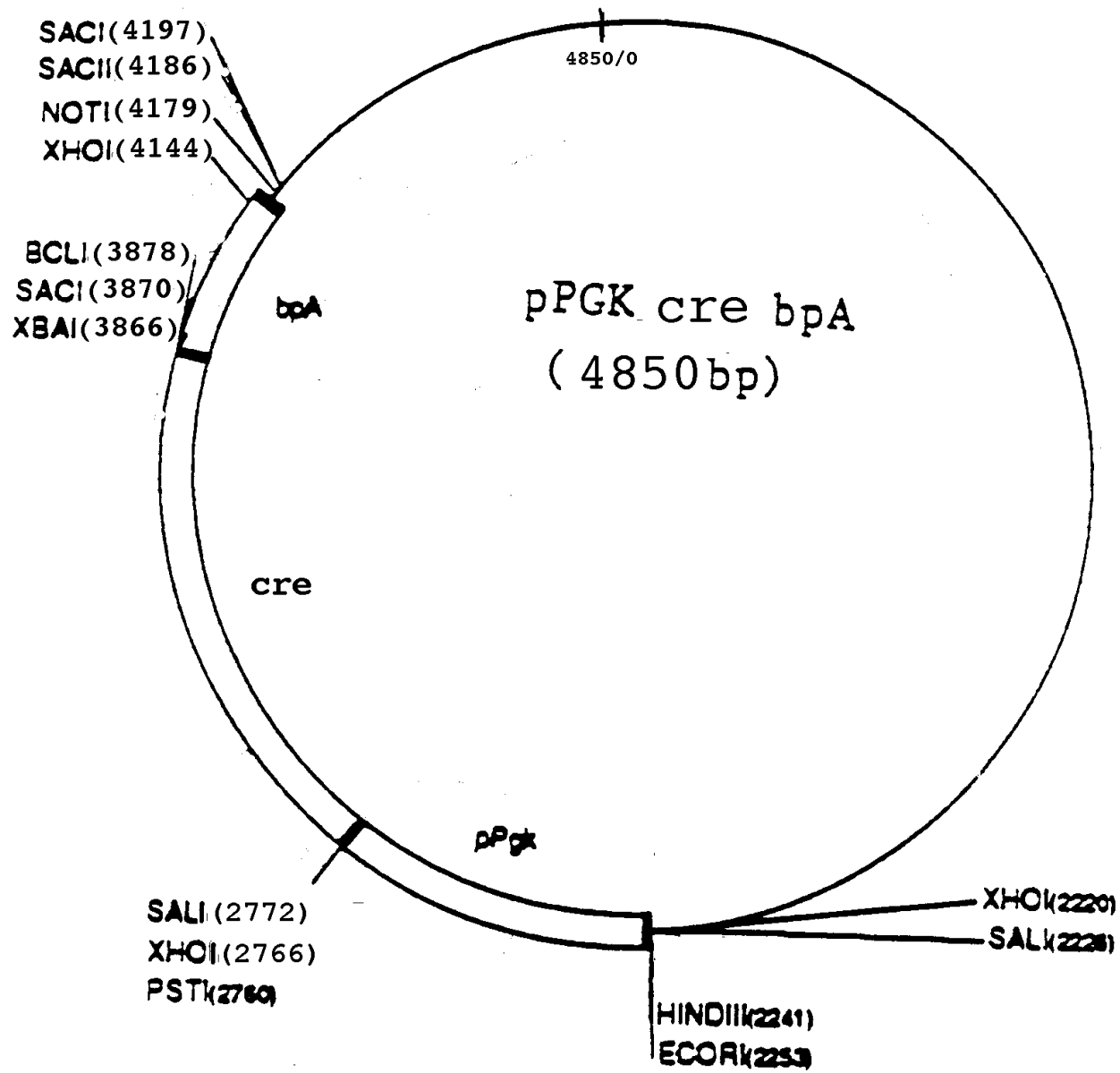# 1. Why high-throughput (HT) data?
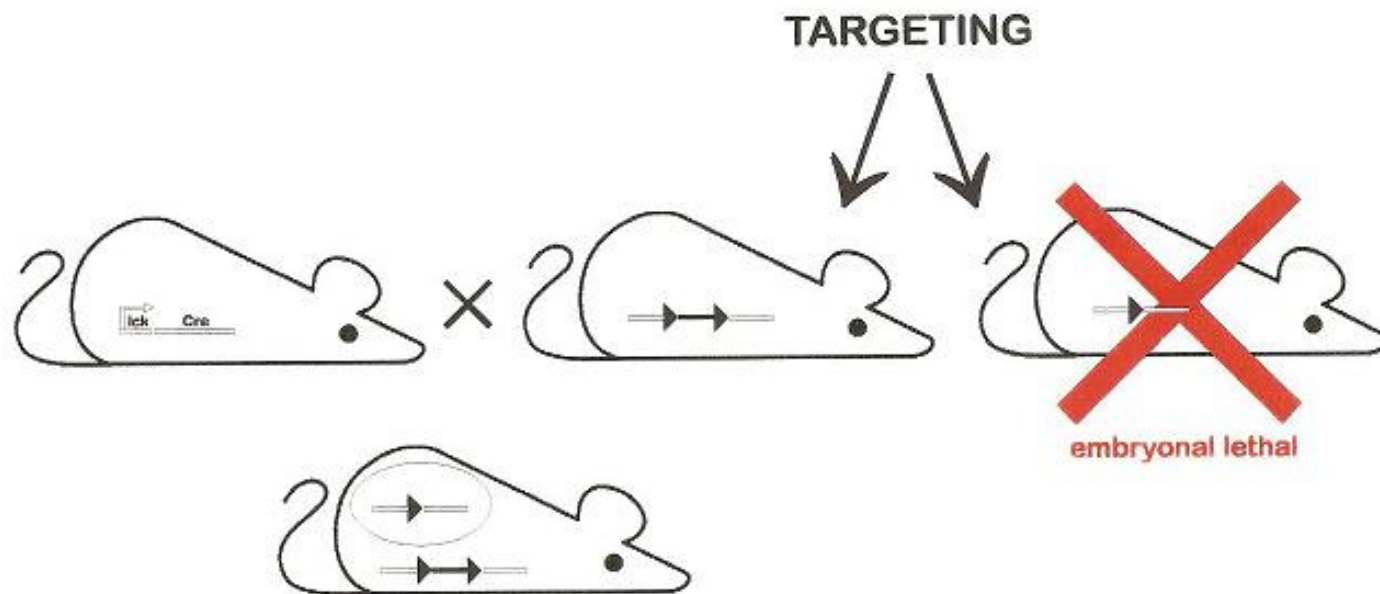
What is it that we cannot achieve by looking at only a few genes at a time?
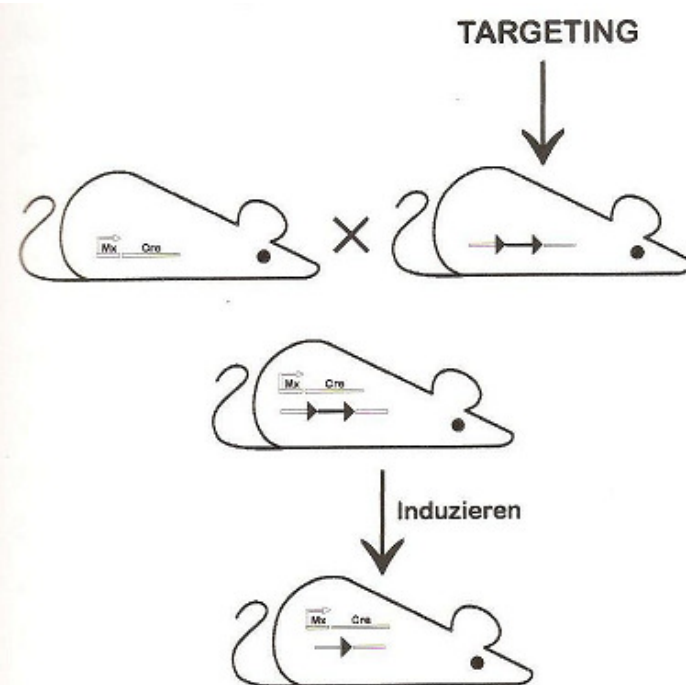
# Why high-throughput (HT) data?

