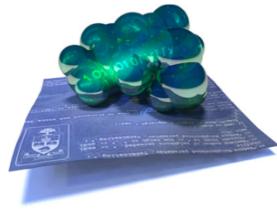


A
BIOINFORMATICS
COURSE

EXPRESSION ANALYSIS



BORIS STEIPE

DEPARTMENT OF BIOCHEMISTRY – DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO

TRANSCRIPTOME INFORMATICS

THE EXPERIMENT: Microarrays, RNA sequencing.

THE DATA: Genes and expression levels stored in databases.
Experimental conditions are important (MAGE, MIAMI).

DATABASES:

N C B I G E O : Repository of experimental information;
select by reference, organism, experiment ...

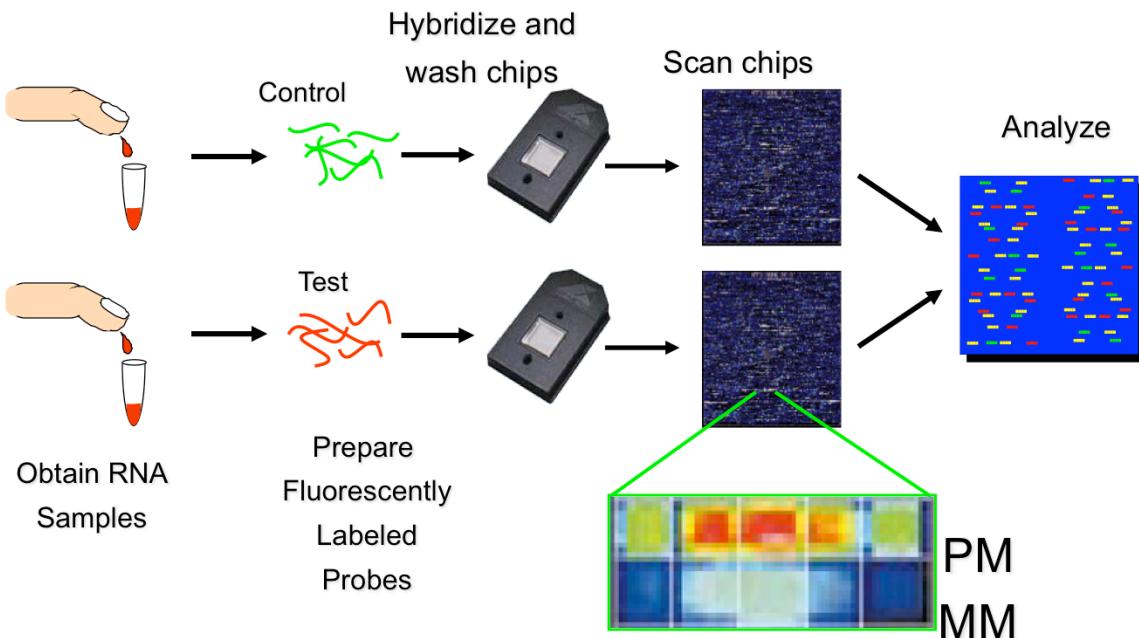
EXPRESSION ATLAS: EMBL–EBI

THE QUESTION: Expression and Differential Expression

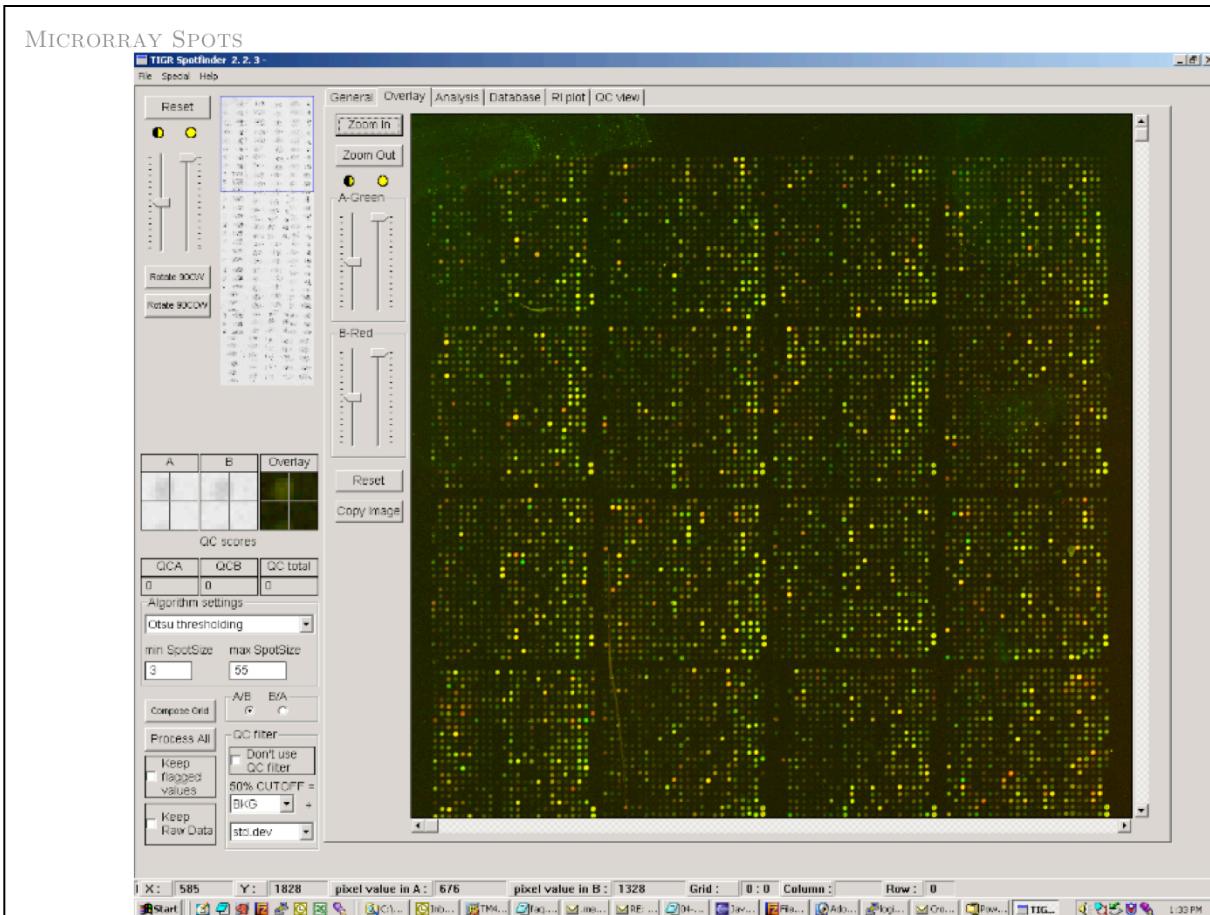
THE WORKFLOW: Collect, normalize, calculate log-ratios,
evaluate statistical significance, annotate, interpret.

