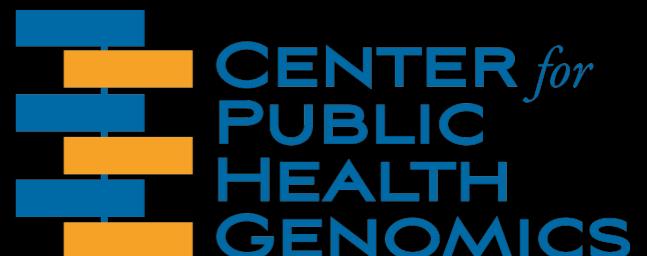


Mining the genome.

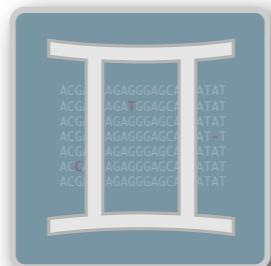
Aaron Quinlan
quinlanlab.org

University of Virginia, Charlottesville VA
Center for Public Health Genomics
Biochemistry and Molecular Genetics



Our research

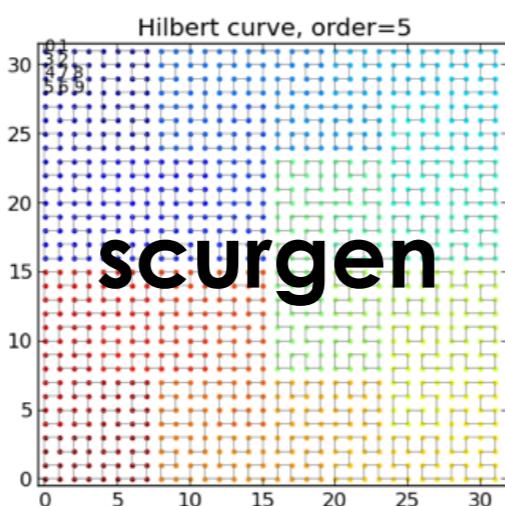
Software for genome
research



hydra

lumpy

cyvcf
pybedtools
grabix
bioawk



Our research

Software for genome research



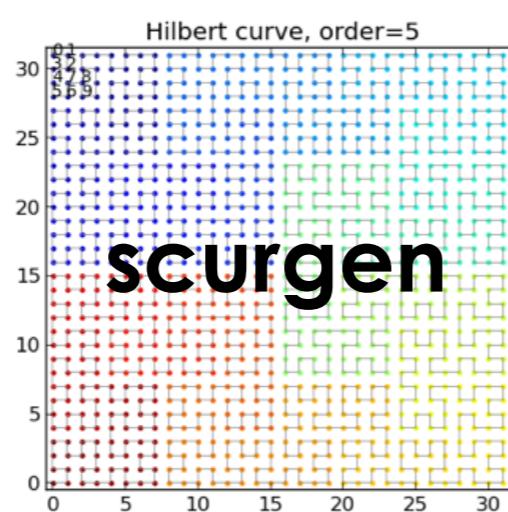
gemini



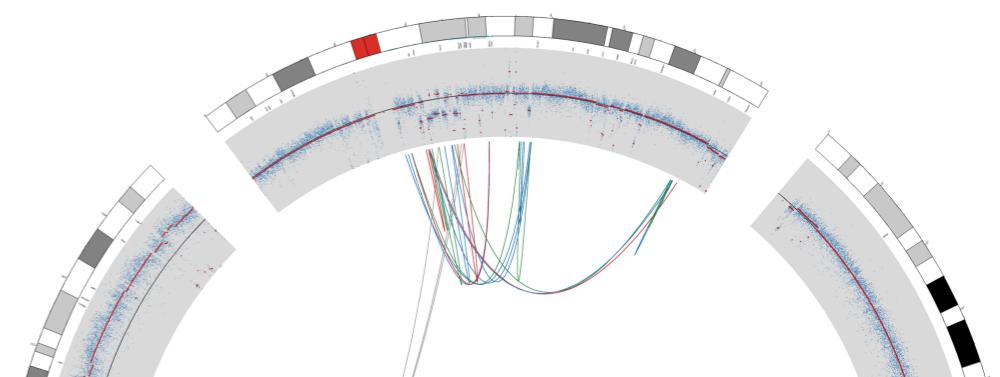
hydra

lumpy

cvcf
pybedtools
grabix
bioawk



Cancer genomics



1. Cellular and molecular origins of GBM
2. Initiating mutations underlying OV?
3. Tumor evolution. **Ira Hall's talk on Thurs.**

Human evolution

Understanding “unexplained” conservation

Disease genetics

Type 1 diabetes
SLE (Lupus)
Unexplained developmental disorders
Radiation hypersensitivity

Key contributors



Uma Paila, Ph.D.

Postdoctoral Research Associate

udp3f @ virginia.edu

Research Projects and Interests: Investigation of the genetic basis of extreme sensitivity to ionizing radiation; development of new analytical tools for exploring genetic variation identified through next-generation sequencing projects.



John Kubinski
Undergraduate
wunderkind



Neil Kindlon, M.S.

Staff Scientist and Software Engineer

nek3d @ virginia.edu

Research Projects and Interests: Software development for genomic analysis. Structural variation discovery and interpretation using DNA sequencing technologies.



Ryan Layer

Graduate student

rl6sf @ virginia.edu

Research Interests: Scalable algorithm development for high-throughput genomic analysis; genome data mining and analysis; structural variation discovery and interpretation.

Outline

1. Exploring genomic variation

- Challenges
- Approaches
- Exploring the genetic basis of radiation sensitivity

2. Break

3. Mining genomic datasets ➤ genome biology

- Algorithms and Tools
- Visualization
- Future

Gemini:

a flexible framework for mining genome variation



Uma Paila, Ph.D.

Postdoctoral Research Associate

udp3f @ virginia.edu

<https://github.com/arq5x/gemini>

Gemini overview

- **Goal:** unified framework for exploring genetic variation for disease and population genetics.
- Samples ➤ Genomes ➤ Variants ➤ VCF ➤ Gemini ➤ Test **H**
- Annotate variants with wealth of genome annotations
- Structured querying interface
- Extensible for new tool development
- Powerful, yet easy to use

Other options

PLINK/SEQ

A library for the analysis of genetic variation data

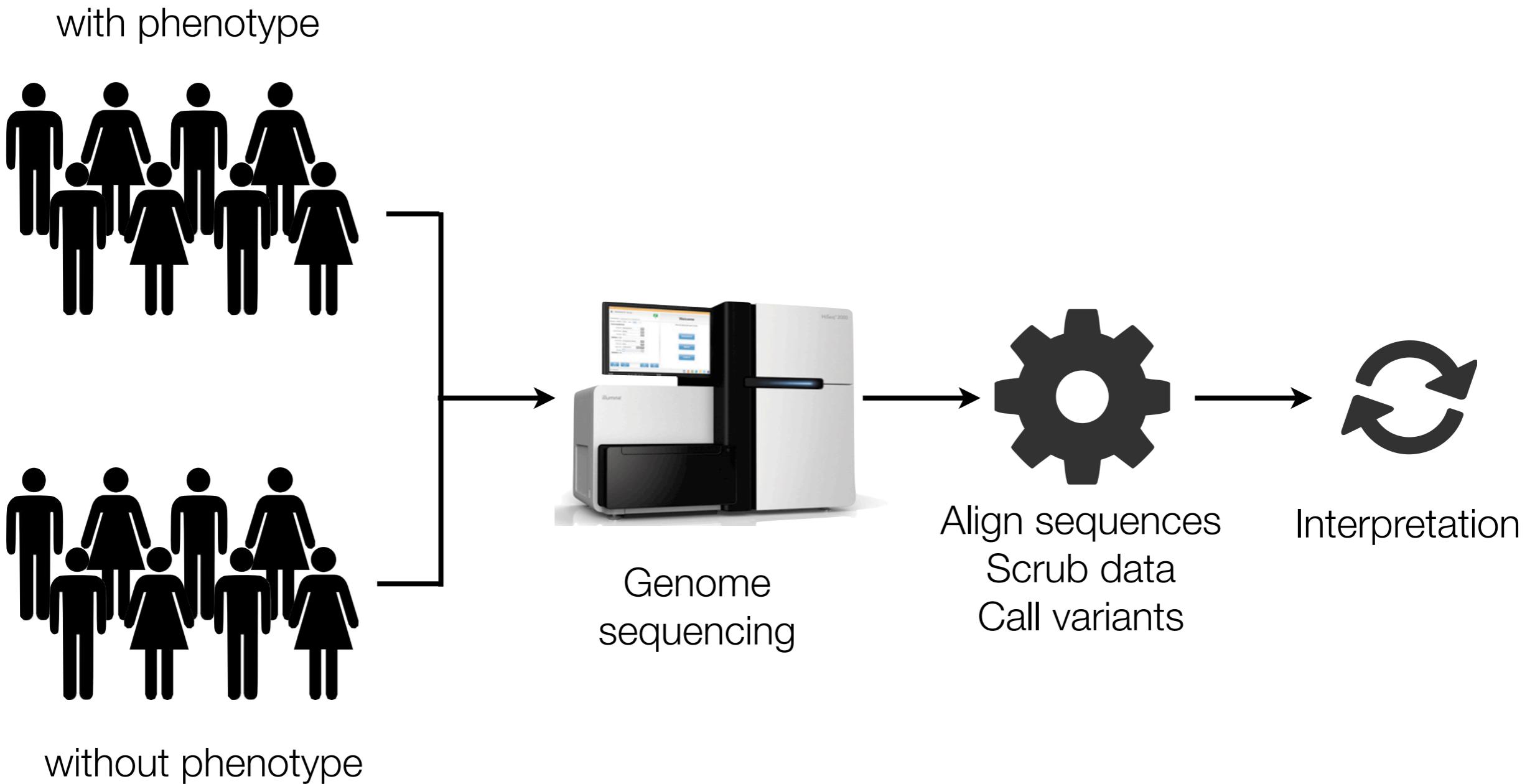
ANNOVAR

gSearch

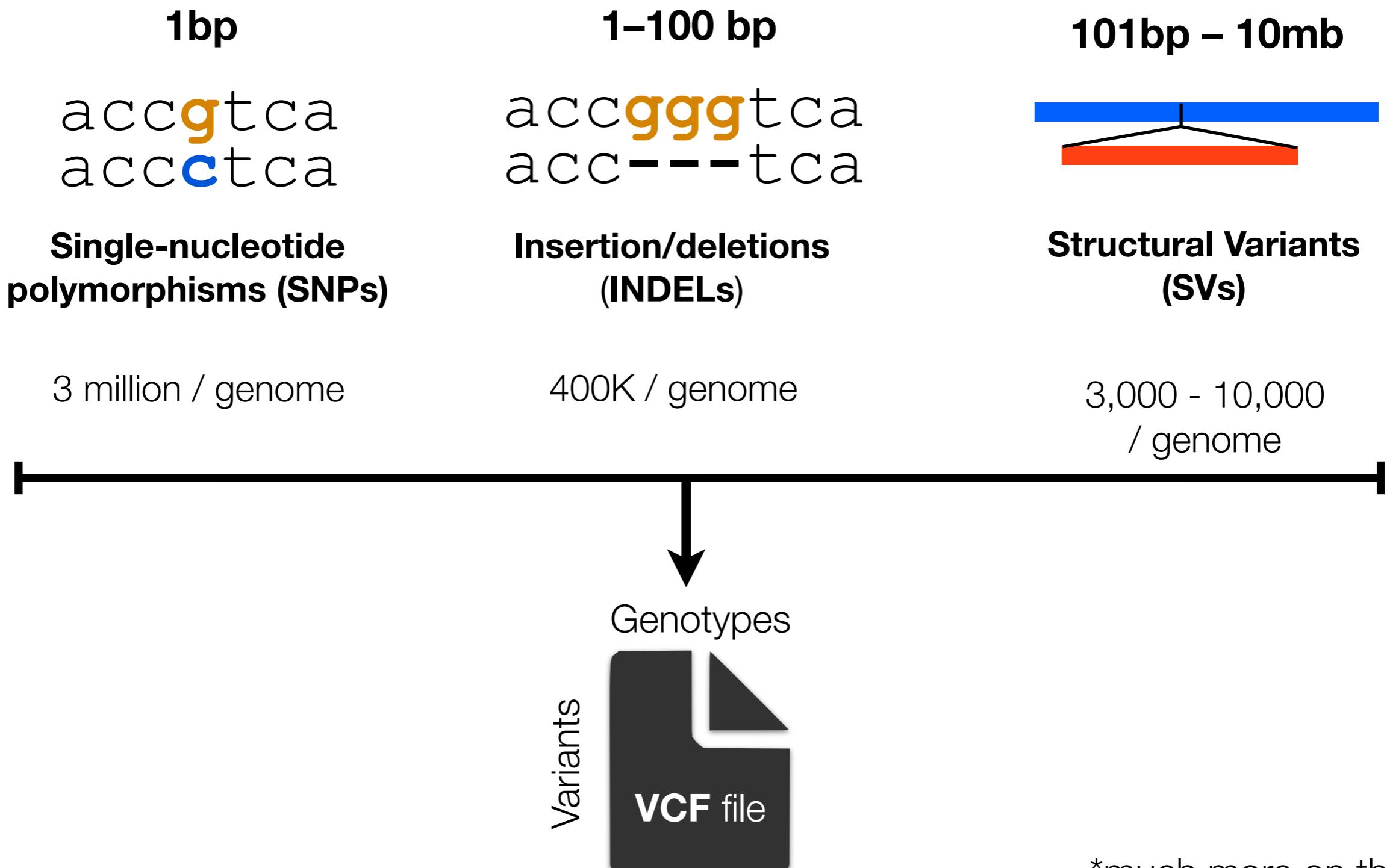
AnnTools

VAAST

A typical study

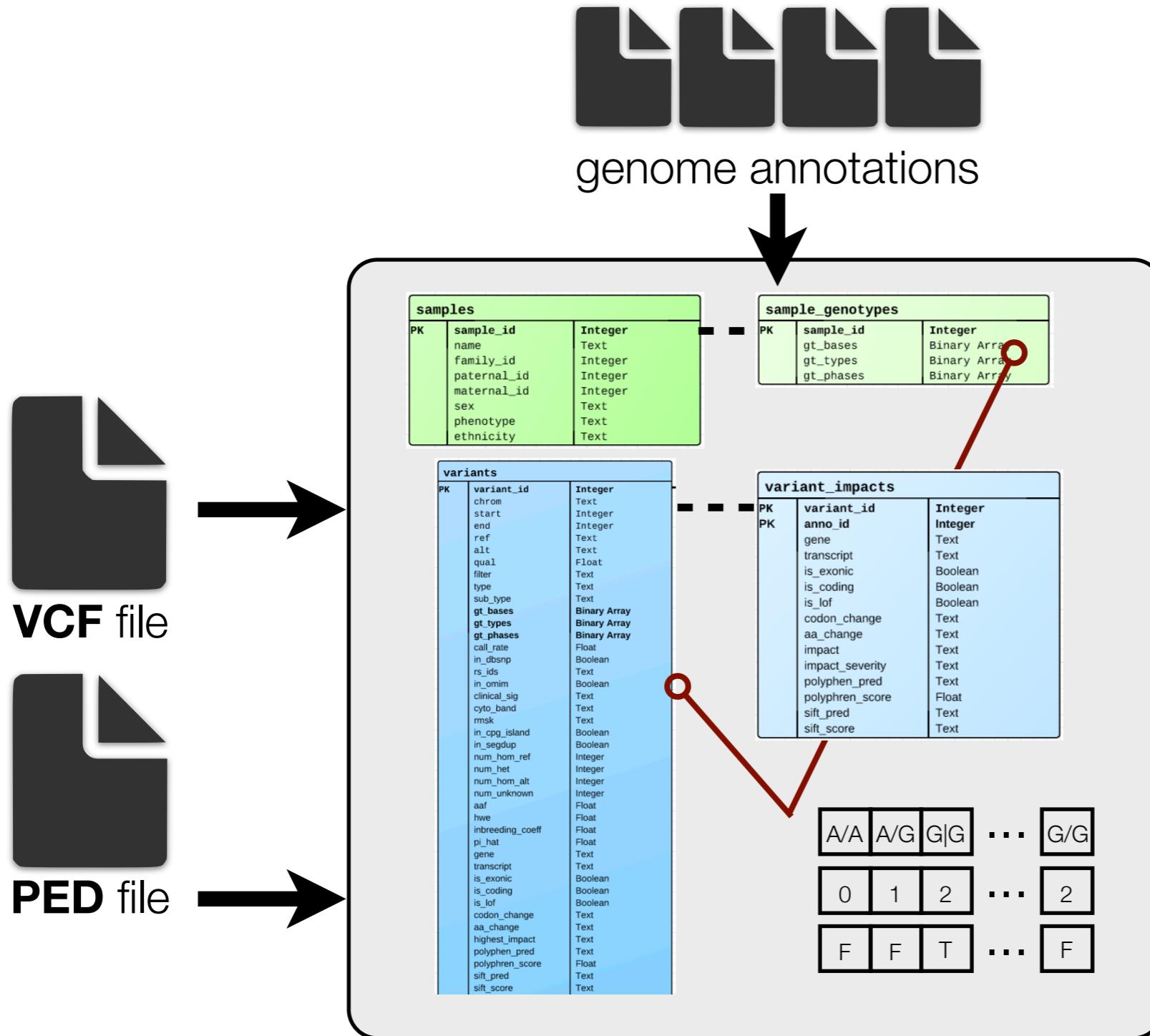


Variants for all samples are stored in a VCF file*



*much more on this later

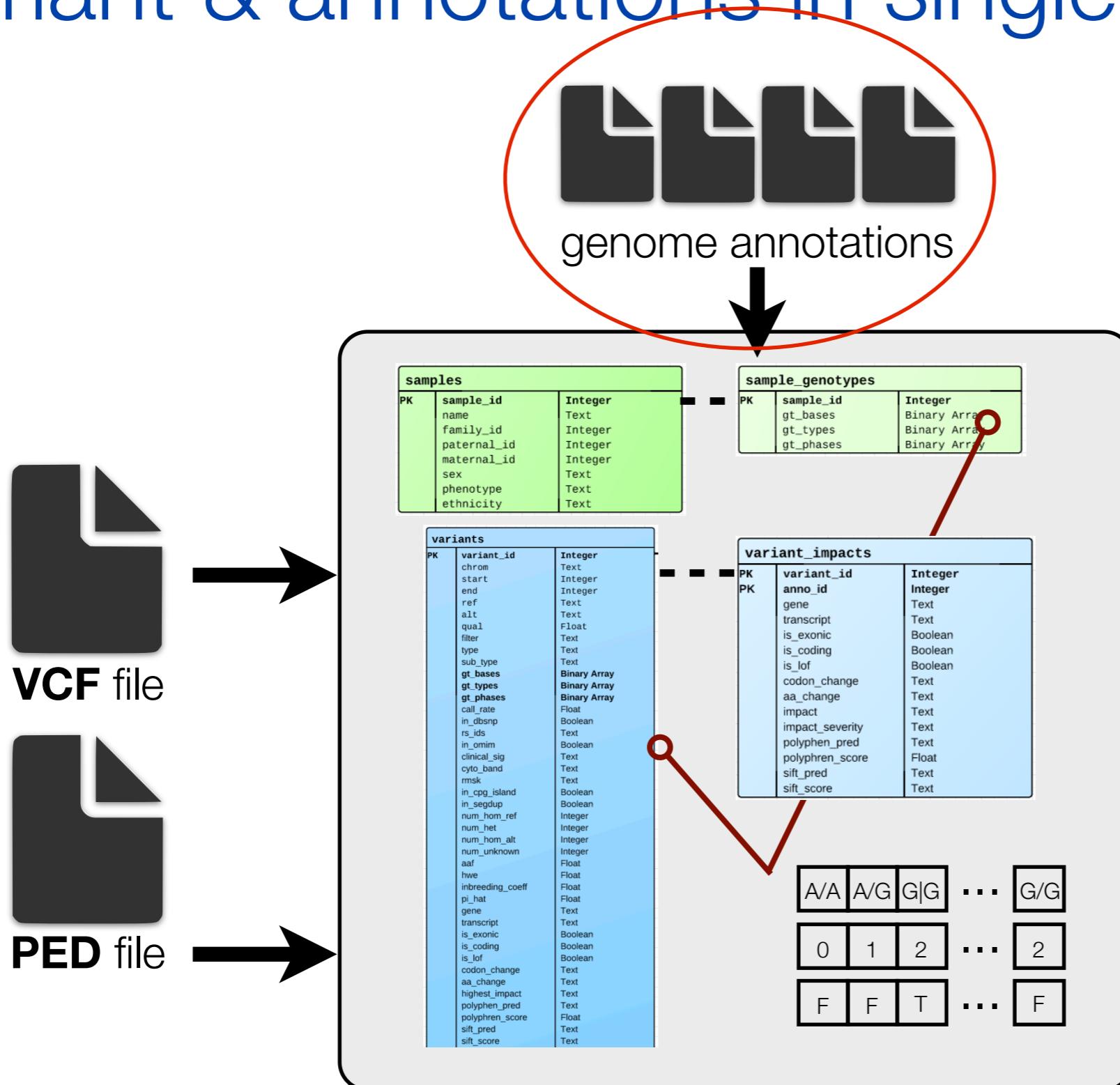
Variant & annotations in single framework



