

Integrated Data Systems for Interpreting Genome-Focused Data in Cancer

Ajay N. Jain, PhD

ajain@cc.ucsf.edu

<http://jainlab.ucsf.edu>

Copyright © 2002, Ajay N. Jain

All Rights Reserved

Biology is shifting from being an observational science to being a quantitative molecular science

Old biology: measure one/two things in two/three conditions

- ◆ High cost per measurement
- ◆ Analysis straightforward
- ◆ Enormously difficult to work out pathways

New biology: measure 10,000 things under many conditions

- ◆ Low cost per measurement
- ◆ Analysis no longer straightforward, but payoff can be bigger
- ◆ Biology as a complex system: Can we work out biological pathways this way?

Cancer biology as a complex system: The marriage of experimental data with annotation information

Phenotype

Proliferation
Measurable phenotypes.

Cell Lines



Apoptosis



Protein

P₁
Protein status for over multiple conditions.



P_n



RNA

G₁
Gene expression levels over multiple conditions.



G_n

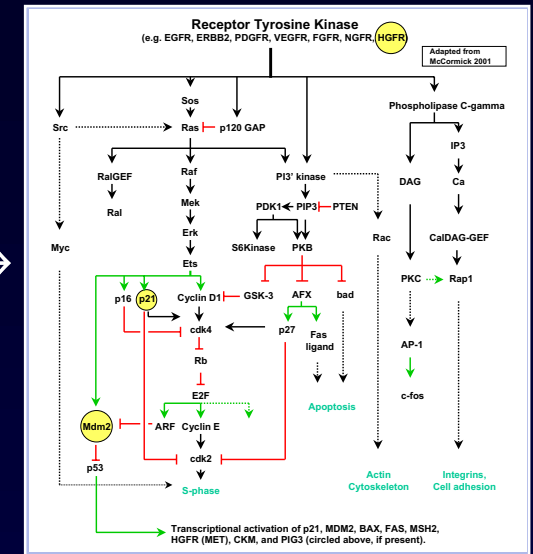


DNA

L₁
DNA copy number over the entire genome.



Pathway Structure →



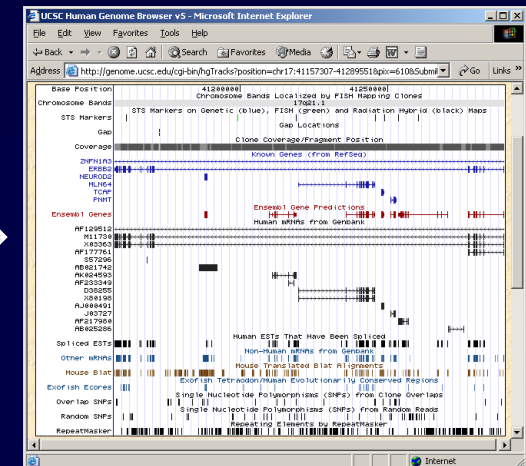
ERBB2:

EC Number: 2.7.1.112

- oncogenesis
- cell proliferation
- Neu/ErbB-2 receptor
- protein phosphorylation
- protein dephosphorylation
- cell growth and maintenance
- receptor signaling tyrosine kinase

← Gene Annotations

Genomic Mapping + Context →



DNA Copy Number Aberrations in Cancer

Some important cancer genes are sometimes present in altered copy number in some tumors

Increases over express oncogenes, decreases help inactivate suppressor genes, dosage changes affect expression

These copy number alterations can be predictive of tumor phenotype and patient outcome

Quantitative analysis of chromosomal CGH in human breast tumors associates copy number abnormalities with p53 status and patient survival

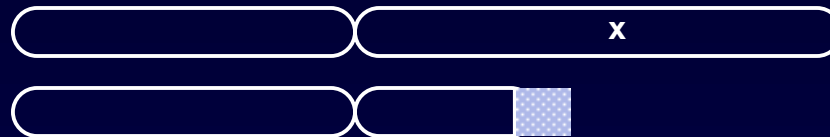
Ajay N. Jain*[†], Koei Chin*[†], Anne-Lise Børresen-Dale[‡], Bjorn K. Erikstein[‡], Per Eystein Lonning[§], Rolf Kaaresen[¶], and Joe W. Gray*^{||}

*UCSF Cancer Center, University of California, San Francisco, Box 0128, San Francisco, CA 94143-0128; [†]Departments of Genetics and Oncology, Institute for Cancer Research, Norwegian Radium Hospital, Montebello, 0310 Oslo, Norway; [‡]Department of Oncology, Haukeland Hospital, 5021 Haukeland Sykehus, Norway; and [¶]Department of Surgery, Ullev Hospital, 0407 Oslo, Norway

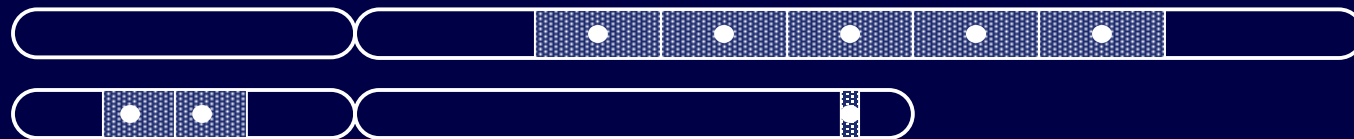
Communicated by James E. Cleaver, University of California, San Francisco, CA, May 15, 2001 (received for review February 5, 2001)

Copy number changes usually involve a DNA segment that is substantially larger than the critical gene(s)

Mutation plus terminal deletion of tumor a suppressor gene



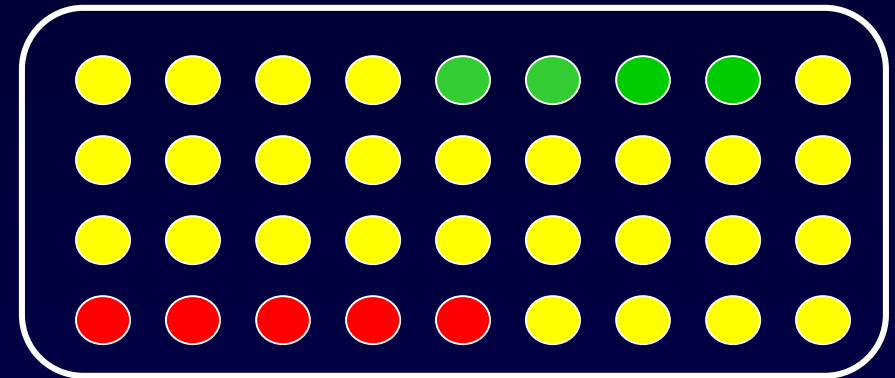
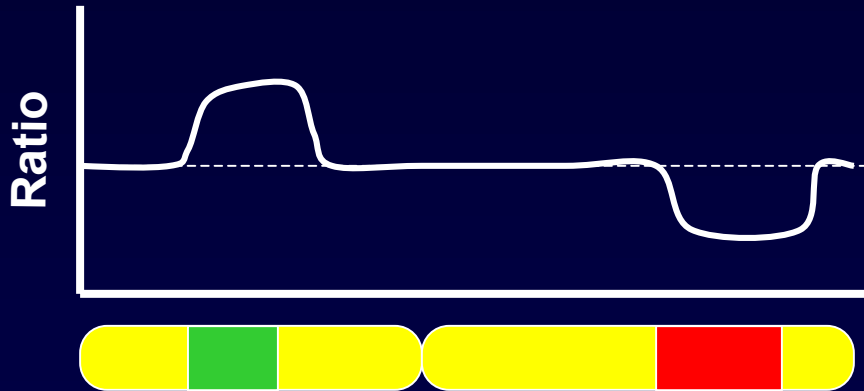
Amplification of an oncogene and surrounding DNA; extra copies can be located anywhere in the genome



Comparative Genomic Hybridization (CGH): Chromosomal targets versus array targets

Test DNA

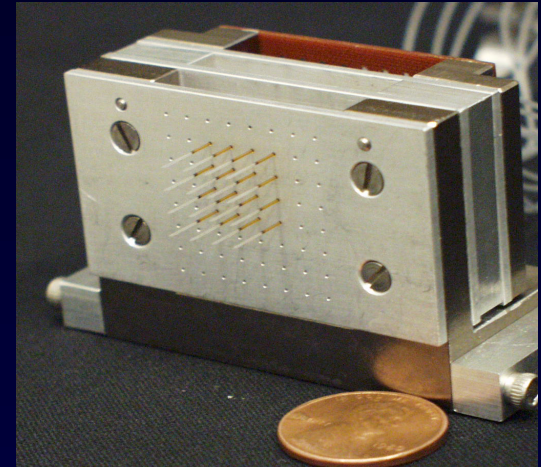
Reference DNA



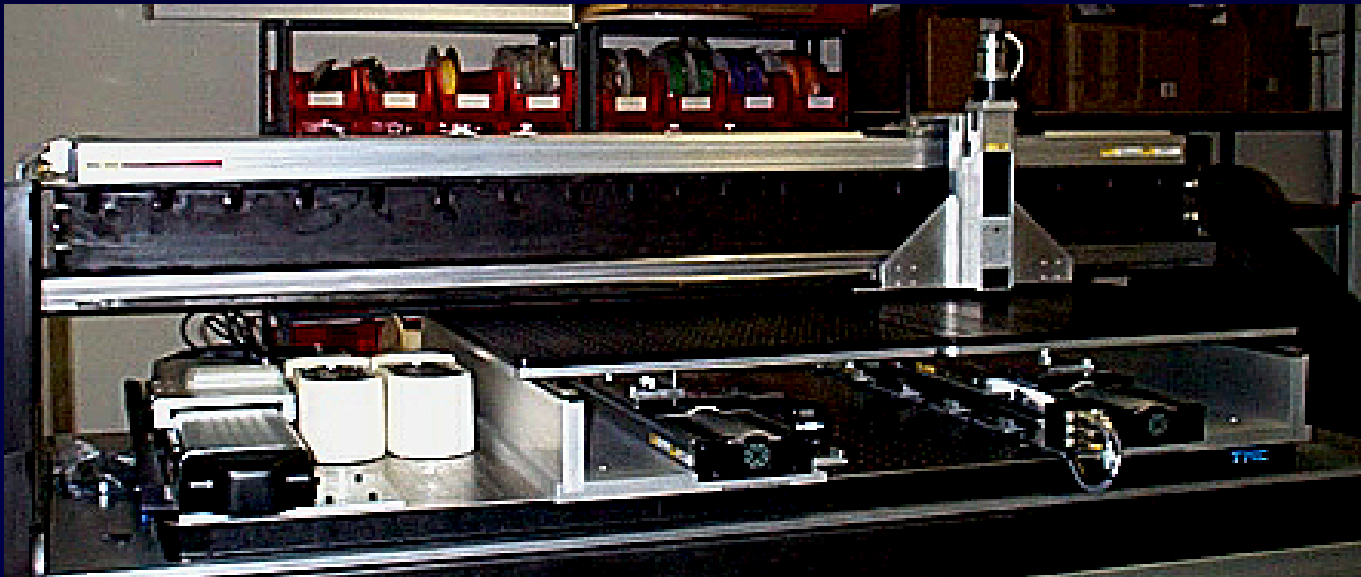
Chromosome CGH provides
"cytogenetic" resolution ~ 10 Mb

