


Differential gene expression


General Introduction



Swiss Institute of Bioinformatics - LF 11.2010

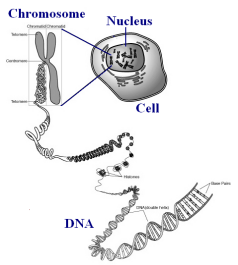

Overview (1)

- Reminder of biology
- Major steps in microarray analysis
 - Microarray preparation design, clone/probe selection
 - RNA extraction, hybridization on chip
 - Scanning, data extraction from image
 - "Low-level" Quality Control
 - Summarization of per-chip information (one number per feature)
 - "High-level" analysis
- High-throughput RNA-level technologies
 - Microarrays
 - Affymetrix Chips
 - SAGE
 - MPSS



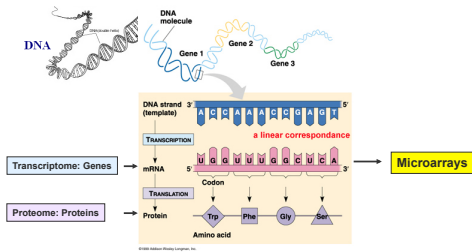
Swiss Institute of Bioinformatics - LF 11.2010

Biology Fundamentals - Genes

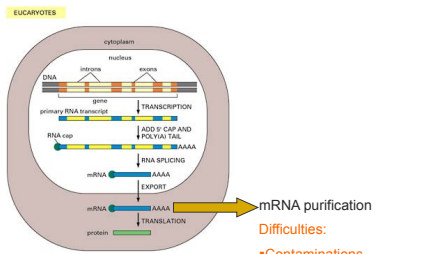
Swiss Institute of Bioinformatics - LF 11.2010

Biology Fundamentals - Expression



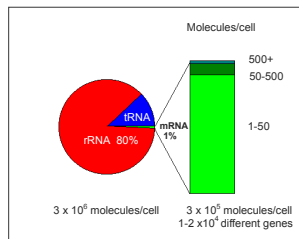
SIB Swiss Institute of Bioinformatics - LF 11 2010

Genomics Fundamentals - Complexity



SIB Swiss Institute of Bioinformatics - LF 11 2010


RNA abundance in mammalian cells



SIB Swiss Institute of Bioinformatics - LF 11 2010

Expression analysis


- Low throughput
 - Northern blot
 - Differential display
 - Quantitative PCR
- High throughput
 - DNA arrays / Chips
 - Spotted arrays (Stanford arrays)
 - Affymetrix (photolithography inspired)
 - Oligo-arrays (Agilent, NimbleGen)
 - Serial Analysis of Gene Expression (SAGE)
 - RNASeq

 Swiss Institute of Bioinformatics - LF 11 2010

What are DNA Microarrays ?


Microarray analysis is a technology that allows scientists to simultaneously detect thousands of genes in a small sample and to analyze the expression of those genes.

Microarrays are simply ordered sets of DNA molecules of known sequence. Usually rectangular shaped, they can consist of a few hundred to hundreds of thousands of sets. Each individual sequence goes on the array at precisely defined location.

 Swiss Institute of Bioinformatics - LF 11 2010

Potential application domains

- Identification of complex genetic diseases
- Drug discovery and toxicology studies
- Mutation/polymorphism detection (SNP' s)
- Pathogen analysis
- Differing expression of genes over time, between tissues, and disease states
- Preventive medicine
- Specific genotype (population) targeted drugs
- More targeted drug treatments – AIDS
- Genetic testing and privacy

 Swiss Institute of Bioinformatics - LF 11 2010

The challenge

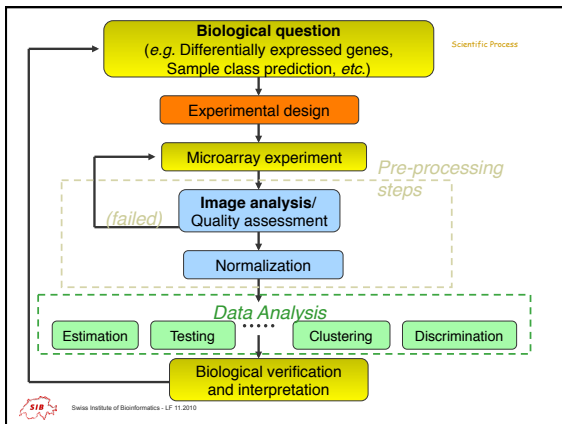
The big revolution here is in the "micro" term. New slides will contain a survey of the human genome on a 2 cm² chip! The use of this large-scale method tends to create phenomenal amounts of data, that have then to be analyzed, processed and stored.

This is a job for... **Bioinformatics!**

General overview

- Making the chip
 - Experiment design, clone/probe selection, collection maintenance, PCR, spotting, printing, synthesis
- Sample hybridization
 - Sample purification, labelling, hybridization, washing
- Scanning and image treatment
 - Fluorescence correction, find spots, background
- Analysing the data
 - Filtering, normalisation
 - Clustering (hierarchical, centroid,...)
- Representation, storage
 - Graphics, databases, web public resources

} wet lab



Question addressed by microarrays

- What are the differences (in gene expression) between two cell lines ?
- What is the difference between knock-out and wild-type mice?
- What is the difference between a tumor and a healthy tissue ?
- Are there different tumor types ?

- Key concept: *Compare* gene expression in two (or more) cell/ tissue types ?
 - Gene expression assessed by measuring the number of RNA transcripts.
 - No absolute measurement.

