij}/g_i)$, of the log of its actual expression level $r_{ij}$ in this sample over its expression level $g_i$ in a control.

É  When this data is visualized $e_{ij}$ is color coded to a mixture of red ($r_{ij} >> g_i$) and green ($r_{ij} << g_i$) and a mixture in between.

# Nylon Membrane and Plastic Arrays (by Clontech)

É A raw intensity and a background value are measured for each gene.

É The analyst is free to choose the raw intensity or can adjust it by subtracting the background intensity.

# Oligonucleotide Silicon Chips (by Affymetrix)

É These arrays produce a variety of numbers derived from 16-20 pairs of perfect match (PM) and mismatch (MM) probes.

É There are several statistics related to gene expression that can be derived from this data. The most commonly used one is the *average difference* (AVD), which is derived from the differences of PM-MM in the 16-20 probe pairs.

É The next most commonly used method is the *log absolute value* (LAV), which comes from the ratios PM/MM in the probe pairs.

É Note: The Affymetrix gene-chip software has a absent/present call for each gene on a chip. According to Jagota, the method is complex and arbitrary so they usually ignore it.

# For What Do We Use Microarray Data?

É  Genes with similar expression patterns over all samples ó We can compare the expression patterns $e_i$ and $e_{i'}$ of two genes $i$ and $i'$ over all samples.

É  If we use cluster analysis, we can separate the genes into groups of genes with similar expression patterns (trees).

É  This will allow us to find what unknown genes have altered expression in a particular disease by comparing the pattern to genes know to be affiliated with a disease.

É  It can also find genes that fit a certain pattern such as a particular pattern of change with time.

É  It can also characterize broad functional classes of new genes from the known classes of genes with similar expression.

# For What Do We Use Microarray Data?

- É  Genes with unusual expression levels in a sample ó In contrast to standard statistical methods where we ignore outliers, here outliers might have particular importance. Hence, we look for genes whose expression levels are very different from the others.
- É  Genes whose expression levels vary across samples ó We can compare gene expression levels of a particular gene or set of genes in different samples.  This can be used to look compare normal and diseased tissues or diseased tissue before and after treatment.

# For What Do We Use Microarray Data?

É  Samples that have similar expression patterns ó We might want to compare the expression patters of all genes between two samples.  We might cluster the genes into gene with similar expression patterns to help with the comparison.  This can be used to look compare normal and diseased tissues or diseased tissue before and after treatment.

É  Tissues that might be cancerous (diseased) ó We can take the gene expression pattern of sample and compare it to library expression patterns that indicate diseased or not diseased tissue.

# Statistical Methods Can Help

É  Experimental Design ó Since using microarrays is costly and time consuming, we want to design experiments to use the minimal number of micorarrays that will give a statistically significant result.

É  Data Pre-processing ó It is sometimes useful to preprocess the data prior to visualization.  An example of this is the log ratio mentioned earlier.  It is often necessary to rescale data from different microarrays so that they can be compared.  This is due to variation in chip to chip intensity.   Another type of preprocessing is subtracting the mean and dividing by the variance.

# Statistical Methods Can Help

É Data Visualization ó Principle component analysis and multidimensional scaling are two useful techniques for reducing multidimensional data to two and three dimensions.  This allows us to visualize it.

É Cluster Analysis ó By associating genes with similar expression patterns, we might be able to draw conclusions about their functional expression.

É Probability Theory ó We can use statistical modeling and inference to analyze our data.  Probability theory is the basis for these.

# Statistical Methods Can Help

É  Statistical Inference ó This is the formulation and statistical testing of a hypothesis and alternative hypothesis.

É  Classifiers for the Data ó We can construct classes from data, such a diseased vs. non-diseased tissue.  We can build a model (such as a hidden Markov model) that fits know data for the different classes.  This can then be used to classify previously unclassified data.

# Preprocessing Microarray Data

É Before microarray data can be analyzed or stored, a number of procedures or transformations must be applied to it.

É In order to analyze the data correctly, it is important to understand what the transformations might be doing to the data.

# Preprocessing Microarray Data

- Ratioing the data
- Log-tranforming ratioed data
- Alternative to ratioing the data
- Differencing the data
- Scaling data across chips to account for chip-to-chip difference
- Zero-centering a gene on a sample expression pattern
- Weighting the components of a gene or sample expression pattern differently
- Handling missing data
- Variation filtering expression patterns
- Discretizing expression data

# Ratioing the data

É This is the most popular transformation.

É The expression level eij of a gene i in sample j is expressed as a ratio, $(r_{ij}/g_i)$, of its actual expression level $r_{ij}$ in this sample over its expression level gi in a control.

É This tells us the level of under- or over- expression of a gene i in the sample j.

É If the control value $g_i$ is very small, it can make the ratio very big.

É This can skew results incorrectly.

# Log-tranforming ratioed data

É This is also a popular transformation.

É The expression level $e_{ij}$ of a gene i in sample j is expressed as a log ratio, $\log(r_{ij}/g_i)$, of the log of its actual expression level $r_{ij}$ in this sample over its expression level $g_i$ in a control.