Resolution of array CGH depends on
spacing and length of clones

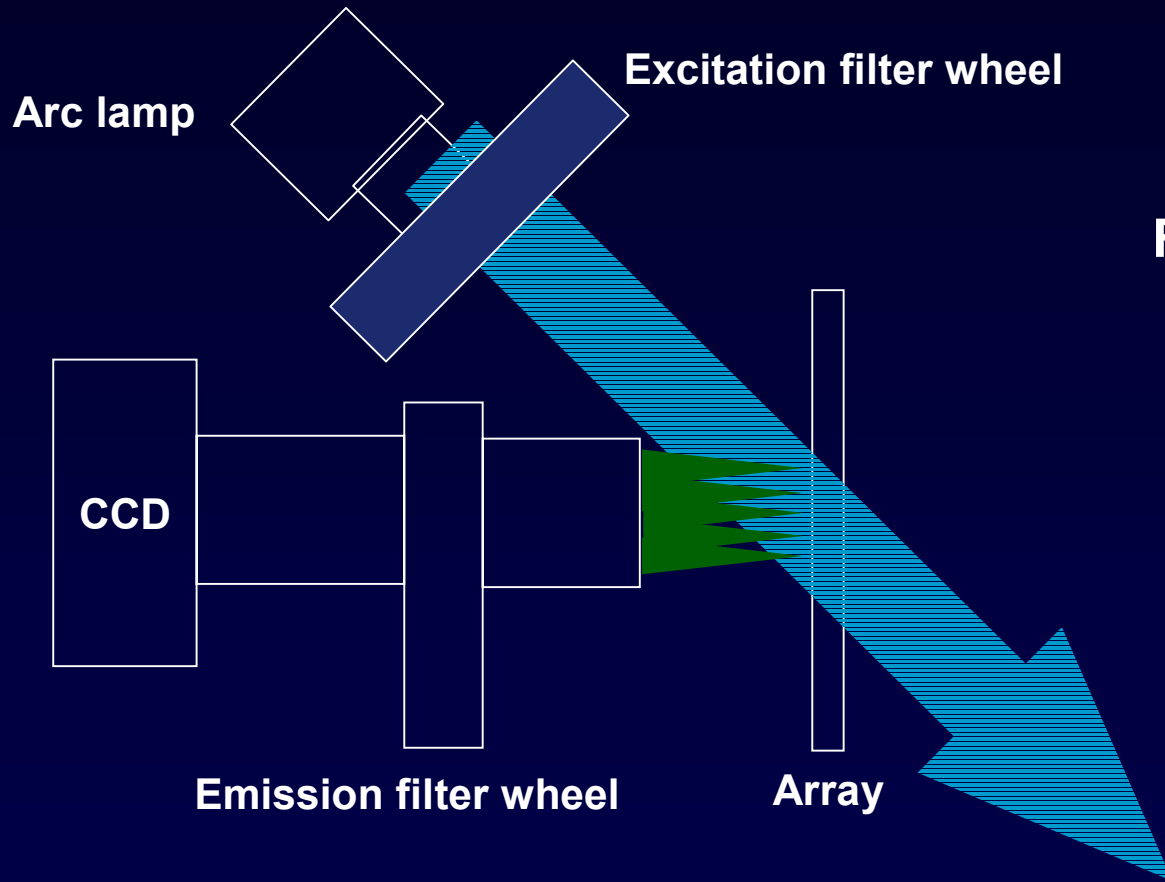
Capillary Print head



Overview of Layout



Custom CCD Imaging System: Allows for imaging multiple distinct nucleic acid species



Filters permit use of:
Dapi etc.
Fluorescein etc.
Cy3 etc.
Texas red etc.
Cy5

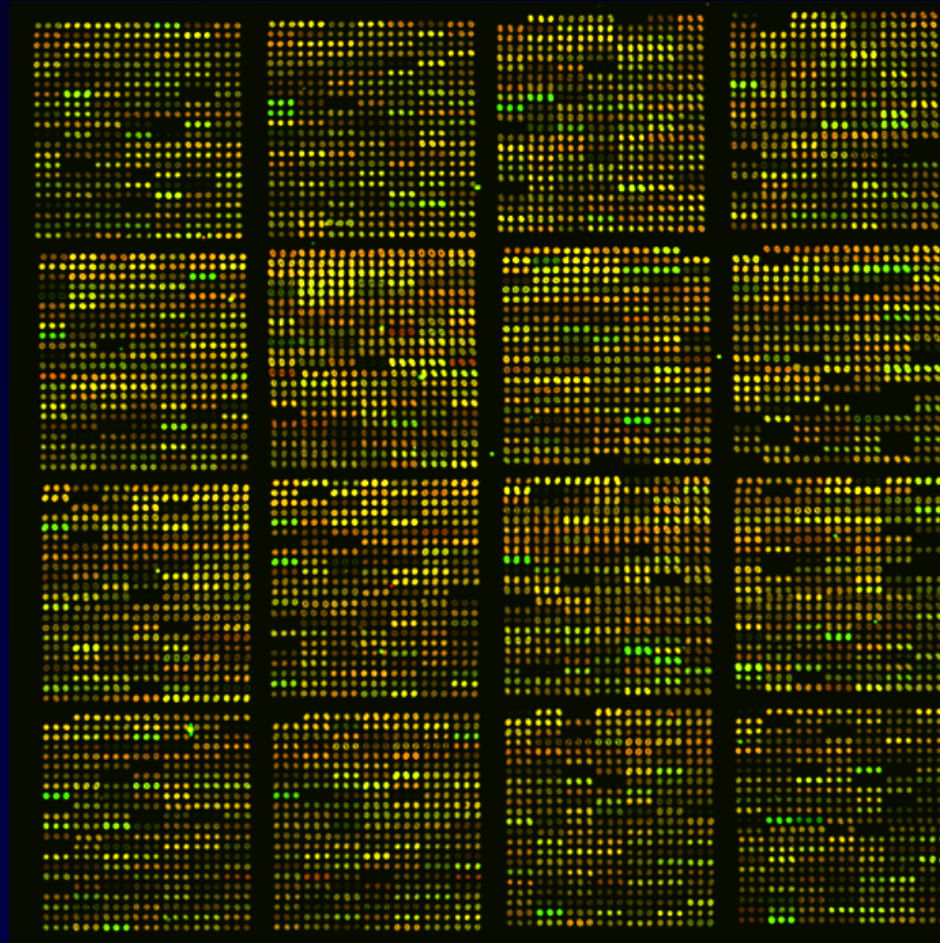
Genome Scanning Array with ~ 1.4 Mb Resolution

DNA obtained from peripheral blood, fresh, frozen or fixed tissue

3 ng to 0.5 μ g input DNA

Random Prime Nick Translation FITC, Cy3, Cy5

16 - 40 hr Hyb.



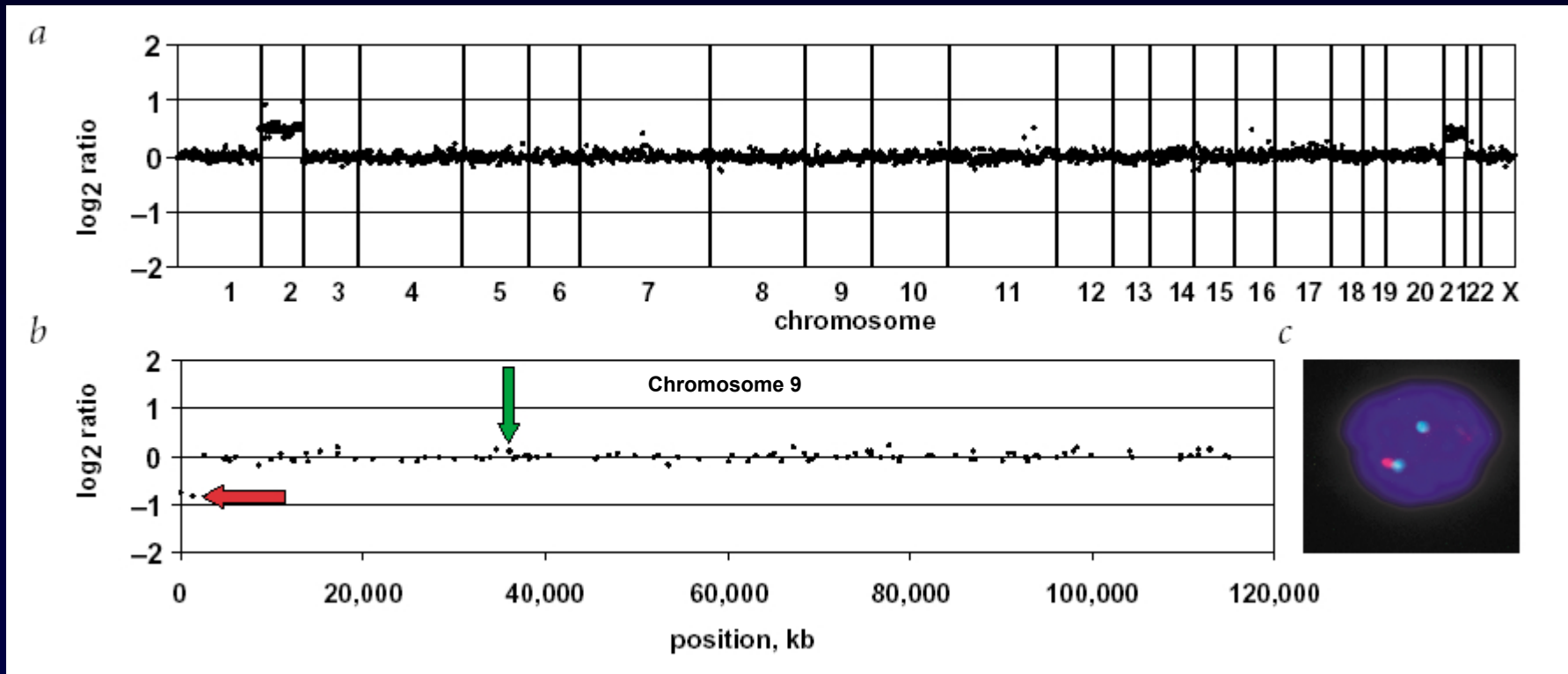
2500 BACs

Triplicate spots

130 μ m centers

864 well plates

12 mm



nature genetics • volume 29 • november 2001

Assembly of microarrays for genome-wide measurement of DNA copy number

Published online: 30 October 2001, DOI: 10.1038/ng754

We have assembled arrays of approximately 2,400 BAC clones for measurement of DNA copy number across the human genome. The arrays provide precise measurement (s.d. of log₂ ratios=0.05–0.10) in cell lines and clinical material, so that we reliably detect and quantify high-level amplifications and single-copy alteration in diploid, polyploid and heterogeneous backgrounds.

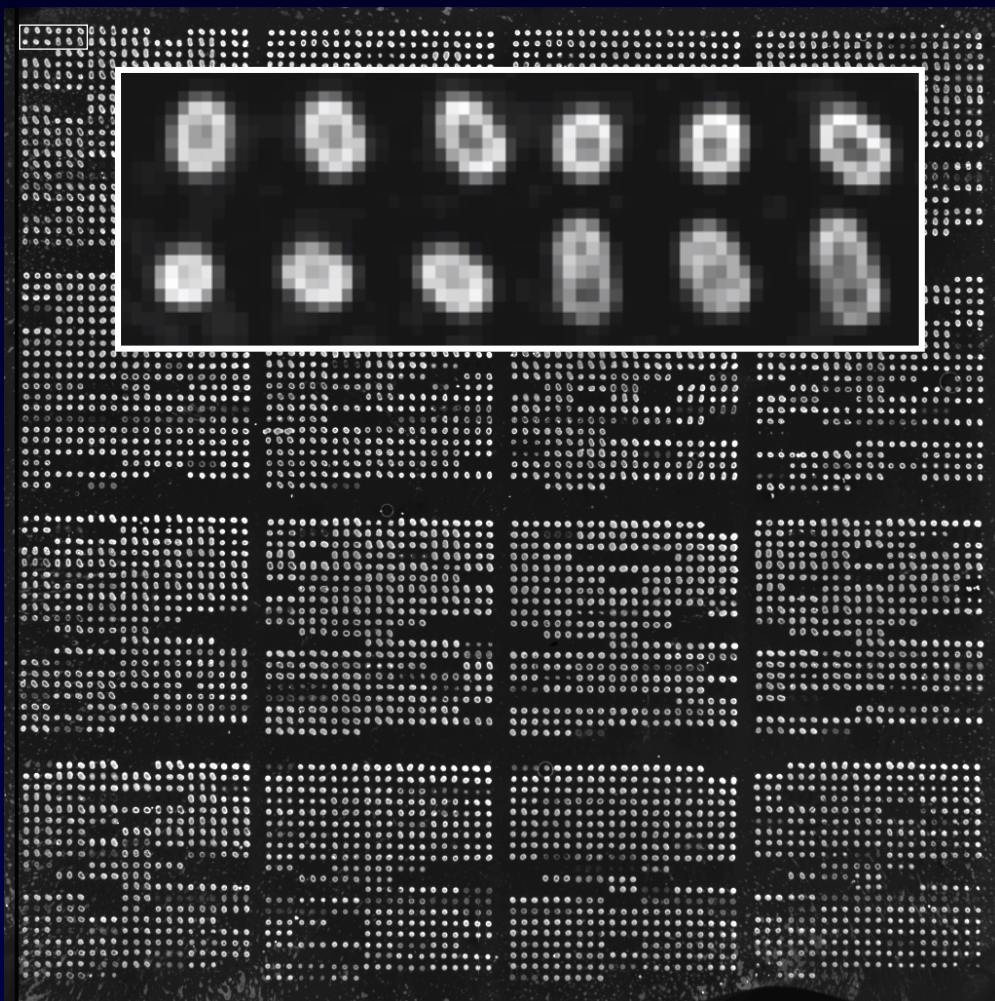
Antoine M. Snijders^{1,2}, Norma Nowak⁴, Richard Seagraves¹, Stephanie Blackwood^{1,2}, Nils Brown¹, Jeffrey Conroy⁴, Greg Hamilton¹, Anna Katherine Hindle^{1,2}, Bing Huey¹, Karen Kimura¹, Sindy Law^{1,2}, Ken Myambo¹, Joel Palmer^{1,2}, Bauke Ylstra^{1,2}, Jingzhu Pearl Yue¹, Joe W. Gray^{1,3}, Ajay N. Jain^{1–3}, Daniel Pinkel^{1,3} & Donna G. Albertson^{1–3}

¹Comprehensive Cancer Center, ²Cancer Research Institute and ³Department of Laboratory Medicine, University of California San Francisco, San Francisco, California 94143. ⁴Roswell Park Cancer Institute, Elm and Carlton Streets, Buffalo, New York 14263. Correspondence should be addressed to D.G.A. (e-mail: albertson@cc.ucsf.edu).