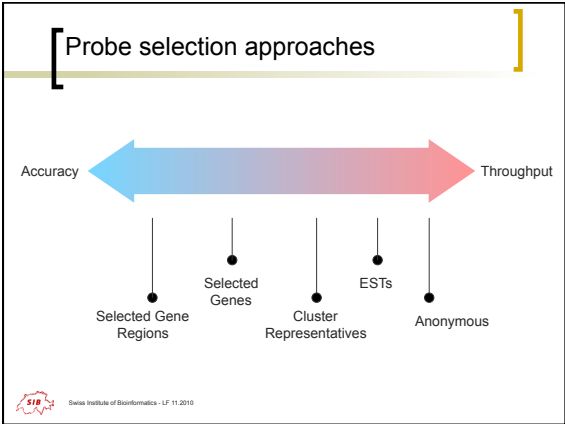
THE EXPERIMENT : making the chip

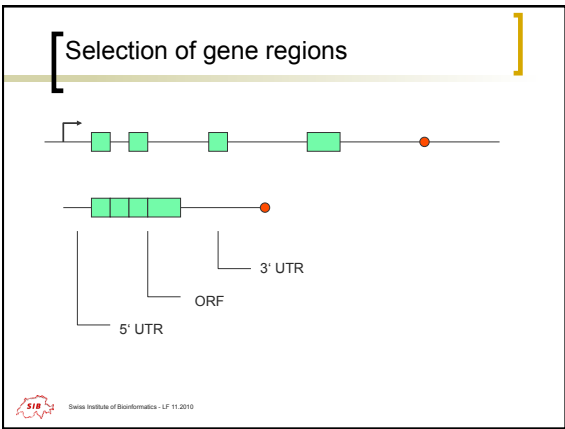
1- Designing the chip : choosing genes of interest for the experiment and/or select the samples

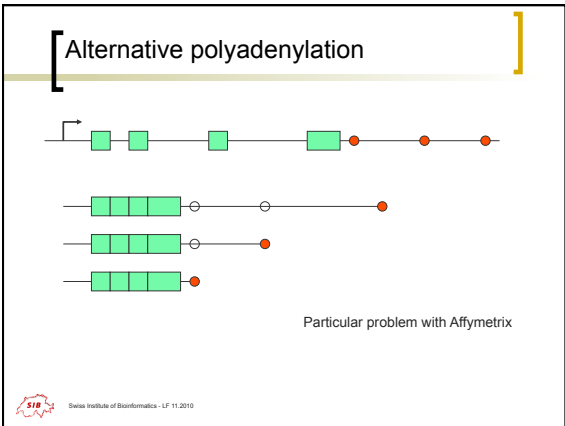
- Selection of sequences that represent the investigated genes.
- Finding sequences, usually in the EST database.
- Problems : sequencing errors, alternative splicing, chimeric sequences, contamination...

Clone/probe selection

- General
 - Not too short (sensitivity, selectivity)
 - Not too long (viscosity, surface properties)
 - Not too heterogeneous (robustness)
 - Degree of importance depends on method
- Single strand methods (Oligos, ss-cDNA)
 - Orientation must be known
 - ss-cDNA methods are not perfect
 - ds-cDNA methods don't care







Alternative splicing

Swiss Institute of Bioinformatics - LF 11 2010

Alternative promoter usage

Swiss Institute of Bioinformatics - LF 11 2010

Selection of gene regions - summary

- Coding region (ORF)
 - Annotation relatively safe
 - No problems with alternative polyA sites
 - No repetitive elements or other funny sequences
 - danger of close isoforms
 - danger of alternative splicing
 - might be missing in short RT products
- 3' untranslated region
 - Annotation less safe
 - danger of alternative polyA sites
 - danger of repetitive elements
 - less likely to cross-hybridize with isoforms
 - little danger of alternative splicing
- 5' untranslated region
 - close linkage to promoter
 - frequently not available

Swiss Institute of Bioinformatics - LF 11 2010

A checklist

- Pick a gene
- Try to get a complete cDNA sequence
- Verify sequence architecture (e.g. cross-species comparison)
- Mask repetitive elements (and vector!)
- If possible, discard 3'-UTR beyond first polyA signal
- Look for alternative splice events
- Use remaining region of interest for similarity searches
- Mask regions that could cross-hybridize

- Use the remaining region for probe amplification or EST selection
- When working with ESTs, use sequence-verified clones



Swiss Institute of Bioinformatics - LF 11 2010

THE EXPERIMENT : making the chip

2- Spotting the sequences on the substrate

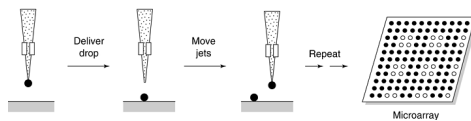
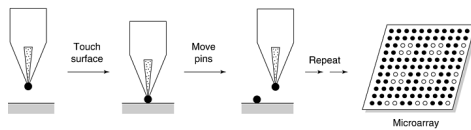
- Substrate : usually glass, but also nylon membranes, plastic, ceramic...
- Sequences : cDNA (500-5000 nucleotides), oligonucleotides (20-80-mer oligos), genomic DNA (~50'000 bases)
- Printing methods : microspotting, ink-jetting or in-situ printing, photolithography



Swiss Institute of Bioinformatics - LF 11 2010

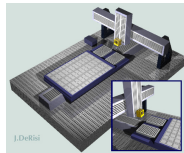
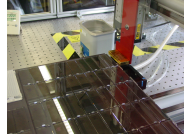
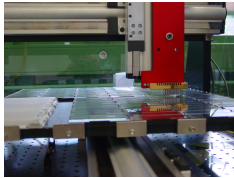
Microarrays: the making of

Microspotting and ink-jetting



Swiss Institute of Bioinformatics - LF 11 2010

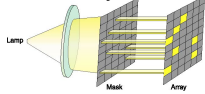
Array Production: Spotting



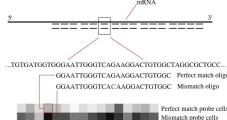
SIB Swiss Institute of Bioinformatics - LF 11 2010

Array Production: "photolithography"

Affymetrix

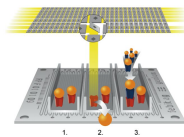


- Each probe 25 bp long
- 22-40 probes per gene
- Perfect Match (PM) as well as Mismatch (MM) probes



