

Discovery & Modeling of Genomic Regulatory Networks with Big Data

Hamid Bolouri

Division of Human Biology

Fred Hutchinson Cancer Research Center

labs.fhcrc.org/bolouri

I have no financial relationships with any commercial entities.



Frontiers in Laboratory Science

Overview:

- (1) What do I mean by Genomic Regulatory Networks (GRNs)?
 - Why go genome-wide?
 - Why focus on regulatory interactions?
- (2) What do I mean by big (molecular biology) data?
- (3) Overview of integrative network discovery/modeling
- (4) Examples of available bioinformatics resources
- (5) Limitations and common pitfalls
- (6) Conclusion: *Proceed With Caution!*

❑ What are Genomic Regulatory Networks?

Molecular interaction networks regulating gene expression

❑ Why *genomic*?

Genome-wide for unbiased discovery of network components

Cost-effective exploration via high-throughput sequencing

❑ Why focus on *regulatory* interactions?

Determinants of cellular identity and state

❑ What are the main challenges?

Complex, intricate (much detail), and large-scale networks

Large-scale, noisy data

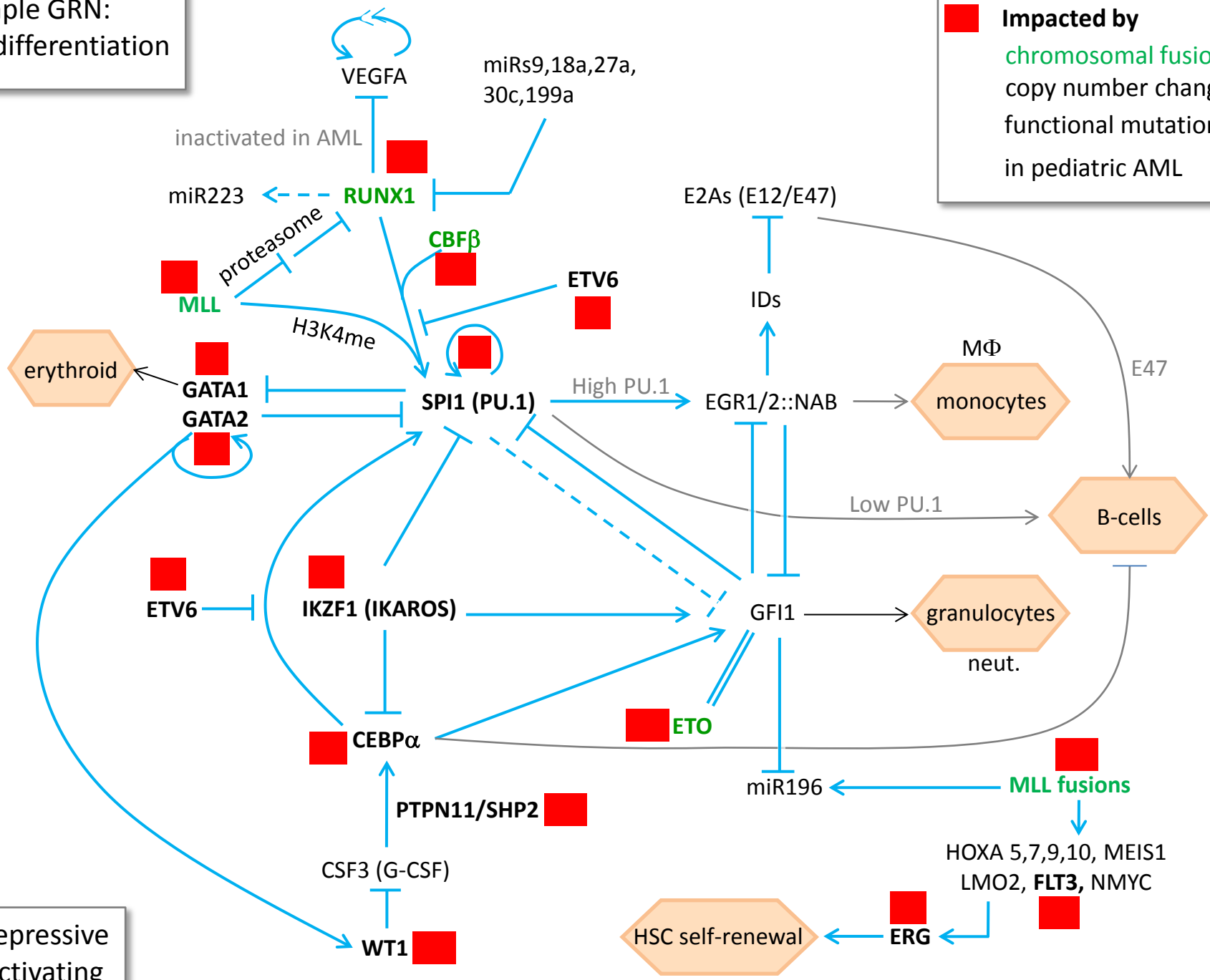
Example regulators of gene expression

- Enhancers, promoters, silencers, insulators/boundary elements
- Transcription factor state, concentration, localization, binding partners
- DNA methylation
- Histone/chromatin state regulation
- Histone variants
- Nucleosome remodeling
- 3D DNA conformation, nuclear localization
- miRNAs, other ncRNAs

Example processes regulating gene expression

- Chromatin accessibility
- TF complex occupancy
- Rate of RNA polymerase complex formation/transcription initiation
- Abortive & paused transcription initiation, elongation
- RNA processing (splicing, capping, editing, transport, etc.)
- Translation initiation and elongation
- Post-translational modifications

Example GRN:
HSC differentiation



Impacted by
chromosomal fusions,
copy number changes,
functional mutations
in pediatric AML

—| repressive
—> activating

Example sources of big data for GRN discovery and analysis

Public datasets:

- ✓ Gene Expression Omnibus (GEO): **1,335,996 samples** (ncbi.nlm.nih.gov/geo/). See also GTExportal.org
- ✓ 1000 Genomes Project: DNA sequence variants in **2,504 individuals** from 26 populations (www.1000genomes.org). See also NHLBI exome seq project: (evs.gs.washington.edu, **6,500 exomes**)
- ✓ Regulatory regions & interactions in DNA: ENCODE (**human 5,040 samples, mouse, 667 samples**) www.encodeproject.org/. See also modENCODE (fly, worm, www.modencode.org/)
- ✓ Roadmap Epigenomics Consortium. DNA methylation, histone mods, chromatin accessibility & small RNA transcripts in 23 stem cells and primary ex vivo tissues (**3,135 samples**, www.roadmapepigenomics.org)
- ✓ Protein localization in tissues (www.proteinatlas.org/), & organelles (locate.imb.uq.edu.au/)
- ✓ The Cancer Cell Line Encyclopedia (CCLE): sequence, expression, drug sensitivity in ~ **1000 cell lines** www.broadinstitute.org/ccle . See also: www.cancerrxgene.org/ , discover.nci.nih.gov/cellminer , cancer.sanger.ac.uk/cancergenome/projects/cell_lines/
- ✓ NCI TCGA, TARGET projects. Comprehensive multi-omics and clinical patient data (see also ICGC).

Your (genome-wide) data:

- ✓ DNA-seq (whole genome / exome / targeted sequencing)
- ✓ RNA-seq (mRNA-seq, miRNA-seq, ribosome profiling)
- ✓ CHIP-seq (TFs, pol2, histone modifications, methylated-DNA seq). Also DNase/MNase seq, FAIRE, etc.
- ✓ Chromatin interactions, location, ...

