## Slide 1

# Mass Spec and MicroArrays

## Applications in Proteomics and Systems Biology



HUPO 6th ANNUAL WORLD CONGRESS, SEOUL 2007
Oct. 6th to Oct.10th, 2007  COEX, Seoul, Korea

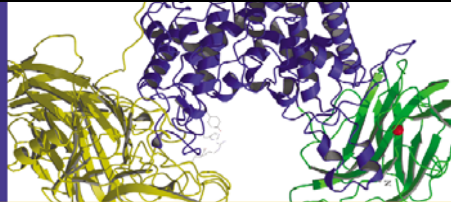Proteomics: *From Technology Development to Biomarker Applications*

Co-sponsored by HUPO, AOHUPO and KHUPO

**CH370 - Hackert**

## Slide 2



**HUPO**
Human Proteome Organisation

Long Beach Convention Center, California, USA

Saturday October 28th through Wednesday November 1st, 2006

**HUPO 5TH ANNUAL WORLD CONGRESS, LONG BEACH 2006**
**TRANSLATING PROTEOMICS FROM BENCH TO BEDSIDE**

Professor Peipei Ping, Congress Co-Chair

Professor John Yates, Congress Co-Chair

## Slide 3

# Report

Practical Proteomics 1-2/2006

**Proteomics Education, an Important Challenge for the Scientific Community: Report on the Activities of the EuPA Education Committee**

EuPA Tutorial Program (preliminary draft)
**Fundamentals and Core Techniques**

**European Proteomics Association (EuPA)**



## Slide 4

**Genomics**

**Proteomics**

**Interactomics**

**Systems Biology** –

None of these fields of research would be possible without **Bioinformatics,** which would not be possible with lots of **computing power!**



THE PENGUIN CURES

HEADACHES

1

## Mass Spec and MicroArrays / Applications

**Genome** – the genome of an organism is its whole hereditary information encoded in its DNA (or, RNA for some viruses) and includes both the coding (genes) and non-coding sequences of the DNA.

**Proteome** – Proteomics is often considered the next step in the study of biological systems, after genomics. It is much more complicated than genomics, mostly because while an organism's genome is rather constant, a proteome differs from cell to cell and constantly changes through its biochemical interactions with the genome and the environment.

**Interactome** – whole set of molecular interactions in cells, in the context of proteomics, it refers to protein-protein interaction network(PPI), or protein network (PN).

**Systems Biology** - seeks to understand how biological systems function. By studying the relationships and interactions between various parts of a biological system (e.g. metabolic pathways, organelles, cells, physiological systems, organisms etc.), it is hoped that eventually a model of the whole system can be developed.
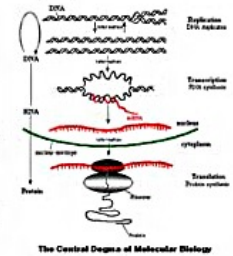
---

### The Proteome

All an organism's cells carry the same Genome, and it is Static. Genomes do not describe function. They are a parts list.

Different cells express different proteins. The type and quantity of this expression changes.

The Proteome is Dynamic. It is the total of all proteins expressed by a particular cell at a given time, under specific conditions.



The Central Dogma of Molecular Biology

A Proteome cannot be studied the way a Genome is sequenced. There has to be a specific biological question behind an experiment. The questions may be either very broad or strictly defined.

---

# Mass spectrometry-based proteomics

Ruedi Aebersold* & Matthias Mann†

*Institute for Systems Biology, 1441 North 34th Street, Seattle, Washington 98103-8904, USA (e-mail: raebersold@systemsbiology.org)
†Center for Experimental BioInformatics (CEBI), Department of Biochemistry and Molecular Biology, University of Southern Denmark, Campusvej 55, DK-5230 Odense M, Denmark (e-mail: mann@bmb.sdu.dk)

Recent successes illustrate the role of mass spectrometry-based proteomics as an indispensable tool for molecular and cellular biology and for the emerging field of systems biology. These include the study of protein–protein interactions via affinity-based isolations on a small and proteome-wide scale, the mapping of numerous organelles, the concurrent description of the malaria parasite genome and proteome, and the generation of quantitative protein profiles from diverse species. The ability of mass spectrometry to identify and, increasingly, to precisely quantify thousands of proteins from complex samples can be expected to impact broadly on biology and medicine.

Note: HT Proteomics is restricted to those species where a sequence database exists!

---

## Generic Mass Spectrometry-based Proteomics

(1) Sample fractionation   SDS–PAGE   Excised proteins   (2) Trypsin digestion (+)   Peptide mixture



(3) Peptide chromatography and ESI   q1   q2

**(4) MS**

Intensity (arbitrary units)

400

200

0

516.27 (2+)

400 600 800

*m/z*

**(5) MS/MS**

200

100

0

LLEAAAQSTK
516.27 (2+)

a2

b2

y3

y4

y5 y6

y7 y8

y9

S Q A A E L L

200 600 1000

*m/z*

---

Ion source — Mass analyser — Detector

Pulsed laser

Liquid chromatography

Nozzle  Sampling cone

Spray needle

Ions

Electrospray ionization (ESI)

Sample plate

Ions

Extraction grid

Matrix-assisted laser desorption/ionization (MALDI)