SIB Swiss Institute of Bioinformatics - LF 11 2010

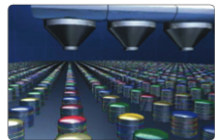
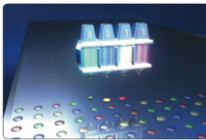
Febit/NimbleGen



- Probe length: 24mer -70mer
- Gene/Array: Up to 38,000
- Probes/Gene: 10-25
- Only perfect match probes

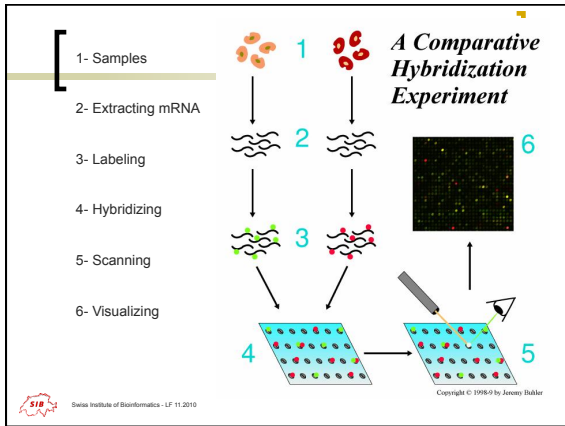
Array Production: "Inkjet"

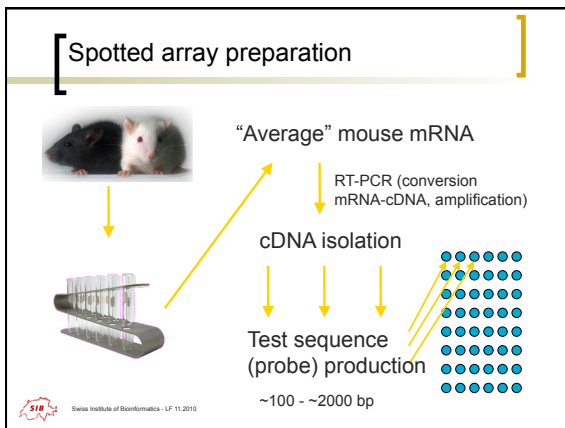
Agilent (HP SurePrint technology)

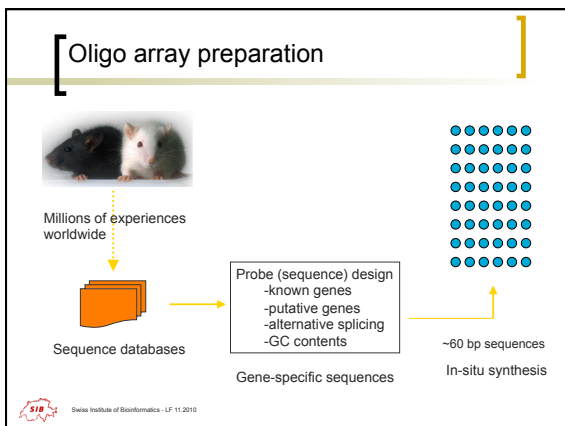


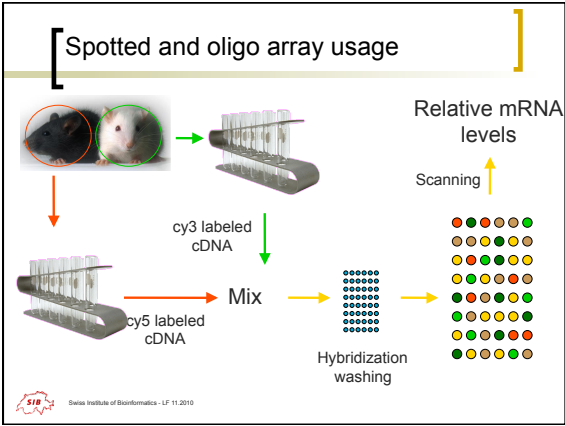
- cDNA printing
- 60bp oligo in-situ synthesis

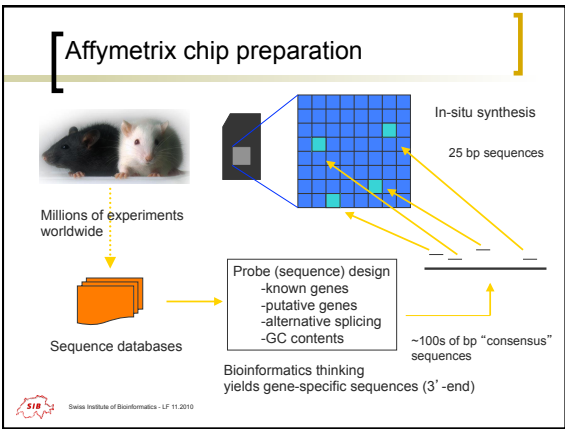
SIB Swiss Institute of Bioinformatics - LF 11 2010

