



Computational Biology Lecture #11: OMICS: Transcriptomics & Proteomics

Bud Mishra
Professor of Computer Science, Mathematics, & Cell Biology
Nov 28 2005

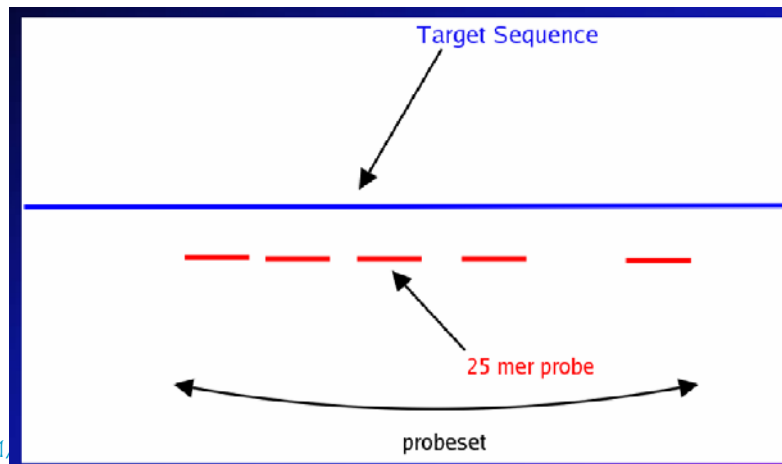
11/28/2005

© Bud Mishra, 2005

L7-1

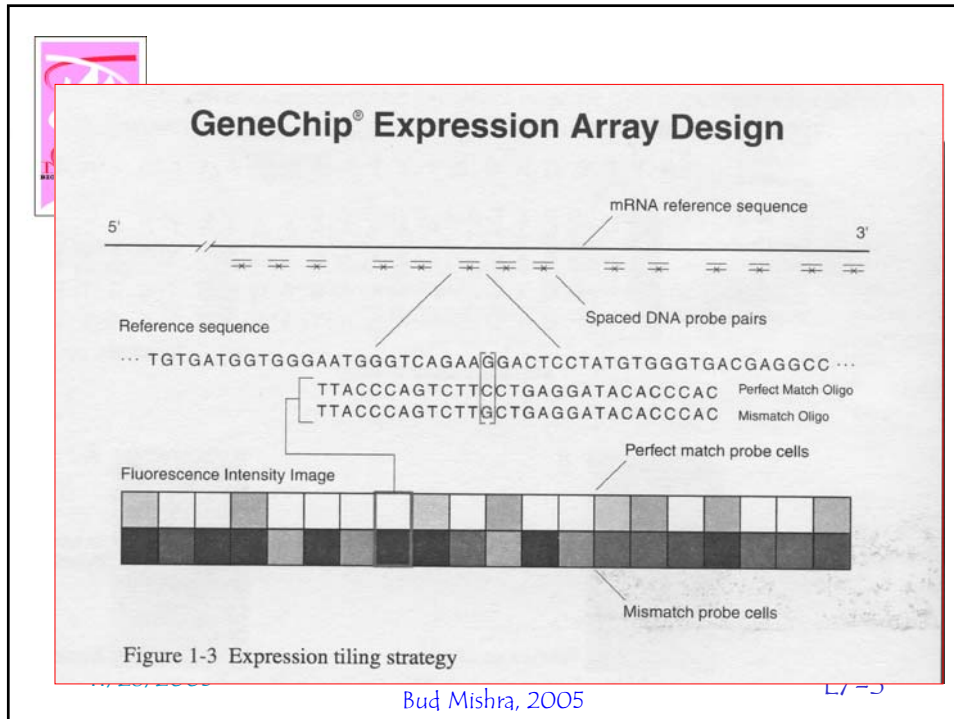



Probes & ProbeSets in Affymetrix Chips



11,

Bud Mishra, 2005

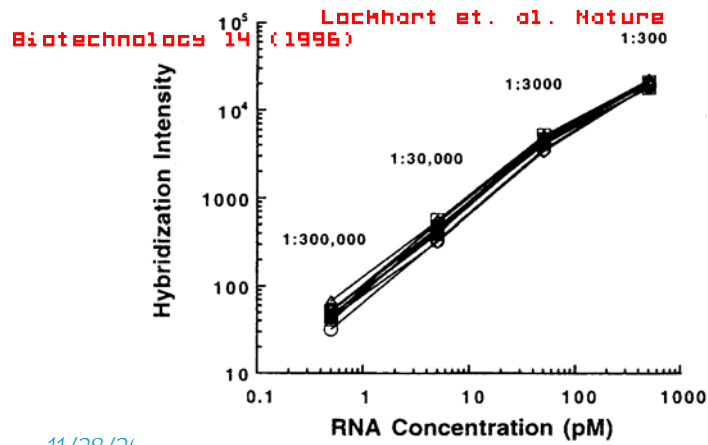
The big picture

- ◇ Summarize 20 PM,MM pairs (probe level data) into one number for each probe set (gene)
 - We call this number an expression measure
 - Affymetrix GeneChip Software has defaults.
- ◇ Does it work? Can it be improved?

11/28/2005 © Bud Mishra, 2005 L7-4



Where is the evidence that it works?



11/28/2005

© Bud Mishra, 2005

L7-5



Comments

- The chips used in Lockhart et. al. contained around 1000 probes per gene
- Current chips contain 11-20 probes per gene
- These are quite different situations
- ◇ We haven't seen a plot like the previous one for current chips

11/28/2005

© Bud Mishra, 2005

L7-6



Data Processing

- ◇ The original GeneChip® software used AvDiff

$$AvDiff = |A|^{-1} \sum_{j \in A} (PM_j - MM_j)$$

- where A is a suitable set of pairs chosen by the software. Here 30%-40% could be <0, which was a major irritant.
- $\log PM_j / MM_j$ was also used in the above.

11/28/2005

© Bud Mishra, 2005

L7-7



Data Processing

- ◇ Li and Wong (dChip) fit the following model to sets of chips

$$PM_{ij} - MM_{ij} = \theta_i \rho_j + \varepsilon_{ij}$$

- where $\varepsilon_{ij} \sim N(0, \sigma^2)$. They consider θ_i to be expression in chip i. Their model is also fitted to PM only, or to both PM and MM. Note that by taking logs, assuming the LHS is ≈ 0 , this is close to an additive model.
- ◇ Efron et al consider $\log PM_j - 0.5 \log MM_j$. It is much less frequently <0.
- ◇ Another summary is the second largest PM, $PM(2)$.

11/28/2005

© Bud Mishra, 2005

L7-8



Data Processing

- ◇ The latest version of GeneChip[®] uses something else, namely

$$\text{Log}\{\text{Signal Intensity}\} = \text{TukeyBiweight}\{\log(\text{PM}_j - \text{MM}_j^*)\}$$

- with MM_j^* a version of MM_j that is never bigger than PM_j .
- ◇ Here TukeyBiweight can be regarded as a kind of robust/resistant mean.

11/28/2005

© Bud Mishra, 2005

L7-9



Tukey Biweight A robust mean

- ◇ Tukey Biweight mean of the dataset
 - Calculate the median (MED) of the data and the mean absolute deviation (MAD)
 - $\text{MED} \pm 5.0 * \text{MAD}$ comprise the limits outside which we consider the data to be outlier. (5.0 is a parameter)
 - $X - \text{MED}$ is used to compute a weight that decays to zero outside the limits of outlier using the bi-square function.
 - Compute the weighted mean to eliminate the outliers.

11/28/2005

© Bud Mishra, 2005

L7-10



Data Processing ...

- ◇ *RMA (Robust Multi-Array Averaging)*
- ◇ **3 Step**
 - Background removal
 - Normalization
 - Summarization

11/28/2005

© Bud Mishra, 2005

L7-11



Data Processing ...