**a** Reflector time-of-flight (TOF)

Pulsed laser

Sample plate   TOF

**b** Time-of-flight time-of-flight (TOF-TOF)

TOF₁   TOF₂

Reflector

Collision cell

**c** Triple quadrupole or linear ion trap

$Q_1$   $q_2$   $Q_3$

**d** Quadrupole time-of-flight

$Q_1$   $q_2$   TOF

Reflector

**e** Ion trap

**f** Fourier transform ion cyclotron resonance mass spectrometer (FT-MS)

$Q_1$

Super conducting magnet

---

**Differential Expression Proteomics**

CONTROL   DISEASE

2D GELS

IMAGE COMPARISON
QUANTITATION
SPOT PICKING

**DIGESTION OF SELECTED GEL SPOTS**

MALDI-TOF/MS MASS MAPPING
DATABASE SEARCHES
PROTEIN IDENTIFICATION

DATA-DEPENDENT LC/MS/MS
DATABASE SEARCHES
PROTEIN IDENTIFICATION

---

Two Dimensional Gel Electrophoresis

Isoelectric focusing is performed on precast gel strips using commercial instruments. Many pH ranges are available. Multiple strips can be run in parallel.

An immobilized pH gradient is created in a polyacrylamide gel strip by incorporating a gradient of acidic and basic buffering groups when the gel is cast.

Resolution is determined by the slope of the pH gradient and the field strength.

Loading capacity depends on gel size and thickness.

In 2D IEF/PAGE, the gel strip from IEF is loaded into a single large well.

Fig. 1. Principle of 2-D electrophoresis. A, pro-B lymphoma cell extract (1 mg) was separated by IEF on a ReadyStrip pH 5–8 IPG strip, and stained with Bio Safe Coomassie stain. B, Equilibrated strip was run in the second dimension by SDS-PAGE (12% acrylamide). The gel was stained with Coomassie Blue.
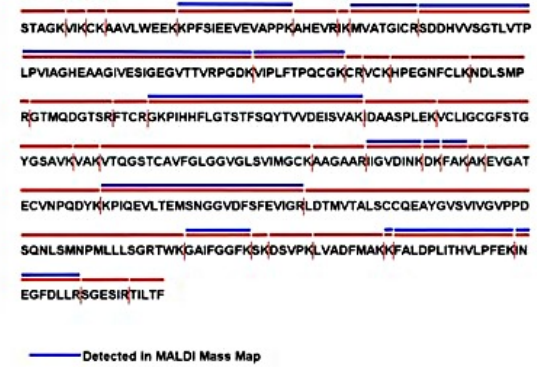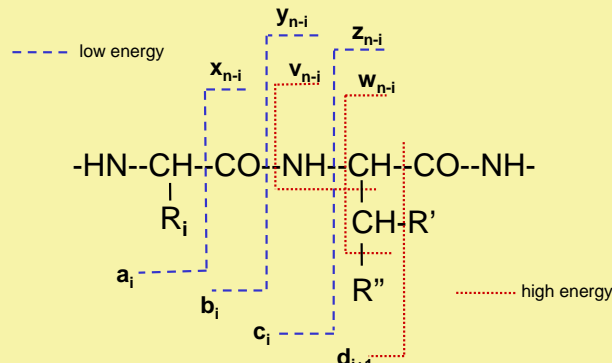
Figure from BioRad Product Literature

3

With the new genomic data bases of model species, such as *Esherichia coli, Saccharomyces cerevisae,* mouse, and human, the sequences of many/most proteins of biological interest will in principle be known, and the problem of characterizing a protein primary structure will be reduced to identifying it in the data base.

Within the past few years research groups have demonstrated how MS can be used for identification of proteins in sequence data bases. One approach is to **cleave the protein** with a **sequence-specific proteolytic enzyme**, **measure molecular weight** values for the resulting peptide mixture by mass spectrometry, and **search a sequence data base for proteins that should yield these values**. **Search algorithms** can utilize low resolution tandem mass spectra of selected peptides (<3 kDa) from the protein degradation. Yates and coworkers compared the MS/MS sequence data to the sequences predicted for each of the peptides that would be generated from each protein in the data base. **In the PEPTIDESEARCH sequence tag approach of Mann and Wilm, a partial sequence of 2–3 amino acids is assigned from the fragment mass differences in the MS/MS spectrum.** This partial sequence and its mass distance from each end of the peptide (based on the masses of the fragment and molecular ions) are used for the data base search. Often, **a single sequence tag retrieved only the correct protein from the data base**.