É This will suppress outliers caused when the control value $g_i$ is very small.

É However, it creates a new outlier when $r_{ij}$ is very small.

# Alternative to ratioing the data

An alternative that eliminated both of the outlier problems above is

$$\frac{r_{ij}}{r_{ij} + g_i}$$

This gives a value in [0,1] and can be interpreted at the probability of gene i is higher in sample j than in control.

# Differencing the data

É  Another transformation is to difference the data ie. $r_{ij}$ ó $g_i$.

É  This is not really appropriate in our previous context.

É  However, this is used by Affymetrix in a different context.

É  In their data $r_{ij}$ is the strength of the match of the target $i$ to a specific probe $j$ and $g_i$ is he strength of the match of the target $i$ to a control for this probe.

# Scaling data across chips to account for chip-to-chip difference

- É  As mentioned previously, different chips might display different intensities.
- É  When comparing different chips the data might need to be scaled so that they are on the same scale.
- É  Alternatively, they can be normalized so that they are between [0,1] and compared.

# Zero-centering a gene on a sample expression pattern

This in effect the same as subtracting the mean expression pattern.

Suppose that **x** is an expression pattern for a particular gene $g_i$ whose components are log-ratios.  Let  where is the average expression pattern or control.  Then **x** indicates whether the gene $g_i$ is induced or repressed relative to control. (Remember the **x**'s are vectors).

Subtracting the mean expression pattern and dividing by the standard deviation can accomplish this.

$$x' = \frac{x - \bar{x}}{\sigma}$$

(Remember the **x**'s are vectors).

# Weighting the components of a gene or sample expression pattern differently

If we have a matrix of weights $\mathbf{W}=\text{diag}(w_1,\acute{\text{i}}\ ,w_n)$, we can weight the expression patterns by

$$\mathbf{x}_w = \mathbf{Wx}$$

In this way, we can weight the contributions from different genes differently.

# Handling missing data

É Sometimes components of an expression pattern **x** are missing.

É To fix this, the missing values can be replaced by the mean over the non-missing values in **x**.

# Variation filtering expression patterns

É  When we are performing cluster analysis on gene or sample expression patterns, patterns with low variance will all seem sufficiently similar to each other and might form a cluster.

É  This cluster will probably not reflect any interesting result.

# Discretizing expression data

Sometimes we might want to convert gene or sample expression pattern into discrete values. For example, if we have log-ratio, we may want to simply look at whether something is up- or down-regulated. To do this, we can do the following:

$$x_b = (x_{b,i}) \quad where \quad x_{b,i} = \begin{cases} +1 & when \; x_i > 0 \\ -1 & when \; x_i < 0 \\ 0 & when \; x_i = 0 \end{cases}$$

In this case +1 would indicate up-regulation, 0 would indicate no change and ó1 would indicate down-regulation.

# Measuring Dissimilarity of Expression Data

É  We might want to compare two or more gene or sample expression patterns.

É  This might be used to differentiate between diseased and normal cells or finding out the genetic similarity of tissues.

É  To do this we need a *distance metric* or a *dissimilarity measure*.

# Example Distance Metric

Euclidean Distance ó This is the most common distance measure.

$$d(x, y) = \sqrt{\sum_i (x_i - y_i)^2}$$

This should not be used if either

1)    Not all components of the vectors being compared have equal weight.

2)    There is missing data.

Preprocessing the data can often alleviate these problems.

We can also use the normalized Euclidean distance

$$d(x, y) = \frac{\sqrt{\sum_i (x_i - y_i)^2}}{\sqrt{n}}$$

# Example Dissimilarity Measures

É  Maximum Coordinate Difference ó The following computes the maximum absolute distance along a coordinate

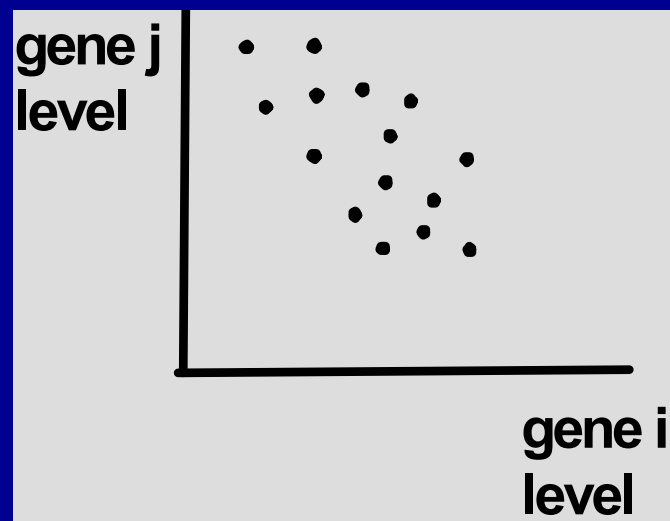$$d(x, y) = \max_i |x_i - y_i|$$

É  Dot Product óThis is a dissimilarity version of the dot product.

$$d(x, y) = -x \circ y$$

# Visualizing Micorarray Data

It is usually easiest to understand data if it can be represented in 2 or 3 dimensions.  For example, a 2-D scatter plot of the expression levels of genes $i$  and  $j$ over a number of samples can show the relationship between these two genes.