- 1st example:

  Studying e.g. immunological coherences

  by knocking out genes

# cre for eukaryotic expression



SACI(4197)
SACII(4186)
NOTI(4179)
XHOI(4144)

BCLI(3878)
SACI(3870)
XBAI(3866)

bpA

cre

pPGK cre bpA
( 4850bp)

4850/0

pPgk

SALI(2772)
XHOI(2766)
PSTI(2760)

HINDIII(2241)
ECORI(2253)

XHOI(2220)
SALI(2223)

TARGETING

lck  Cre

×

embryonal lethal

**TARGETING**

Der interferonabhängige Mx-Promotor dient hier zur induzierbaren Ausprägung von Cre (Kuhn et al., 1995). Bei anderen auf Transkriptionsebene kontrollierenden Systemen kann Cre sowohl induzierbar als auch zelltypspezifisch ausgeprägt werden.

Induzieren

Abb. 1–3: Mx-cre System (Kuhn et al., 1995).

# Often:

✓ make construct ("flox" gene)
✓ transfect ES cells
✓ transfer to blastocyst
✓ transfer to foster mother
✓ get germline transmission
✓ cross with cre mice
✓ induce
✓ get considerable knockout rate

at least
½ year
of work

☹ observing no phenotype

That sucks!

# Why is that?

- Function of important genes is backuped by other pathways!

- In order to understand how that works
  (→ systems biology)
  one needs to know
  the status of **<u>many</u>** genes

# Why high-throughput (HT)?

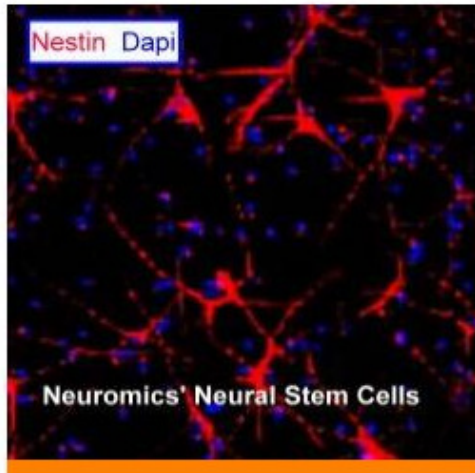- 2nd example:

  cancer

- same problem:

- Cancer is a genetic disease
  - Not monogenic like MD or CF, but multigenic
- Cancer is caused by mutations in somatic cells
- Cancer can be caused by mutagens, chemicals that damage DNA, or viruses
- Cancer is caused by an accumulation of mutations in different genes in a single cell
- Cancer is caused by altered expression of genes or by accumulation of mutations in a single cell

- There are five major pathways that must be activated or inactivated in a cell for the cell to become a cancer cell
  - Growth stimulatory signals
  - Growth inhibitory signals
  - Apoptosis resistance
  - Infinite proliferative capacity
  - Angiogenic potential

# Why is that? <span style="color:red">Again:</span>

- Function of important genes is backuped by other pathways!

- In order to understand how that works
  ($\rightarrow$ systems biology)
  one needs to know
  the status of **<u>many</u>** genes

Neuromics' Neural Stem Cells

**Danish and Belgian researchers have found a computer key that maps genes underlying heritable disorders, such as breast cancer, multiple sclerosis, and Alzheimer's disease. These results will possibly ease the discovery of new medicines and improve treatment in various disorders.**

The results - which are published in the current issue of Nature Biotechnology - show that genes important for the development of diseases like Alzheimer's follow the same cellular rules as genes involved in fundamentally different disorders, such as heart disorders, multiple sclerosis, breast cancer, and Type 2 diabetes.

"Many disorders manifest themselves in fundamentally different ways, but the new surprising discovery is that the underlying genes play together after the same rules. Our results show that the genes that trigger diseases, regardless of the type of disease in question, are social team players who cooperate according to highly specific rules. These rules have now been mapped, and we have pointed at hundreds of new genes that are likely to be involved in disorders including multiple sclerosis, Parkinson, heart disorders, and diabetes", says Kasper Lage from Technical University of Denmark, who is the project coordinator on this work.

Heritable disorders will be easier to interpret for clinicians using the new results. Furthermore, the identification of new genes likely to be involved in disorders will help patients with defects in these genes. For example, if you are a high risk carrier of a gene that underlies a disease such as Type 2 diabetes, physicians could prevent or delay the manifestations of the disease by dietary guidance early in life.