Examples of combining in-lab and 3rd party data to discover GRNs

Use data from related cell types to augment/test in-lab findings, e.g.:

- ❑ Correlate mRNA data with 3rd party
 - miRNA expression
 - miR binding motifs in UTR sequences
 - DNA-methylation

- ❑ Match up-regulated transcription factors (TFs) to
 - TF ChIP-seq peaks
 - TF binding motifs in open chromatin (e.g. DNase1 HS / footprints) regions

- ❑ Match DNA-sequence variations to
 - Changes in mRNA isoforms
 - mRNA/miRNA expression
 - Changes in correlated gene expression

GRN discovery/modeling **example approach (1):**

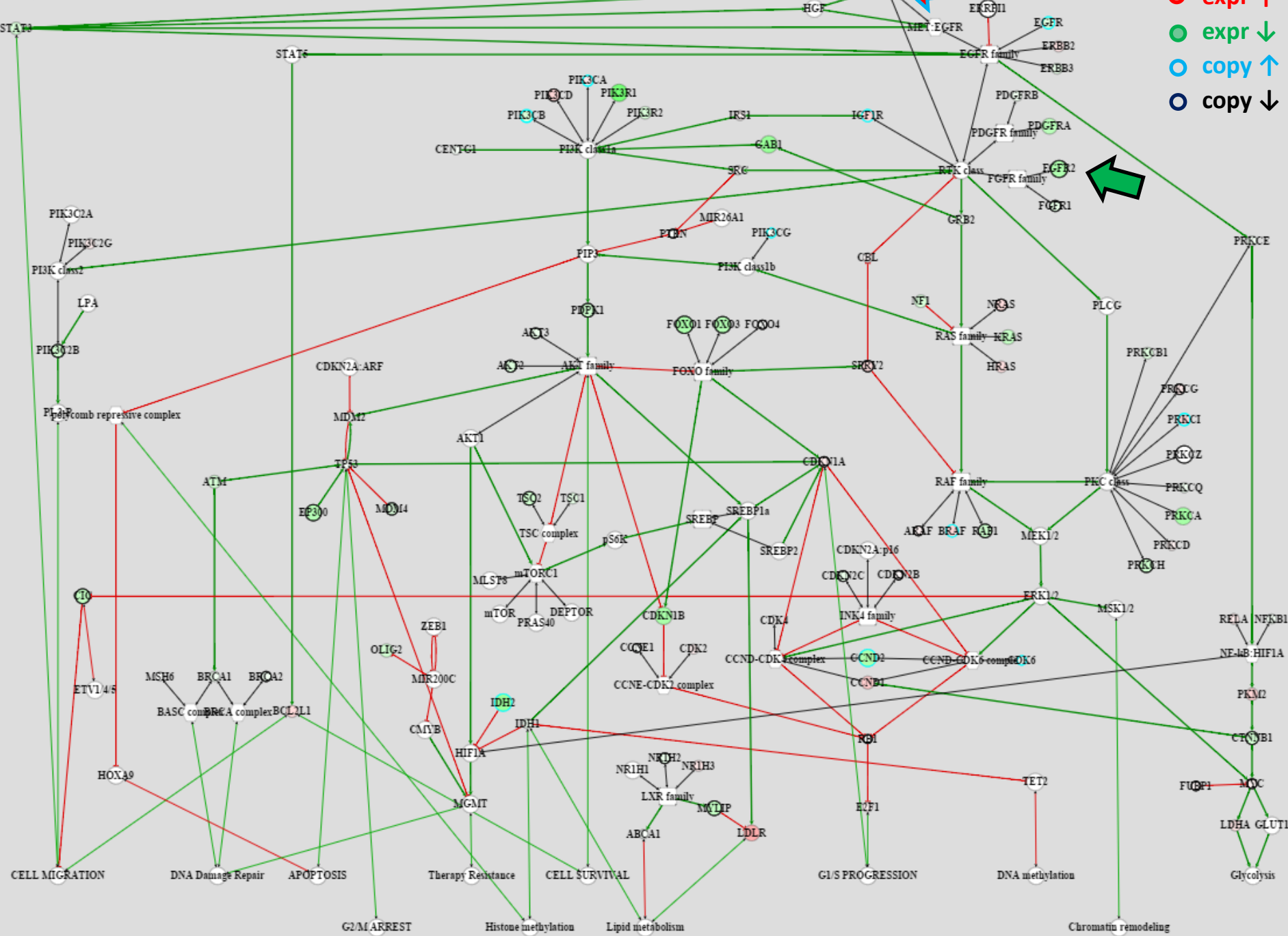
visual exploration of 'omics data super-imposed on an interaction/pathway map

Example software:

Cytoscape: cytoscape.org (see also <http://js.cytoscape.org>)

BipTapestry: <http://www.biotapestry.org/#download>

- expr ↑
- expr ↓
- copy ↑
- copy ↓



CELL MIGRATION DNA Damage Repair APOPTOSIS Therapy Resistance CELL SURVIVAL G1/S PROGRESSION DNA methylation Chromatin remodeling Lipid metabolism Histone methylation G2/M ARREST

GRN discovery/modeling **example approach (2):**

GRN inference by integration of complementary 'omics data sets

Example software:

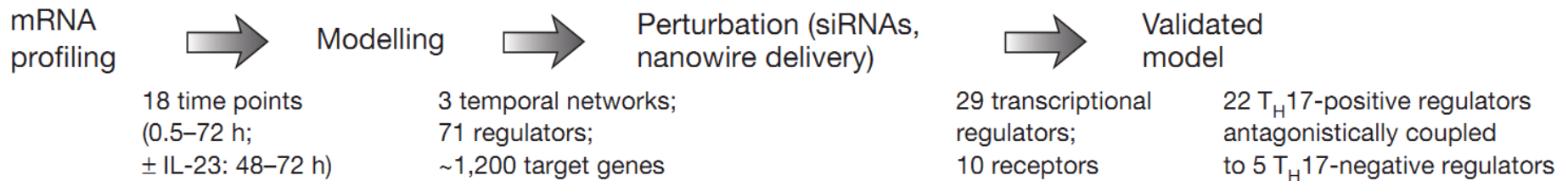
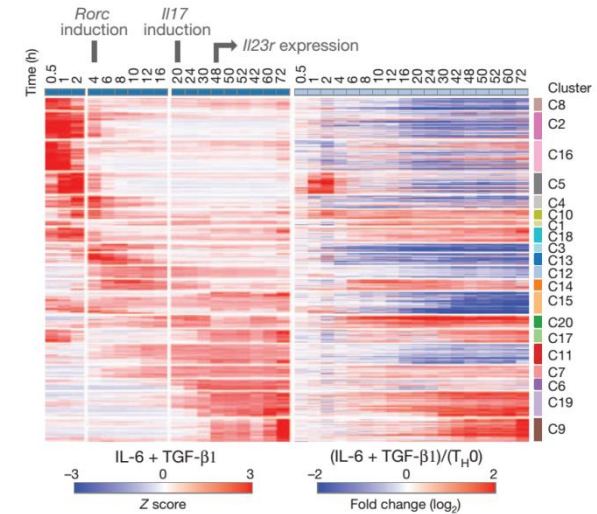
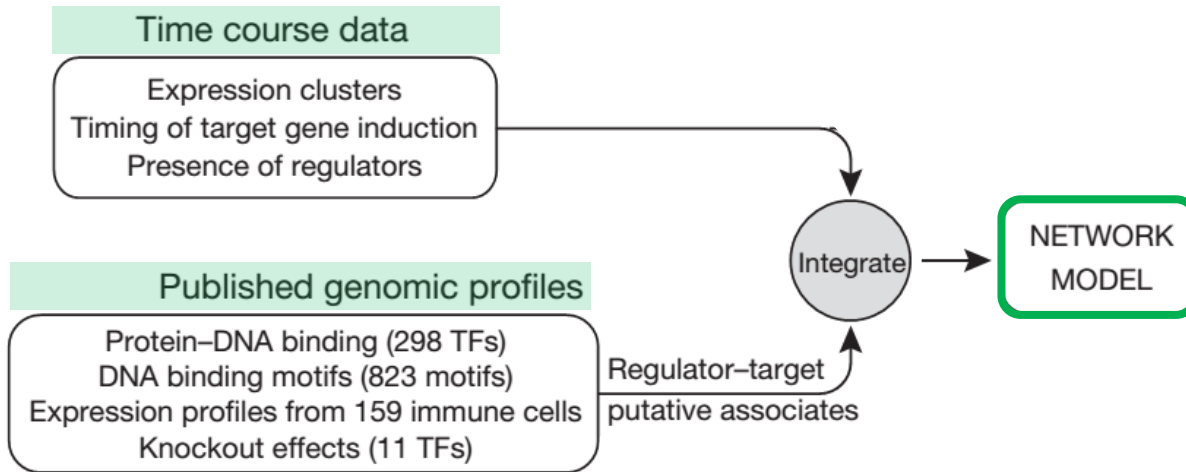
PARADIGM,

