

Data Visualization & Annotation

8th International Biocuration Conference

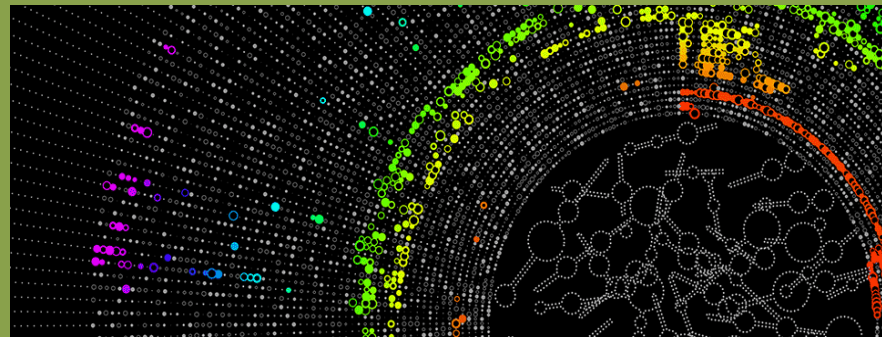
24 April 2015 | Beijing, China

Rama Balakrishnan

Saccharomyces Genome Database
Gene Ontology Consortium
Stanford University, CA, USA

Monica Munoz-Torres

Berkeley Bioinformatics Open-Source Projects
Lawrence Berkeley National Lab, CA, USA



Outline

1. Introduction

- Goals
- Examples of genome visualization tools

2. Panelists

- **Lorna Richardson:**
eMouseAtlas and Image Informatics
- **Justyna Szostak:**
Curated Causal Biological Network Models

3. Discussion

- Featuring you!

Goals of the workshop

1. To learn about tools available for human interpretation of genomic data, specifically in the context of annotation.
2. To open a space for discussion: genomic data are ever more abundant and heterogeneous, with widely varied sources, production techniques, and intrinsic experimental error.
 - How do we analyze these data?
 - What is the best way to interpret the stories the data are telling us?
 - How to put these together (overlay) visually?
 - Developers: what is the best way to disseminate and contribute code to make tool development easier?

Then and Now

Figures 1
Tables 0
References 6

Figures 49
Tables 27
References 452



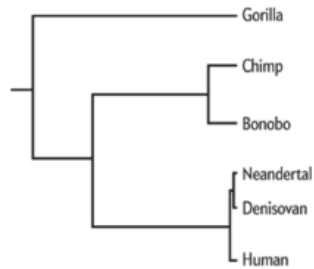
Genomic Data:

Heterogeneous & Abundant

- Structural: gene models, transcriptomes, RNAseq, differential expression, etc.
- Functional: gene ontology, interactions, phenotypes, SNPs, complexes, protein abundance, diseases, images, etc.
- Some examples . . .

Genomic differences

Family Tree



How to Read This Graphic

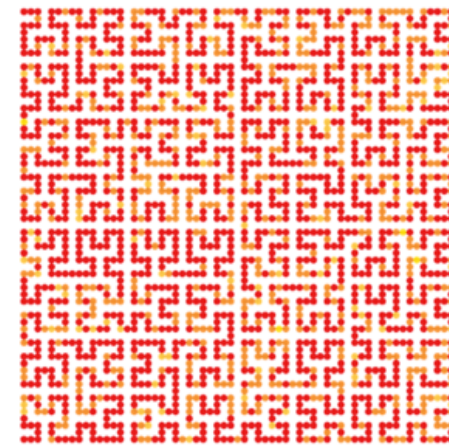
Each dot represents a sequence of about 500,000 pairs of chemical bases—the A, T, C and G of our genetic code—in the protein-coding portion of the human genome in the order that they appear on our chromosomes.

ATGCCCGTTCTGAA ...



The color of the dot indicates how well the human sequence matches up with the corresponding sequence in the comparison species, with red signifying a greater difference between the two.

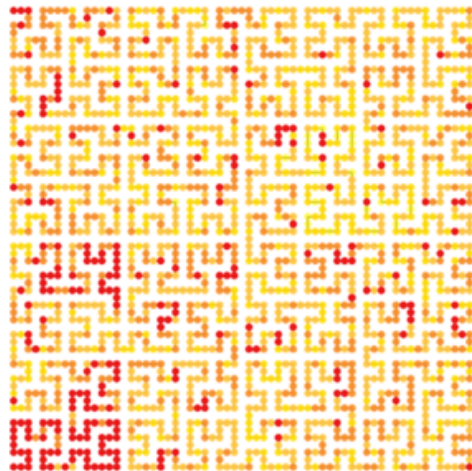
Fraction of different or unaligned bases (%)



Gorilla



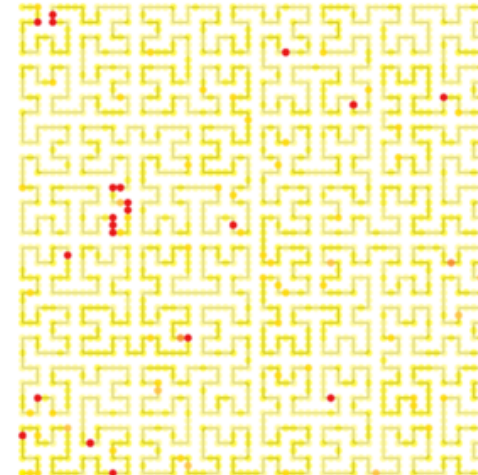
On the whole, our coding genome differs more from the gorilla's than from the chimp's or the bonobo's, reflecting the fact that we have been evolving along separate trajectories for a longer period. But about 15 percent of the human genome looks more like the gorilla's than the chimp's or the bonobo's.



Chimp



Researchers have traditionally considered the chimpanzee, which lives in patriarchal societies, to be our closest living relative and thus the best model for reconstructing the lives of ancient human ancestors. The recent genome-sequencing work calls that view into question, however.



