# Information Systems for Genetics

Erwin M. Bakker

ISLS 23-9 2004

---

## Overview

- DIAL
- Sequence Alignment
  - BLAST, FASTA, etc.
- Genebrowsing
  - TIGR, MeV
  - Ensembl
  - http://www.genebrowser.com/

---

## DIAL

CMSB

- CMSB (Centre for Medical Systems Biology, www.cmsb.nl) is one of the Genomics Centre of Excellence
  - Genomics for identifying hidden connections between diseases: Improving diagnosis, treatment and prevention of common diseases *such as Alzheimer's, cardiovascular disease, diabetes and rheumatism.*
- DIAL (Data Integration, Analysis and Logistics) The project works on the data of the experimental projects of the CMSB.

---

## DIAL Institutes

- Leiden University Medical Center
- Leiden University
- VU University medical center
- Vrije Universiteit Amsterdam
- TNO Pharma Leiden
- Erasmus MC Rotterdam

## DIAL Example Study Groups

- RotterdamStudy (ERGO: Hofman, van Duijn, a.o.)
  - population-based cohort study of 12,000 subjects aged 55+ years. Patients have been followed for over 10 years now.

- Grip Cohort Study (Rotterdam: van Duijn, Oostra)
  - population-based cohort study of 3 generation families (2500 subjects). They are screened for the presence of multiple diseases.

- Netherlands Twin Register (Boomsma VU Amsterdam)
  - number of twins (60,670) and siblings (3,175)



---

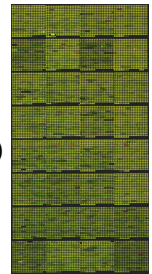## CMSB Center for Medical Systems Biology

Six Projects

- Epidemiology: cohorts & genotyping
- Systems Biology: transcriptomics/arraying proteomics metabolomics
- Technology: magnetic resonance microscopy other imaging molecular interactions
- Model Systems: animal models (mouse, zebra fish etc).
- Clinical Applications : translation (cells, vaccines, viral, pharmaceutical)
- DIAL: Data Integration, Analysis and Logistics

---

## DIAL Data Integration, Analysis and Logistics

- database/data mining group
- genomical statistics
- biomolecular bioinformatics
- computational systems biology

---

## DIAL Micro array CGH

- CGH (Comparative Genomic Hybridization) micro array experiments
- Data collection
- Data processing (Normalization)
- Data storage
- Data selection
- Data integration

## DIAL Micro Array CGH
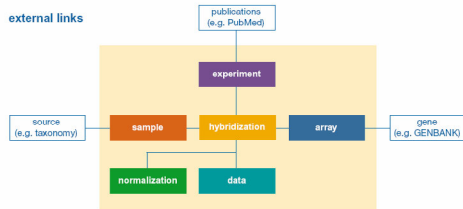MIAME minimum information about a micro array experiment



**Fig. 3** A schematic representation of six components of a microarray experiment.

## DIAL Interfaces

- GenePix, ImageGene
- Normalization Routines
- Smoothing
- BAC (Bacterial Artificial Chromosomes) updates
- CGHAnalyzer, MeV
- Rosetta, Ensembl, etc.
- Data mining

## Sequence Alignment

- sequence similarity <-> homology
- database queries; comparative genomics
- different possible alignments
- different quality of alignments

```
ACTGTGACCATATCG        ACTGTGACCATATCG
   ||X||  |             |||X|X
   ACTAT-T              ACTATT
```

## Sequence Alignment

- different possible alignments
- different quality of alignments -> score/metric
- match/mismatch: PAM, BLOSUM matrices
- insertions and deletions =>gaps: gap-penalty, gap-extension-penalty

```
ACTGTGACCATATCG        ACTGTGACCATATCG
   ||X||  |             |||X|X
   ACTAT-T              ACTATT
```

## Sequence Alignment

- what kind of alignments are valid
- scoring system for ranking the alignments
- the algorithm to find optimal or good alignments using the scoring mechanism
- statistical evaluation of the significance of an alignment score

## Sequence Alignment

- There are ~$2^{2n}/sqrt(pi.n)$ different alignments between to sequences of length n (too much)
- Dynamic Programming
  - O(nm) (n<=m) (time and space)
  - O(km) (differences are bounded by k <<n)
- Exclusion Algorithms: fast expected time
- FASTA
- BLAST
  - sequence alignment
  - Heuristics
- Hidden Markov Models (model a family of sequences)
- Chomsky Hierarchy
- Stochastic Context Free Grammars (structure prediction)

## Exclusion Algorithms

- T text, P pattern
- **Partition** P in consecutive regions of length r
- **Search** T to find length-r intervals R that could be contained in an approximate occurrence of P (surviving intervals)
- **Check** for each surviving interval R if there is an approximate occurrence of P in a larger interval around R

It is expected that already a large part of T is excluded from the **check** phase so that only (sub) linear time is needed here.