- ◇ For example, dChip Background Removal was PM – MM, MAS-5 was somewhat similar
- ◇ RMA bg.correct uses a signal + plus noise model and uses the posterior mean to detect the signal.
- ◇ Works only on the PM values. The MM values serve in parameter estimation for this and normalization steps

11/28/2005

© Bud Mishra, 2005

L7-12



RMA bg.correct ...

- ◇ Signal: exponentially distributed
- ◇ Observed PM probe value: $X = Y + \text{Noise}$
- ◇ Noise: independent, mean μ , std dev = σ
- ◇ μ , σ , α (for the exponential distribution) are the three parameters to be estimated.
- ◇ Different methods for this.
 - All PM's
 - All MM's
 - Alpha from PM's mu and sigma from MM's

11/28/2005

© Bud Mishra, 2005

L7-13



RMA bg.correct

- ◇ The last might be problematic
 - The MM's have a strong signal components and lead to mis-estimation of μ and σ
 - Result sensitive to mis-estimation of σ .
 - α is usually very small 0.001 – 0.002
 - We are looking at an improper flat prior being approximated by a slowly decaying exponential
 - We can take $\alpha = 0.0$ in the final formula and formulate the estimation problem as estimating from an improper prior by taking limits.

11/28/2005

© Bud Mishra, 2005

L7-14



Normalization

- ◇ Goal:
 - Remove unwanted variability between chips/experiments
 - Combined with scaling to get the values between certain pre-fixed limits (MAS-5)
 - RMA: quantile normalization. Tries to achieve a linear relation between gene expression rank and response.

11/28/2005

© Bud Mishra, 2005

L7-15



Summarization

- ◇ Combining the responses of the probes in the probeset to generate one value for the probeset.
- ◇ A form of mean.
 - Usually robustified
- ◇ RMA: median polish on the logged expression values
- ◇ MAS-5: Tukey Biweight (as explained earlier)
- ◇ dChip: Model based (see earlier)

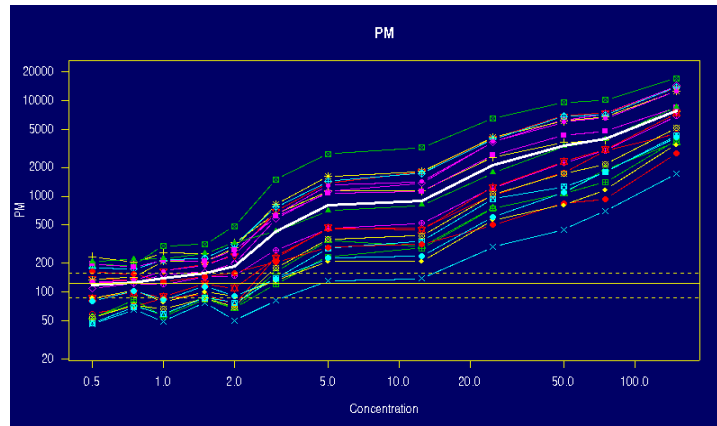
11/28/2005

© Bud Mishra, 2005

L7-16



Probe level data exhibiting parallel behaviour on the log scale



11/28/2005

© Bud Mishra, 2005

L7-17



RMA in summary

- ◇ We *background correct* PM on original scale
- ◇ We carry out *quantile normalization*
- ◇ We take *log₂*
- ◇
- ◇ Under the *additive model*
- ◇ $\log_2 n(\text{PM}_{ij} - \text{BG}) = m + a_i + b_j + \epsilon_{ij}$
- ◇ We estimate chip effects a_i and probe effects b_j using a *robust/resistant method*.

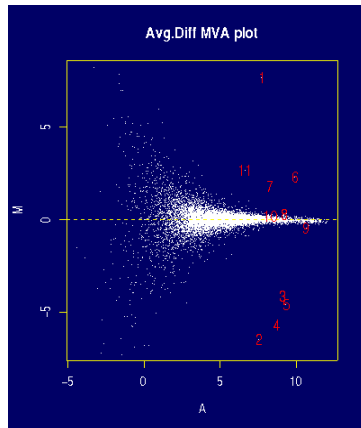
11/28/2005

© Bud Mishra, 2005

L7-18

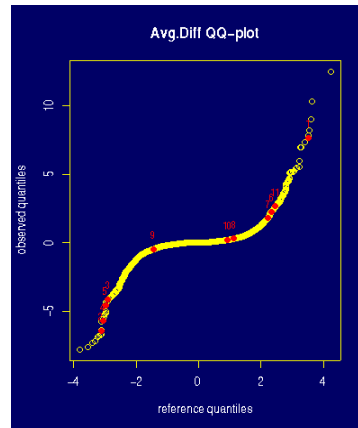


Performance (AvDiff)



11/28/2005

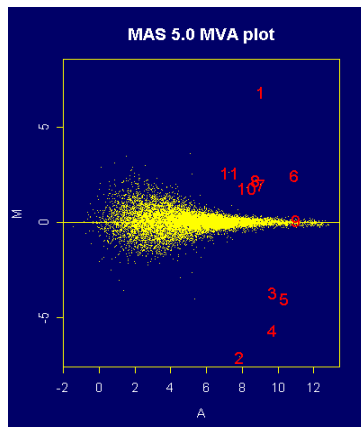
© Bud Mishra, 2005



L7-19

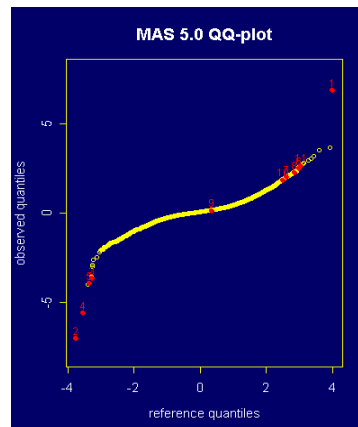


MAS-5



11/28/2005

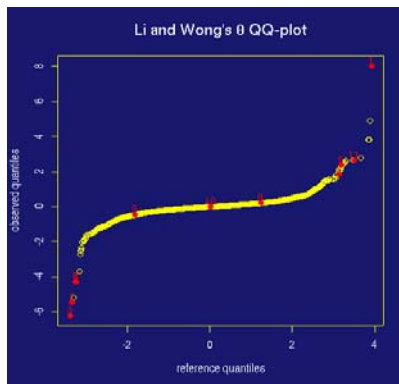
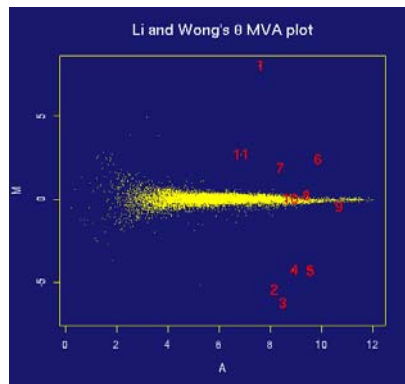
© Bud Mishra, 2005



L7-20



dChip (Li & Wong)



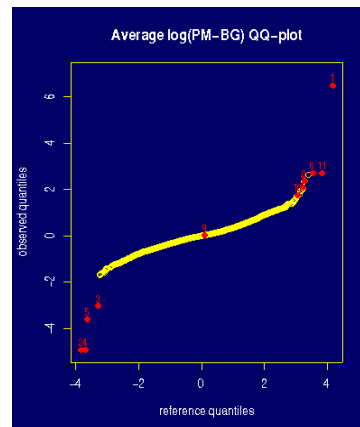
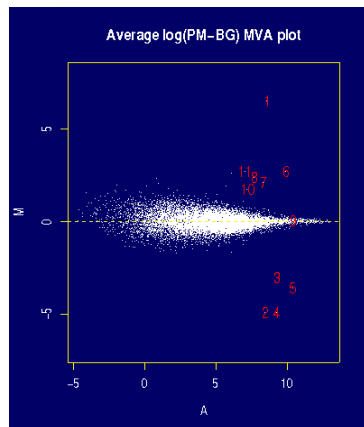
11/28/2005