SPOTTED ARRAYS AND GENE CHIPS

Gene chip experiment

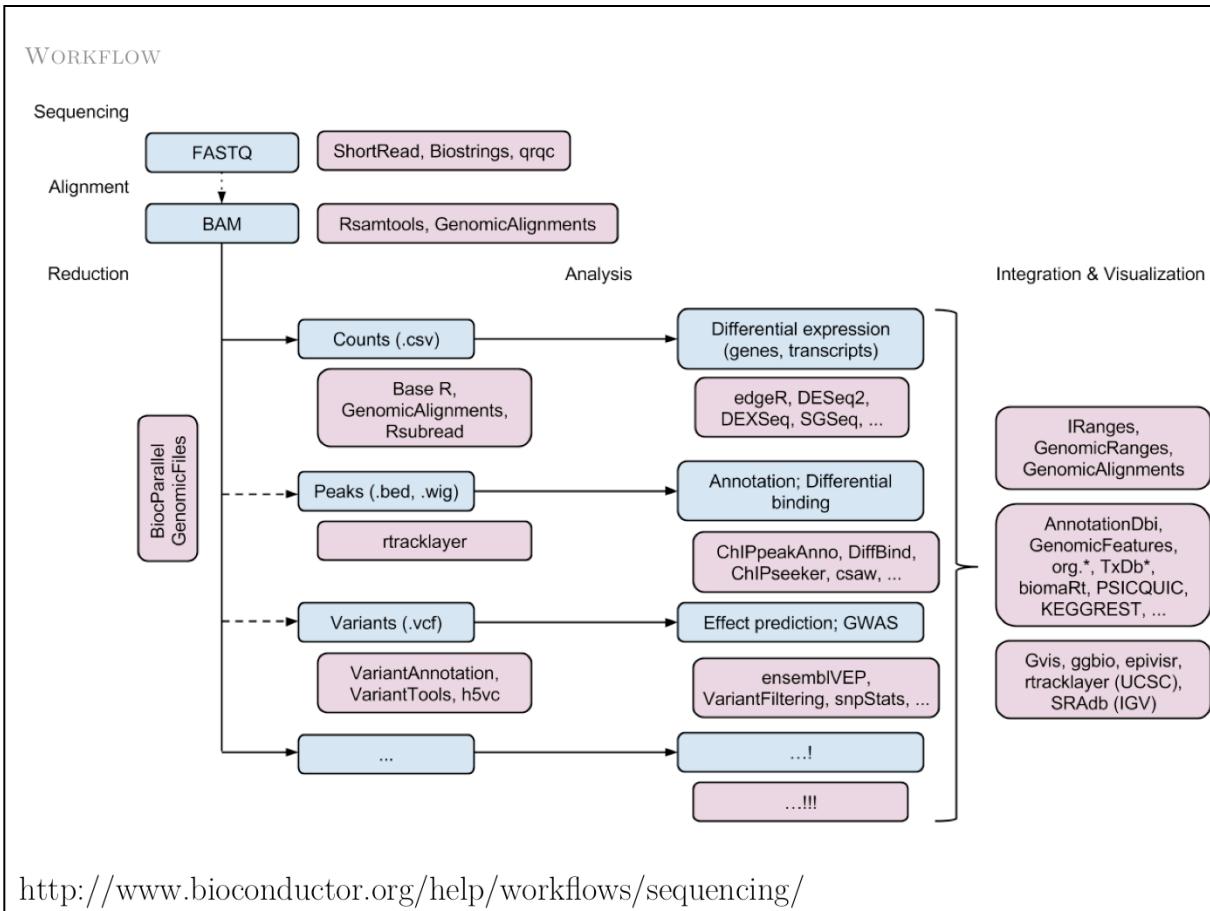


Chips, or microarrays, are solid supports that have crosslinked probes of oligonucleotides in defined locations. The oligonucleotides hybridize (more or less) specifically with mRNA molecules from the fluorescent-labelled samples and thus have fluorescent spots. Since the location of the specific probes is known, the identity of the hybridized mRNA molecule can be inferred from the sequence of the probe. In this way spots are associated with genes. The intensity and color of fluorescence depends on the absolute and relative amounts of mRNA in the sample. In our example, a spot that contains more mRNA in the test sample (gene has been upregulated) will fluoresce red.



View of a “spotted array” image. This is the raw data from which microarray expression data are derived.

Each spot provides an averaged view of the amount of mRNA that is present in the sample.



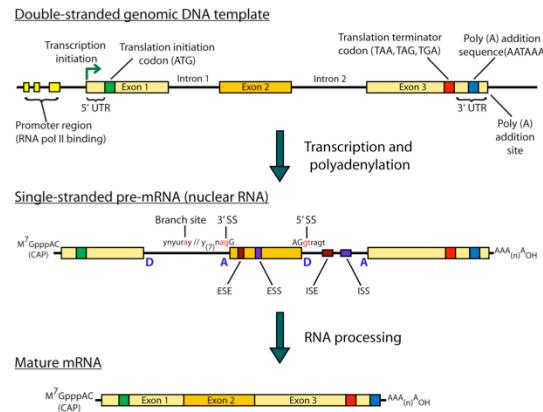
Microarray experiments have been almost completely replaced by RNAseq experiments. But these are **very** different types of data and they raise unique challenges.