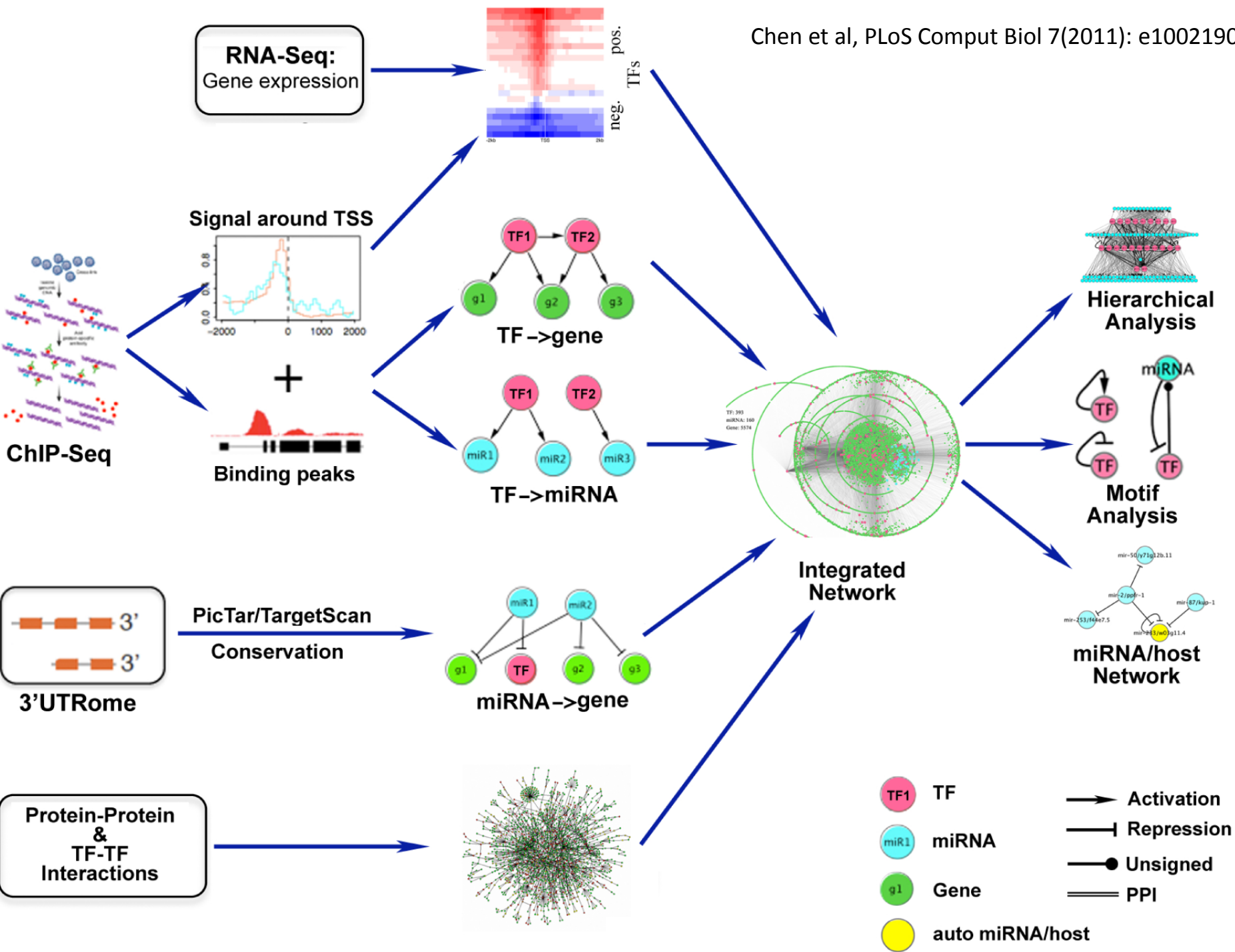
Vaske, Benz et al, Bioinformatics (2010) 26 (12):i237-i245.

sbenz.github.com/Paradigm

Dynamic regulatory network controlling T_H17 cell differentiation

Nir Yosef, Aviv Regev & colleagues





Examples of available bioinformatics resources

... for people who do not write Perl, Python, R (Bioconductor), etc.

- Tools
- NGS: QC and manipulation
- NGS: Mapping
- NGS: SAM Tools
- snpEff
- BEDTools
- NGS: Variant Analysis
- Genome Diversity
- NGS: RNA-seq
 - Cuffcompare compare assembled transcripts to a reference annotation and track Cufflinks transcripts across multiple experiments
 - Cufflinks transcript assembly and FPKM (RPKM) estimates for RNA-Seq data
 - Cuffmerge merge together several Cufflinks assemblies
 - Filter_Combined_Transcripts using tracking file
 - Tophat for Illumina Find splice junctions using RNA-seq data
 - Tophat2 Gapped-read mapper for RNA-seq data
 - Cuffdiff find significant changes in transcript expression, splicing, and promoter use
 - Tophat Fusion Post post-processing to identify fusion genes
- Operate on Genomic Intervals
- NGS: GATK Tools (beta)
- NGS: VCF Manipulation
- EMBOSS
- Regional Variation
- FASTA manipulation
- Evolution
- NGS: Picard (beta)
- Multiple Alignments
- Metagenomic analyses
- NGS: Peak Calling
- Motif Tools
- Workflows
 - All workflows

Cufflinks version 0.0.7 Help from Biostar

SAM or BAM file of aligned RNA-Seq reads:

Max Intron Length:

Min Isoform Fraction:

Pre mRNA Fraction:

Perform quartile normalization:

Removes top 25% of genes from FPKM denominator to improve accuracy of differential expression calls for low abundance transcripts.

Use Reference Annotation:

Perform Bias Correction:

Bias detection and correction can significantly improve accuracy of transcript abundance estimates.

Use multi-read correct:


Tells Cufflinks to do an initial estimation procedure to more accurately weight reads mapping to multiple locations in the genome.

Use effective length correction:

Cufflinks will not employ its 'effective' length normalization to transcript FPKM.

Cufflinks Overview

Cufflinks assembles transcripts, estimates their abundances, and tests for differential expression and regulation in RNA-Seq samples. It accepts aligned RNA-Seq reads and assembles the alignments into a parsimonious set of transcripts. Cufflinks then estimates the relative abundances of these transcripts based on how many reads support each one. Please cite: Trapnell C, Williams BA, Pertea G, Mortazavi AM, Kwan G, van Baren MJ, Salzberg SL, Wold B, Pachter L. Transcript assembly and abundance estimation from RNA-Seq reveals thousands of new transcripts and switching among isoforms. Nature Biotechnology doi:10.1038/nbt.1621

Know what you are doing  **NB!**

⚠ There is no such thing (yet) as an automated gearshift in expression analysis. It is all like stick-shift driving in San Francisco. In other words, running this tool with default parameters will probably not give you meaningful results. A way to deal with this is to **understand** the parameters by carefully reading the [documentation](#) and experimenting. Fortunately, Galaxy makes experimenting easy.

Input formats

Cufflinks takes a text file of SAM alignments as input. The RNA-Seq read mapper TopHat produces output in this format, and is recommended for use with Cufflinks. However Cufflinks will accept SAM alignments generated by any read mapper. Here's an example of an alignment Cufflinks will accept:

```
s6.25mer.txt-913508 16 chr1 4482736 255 14M431N11M * 0 0 CAAGATGCTAGGCAAGTCTGGAAG IIIIIIIIIIIIIIIIIIIIIIIIIIIII NM:i:0 XS:A:-
```

History

search datasets

Unnamed history

0 bytes

This history is empty. You can load your own data or get data from an external source

Favorite Modules

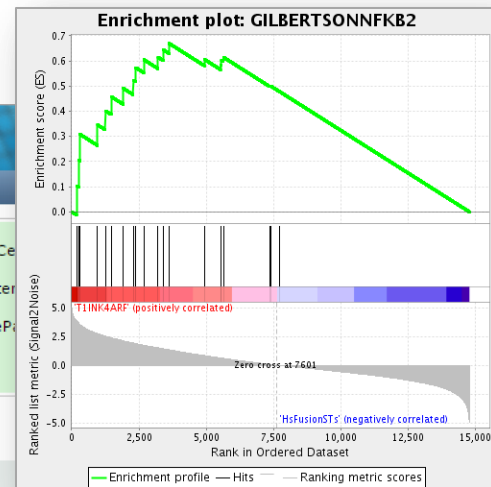
ssGSEAProjection

Recent Modules

GSEA

GSEALeadingEdgeViewer

-- Feb 11 -- **FCSNormalization** v2 now available (via GParc and BC Cancer Research Ce
-- Jan 20 -- **IGV** updated: addresses previous issues with launching IGV from GenePatter
-- 9:10am Jan 20 -- The software update to 3.9.1 is complete. The "what's new in GeneP
gp-help with your questions and comments.
Thanks!