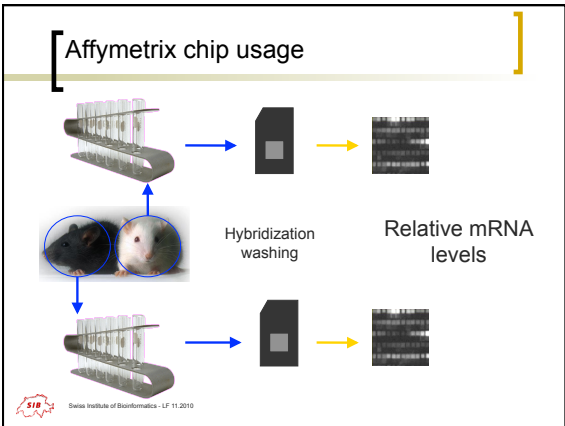


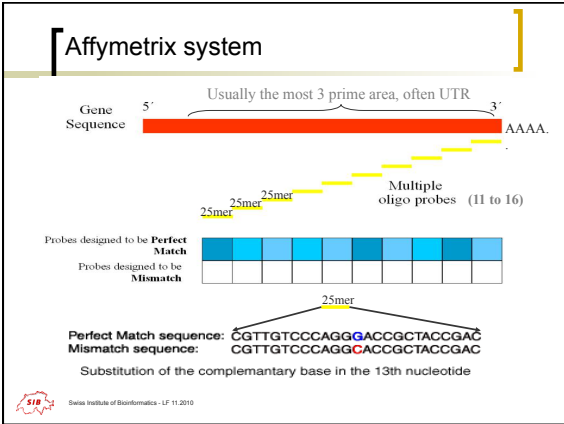


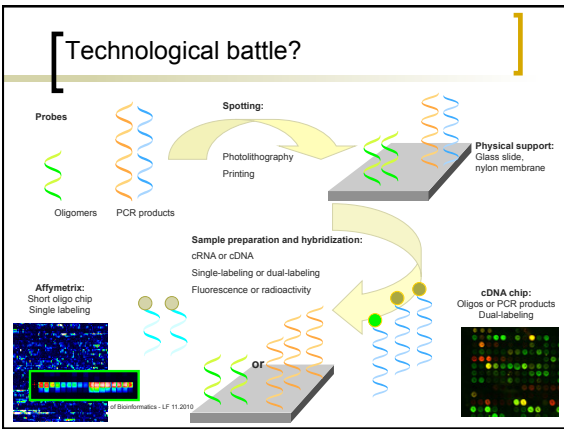












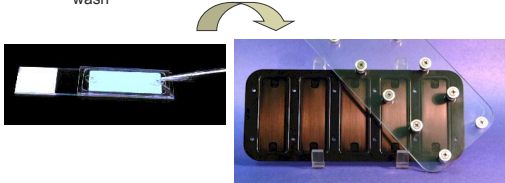
Comparison of techniques

In-situ Synthesis / Oligos	PCR Products / cDNA Probes
<p>Advantages</p> <ul style="list-style-type: none"> No need to isolate and purify cDNAs because oligonucleotides can be synthesized. Short oligonucleotides are less likely to have cross-reactivity with other sequences in the target DNA. Density of chips is higher than with cDNAs. 	<p>Advantages</p> <ul style="list-style-type: none"> Flexibility to study cDNAs from any source. cDNAs do not require any a priori information about the corresponding genes. Longer sequences increase hybridization specificity, which reduces false positives.
<p>Limitations</p> <ul style="list-style-type: none"> The sequence has to be known. Synthesis can be expensive and time-consuming. The short sequences are not as specific for target DNA, so appropriate controls must be added. 	<p>Limitations</p> <ul style="list-style-type: none"> Isolation of individual cDNAs to immobilize on each spot can be cumbersome. Density is lower than synthesizing oligonucleotides on the surface of the chip. cDNAs are longer sequences and are more likely to randomly contain sequences found in target DNA, which results in cross-reactivity.

SIB Swiss Institute of Bioinformatics - LF 11.2010

Probe preparation & hybridization

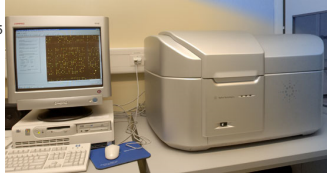
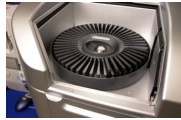
- Extract mRNA or total RNA
- RT, add 5' anchor
- PCR with labelled nucleotide (Cy3, Cy5, DIG or radiolabelling)
- Overlay probe on the chip, put in the hybridization chamber, wash



 Swiss Institute of Bioinformatics - LF 11 2010

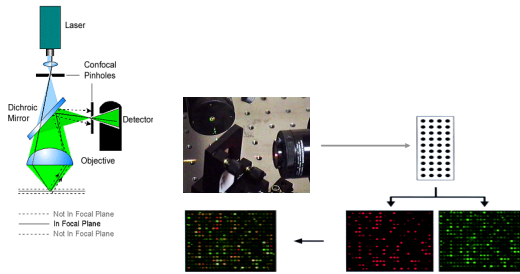
Scanner basics

- Based on fluorescence
 - 1 or 2 lasers: cy3 cy5 (seldom more)
- Most scanners are confocal
 - Target a very limited volume of space (signal only from focal plane)
 - Need to "scan" the surface
- 16-bits ADC converters
 - Range of values: 0-65535
 - Log2 range: 0 - 16
- Scan various supports
 - Glass Slide (e.g. Agilent, PerkinElmer)
 - Affymetrix



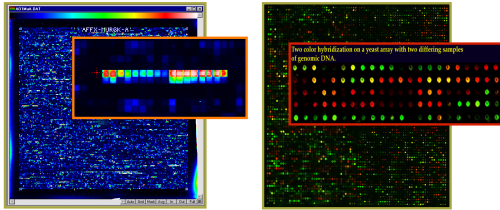
 Swiss Institute of Bioinformatics - LF 11 2010

Confocal scanner



 Swiss Institute of Bioinformatics - LF 11 2010

Scanner output: image(s)



Affymetrix chip
1 channel, false colors

Stanford array
dual-channel, color addition



Statistik-Institut der Universität Wien - 11.11.2010

Image analysis (scanner variability)

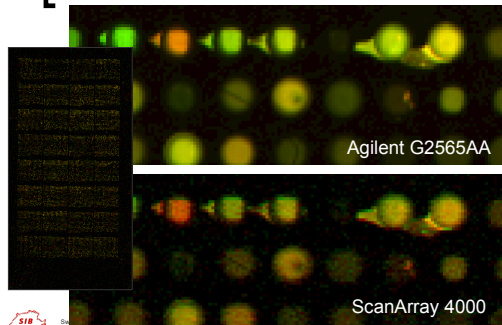
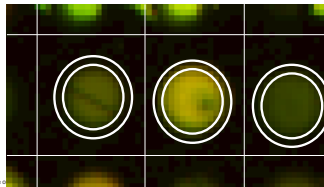