© Bud Mishra, 2005

L7-21



RMA (no median polish)



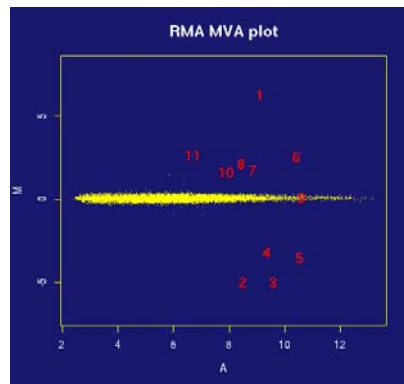
11/28/2005

© Bud Mishra, 2005

L7-22

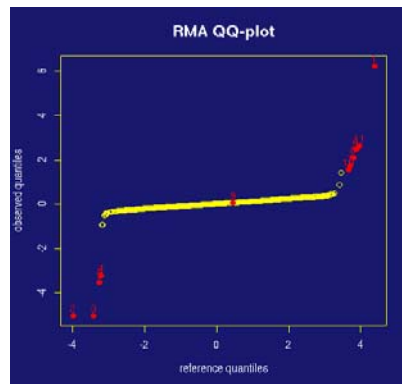


RMA (median Polish)



11/28/2005

© Bud Mishra, 2005



L7-23



RMA background correction

- ◇ It assumes a model $O = S + N$,
- ◇ where S is an exponentially distributed (parameter α) signal, and N is a Gaussian noise with mean μ and standard deviation σ .
- ◇ Various truncation possibilities have also been suggested.
- ◇ The estimator used in all these papers is
- ◇ $E[S | O=o]$.

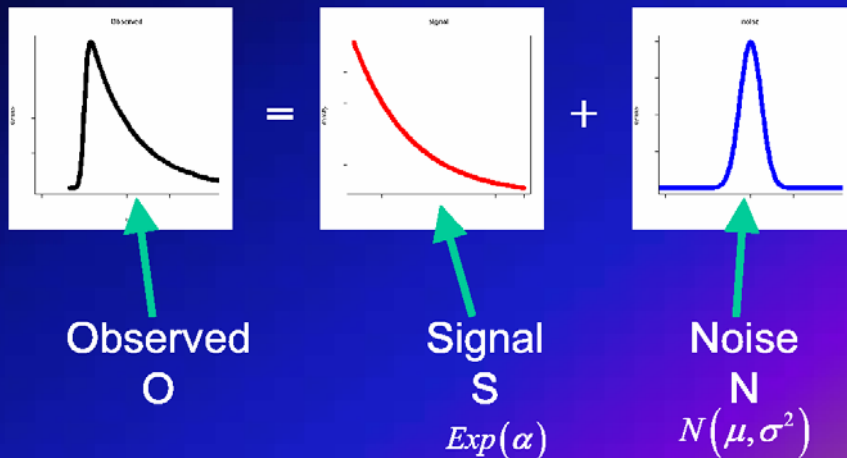
11/28/2005

© Bud Mishra, 2005

L7-24

RMA Background Approach

- Convolution Model



- ◇ $P(S=s) = \alpha \exp(-\alpha s)$, for s, O
- ◇ $P(O = o | S = s) = \phi_{\sigma}(o - s - \mu)$
- ◇ The posterior is computed by
- ◇ $P(S = s | O = o)$
 $= P(S=s) P(O=o|S=s) / \int_{s_0}^{\infty} P(S=s) P(O=o|S=s) ds$
- ◇ Numerator = $\alpha \exp(-\alpha s) \phi_{\sigma}(o - s - \mu)$
 $= \alpha / (\sqrt{2\pi} \sigma) \int_{s_0}^{\infty} \exp(-((s - (o - \mu))^2 + 2\sigma^2 \alpha s) / 2\sigma^2) ds$
- ◇ Denominator = $\int_{s_0}^{\infty}$ Numerator

11/28/2005

© Bud Mishra, 2005

L7-26



$$\begin{aligned} \diamond \text{ Thus, } P(S = s \mid O = o) &= \phi_{\sigma}(s - a) / \Phi_{\sigma}(a) \\ &= \phi_{\sigma}(s - \mu - \sigma^2 \alpha) \\ &+ \sigma \phi((s - \mu - \sigma^2 \alpha) / \sigma) / \Phi((s - \mu - \sigma^2 \alpha) / \sigma) \\ &= s - *BG(\mu, \sigma) \end{aligned}$$

11/28/2005

© Bud Mishra, 2005

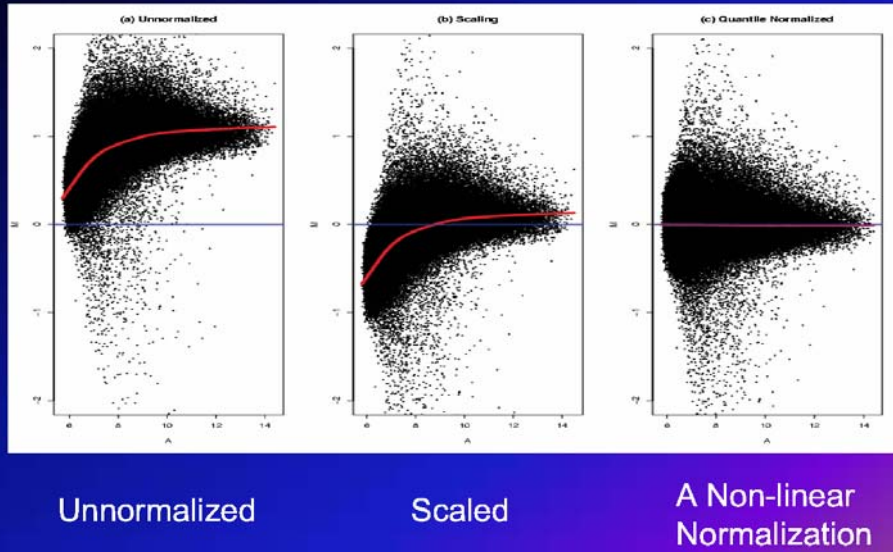
L7-27

Correction is given by

$$E(S \mid O = o) = a + b \frac{\phi\left(\frac{a}{b}\right) - \phi\left(\frac{o - a}{b}\right)}{\Phi\left(\frac{a}{b}\right) + \Phi\left(\frac{o - a}{b}\right) - 1}$$

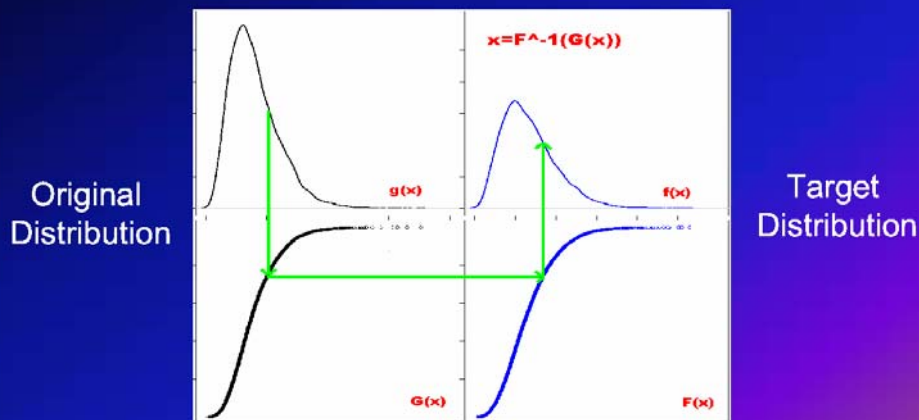
$$a = o - \mu - \sigma^2 \alpha, b = \sigma$$