Tryptic Digest of ADH: Expected Peptides vs. Those Detected

STAGKVIKCKAAVLWEEKKPFSIEEVEVAPPKAHEVRIKMVATGICRSDDHVVSGTLVTP

LPVIAGHEAAGIVESIGEGVTTVRPGDKVIPLFTPQCGKCRVCKHPEGNFCLKNDLSMP

RGTMQDGTSRFTCRGKPIHHFLGTSTFSQYTVVDEISVAKIDAASPLEKVCLIGCGFSTG

YGSAVKVAKVTQGSTCAVFGLGGVGLSVIMGCKAAGAARIIGVDINKDKFAKAKEVGAT

ECVNPQDYKKPIQEVLTEMSNGGVDFSFEVIGRLDTMVTALSCCQEAYGVSVIVGVPPD

SQNLSMNPMLLLSGRTWKGAIFGGFKSKDSVPKLVADFMAKKFALDPLITHVLPFEKIN
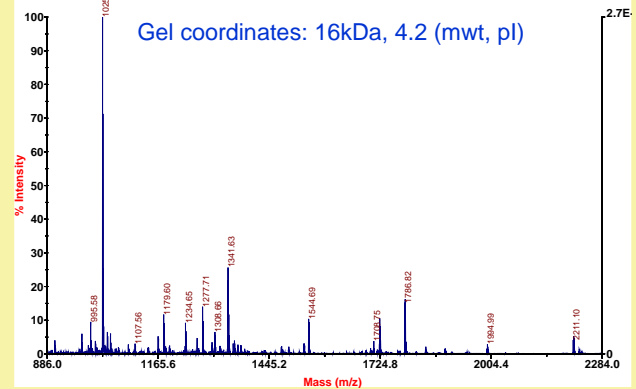
EGFDLLRSGESIRTILTF

——— Detected in MALDI Mass Map

## Cleavages Observed in MS/MS of Peptides



**CID (Collision InDuced) Spectra** – adds **sequence data** to **mass mapping** for improved database identification!

Peptide mass fingerprint of Spot A

Gel coordinates: 16kDa, 4.2 (mwt, pI)



4

MS-Fit (by Peter Baker and Karl Clauser) Instructions

Mass accuracy tolerance = 15 ppm

This means that the mass is within 0.015 Da at m/z 1000

MS-Fit Search Results

## Stable-isotope Protein Labelling for Proteomics

Metabolic stable-isotope labelling

Isotope tagging by chemical reaction

Stable-isotope incorporation via enzyme reaction

Protein labelling

Digest

Label

Digest

Digest

Data collection

Mass spectrometry

Data analysis

Intensity

m/z

Light
Heavy

Organellar
Proteomics:
Combined
MS and
Imaging
Methods



**a**

**b**

YFP-PSP1 actinomycin D — YFP-PSP1 untreated

**c**

Others 12%
Chaperones 6%
Dead box protein 5%
RNA-modifying enzymes and related proteins 8%
Other translation factors 3%
Ribosomal proteins 13%
Novel 32%
Nucleotide-binding and nucleic acid-binding proteins 21%

---

Summary of the functions of various proteins identified in specific tissues of *M. truncatula*.



---

# A Mammalian Organelle Map by Protein Correlation Profiling

Leonard J. Foster,[1,2] Carmen L. de Hoog,[1,2] Yanling Zhang,[3,4] Yong Zhang,[3,4] Xiaohui Xie,[5] Vamsi K. Mootha,[5,6] and Matthias Mann[1,3,*]

[1] Center for Experimental BioInformatics (CEBI), Department of Biochemistry and Molecular Biology, University of Southern Denmark, Campusvej 55, DK-5230 Odense M, Denmark
[2] Centre for Proteomics, Department of Biochemistry and Molecular Biology, University of British Columbia, Vancouver, BC V6T 1Z4, Canada
[3] Department of Proteomics and Signal Transduction, Max-Planck Institute for Biochemistry, Martinsried, Germany D-82152
[4] Beijing Institute of Genomics, Chinese Academy of Sciences, Beijing 101300, China
[5] Broad Institute of Harvard and MIT, Cambridge, MA 02139, USA
[6] Department of Systems Biology, Massachusetts General Hospital, Harvard Medical School, Boston, MA 02115, USA
*Contact: mmann@biochem.mpg.de
DOI 10.1016/j.cell.2006.03.022

**SUMMARY**

Protein localization to membrane-enclosed organelles is a central feature of cellular organization. Using protein correlation profiling, we have mapped 1,404 proteins to ten subcellular locations in mouse liver, and these correspond with enzymatic assays, marker protein profiles, and confocal microscopy. These localizations allowed assessment of the specificity in published organellar proteomic inventories and demonstrate multiple locations for 39% of all organellar proteins. Integration of proteomic and genomic data enabled us to identify networks of coexpressed genes, *cis*-regulatory motifs, and putative transcriptional regulators involved in organelle biogenesis. Our analysis ties biochemistry, cell biology, and genomics into a common framework for organelle analysis.

microscopic examination of an organelle, certain proteins or enzyme activities that appear to localize exclusively to that organelle are considered markers, essentially defining that compartment.