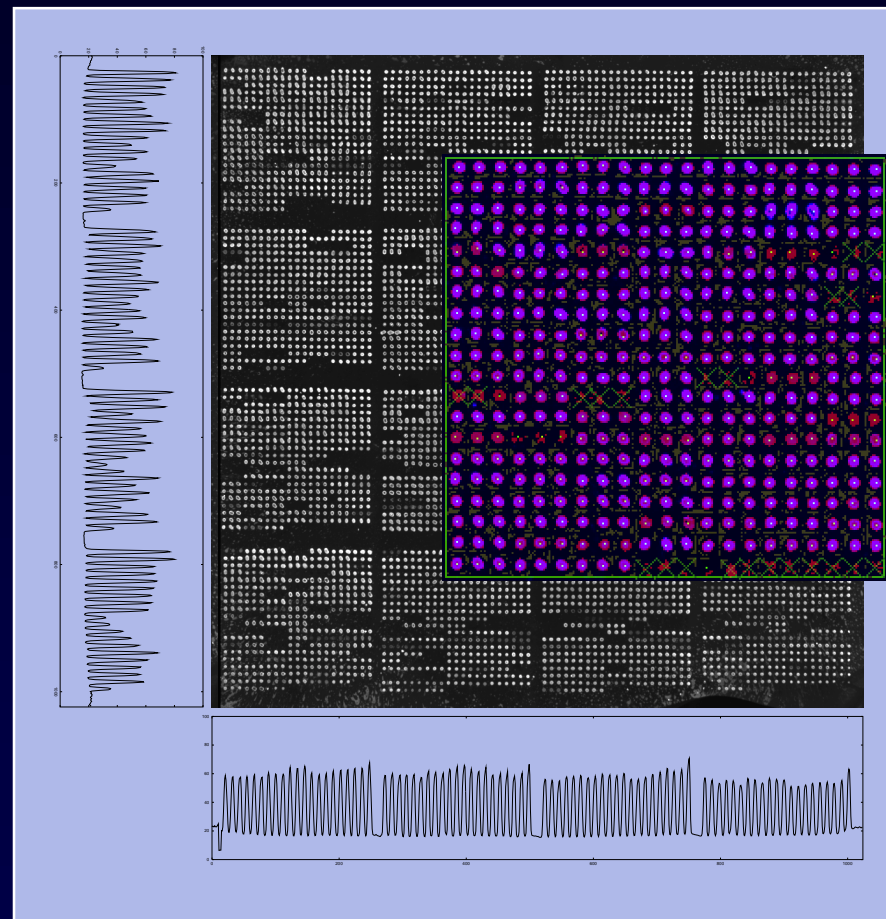
ArrayCGH can detect single copy gains and losses. Top shows genome-wide copy number for cell strain GM03576.

Bottom left shows copy number for cell strain GM03563 on chromosome 9. Bottom right: single copy deletion verified by FISH.

UCSF Spot processes typical arrayCGH hybridizations in less than 20 seconds, fully automatically



Dapi image, ~6000 spots, 135 micron centers
Spots are not perfect filled circles.



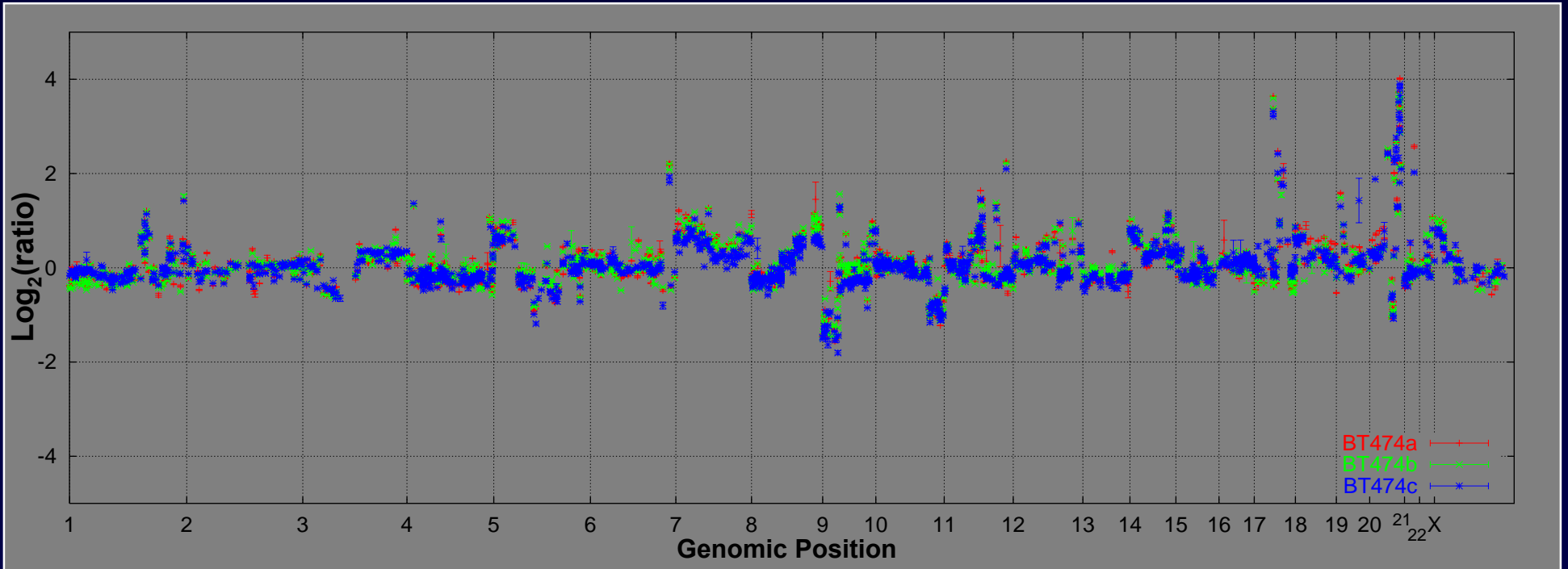
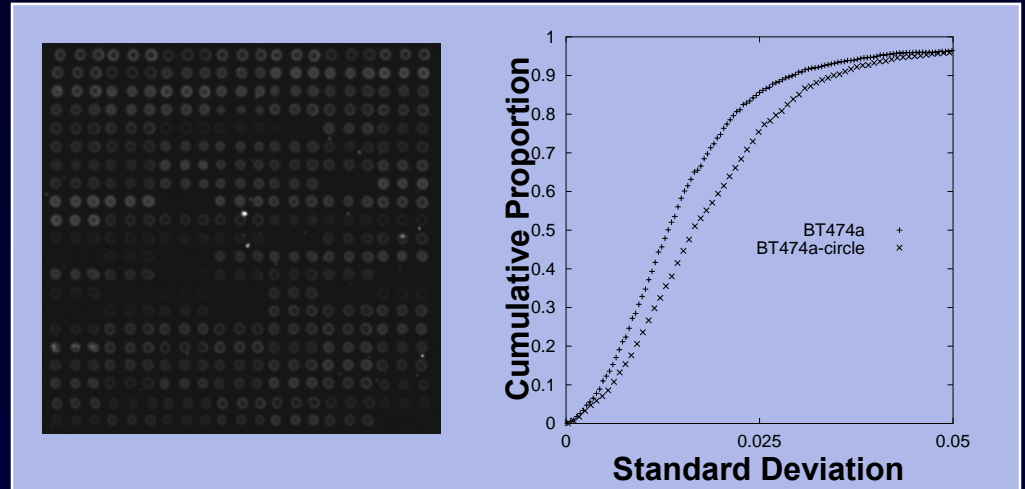
A. N. Jain, T. A. Tokuyasu, A. M. Snijders, R. Segraves, D. G. Albertson, and D. Pinkel. Fully automatic quantification of microarray image data. *Genome Research* 12: 325-332, 2002.

Explicit spot segmentation yields accurate results

Complex patterns of amplification/deletion are common

Circularity assumption leads to poorer spot replicate stddev distributions. Higher signal reduces the difference.

Experimental replicates have larger deviation than that due to image quantification noise (BT474x3, Albertson and Pinkel Labs)



Cancer biology as a complex system: The marriage of experimental data with annotation information

Phenotype

Proliferation
Measurable phenotypes.

Cell Lines



Apoptosis



Protein

P₁
Protein status for over multiple conditions.



P_n



RNA

G₁
Gene expression levels over multiple conditions.



G_n



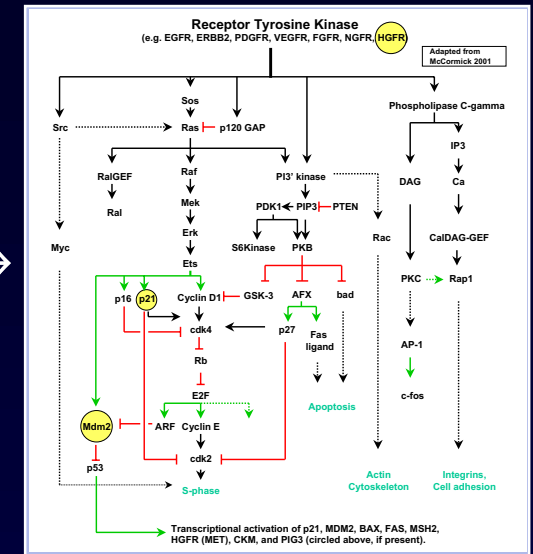
DNA

L₁
DNA copy number over the entire genome.