Compressed arrays of genotype information accommodates 1000s of genotypes per variant.

```
$ gemini load -v my.vcf -t VEP my.db
```

Variant & annotations in single framework



Compressed arrays of genotype information accommodates 1000s of genotypes per variant.

```
$ gemini load -v my.vcf -t VEP my.db
```

Genome annotations place variants in context



OMIM



COSMIC



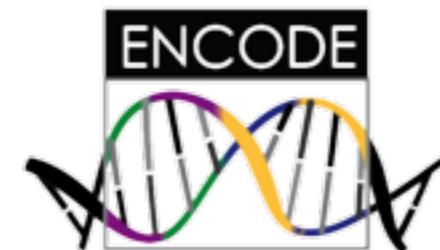
dbSNP



pathways



CpG
Repeats
SegDups
Conservation
(more)



DNase
FAIRE
TF Binding
chromatin states



Recombination



ESP



protein interactions

Variants, genotypes, & sample info all in a single database

samples		
PK	sample_id	Integer
	name	Text
	family_id	Integer
	paternal_id	Integer
	maternal_id	Integer
	sex	Text
	phenotype	Text
	ethnicity	Text

sample_genotypes		
PK	sample_id	Integer
	gt_bases	Binary Array
	gt_types	Binary Array
	gt_phases	Binary Array

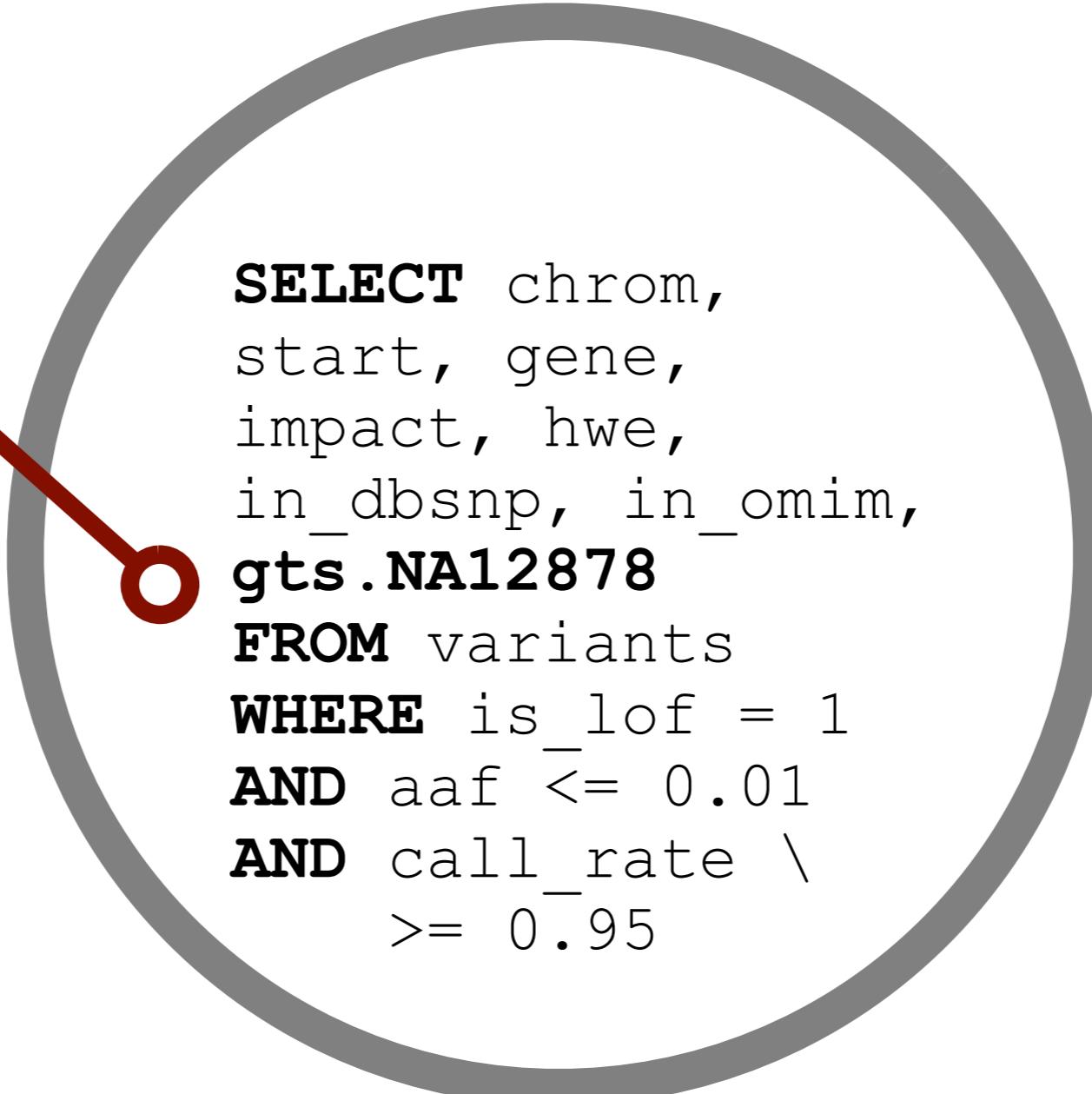
variants		
PK	variant_id	Integer
	chrom	Text
	start	Integer
	end	Integer
	ref	Text
	alt	Text
	qual	Float
	filter	Text
	type	Text
	sub_type	Text
	gt_bases	Binary Array
	gt_types	Binary Array
	gt_phases	Binary Array
	call_rate	Float
	in_dbsnp	Boolean
	rs_ids	Text
	in_omim	Boolean
	clinical_sig	Text
	cyto_band	Text
	rmsk	Text
	in_cpg_island	Boolean
	in_segdup	Boolean
	num_hem_rf	Integer

variant_impacts		
PK	variant_id	Integer
	anno_id	Integer
	gene	Text
	transcript	Text
	is_exonic	Boolean
	is_coding	Boolean
	is_lof	Boolean
	codon_change	Text
	aa_change	Text
	impact	Text
	impact_severity	Text
	polyphen_pred	Text
	polyphren_score	Float
	sift_pred	Text
	sift_score	Text

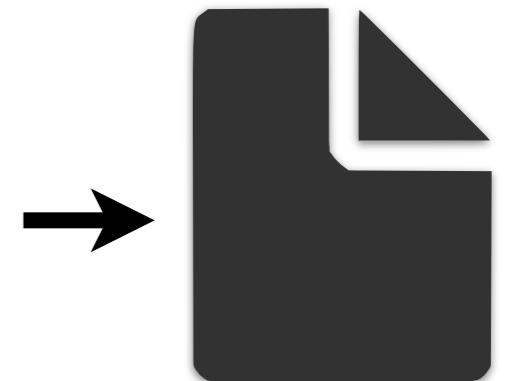
Ad hoc data exploration

Enhanced SQL engine
allows **selection** and
filtering based on
individual genotypes.

**Scales to 1000s of
samples.**



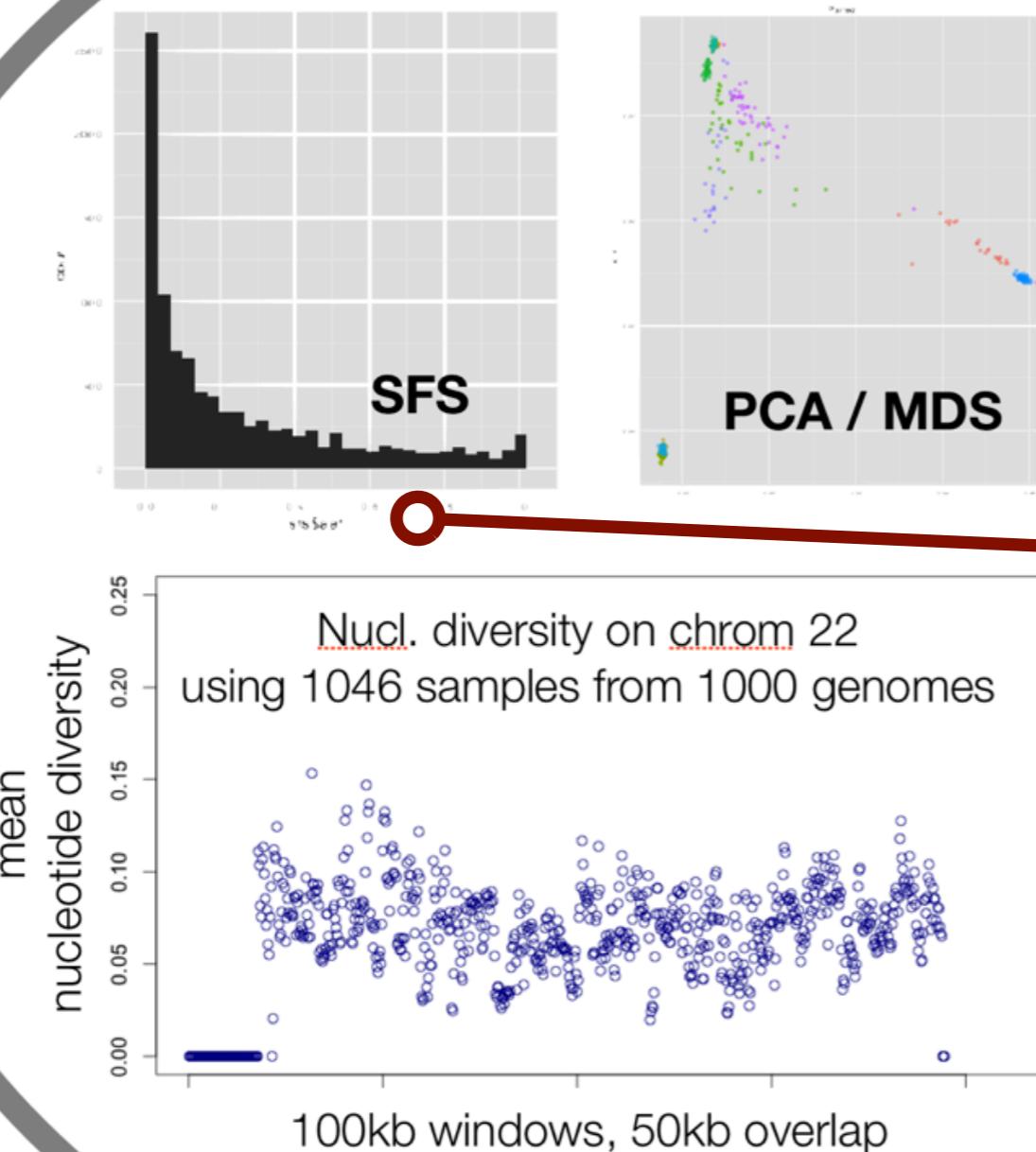
```
SELECT chrom,  
start, gene,  
impact, hwe,  
in_dbsnp, in_omim,  
gts.NA12878  
FROM variants  
WHERE is_lof = 1  
AND aaf <= 0.01  
AND call_rate \  
>= 0.95
```



Output files in
standard
formats (BED, etc.)

```
> gemini query -q [QUERY] my.db
```

Population genetics



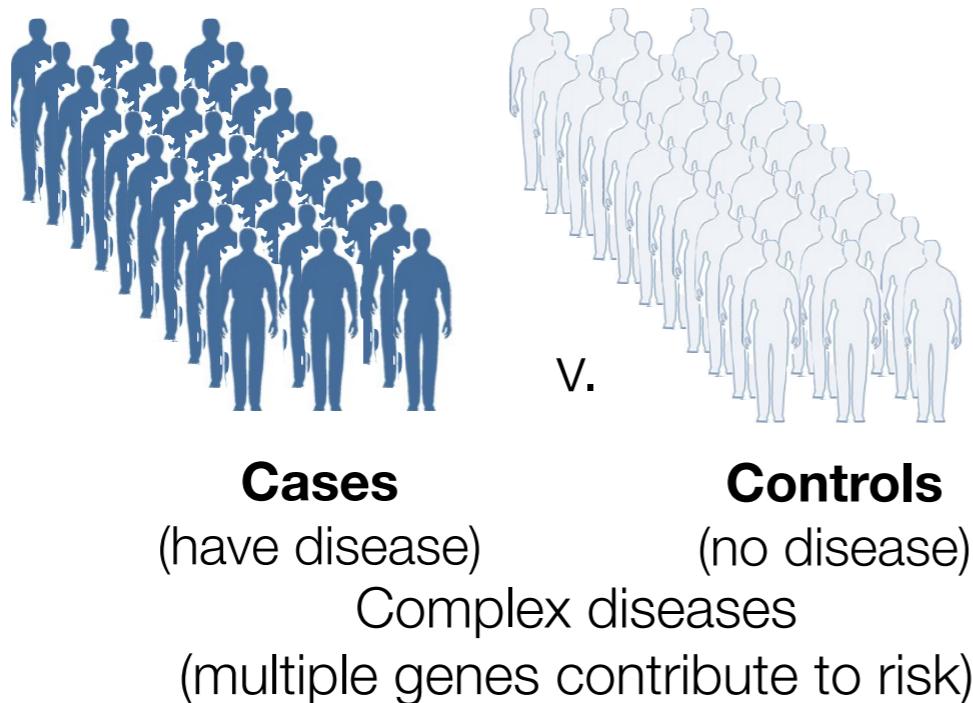
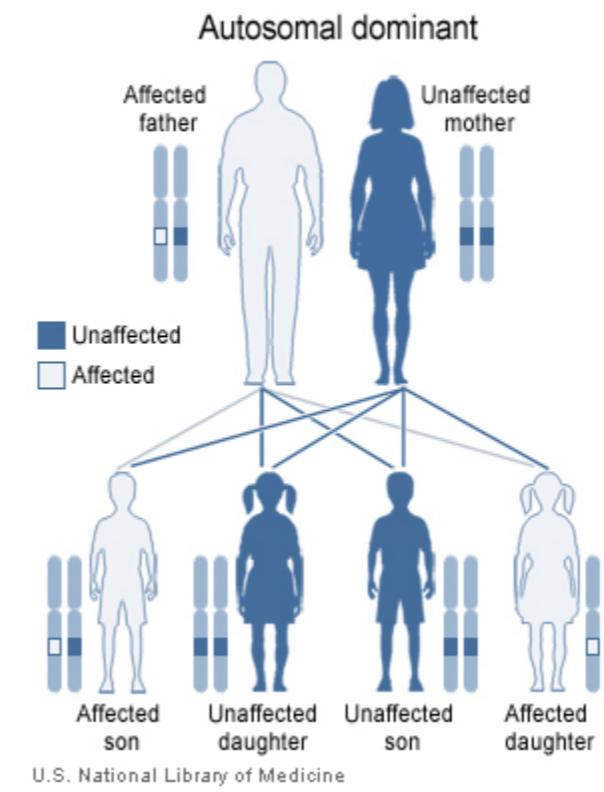
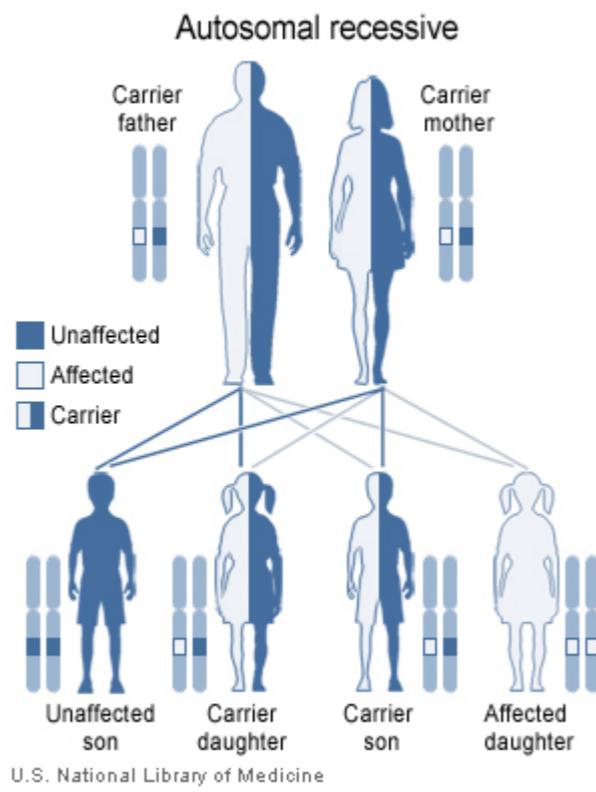
Multidimensional scaling,
Kinship coefficients
PCA

Variant QC and profiling

“windowing” tools

```
> gemini stats --sfs my.db
> gemini stats --mds my.db
> gemini stats --tstv my.db
> gemini stats --vars-by-sample my.db
> gemini stats --gts-by-sample my.db
```

Disease genetics



Tumor / Normal comparisons

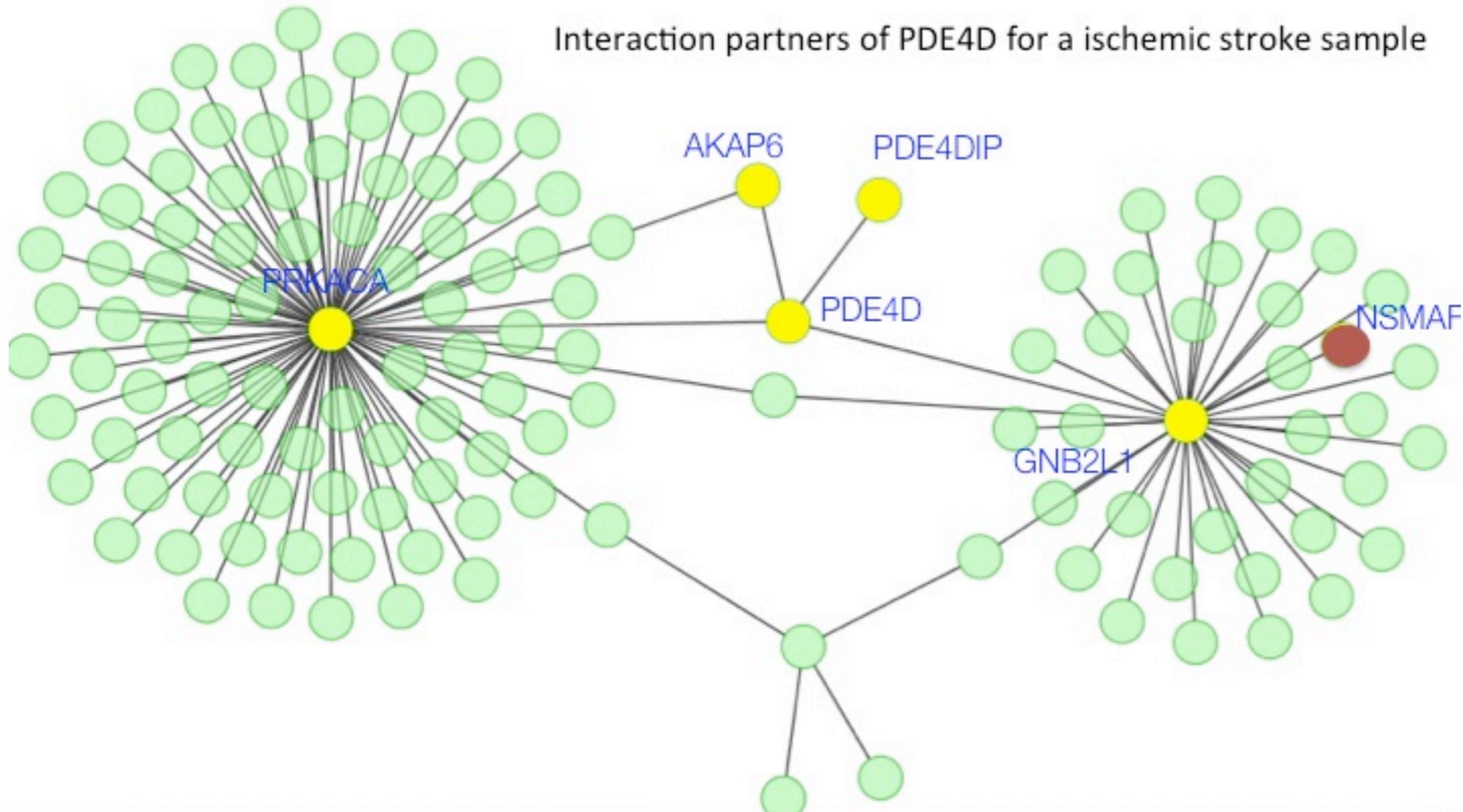
Trio sequencing
(pilot study at UVA
for undiagnosed
dev. disorders)