Recently, proteomics (de Hoog and Mann, 2004) has been applied to study organelle composition. The genetic tractability of *Saccharomyces cerevisiae* has allowed a large fraction of yeast ORFs to be tagged for localization studies (Ross-Macdonald et al., 1999; Kumar et al., 2002; Huh et al., 2003), but such an approach is more challenging in mammalian systems due, in part, to artifacts from overexpression (Simpson et al., 2000). Mass spectrometry-based proteomics (Aebersold and Mann, 2003) is often employed to characterize the protein composition of organelle-enriched fractions. Indeed, protein catalogs are now available for virtually all cytoplasmic organelles as well as most of the major nuclear ones (reviewed in Yates et al., 2005). However, due to the high sensitivity of mass spectrometers and the difficulties inherent in purifying organelles to homogeneity, it has been challenging
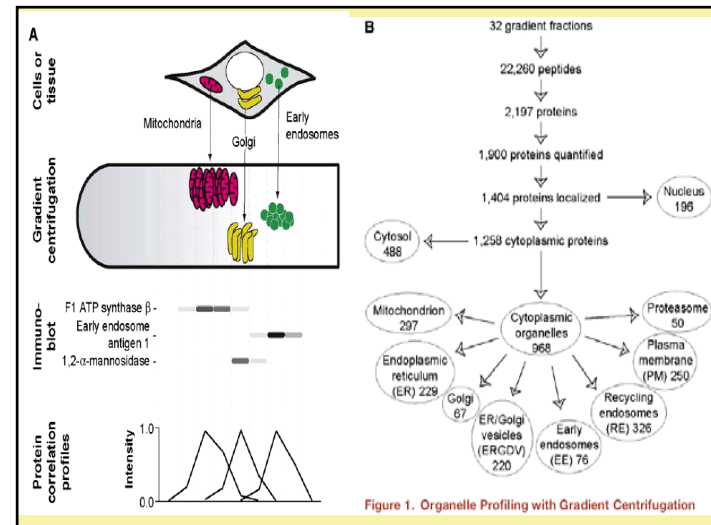
---



Figure 1. Organelle Profiling with Gradient Centrifugation

**The Birth of Molecular Biology: DNA Structure**

inside, linked together by hydrogen bonds. This structure as described is rather ill-defined, and for this reason we shall not comment on it.

We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate di-ester groups joining β-D-deoxy-ribofuranose residues with 3',5' linkages. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow right-handed helices, but owing to the dyad the sequences of the atoms in the two chains run in opposite directions. Each chain loosely resembles Fur-berg's² model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside. The configuration of the sugar and the atoms near it is close to Furberg's 'standard configuration', the sugar being roughly perpendi-cular to the attached base. There

This figure is purely diagrammatic. The two ribbons symbolise the two phosphate—sugar chains, and the horizontal rods the pairs of bases holding the chains together. The vertical line marks the fibre axis

*nature*
the **human** genome

*Nature* – 1953    *Nature* – 2001

---

Dideoxy sequencing

(a) Chain-terminating (dideoxy) nucleotides — Cannot form a phosphodiester bond with next incoming dNTP

5' DNA strand 3'
TTAGACCCGATAAGCCCGCA

DNA polymerase I + 4 dNTPs + ddATP

TTAGACCCGATAAGCCCGCA
+
ATTCGGGCGT
+
ATCTGGGCTATTCGGGCGT
+
AATCTGGGCTATTCGGGCGT

(b) DNA polymerase I + 4 dNTPs + Labeled primer
ddATP  ddTTP  ddCTP  ddGTP

Acrylamide gel

DNA sequence of original strand

---

Automated dye-terminator sequencing

4-fluorescently labelled dideoxy dye terminators
ddATP
ddGTP    pool and load in a single well or capillary
ddCTP      • scan with laser + detector specific for each dye
ddTTP      • automated base calling
           • very long reads (~ 1000 bases)/run

---

Physical mapping and sequencing of the human genome

Genomic DNA

BAC library (bacterial artificial chromosome)

Organized mapped large clone contigs

BAC to be sequenced

Shotgun clones

Shotgun sequence    ...ACCGTAAATGGGCTGATCATGCTTAAA
                       TGATCATGCTTAAACCCTGTGCATCCTACTG...

Assembly    ...ACCGTAAATGGGCTGATCATGCTTAAACCCTGTGCATCCTACTG...

*Nature* (2001) **409** p. 860-921

**Jim Kent** is a research scientist at UC Santa Cruz.

The human genome project was ultimately a race between Celera Genomics and the public effort, with the final push being a bioinformatics problem to put all of the sequence reads together into a draft genome sequence. **Jim Kent was a grad student at UCSC**, who worked for weeks developing the algorithm to put all of this together, **beating Celera by 3 days** to an assembled human genome sequence.

His efforts ensured that the human genome data remained in the public domain and were not patented into private intellectual property.

Kent built a grid of cheap, commodity PC's running the Linux operating system and other Freeware to beat Celera's, what was thought of then as the, world's most powerful civilian computer. In **June 2000**, thanks to the work done by Kent and several others, the Human Genome Project was able to publish its data in the Public Domain just hours ahead of Celera.
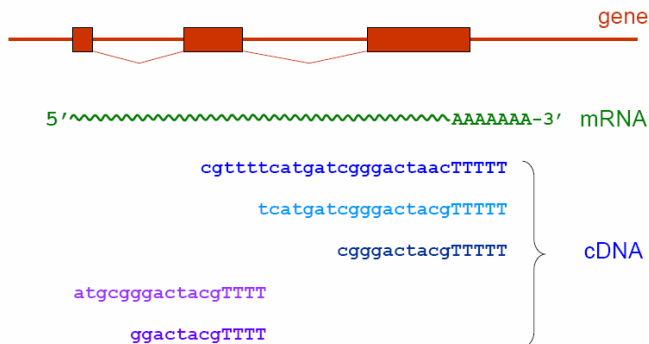
Kent went on to write BLAT and the UCSC Human Genome Browser to help analyze important genome data, receiving his PhD in biology in 2002. Today at UCSC he works primarily on web tools to help understand the human genome. He helps maintain and upgrade the browser, and has worked on recent projects such as comparative genomics and Parasol.