L_n

Pathway Structure →



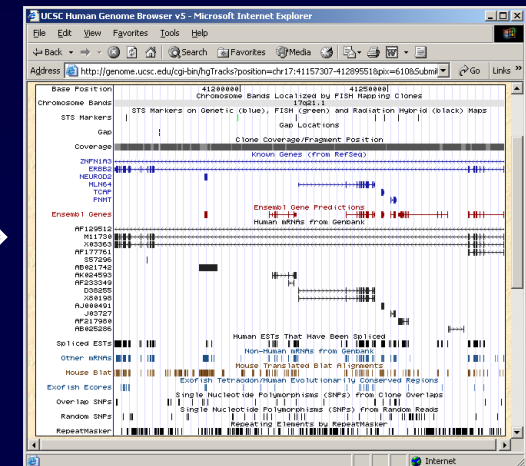
ERBB2:

EC Number: 2.7.1.112

- oncogenesis
- cell proliferation
- Neu/ErbB-2 receptor
- protein phosphorylation
- protein dephosphorylation
- cell growth and maintenance
- receptor signaling tyrosine kinase

← Gene Annotations

Genomic Mapping + Context →



Biological Pathways: The receptor tyrosine kinase signaling pathways

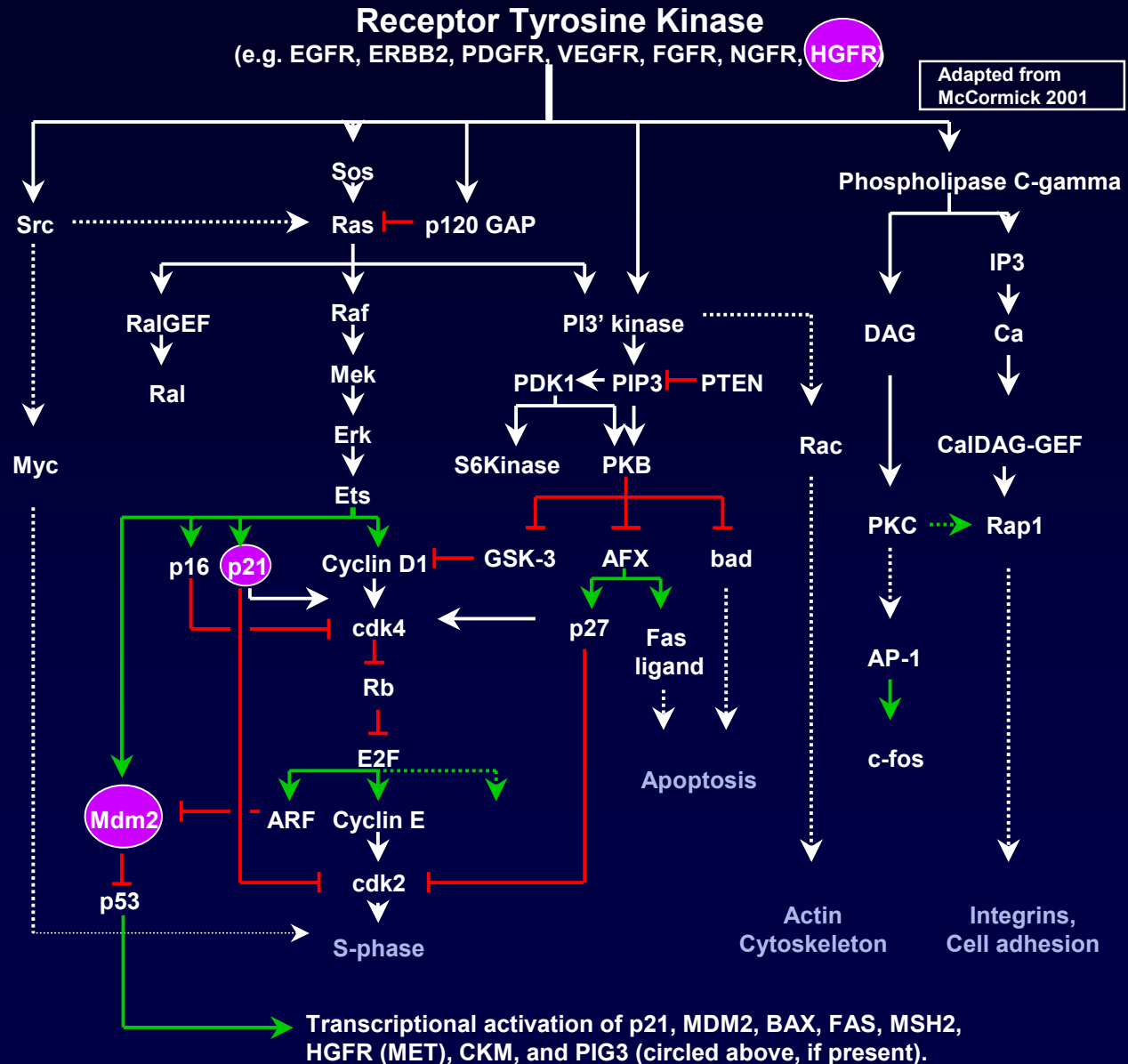
RTK signaling:

- ◆ Cell cycle control
- ◆ Apoptosis
- ◆ Cell adhesion

Therapeutics:

- ◆ Herceptin
- ◆ PI3-kinase inhibitors
- ◆ Onyx015

Pathway derived through hard molecular biology and biochemistry

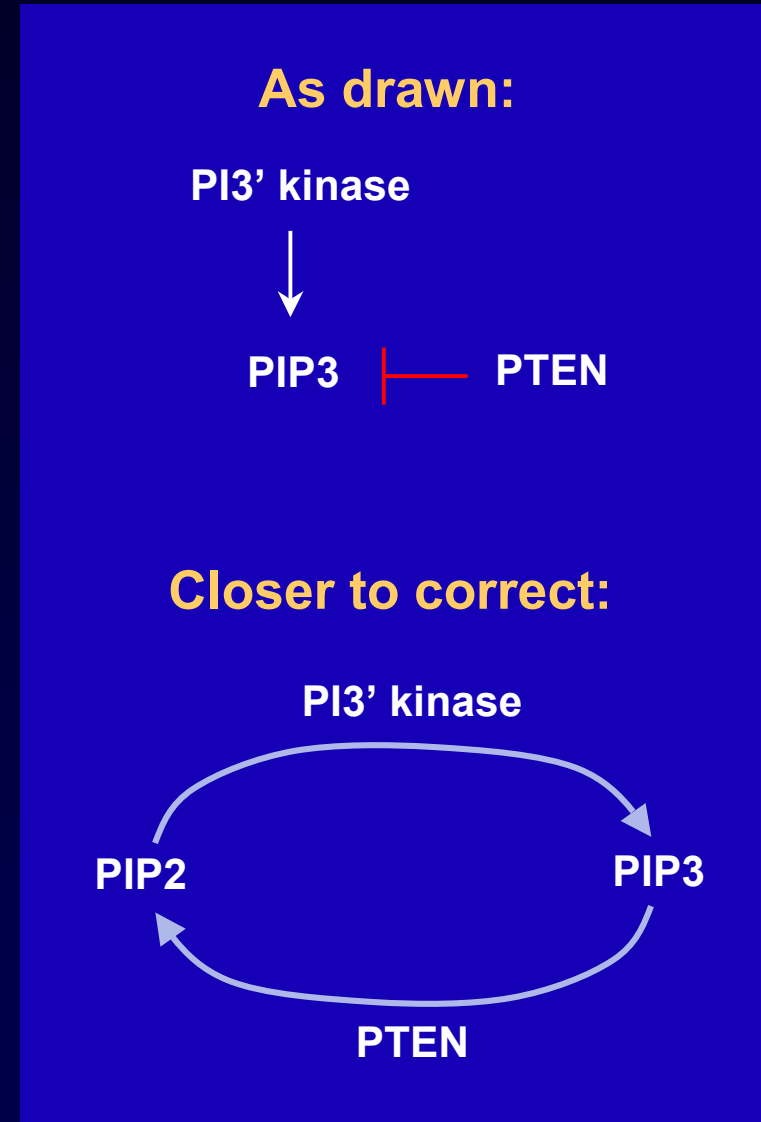


Problem 1: symbols are overloaded

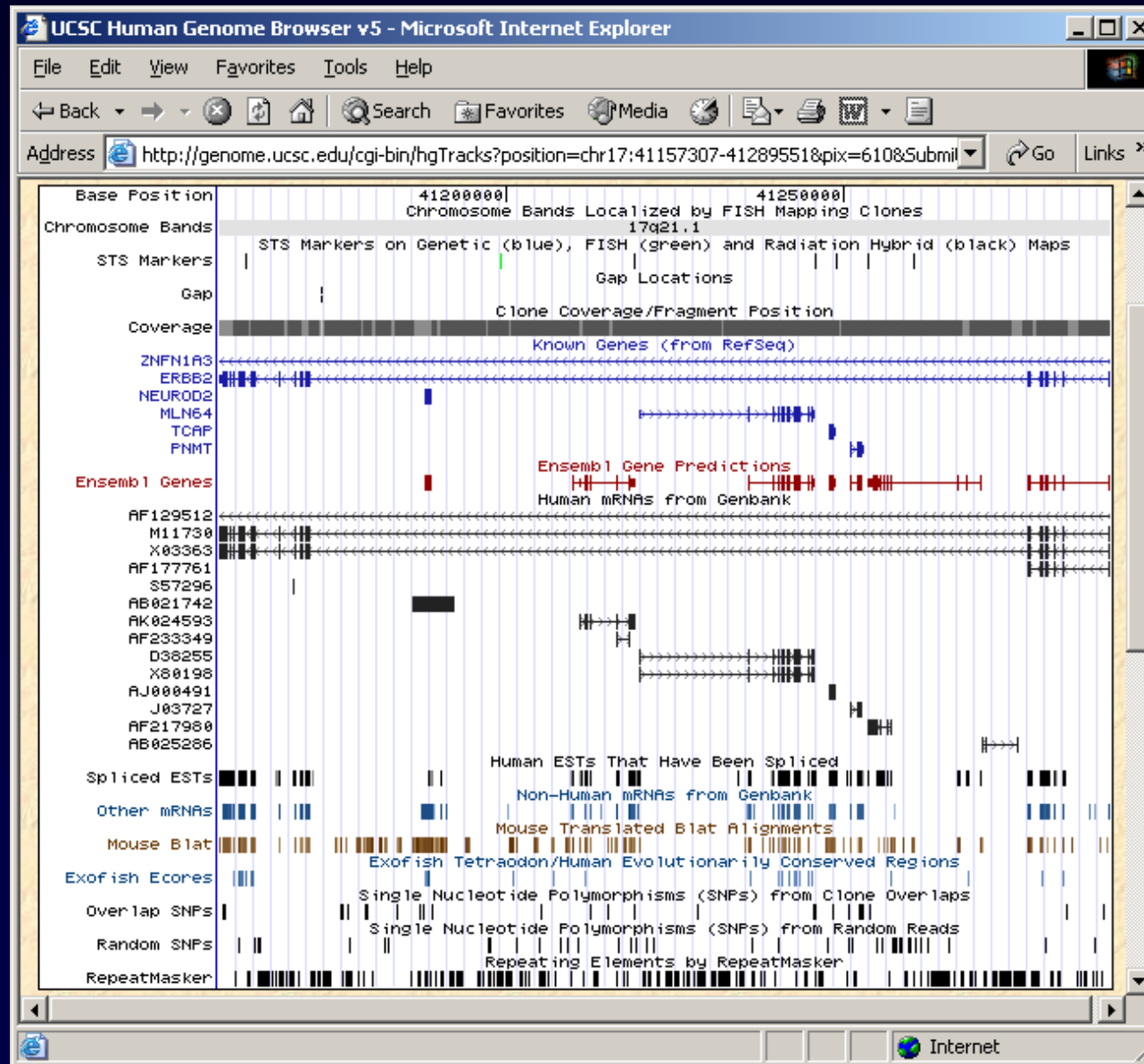
Problem 2: shorthand is used

Problem 3: knowledge is incomplete

But we must still try to represent this information



The human genome contains a great deal of information



Large number of measurements

- ◆ ~3000 for genome-wide array-CGH
- ◆ >30,000 for expression arrays

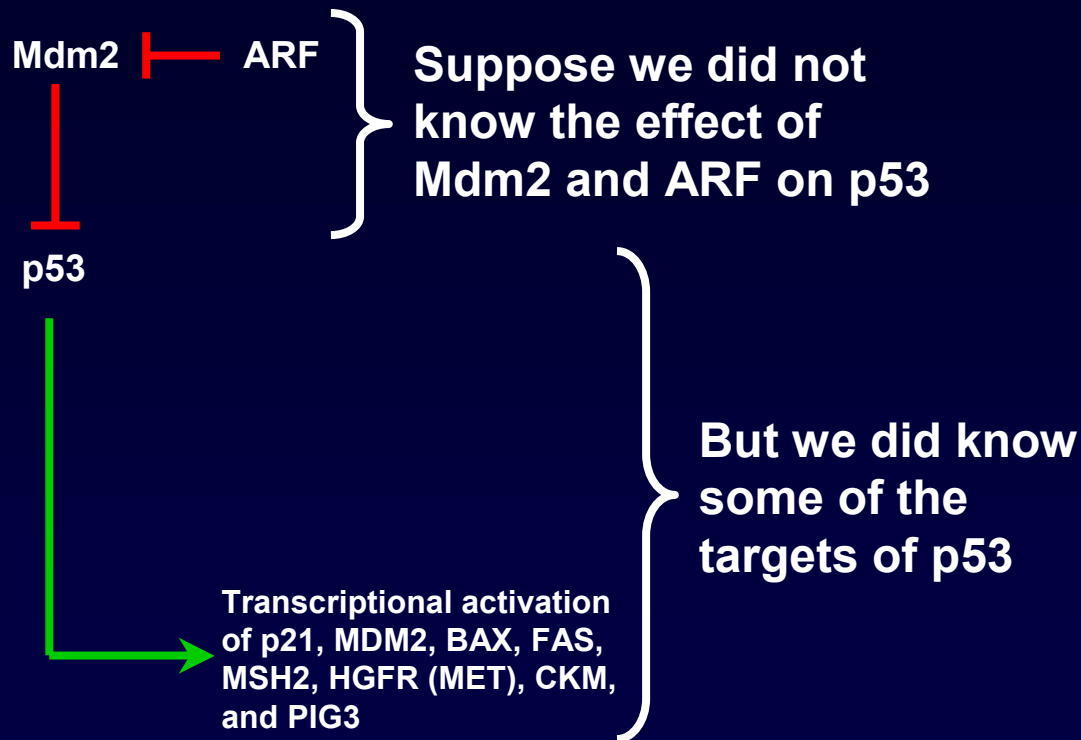
Small number of samples

- ◆ Typically 5 to 50 cell lines, time points, or tissue samples
- ◆ Ratio of measurements to samples can be as bad as 10^3