Pathway and interaction analysis

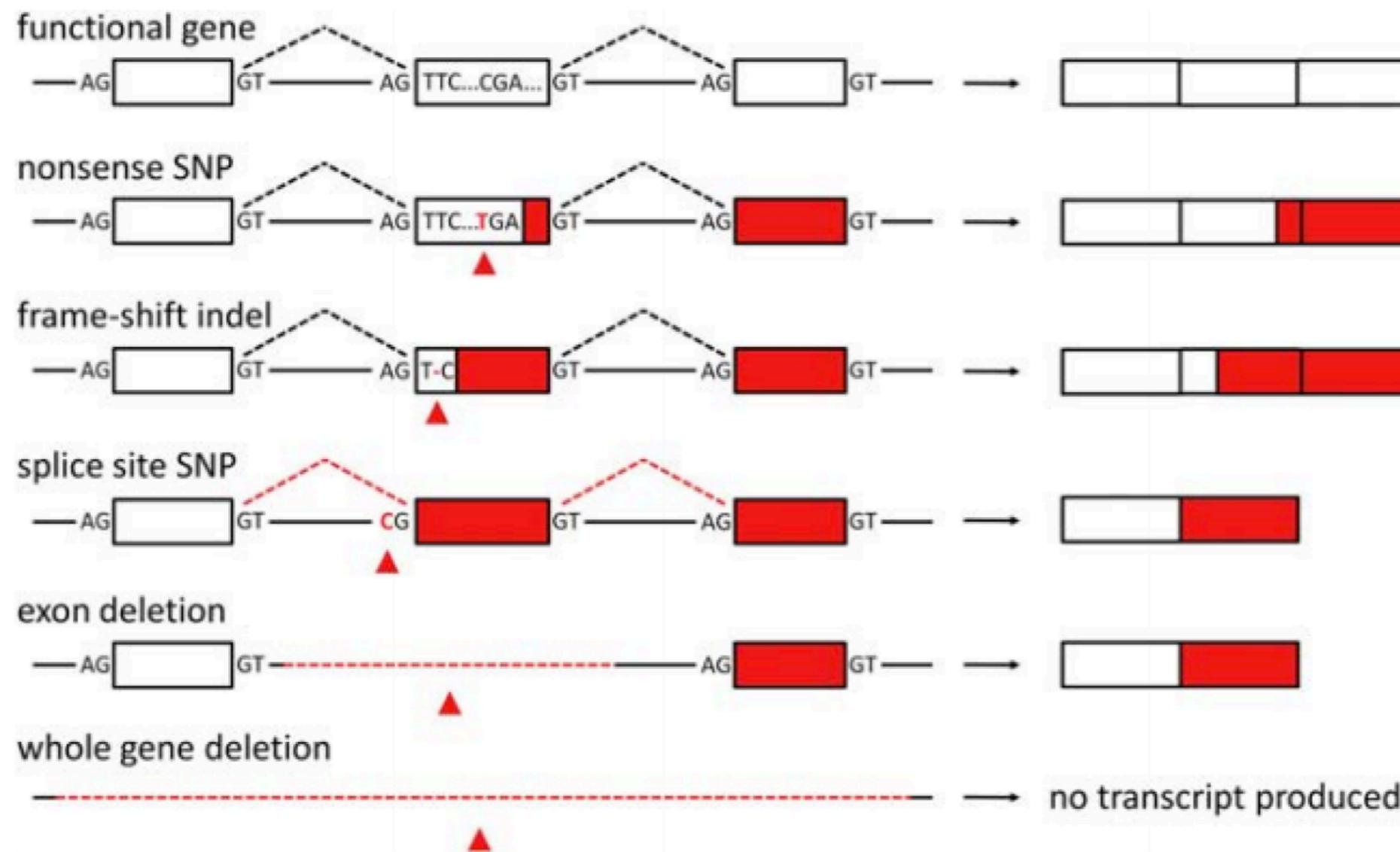
pathways



protein
interactions



Loss-of-function (LoF) variants



**Loss-of-function variants in the genomes
of healthy humans**

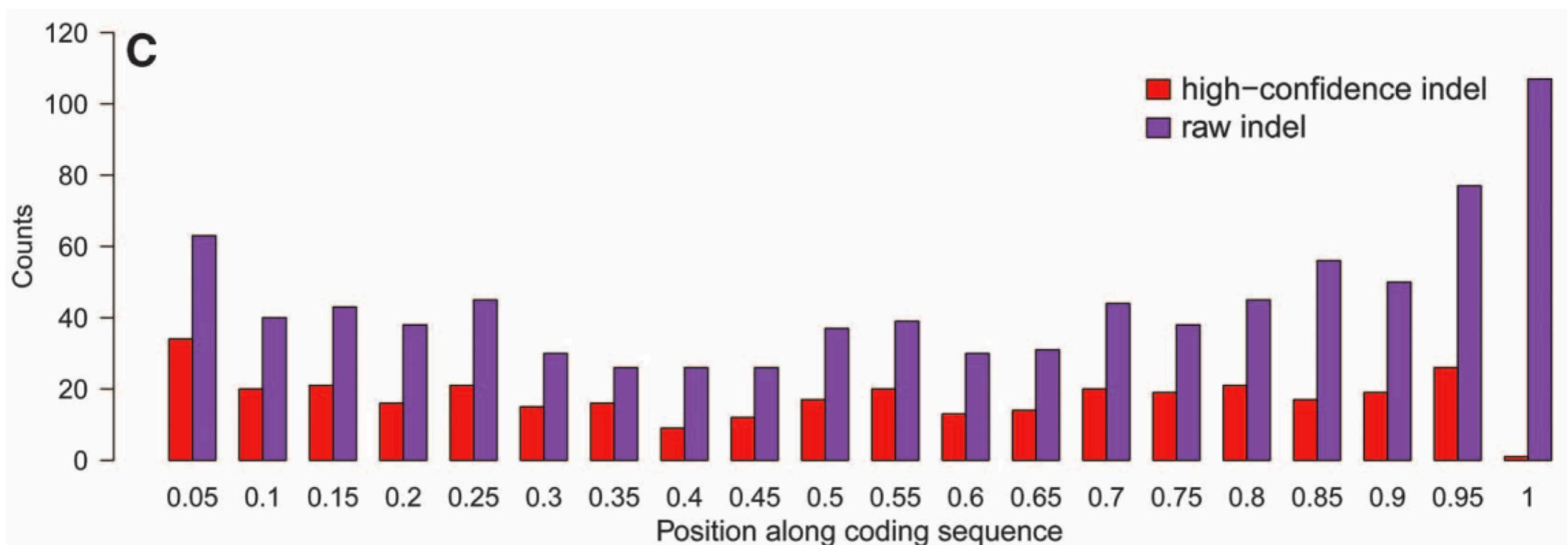
Daniel G. MacArthur* and Chris Tyler-Smith

Human Molecular Genetics, 2010, Vol. 19, Review Issue 2
doi:10.1093/hmg/ddq365
Advance Access published on August 30, 2010

R125–R130

Loss-of-function (LoF) variants

Not all LoF variants are created equal.



1. Consider the position of the LoF variant in the polypeptide.
2. Test whether a second frameshift restores reading frame.

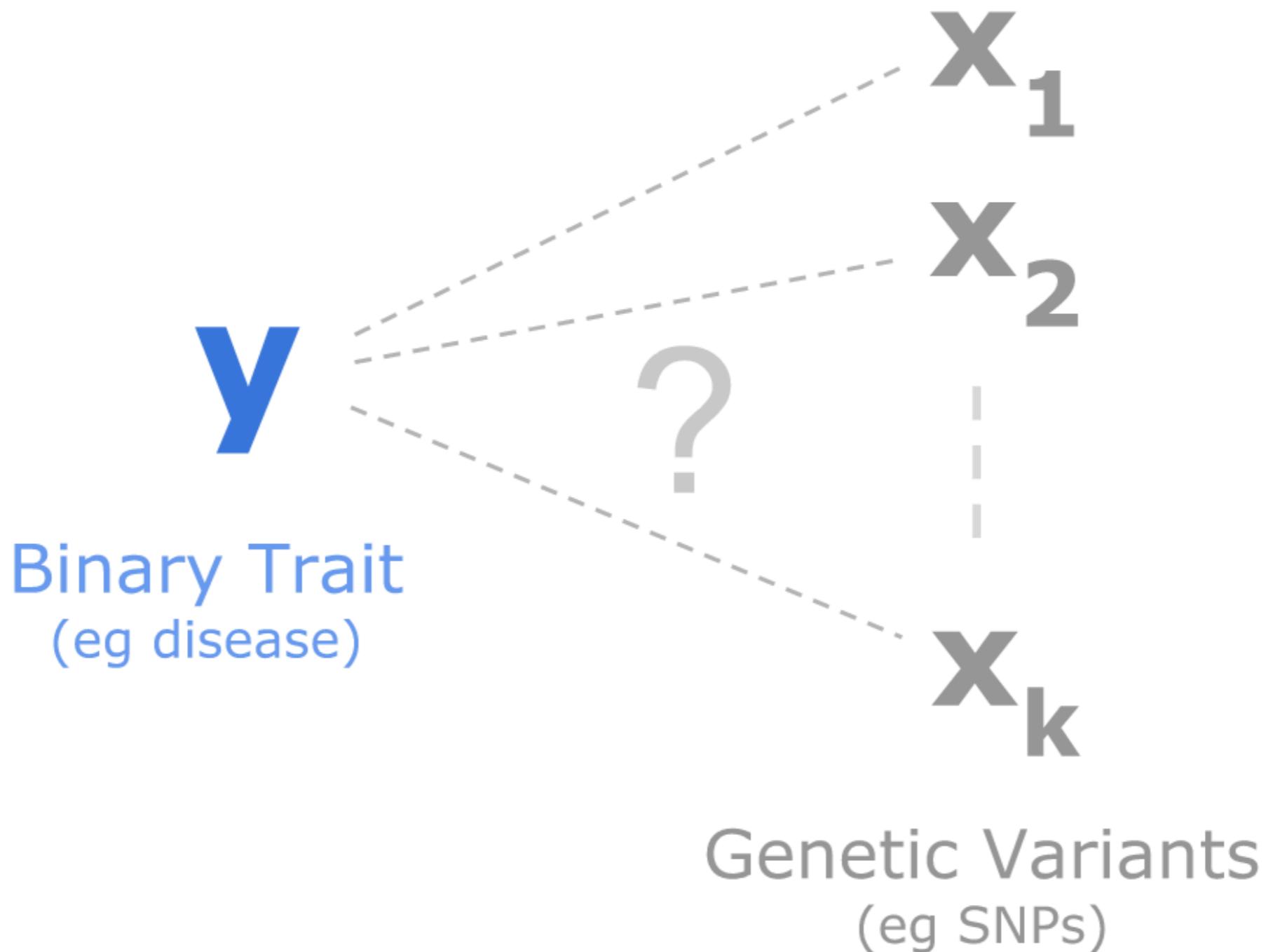
Burden tests for rare variants

- GWAS based on the notion that common variants underly common disease.
- There are many, many more rare variants than common
- Standard, GWAS-like single-locus tests will be underpowered unless there is a huge effect size.
- Burden tests and Collapsing methods.
 - C-alpha, SKAT, KBAT, Morris-Zeggini, Hotelling's T, ...
 - Active area of research. Gemini is a common framework.

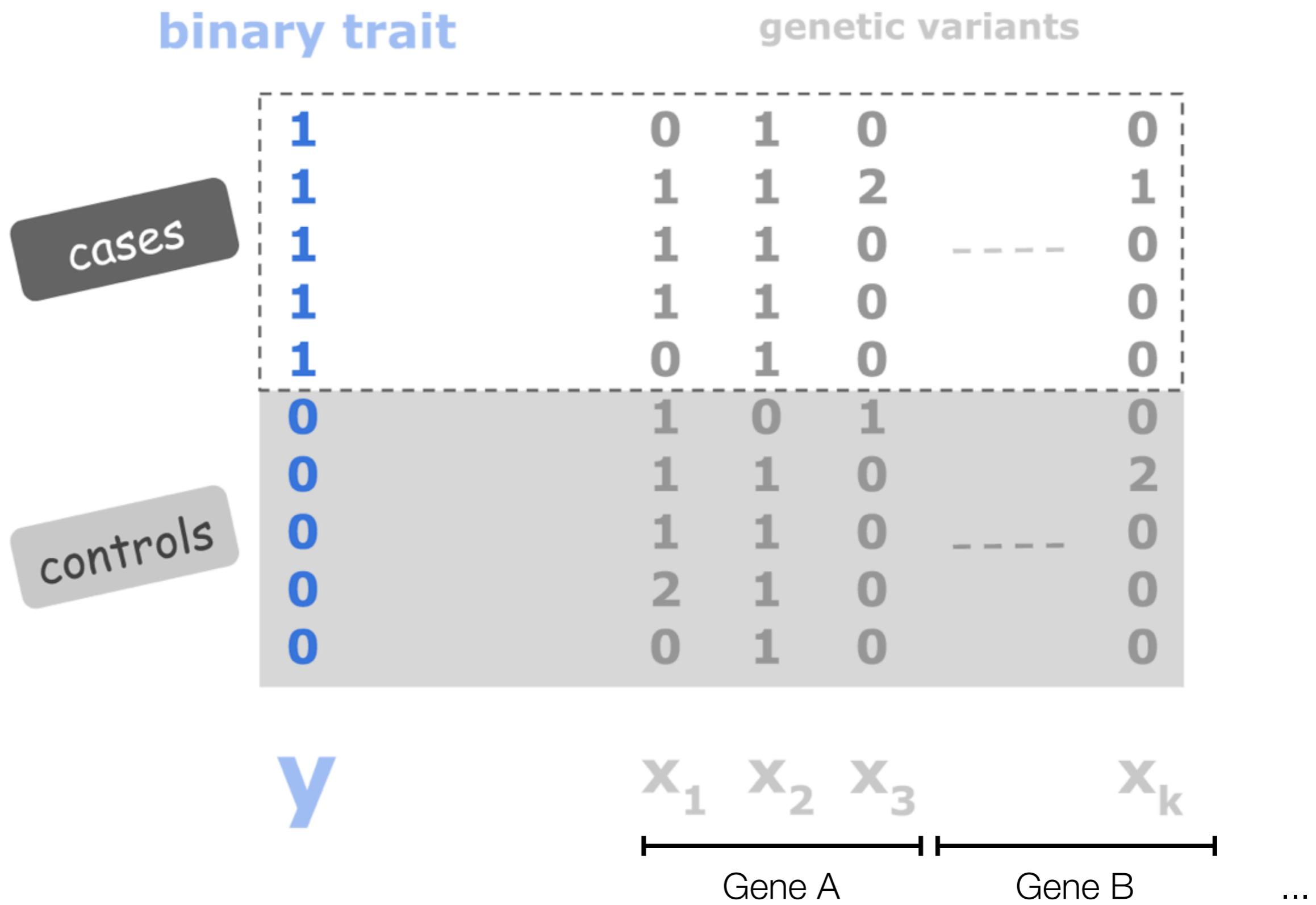
Burden tests for rare variants

		binary trait			genetic variants		
		0	1	0	0	1	0
cases	1	0	1	0	0	1	0
	1	1	1	2	0	0	0
controls	0	1	1	0	0	2	0
	0	0	1	0	0	0	0
	0	1	1	0	0	0	0
	0	1	1	0	0	0	0
	0	2	1	0	0	0	0
	0	0	1	0	0	0	0
	Y	x_1	x_2	x_3	x_k		

Burden tests for rare variants



Burden tests for rare variants



~Simple and extensible.

Framework for new tool development

```
# gemini imports
import gemini_utils as util

def my_tool(c, args):
    """
    Execute a query against the gemini database and
    conduct a custom analysis on the results
    """
    # build and execute the relevant query against the database
    query = "SELECT * FROM variants \
              WHERE is_coding = 1"
    c.execute(query)

    # loop through the results.
    for row in c:
        gt_types = np.array(unpack(row['gt_types'])))
        gt_phases = np.array(unpack(row['gt_phases'])))
        gt_bases = np.array(unpack(row['gts'])))

        # MAGIC HAPPENS HERE

def run(parser, args):
    if os.path.exists(args.db):
        conn = sqlite3.connect(args.db)
        conn.isolation_level = None
        conn.row_factory = sqlite3.Row
        conn.cursor()
```

Development in Python
with C and C++ code
for heavy lifting.

1. Issue a query against DB

2. Iterate through results

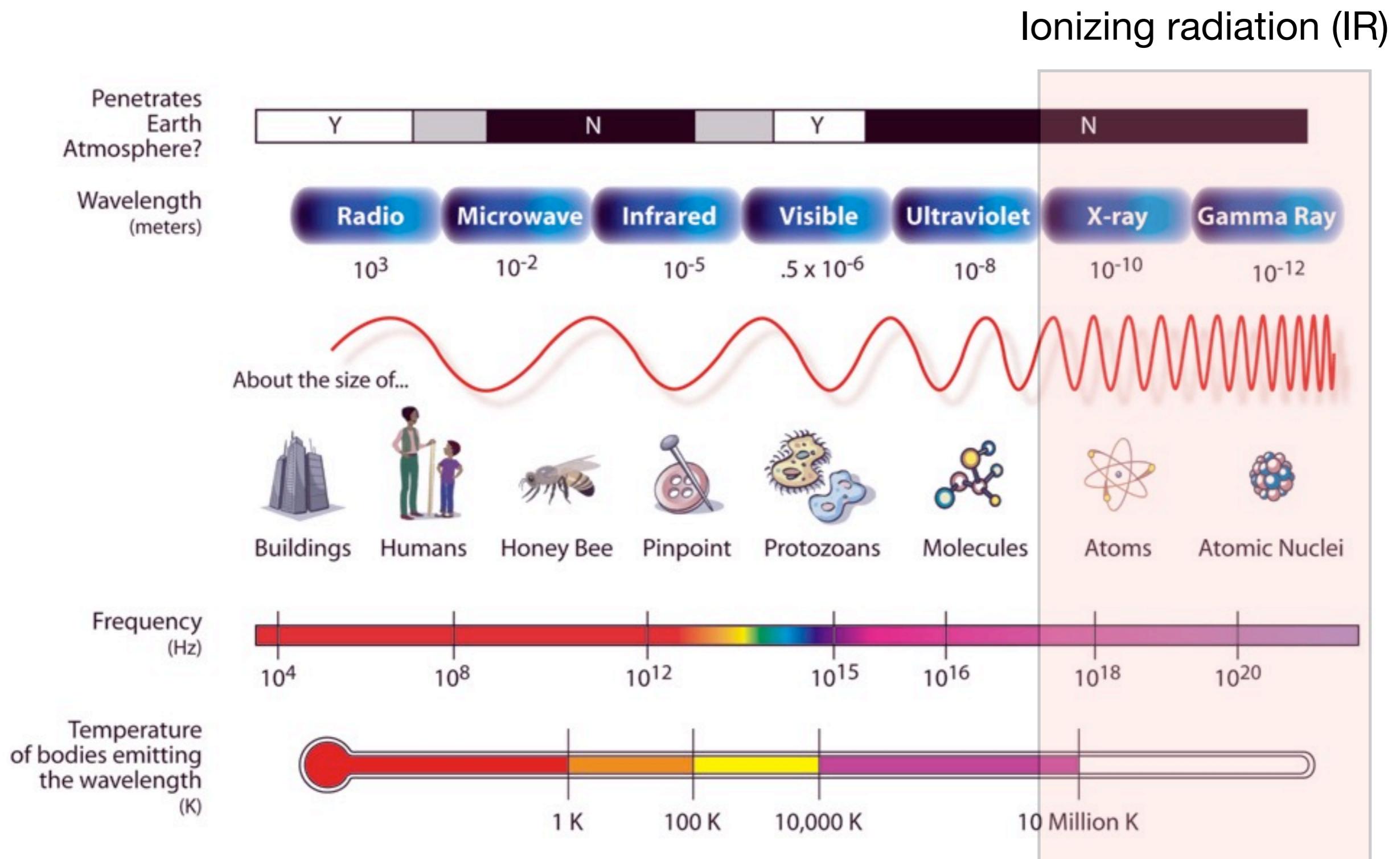
3. Work with genotypes

4. Your genius.

> **gemini my_tool -arg1 -arg2 my.db**

**What is the genetic basis of
extreme, unexplained
sensitivity to ionizing radiation?**

The electromagnetic spectrum



Beneficial applications of IR

1. Cancer therapy
2. Medical imaging
3. Food sterilization
4. Smoke detectors
5. QC for industrial production processes

Cellular consequences of IR

Direct effects

1. Causes direct damage to DNA and proteins in cell
2. More likely when the beam of charged particles consists of alpha particles, protons and electrons

Indirect effects

1. Causes damage by interacting with the cellular medium producing free radicals which, in turn, can damage DNA
2. Typical effect of X-rays or gamma rays

Damage to DNA can include base loss or modification, single strand gaps or double-strand breaks (DSBs)

1. Often, damage is complex with multiple lesions occurring
2. Of these types of damage, the most dangerous is DSBs

Impact of standard radiation therapy in an undiagnosed ataxia-telangiectasia (A-T) patient



Background

Several rare Mendelian disorders have been described sharing the phenotype of radiation hypersensitivity with a constellation of other phenotypes including increased cancer incidence, immunodeficiency, neurologic defects, and developmental delay.