---

## Finding genes in genomes

- compare to EST or cDNA sequence

- look for open reading frames

- similarity to other genes and proteins

- Gene prediction algorithms (identifying splice sites, coding sequence bias, etc.)

---

Genes can also be identified by sequencing cDNAs at random. The sequenced cDNAs are called ESTs (expressed sequence tags)



---

## The BIG QUESTION:

### Why do we have so few genes?

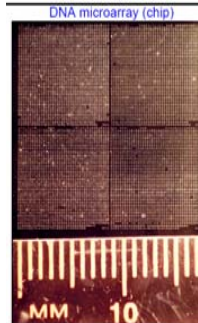| Species | Genome size | Number of genes |
|---|---|---|
| Human (*Homo sapiens*) | 2.9 billion base pairs | 25,000 - 30,000 |
| Fruit fly (*Drosophila melanogaster*) | 120 million base pairs | 13,600 |
| Worm (*Caenorhabditis elegans*) | 97 million base pairs | 19,000 |
| Budding yeast (*Saccharomyces cerevisiae*) | 12 million base pairs | 6,000 |
| *E. coli* | 4.1 million base pairs | 4,800 |

## Genomics   vs.   Proteomics

**With the completion of a rough draft of the human genome, many researchers are looking at how genes and proteins interact to form other proteins.  A surprising finding of the Human Genome Project is that there are far fewer protein-coding genes in the human genome than proteins in the human proteome (20,000 to 25,000 genes vs. about 1,000,000 proteins).  The human body may contain more than 2 million proteins, each having different functions. The protein diversity is thought to be due to alternative splicing and post-translational modification of proteins. The discrepancy implies that protein diversity cannot be fully characterized by gene expression analysis, thus proteomics is useful for characterizing cells and tissues.**

## Functional genomics and proteomics

- Identify genes and proteins encoded in the genome (Gene finding)

- Measure gene expression on a genome-wide scale (microarrays)

- Identify protein function
  30-50% of the genes in a genome are of unknown function

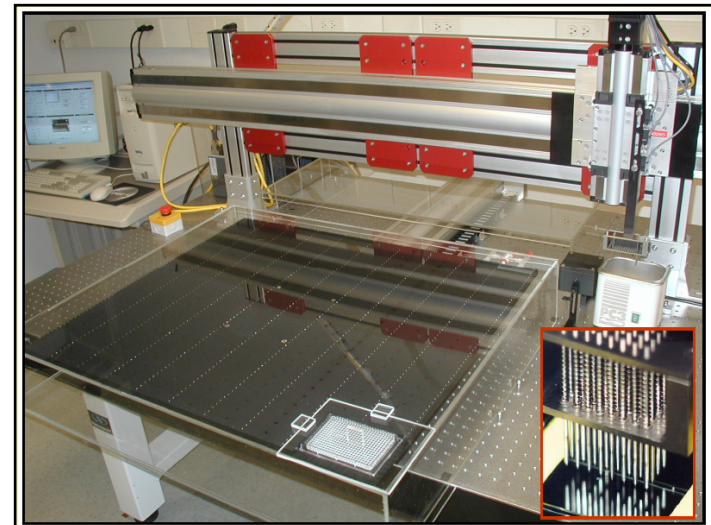- Identify protein interactions, biochemical pathways, gene interaction networks inside cells

## Methods of making microarrays

- Robotic spotting
  - using a printing tip
  - using inkjets

- Synthesis of oligonucleotides
  - photolithography (Affymetrix)
  - using inkjets
  - Digital Light Processor (DLP) or Digital Micromirror Device (DMD)



DNA microarray (chip)

Microarrays can be used to study gene expression, DNA-protein interactions, mutations, protein-protein interactions, etc., all on a genome-wide scale

*Note: Thanks to Prof. Vishy Iyer for many of these slides on microarrrays.*

Affymetrix GeneChip

courtesy: www.affymetrix.com



In cell B, relative to cell A,

Gene 1 is equally expressed

Gene 2 is overexpressed

Gene 3 is underexpressed



DNA microarray after hybridization of fluorescent probes



Original microarray image

Colour representation of differential gene expression

| Green | Red | Red/Green | |
|---|---|---|---|
| 200 | 10000 | 50.00 | Gene 1 |
| 4800 | 4800 | 1.00 | Gene 2 |
| 9000 | 300 | 0.03 | Gene 3 |

- Large amounts of data can be displayed in this manner
- Gene expression data can be computationally analyzed and organized to reveal patterns

**Data after hierarchical clustering**

Experiments

Original data

Genes

---

Functionally related genes are often co-expressed

A cluster of co-expressed genes

- Ribosomal protein 1
- Ribosomal protein 2
- Ribosomal protein 3
- Ribosomal protein 4
- Ribosomal protein 5
- Ribosomal protein 6
- Unknown Gene X
- Ribosomal protein 7
- Ribosomal protein 8