GSEA version 14

Gene Set Enrichment Analysis

* required field

Reset Run

expression dataset*

Upload File...

Add Path or URL...

Drag Files Here

Batch

2GB file upload limit using the Upload File... button. For files > 2GB upload from the Files tab.

Dataset file - .res, .gct

gene sets database

Gene sets database from GSEA website.

gene sets database file

Upload File...

Add Path or URL...

Drag Files Here

Batch

2GB file upload limit using the Upload File... button. For files > 2GB upload from the Files tab.

Gene sets database - .gmt, .gmx, .grp. Upload a gene set if your gene set is not listed as a choice for the gene sets database parameter.

number of permutations*

1000

Number of permutations to perform

phenotype labels*

Upload File...

Add Path or URL...

Drag Files Here

Batch

2GB file upload limit using the Upload File... button. For files > 2GB upload from the Files tab.



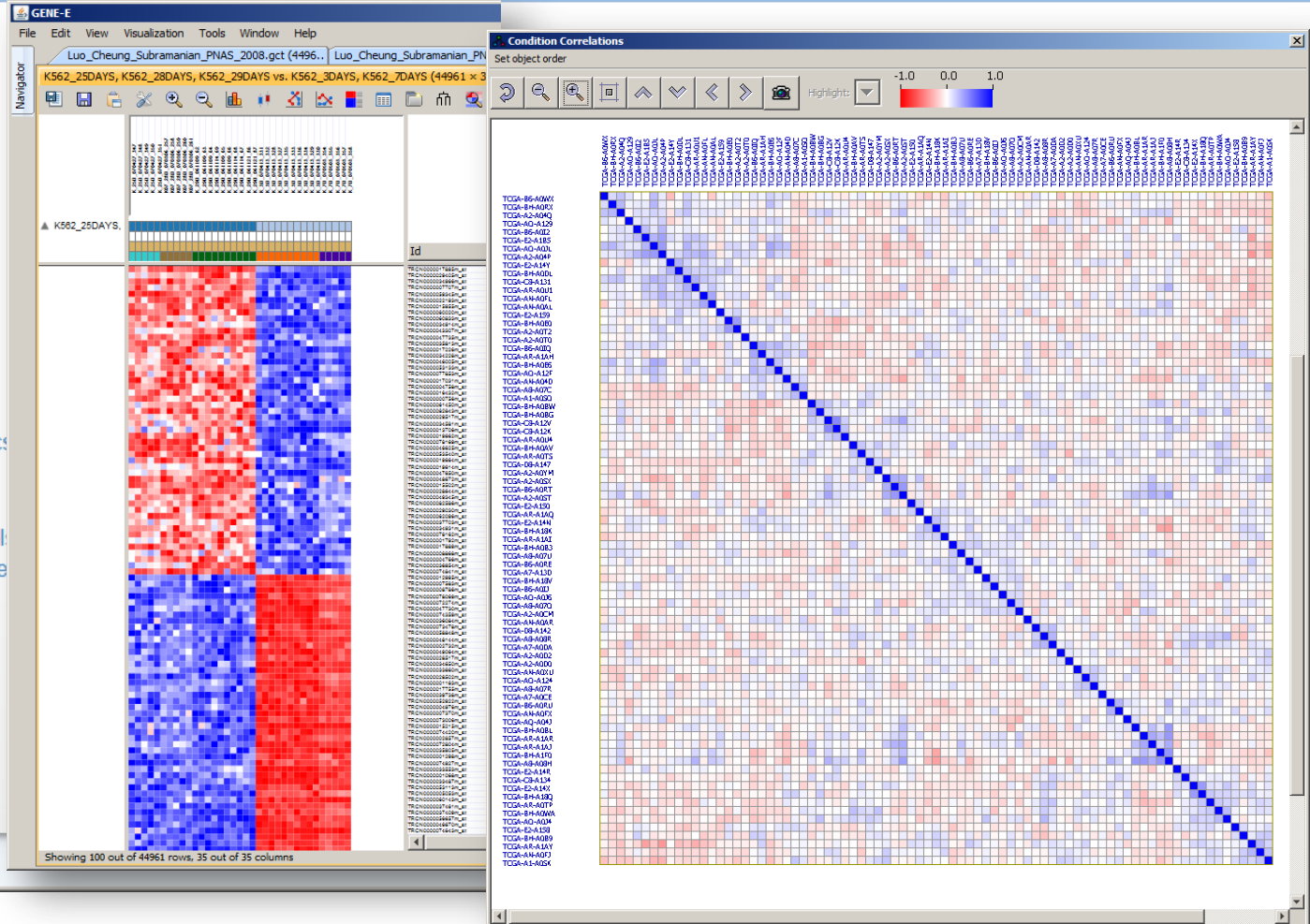
User Login

What is GenomeSpace? Tools Recipes Documentation Developers Support About

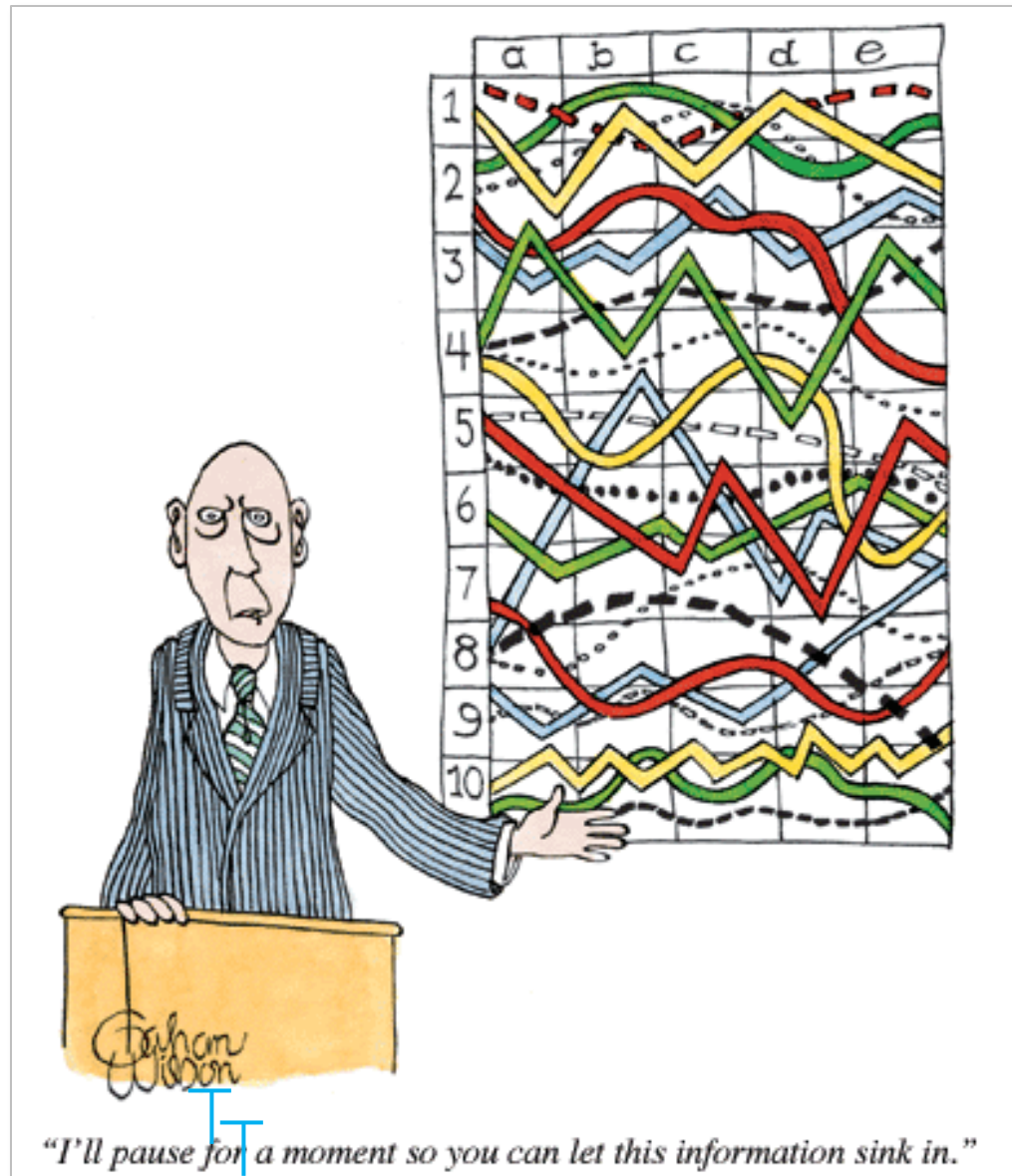


Tool Guide

- Introduction
- ArrayExpress
- Cistrome
- Cytoscape
- Galaxy
- GenePattern
- Genomica
- geWorkbench
- Gitools
- InSilico DB
- Integrative Genomics Viewer (IGV)
- ISAcreeator
- MSigDB Online Tool
- UCSC Table Browse
- Tool Term Glossary
- Cytoscape 3