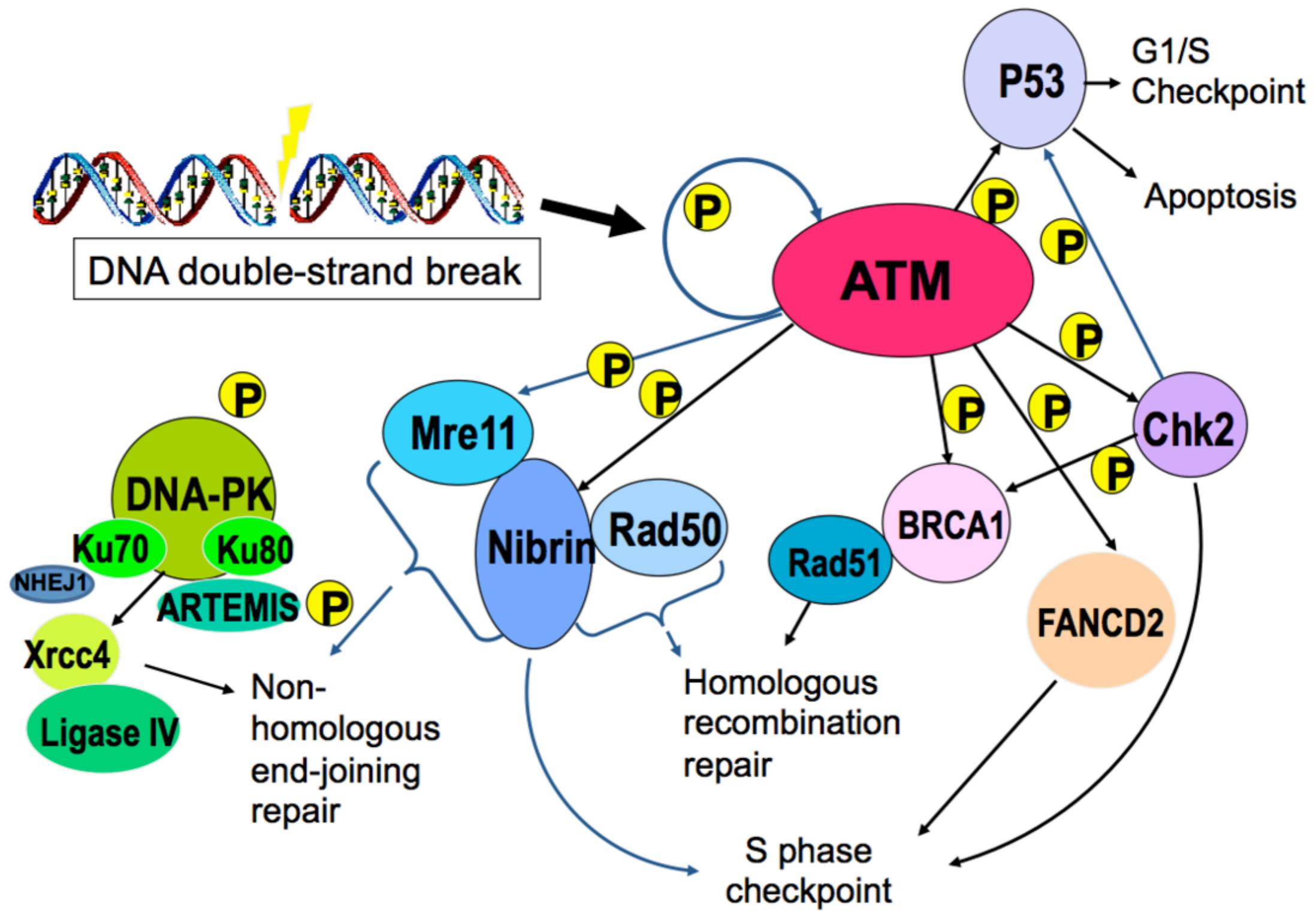
Studies of these rare, autosomal recessive genetic disorders has **identified key molecules in the human DNA damage response pathways**.

- Studies of Ataxia-Telangiectasia (A-T) identified ATM
- Studies of Nijmegen Breakage Syndrome (NBS), identified NBN as part of the MRE11-RAD50-NBN complex

Known radiosensitivity disorders

	A-T	NBS	ATLD	LIG4S
Incidence	$1/10^5$	>100 cases	40 cases	18 cases
Ataxia	+		+	
Telangiectasia	+			
Radiation sensitivity	+	+	+	+
Immunodeficiency	+	+	+	pancytopenia
Chromosomal instability	+ (7,14)	+ (7,14)	+ (7,14)	+
Cell cycle defects	+	+	+	
Lymphoid cancers	+	+	ND	+
Microcephaly		+		+
Growth retardation		+		+
Gene mutated:	ATM	NBN	MRE11	LIG4

ATM regulates the mammalian cellular response to DNA DSBs



RNF168 ubiquitin ligase: precedent for finding new *rads* genes among *rads* cell lines

- 2009: Discovery of radiosensitivity disorder RIDDLE syndrome caused by mutations in ubiquitin ligase RNF168
 - *Stewart et al. Cell 136:420*
- RS66 found to have mutations in RNF168, leading to lack of the protein.
 - *Devgan et al. 2011*

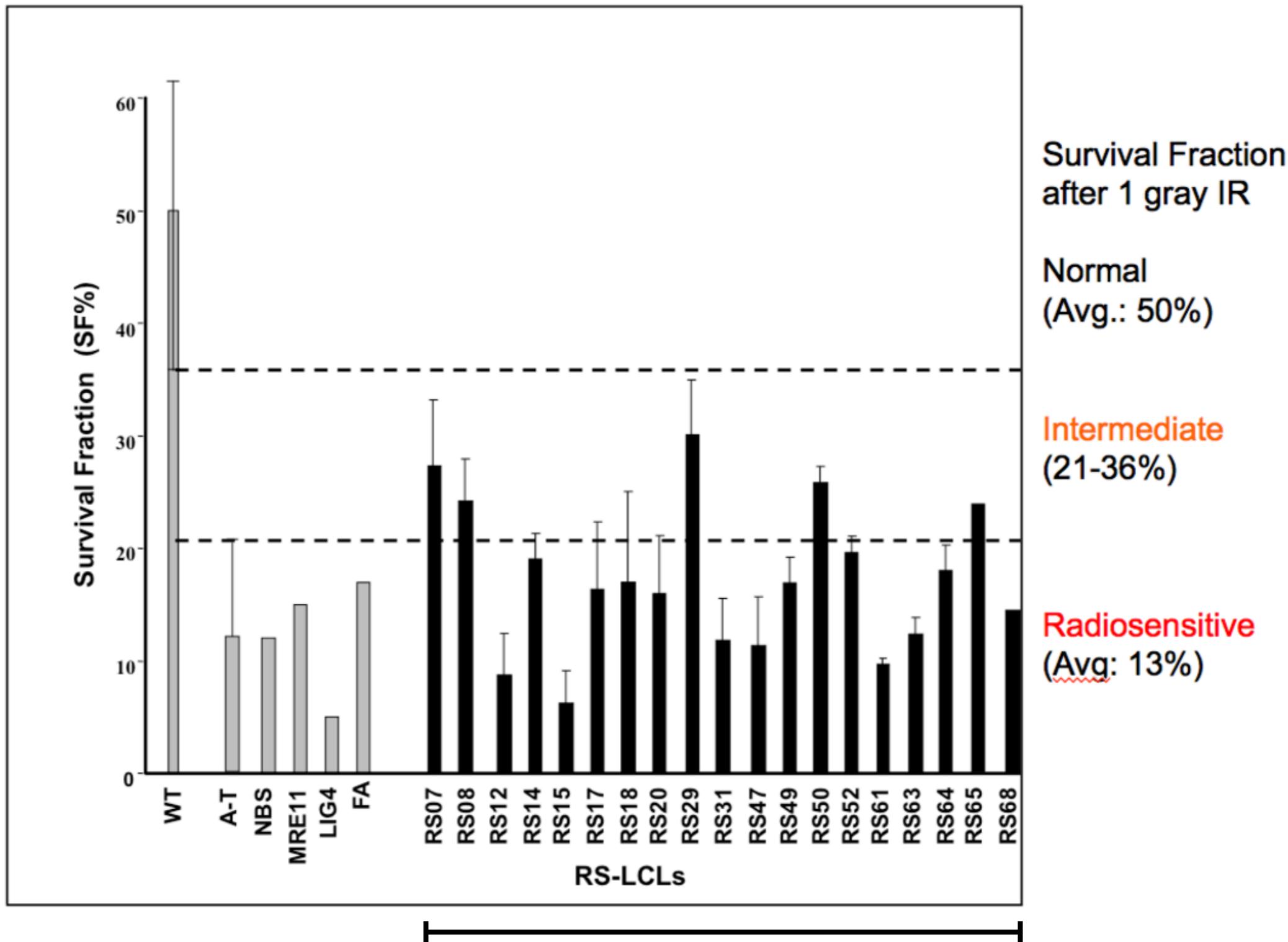
96 patients with *unexplained* radiosensitivity

- Collaboration w Pat Concannon & Richard Gatti, UCLA
- All patients have been previously screened for dysfunction in known radiosensitivity genes (e.g., ATM and NBN).
- Thus, opportunity to **discover new genes underlying response to DNA damage.**
- We hypothesize that each patient has a single gene disorder - yet the phenotype is only observed when they receive IR.
- We have cell lines - functional assays to confirm variants via complementation assays

Rationale for exome sequencing

- The RS cell lines derive from patients who display multiple clinical phenotypes known to be associated with A-T and NBS, autosomal recessive, monogenic, DNA damage response disorders
- Prior candidate gene approaches yielded biallelic mutations in single genes LIG4 or RNF168
- >84% of disease-causing mutations in ATM are in the coding regions
- The majority of disease-causing mutations in ATM and NBS were truncating mutations, readily recognizable (though not totally solved)

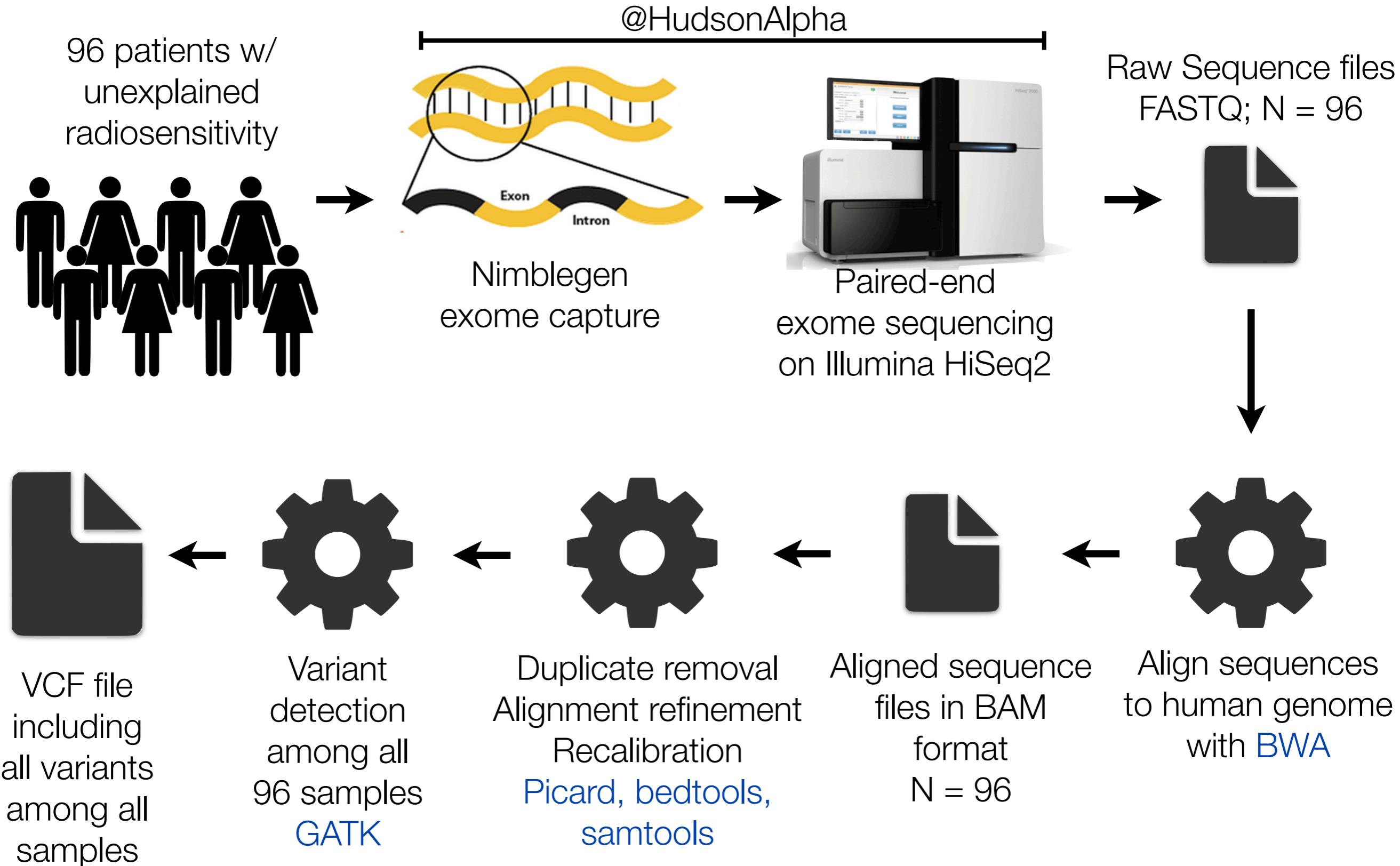
Survival after Ionizing Radiation



19 of 96 IR-sensitive patients

The pipeline

Project Title: *Identification of radiation sensitivity alleles by whole exome sequencing.*
PI: Pat Concannon
Co-investigator: Aaron Quinlan
Source: NIH/NIEHS (R21 ES020521-01)

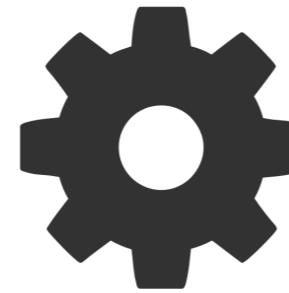


The pipeline

Project Title: *Identification of radiation sensitivity alleles by whole exome sequencing.*
PI: Pat Concannon
Co-investigator: Aaron Quinlan
Source: NIH/NIEHS (R21 ES020521-01)



VCF file
including
all variants
among all
96 samples



Predict functional
impact of variants
on coding genes



Annotate with **gemini**.
Explore variants
and screen for
candidates underlying
radiosensitivity.

Preliminary variant filters

- We are looking for rare, highly penetrant alleles.
- **Exclude if:**
 - appears in dbSNP, except if has clin. sig.
 - in 1000 Genomes with MAF > 1%
 - in Exome Sequencing Project w/ MAF > 1%
- **Looking for:**
 - LoF variants fitting recessive model
 - LoF compound heterozygotes
 - Genes with 2 het LoF mutations
 - Genes known to be phosphorylation targets of ATM/ATR
 - Genes with functions consistent with a role in DNA damage response.

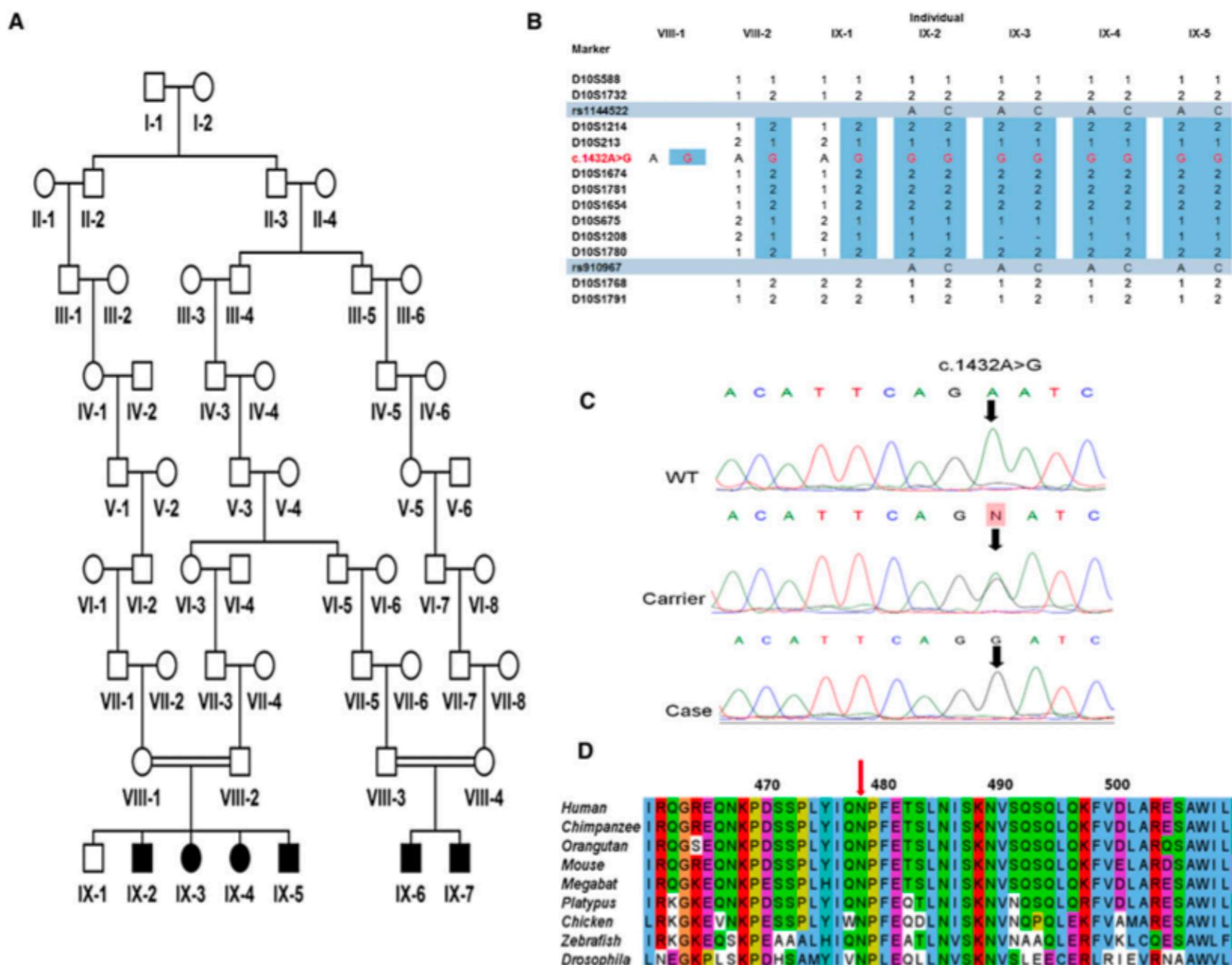
A quick, unlikely success

Sample A21: chr10, MTPAP, exon9, N478D homozygote



Defective Mitochondrial mRNA Maturation Is Associated with Spastic Ataxia

Andrew H. Crosby,^{1,7,*} Heema Patel,^{1,7} Barry A. Chioza,¹ Christos Proukakis,^{1,2} Kay Gurtz,³ Michael A. Patton,¹ Reza Sharifi,¹ Gaurav Harlalka,¹ Michael A. Simpson,¹ Katherine Dick,¹ Johanna A. Reed,¹ Ali Al-Memar,⁴ Zofia M.A. Chrzanowska-Lightowlers,⁵ Harold E. Cross,⁶ and Robert N. Lightowlers⁵



Role for MTPAP in stability of miRNA

RS cell lines tested with miRNA expression array:
MTPAP mutant cell line had global decrease in
miRNA compared to others—Hailang Hu, Gatti lab

Further evidence: MTPAP adds nucleotides to the 3'
end of miRNA, increasing their stability
Wyman et al. Genome Res. 2011 21: 1450-1461

Current: complementation assays in patient's cell
lines

coffee.