Image processing

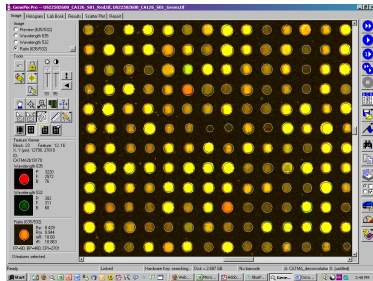
- Align channels
- Identify spot pixels
- Identify background pixels
- Compute representative value, e.g.
 - Mean foreground value
 - Median background value



Statistik-Institut der Universität Wien

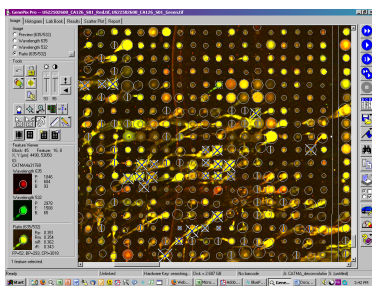
2-color Arrays Image Processing

GenePix



Swiss Institute of Bioinformatics - LF 11.2010

2-color Arrays Image Processing



A difficult case... 😊

Swiss Institute of Bioinformatics - LF 11.2010

Other high-throughput techniques: sequencing

- EST counting
- SAGE (Serial Analysis of Gene Expression)
- MPSS (Massively Parallel Signature Sequencing)
- RNASeq

Swiss Institute of Bioinformatics - LF 11.2010

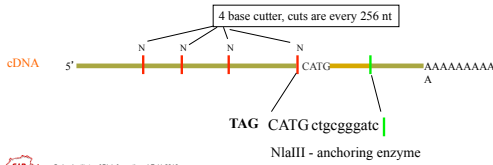
Comparison of techniques

Method	Microarrays	EST counts	SAGE/MPSS
Genes	Selected genes	Highly expressed	Almost all genes
Sampling	Analog	Digital	Digital
Statistics	Robust	Robust	Robust
Molecules sampled	High	Low-medium	Medium-high
Duplicates	Required	Desirable	Not required
Costs	Medium	High	High
Sharing	Variable	Easy	Easy

 Swiss Institute of Bioinformatics - LF 11 2010

SAGE (Principle)

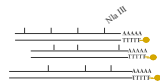
- A short nucleotide sequence "TAG" (9 to 10 bp) from a defined position within the transcript contains sufficient information to uniquely identify a transcript.
 - e.g. a sequence of 10 bp can distinguish 1,048,576 transcripts (4^{10}) given random nucleotide distribution at the tag site.
 - current estimates suggest the human genome only encodes ~ 100,000 transcripts (from 35,000 genes)



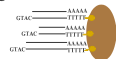
 Swiss Institute of Bioinformatics - LF 11 2010

SAGE Protocol (I)

- ◆ cDNA synthesis using biotinylated oligo-dT nucleotide



- ◆ Digestion of biotinylated cDNA with Nla III (anchoring enzyme) and binding to magnetic beads



 Swiss Institute of Bioinformatics - LF 11 2010

SAGE Protocol (II)

- Divide in half and ligate to linkers A and B

- Digest with BsmFI enzyme (tagging enzyme, TE), keep supernatant
- Blunt end (fill in)

SIB Swiss Institute of Bioinformatics - LF 11.2010

SAGE Protocol (III)

- Ligated and amplified with primers A & B

- PCR amplification using primers A & B

SIB Swiss Institute of Bioinformatics - LF 11.2010

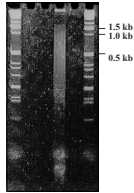
SAGE Protocol (IV)

- Purification of 102 bp "ditag" band and Nla III digestion

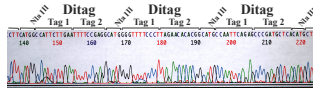
SIB Swiss Institute of Bioinformatics - LF 11.2010

SAGE Protocol (V)

- Purification of 26 bp "ditag" band and concatemerization



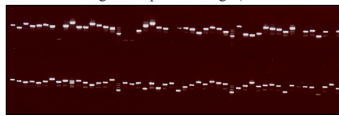
- Purification of 0.7-1.5 kb long concatemers subcloning and sequencing



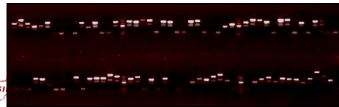
 Swiss Institute of Bioinformatics - LF 11.2010

SAGE Protocol (VI)

- Library verification
 - Insert size
 - Percentage inserts
 - Percentage of ditags derived from linkers (between 2 and 15%)
 - Percentage of duplicate ditags (should correlate to TAG abundance)



- 89/96 inserts (93%)
- Ave size ~ 900bp



- 67/96 inserts (70%)
- Ave size ~ 600bp



DiTAG and TAG extraction from sequence data

SAGE300: (version 3.03, 1998)

Kinzler: kinzlk@welchlink.welch.jhu.edu

SAGE
Serial Analysis of Gene Expression

```
12TC7T7
NCNTTCCCAGGGGCTTGGCCTCGATATGCATGCCCATCGTCCTAGTAACAAGATCATGAACAGTA
TGTGGGATGCTTGCATGTCAGCTCCCATATTTCCGAGGCATGTCCTATTAAGCAGAGGACCA
ACATGCTGCAACCTATGAAGCCTTATCATGAAGCAGTTACAAAACAGCCAGCCATGCACCTAATT
GGAGGCTGTGATCATGTACATAAATTACTGGGGTTTCGACATGTGAGAGACATCTANACTTTTA
CCATGCTCGAGCGGCCCGCCAGTGTGATGGATATCTGCAGAAATTCGGCACACTGGGGGGG
```