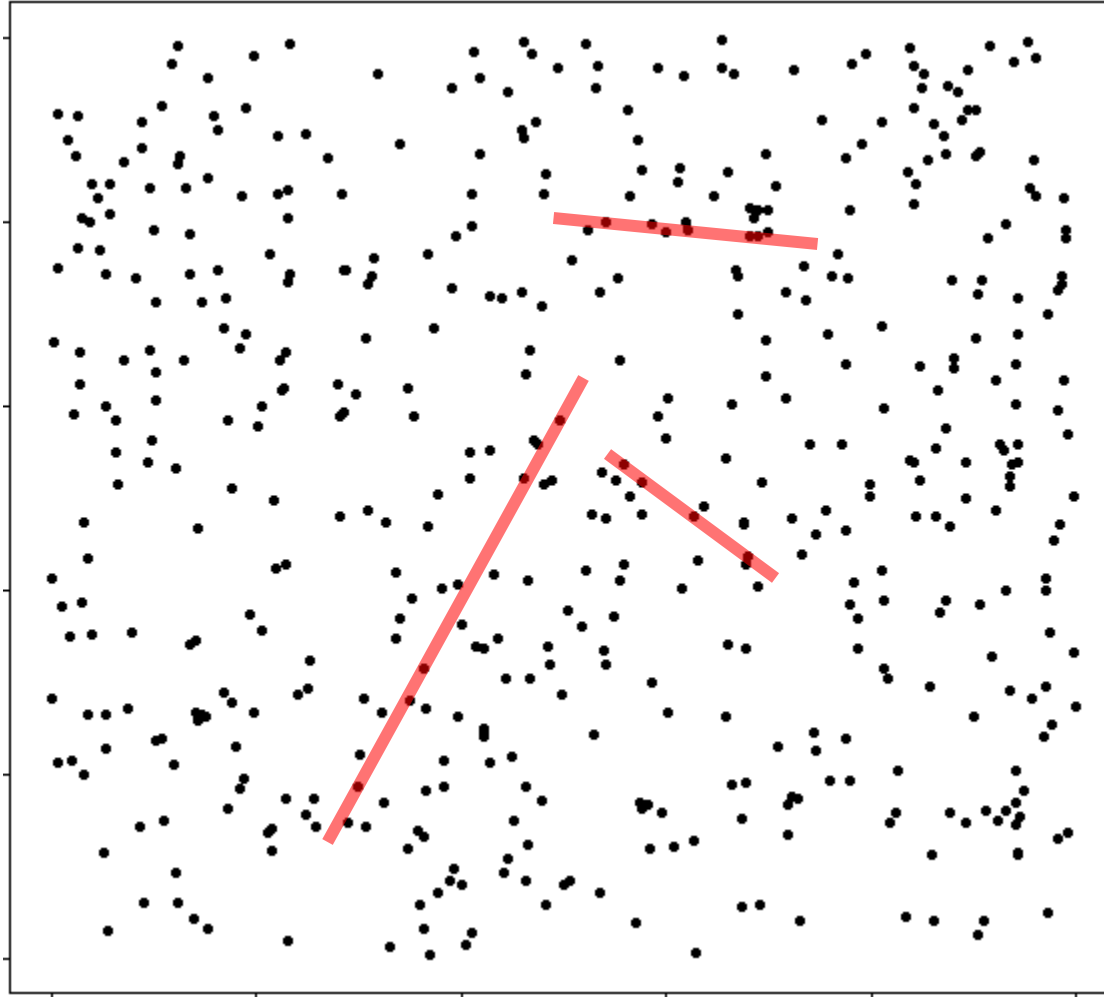


Limitations and common pitfalls



In large-scale (e.g. genome-wide) data, patterns will occur by chance

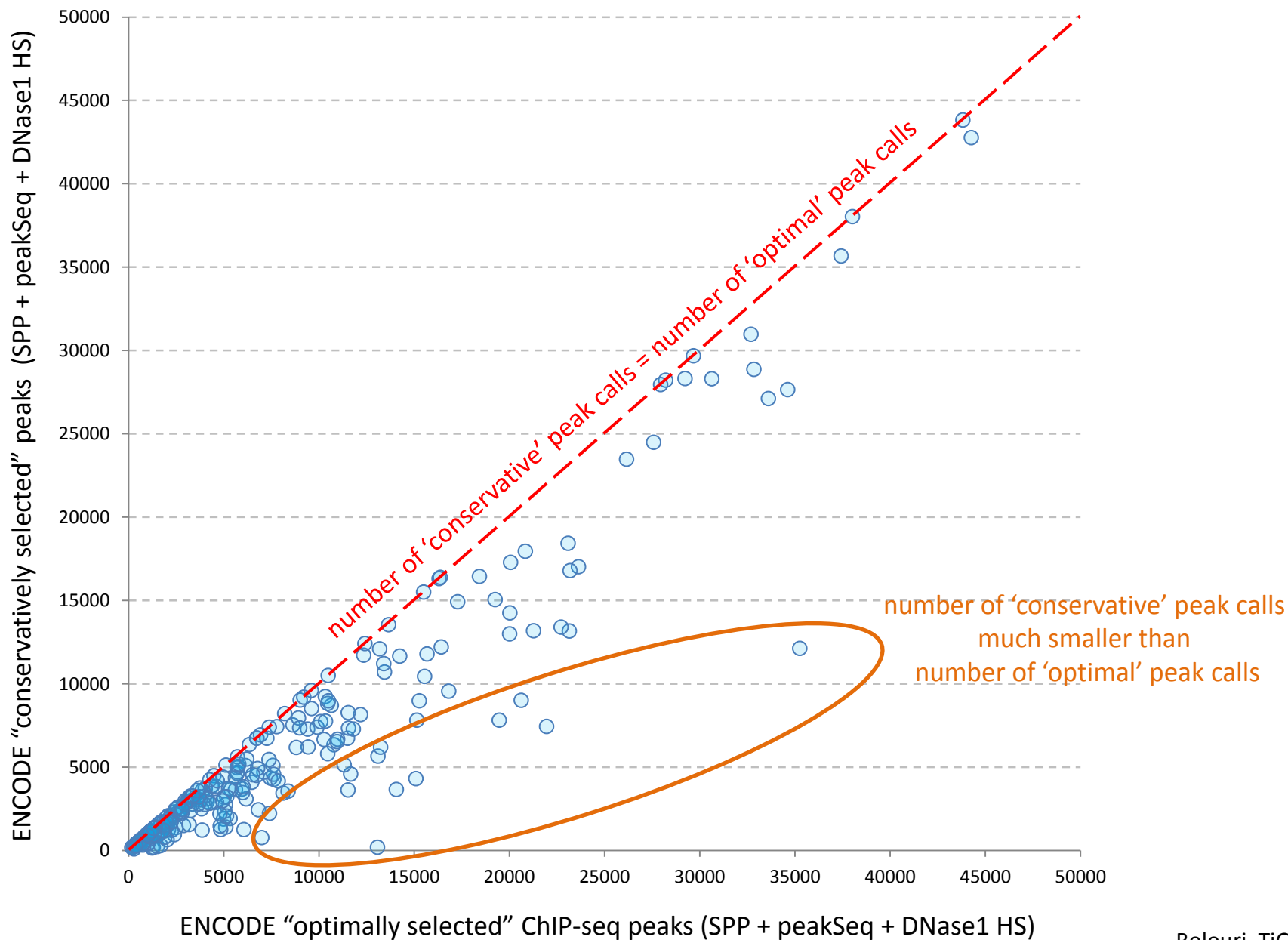
Always test the statistical significance of finding a pattern