"This is a crucial breakthrough for our understanding of heritable disorders, and a breakthrough for systems biology as a research strategy in the field genetics and disease", says S?Brunak leader of Center for Biological Sequence analysis at the Technical University of Denmark. "We work with genes and proteins, but also with clinical literature describing the characteristics of different disorders. Then we let the computer integrate all of these data, and extract the pattern", he adds.

# What HT data?

What is measured?

How?

- Genomics (DNA, → genome)
  (e.g. by sequencers, microarrays, …)
- Epigenetics (→ e.g. methylome)
  (e.g. by microarrays)
- Transcriptomics (→ transcriptome)
  (e.g. deep sequencing, microarrays, …)
- Proteomics (→ proteome)
  (e.g. 2D-gel electrophoresis, mass spectrometry, guess what – microarrays (AB or peptide chips)

# First microarrays:

- cDNA

- on a nylon membrane

- prepared RNA reversely transcribed into cDNA (like today) ...

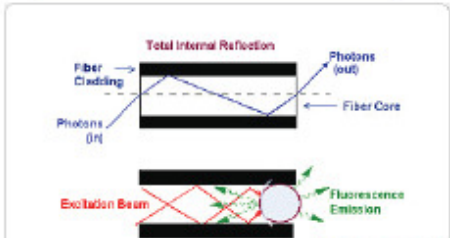- ... using radioactively labelled nucleotides (today: mostly fluorescence labelling)

**SAGE**

2003: Illumina Bead Arrays

1975: Southern Blotting Technology (Edward Southern)

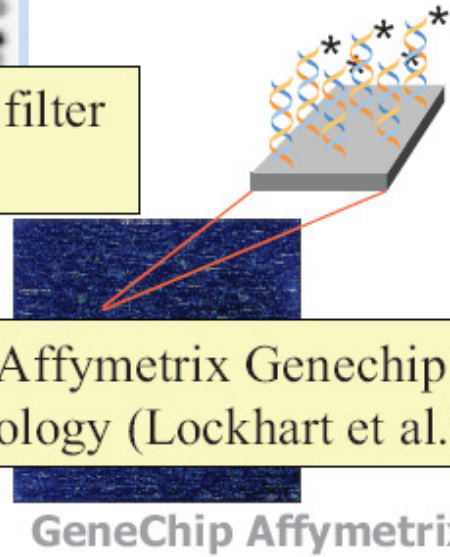1991: First high-density Nylon filter Arrays (Lennon, Lehrach)

**Illumina Bead Array**

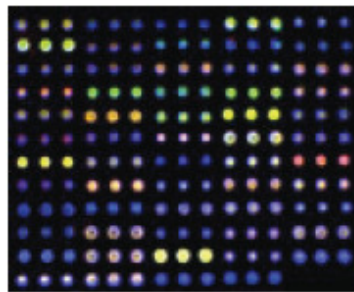Different Technologies for Measuring Gene Expression
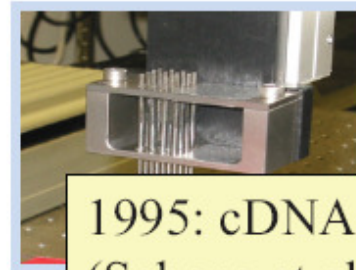
1996: Affymetrix Genechip Technology (Lockhart et al.)
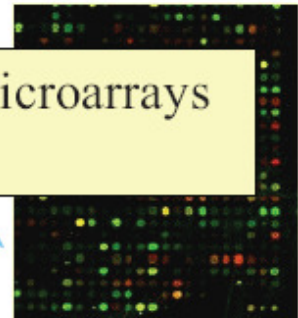
**GeneChip Affymetrix**

**Agilent: Long oligo Ink Jet**

**CGH**

**cDNA microarray**

1995: cDNA-Microarrays (Schena et al.)

two-channel

single channel



cDNA-microarrays

high-density oligonucleotide arrays

probe preparation

cDNA collection

insert amplification by PCR
vector specific primers
gene specfic primers

printing
coupling
denaturing

ratio Cy5/Cy3

mRNA reference
sequence

perfect match
mismatch    probe set

in situ synthesis
by photolithography

array 2
array 1

ratio array 1/array 2

target preparation

hybridization
mixing

Cy3          Cy5

Cy3 or Cy5
labeled cDNA

modified oligodT    cDNA synthesis

total RNA

cells/tissue

staining    fragmented
hybridization

biotin labeled
cRNA
(generated using T7 promoter)

in vitro transcription

double-stranded
cDNA
(generated using T7 primer)

cDNA synthesis

total RNA

cells/tissue

# Microarray Experiment

**Animation:**

http://www.bio.davidson.edu/Courses/genomics/chip/chip.html

# What should it do?

M-CHIPS

TO BIG JOURNALS

**Experimental Cycle**

Biological question
(hypothesis-driven or explorative)