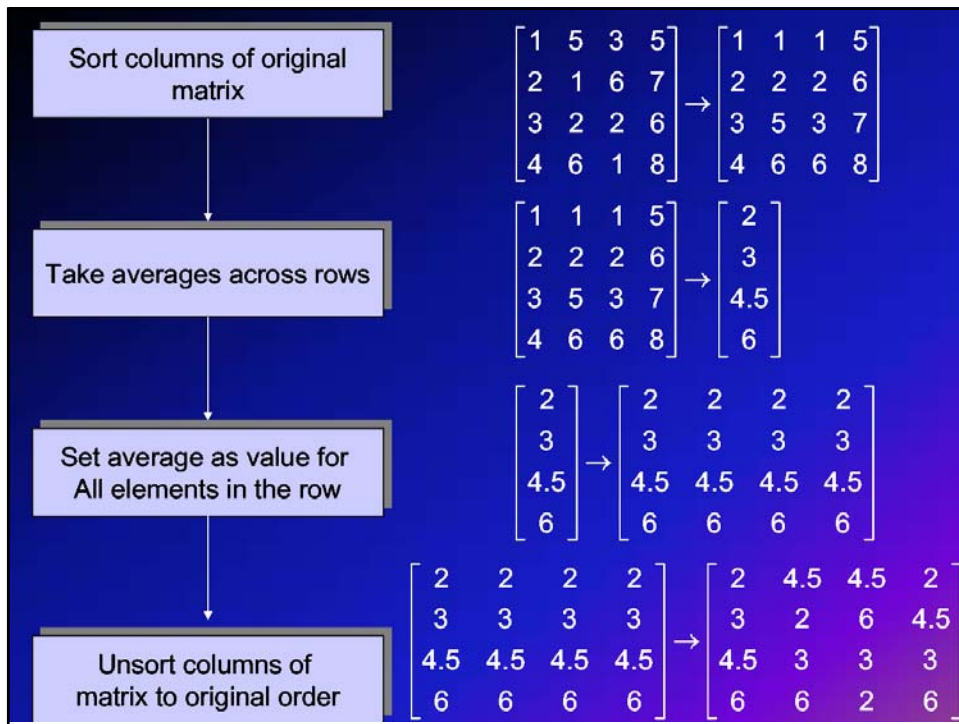
Non-linear normalization needed



Quantile Normalization

- Normalize so that the quantiles of each chip are equal. Simple and fast algorithm. Goal is to give same distribution to each chip.

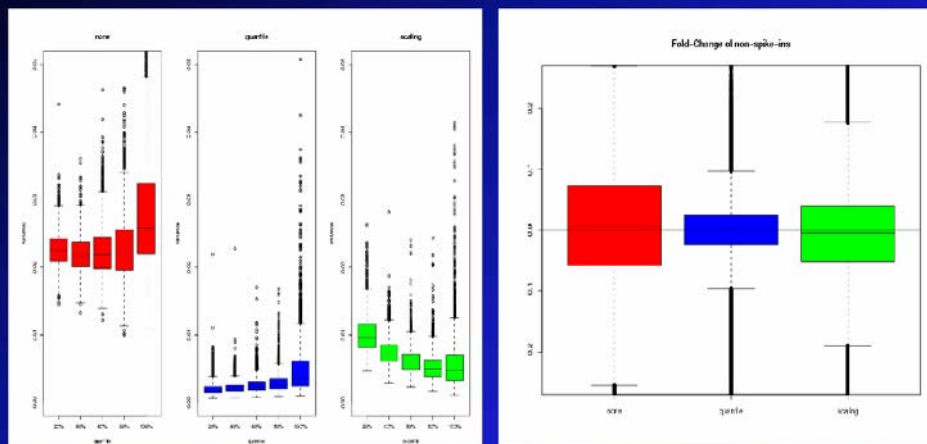




It Reduces Variability

Expression Values

Fold change



Also no serious bias effects. For more see Bolstad et al (2003)

General Probe Level Model

$$y_{ij} = f(\mathbf{X}) + \varepsilon_{ij}$$

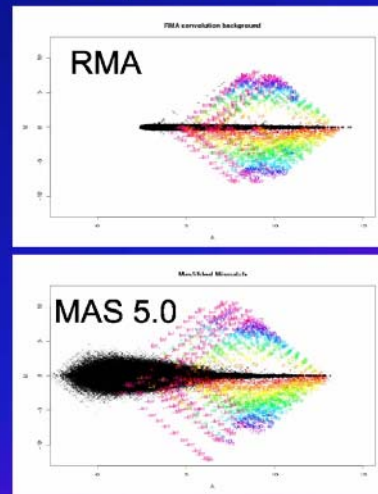
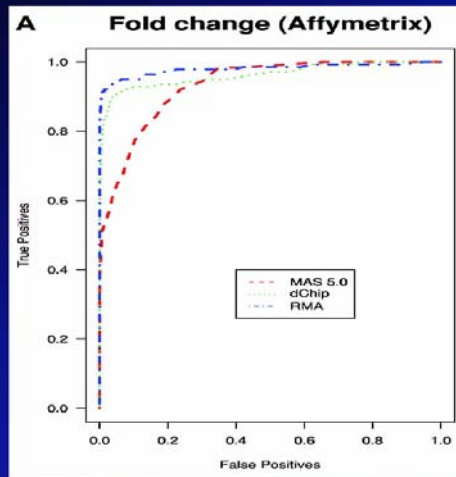
- Where $f(\mathbf{X})$ is function of factor (and possibly covariate) variables (our interest will be in linear functions)
- y_{ij} is a pre-processed probe intensity (usually log scale)
- Assume that $E[\varepsilon_{ij}] = 0$
 $\text{Var}[\varepsilon_{ij}] = \sigma^2$

The Three Steps of RMA

1. Convolution Background
 2. Quantile Normalization
 3. Linear model on the log2 scale fit robustly.
- Software for implementing RMA is in the Bioconductor *affy* package

RMA mostly does well in practice

Detecting Differential Expression Not noisy in low intensities



The RMA model

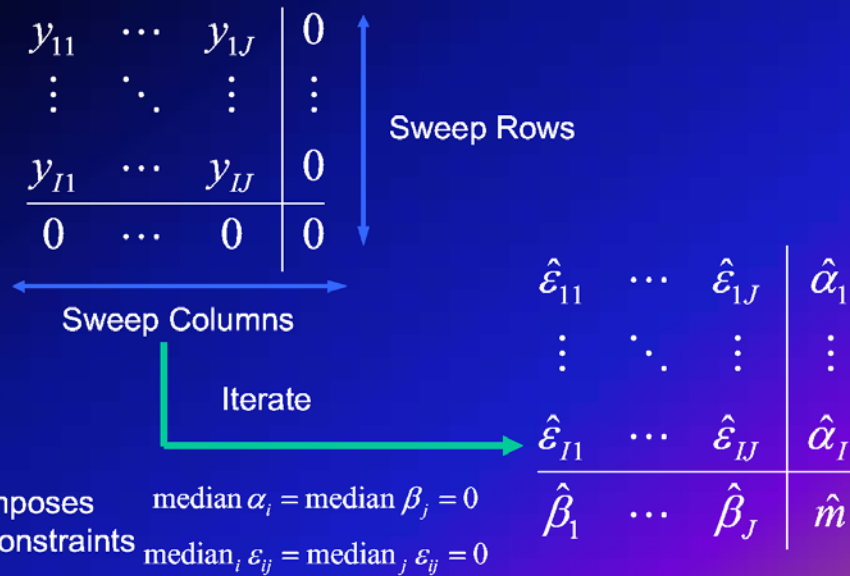
$$y_{ij} = m + \alpha_i + \beta_j + \varepsilon_{ij}$$

where $y_{ij} = \log_2 N(B(PM_{ij}))$

α_i is a probe-effect $i = 1, \dots, I$

β_j is chip-effect ($m + \beta_j$ is log2 gene expression on array j) $j = 1, \dots, J$

Median Polish Algorithm



Median Polish

- Advantages
 - Fast
 - Very robust
- Disadvantages
 - No algorithmic flexibility to fit alternative models
 - No standard error estimates