It is often impossible to make a rigorous quantitative conclusion in such cases.

Explicit use of orthogonal knowledge sources to constrain your questions makes it possible to derive quantitative conclusions.

Array-based data can potentially accelerate the derivation of biological pathways



We have a collection of cell lines and collect the following information

- ◆ P53 mutational status
- ◆ Temporal expression data

Why are the targets of p53 non-responsive to wildtype p53 in some cell lines?

Look for genes whose expression is positively or negatively correlated with expression of p53 targets in the subset of cell lines with WT p53.

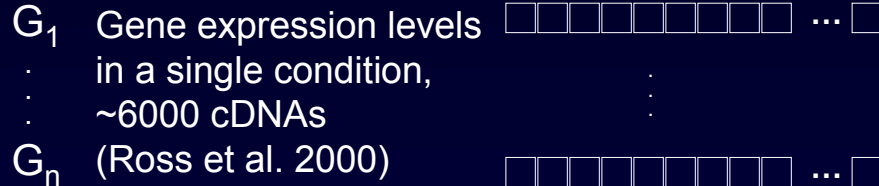
Find that mdm2 is negatively correlated with p53 targets and ARF is positively correlated

In order to do this, we need to integrate experimental data with systematized biological knowledge.

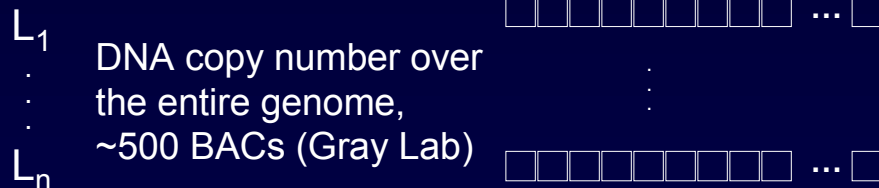
The NCI60 cell-lines as a complex system: Each cell line is a (large) perturbation

RNA

60 tumor-derived cell lines



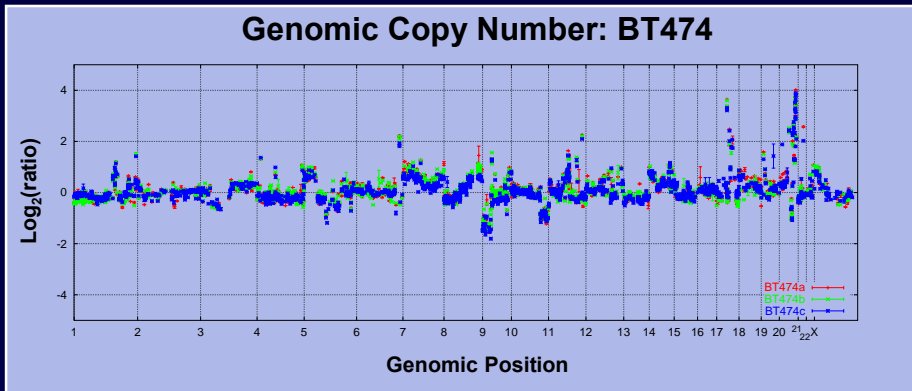
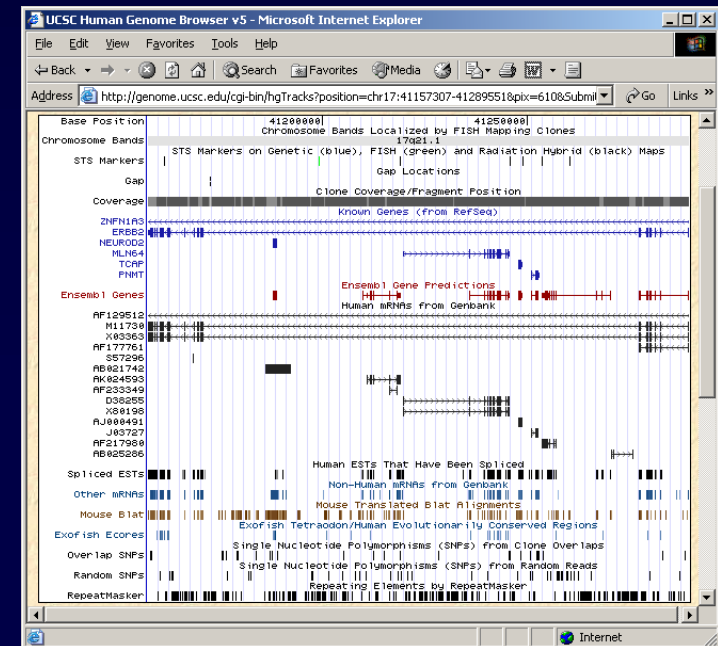
DNA



Gene Annotations

ERBB2:
EC Number: 2.7.1.112
oncogenesis
cell proliferation
Neu/ErbB-2 receptor
protein phosphorylation
protein dephosphorylation
cell growth and maintenance
receptor signaling tyrosine kinase

Genomic Mapping + Context



The signals are clearly different in the CGH and expression data

Ovary
Leukemia

Melanoma

CNS

Renal

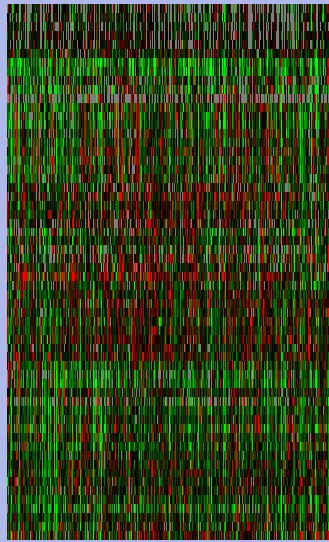
Colon

Lung

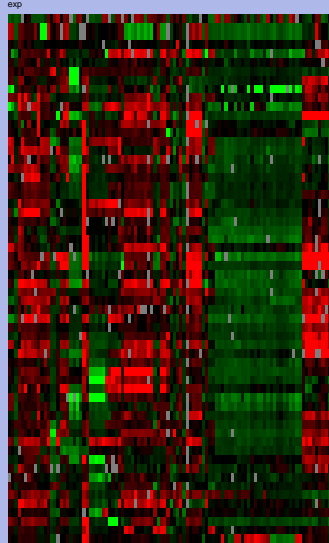
Melanoma

Leukemia

Leukemia_K562_D280aM
Leukemia_HL60_D437z
Melanoma_LOXIMV1_D308a
Prostate_PC3_D338a
Ovary_OVCAR4_D322z
Ovary_OVCAR3_D322z
Leukemia_MOLT4_D437z
Leukemia_CCRF_D280aM
Leukemia_RPM8226_D281aM
Leukemia_SR_D374a
Lung_H460_D374b
Melanoma_SK-MEL5_D314c
Breast_MDA435_D266a
Breast_MDAMB231_D266a
Melanoma_M14_D308b
Melanoma_UACC267_D315a
Melanoma_MALME3M_D316b
Melanoma_SK-MEL28_D316a
Melanoma_SK-MEL2_D314c
Melanoma_UACC262_D315b
Breast_BT549_D439b
CNS_SNB75_D337b
CNS_SNB19_D337b
CNS_U251_D337a
CNS_SF539_D335a
Renal_SNH2C_D376a
Ovary_OVCAR6_D214b
Breast_MDA231_D213c
Lung_HOP92_D414a
CNS_SF539_D335a
Breast_HS578T_D268a
CNS_SF268_D334a
Lung_HOP62_D414a
Renal_UO31_D376b
Renal_ACHN_D365a
Renal_TK10_D367a
Renal_786-O_D364a
Renal_RXF393_D375b
Renal_CAK1_D365a
Renal_A498_D364b
Breast_T47D_D212z
Breast_MCF7-A_D212a
Breast_MCF7_D213a
Lung_H322M_D318aM
Colon_H129_D211b
Colon_HCT15_D263b
Colon_HCC2998_D263a
Colon_Colo205_D265a
Colon_KM12_D265b
Colon_SW620
Colon_HCT166_D211a
Ovary_OVCAR5_D317b
Prostate_DU145_D338b
Lung_A549_D438a
Lung_EKVX_D377b
Lung_H23_D438b
Lung_H522_D211aM
Ovary_IGROV1_D317a
Ovary_SKOV3_D439b
Lung_H226_D375a



Expression



Array-CGH

Is gene expression related to DNA copy number?

Are genes within similar pathways expressed similarly?

Are genes that are co-expressed similar in their regulatory sequences?

The dominant signal in the expression data is tissue of origin.
The dominant signal in the copy number data is not.

We expect that there will be effects of gene dosage on expression

Direct effect

- ◆ A gene that is on a genomic region represented 8 times in the genome of a particular cell line should have higher expression than in a cell line with normal copy number
- ◆ A gene that is homozygously deleted should not be expressed

Indirect effects

- ◆ A genomic locus that is amplified as a result of selective pressure might have specific effects on many genes' expression

Cannot find this using direct correlation and permutation analysis

- ◆ 1.5 million correlations
- ◆ Nothing is nominally significant

We can show that there is a relationship between gene expression and copy number by considering sample/sample distances (Mantel statistic)

- ◆ Compute 60x60 sample/sample distances using expression vectors
- ◆ Compute 60x60 sample/sample distances using copy number vectors
- ◆ There is a correlation, but it is not stunning

How can we exploit annotation information to look at the direct effects? **Subsets!**

We can look at the direct effect by considering subsets

- ◆ Consider the set of genes that map to a particular genomic position
- ◆ Consider the set of BACs that map to the same place
- ◆ Are those genes' expression correlated with copy number at those loci?