## FASTA

- A heuristic exclusion method using dynamic programming
- D.J. Lipman, W.R. Pearson Rapid and sensitive protein similarity searches. Science, 227:1435-41, 1985
- W.R. Pearson, D.J. Lipman Improved tools for biological sequence comparison. Proc. Natl. Academy Science, 85, 2444-48, 1988

# FAST A (Fast All)

- query string S; text string T
- k: is the length of the hot-spots in the dynamic programming table (k=6 for DNA, k=2 for Protein)

1. find pairs (i,j) such that the substring of length k starting at position i in S exactly matches the substring of length k at position j in T (using lookup table), and then look for diagonals with many supporting word matches (sort the matches on the difference (i-j))
2. The best diagonals from 1) are pursued: extend with ungapped alignment to find maximal scoring ungapped regions
3. join ungapped regions with gaps and find highest scoring matches using dynamic programming

- (full local dynamic programing, if k=1)

# BLAST (Basic Local Alignment Search Tool)

- Altschul et al., 1990
- O(mn) is too slow
- Probabilistic approach to searching
- True alignments will have short stretches of perfect match

# BLAST (Basic Local Alignment Search Tool)

- make a list of words W in P (proteins: length 3; nucleic acids: 11) that match T with score > threshold
- scan T for occurrences of W
- if a hit is obtained extend the match in both directions (no gaps), stopping at maximum scoring extension

=> finds ungapped alignments only

newer versions of Blast => gapped alignments

# BLOSOM62

|   | C | S | T | P | A | G | N | D | E | Q | H | R | K | M | I | L | V | F | Y | W |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| C | **9** | -1 | -1 | -3 | 0 | -3 | -3 | -3 | -4 | -3 | -3 | -3 | -3 | -1 | -1 | -1 | -1 | -2 | -2 | -2 |
| S | -1 | **4** | 1 | -1 | 1 | 0 | 1 | 0 | 0 | 0 | -1 | -1 | 0 | -1 | -2 | -2 | -2 | -2 | -2 | -3 |
| T | -1 | 1 | **4** | 1 | -1 | 1 | 0 | 1 | 0 | 0 | -1 | 0 | -1 | -2 | -2 | -2 | -2 | -2 | -2 | -3 |
| P | -3 | -1 | 1 | **7** | -1 | -2 | -1 | -1 | -1 | -1 | -2 | -2 | -1 | -2 | -3 | -3 | -2 | -4 | -3 | -4 |
| A | 0 | 1 | -1 | -1 | **4** | 0 | -1 | -2 | -1 | -1 | -2 | -1 | -1 | -1 | -1 | -1 | -2 | -2 | -2 | -3 |
| G | -3 | 0 | 1 | -2 | 0 | **6** | -2 | -1 | -2 | -2 | -2 | -2 | -2 | -3 | -4 | -4 | 0 | -3 | -3 | -2 |
| N | -3 | 1 | 0 | -2 | -2 | 0 | **6** | 1 | 0 | 0 | -1 | 0 | 0 | -2 | -3 | -3 | -3 | -3 | -2 | -4 |
| D | -3 | 0 | 1 | -1 | -2 | -1 | 1 | **6** | 2 | 0 | -1 | -2 | -1 | -3 | -3 | -4 | -3 | -3 | -3 | -4 |
| E | -4 | 0 | 0 | -1 | -1 | -2 | 0 | 2 | **5** | 2 | 0 | 0 | 1 | -2 | -3 | -3 | -2 | -3 | -2 | -3 |
| Q | -3 | 0 | 0 | -1 | -1 | -2 | 0 | 0 | 2 | **5** | 0 | 1 | 1 | 0 | -3 | -2 | -2 | -3 | -1 | -2 |
| H | -3 | -1 | 0 | -2 | -2 | -2 | 1 | 1 | 0 | 0 | **8** | 0 | -1 | -2 | -3 | -3 | -2 | -1 | 2 | -2 |
| R | -3 | -1 | -1 | -2 | -1 | -2 | 0 | -2 | 0 | 1 | 0 | **5** | 2 | -1 | -3 | -2 | -3 | -3 | -2 | -3 |
| K | -3 | 0 | 0 | -1 | -1 | -2 | 0 | -1 | 1 | 1 | -1 | 2 | **5** | -1 | -3 | -2 | -3 | -3 | -2 | -3 |
| M | -1 | -1 | -1 | -2 | -1 | -3 | -2 | -3 | -2 | 0 | -2 | -1 | -1 | **5** | 1 | 2 | -2 | 0 | -1 | -1 |
| I | -1 | -2 | -2 | -3 | -1 | -4 | -3 | -3 | -3 | -3 | -3 | -3 | -3 | 1 | **4** | 2 | 1 | 0 | -1 | -3 |
| L | -1 | -2 | -2 | -3 | -1 | -4 | -3 | -4 | -3 | -2 | -3 | -2 | -2 | 2 | 2 | **4** | 3 | 0 | -1 | -2 |
| V | -1 | -2 | -2 | -2 | 0 | -3 | -3 | -3 | -2 | -2 | -3 | -3 | -2 | 1 | 3 | 1 | **4** | -1 | -1 | -3 |
| F | -2 | -2 | -2 | -4 | -2 | -3 | -3 | -3 | -3 | -3 | -1 | -3 | -3 | 0 | 0 | 0 | -1 | **6** | 3 | 1 |
| Y | -2 | -2 | -2 | -3 | -2 | -3 | -2 | -3 | -2 | -1 | 2 | -2 | -2 | -1 | -1 | -1 | -1 | 3 | **7** | 2 |
| W | -2 | -3 | -3 | -4 | -3 | -2 | -4 | -4 | -3 | -2 | -2 | -3 | -3 | -1 | -3 | -2 | -3 | 1 | 2 | **11** |