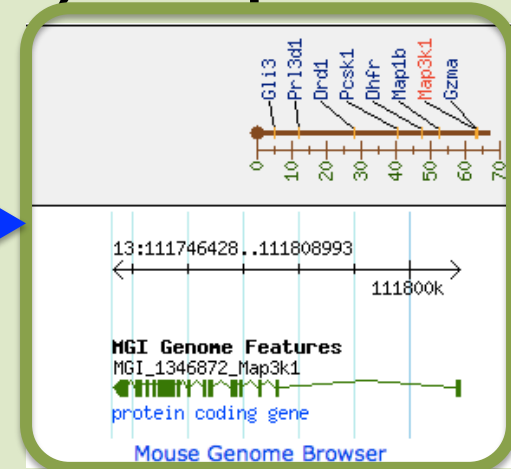
Denisovan



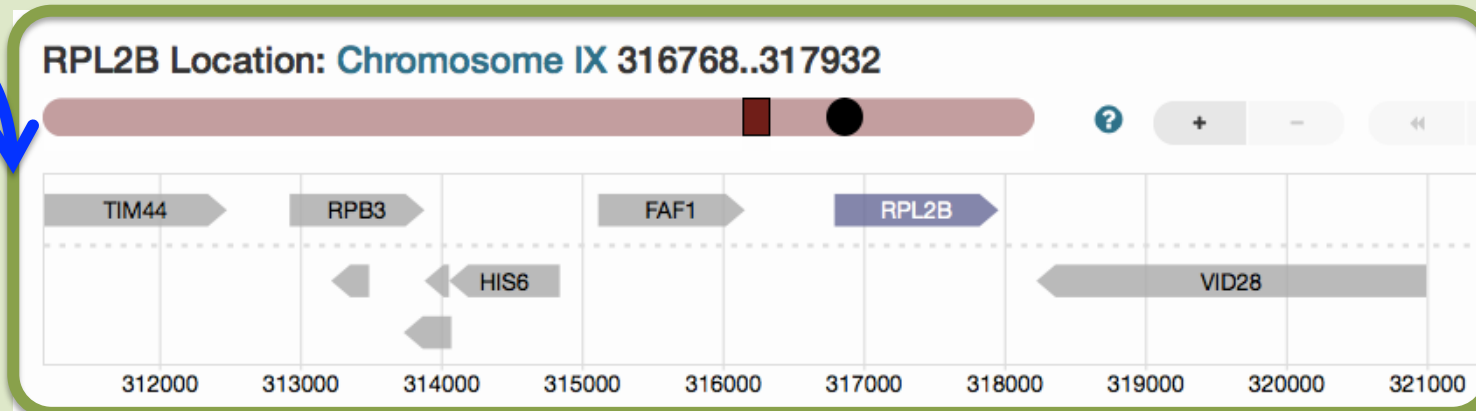
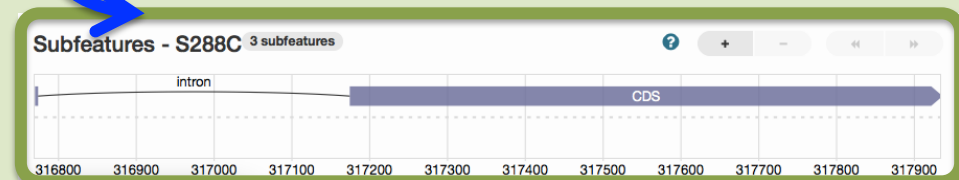
The Denisovans—a group of archaic humans closely related to the Neanderthals—show far fewer sequence differences from us than any of the African apes do, having shared a common ancestor with *H. sapiens* in the much more recent past, around 400,000 years ago.

Gene structure, ideograms, maps

Genetic Map	Chromosome 13 63.36 cM Detailed Genetic Map \pm 1 cM Mapping data(3)
Sequence Map	Chr13:111746428-111808993 bp, - strand From VEGA annotation of GRCm38



Feature	Relative Coordinates	Coordinates
CDS	1..4	chrIX:316768..316771
intron	5..404	chrIX:316772..317171
CDS	405..1165	chrIX:317172..317932



Most of the curated data is text

Gene Ontology ⁱ

[Gene Ontology Details](#)

Summary: Subunit of the PeBoW complex and the preribosomal complex; binds to the large ribosomal subunit rRNA and is involved in the processing of the rRNA and the biogenesis of the large ribosomal subunit

[View computational annotations](#)

Molecular Function

Manually Curated: • [large ribosomal subunit rRNA binding \(IDA\)](#)

Biological Process

Manually Curated: • [ribosomal large subunit biogenesis \(IPI\)](#)
• [rRNA processing \(IMP\)](#)

Cellular Component

Manually Curated: • [PeBoW complex \(IDA, IPI\)](#)
• [preribosome, large subunit precursor \(IDA\)](#)

High-Throughput: • [nucleolus \(IDA\)](#)
• [nucleus \(IDA\)](#)
• [preribosome, large subunit precursor \(IDA\)](#)

Phenotype ⁱ

[Phenotype Details](#)

Classical Genetics

conditional: • [cell cycle passage through the metaphase-anaphase transition: delayed](#)
• [resistance to hydroxyurea: decreased](#)
• [spindle morphology: abnormal](#)

Large-scale Survey

reduction of • [competitive fitness: decreased](#)
function: • [chromosome/plasmid maintenance: abnormal](#)
conditional: • [colony sectoring: increased](#)
null: • [haploinsufficient](#)
• [inviable](#)

Visualizing interaction data

Interaction i

249 total interactions for 129 unique genes

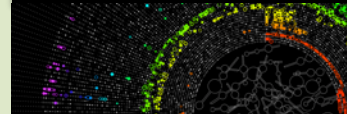
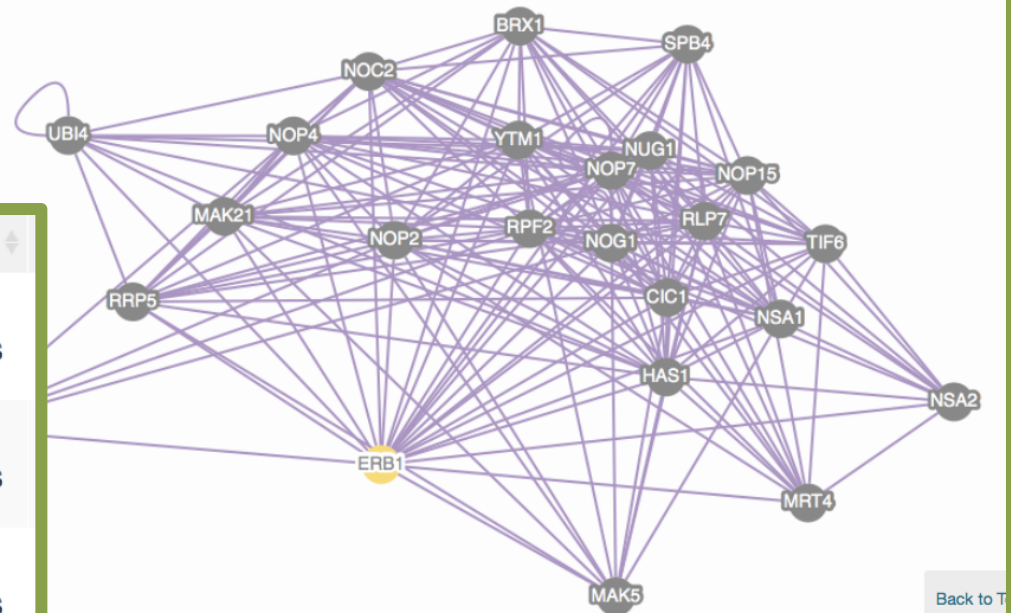
Physical Interactions

- Affinity Capture-MS: 220
- Affinity Capture-RNA: 6
- Affinity Capture-Western: 10
- Biochemical Activity: 4
- Co-purification: 2
- Reconstituted Co
- Two-hybrid: 2