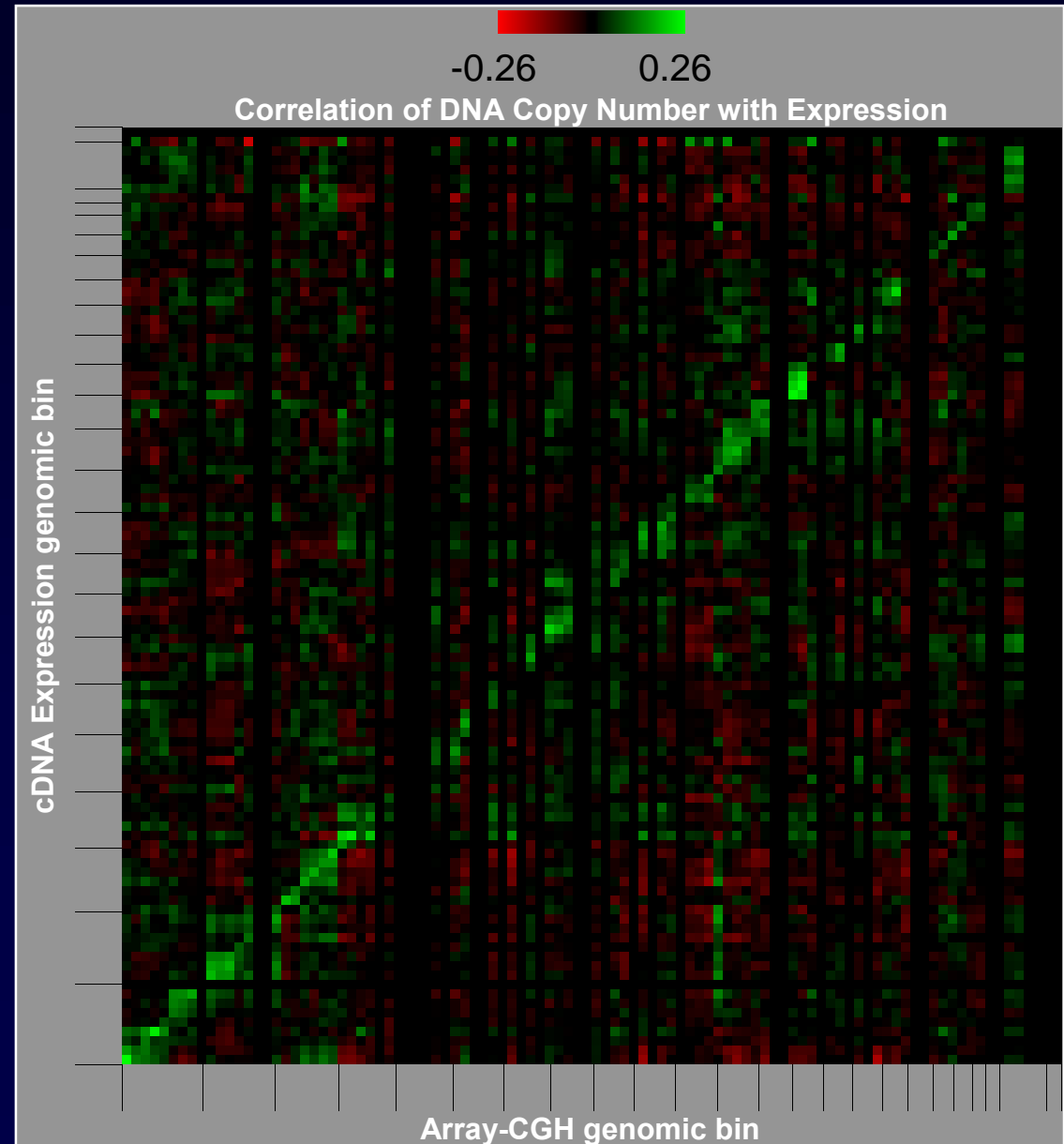
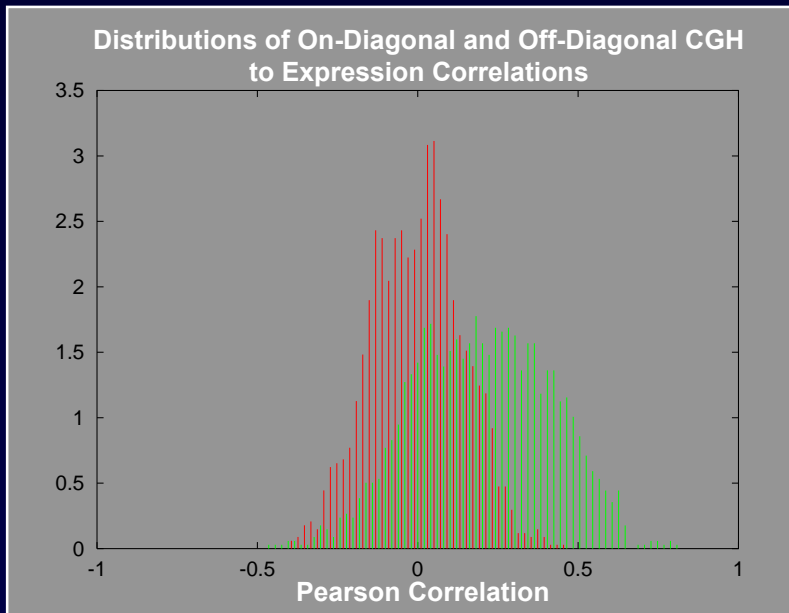
Comparing subsets completes the argument

- ◆ Consider the set of gene/locus pairs that map within 1 Mb of one another
- ◆ Consider the set of gene/locus pairs that map greater than 50 Mb apart
- ◆ Are the correlations from (1) higher than from (2)?

In order to do this, one must accurately map both cDNAs and BACs to genomic sequence.

Subsets based on mapping: Genome-wide gene expression, on average, correlates with genomic copy number

The close-mapping pairs have significantly higher correlations than the distant-mapping pairs



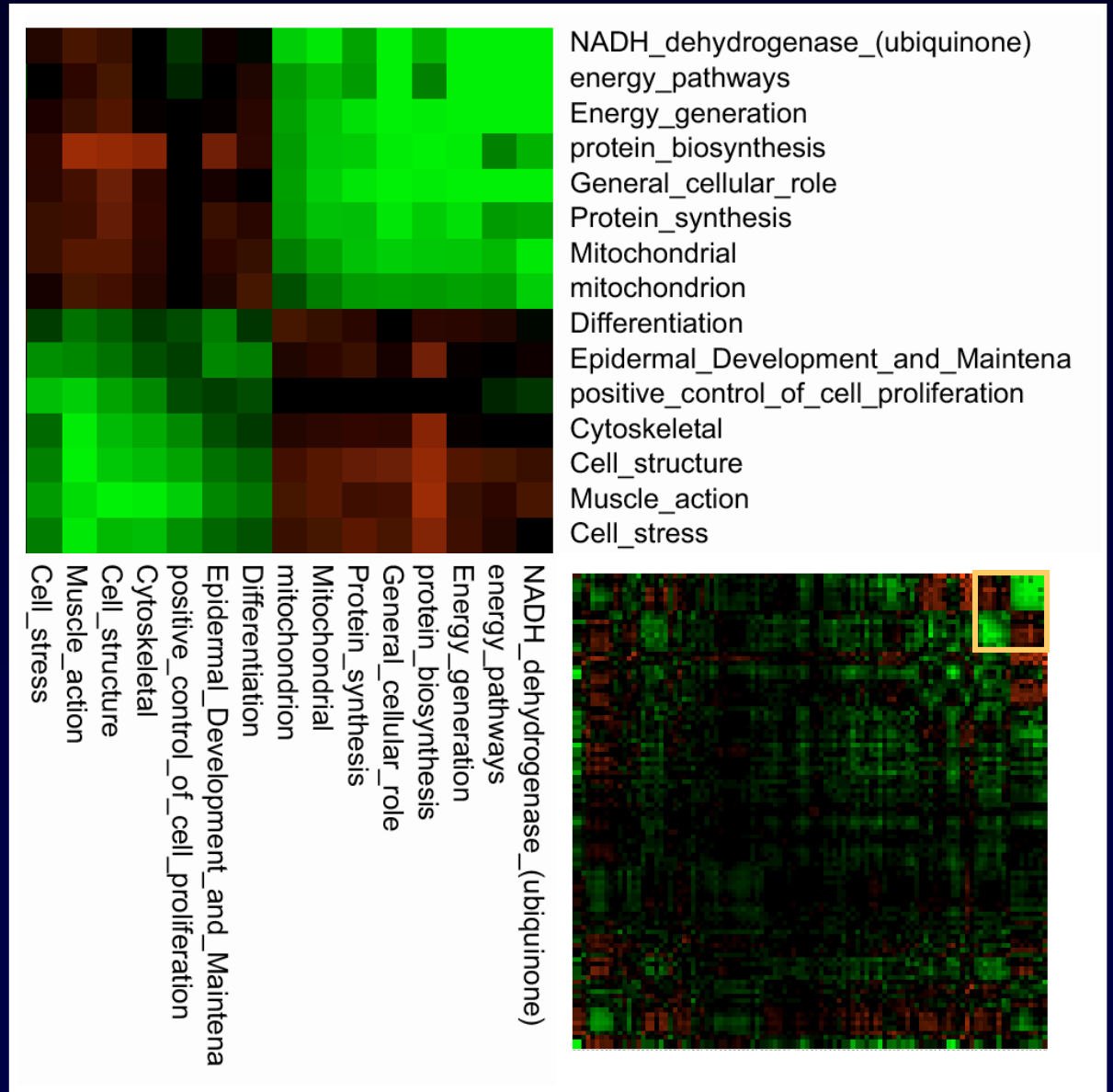
Subsets based on gene function: Gene Ontology information as a surrogate for pathway reveals enrichment of co-expression

Map cDNAs to curated NCBI RefSeqs

Use Gene Ontology and other controlled vocabularies as annotations of the genes

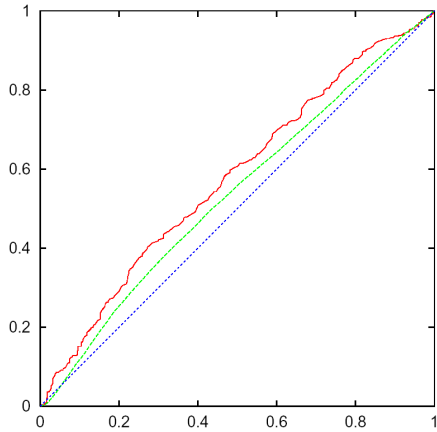
Expression of sets of genes that share annotations is correlated

- ◆ Particularly in basic cellular metabolic systems
- ◆ Note: we have eliminated all identity correlations. Each bin at right contains the average of at least 50 non-degenerate correlations.



Subsets based on gene regulatory sequence similarity:

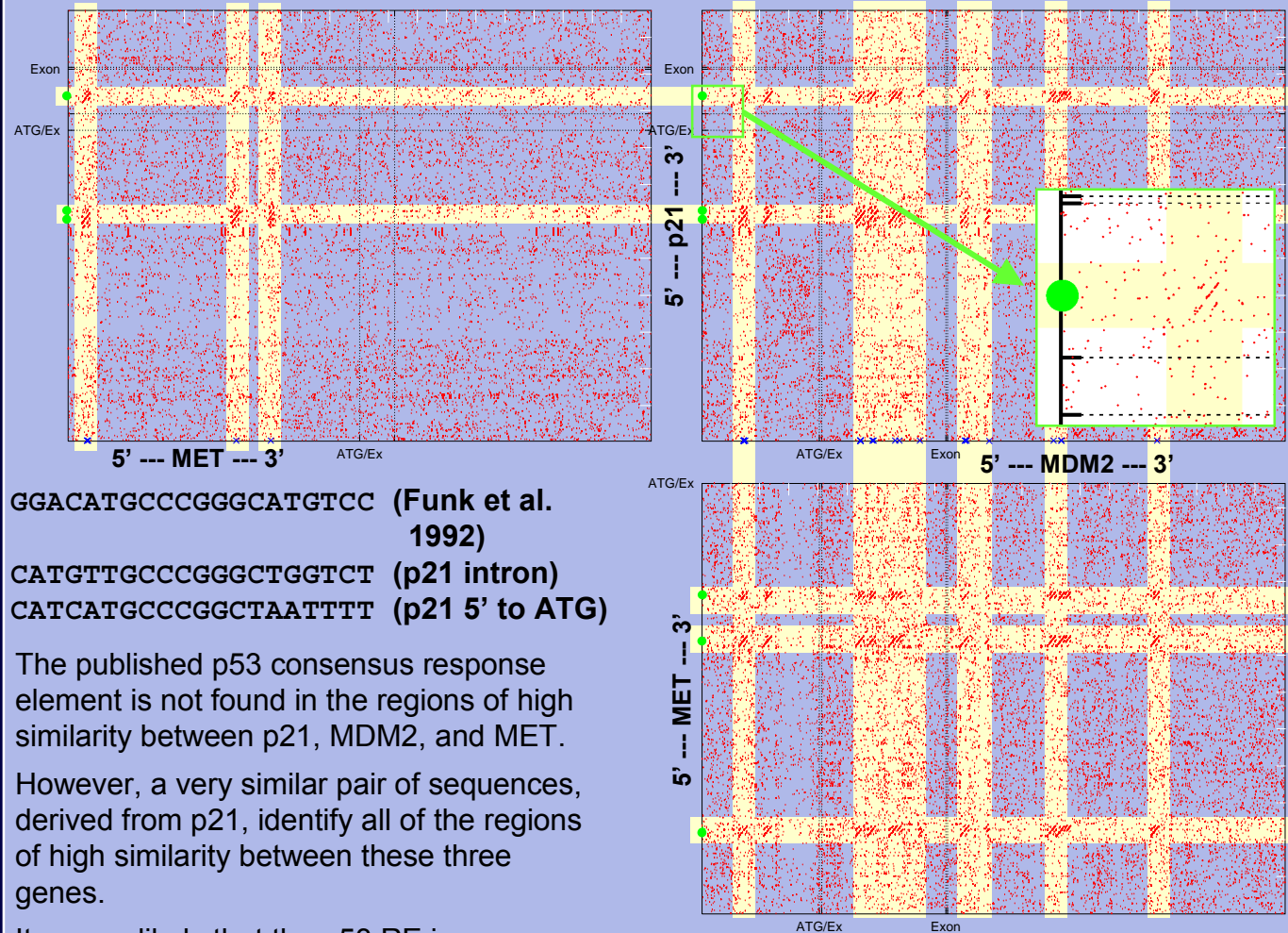
Genes pairs that are co-expressed have more similar regulatory regions than other pairs