# BLAST Queries

- NCBI  http://www.ncbi.nlm.nih.gov/
- Ensembl  http://www.ensembl.org/

## BLAST Results

- NCBI  http://www.ncbi.nlm.nih.gov/
  - Different types of BLAST
  - BLAST Results
- Ensembl  http://www.ensembl.org/
  - BLAST Result
  - Detailed View

## Genebrowsers and Experiment Visualization Tools

- TIGR TM4, MeV
  - FISH, Micro array CGH Experiments
- Ensembl
- http://www.genebrowser.com/

## CGH Experiments

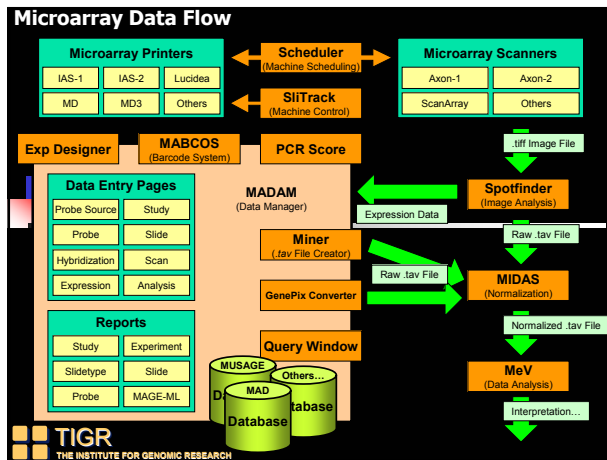*Determination of Genomic Imbalances by Genome-wide Screening Approaches*
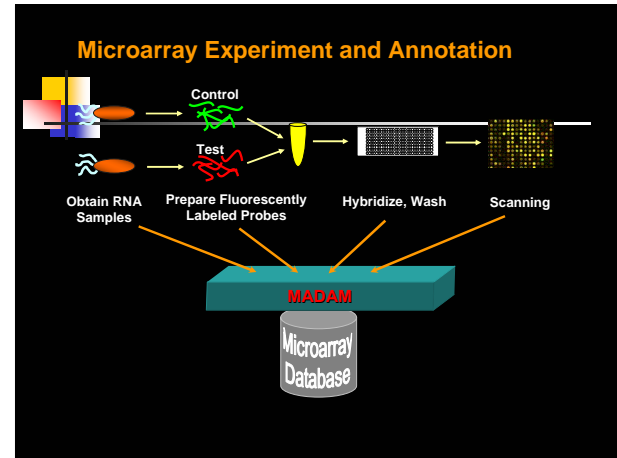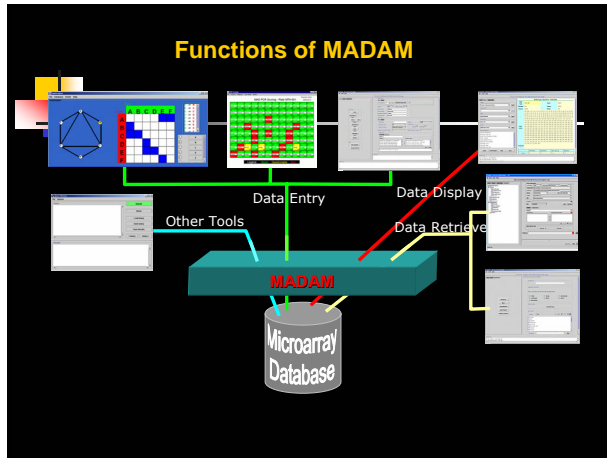
By