Making sense of complex datasets

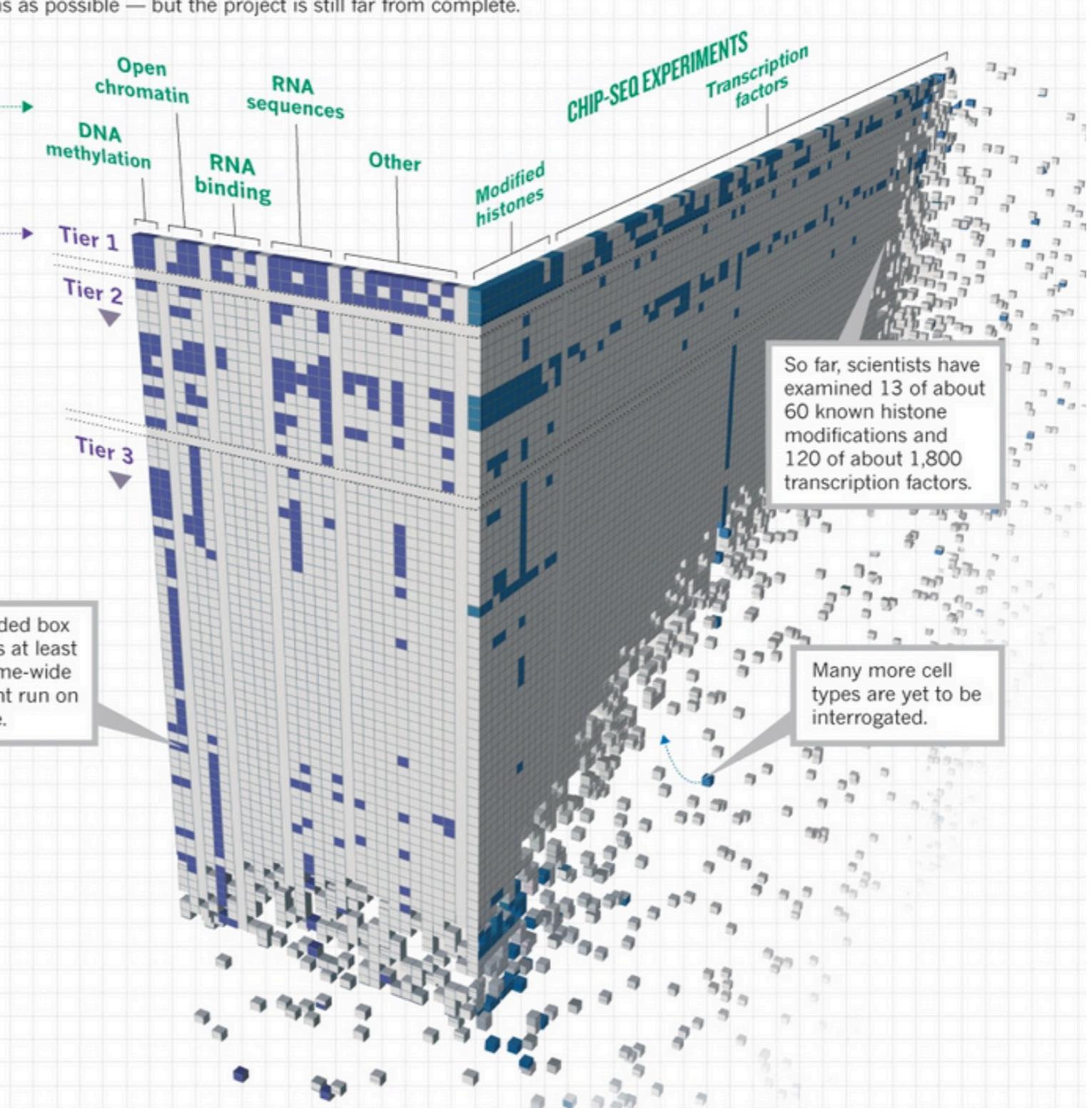
MAKING A GENOME MANUAL

Scientists in the Encyclopedia of DNA Elements Consortium have applied 24 experiment types (across) to more than 150 cell lines (down) to assign functions to as many DNA regions as possible — but the project is still far from complete.

EXPERIMENTAL TARGETS	
DNA methylation:	regions layered with chemical methyl groups, which regulate gene expression.
Open chromatin:	areas in which the DNA and proteins that make up chromatin are accessible to regulatory proteins.
RNA binding:	positions where regulatory proteins attach to RNA.
RNA sequences:	regions that are transcribed into RNA.
ChIP-seq:	technique that reveals where proteins bind to DNA.
Modified histones:	histone proteins, which package DNA into chromosomes, modified by chemical marks.
Transcription factors:	proteins that bind to DNA and regulate transcription.

CELL LINES

- Tiers 1 and 2: widely used cell lines that were given priority.
- Tier 3: all other cell types.



This is a hard and important problem.

- Understand the function (or lack) of every base pair in different cell types and contexts.
- **Challenges** (among many):
 - Basic exploratory data analysis: slicing and dicing.
 - Testing for significance
 - Visualization: unbiased exploration; let the data tell its story.
- Necessary for new discovery and understanding of genome biology

bedtools: genome arithmetic



Neil Kindlon, M.S.

Staff Scientist and Software Engineer

nek3d @ virginia.edu

Research Projects and Interests: Software variation discovery and interpretation



Ryan Layer

Graduate student

r16sf @ virginia.edu

Research Interests: Scalable algorithm analysis; genome data mining and analysis interpretation.

The screenshot shows the Google Code project page for 'bedtools'. The page title is 'bedtools - bedtools: a flexible suite of utilities for comparing genomic features.' The navigation bar includes links for Project Home, Downloads, Wiki, Issues, Source, and Administer. The 'Project Home' tab is selected. Below it, there are tabs for Summary and People. The 'Project Information' section on the left lists the following details:

- Starred by 96 users ([Project feeds](#))
- Code license: GNU GPL v2
- Labels: bioinformatics, genomics, bed, sam, bam, overlap, features, sequencing, intersect, coverage, gff, vcf, bedgraph, intervals, genome arithmetic
- Members: [aaronquinlan](#)
- Your role: Owner
- Featured: Downloads

The 'Downloads' section contains links to 'BEDTools-User-Manual.v4.pdf' and 'BEDTools.v2.16.2.tar.gz'. There is also a link to 'Show all >'. The main content area on the right is titled 'bedtools' and describes the project as a flexible suite of utilities for comparing genomic features. It includes a sidebar with links to Citation, Latest news (Version 2.16.2, 30-March-2012), BEDTools Summary, User base, Brief example, Table of supported utilities, Documentation, Notes regarding usage, Installation, Source Repository, Package Managers, and Contact. A section for 'Citation' asks users to cite a specific article if they use BEDTools in their research. Another note mentions 'pybedtools', a Python extension of BEDTools, which provides a powerful and flexible Python interface for manipulating and comparing genomic features in BED/VCF/GFF/GTF/SAM/BAM format. A final list of references includes Quinlan AR and Hall IM (2010) and Dale RK, Pedersen BS, and Quinlan AR (2011).

Software <http://code.google.com/p/bedtools/>

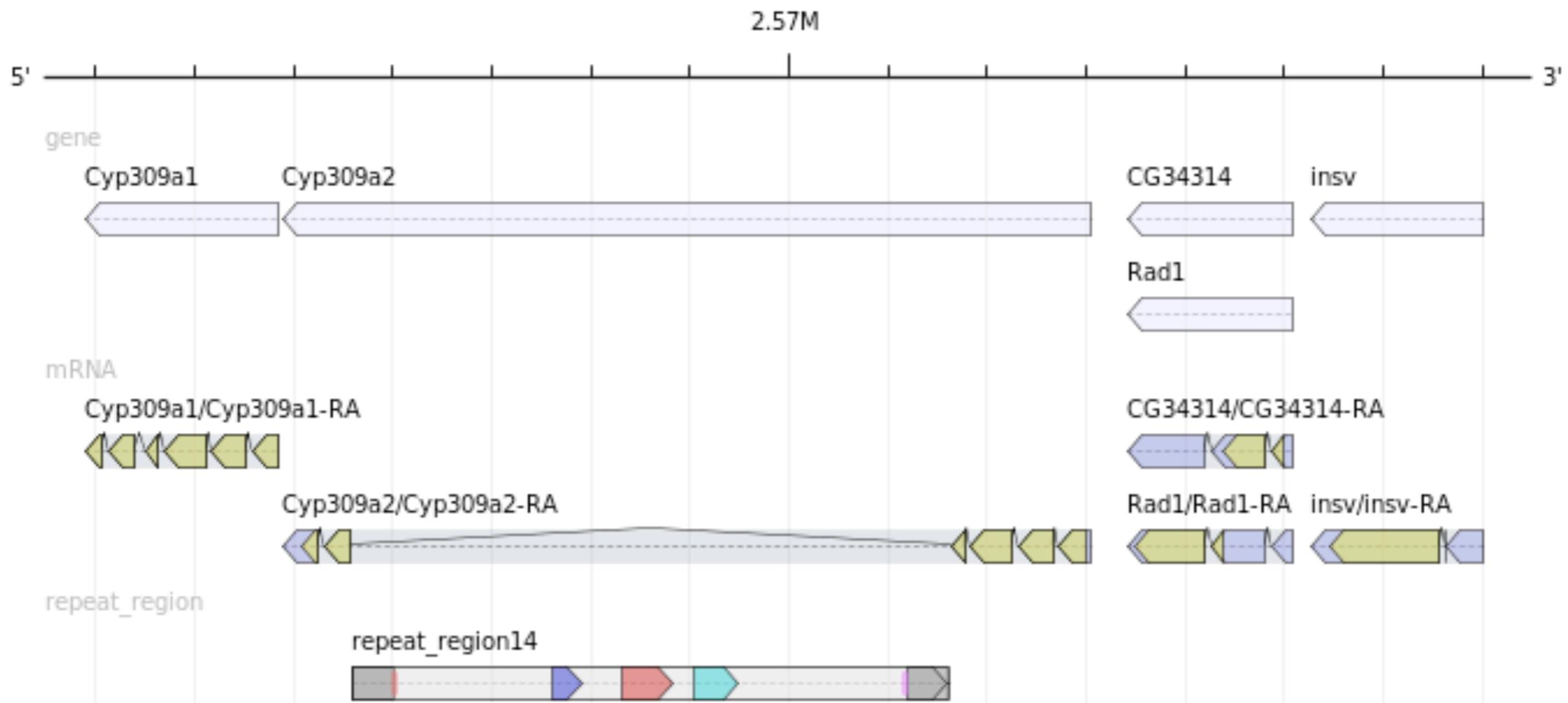
Docs <https://bedtools.googlecode.com/files/BEDTools-User-Manual.v4.pdf>

Tools for comparing and exploring genome “intervals”

What is a genome “interval”?

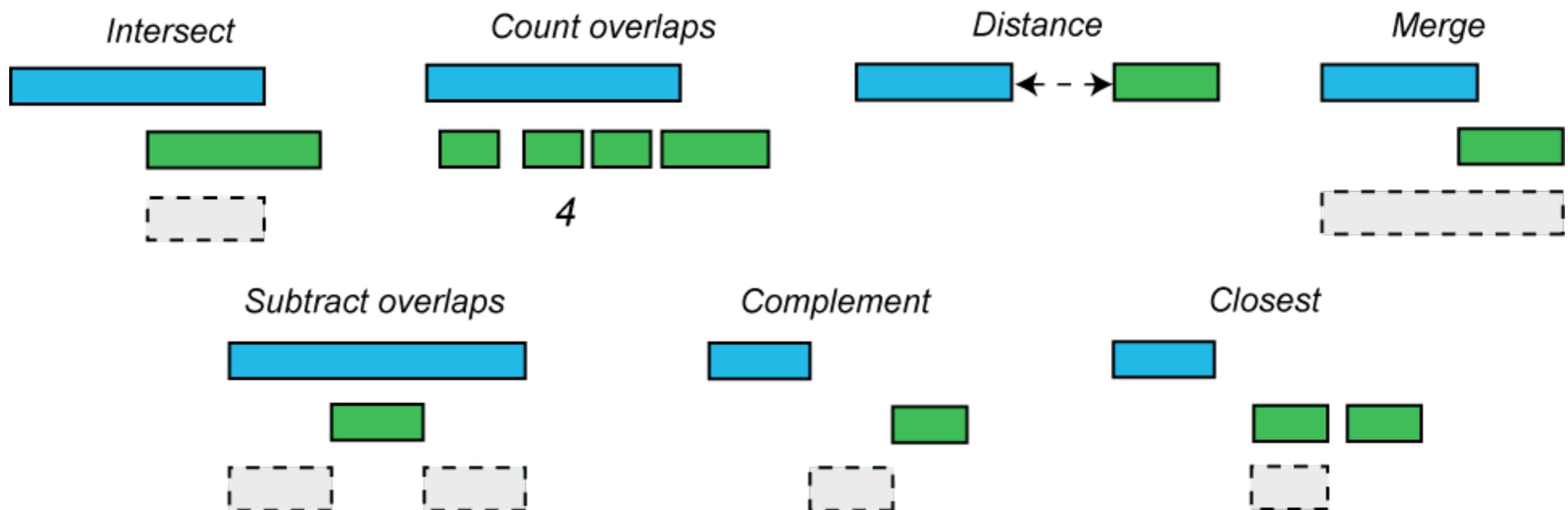
- Genes: exons, introns, UTRs, promoters
- Conservation
- Genetic variation
- Transposons
- Origins of replication
- TF binding sites
- CpG islands
- Segmental duplications
- Sequence alignments
- Chromatin annotations
- Gene expression data
- ...

Genome intervals



What is genome arithmetic?

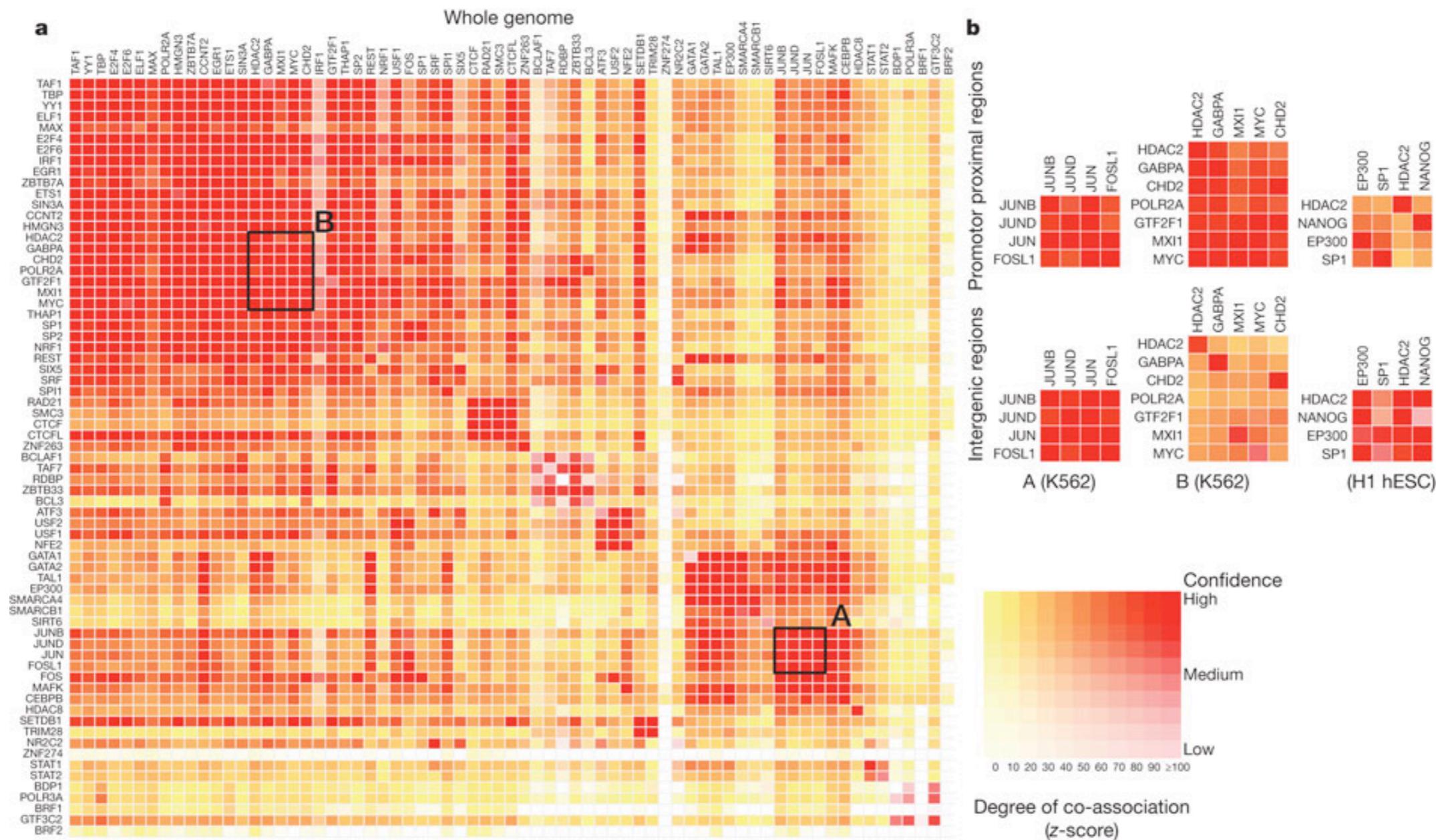
set theory on genome intervals



Answerable questions

- Closest gene to a ChIP-seq peak.
- Is my latest discovery novel w.r.t. other datasets?
- Is there strand bias in my data?
- How many genes does this mutation affect?
- Where did I fail to collect sequence coverage?
- Is my favorite feature significantly correlated with some other feature?

How do we measure significance?



Co-association between transcription factors.

Monte-Carlo simulation

Are the observed feature overlaps more common than what you would expect by chance?

1. **Observed.** Count the number of overlaps between the two sets of features (e.g. ChIP peaks for TF 1 v. TF 2)
2. **Expected.**
 - 2.1. Randomly reassign the features in each set to new genomic locations.
 - 2.2. Count the number of overlaps.
 - 2.3. Repeat 1000 or more times. **SLOW!**
3. **Return either:**
 - 3.1. **P-value:** how many times were the shuffled overlaps > observed?
 - 3.1.1. e.g., if answer is 3 and there were 1000 simulations, $P = 3e-3$
 - 3.2. **Enrichment score:** $\log_2(\text{observed} / \text{median expected})$

How do we speed this up?

We invent a new algorithm, of course!

PREPRINT

Vol. 00 no. 00 2012
Pages 1–8

Binary Interval Search (BITS): A Scalable Algorithm for Counting Interval Intersections

Ryan M. Layer¹, Kevin Skadron¹, Gabriel Robins¹, Ira M. Hall², and Aaron R. Quinlan^{3*}

¹Department of Computer Science, University of Virginia, Charlottesville, VA

²Department of Biochemistry and Molecular Genetics, University of Virginia, Charlottesville, VA

³Department of Public Health Sciences and Center for Public Health Genomics, University of Virginia, Charlottesville, VA



Ryan Layer

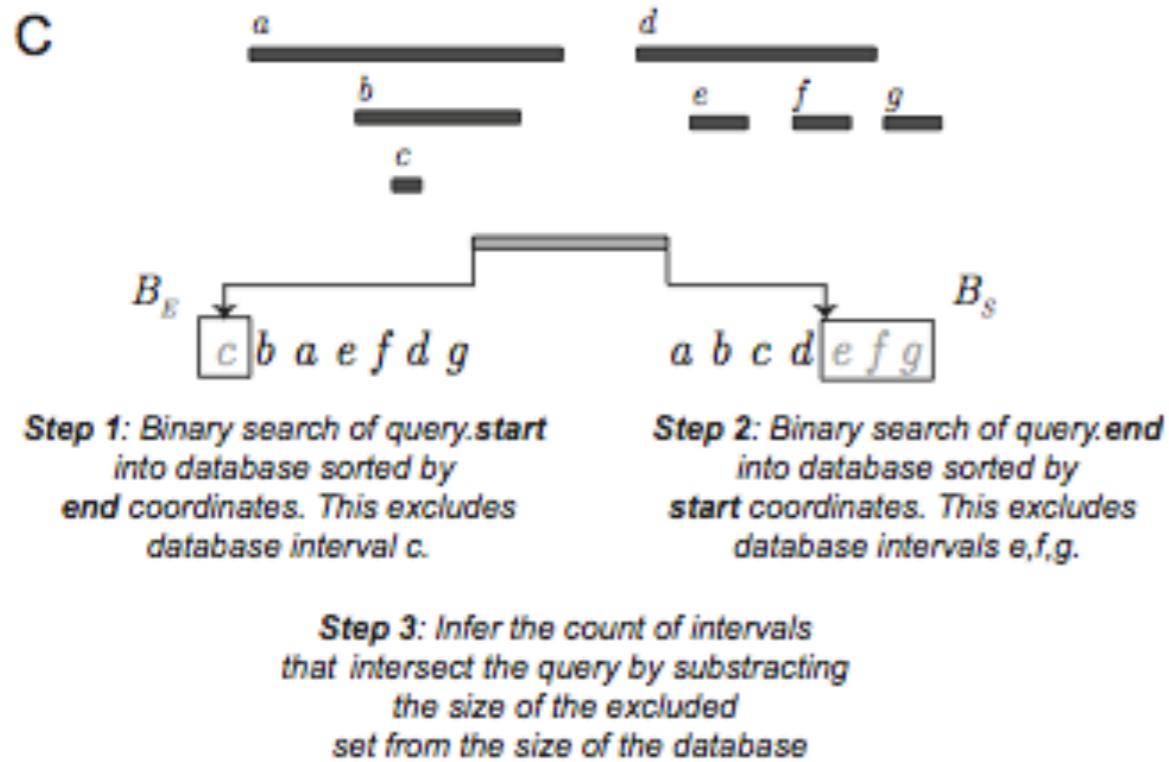
Graduate student

rl6sf @ virginia.edu

Research Interests: Scalable algorithm analysis; genome data mining and analysis interpretation.