Within-slide Normalizations

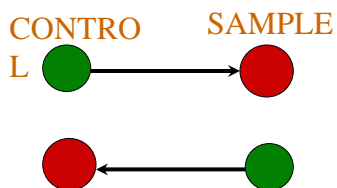
11/28/2005

© Bud Mishra, 2005

L7-39



Dye Bias



- ◇ **Dye Bias**
 - Two-channel microarrays
 - Intensity in one channel is higher than other
 - Dye swapping experiments
- ◇ **Additionally, can be normalized**
 - Take sum intensities for each signal
 - Normalize sums: Assumes most genes regulated at same level

11/28/2005

© Bud Mishra, 2005

L7-40



Spatial Normalization

- ◇ Signal varies according to spot location
 - Particularly, corners
 - ❖ Less hybridization solution
 - ❖ Susceptible to desiccation
 - Chip design
 - ❖ DO NOT cluster genes with similar expression profiles

11/28/2005

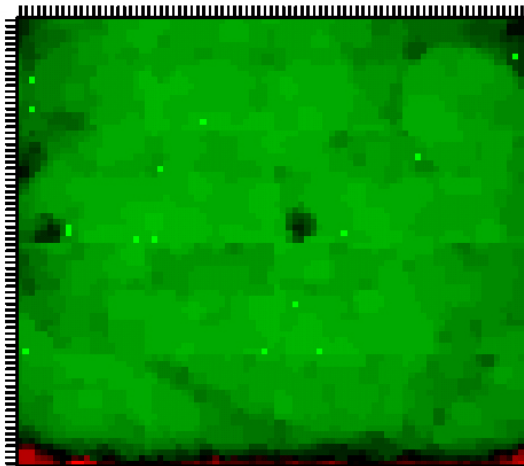
© Bud Mishra, 2005

L7-41



Spatial Bias

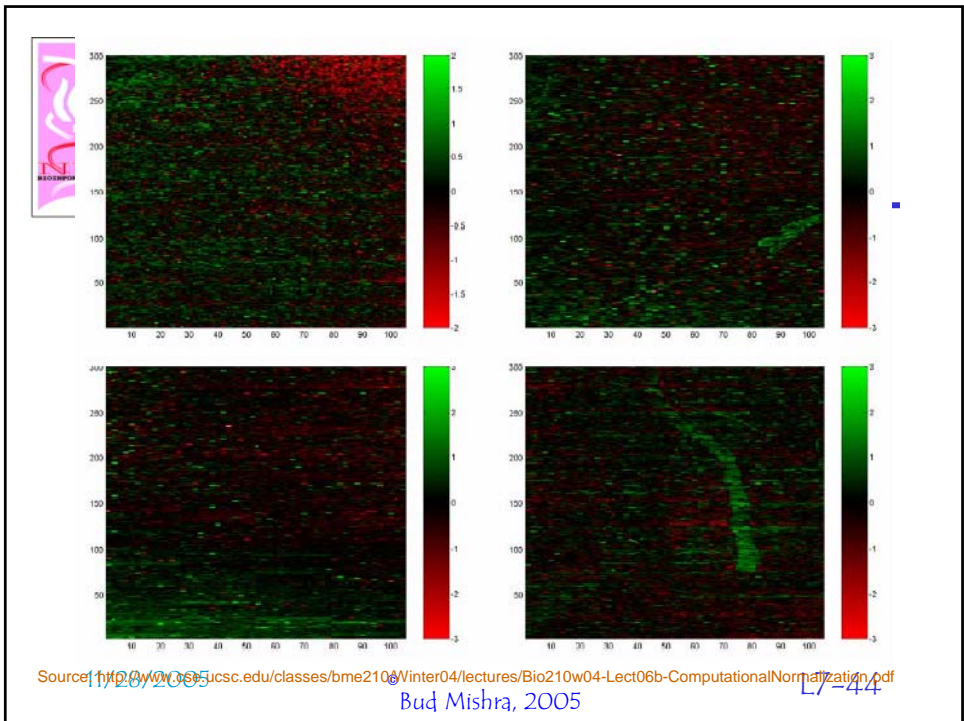
Source: <http://www.csc.fi/oppaat/siru/sirupart11.pdf>



11/28/2005

Bud Mishra, 2005

L7-42





Intensity Dependent Biases

- ◇ Low intensities have much greater variation

11/28/2005

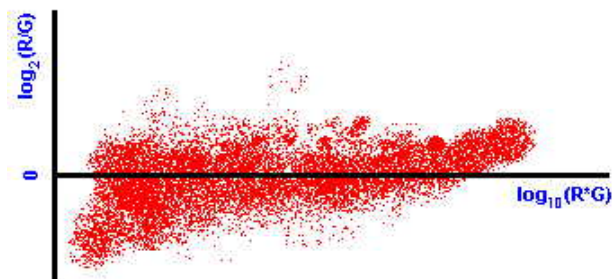
© Bud Mishra, 2005

L7-45



RI Plots

- ◇ Ratio-Intensity
- ◇ R: $\log_2(R/G)$
- ◇ I: $\log_{10}(R^*G)$



http://www.ucl.ac.uk/oncology/MicroCore/HTML_resource/
11/28/2005

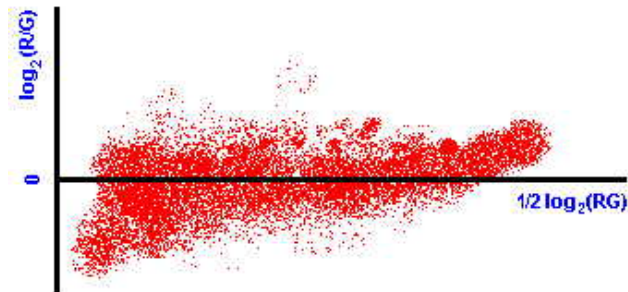
© Bud Mishra, 2005

L7-46



MA Plots

- ◇ $M: \log_2(R/G)$
- ◇ $A: \log_2\text{SQRT}(R*G) = 1/2 \log_2(R*G)$



http://www.ucl.ac.uk/oncology/MicroCore/HTML_resource/MA_plots_popup.htm

11/28/2005

© Bud Mishra, 2005

L7-47



Lowess (Loess) Normalization

- ◇ Locally Weighted Linear Regression
- ◇ Linearises Data

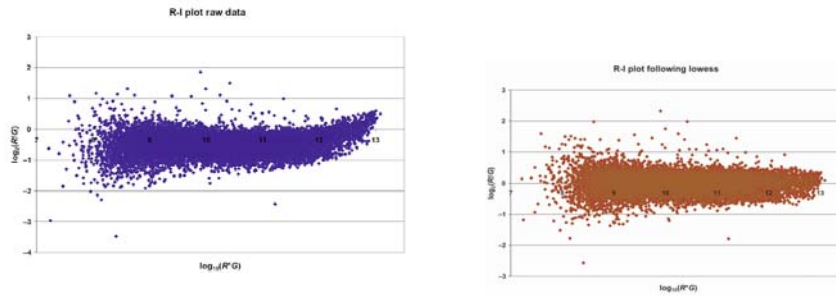
11/28/2005

© Bud Mishra, 2005

L7-48



Loess Normalization



Figures from Quackenbush, 2002

11/28/2005

© Bud Mishra, 2005

L7-49



Cross-slide Normalizations

- ◇ Comparisons between chips needed
- ◇ Slides normalized so comparisons can be made

11/28/2005

© Bud Mishra, 2005

L7-50



Per-Chip Normalization

- ◇ Mean/Median centering – mean/median intensity of every chip brought to same level
- ◇ Total intensity normalization – scaling factor determined by summing intensities
- ◇ Spiked-control, housekeeping normalization

11/28/2005

© Bud Mishra, 2005

L7-51



Differential Expression

- ◇ Crude filter
 - Genes over/underexpressed by a factor of two
 - \log_2 values of 1 and -1
 - Plus: Calculation very easy
 - Minus: Does not consider reliability of data