After sequencing samples and controls of suitably fragmented mRNA, reads are aligned to the reference genome. But to convert reads to counts of mRNA molecules, the probabilistic nature of the experiment has to be taken into account because reads are very much shorter than mRNA molecules. If we get two reads from different regions of an mRNA, does that mean there were two copies of the mRNA, or just one copy that was seen twice? And since the sampling is stochastic, most reads may come from a small number of highly expressed mRNAs – after all, the range of concentrations of individual mRNA molecules in the cell spans 5 to 7 orders of magnitude! Moreover the length of mRNA molecules can vary by 2 orders of magnitude and that significantly influences the probability of observing reads.

Bottom line: technology is needed to convert reads to counts.

ALIGNMENT TO THE REFERENCE GENOME MUST BE SPlice-AWARE

- RNA-seq reads are derived from mRNA after splicing
- They may span large introns
- To align them back to the reference genome DNA sequence requires splice-aware aligners
 - TopHat, STAR, MapSplice, etc.



The first step of converting reads to counts is to align them to the genome and identify the gene they came from.

SAM/BAM FILES

SAM: Sequence Alignment Map

These files map reads to the genome. They are typically several GB in size. BAM files are compressed relative to SAM files but need special software to read and process.

```

mriffiti@linus27 ~ $ samtools view -H /gscmnt/gc13001/info/model_data/281932684/build136494552/alignments/13688019.bam | grep -P "SN:i;22;HD/[RG/PC"
@HD VN:1.4 SD:coordinate
@SN:22 LN:1358456
@RG ID:linus27_1299 FN:Pluimenna_PU1D844ACXX_1LH_BA-K42198-0817007-[C-3]_BA-1191 PI:365 DS:paired end DT:2012-01-31T00:00:05Z SM:H_KA-K42198-0817007 ON:WUGSC
@PG ID:2888721359 VN:2.0.8 CL:library --library-type=secondstrand --barcode-version=2.1.0
@PG ID:MarkDuplicates PP:2888721359 WI:1.85 exported C:\net\sf.picard\MarkDuplicates INPUT=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blad1e18-2-5,gsc.wustl.edu-jwukar-15434-13688019/search-1lgv-H_KA-K42198-0817007-[C-3]-BA-1191-2888363803.post_duB METRICS_FILE=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blad1e18-2-5,gsc.wustl.edu-jwukar-15434-13688019/search-1lgv-H_KA-K42198-0817007-[C-3]-BA-1191-2888363803.DUMP REMOVE_DUPLICATES=false ALLOW_SORTED=true FILENAMES FOR_DNA_IS_IN_FASTA=false VALIDATION_STRICTNESS=SILENT MAX_RECORDS_IN_RAM=5000000 PROGRAM_RECORD_INDEX=1 PROGRAM_GROUP_NAME=DUPликаты Duplicates MAX_SEQUENCE_FOR_DISK_READ_ENDS_MAP=500000 SORTING_ALGORITHM=SIZE_RATIO=0.25 READ_NAME_REGEX=[a-zA-Z0-9]+:[0-9]+:[0-9]+.* OPTICAL_DUPLICATE_PIXEL_DISTANCE=100 VERBOSITY=INFO QUIET=false COMPRESSION_LEVEL=5 CREATE_INDEX=false CREATE_MSI=false
mriffiti@linus27 ~ $