<https://github.com/arg5x/bits>

Binary InTerval Search BITS



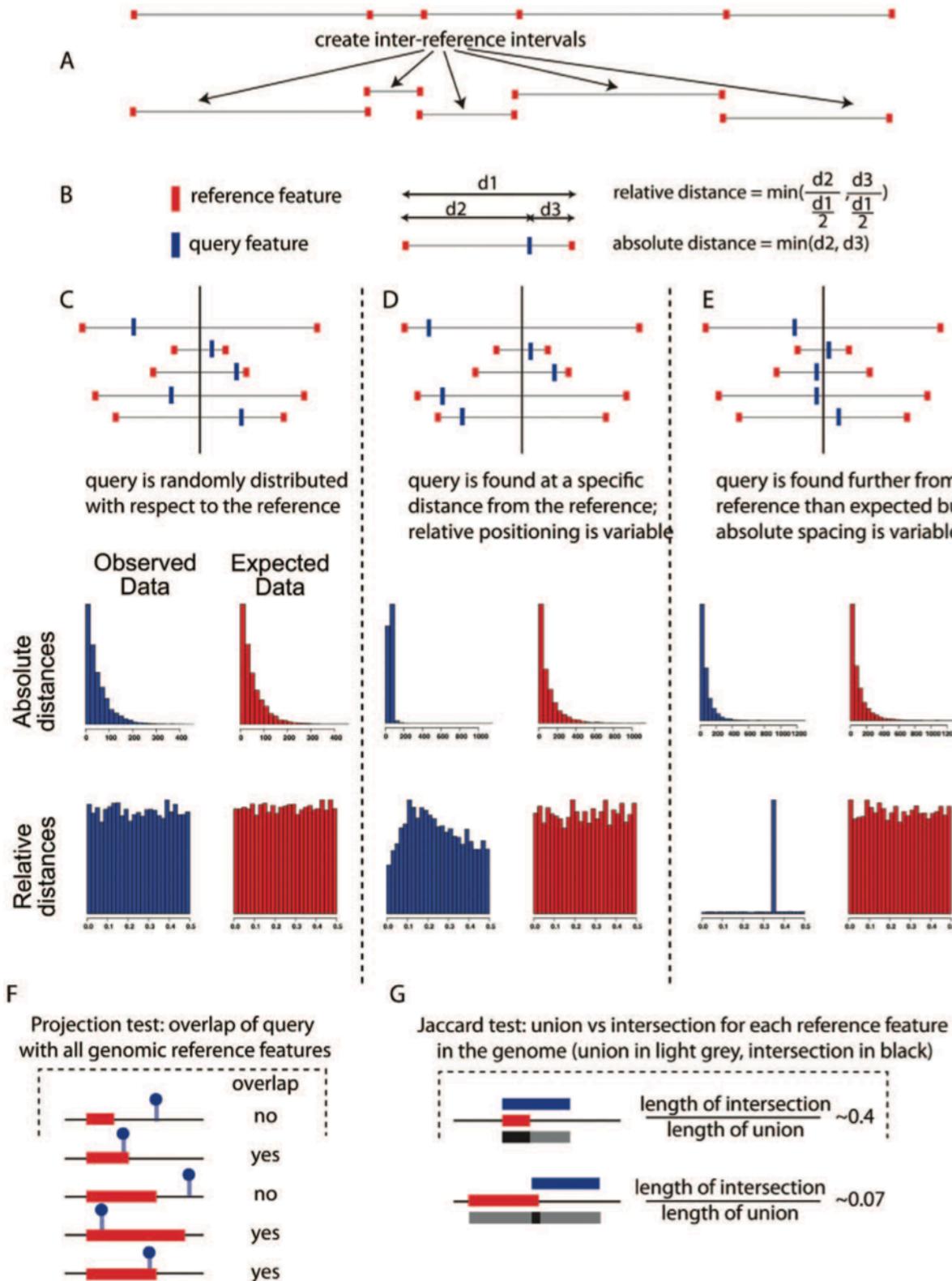
Algorithm 1: Single interval intersection counter

Input: Sorted interval starts and ends B_S and B_E , query interval a
Output: Number of intervals c intersecting a

```
Function ICOUNT( $B_S, B_E, a$ ) begin
    first  $\leftarrow$  BINARYSEARCH( $B_S, a.end$ )
    last  $\leftarrow$  BINARYSEARCH( $B_E, a.start$ )
     $c \leftarrow first - last$  /*  $= |B| - (last + (|B| - first))$  */
    return  $c$ 
```

- Novel algorithm for detecting genomic interval intersections.
- Clever aspect: unlike any other algorithm, it can deduce the **count** of overlaps without having to **enumerate** each individual intersection.
- Uses two binary searches. Very, very fast.
- If you haven't heard, speed is good.

Other approaches



Exploring Massive, Genome Scale Datasets with the GenometriCorr Package

Alexander Favorov^{1,2,3*}, Loris Mularoni^{1,9#}, Leslie M. Cope¹, Yulia Medvedeva^{2,3#b}, Andrey A. Mironov^{4,5}, Vsevolod J. Makeev^{2,3}, Sarah J. Wheelan^{1*}

¹ Department of Oncology, Division of Biostatistics and Bioinformatics, Johns Hopkins University School of Medicine, Baltimore, Maryland, United States of America, ²Vavilov Institute of General Genetics, Russian Academy of Sciences, Moscow, Russia, ³ Research Institute of Genetics and Selection of Industrial Microorganisms, Moscow, Russia, ⁴ Department of Bioengineering and Bioinformatics, Moscow State University, Moscow, Russia, ⁵ Institute for Information Transmission Problems, Russian Academy of Sciences, Moscow, Russia

NON PARAMETRIC METHODS FOR GENOMIC INFERENCE

BY PETER J. BICKEL*,†, NATHAN BOLEY*,†, JAMES B. BROWN*,†, HAIYAN HUANG*,†, NANCY R. ZHANG*,†

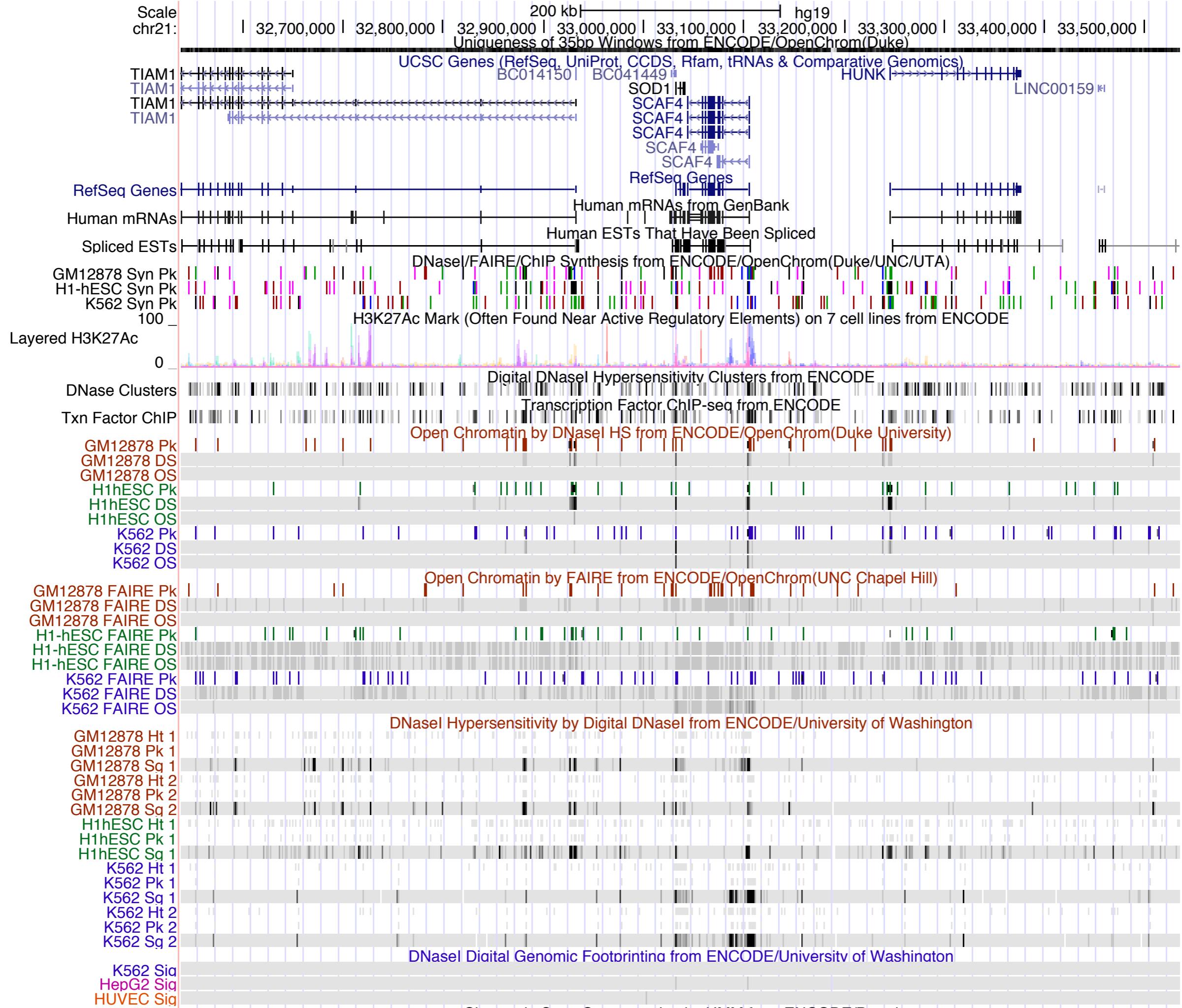
Statistics, University of California at Berkeley † and Statistics, Stanford University†

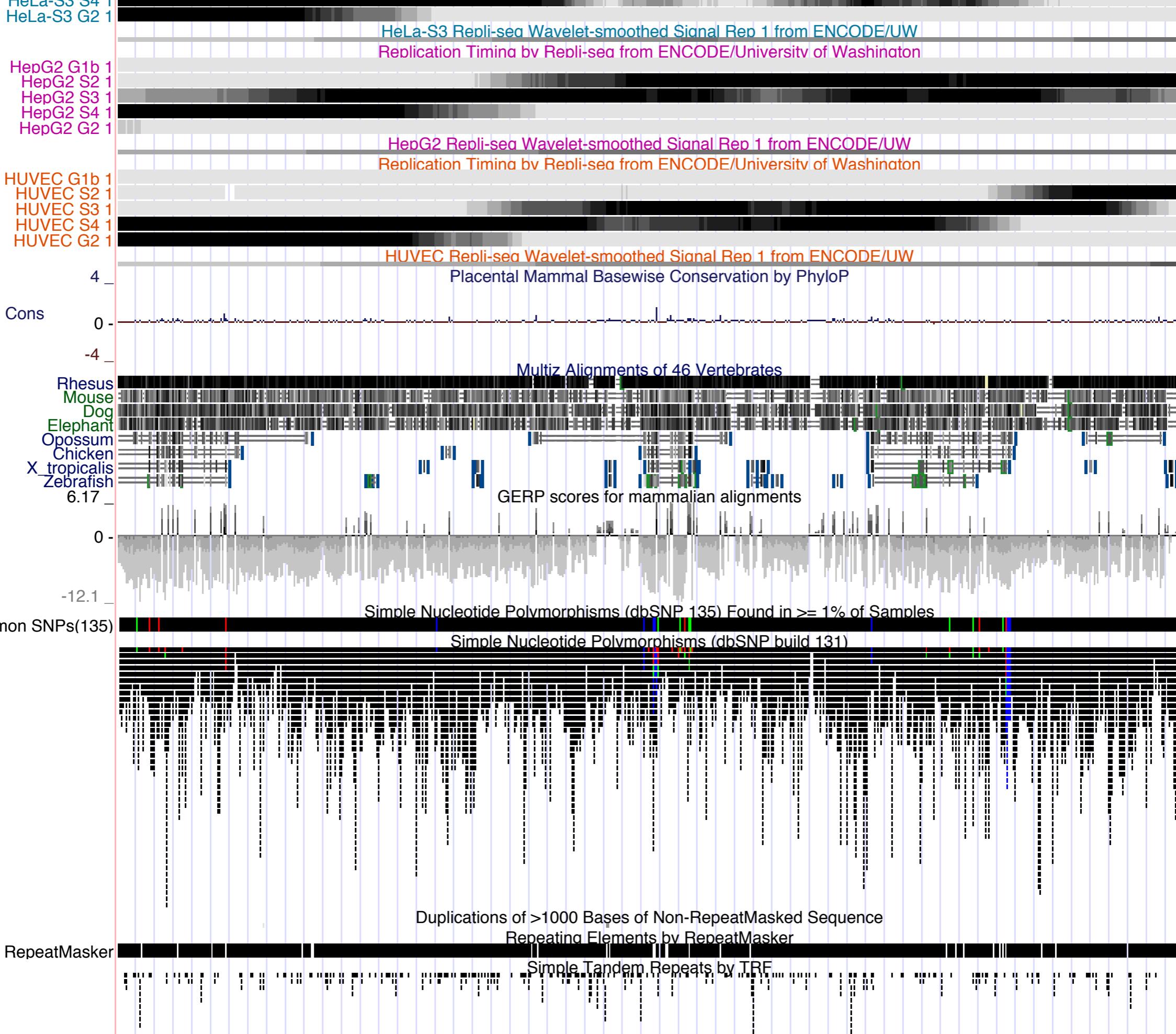
Large-scale statistical analysis of data sets associated with genome sequences plays an important role in modern biology. A key component of such statistical analyses is the computation of p-values and confidence bounds for statistics defined on the genome. Currently such computation is commonly achieved through ad hoc simulation measures. The method of randomization, which is at the heart of these simulation procedures, can significantly affect the resulting statistical conclusions. Most simulation schemes introduce a variety of hidden assumptions regarding the nature of the randomness in the data, resulting in a failure to capture biologically meaningful relationships. To address the need for a method of assessing the significance of observations within large scale genomic studies, where there often exists a complex dependency structure between observations, we propose a unified solution built upon a data subsampling approach. We propose a piecewise stationary model for genome sequences and show that the subsampling approach gives correct answers under this model. We illustrate the method on three simulation studies and two real data examples.

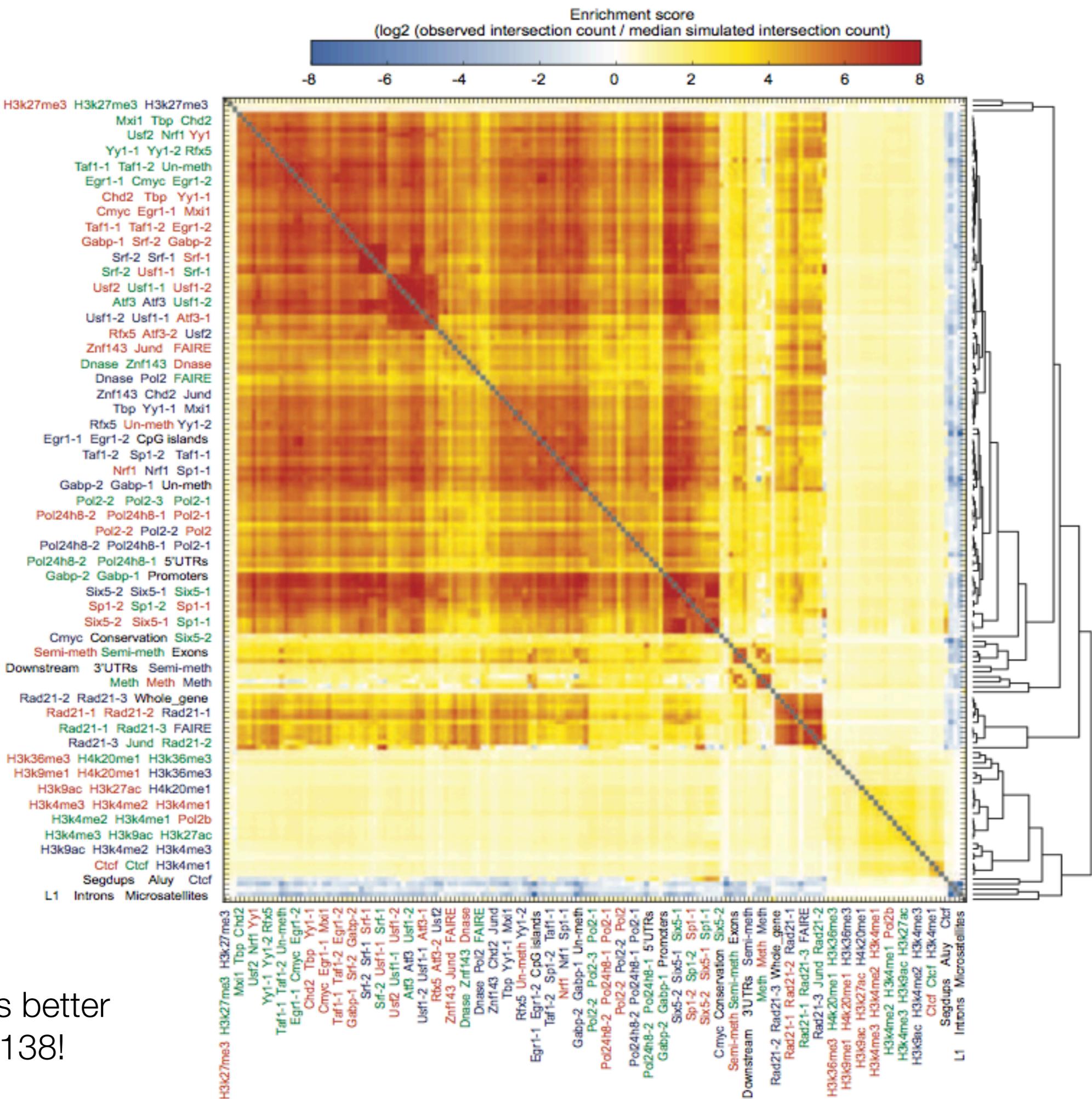
Used by ENCODE project

So, what can we do?

*Large, unbiased screens for novel
genome biology*







6 days is better
than 138!

How do we understand the details of these signals?

Problem: *The genome is big. Patterns are hard to detect and quantify.*

We need better tools for data visualization.



Space-filling curves for genomic data

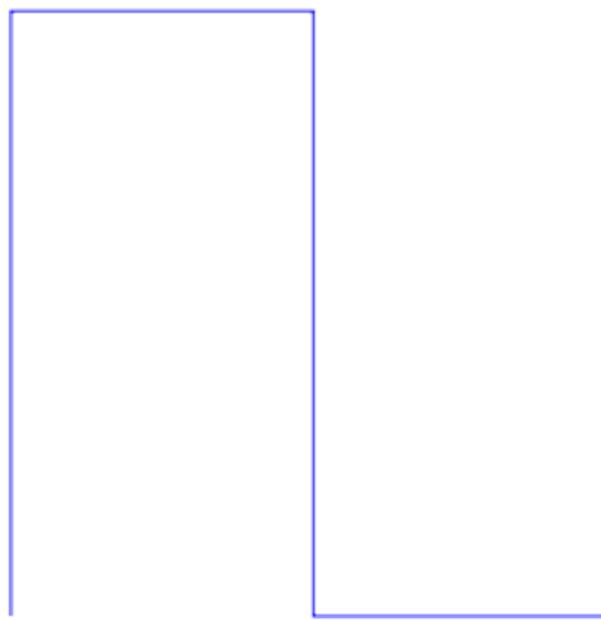
What is a space-filling curve?

Fractal algorithm to convert a 1D line to a **curve** that fills a 2 (or more) dimension space.