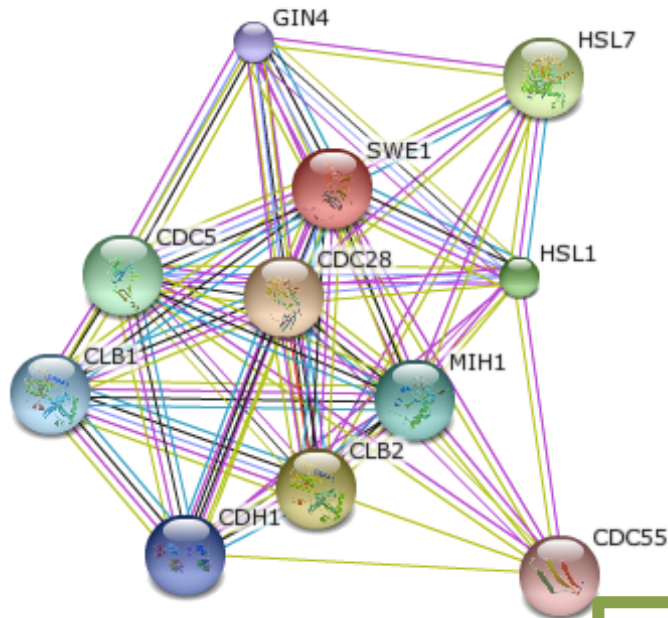
Interactor ▲	Type ▲	Assay ▲
ACO1	Physical	Affinity Capture-MS
ARD1	Physical	Affinity Capture-MS
ASK10	Physical	Affinity Capture-MS
i ATG1	Physical	Biochemical Activity
i BRE5	Physical	Affinity Capture-MS

Interaction Network i

Reset



Overlaying curated data

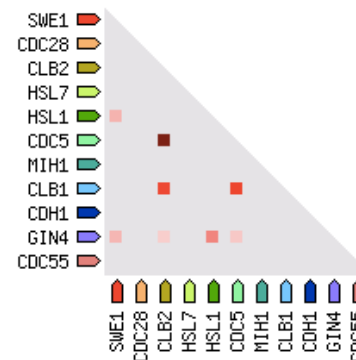


<http://string-db.org>

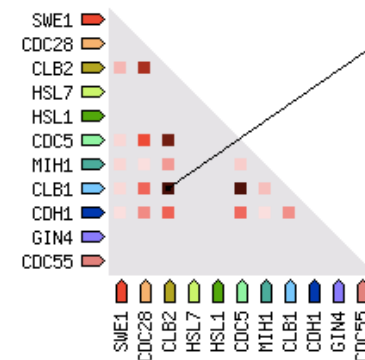
- Coexpression -

association score 0.0 1.0

... from Coexpression in
Saccharomyces cerevisiae:



... from Coexpression in
other species (transferred):



from *D. rerio*
from *A. thaliana*
from *M. musculus*
from *R. norvegicus*
from *B. taurus*
from *D. melanogaster*
from *S. scrofa*
from *M. mulatta*
from *C. elegans*
from *G. gallus*

Show

Complexes, interactions, and more

Complex SGD_GO:0005955
calcineurin complex subunits



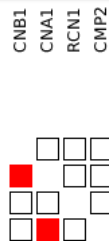
Crosstalk (Help)

SGD GO:0008287
Krogan 57



Interactions

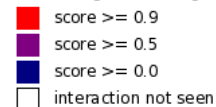
CNB1
CNA1
RCN1
CMP2



Upper diagonal - Socio-affinities

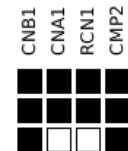


Lower diagonal - String experimental interactions



Subcellular Locations (Help)

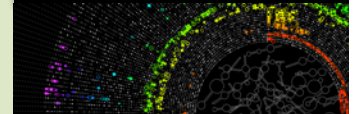
calcineurin complex
protein serine/threonine phosphatase complex
cytoplasm



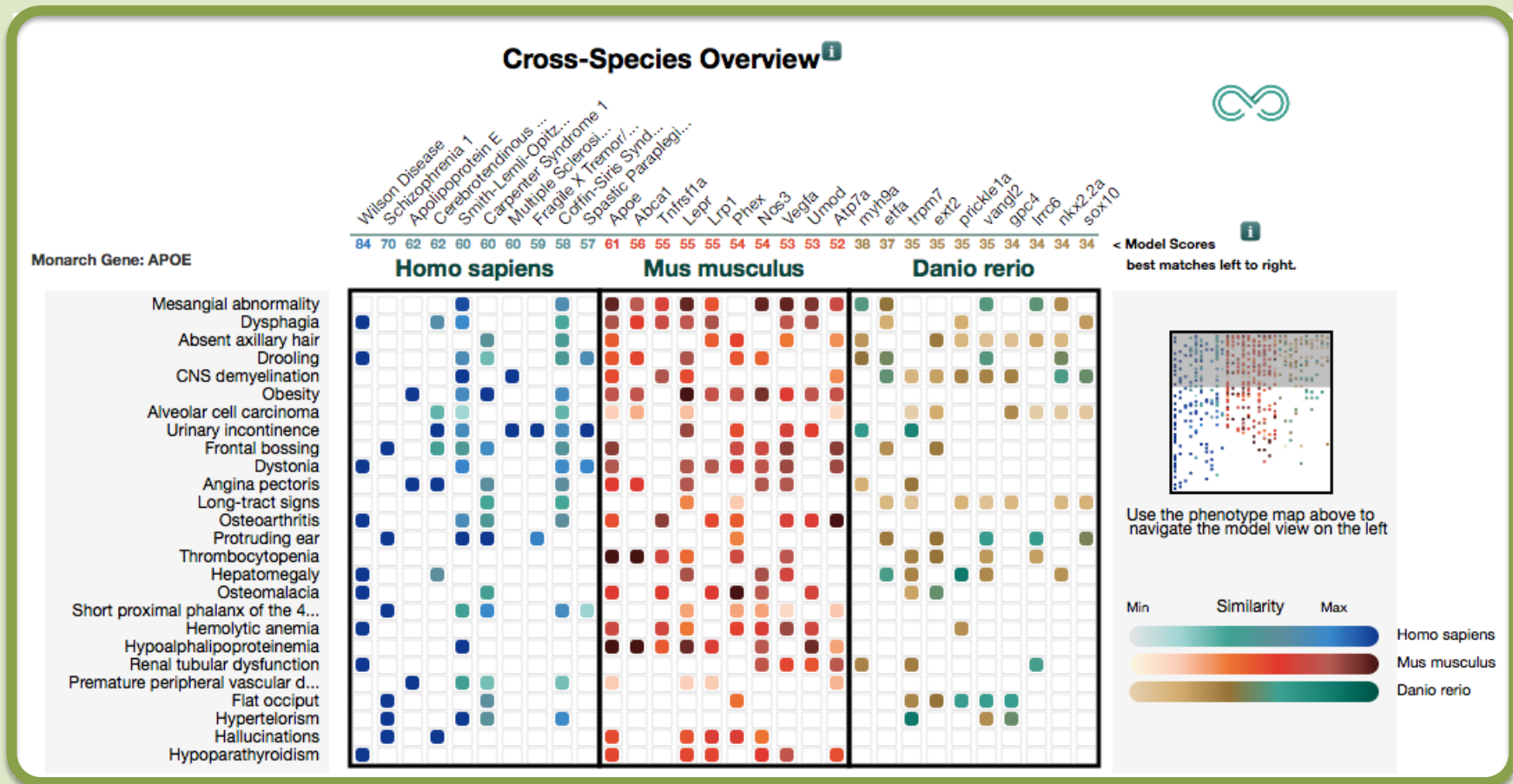
Copy Numbers (Help)