```

Samples from M. Griffith RNAseq course at bioinformatics.ca

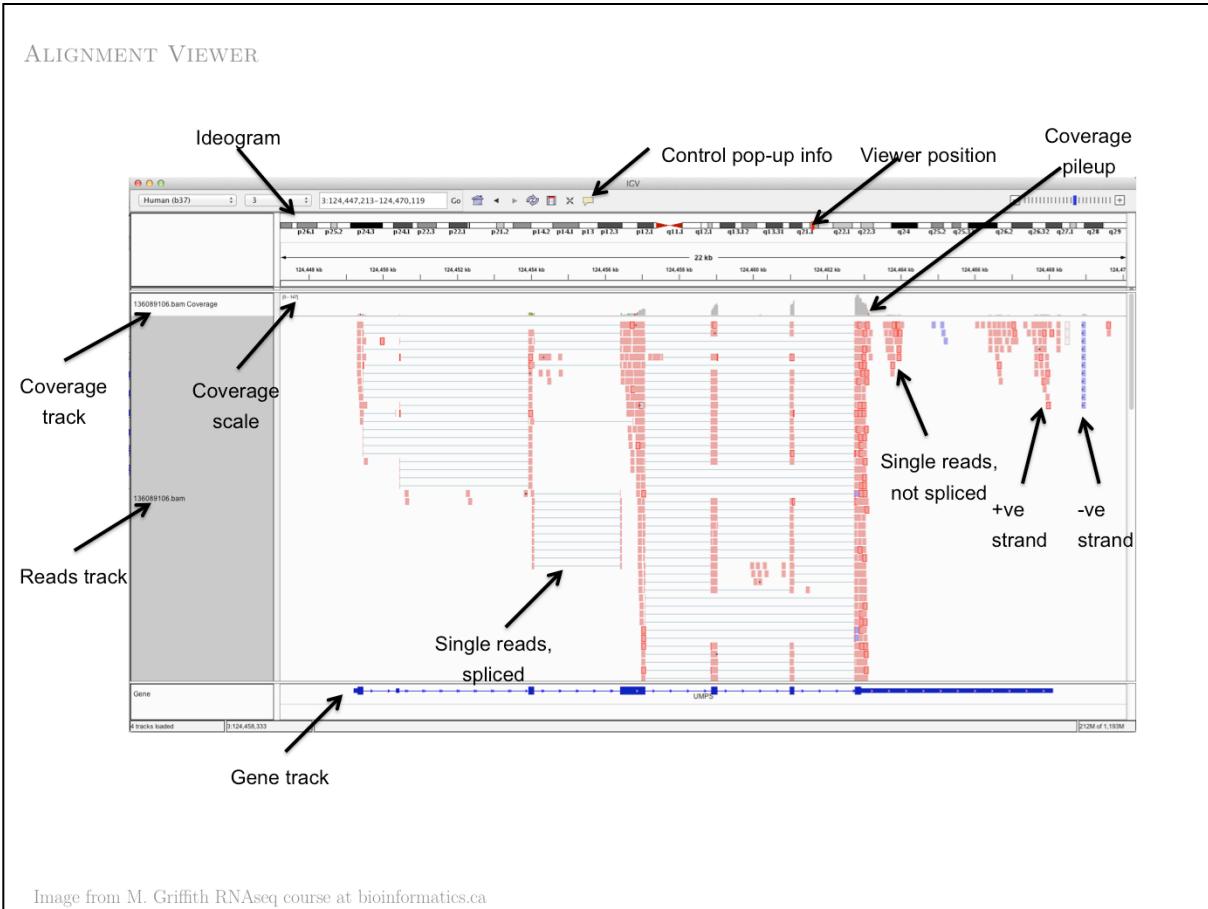
Header section

Alignment section

Aligned reads are stored in SAM (or BAM) files.

BED files are often used to specify regions of interest in a SAM/BAM file.

They contain Chromosome name, start, end and other annotation. and thus can be used to define complex subsets of a genome.



Alignment viewers give an overview of the alignment process.

READS TO COUNTS: THE QUESTION OF UNITS

- RPKM: Reads Per Kilobase of transcript per Million mapped reads.
- FPKM: Fragments Per Kilobase of transcript per Million mapped reads.
- In RNA-Seq, the relative expression of a transcript is proportional to the number of cDNA fragments that originate from it. However:
 - The number of fragments is also biased towards larger genes
 - The total number of fragments is related to total library depth
- FPKM (or RPKM) attempt to normalize for gene size and library depth
- $\text{RPKM} \text{ (or FPKM)} = (10^9 * C) / (N * L)$
 - C = number of mappable reads/fragments for a gene/transcript/exon/etc
 - N = total number of mappable reads/fragments in the library
 - L = number of base pairs in the gene/transcript/exon/etc

Cf. M. Griffith RNaseq course at bioinformatics.ca

After alignment, software exists to convert the mapped reads to counts. Different units are in use:

FPKM incorporates a probabilistic model to reduce reads to transcripts. This is good for calculating fold- changes of expression levels. It is a robust, widely accepted measure. The cufflinks program and other tools in the widely used Tuxedo suite use FPKM.