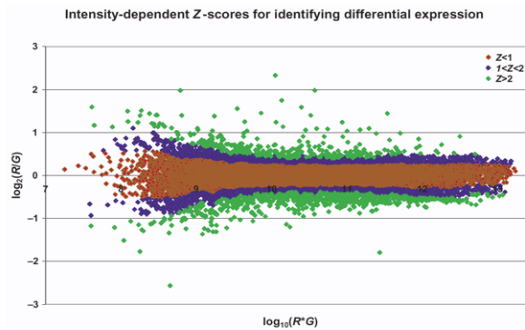
11/28/2005

© Bud Mishra, 2005

L7-52



Localized Z-Scores



Figures from Quackenbush, 2002
11/28/2005 © Bud Mishra, 2005

L7-53



Two sample t-test

- ◇ Calculates probability values sampled from same distribution
 - Considers mean, variance

11/28/2005

© Bud Mishra, 2005

L7-54



Analysis Tools

- ◇ R, Bioconductor
- ◇ S+, ArrayAnalyzer
- ◇ Affymetrix Tools
- ◇ GeneSpring

11/28/2005

© Bud Mishra, 2005

L7-55



MIAME Data

- ◇ **MIAME: "Minimum Information About a Microarray Experiment"**
 - Specifies content, not format
 - Specifies type of data to be published
- ◇ **MIAME Checklist**
 - Experiment design
 - Samples used; extract preparation and labeling
 - Hybridization procedures and parameters
 - Measurement data and specifications
 - Array design

11/28/2005

© Bud Mishra, 2005

L7-56



MAGE-ML

- ◇ MicroArray Gene Expression Markup Language
- ◇ XML based
 - Object model
 - Document exchange format
 - Software toolkits

11/28/2005

© Bud Mishra, 2005

L7-57



Repositories

- ◇ ArrayExpress (EBI)
- ◇ Gene Expression Omnibus (NCBI)
- ◇ Stanford Microarray Database (SMD)

- ◇ Microarray databases
 - AMAD: Another Microarray Database
 - LONGHORN:
 - MIDAS:
 - BASE:

11/28/2005

© Bud Mishra, 2005

L7-58



Proteomics

11/28/2005

© Bud Mishra, 2005

L7-59



Beyond Genomics

- ◇ **Human Genome**
 - 30,000 to 60,000 genes
- ◇ **Human Proteome**
 - 300,000 to 1,200,000 protein variants
- ◇ **Human Metabolome**
 - Metabolic products of the organism (lipids, carbohydrates, amino acids, peptides, prostaglandins, etc)

11/28/2005

© Bud Mishra, 2005

L7-60



Proteome

- ◇ **Proteome: The entire protein complement in a given cell, tissue or organism.**
 - Protein Activities
 - 3D Structure
 - Modifications and Localization
 - Protein-Protein Interaction: Proteins in Complexes
 - Protein Profile: Global patterns of protein content and activity (particularly in response to a disease state.)
 - Understanding system-level cellular behavior

11/28/2005

© Bud Mishra, 2005

L7-61



Technology & Databases

- ◇ **Identify proteins and protein complexes in biological samples comprehensively and quantitatively with both high sensitivity and fidelity.**
 - Develop new diagnostic markers
 - Identification of new drug-target
- ◇ **HUPO (Human Proteome Organization)**
 - Coordinating proteomics projects worldwide.

11/28/2005

© Bud Mishra, 2005

L7-62



Integration

- ◇ **Complementary to other functional genomic approaches:**
 - Micro-array based expression profiles
 - Systematic phenotypic profiles at the cell and organism level
 - Systematic genetics
 - Small-molecule-based arrays

11/28/2005

© Bud Mishra, 2005

L7-63



Applications of Proteomics

- ◇ **Protein Mining**
 - Catalog all the proteins present in a tissue, cell, organelle, etc.
- ◇ **Differential Expression Profiling**
 - Identification of proteins in a sample as a function of a particular state: differentiation, stage of development, disease state, response to drug or stimulus
- ◇ **Network Mapping**
 - Identification of proteins in functional networks: biosynthetic pathways, signal transduction pathways, multiprotein complexes
- ◇ **Mapping Protein Modifications**
 - Characterization of posttranslational modifications: phosphorylation, glycosylation, oxidation, etc.

11/28/2005

© Bud Mishra, 2005

L7-64



Challenges of Proteomics

- ◇ Limited and Variable Sample Material
- ◇ Sample Degradation
- ◇ Vast Dynamic Range
 - (more than 10^6 -fold for protein abundance)
- ◇ Post-translational Modifications
- ◇ Unlimited tissue, developmental and temporal specificity
- ◇ Disease and drug perturbation.

11/28/2005

© Bud Mishra, 2005

L7-65



Proteomics Technologies

- ◇ Development of genome and protein sequence databases
 - Bioinformatics and Data mining software
- ◇ Development of mass spectrometry instrumentation suitable to analyze biomolecules
 - Protein mass, Peptide mass, Peptide sequence
- ◇ Development of analytical protein separation technology
 - IEF, 2D-SDS-PAGE, HPLC, Capillary Electrophoresis, Affinity Chromatography

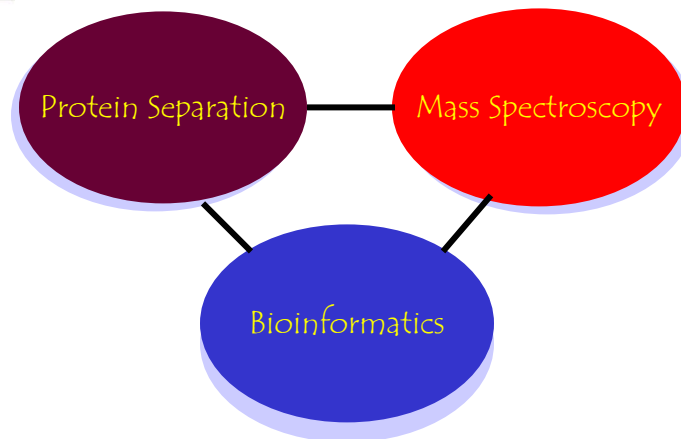
11/28/2005

© Bud Mishra, 2005

L7-66



Components of Proteomics



11/28/2005

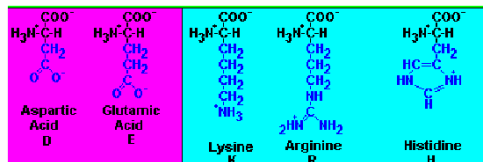
© Bud Mishra, 2005

L7-67



2 D Electrophoresis

- ◇ **Property of proteins**
 - Some amino acids are acidic / basic (donate / accept H⁺)
 - Collection of amino acids in protein determines its pI value
 - ❖ pI = pH at which molecular charge = zero
- ◇ **2D electrophoresis**
 - Separate proteins according to both pI & molecular weight



11/28/2005

© Bud Mishra, 2005

L7-68



2 D Electrophoresis

- ◇ Method
- ◇ 1. Extract & prepare protein sample in solution
- ◇ 2. Separate proteins (in each dimension)
 - I. Based on pH
 - ❖ Using isoelectric focusing (IEF)
 - ❖ Using immobilized pH gradient (IPG) strips
 - II. Based on molecular weight (size)
 - ❖ Using gel electrophoresis