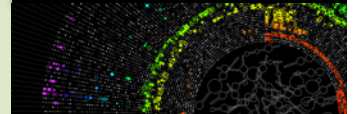
<http://3drepertoire.russelllab.org/>



Phenotypes and diseases

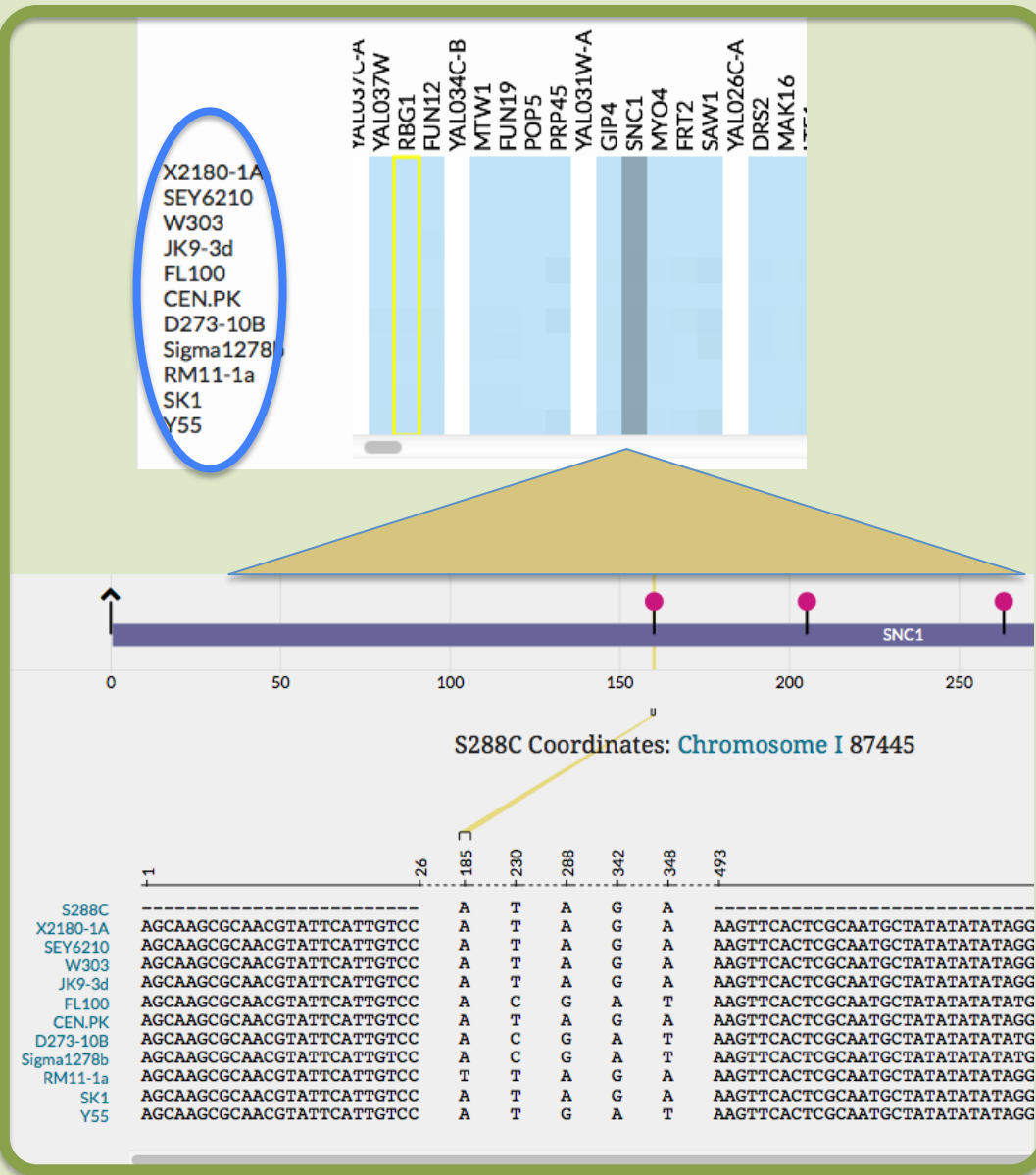


<http://monarchinitiative.org/>



Sequence variations

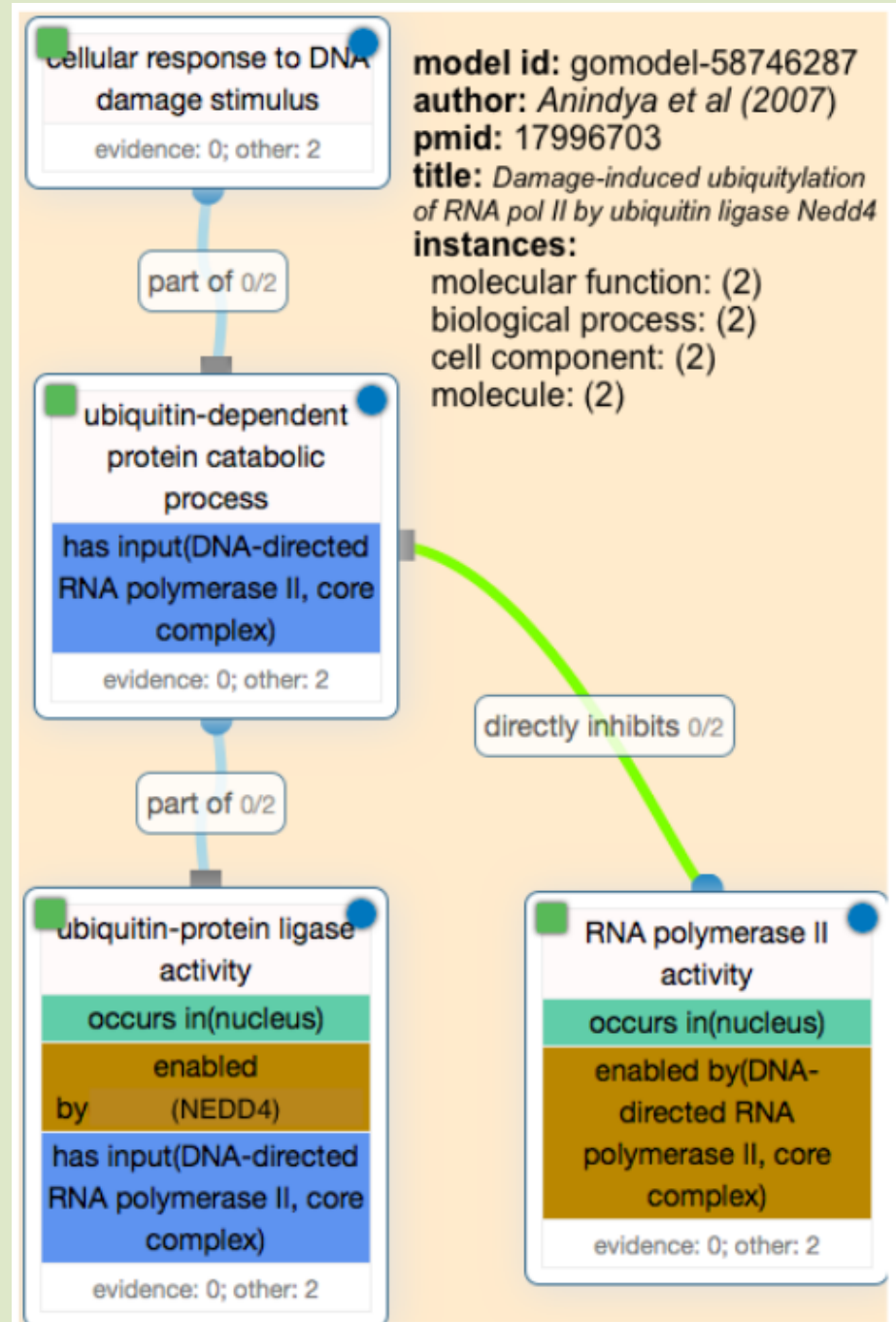
Sequence variation in various strains of *S. cerevisiae*