*Károly Szuhai*



## CGH Experiments TIGR

- Spotfinder
- Madam
- Midas
- MeV (Demo)

**Microarray Data Flow**

Printer | Scanner | .tiff Image File | TIGR Spotfinder | Image Analysis

Raw Gene Expression Data

Gene Annotation | TIGR MIDAS | Normalization / Filtering

Normalized Data with Gene Annotation

AGED | MAD | Others… | Database | TIGR MeV | Expression Analysis

Data Entry / Management

Interpretation of Analysis Results

---



**Microarray Experiment and Annotation**

Control
Test

Obtain RNA Samples | Prepare Fluorescently Labeled Probes | Hybridize, Wash | Scanning

MADAM

Microarray Database

---



**Microarray Data Flow**

Microarray Printers: IAS-1 | IAS-2 | Lucidea | MD | MD3 | Others

Scheduler (Machine Scheduling)
SliTrack (Machine Control)

Microarray Scanners: Axon-1 | Axon-2 | ScanArray | Others

Exp Designer | MABCOS (Barcode System) | PCR Score

Data Entry Pages: Probe Source | Study | Probe | Slide | Hybridization | Scan | Expression | Analysis

MADAM (Data Manager)

.tiff Image File
Spotfinder (Image Analysis)
Expression Data
Raw .tav File

Miner (.tav File Creator)
GenePix Converter
Raw .tav File
MIDAS (Normalization)

Reports: Study | Experiment | Slidetype | Slide | Probe | MAGE-ML

Query Window

MUSAGE | MAD | Others… | Database

Normalized .tav File
MeV (Data Analysis)
Interpretation…
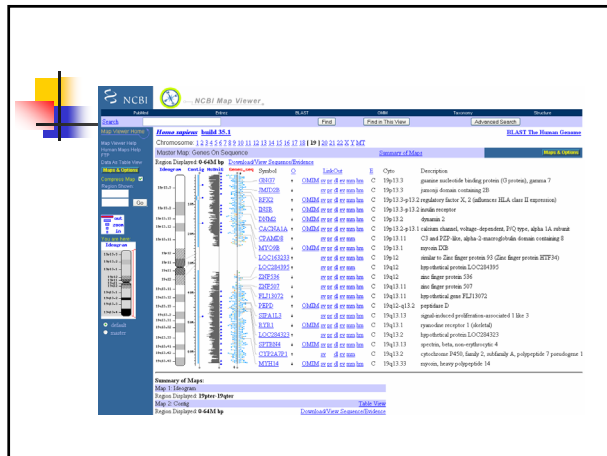
**TIGR** THE INSTITUTE FOR GENOMIC RESEARCH

---

**Requirements of MADAM**

- Communicates with a DBMS.
- Uploads microarray information to database with a convenient way.
- Allows users to view the information.
- Retrieves the information and converts it to a data format for analysis.
- Checks errors for the data entries.
- Starts other tools to manipulate the data.

Functions of MADAM

---

# Genebrowsers and Experiment Visualization Tools

- TIGR TM4, MeV
  - FISH, Micro array CGH Experiments
- Ensembl
- http://www.genebrowser.com/

---



---