 Swiss Institute of Bioinformatics - LF 11.2010

DiTAG and TAG Extraction From Sequence Data

Input File: X:\sage\XXX12TC7T7.SEQ

```
1 ) CCCATCCTCCTAGTAACAAGAT - 347574 - 229106
2 ) AAACAGTATGTGGGGATGCTTG - 4815 - 271574
3 ) TCAGCTCCCATATTTCCGAGG - 862037 - 383489
4 ) TCCCTATTAGCAGAGACCAC - 875761 - 1027448
5 ) CTGCCAACCTATGAGCCTTAT - 495709 - 199294
6 ) AAGCAGTTACAAAACAGCCAGC - 37618 - 649712
7 ) CACCTAATTGGAGGCTGTGAT - 285759 - 214182
8 ) TACATAAATTAAGGGTTTCGA - 805949 - 884822
9 ) TGAGAGACATCTAARCTTTTAC - 926228 - 0
```

Total Dimers: 9
Short Dimers: 0
Long Dimers: 0
Duplicate Dimers: 0
Good Tags: 17

 Swiss Institute of Bioinformatics - LF 11 2010

Comparisons of SAGE Projects

Lib # 1	Lib # 2	Total	Tag Sequence
962	350	1312	GTGGCTCACA
614	687	1301	GCTGCCCTCC
288	140	428	AGCAGTCCCC
262	135	397	GCTTCGTCCA
221	136	357	TCAGGCTGCC
233	119	352	ATACTGACAT
202	102	304	AAAAAAAAAA
239	65	304	GCCTCCAAGG
212	90	302	GTGACCTGGC
233	64	297	GCGGGTCCGC
160	97	257	GCAACTCTTG
180	70	250	CATCGCCAGT
163	71	234	GAGCGTTTIG

 Swiss Institute of Bioinformatics - LF 11 2010

Massively Parallel Signature Sequencing

- Alternative to SAGE: generate 13-nt tags from a large ($>10^5$) sample of cDNAs
- Longer tags means higher specificity
- Solid-state technology gives high throughput
- Cost comparable to SAGE, but much larger number of tags obtained



 Swiss Institute of Bioinformatics - LF 11 2010

Overall strategy of MPSS

- “Megacloning”
 - Generate cDNA population where each molecule is attached to a different tag
 - Amplify this population
 - Attach to beads carrying anti-tags
 - Purify beads that have captured cDNAs
- Massively parallel sequencing
 - Use cycles of Rx enzyme cleavage, ligation and hybridization to read blocks of 4 nucleotides

 Swiss Institute of Bioinformatics - LF 11 2010

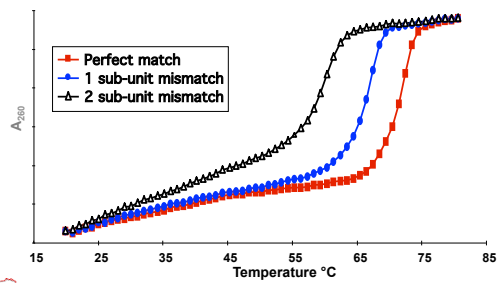
“Words” for tag construction

- TTAC, AATC, TACT, ATCA, ACAT, TCTA, CTTT, and CAAA
 - No restriction sites
 - Isothermal denaturation
 - Large ΔT_m for match/mismatch pairs
 - Large total repertoire (8^8 or 16,777,216)
 - Example:

5' -TACT . TTAC . ACAT . ATCA . CTTT . CTTT . CAAA . AATC- 3'
3' -ATGA . AATG . TGTA . TAGT . GAAA . GAAA . GTTT . TTAG- 5'

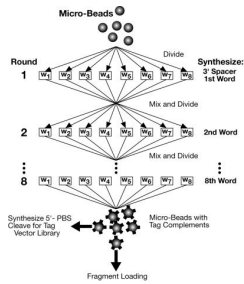
 Swiss Institute of Bioinformatics - LF 11 2010

Tag Hybridization Specificity



 Swiss Institute of Bioinformatics - LF 11 2010

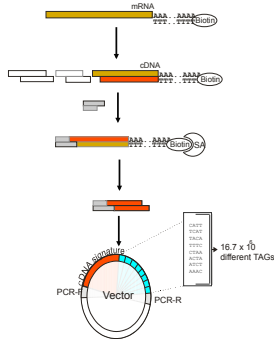
Generation of tags from "words"



SIB Swiss Institute of Bioinformatics - LF 11 2010

Signature Capture & Tagging

cDNA synthesis & DpnII digestion



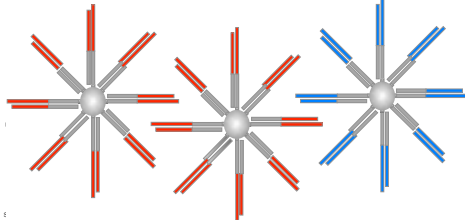
SIB Swiss Institute of Bioinformatics - LF 11 2010

Loading DNA on Microbeads

Each bead collects $\sim 10^5$ copies of a given fragment.

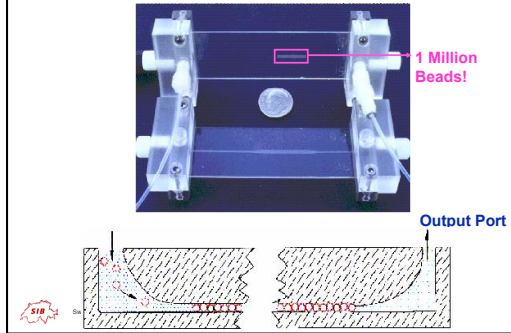
After loading one strand is covalently ligated to the bead

Clones of X_1 Clones of X_2 Clones of Y



SIB Swiss Institute of Bioinformatics - LF 11 2010

Beads Immobilized in a Flow Cell



Flowchart for the sequencing reactions

Table 1. Sequences of encoded adaptors*

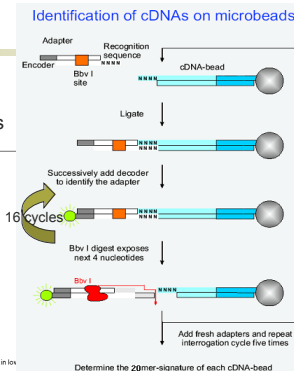
Common strand:
5'-GACTGTSAGACTCGT

Encoded adaptors for detecting base 1:
5'-NNNACGAGCTTCCAGTTCatflapagg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapacg
5'-NNNACGAGCTTCCAGTTCatflapaccg