Molecular Model Editing Environment

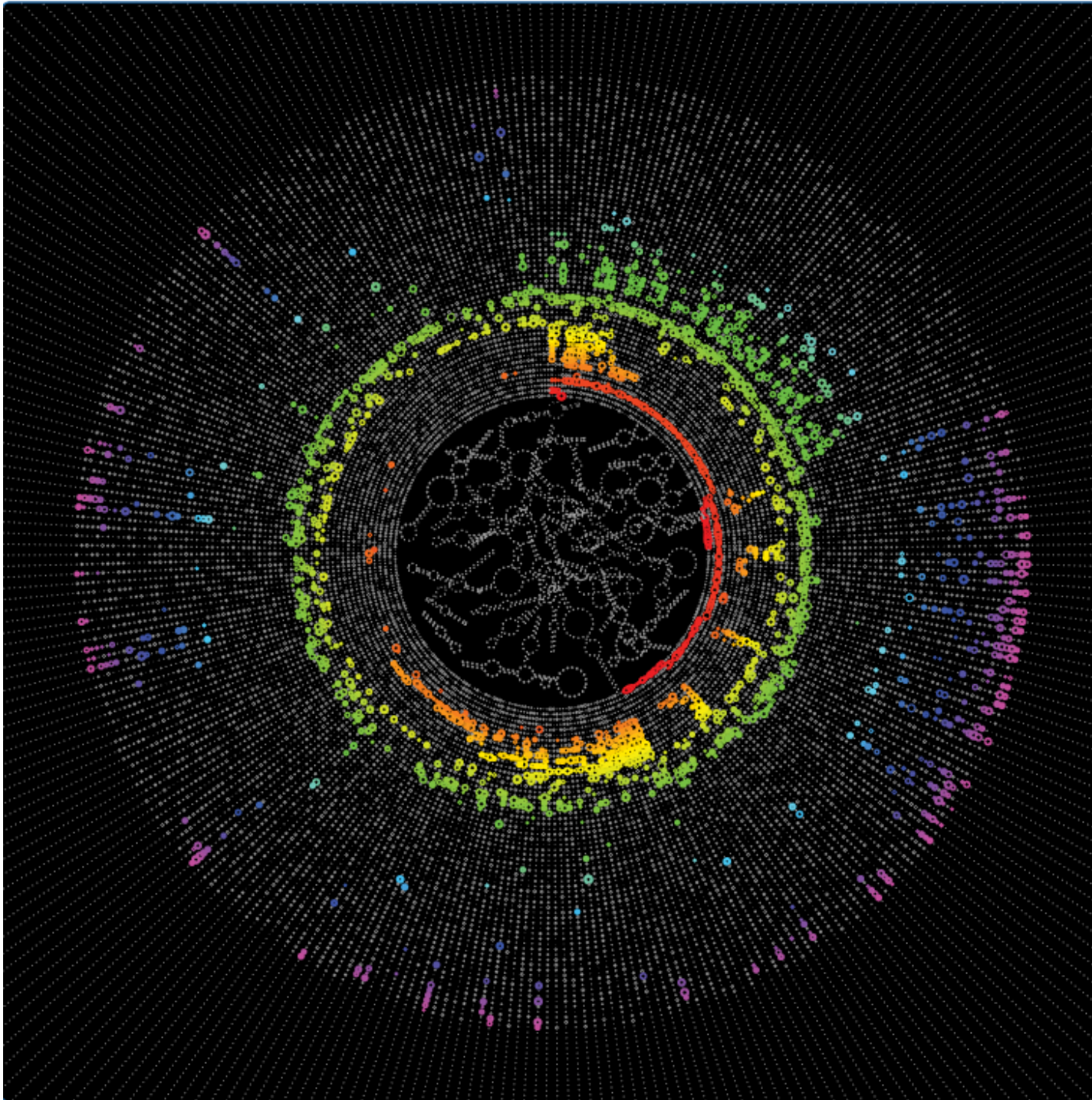
Noctua – prototype from GOC

- Each node (box) is a function or process.
- Other nodes are folded in as OWL expressions.
- Users may add and drag elements
- Supports real time collaboration



Understanding the Data

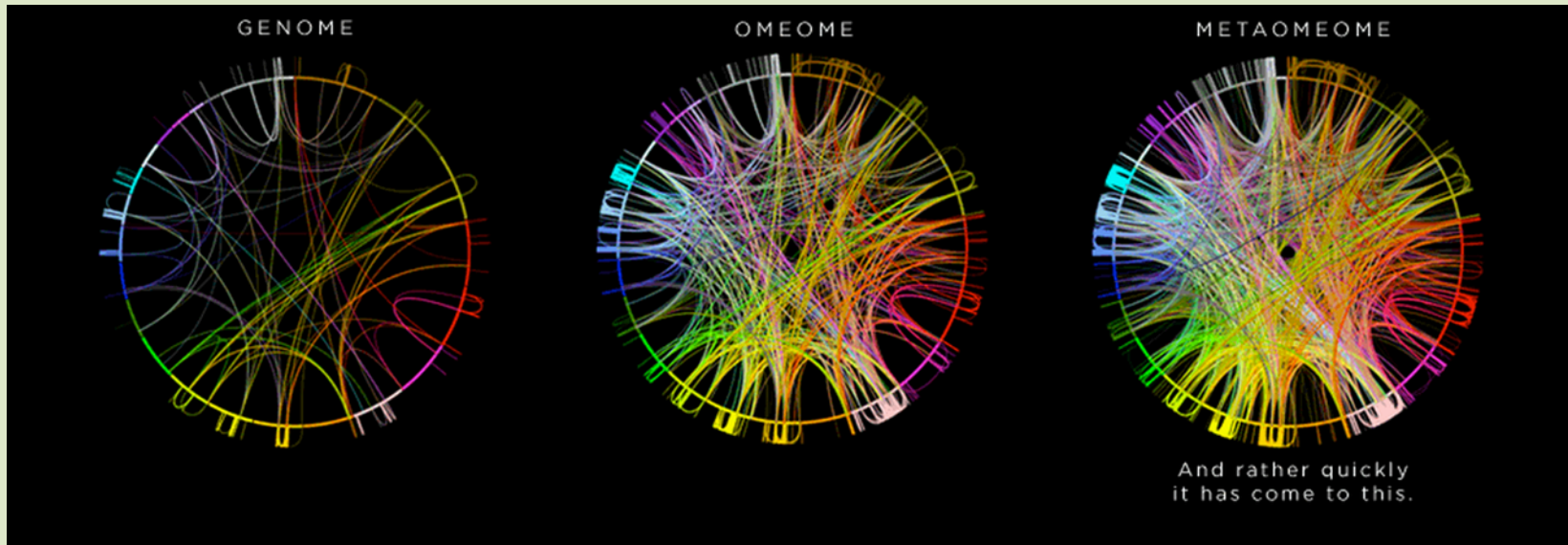
Much of the interpretation requires human judgment. Visualization improves our understanding and increases our chances of extracting meaningful conclusions.