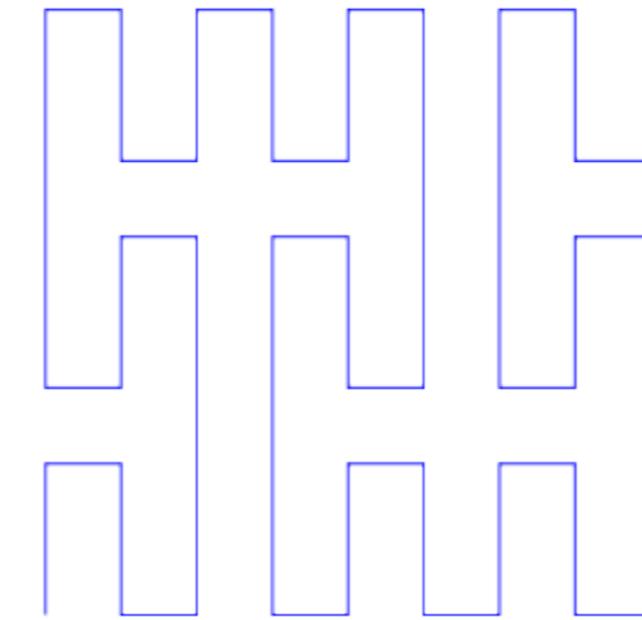
1D



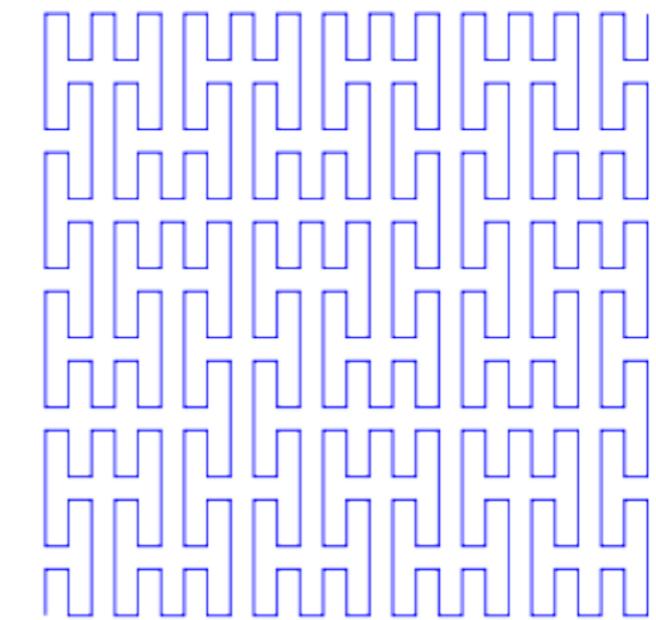
first iteration



second iteration

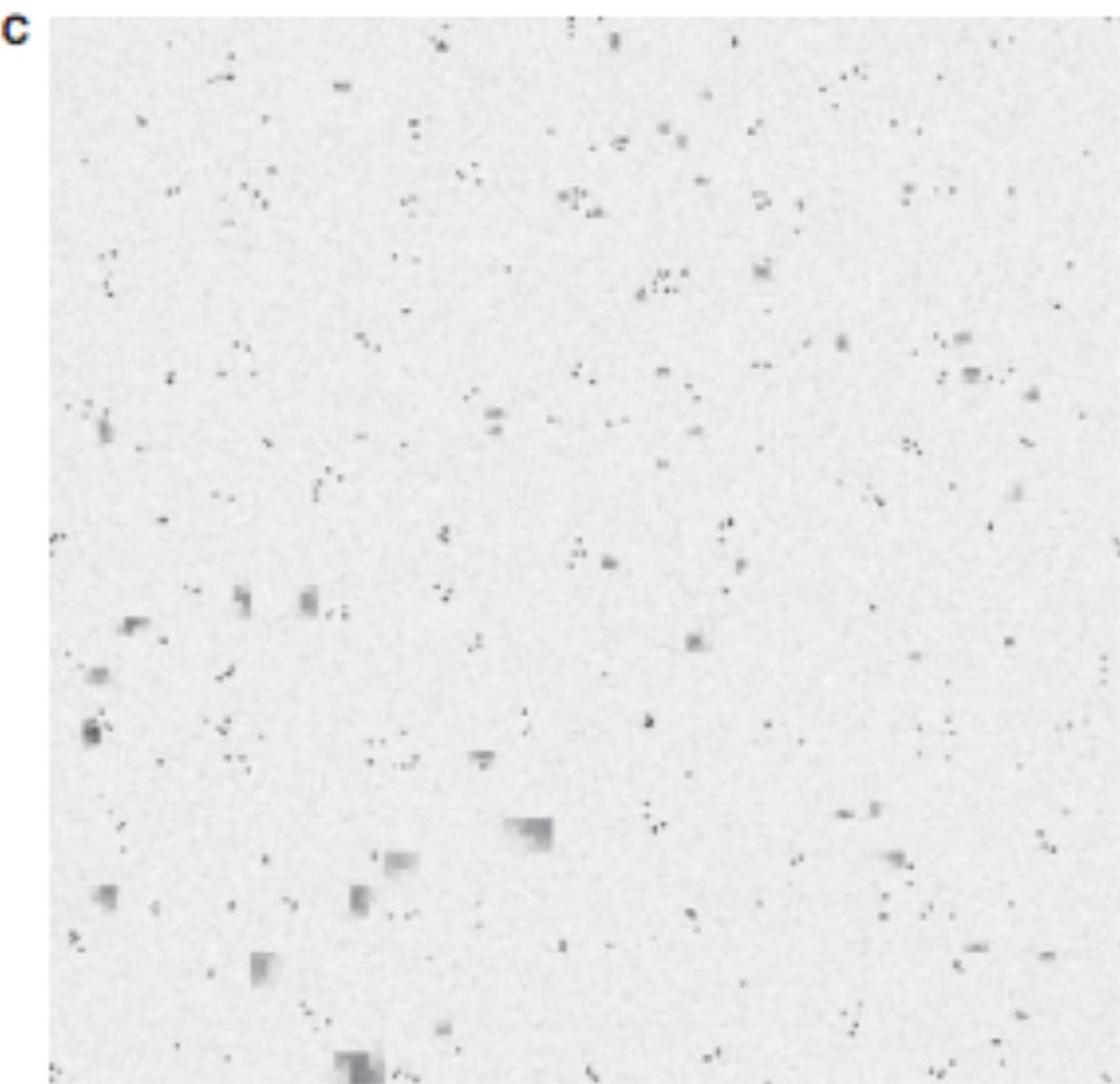
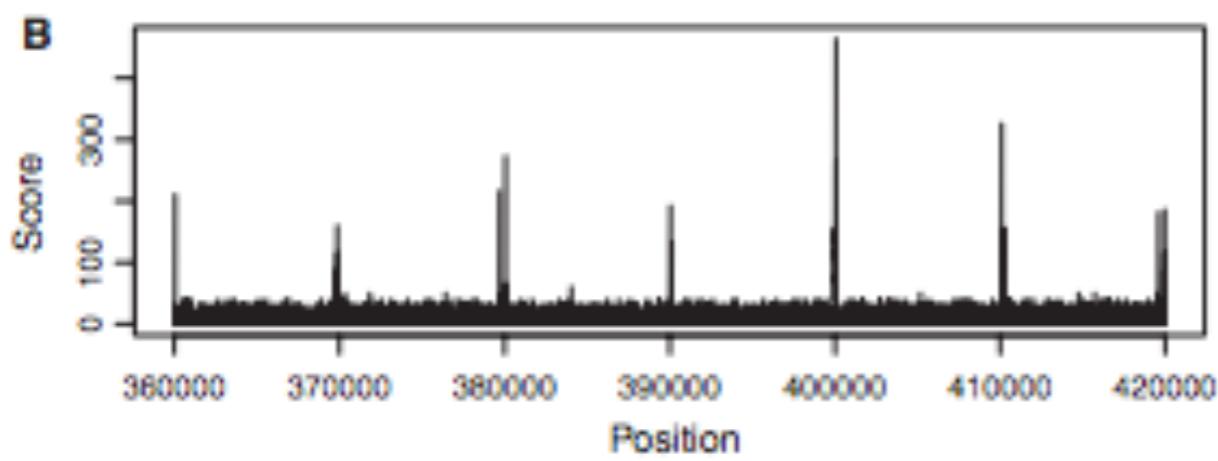
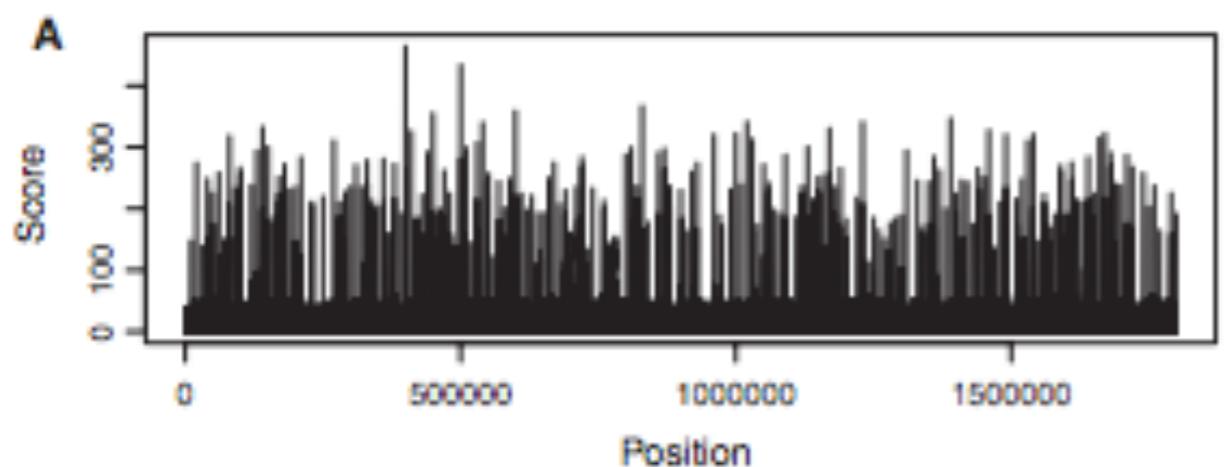


third iteration

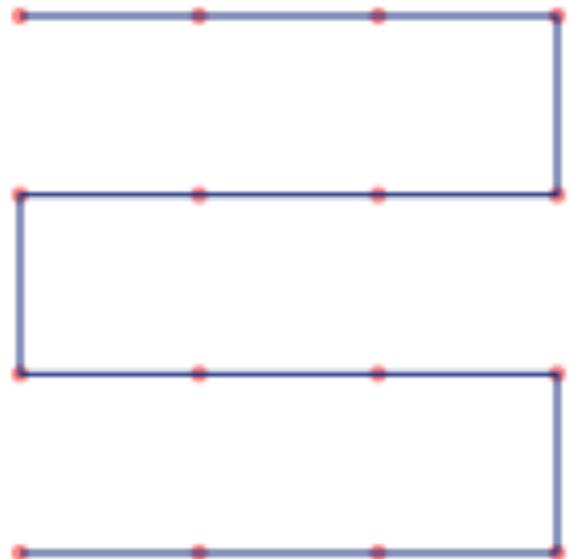


2D

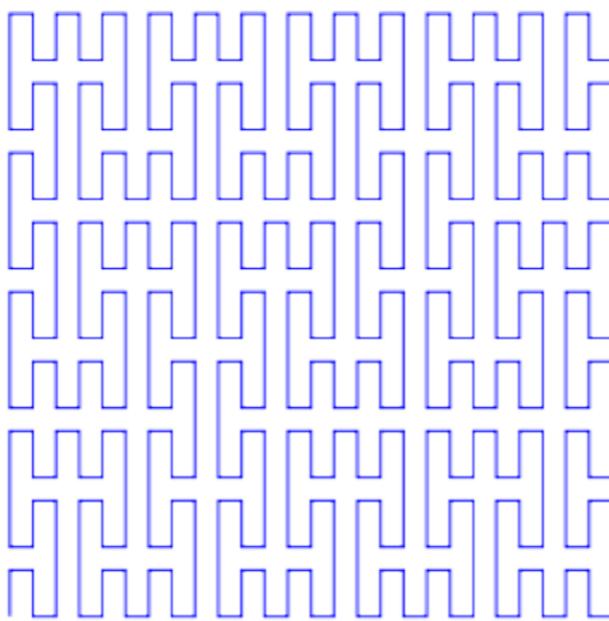
Hilbert curves use 2D. Patterns revealed.



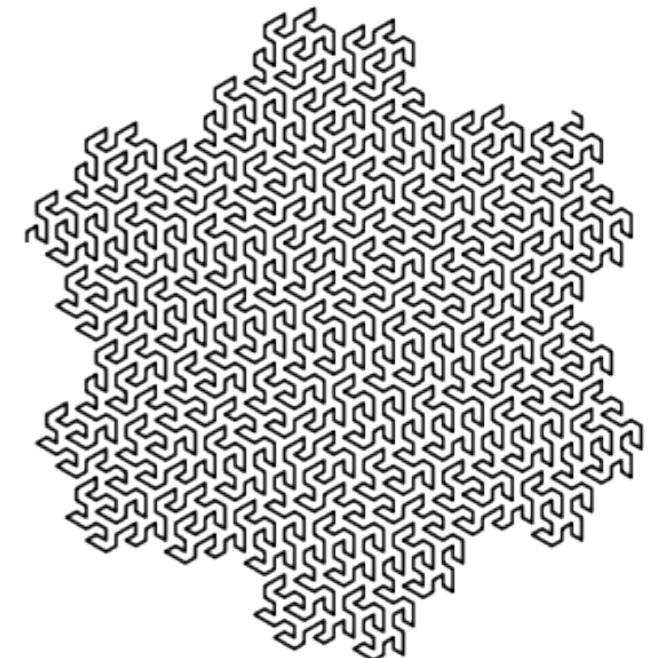
Types of SFCs



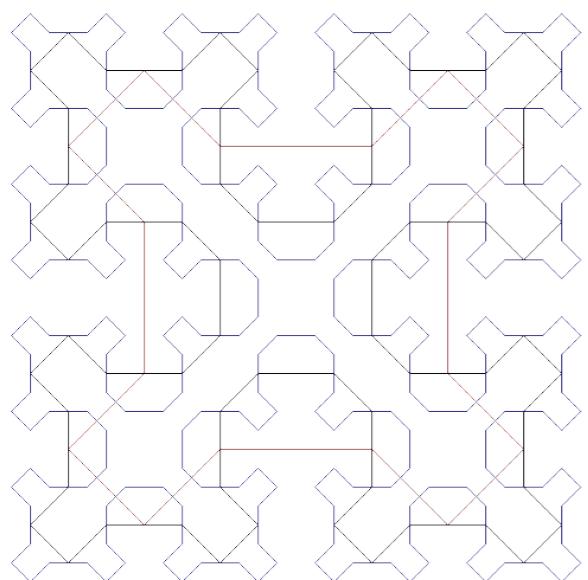
Zig-zag



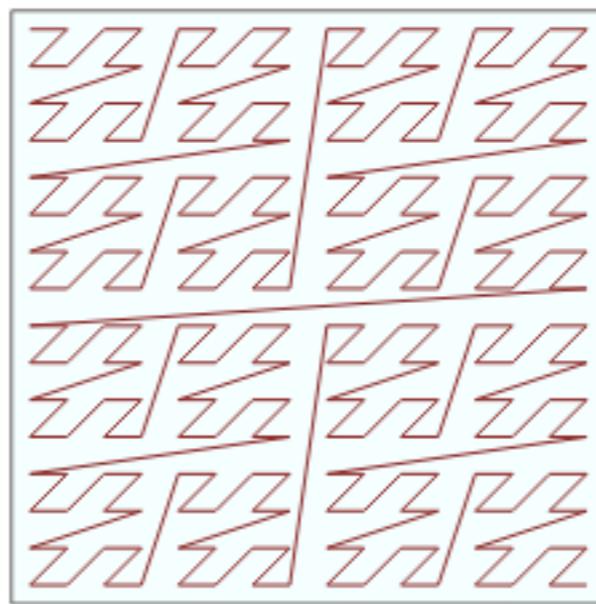
Peano



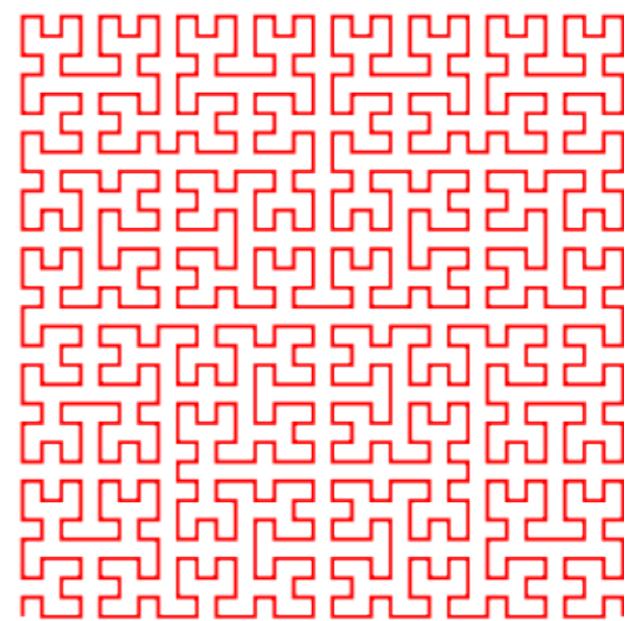
Gosper



Sierpinski



Z-order



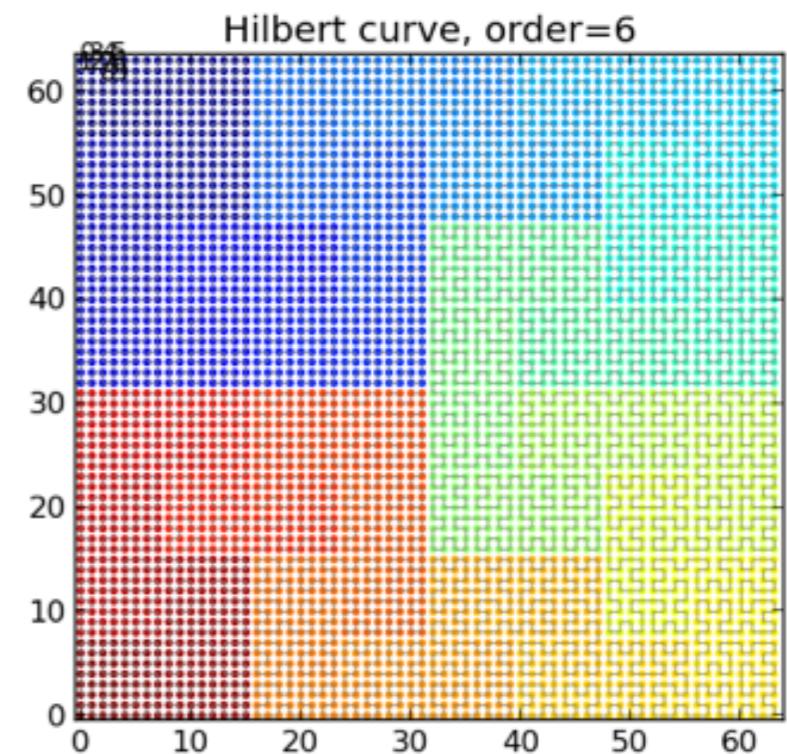
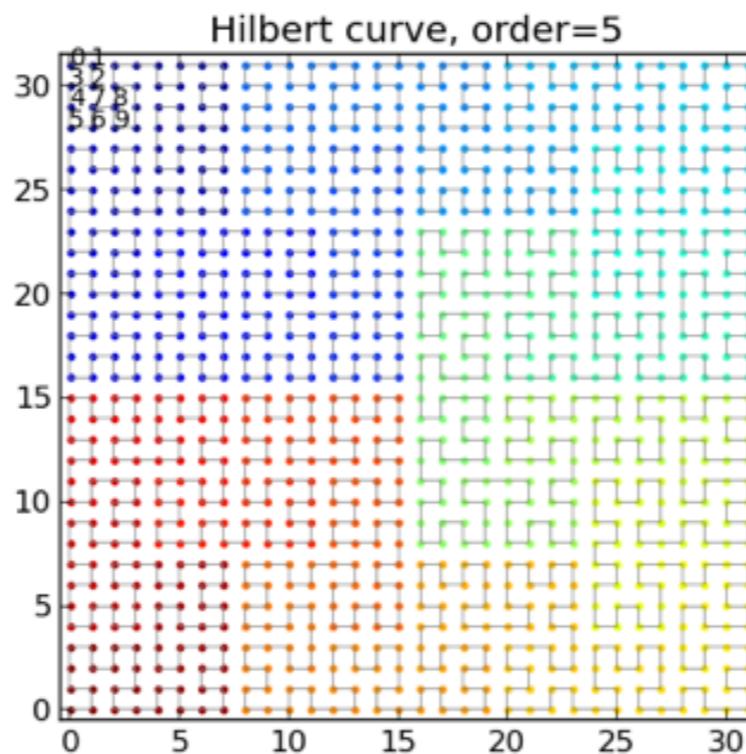
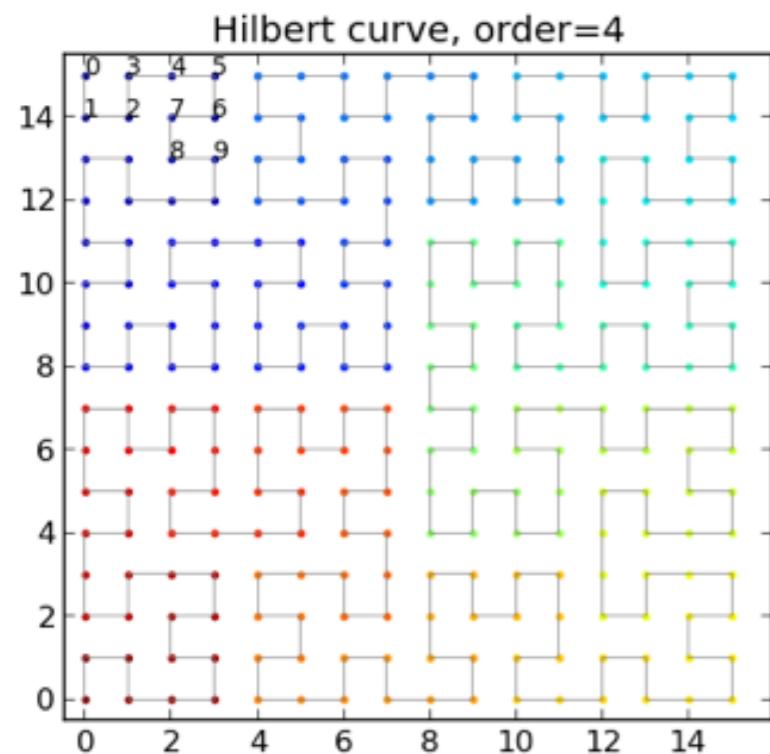
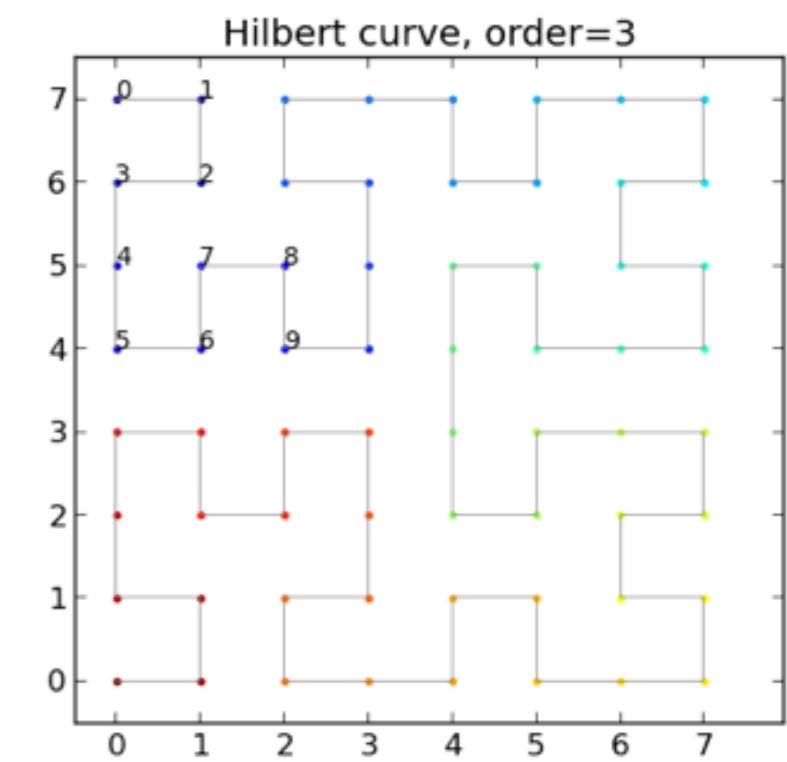
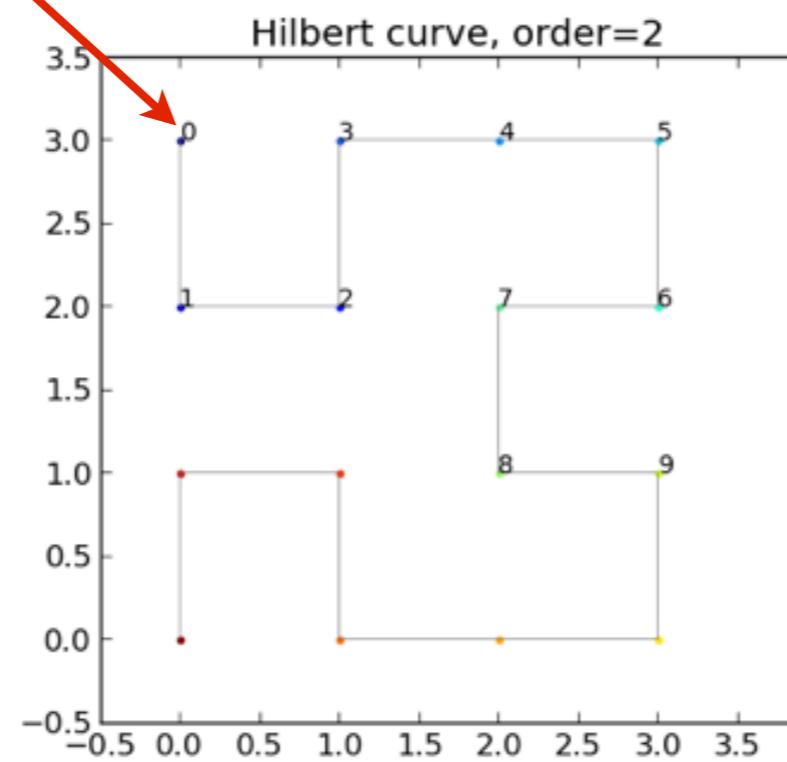
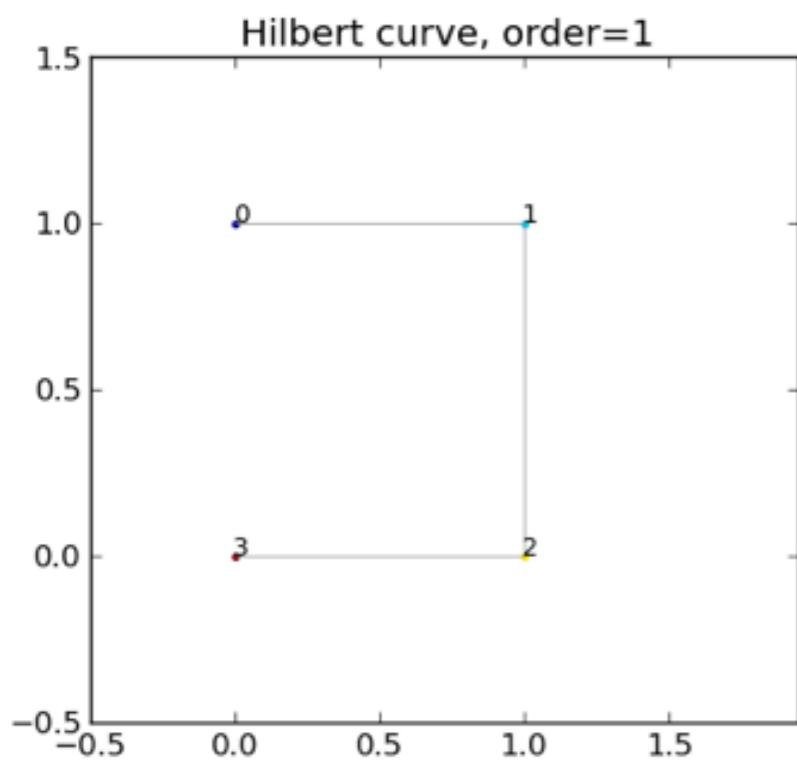
Hilbert