Thus, unknown Gene X may also be a ribosomal protein

---

*S. cerevisiae* mitotic cell-cycle

α    cdc 15    cdc28    elu

M/G1
G1
S
G2
M

G1
M    S
G2

Brenda Andrews lab, University of Toronto

Spellman *et al*, (1998) *Mol. Biol. Cell*

---

# Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling
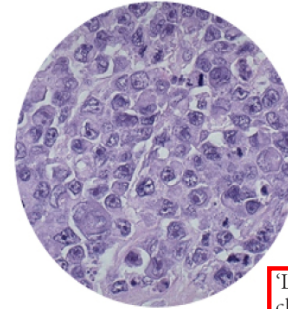
Ash A. Alizadeh[1,2], Michael B. Eisen[2,3,4], R. Eric Davis[5], Chi Ma[5], Izidore S. Lossos[6], Andreas Rosenwald[5], Jennifer C. Boldrick[1], Hajeer Sabet[5], Truc Tran[5], Xin Yu[5], John I. Powell[7], Liming Yang[7], Gerald E. Marti[8], Troy Moore[9], James Hudson Jr[9], Lisheng Lu[10], David B. Lewis[10], Robert Tibshirani[11], Gavin Sherlock[4], Wing C. Chan[12], Timothy C. Greiner[12], Dennis D. Weisenburger[12], James O. Armitage[13], Roger Warnke[14], Ronald Levy[6], Wyndham Wilson[15], Michael R. Grever[16], John C. Byrd[17], David Botstein[4], Patrick O. Brown[1,18] & Louis M. Staudt[5]

Diffuse large B-cell lymphoma (DLBCL), the most common subtype of non-Hodgkin's lymphoma, is clinically heterogeneous: 40% of patients respond well to current therapy and have prolonged survival, whereas the remainder succumb to the disease. We proposed that this variability in natural history reflects unrecognized molecular heterogeneity in the tumours. Using DNA microarrays, we have conducted a systematic characterization of gene expression in B-cell malignancies. Here we show that there is diversity in gene expression among the tumours of DLBCL patients, apparently reflecting the variation in tumour proliferation rate, host response and differentiation state of the tumour. We identified two molecularly distinct forms of DLBCL which had gene expression patterns indicative of different stages of B-cell differentiation. One type expressed genes characteristic of germinal centre B cells ('germinal centre B-like DLBCL'); the second type expressed genes normally induced during *in vitro* activation of peripheral blood B cells ('activated B-like DLBCL'). Patients with germinal centre B-like DLBCL had a significantly better overall survival than those with activated B-like DLBCL. The molecular classification of tumours on the basis of gene expression can thus identify previously undetected and clinically significant subtypes of cancer.

Despite the variety of clinical, morphological and molecular parameters used to classify human malignancies today, patients receiving the same diagnosis can have markedly different clinical courses and treatment responses. The history of cancer diagnosis has been punctuated by reassortments and subdivisions of diagnostic categories. There is little doubt that our current taxonomy of cancer still lumps together molecularly distinct diseases with distinct clinical phenotypes. Molecular heterogeneity within individual cancer diagnostic categories is already evident in the variable presence of chromosomal translocations, deletions of tumour suppressor genes and numerical chromosomal abnormalities. The classification of human cancer is likely to become increasingly more informative and clinically useful as more detailed molecular analyses of the tumours are conducted.

---

## The challenge of cancer diagnosis



**Diffuse large B-cell lymphoma** is the most common subtype of non-Hodgkin's lymphoma. With current treatments, long-term survival can be achieved in only 40% of patients. There are no reliable indicators — morphological, clinical, immunohistochemical or genetic — that can be used to recognize subclasses of **DLBCL** and point to a differential therapeutic approach to patients.

'Lymphochip', a microarray carrying 18,000 clones of complementary DNA designed to monitor genes involved in normal and abnormal lymphocyte development.

What type of cancer?

What is the underlying molecular basis?

What is the optimal treatment?

---

# Box 1: Gene-expression profiling with microarrays

Imagine a 1-cm² chessboard. Instead of 64 squares, it has thousands, each containing DNA from a specific gene. This is a DNA microarray. The activity of each gene on the microarray can be compared in two populations of cells (A and B).