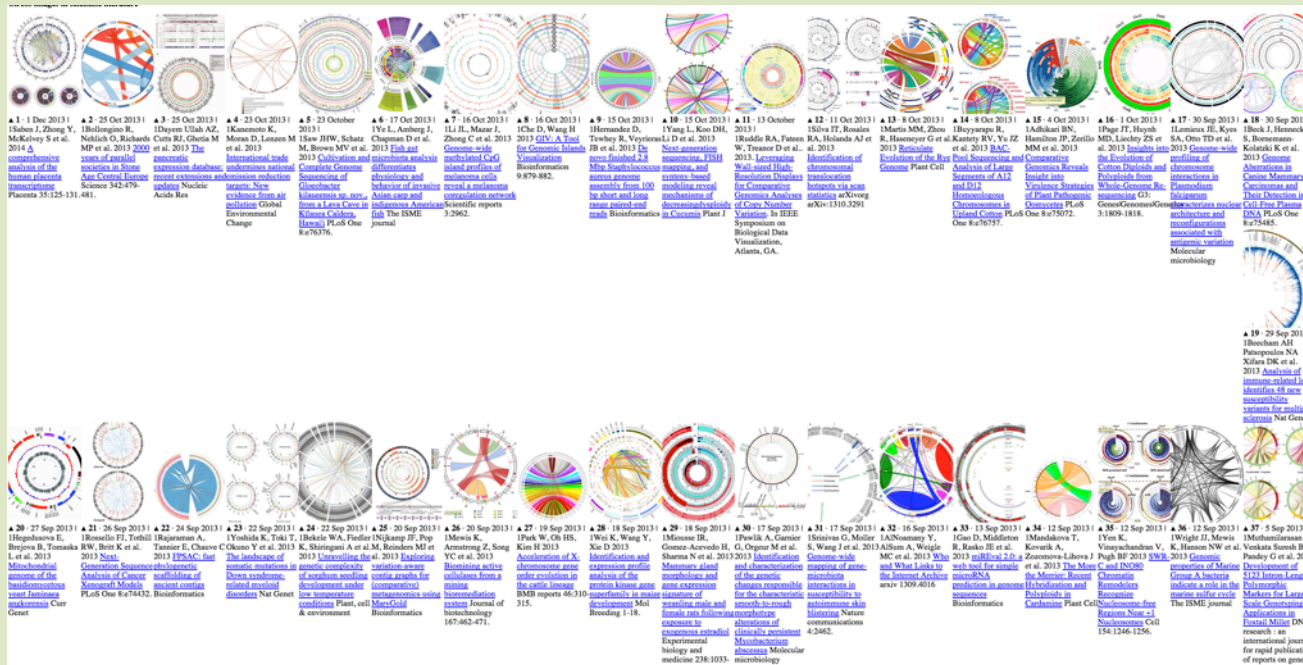
Cancer
miRNome
revealed survival
differences in
diffuse large B-
cell lymphoma
patients

Lim et al. *Genome Biol*
16:18 (2015)

Circos



ENCODE



Circular Genome Data Visualization

- Human placenta transcriptome
- Pancreatic expression db
- Wall-sized High-res display for comparative analys. of CNV
- Chromosomal translocations
- Variant identification in multiple sclerosis
- Sorghum seedling development under Low Temp conditions
- Etc., etc., etc...

READING THE G

The chart shows variety of bacter that cause huma

The x-axis represent burden, the axis worldwide death

The y-axis depict percentage of result in death.

Each colored line genome of the b that causes the d the genome and quinine or cycles (GC content) are viruses, the gene strains are show

THE GENOME A

The genome of e virus is represent line or path. The is proportional to genome. At each path, color is use content near the



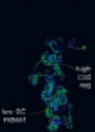
Individual genes i scale, but red reg indicate breaks b Path direction is a repeat and GC co location. Regions are straighter the

REPEAT



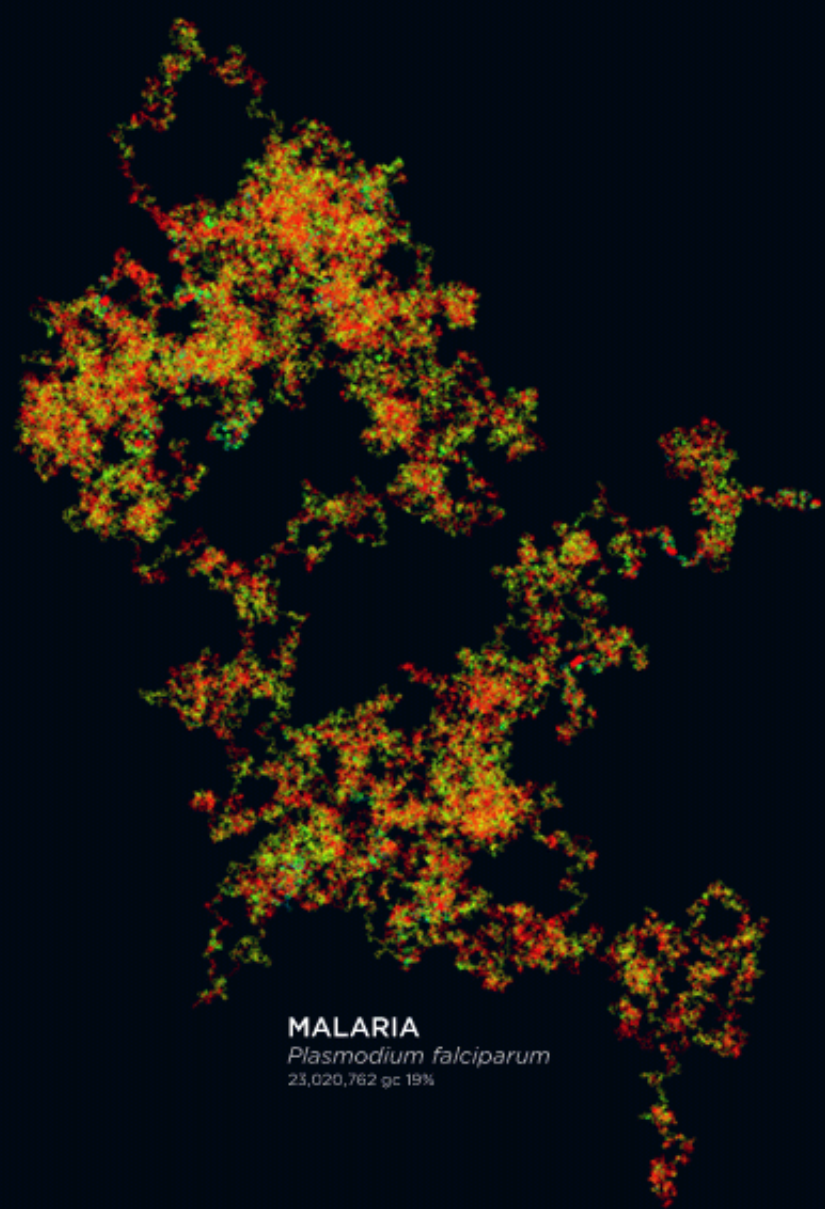
INTERPRETING E

Bacterial genome gene rich. The lo intergenic region can be seen by li the genome path