Experimental design

To call in the statistician after the experiment is done may be no
more than asking him to perform a post-mortem examination:
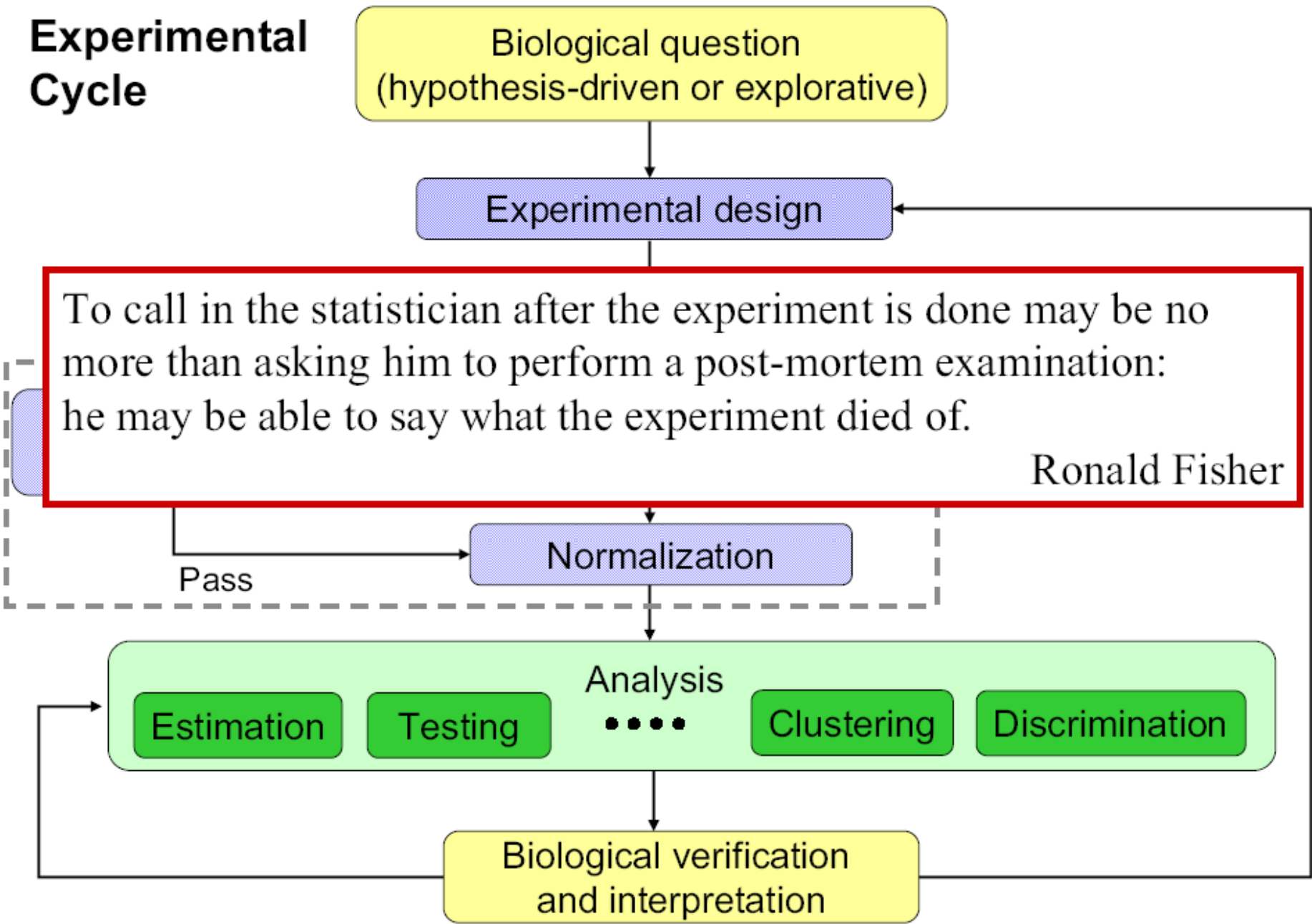he may be able to say what the experiment died of.
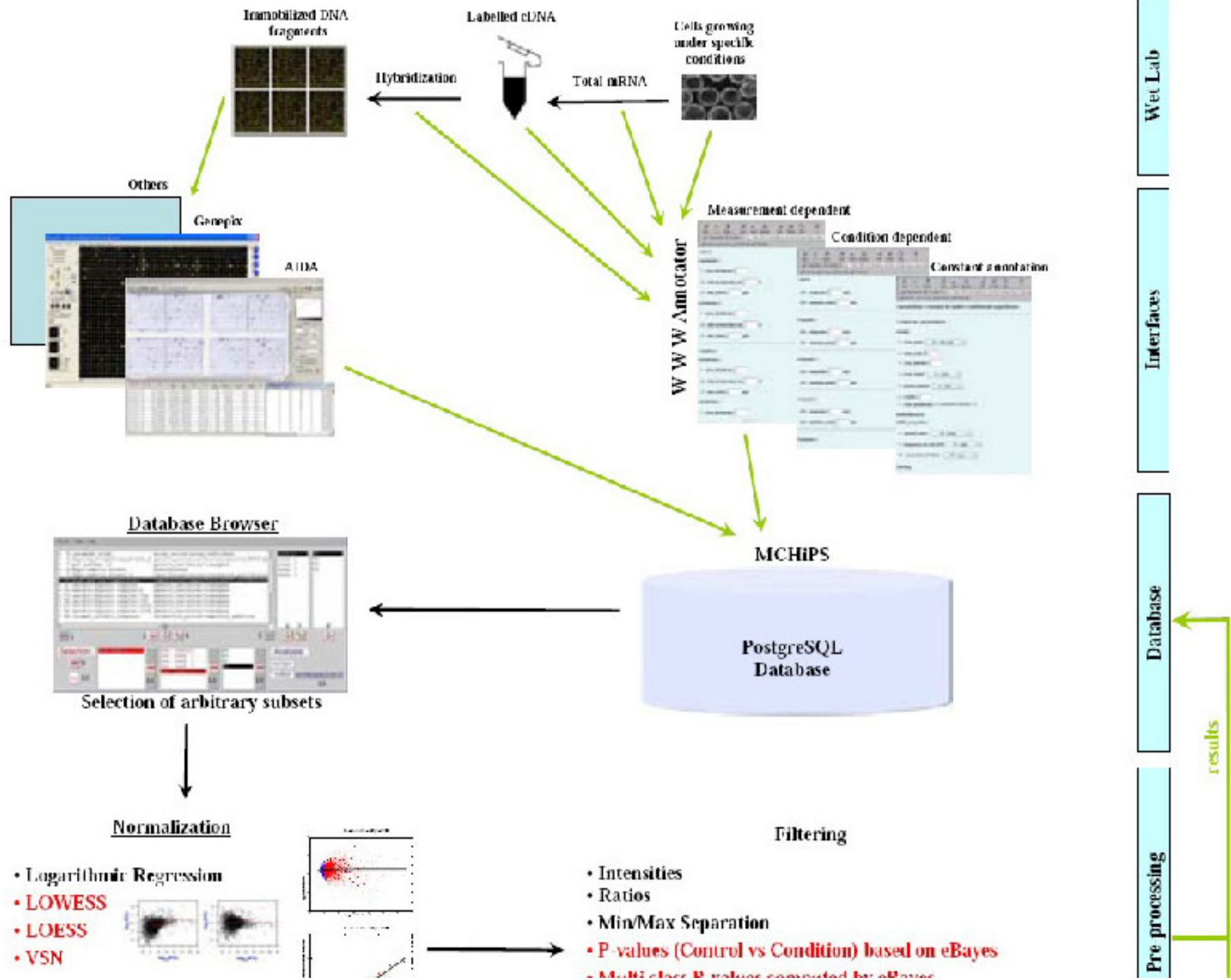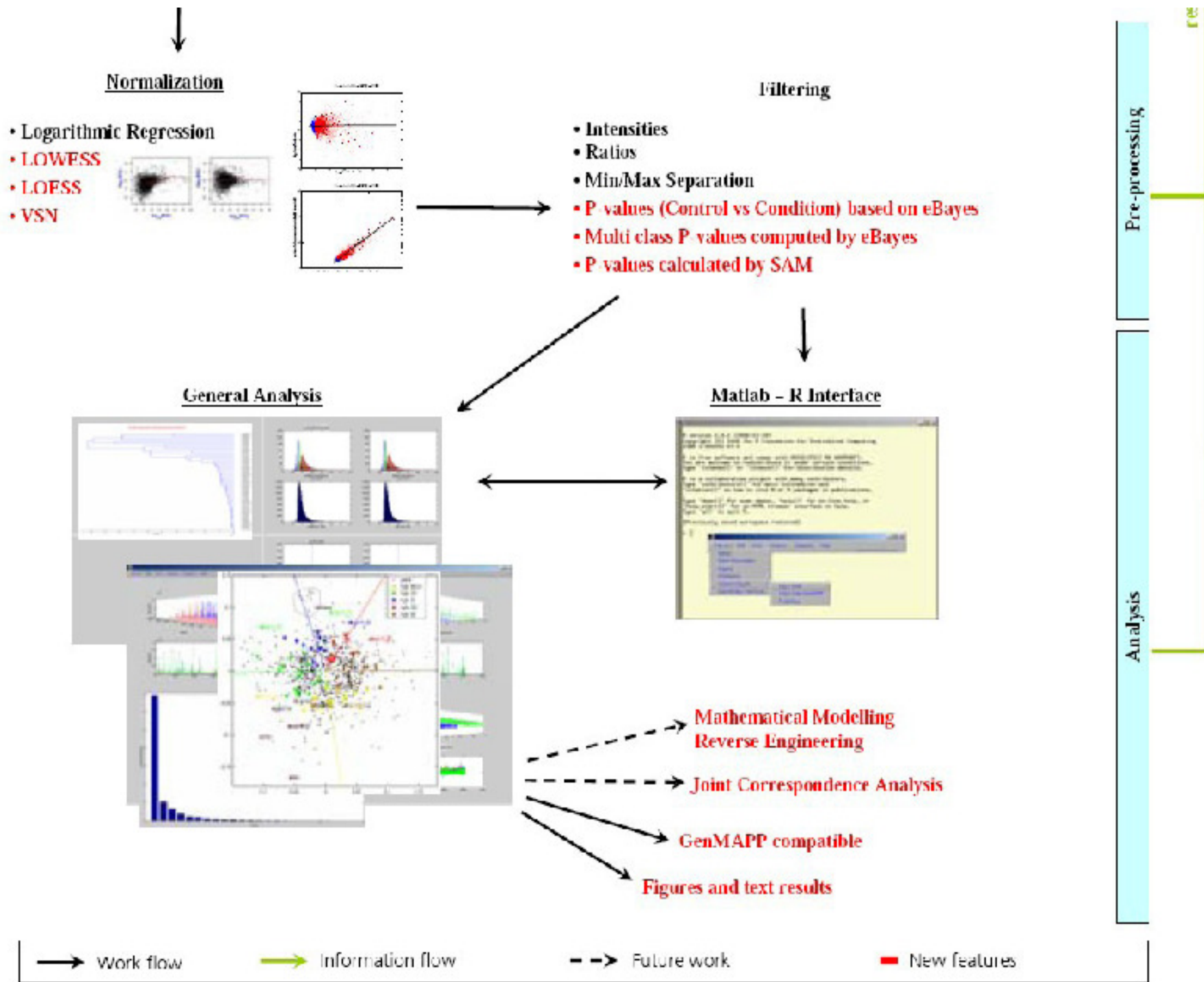
Ronald Fisher