Raw counts can feed into a number of different statistical procedures, which is important for the expert. They allow for more sophisticated experimental design and analysis. Raw counts are used by DEseq and edgeR.

NCBI: GEO

The screenshot shows the NCBI GEO DataSet Browser interface. At the top, there's a search bar with the query "GDS2347[ACCN]". Below it, the dataset record for GDS2347 is displayed, including fields like Title, Summary, Organism, Platform, Citation, Reference Series, and Value type. To the right of the main content area, there are two boxes: one for "Cluster Analysis" showing a heatmap, and another for "Download" listing various file formats. At the bottom, there are links for NLM, NIH, GEO Help, Disclaimer, and Accessibility.

DataSet Record GDS2347: Expression Profiles Data Analysis Tools Sample Subsets			
Title:	Wild type strain across two cell cycles (I)		
Summary:	Analysis of wild type W303 cells across two cell cycles, a length of 2 hours after synchronization with alpha factor. Results compared to those from an experiment using a yox1 yhp1 double mutant strain (GDS2318).		
Organism:	<i>Saccharomyces cerevisiae</i>		
Platform:	GPL1914: FHCRC Yeast Amplicon v1.1		
Citation:	Pramila T, Miles S, GuhaThakurta D, Jemioło D et al. Conserved homeodomain proteins interact with MADS box protein Mcm1 to restrict ECB-dependent transcription to the M/G1 phase of the cell cycle. <i>Genes Dev</i> 2002 Dec 1;16(23):3034-45. PMID: 12464633		
Reference Series:	GSE3635	Sample count:	13
Value type:	log ratio	Series published:	2006/06/01

Data Analysis Tools

- Find genes
- Compare 2 sets of samples
- Cluster heatmaps
- Experiment design and value distribution

Find gene name or symbol:

Find genes that are up/down for this condition(s): time

NLM NIH GEO Help Disclaimer Accessibility

The GEO database at the NCBI hosts Microarray and RNA seq expression data.

EBI: EXPRESSION ATLAS

The screenshot shows the EBI Expression Atlas homepage. At the top, there's a search bar with placeholder text "Enter gene query..." and examples like "ASPM, Apoptosis, ENSMUSG0000021789, zinc finger". Below the search bar are tabs for Home, Release notes, FAQ, Download, Help (which is selected), Licence, and About.

On the left, there's a sidebar with a search input for "Search with gene attributes, e.g. gene symbols or annotations" and a "Release search" button. It also includes a section for "Select comparisons to see results for" and "Choose an adjusted p-value cutoff" (set to 0.05) and "Choose a log₂ fold-change cutoff" (set to 1.0).

The main content area has several sections:

- Heatmap:** A heatmap showing differential expression across various cell types. Colored boxes indicate differentially expressed genes. A tooltip for a specific gene shows its ID (ENSG0000021789) and statistics: Adjusted p-value: 2.47, Log₂ fold change: -1.08.
- Ensembl Gene Browser:** A table showing gene IDs and names, with a "Select a gene and a comparison to visualise at Ensembl" link.
- MA Plot and Gene Set Overlap Summaries:** A scatter plot showing MA plots and gene set overlap summaries.
- Network Analysis:** A network graph showing interactions between genes.
- GO and Interpro Terms:** A table showing GO and Interpro terms for a selected gene.

Annotations with arrows explain various features: "Greater colour intensity means larger absolute log₂ fold-change", "Download all statistics", "Click to see MA plot and gene set overlap summaries", "Click to see exact log₂ fold-changes in heatmap", "Select a gene and a comparison to visualise at Ensembl", "Coloured boxes mean genes differentially expressed. Mouseover a box to see statistics for a gene.", and "Mouseover a gene to see GO and Interpro terms. Click a gene for more detailed information."

<https://www.ebi.ac.uk/gxa/home>

The European source for expression data is the EBI Expression Atlas.

<http://steipe.biochemistry.utoronto.ca/abc>

B O R I S . S T E I P E @ U T O R O N T O . C A

DEPARTMENT OF BIOCHEMISTRY & DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO, CANADA