10 Data of human chr 2 near H200 on the cluster

Data Source
Genomes and morphology data is taken from H200, depending the distribution of genes and the data are largely limited to Africa



MALARIA
Plasmodium falciparum
23,020,762 gc 19%



SYPHILIS
Treponema pallidum
1,134,371 gc 53%

SARS
SARS coronavirus TOR2
29,751 gc 41%



BIOGRAPHY · FAST-FORWARD SLIDE



Martin Krzywinski
Scientist, Bioinformatics
[Genome Sciences Centre](#)
BC Cancer Agency
570 W 7th Avenue
Vancouver BC V5Z 4S6
Canada

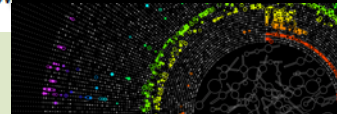
1.604.877.6000 x 673262 martink@bcgsc.ca @mkrzywinski

Visualizing sequencing data

Table 1 | Tools for visualizing sequencing data

Name	Cost	OS	Description	URL
Stand-alone tools				
ABYSS-Explorer ²⁵	Free	Win, Mac, Linux	Interactive assembly structure visualization tool	http://tinyurl.com/abyss-explorer/
CLC Genomics Workbench	\$	Win, Mac, Linux	Integrates NGS data visualization with analysis tools; user friendly	http://www.clcbio.com/
Consed ^{3*}	Free	Mac, Linux	Widely used; assembly finishing package; NGS compatible	http://www.phrap.org/
DNASTAR Lasergene ¹⁴	\$	Win, Mac	Analysis suite with an assembly finishing package; NGS compatible	http://www.dnastar.com/
EagleView ¹⁷	Free	Win, Mac, Linux	Assembly viewer; compatible with single-end NGS	http://tinyurl.com/eagleview/
Gap ^{12,13}	Free	Linux	Widely used; assembly finishing package; Gap5 is NGS compatible	http://staden.sourceforge.net/
Hawkeye ⁶	Free	Win, Mac, Linux (S)	Sanger sequencing assembly viewer	http://amos.sourceforge.net/hawkeye/
Integrative Genomics Viewer (IGV)*	Free	Win, Mac, Linux	Genome browser with alignment view support (Table 2); NGS compatible	http://www.broadinstitute.org/igv/
MapView ¹⁸	Free	Win, Linux	Read alignment viewer; custom file format for fast NGS data loading	http://evolution.sysu.edu.cn/mapview/
MaqView	Free	Mac, Linux	Read alignment viewer; fast NGS data loading from Maq alignment files	http://maq.sourceforge.net/
Orchid	Free	Linux (S)	Assembly viewer customized to display paired-end relationships	http://tinyurl.com/orchid-view/
Sequencher	\$	Win, Mac	Assembly finishing package	http://www.genecodes.com/
SAMtools tview ⁸	Free	Win, Mac, Linux	Simple and fast text alignment viewer; NGS compatible	http://samtools.sourceforge.net/
Web-based tools				
LookSeq ¹⁹	Free		Uses AJAX; y axis for insert size; user configures data resources; NGS compatible	http://lookseq.sourceforge.net/
NCBI Assembly Archive Viewer ⁷	Free		Graphical interface to contig and trace data in NCBI's Assembly Archive	http://tinyurl.com/assmbrowser/

Free means the tool is free for academic use; \$ means there is a cost. OS, operating system: Win, Microsoft Windows; Mac, Macintosh OS X. Tools running on Linux usually also run on other Unix-like systems. (S) indicates that compilation from source is required. "Assembly finishing package" enables interactive sequence editing and/or integration with tools for automated assembly improvement. *Our recommendation

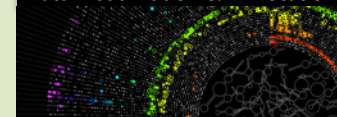


Genome Browsers

Table 2 | Genome browsers

Name	Description	URL
Stand-alone browsers		
Argo	Supports manual annotation of whole genomes	http://tinyurl.com/argo-combo
CGView ⁸²	Circular genome visualization	http://wishart.biology.ualberta.ca/cgview/
Gaggle ⁸³	Genome browser within an analysis framework; good microarray support	http://gaggle.systemsbiology.net/
Integrative Genomics Viewer (IGV)*	Flexible user interface; can integrate metadata as heat maps	http://www.broadinstitute.org/igv/
Integrated Genome Browser (IGB) ⁸⁴	GenoViz project genome browser; reusable visualization components	http://genoviz.sourceforge.net/
NCBI Genome Workbench	Displays sequence data in many views; integrated with BLAST	http://tinyurl.com/gbench/
Web-based browsers		
AnnoJ	Designed for NGS data; uses AJAX; assemble by html configuration	http://www.annoj.org/
Cancer Molecular Analysis Portal	Integrates clinical data; designed for TCGA project	https://cma.nci.nih.gov/cma-tcga/
Ensembl ^{31,32*}	Comprehensive genome browser and database; strong user support	http://www.ensembl.org/
GBrowse ^{28*}	GMOD ^{28*} component; back end of WormBase, FlyBase; v2.0 uses AJAX	http://gmod.org/wiki/Gbrowse
Genome Projector ⁴²	Offers circular and pathway views; user configures data resources	http://tinyurl.com/gprojector/
JBrowse ³⁹	Component of GMOD ^{28*} ; AJAX interface; user configures data resources	http://jbrowse.org/
JGI	Supports live annotation; primary portal for JGI genome projects	http://genome.jgi-psf.org/
NCBI Map Viewer ³³	Vertically oriented viewer; integrated with NCBI resources and tools	http://tinyurl.com/mapview1/
UCSC Genome Browser ^{30*}	Comprehensive genome browser and database; strong user support	http://tinyurl.com/ucscbrowser/
UCSC Cancer Genomics Browser ⁴³	Integrates clinical data; offers a pathway view; portal for TCGA data	http://genome-cancer.ucsc.edu/
UTGB	Toolkit to construct personalized browser; uses AJAX; user configures data resources	http://utgenome.org/
X:map ⁴¹	Customized to view Affymetrix exon arrays	http://xmap.picr.man.ac.uk/

All listed tools are free for academic use, and all are available for Microsoft Windows, Macintosh OS X and Linux. Tools running on Linux usually also run on other versions of Unix.