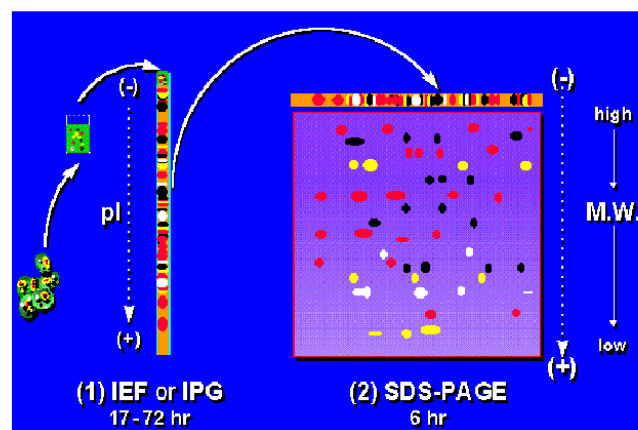
11/28/2005

© Bud Mishra, 2005

L7-69



2 D Electrophoresis



11/28/2005

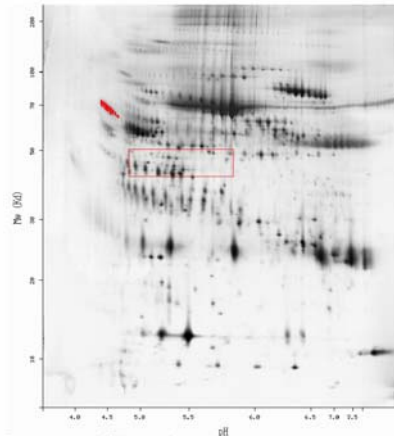
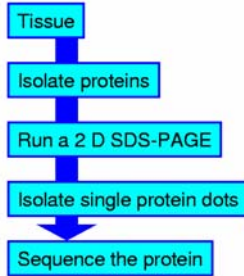
© Bud Mishra, 2005

L7-70



Proteomic Technology

Sequence the proteome



11/28/2005

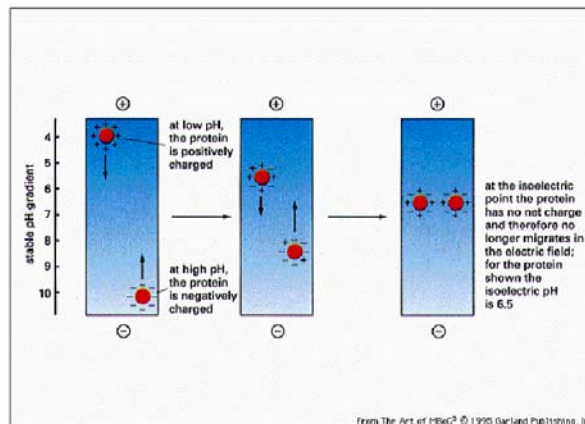
© Bud Mishra, 2005

L/-/1



2-D SDS PAGE

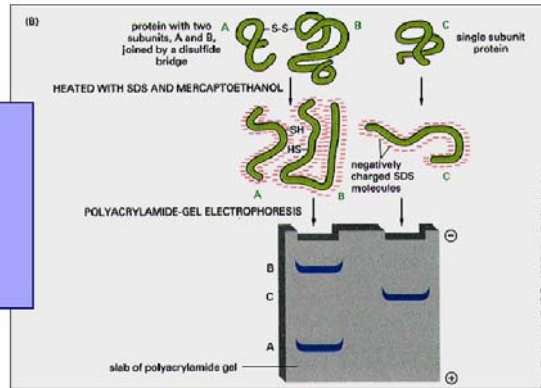
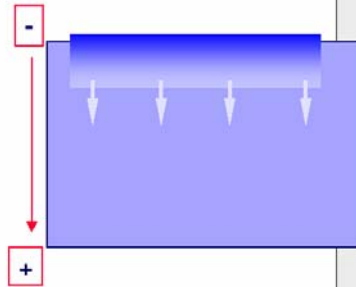
1. Step:



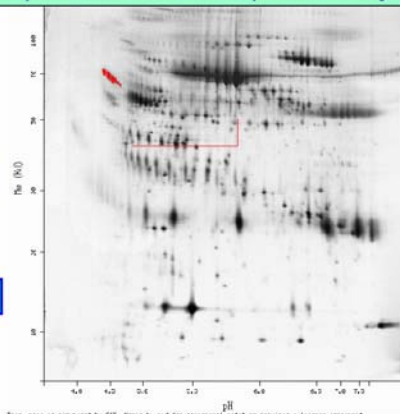
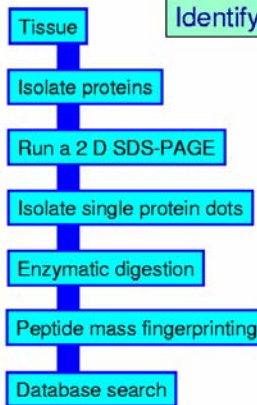


2D SDS PAGE

2. Step:



Integrating with Mass-Spec



11/28/20



Mass Spectrometry

♦ Method

- 1. Excise individual dots from 2D electrophoresis
- 2. Digest protein into fragments with enzyme (e.g., trypsin)
- 3. Ionize protein fragments (without breaking)
 - ❖ Matrix Assisted Laser Desorption Ionization (MALDI)
 - ❖ Electrospray Ionization (ESI)
- 4. Accelerate through mass spectrometer
- 5. Produces peptide mass fingerprint

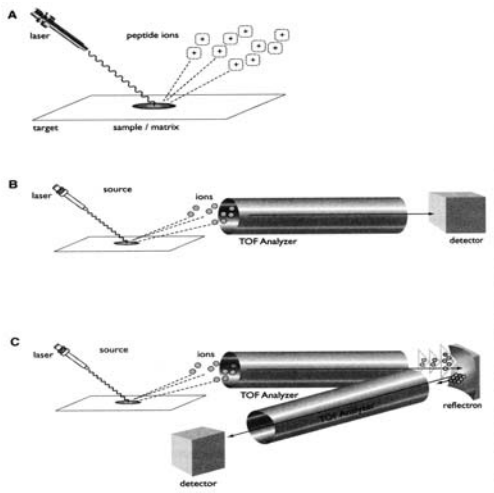
11/28/2005

© Bud Mishra, 2005

L7-75



Principles of MALDI-TOF Mass Spectroscopy



11/28/2005

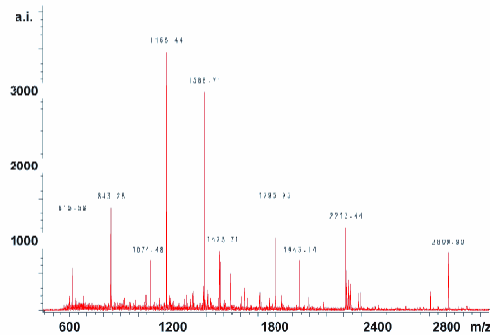
© Bud Mishra, 2005

L7-76



Fingerprint

◆ Peptide mass fingerprint



11/28/2005

© Bud Mishra, 2005

L7-77



Identifying peptide mass fingerprint

- ◆ **Compare with Fingerprint for actual protein in database**
 - Predicted fingerprint for predicted / hypothetical protein (Precompute for efficiency)
 - May fail to distinguish Post-translation modifications to protein
- ◆ **Protein databases / web servers (e.g., SWISS-2D PAGE)**
 - For each protein, record its
 - (1) Protein pI, molecular weight, peptide mass fingerprint...
 - (2) Experimentally determined location in 2D gel

11/28/2005

© Bud Mishra, 2005

L7-78



Peptide Chips

- ◇ Protein-protein interaction,
- ◇ Unravel signal transduction pathways,
- ◇ Perform multi-parameter diagnosis,
- ◇ Study individual immunological repertoires
 - e.g. autoimmune reactions.

11/28/2005

© Bud Mishra, 2005

L7-79



Peptide Chips

- ◇ **Goal**
 - High-throughput analysis of protein expression / interaction
 - Adapt approach similar to DNA microarrays
 - Improves on speed vs. 2D electrophoresis
- ◇ **Approach**
 - No equivalent of hybridization for proteins
 - Exploit other biochemical binding reactions
 - ❖ Antibody-antigen
 - ❖ Receptor-ligand
 - ❖ DNA-protein...