Simpson's paradox

Results are often statistically reliable, but individually noisy

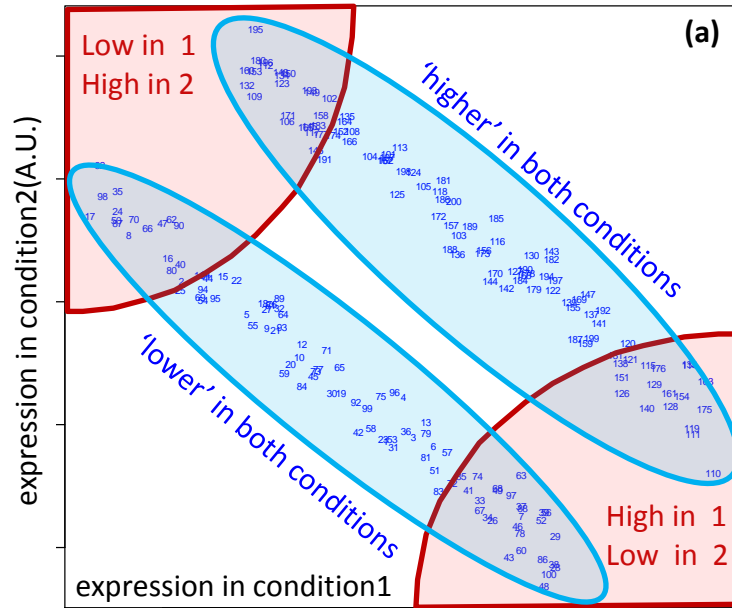
Per sample variability in ENCODE ChIP-seq peak calls (254 transcription factors)



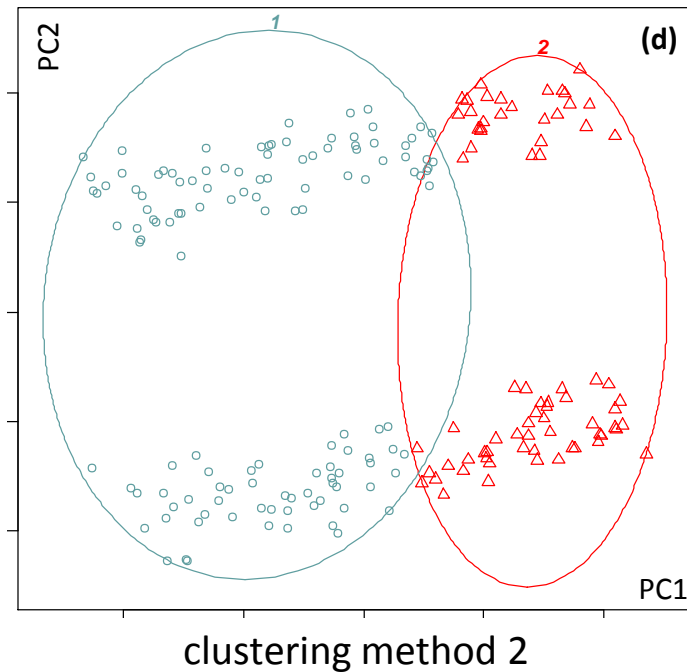
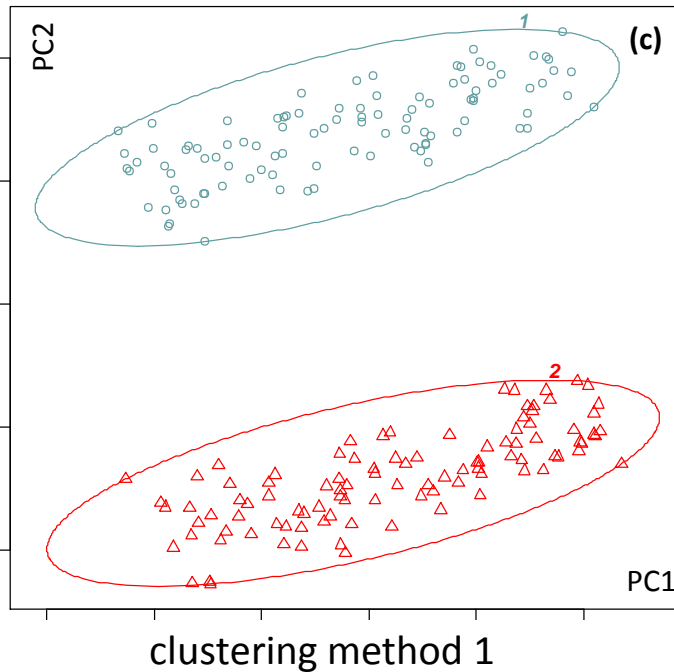
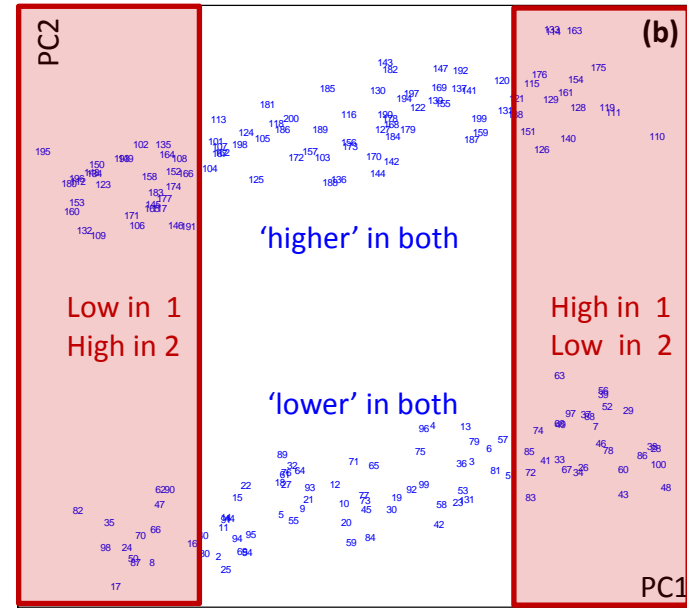
The importance of knowing your tools:

choices of statistical methods & parameters affect results

“original data”

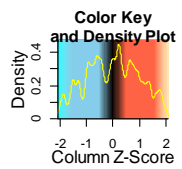


Principal Component Analysis

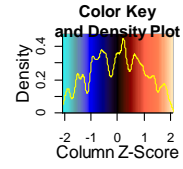


The importance of full disclosure

Example: how a missing color-scale key can mislead



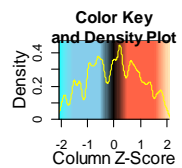
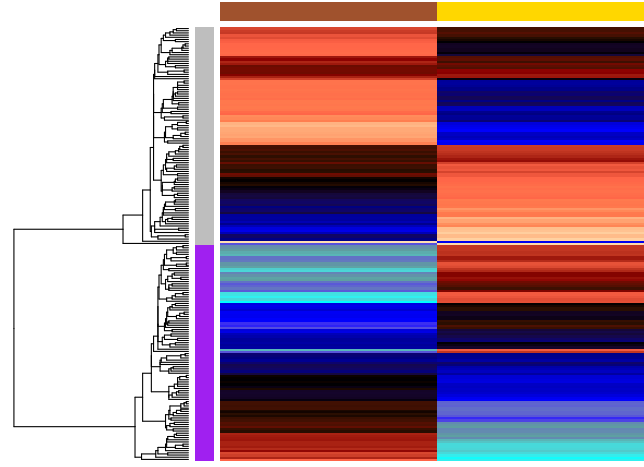
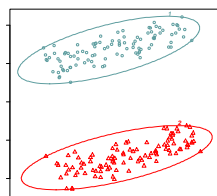
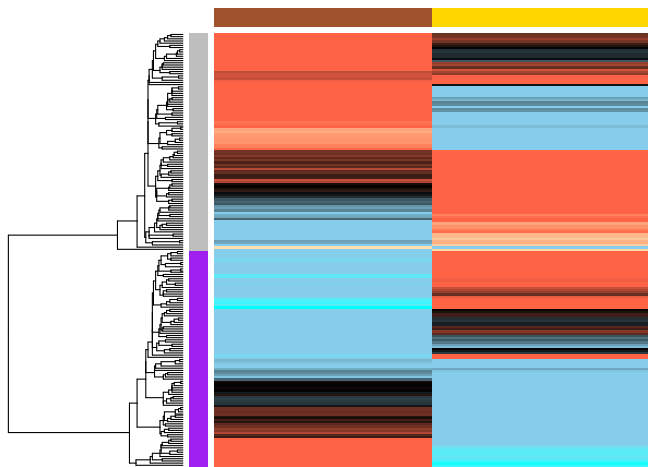
conditions