ROC curve depicting separation of compositional similarity of co-expressed gene pairs from non-co-expressed gene pairs (green)

ROC curve depicting separation of expression correlation for gene pairs with high compositional similarity versus low similarity (red)

p21, MDM2, and MET: p53 Response Elements



GGACATGCCCGGGCATGTCC (Funk et al. 1992)

CATGTTGCCCGGGCTGGTCT (p21 intron)

CATCATGCCCGGCTAATTTT (p21 5' to ATG)

The published p53 consensus response element is not found in the regions of high similarity between p21, MDM2, and MET.

However, a very similar pair of sequences, derived from p21, identify all of the regions of high similarity between these three genes.

It seems likely that the p53 RE is more complex than has been proposed. We are exploring the specificity of these sequences.

Do specific DNA copy number abnormalities bear on patient outcome or tumor phenotype? **Yes**

Is gene expression, genome wide, on average, quantitatively related to genomic copy number? **Yes**

Are the regulatory regions of gene pairs that share patterns of expression more similar than those that don't? **Yes**

Are the patterns of expression of gene pairs that have high similarity in regulatory regions more concordant than those with low similarity? **Yes**

We want to be able to answer complicated questions about biology

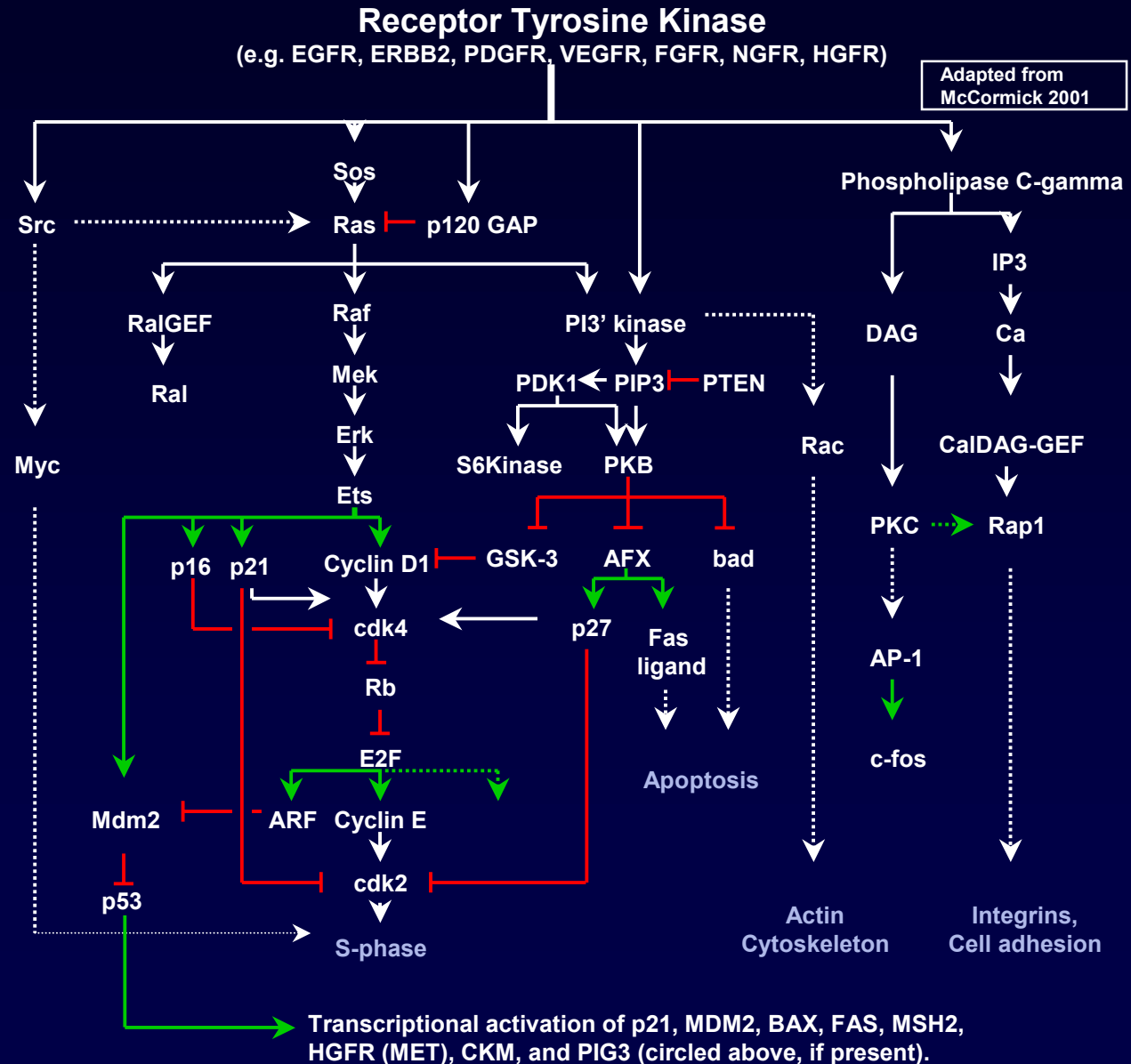
Are the genome copy number patterns of genes that impinge on S-phase checkpoint control quantitatively related?

Are other genes related in their pattern of aberration?

Can the context of the RTK pathway help in the analysis?

Bladder tumor data

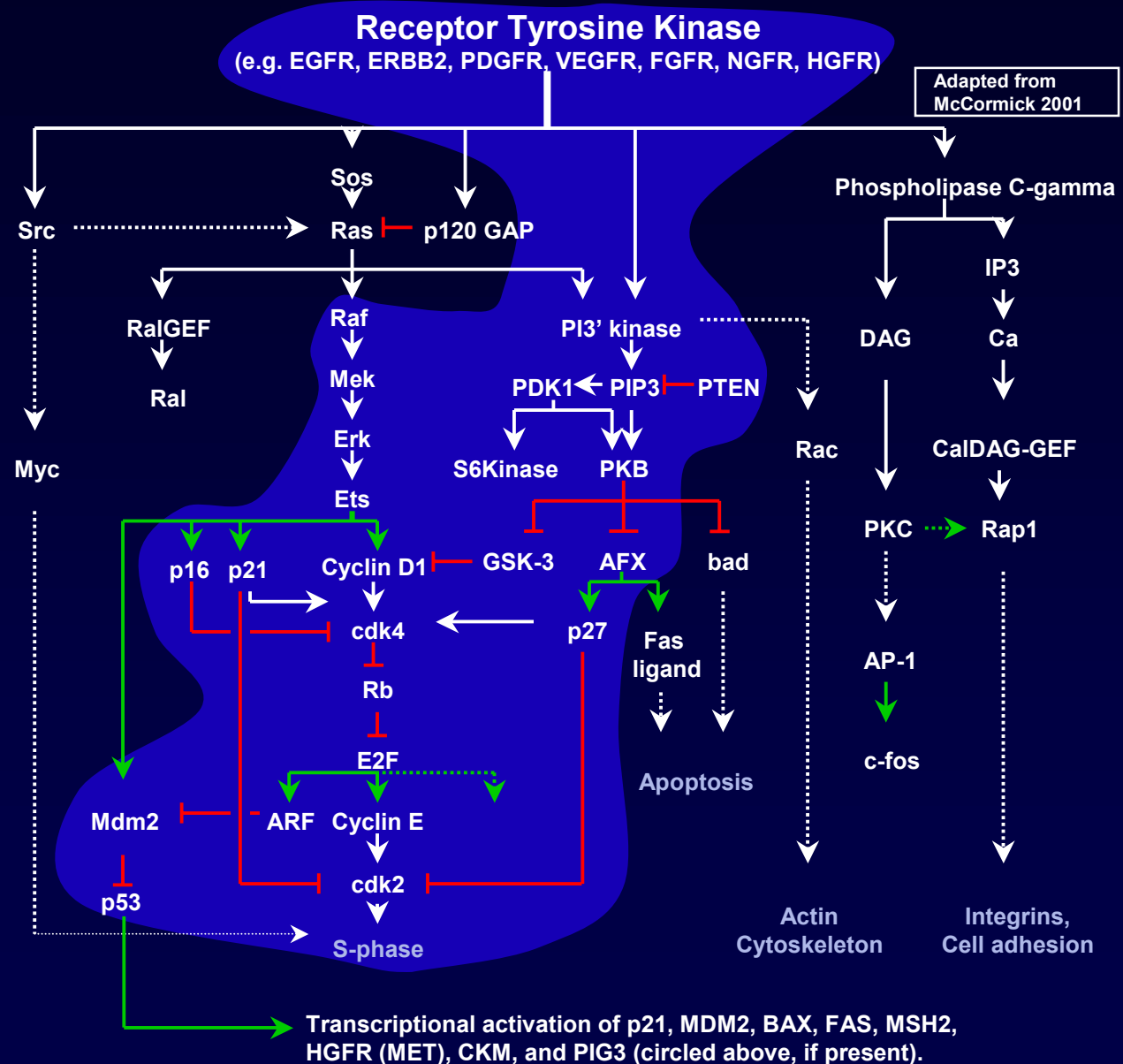
- ◆ Waldman Lab (Joris Veltman)
- ◆ 41 tumors (9 Ta, 7 T1, 25 T2-4)
- ◆ ArrayCGH, both high-resolution (2000 clones) and oncogene focused arrays (500 clones).