11/28/2005

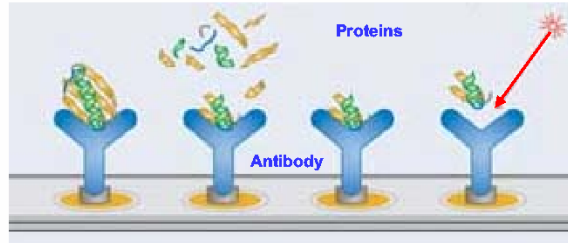
© Bud Mishra, 2005

L7-80



Ciphergen

Ciphergen Antibody Capture Protein Chip



- 1) Protein with antigen bound to antibody probe
- 2) Remainder of protein digested enzyme, leaving peptide antigen
- 3) Wash away protein fragments
- 4) SELDI laser ionizes & desorbs epitope binding peptide, sends to mass spectrometer

11/28/2005

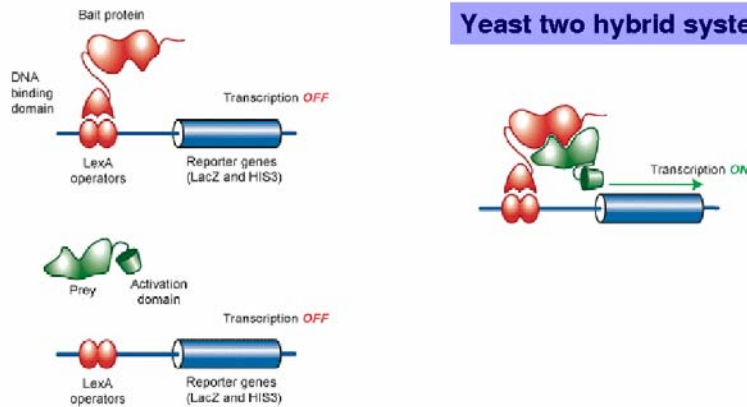
© Bud Mishra, 2005

L7-81



Protein-Protein Interaction

Yeast two hybrid system



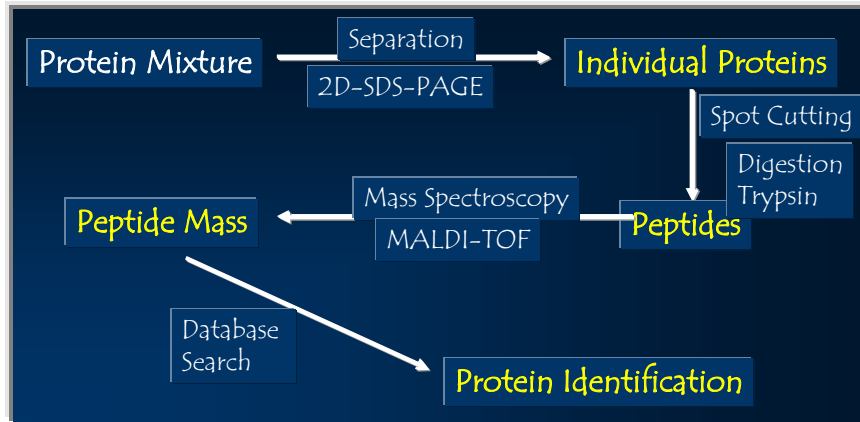
11/28/2005

© Bud Mishra, 2005

L7-82



Basic Proteomic Analysis Scheme



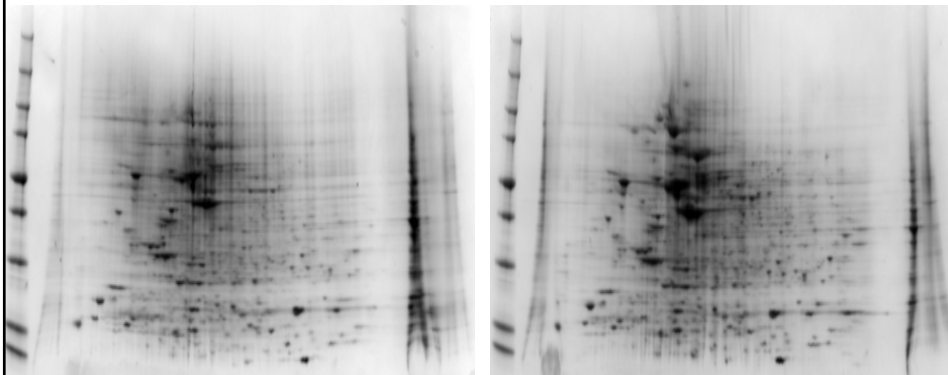
11/28/2005

© Bud Mishra, 2005

L7-83



2D-SDS-PAGE of 2 Types of Cells



Cell Type A

11/28/2005

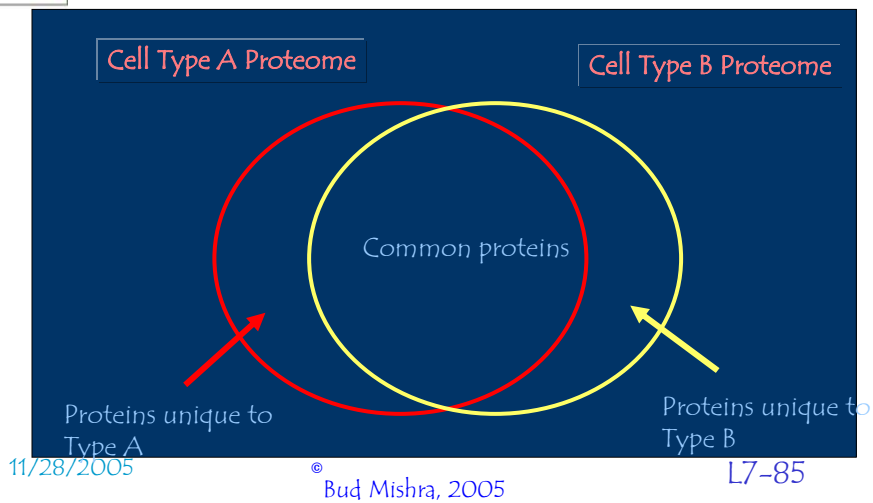
Cell Type B

© Bud Mishra, 2005

L7-84

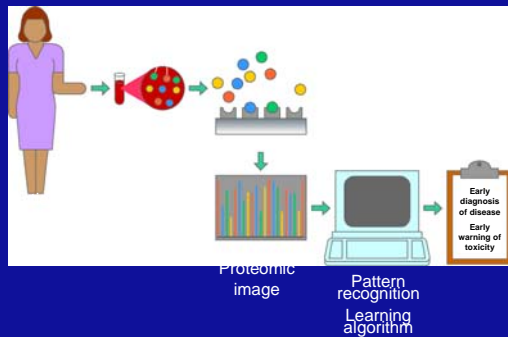


Differential Expression



Clinical Diagnostics Proteomics: Protein Profiling

Serum Protein Pattern Diagnostics



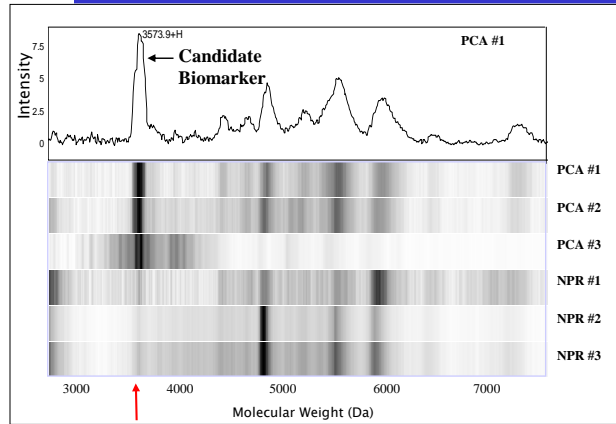
11/28/2005

© Bud Mishra, 2005

L7-86



Protein Profiles of (3 patients)



11/28/2005

© Bud Mishra, 2005

L7-87



To be continued...

...

11/28/2005

© Bud Mishra, 2005

L7-88