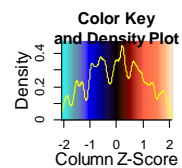
conditions

These 2 plots show identical data

genes



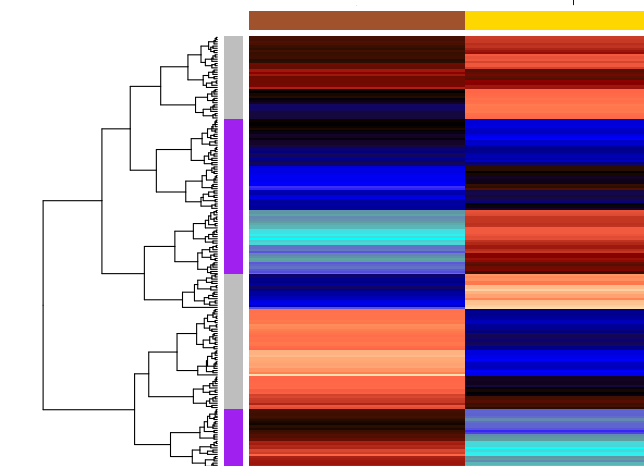
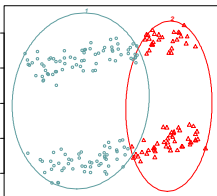
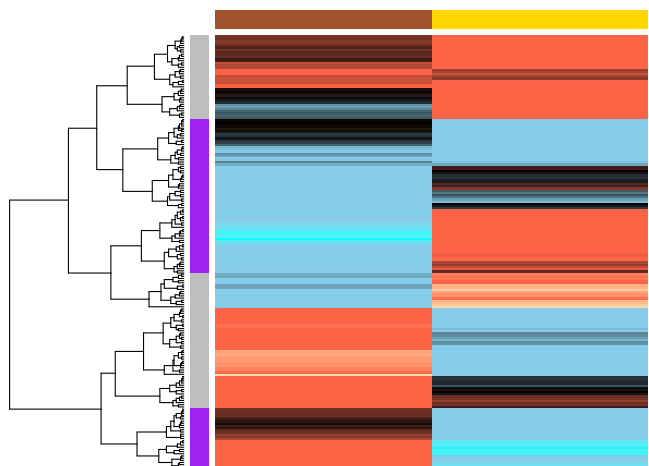
conditions



conditions

These 2 plots show identical data

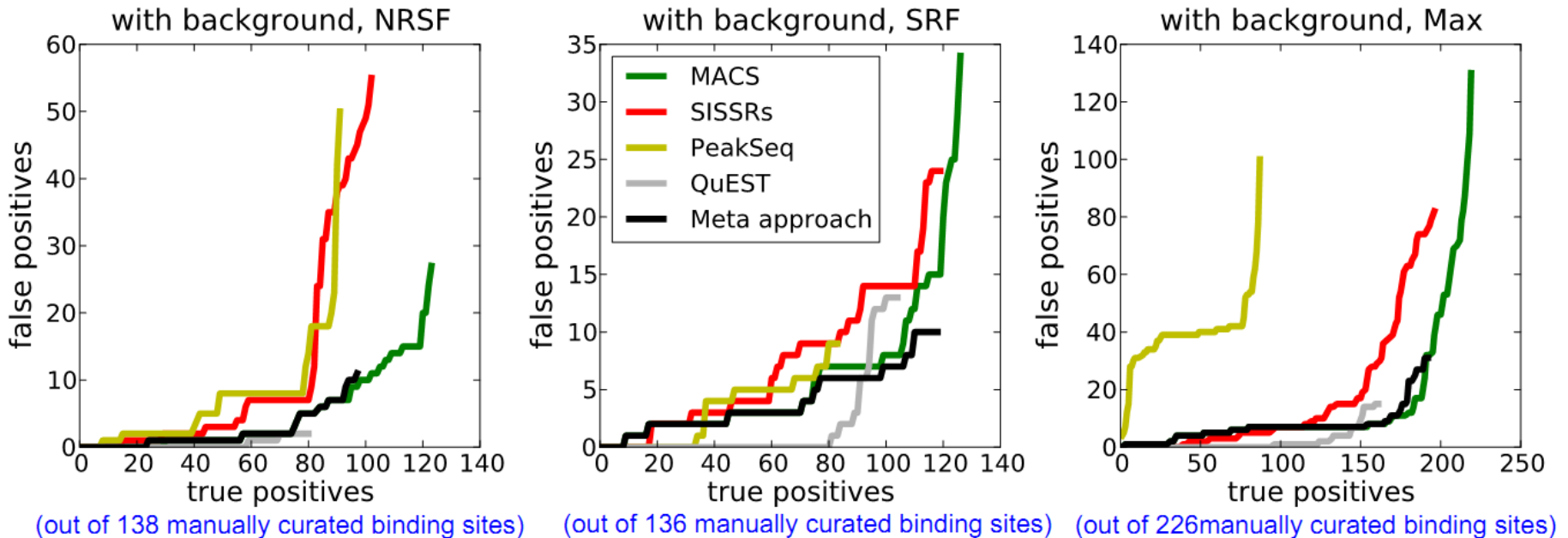
genes



- ❑ Process data with multiple algorithms and parameters
- ❑ When possible use matched controls

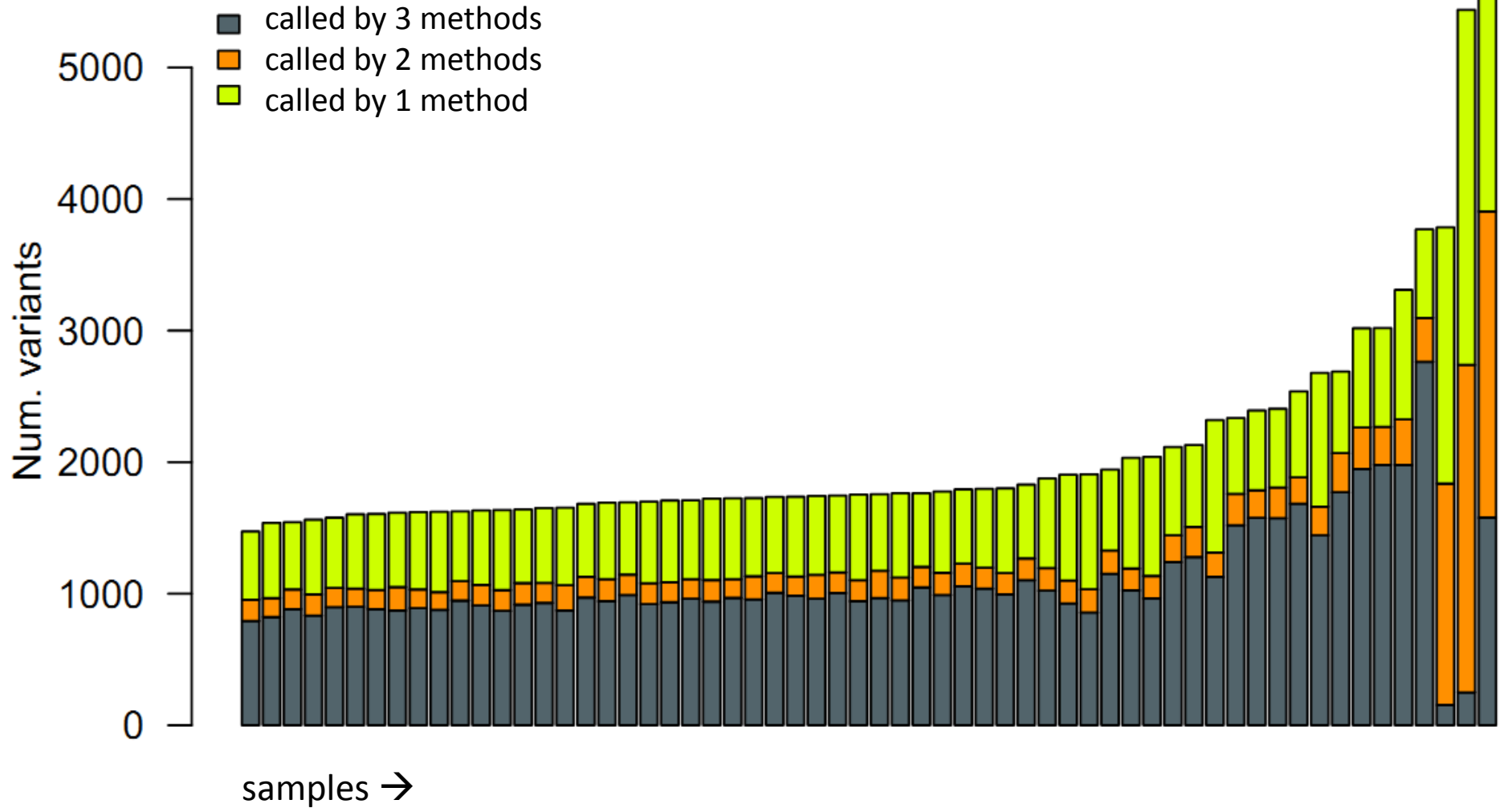
A manually curated ChIP-seq benchmark demonstrates room for improvement in current peak-finder programs