Encoded adaptors for detecting base 2:
5'-NNNTACGAGCTTCCAGTTCatflapagg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg

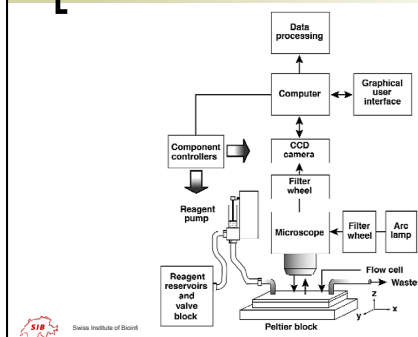
Encoded adaptors for detecting base 3:
5'-NNNNACGAGCTTCCAGTTCatflapagg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg

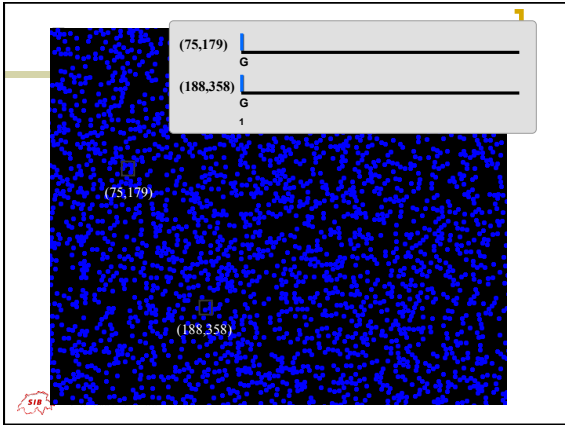
Encoded adaptors for detecting base 4:
5'-ANNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg
5'-NNNACGAGCTTCCAGTTCatflapaccg

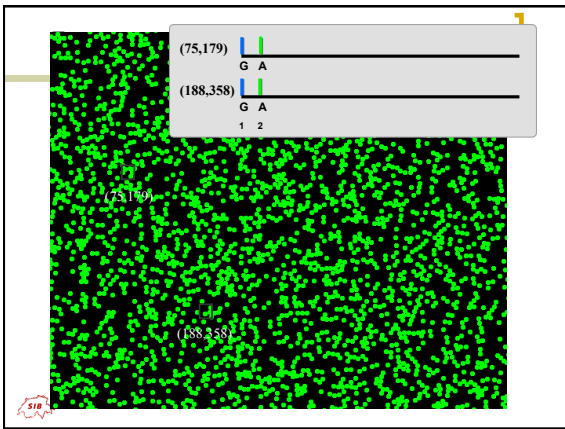


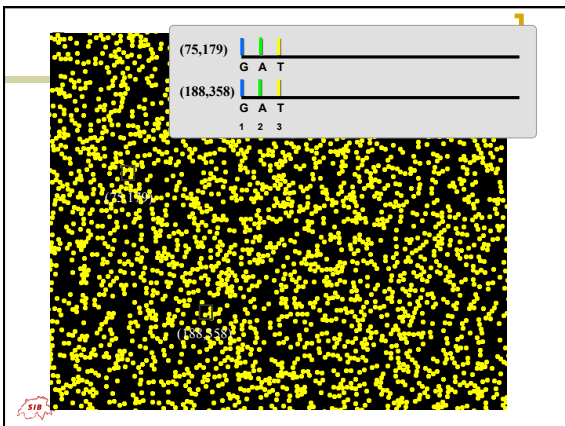
*Four-base overhangs in bold and decoder binding sites in low