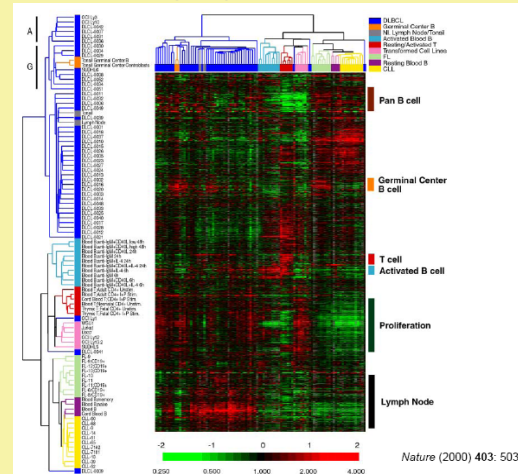When a gene is expressed it makes a transcript, and the whole population of these products from a cell can be tagged with a fluorescent dye (say, red for the A cells, green for the B cells). The microarray is bathed in a mixture of the red and green transcripts. Those that originate from a specific gene will bind to that gene on the microarray, turning red, green or somewhere in between, depending on the relative numbers of transcripts in the two cell types.
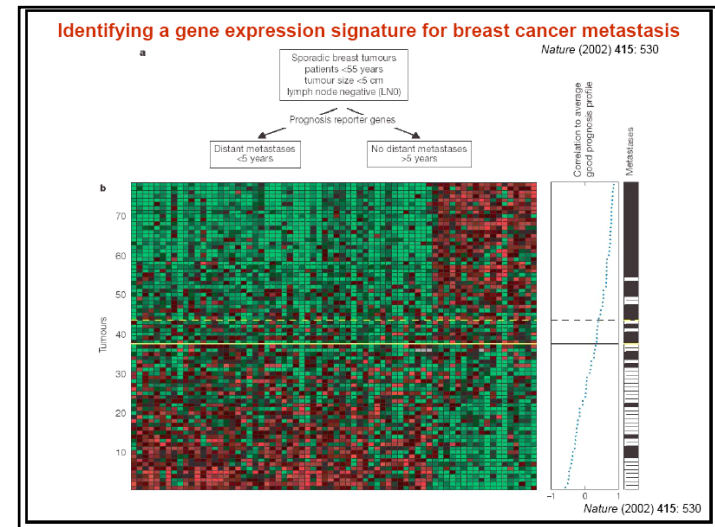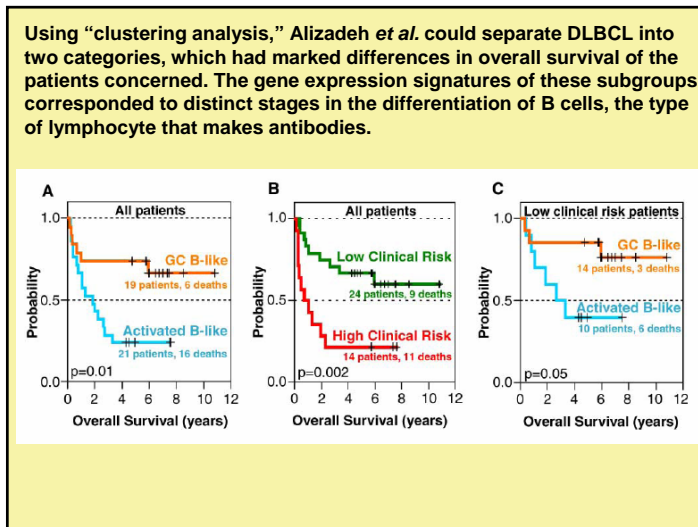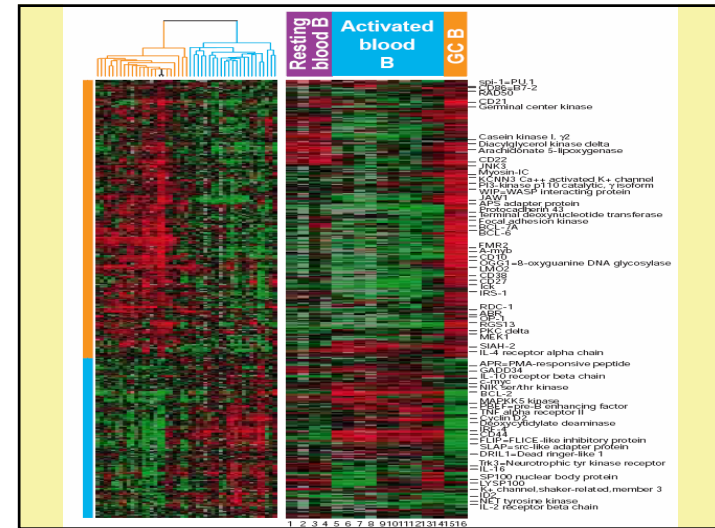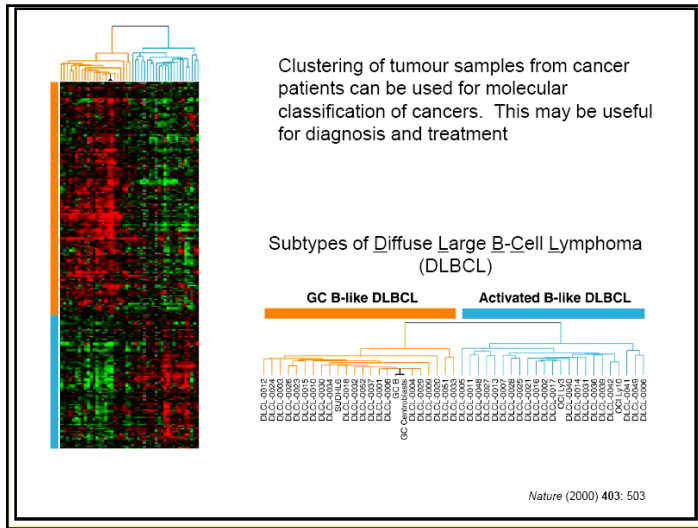
So the microarray provides a snapshot of gene activity for thousands of genes. Data from many experiments can be compared and genes that have consistent patterns of activity can be grouped or clustered. In this way, genes that characterize a particular cell state, such as malignancy, can be identified — so providing new information about the biology of the cell state.