Morten Beck Rye^{1,*}, Pål Sætrom^{1,2} and Finn Drabløs¹



SNAs detected by 3 methods applied to the same 80X tumor exomes

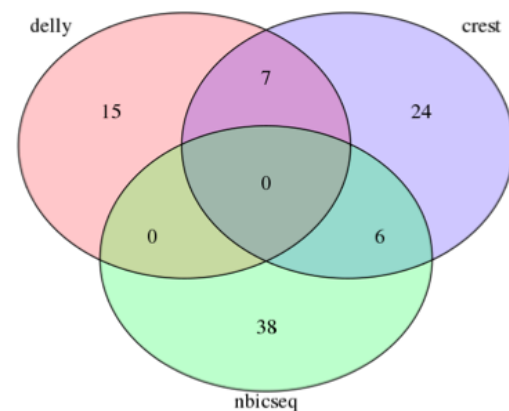
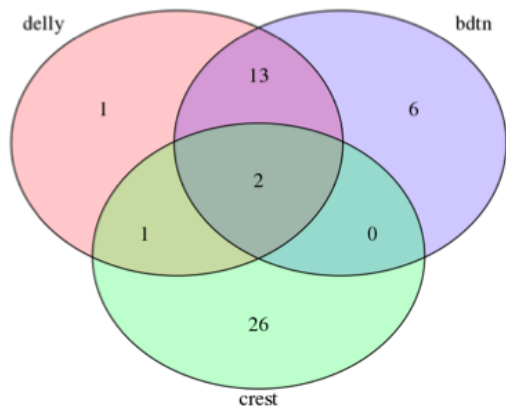
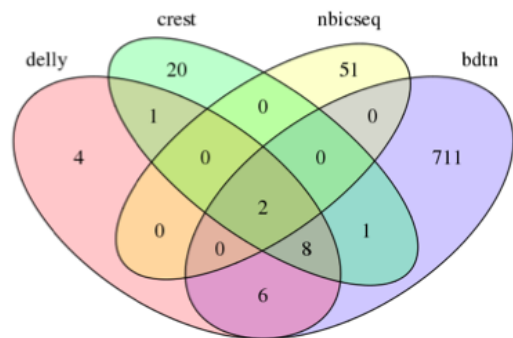
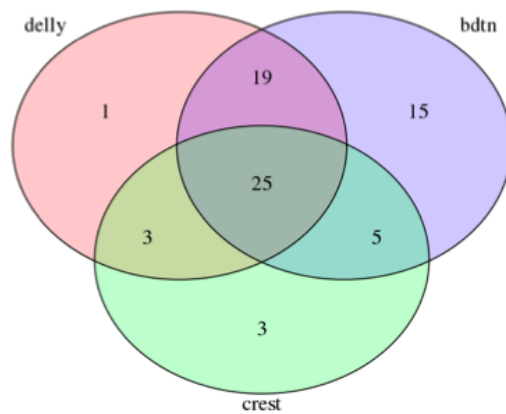
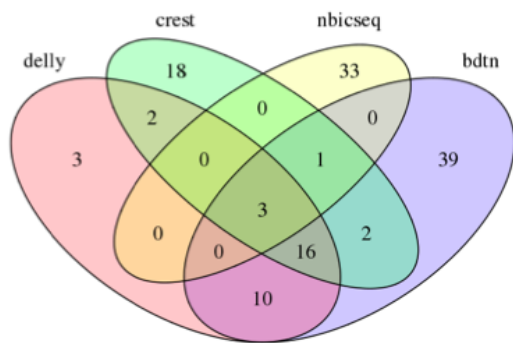
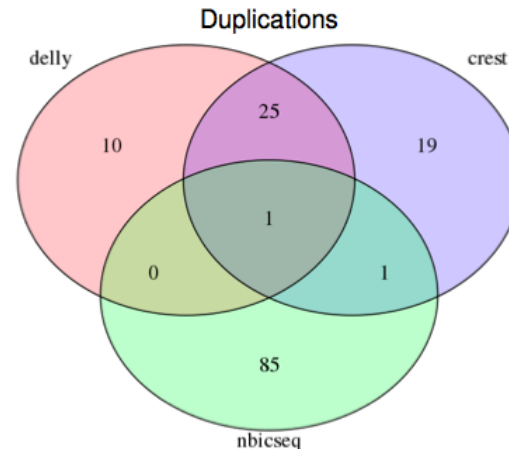
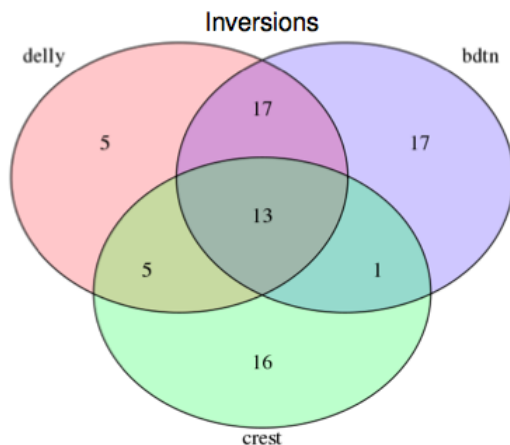
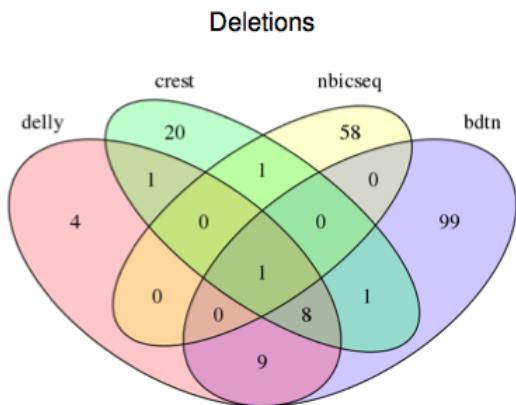
(methods: MuTect, Strelka, Virmid)



2x125bp HiSeq 2500 (high output mode),
Library Prep: Agilent SureSelectXT V4 71mb+UTR Exome

Overlap among four structural variation callers for three 30X WGS tumor samples

2x150bp HiSeq X, Library Prep: TruSeqDNA Nano 350bp



Overview:

- (1) What do I mean by Genomic Regulatory Networks (GRNs)?
 - Why go genome-wide?
 - Why focus on regulatory interactions?
- (2) What do I mean by big (molecular biology) data?
- (3) Overview of integrative network discovery/modeling
- (4) Examples of available bioinformatics resources
- (5) Limitations and common pitfalls
- (6) Conclusion: ***Proceed With Caution!***

Indels detected by 2 methods applied to the same 80X tumor exomes

(methods: SID, Strelka)