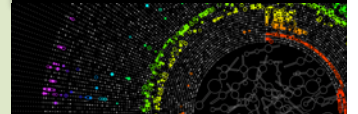
Comparative Genomics Visualization

Table 3 | Tools for comparative genomics visualization

Name	Description	Data	URL
Web-based tools			
Cinteny ⁶⁷	Three-scale view of synteny calculated from user-specified markers	H	http://cinteny.cchmc.org/
CoGe SynMap ⁸⁵	Dot plots from DAGChainer ⁶¹ alignments; histograms of synonymous substitutions	H	http://tinyurl.com/synmap/
GenomeMatcher ⁶³	A rich, mostly dot plot-based viewer displaying alignments and annotation	F,E,G	http://tinyurl.com/genomematcher/
MEDEA*	A Flash-based suite of linked-track, dot-plot and global-synteny viewing tools	C	http://tinyurl.com/broadmedea/
MultiPipMaker ⁸⁶	Vertically arranged display of user-supplied multiple alignments	F	http://pipmaker.bx.psu.edu/pipmaker/
PhIGs ⁶⁹	Ideogram-style interactive display of orthologs across >75 genomes	H	
UCSC Genome Browser ^{72*}	Conservation tracks within popular UCSC genome browser	H,F,G	http://genome.ucsc.edu/cgi-bin/hgGateway/
VISTA ^{87*}	Conservation tracks connected to a variety of analysis tools	H	http://genome.lbl.gov/vista/index.shtml
VSV, VISTA-Dot*	Three-scale viewer for synteny and dynamic, interactive dot plots for whole-genome DNA alignments	H	http://genome.jgi-psf.org/synteny/
Stand-alone tools			
ACT ⁷⁶	Linked-track views; annotation track search; stacking of multiple genomes	E,GF,D	http://www.sanger.ac.uk/Software/ACT/
Circos ⁷⁰	Circle-graph presentation of synteny; animations for increased dimensionality	C	http://mkweb.bcgsc.ca/circos
CMap ⁸⁸	Stacked vertical depictions of arbitrary relations among DNA segments	D,S	http://gmod.org/wiki/CMap
Combo ⁷⁷	Dot-plot and linked-track views; integration of annotation in both views	G,F,C	http://tinyurl.com/argo-combo
GBrowse_syn	GMOD ^{28*} component; highly customizable linked-track view of synteny	D,S	http://gmod.org/wiki/GBrowse_syn
MizBee ⁷¹	Synteny visualized using circular and linked-track views at multiple scales	C	http://mizbee.org/
Sybil ⁷⁸	Local and global views of synteny based on BlastP and protein clustering	D	http://sybil.sourceforge.net/
SynBrowse ⁷⁵	GMOD ²⁸ component; local synteny based on gene order, orthology or structure	D	http://www.synbrowse.org/
SynView ⁷⁹	GMOD ²⁸ component; synteny at different scales with multiple feature tracks	D	http://gmod.org/wiki/SynView

All tools listed are free and are either web-based or available for all three operating systems. The Data column describes the formats accepted for display within each tool: H, only alignment data hosted at the tool's website; F, FASTA format; E, EMBL/GenBank/DDJB format; G, gff format; C, a custom text-based format; D, designed for use with a user-hosted database; S, requires hosting from a user-supplied web server.

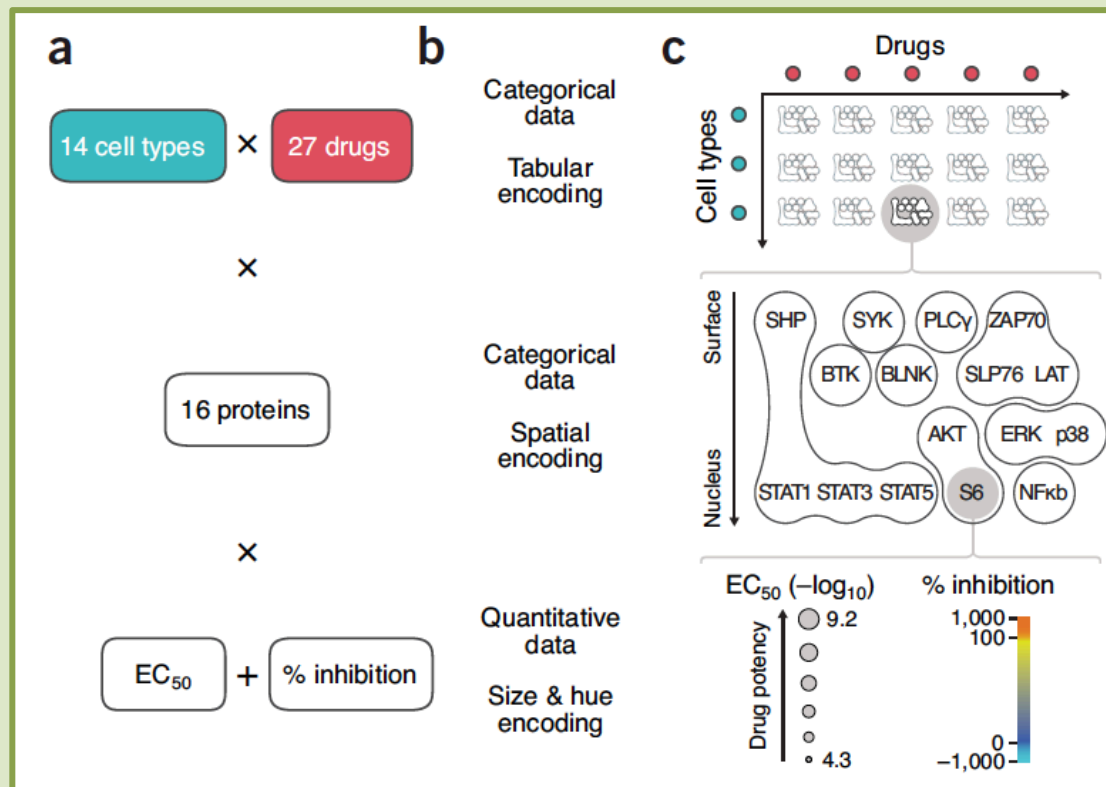
*Our recommendations



Communicating Complex Data

Focus on meaning instead of structure—anchor the figure to relevant biology rather than to methodological details.

- 1) What are the interesting findings, and what representation would communicate them clearly?



- 2) Forgo conventional approaches to displaying multidimensional data. Better to project the data onto familiar visual paradigms, such as a protein network or pathway, to saliently show biological effects in a functional context.

Krzywinski and Savig.
Nature Methods **10**:7, 595 (2013)

Storytelling

- Relate your data using the age-old custom of telling a story.
 - Stories have the capacity to delight and surprise and to spark creativity by making meaningful connections between data and the ideas, interests and lives of your readers.

Open-source: dissemination & contributions

- Genetic & genomic information is more valuable when shared
- Promote and encourage Open Science: transparency, reproducibility, data provenance. E.g. Open Bioinformatics Foundation <http://open-bio.org>
- Public repositories make software easily accessible and allow collaborative efforts, e.g. GitHub

<https://github.com/>



Our Panelists

1. Lorna Richardson:

eMouseAtlas and Image Informatics

2. Justyna Szostak:

Curated Causal Biological Network Models