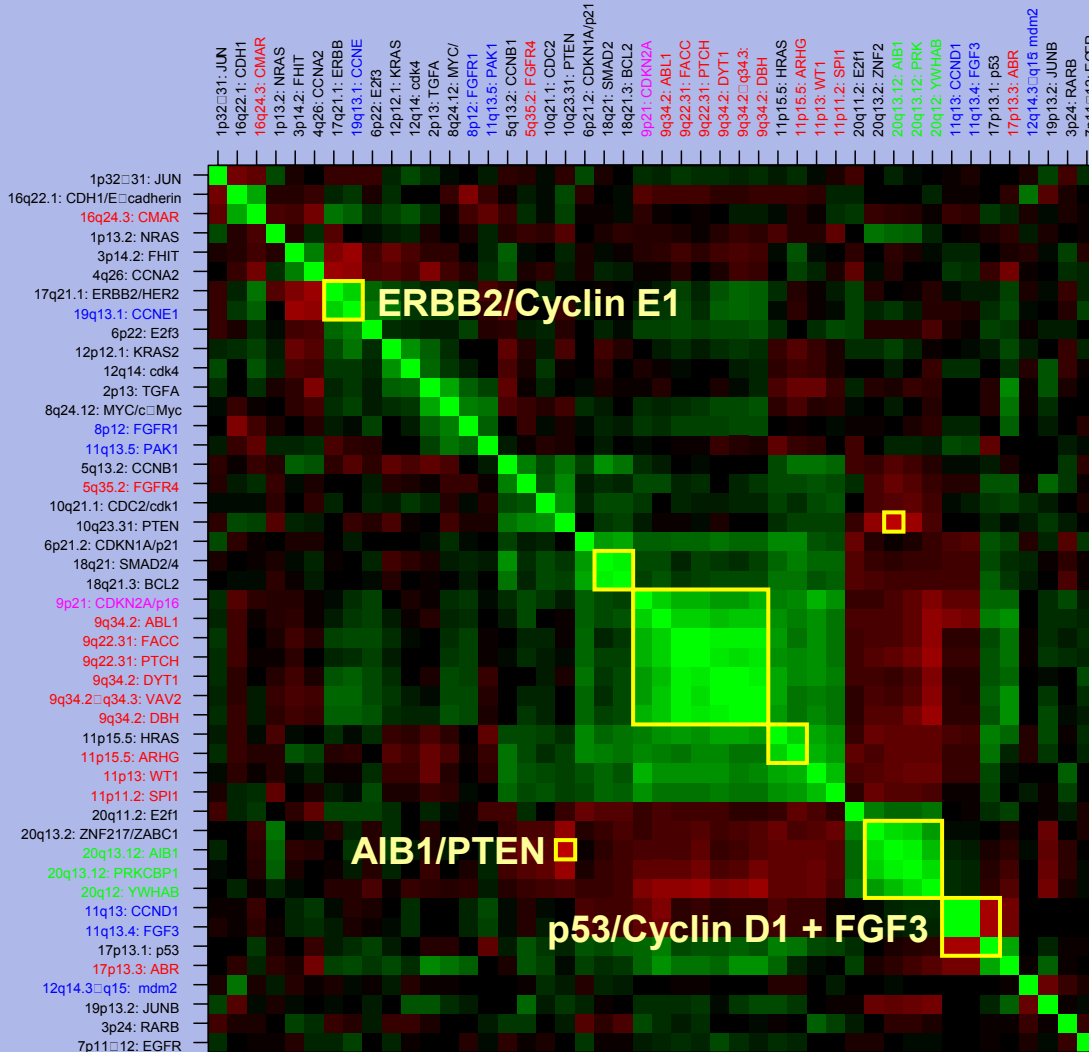


We want to be able to answer complicated questions about biology

Gene/gene relationships in S-phase checkpoint control:

- ◆ Set 1: all genomic clones in the array-CGH experiment that are connected within 7 steps to S-phase control
- ◆ Set 2: all genomic clones that are frequently amplified or deleted
- ◆ Compute the correlation between copy number patterns of all gene pairs.
- ◆ Significance quantified by permutation analysis.





Gain of ERBB2 (17q12) and gain of CCNE (19q13.11)

Gain of AIB1 (20q13) and loss of PTEN (10q23)

Loss of p53 (17p13.3) and gain of CCND1 (11q13)

Loss of p53 (17p13.3) and gain of FGF3 (11q13)

Bladder tumor data

- ◆ Waldman Lab (Joris Veltman)
- ◆ ArrayCGH, both high-resolution (2000 clones) and oncogene focused arrays (500 clones).

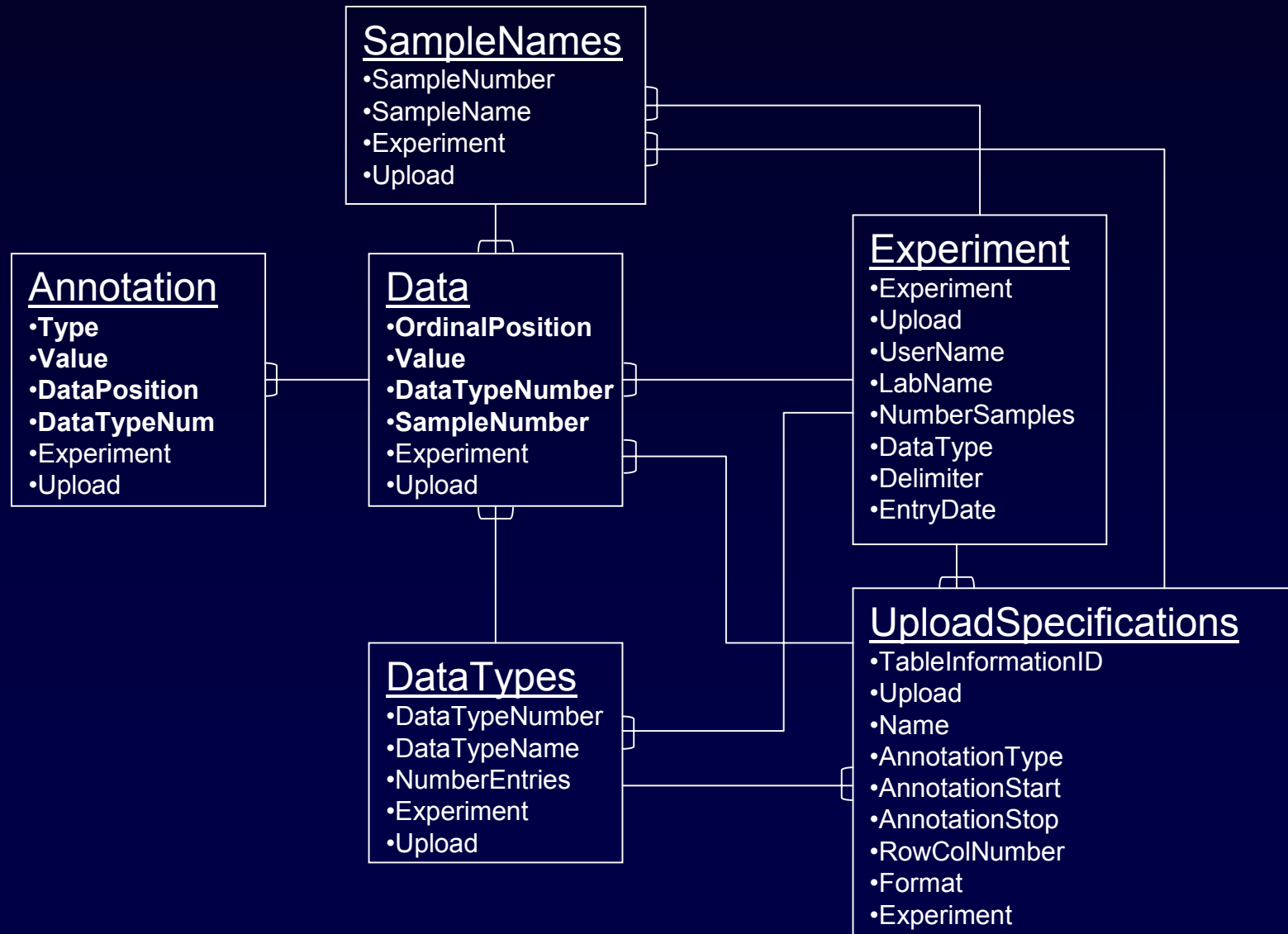
We have adopted a very flexible data model: Entity attribute value

User can define essentially any data type

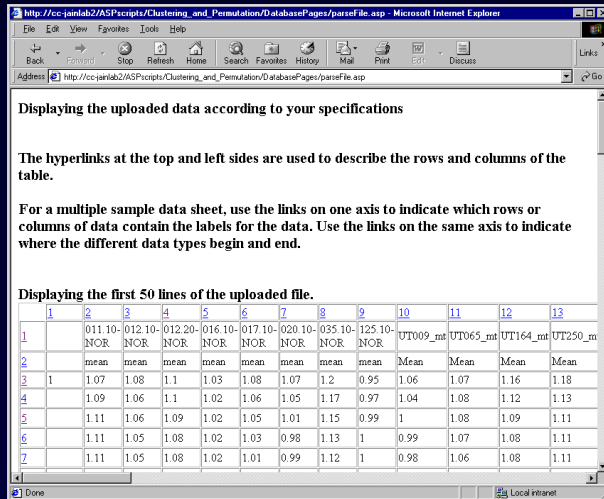
No restriction on quantified biological data

Same basic model can be used for representing free-form user annotation data

However, for the highly structured curated information, we are adopting more specialized models



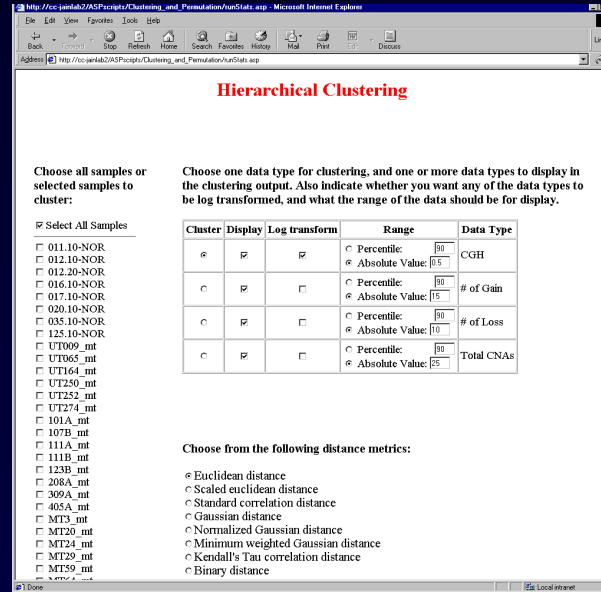
We are building a web-based data system that embeds these analysis and visualization tools



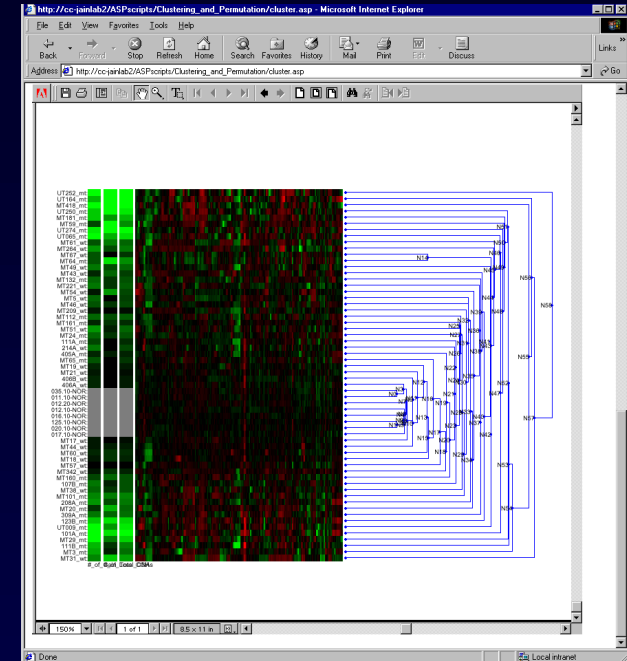
Choose a file to upload to the Web server.

Identify the data types and their positions within the file.

Data can be from any experiment of any type. The only restriction is that it can be represented as tab-delimited text for import to the system.



Several methods of analysis are available for selection. We have displayed hierarchical clustering. The user can select samples, distance metrics, and additional information to display.



The system generates a PDF file that is displayed to the user's browser.

Biology's shift towards being a quantitative molecular science requires fundamentally new analytical methods

The problems inherent in array-based data are not insurmountable: it is often possible to derive quantitatively supportable conclusions

Integration of experimental data with systematically represented biological knowledge (annotations) can reveal relationships otherwise impossible to see by supporting definition of meaningful subsets

Experimental collaborators

- ◆ Albertson Lab
- ◆ Collins Lab
- ◆ Gray Lab
- ◆ Pinkel Lab
- ◆ Waldman Lab

Jain Lab

- ◆ Jane Fridlyand, PhD
 - ◆ Lawrence Hon
 - ◆ Chris Kingsley
 - ◆ Barbara Novak
 - ◆ Taku Tokuyasu, PhD
- UCSF Biological and Medical Informatics (BMI) PhD Students
- ◆ Adam Olshen, PhD
[Now faculty at Sloan-Kettering]