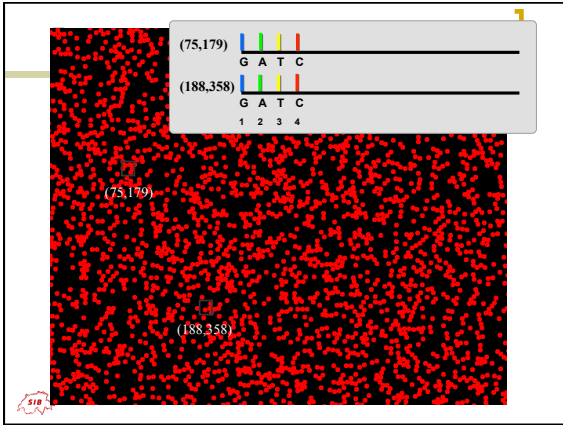
Setup for acquiring signature sequences

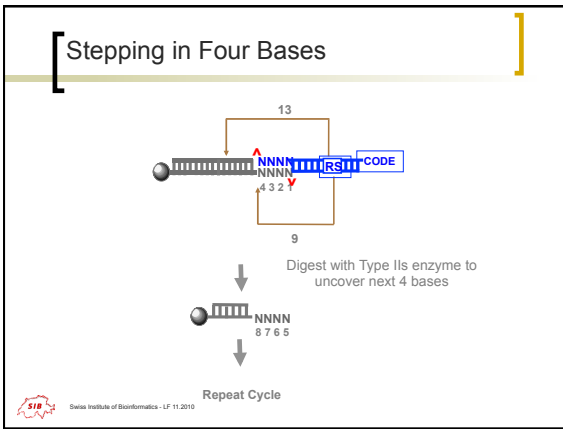


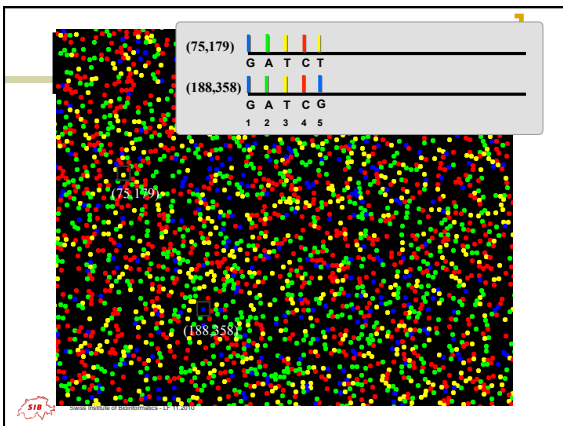


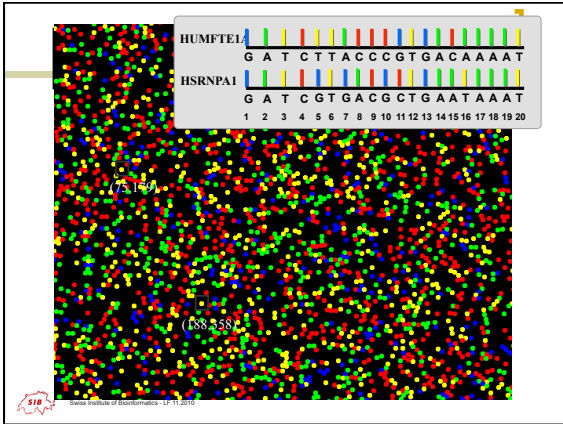


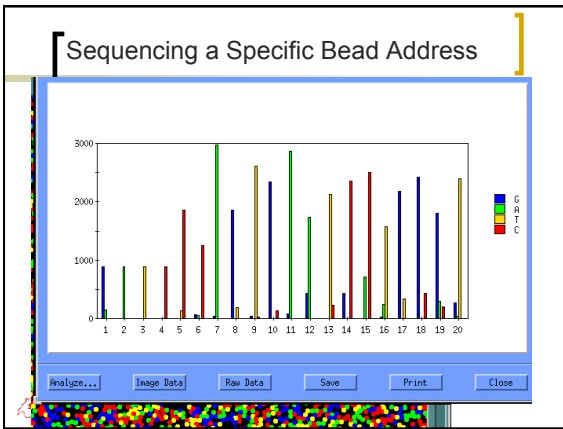












From signatures to genes

- Database searching
- RT-PCR
- cDNA library screening.

- Functional hypotheses and analyses

Problems with SAGE/MPSS data

- Sequencing errors in the libraries
- Sequencing errors in the ESTs used to derive the signatures
- Incomplete digestion by Rx enzymes
- Ambiguities in signature to gene maps

Unreliable sequences

- 1% error rate in sequence means that there is 10% chance a signature is wrong in either library or EST
- Correction in libraries by elimination of low-frequency signatures (singlets) and by merging of neighbours of abundant tags
- Correction in ESTs and detection of SNPs by aligning to genome data

Getting signatures from the transcriptome (NCBI)

- Separate out individual species (e.g., human) sequences from GenBank submission records.
- Assign a SAGE tag to each sequence, by:
 - assigning sequence orientation through a combination of identification poly-adenylation signal (ATTAAA or AATAAA), poly-adenylation tail, and sequence label, and
 - extracting a 10 base tag 3'-adjacent to the 3'-most NlaIII site (CATG).
- Use information from NCBI's UniGene project, assigning an UniGene identifier to each species sequence with a SAGE tag.

The problem of poly(A)

- There is one chance in 256 that the first four nt upstream of the poly(A) are CATG
- There is one chance in 25 that CATG is found within 10 nt of the poly(A)
- Therefore, tags containing multiple A's at their 3' end may not be mapped correctly
- In fact, tags consisting of only A's are found very commonly in SAGE/MPSS libraries

Multiple tags per gene and genes per signature

- Many (probably most) genes have more than one polyadenylation site, and may be associated with multiple UniGene clusters
- About 1% of the tags originate from partially cleaved cDNA (i.e. from the 2nd restriction site)
- The same tag can appear in more than one mRNA, and have a different probability of being generated from each

Getting signatures from the genome

- Extract poly(A) proximal 3' tags from EST trace files and map to the genome
- Map exons on the genome from the transcriptome
- Follow the exons from the 3' tag to find NlaIII site and SAGE tag
- This identifies 120' 000 reliable tags on the human genome

Efficiency of tag annotation

	HB4a	HCT-116	Combined
Total tags	17354	24065	27965
Contaminants ^a	160	264	276
Match virtual transcripts ^b	12109 (70%)	14699 (62%)	17992 (65%)
Match NCBI models ^c	9476 (55%)	10883 (46%)	12326 (44%)
Match Ensembl transcripts	8561 (50%)	9842 (41%)	11105 (40%)

- a) Contaminants include mitochondrial and ribosomal RNAs and repetitive elements.
 b) Percentages are calculated relative to total tags minus contaminants.
 c) Combination of experimental and predicted transcripts (NM and XM identifiers) in RefSeq (November 8, 2002).

 Swiss Institute of Bioinformatics - LF 11 2010

Distribution of tag abundance

Abundance (tpm) ^a	# of tags	# identified ^b	# of genes
>10000	7	7	3
>5000	25	24	14
>1000	154	149	120
>500	298	280	229
>100	1719	1600	1397
>50	3261	3060	2631
>10	10519	9608	8018
>5	15145	13517	10876
>1	27965	25779	17992

Most of the tags that cannot be identified are derived from lowly expressed genes

- a) Mean between the HB4a and HCT-116 libraries.
 b) Identification includes a match to a gene, to a known contaminant, or to a potential sequencing error or polymorphism; it does not include a match to the genome.

Jongeneel CV et al. Comprehensive sampling of gene expression in human cell lines with massively parallel signature sequencing. Proc Natl Acad Sci U S A. 2003 Apr 15;100(8):4702-5.

 Swiss Institute of Bioinformatics - LF 11 2010

Conclusion

- What is your biological question?
 - Choose the right technique
 - Compare the prices
 - Do the right experiment
 - Analysis the results, does it answer your question?
 - Always ask for help - before - doing a mistake
- Questions ?

 Swiss Institute of Bioinformatics - LF 11 2010