# Principal Components Analysis

É In principal components analysis n-dimensional data is converted to d-dimensional data (d<<n) such that the components in the new space are uncorrelated and axis or dimensions of the new space are ordered with respect to the amount of variance they explain.

É The first component explains the most about the data.

É The second component is orthogonal to the first and explains more about the data and so on.

# Principal Components Analysis

Example ó Consider height and weight data for a group of individuals.  This is 2-D data, but there is a correlation between height and weight.  We can use this property to reduce the data to 1D.  PC1 explains most of the data and PC2 explains the rest.

# PCA and Microarrays

Sample Application 1 ó If we want to compare the sample expression patterns from two groups (diseased vs normal, experimental vs control). If we have n genes, the each pattern is a point in n-dimensional space. Suppose we want to see if the sample expression patterns for these two groups cluster by group. We might want to perform PCA analysis and perform cluster analysis at the top three components.

# PCA and Microarrays

Sample Application 2 ó On gene chips (such as the one made by Affymetrix), the same gene occupies multiple cells. In theory, the expression level of all cells with the same gene should be perfectly correlated. However, in practice, this is often not the case due to imperfections in the technology or hybridization of the sequence fragment to other genes in the target.

# PCA and Microarrays

PCA allows us to see how good the correlation among these cells is. To use PCA, we would hybridize k different samples on the same chip. For each sample, the expression levels of a gene x in the n cells is an n-dimensional vector. Hence, there are k points in n-dimensional space. Using PCA, if most of the variance is explained by the first principal component, the effective dimensionality of the data is 1 and there cells are highly correlated.

# PCA Limitations

É Clustering by PCA effectively yields clusters as if the Euclidean distance metric had been used. Hence, it is possible that it might miss clusters.

É The reduction of dimensionality uses all coordinates. If only a few genes out of a thousand differ between two samples (Application 1), clustering by PCA might not yield any meaningful results.

# Cluster Analysis of Microarray Data

Recall that microarray data can be thought of as gene expression patterns or sample expression patterns. These can be each considered to be vectors. The first thing we have to do before applying cluster analysis is to find a distance between the various expression pattern vectors. This is done using similarity/dissimilarity measures such as Euclidean distance, Mahalonobis distance, or linear correlation coefficients. Once a distance matrix is computed, the following clustering algorithms can be used. The clusters formed can differ significantly depending upon the distance measure used.

# Cluster Analysis of Microarray Data

Hierarchical Clustering ó Assume each data point is in a singleton cluster.

Find the two clusters that are closest together. Combine these to form a new cluster.

Compute the distance from all clusters to the new cluster using some form of averaging.

Find the two closest clusters and repeat.

# Cluster Analysis of Microarray Data

k-Means Clustering – An alternate method of clustering called k-means clustering, partitions the data into k clusters and finds cluster means $\mu_i$ for each cluster.  In our case, the means will be vectors also.

Usually, the number of clusters k is fixed in advance.  To choose k something must be know about the data.  There might be a range of possible k values.

To decide which is best, optimization of a quantity that maximizes cluster tightness ie. minimizes distances between points in a cluster.

# Cluster Analysis of Microarray Data

Self-organizing Maps ó This is basically an application of neural networks to microarray data. Assume that there is a 2-dimensional grid of cells and a map from a given set of expression data vectors in $R^n$, ie, there are n nodes in the input layer and a connection neuron from each of these to each cell. Each cell (i, j) gets it own weight from n input neurons. The weight vector $\mu_{ij}$ is the mean of the cluster associated with cell (i, j). Each data vector **d** gets mapped to the cell (i, j) that is closest to **d** using Euclidean distance.In order to train the network, the mean vectors $\mu_{ij}$ for the cells (i, j) must be learned.

# Hidden Markov Models and Microarray Data

We can use Hidden Markov models for pattern recognition in the study of micorarray data. Suppose that we want to consider gene expression data from a tissue sample and want to know if it is control or different from the control (diseased, experimentally altered, responding to drug, etc.). Consider the gene expression data vector as a set of emissions, one for each vector coordinate. Each emission has a value that is defined by some probability distribution function. This can be continuous, or can even discrete. To make it discrete, the data should be preprocessed to indicate, up-regulation, down-regulation, or no significant change.

# Finding Genes Expressed Unusually Different in a Population

The following section addresses the question: Is gene g expressed unusually in the sample?

The first thing to do is to come up with a formal mathematical definition for what unusual is. Assume that the microarray data is log transformed ratio data. If a histogram is constructed of the data, it should yield roughly a normal distribution. Anything that is out near either tail can be considered to be unusually expressed. Note that this can be either a high or low expression level.

# Finding Genes Expressed Unusually Different in a Population

Calculate the Z-score for the data point considered

$$Z = \frac{e_g - \mu}{\sigma}$$

where $e_g$ is the expression level, $\mu$ is the mean and $\sigma$ is the standard deviation. The Z value will give an indication of the how far the data is toward the tail ($\alpha$ - level).

Use statistical inference (hypothesis testing).