Pass

Normalization

Analysis

Estimation   Testing   • • • •   Clustering   Discrimination

Biological verification
and interpretation

# What is the biggest challenge?

Immobilized DNA
fragments

Labelled cDNA

Cells growing
under specific
conditions

Hybridization

Total mRNA

Others

Genepix

AIDA

Measurement dependent

Condition dependent

Constant annotation

W W W Annotator

Database Browser

MCHiPS

PostgreSQL
Database

Selection of arbitrary subsets

**Normalization**

- Logarithmic Regression
- LOWESS
- LOESS
- VSN

**Filtering**

- Intensities
- Ratios
- Min/Max Separation
- P-values (Control vs Condition) based on eBayes
- Multi class P-values computed by eBayes

Wet Lab

Interfaces

Database

results

Pre processing

**Normalization**

- Logarithmic Regression
- LOWESS
- LOESS
- VSN

**Filtering**

- Intensities
- Ratios
- Min/Max Separation
- P-values (Control vs Condition) based on eBayes
- Multi class P-values computed by eBayes
- P-values calculated by SAM

**General Analysis**

**Matlab – R Interface**

Mathematical Modelling
Reverse Engineering

Joint Correspondence Analysis

GenMAPP compatible

Figures and text results

Pre-processing

Analysis

→ Work flow     → Information flow     --→ Future work     ▬ New features

hybridization

cells growing under
specific conditions

immobilized DNA fragments

labelled cDNA

mRNA

convert into numerical
values

**Intensity Table**

| | control condition | | | condition 1 | | | condition 2 | | |
|---|---|---|---|---|---|---|---|---|---|
| | hybr. 1 | hybr. 2 | hybr. 3 | hybr. 4 | hybr. 5 | hybr. 6 | hybr. 7 | hybr. 8 | hybr. 9 |
| gene 1 | 14,243 | 12,154 | | | | | | | |
| gene 2 | 5,323 | 27,152 | | | | | | | |
| gene 3 | 10,300 | 1,407 | ... | | | | | | |
| gene 4 | 1,007 | 3,101 | | | | | | | |
| gene 5 | 100,232 | 120,993 | | | | | | | |
| gene 6 | | | | | | | | | |
| gene 7 | ⋮ | | | | | | | | |
| gene 8 | | | | | | | | | |
| gene 9 | | | | | | | | | |
| ⋮ | | | | | | | | | |