**Mark Patterson**

---

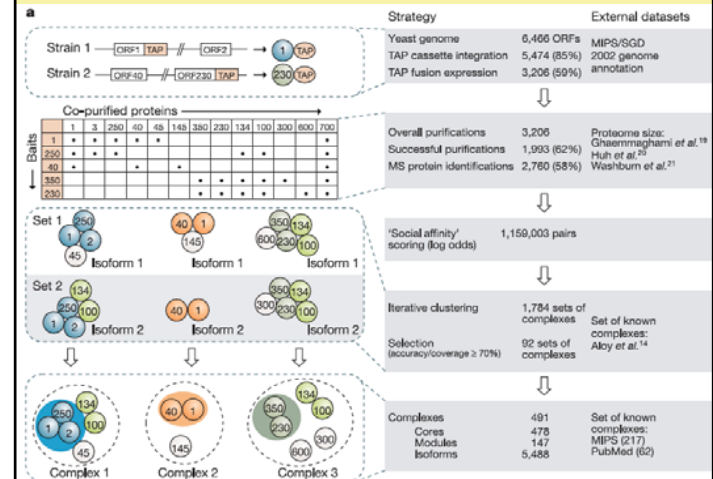## Hierarchical clustering of gene expression data (as ratios).



Nature (2000) **403**: 503

---

Clustering of tumour samples from cancer patients can be used for molecular classification of cancers. This may be useful for diagnosis and treatment

Subtypes of <u>D</u>iffuse <u>L</u>arge <u>B</u>-<u>C</u>ell Lymphoma (DLBCL)

GC B-like DLBCL    Activated B-like DLBCL

*Nature* (2000) **403**: 503



**Using "clustering analysis," Alizadeh *et al.* could separate DLBCL into two categories, which had marked differences in overall survival of the patients concerned. The gene expression signatures of these subgroups corresponded to distinct stages in the differentiation of B cells, the type of lymphocyte that makes antibodies.**



A — All patients
GC B-like — 19 patients, 6 deaths
Activated B-like — 21 patients, 16 deaths
p=0.01

B — All patients
Low Clinical Risk — 24 patients, 9 deaths
High Clinical Risk — 14 patients, 11 deaths
p=0.002

C — Low clinical risk patients
GC B-like — 14 patients, 3 deaths
Activated B-like — 10 patients, 6 deaths
p=0.05

Overall Survival (years)

Identifying a gene expression signature for breast cancer metastasis

*Nature* (2002) **415**: 530



Sporadic breast tumours
patients <55 years
tumour size <5 cm
lymph node negative (LN0)

Prognosis reporter genes

Distant metastases <5 years
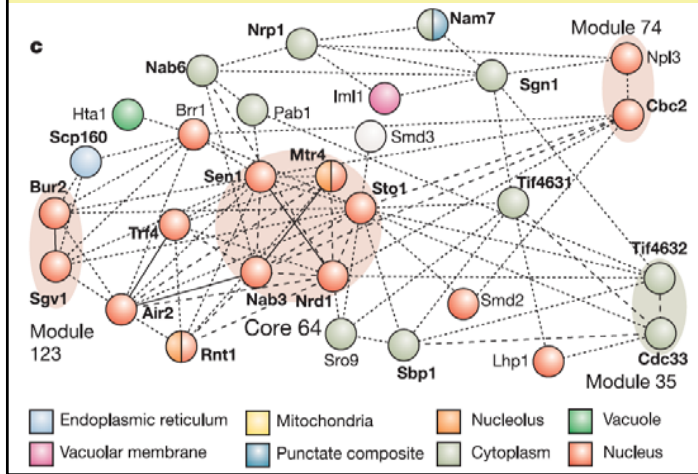No distant metastases >5 years

*Nature* (2002) **415**: 530
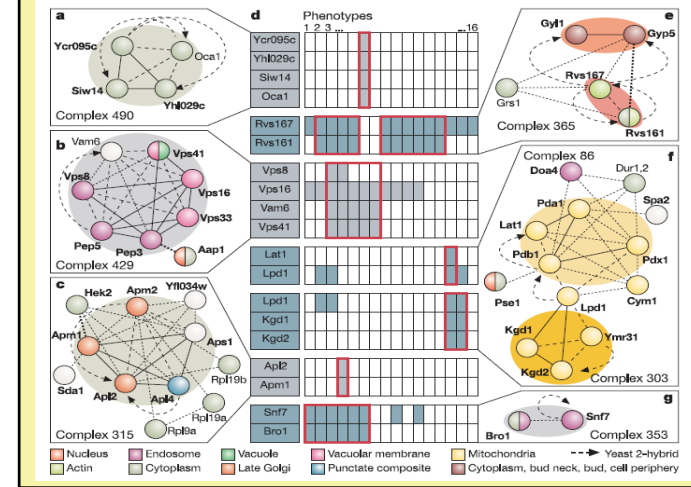
Architecture and Modularity of Complexes

## Architecture and Modularity of Complexes



## Phenotypic Data Mapped to Complexes



## Systems Biology Approach



15