**hybridization**

cells growing under specific conditions

mRNA

labelled cDNA

immobilized DNA fragments

convert into numerical values

**Intensity Table**

| | control condition | | | condition 1 | | | condition 2 | | |
|---|---|---|---|---|---|---|---|---|---|
| | hybr. 1 | hybr. 2 | hybr. 3 | hybr. 4 | hybr. 5 | hybr. 6 | hybr. 7 | hybr. 8 | hybr. 9 |
| gene 1 | 14,243 | 12,154 | | | | | | | |
| gene 2 | 5,323 | 27,152 | | | | | | | |
| gene 3 | 10,300 | 1,407 | ... | | | | | | |
| gene 4 | 1,007 | 3,101 | | | | | | | |
| gene 5 | 100,232 | 120,993 | | | | | | | |
| gene 6 | | | | | | | | | |
| gene 7 | : | | | | | | | | |
| gene 8 | | | | | | | | | |
| gene 9 | | | | | | | | | |
| : | | | | | | | | | |

# Image Analysis – Spot Identification

- The grid structure is provided by the manufacturer or generated individually for custom-made microarrays (e.g. GAL-files)

- The grid is overlaid by hand or automatically onto the image (beware of column/row displacement errors!)

Columns

Rows



Blocks





GAL-file contains Clone-IDs and defines their position on the grid
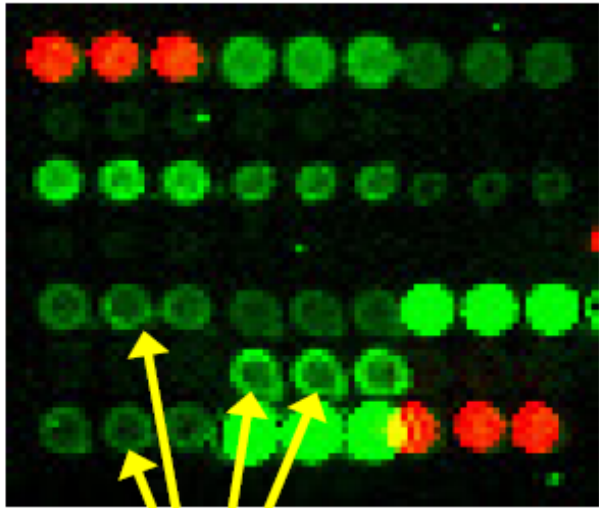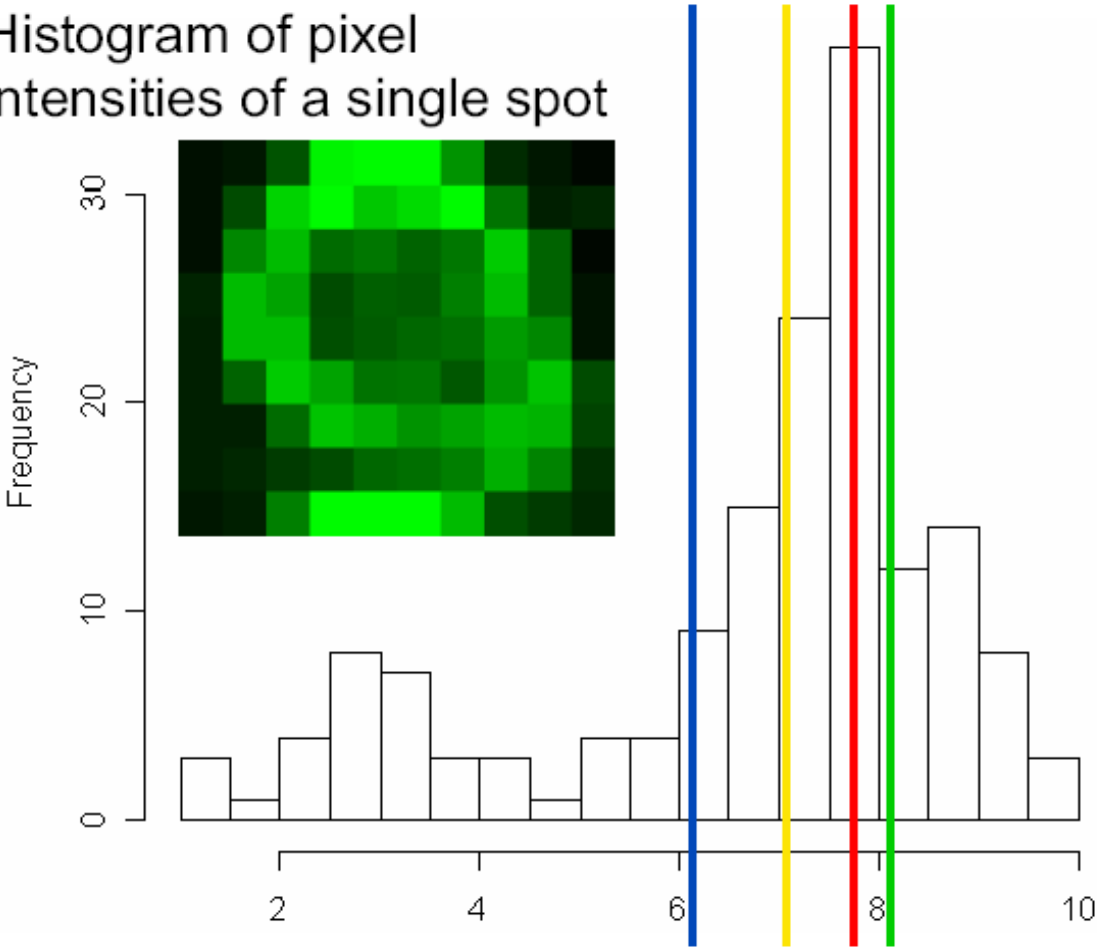
# Spot identification

- The signal of the spots is quantified.
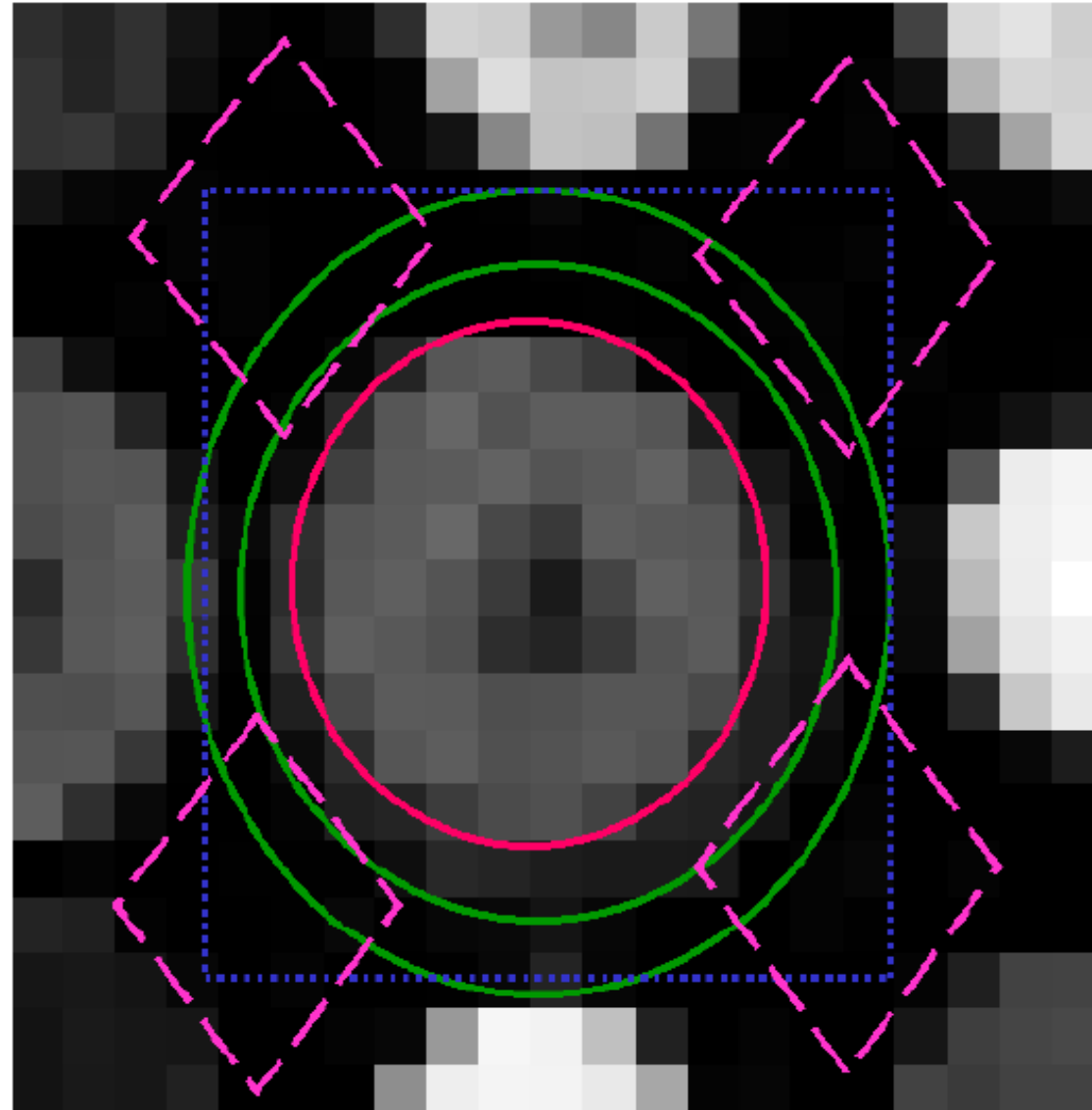


Histogram of pixel intensities of a single spot

Mean / Median / Mode / 75% quantile

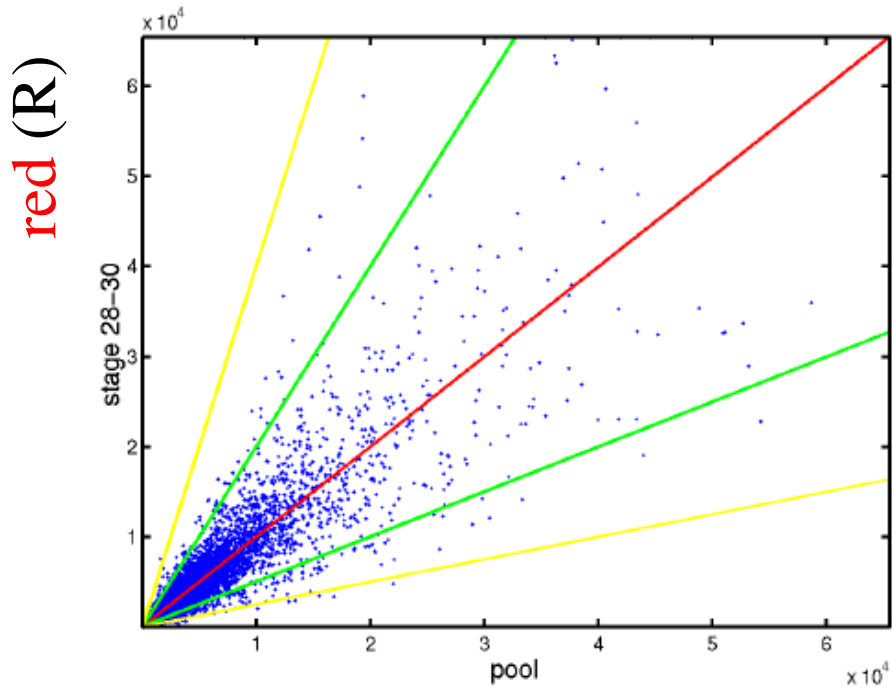"Donuts"

# Different Regions around the spot are quantified to measure local background.
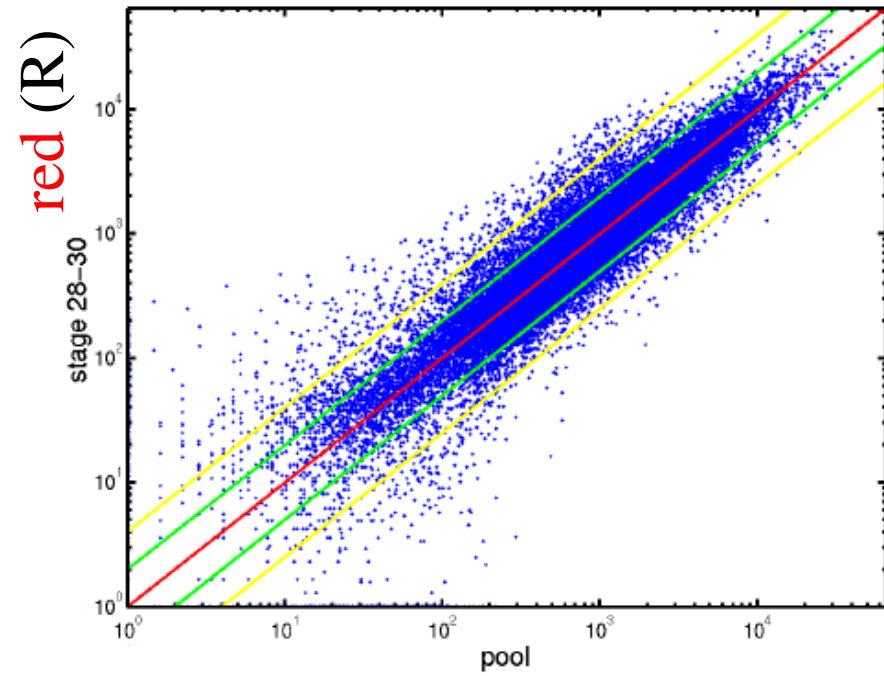


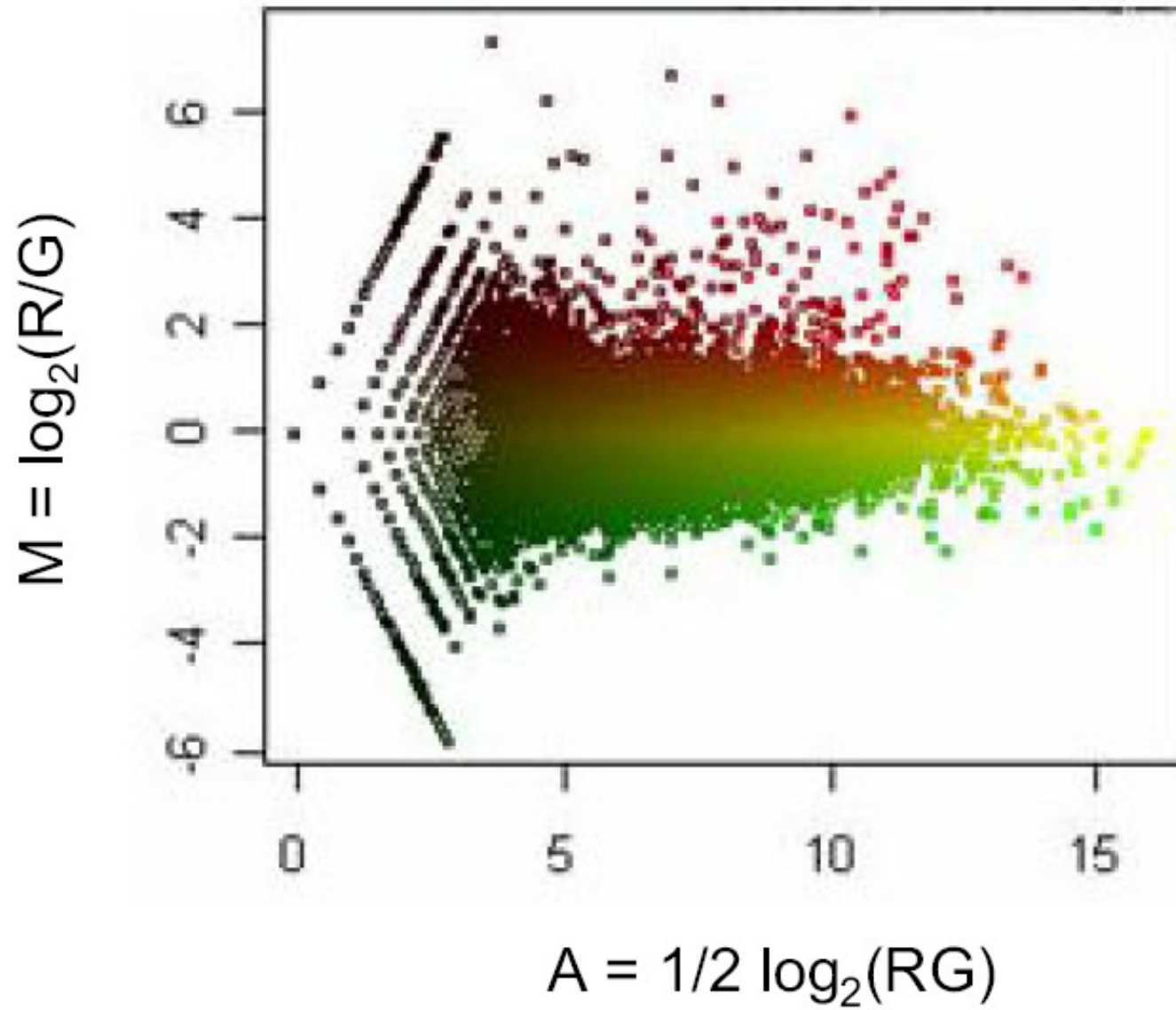**GenePix**

**QuantArray**

**ScanAlyse**

**hybridization**

cells growing under
specific conditions

mRNA

immobilized DNA fragments

labelled cDNA

**convert into numerical
values**

| Intensity Table | control condition | | | condition 1 | | | condition 2 | | |
|---|---|---|---|---|---|---|---|---|---|
| | hybr. 1 | hybr. 2 | hybr. 3 | hybr. 4 | hybr. 5 | hybr. 6 | hybr. 7 | hybr. 8 | hybr. 9 |
| gene 1 | 14,243 | 12,154 | | | | | | | |
| gene 2 | 5,323 | 27,152 | | | | | | | |
| gene 3 | 10,300 | 1,407 | ... | | | | | | |
| gene 4 | 1,007 | 3,101 | | | | | | | |
| gene 5 | 100,232 | 120,993 | | | | | | | |
| gene 6 | | | | | | | | | |
| gene 7 | : | | | | | | | | |
| gene 8 | | | | | | | | | |
| gene 9 | | | | | | | | | |
| : | | | | | | | | | |

## Scatterplot



Data

Data (log scale)

**MA Plot**

$M = \log_2(R/G)$

$A = 1/2 \log_2(RG)$

# Computing

- Microarray data analysis does not need much processor time (interactive instead of batch processing)

- However, it needs considerable amounts of memory (RAM)

# Computing, cont.

- Imaging or scatterplots

  comprise one hybridization at a time

  → often done on PCs
  → mostly running Windows

# Computing, cont.

- High level analyses
  (classification, clustering, projection,
  modelling, ...)
  may comprise hundreds of hybridizations

  → often done on Servers
  → mostly running Unix