Why Hilbert curves for genomics?

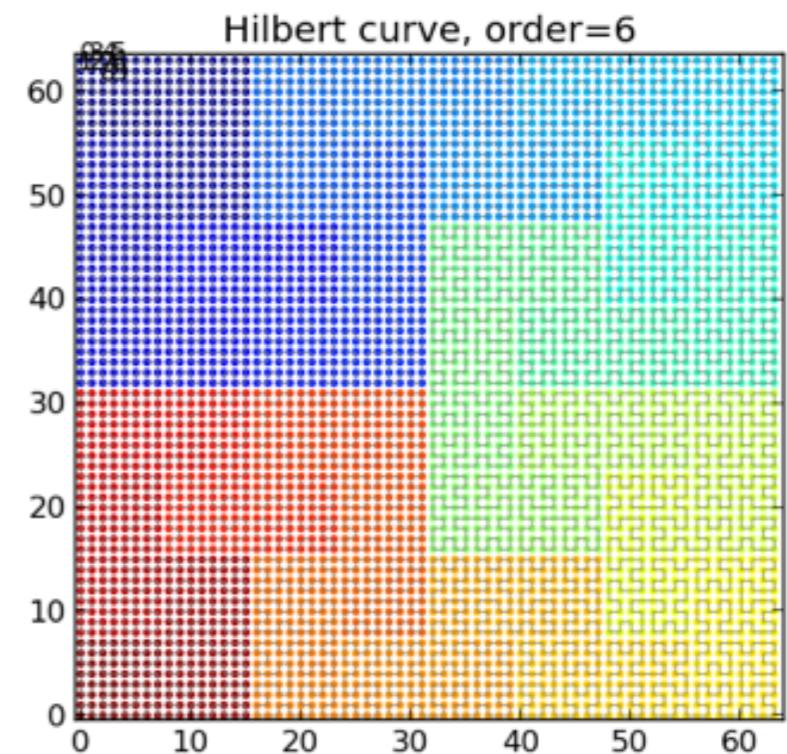
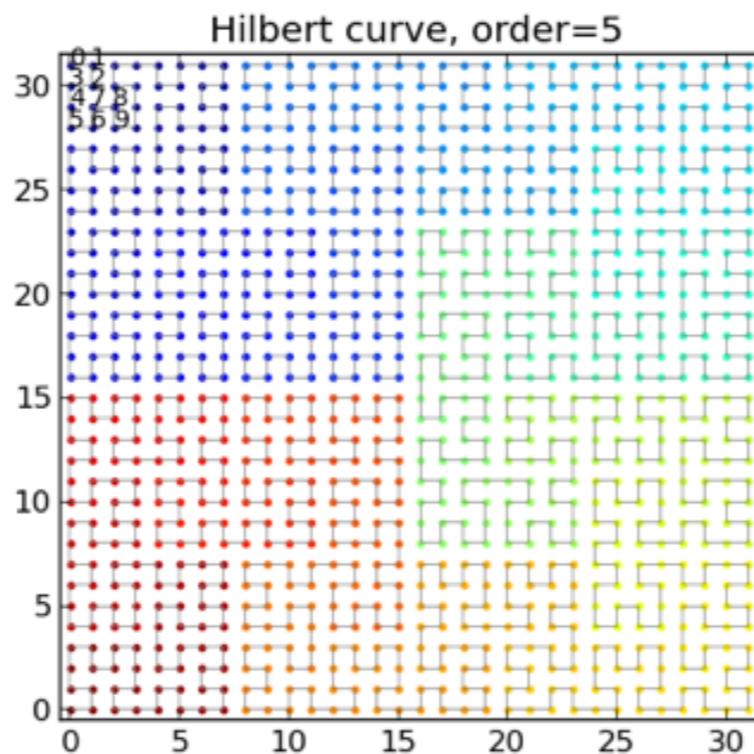
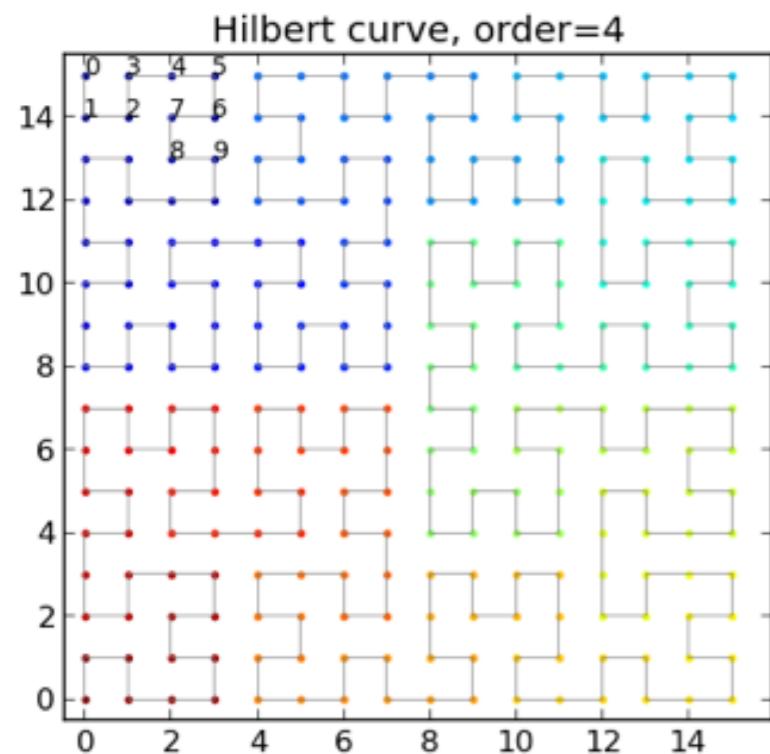
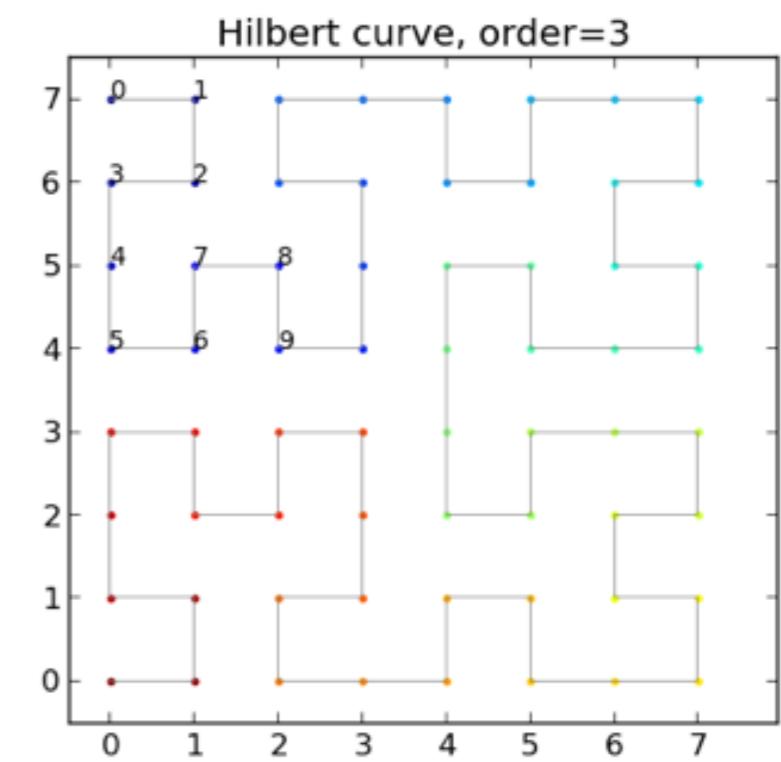
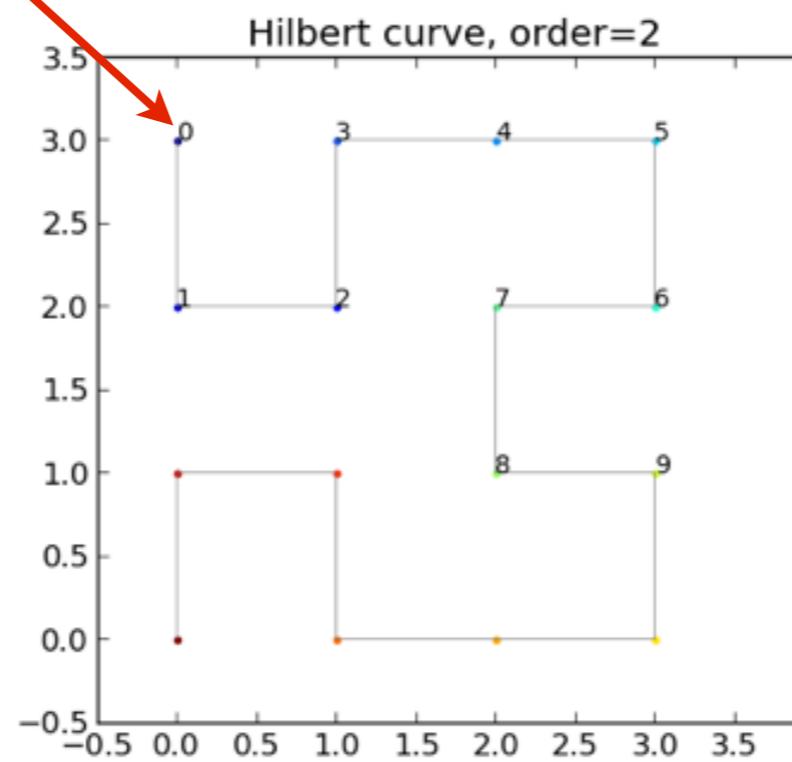
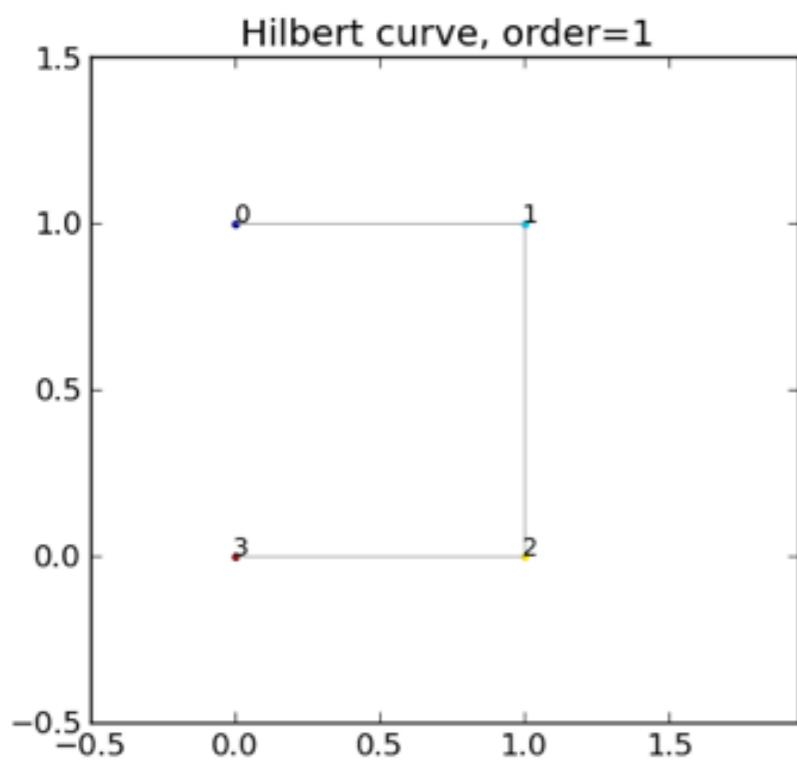
- The human genome is enormous. 3 billion base pairs.
- Difficult to see patterns in 1D plots when the coordinate space is so gargantuan.



Origin (start of chrom)



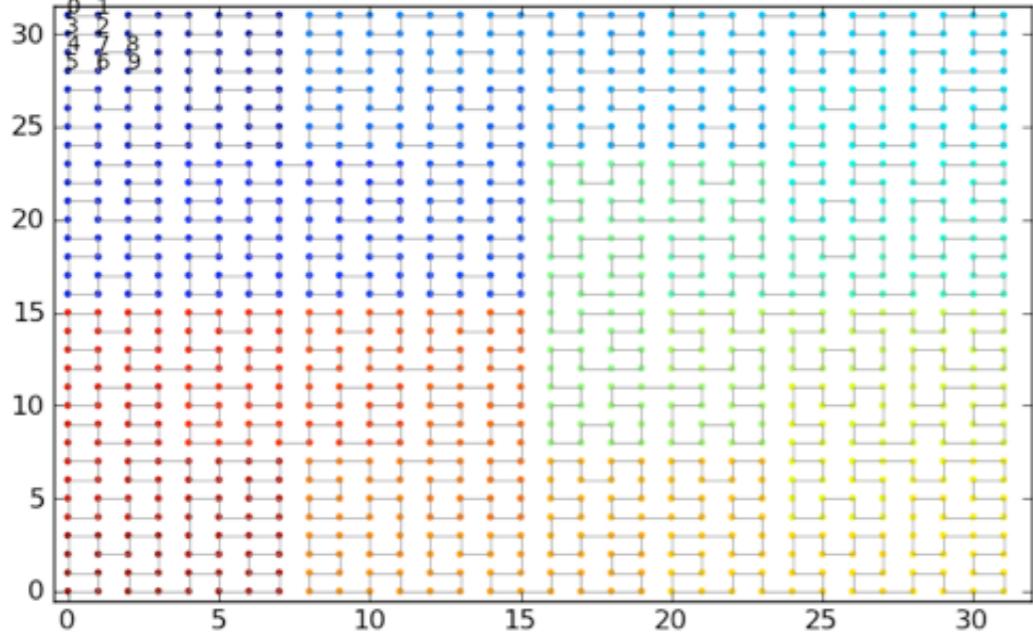
Origin (start of chrom)



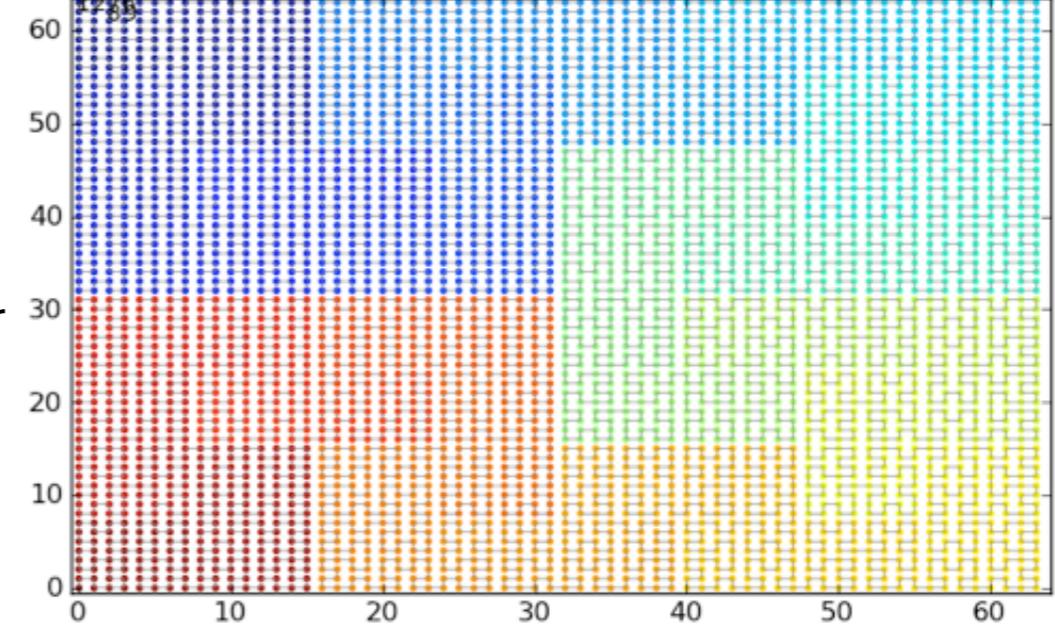
The higher the “order” the higher the resolution - that is,
each cell represents less genome space

Human chromosome 1: 249,250,621 nucleotides

order=5;
 2^5 ($32 \times 32 = 1024$) cells
each cell represents a bin of 243,408bp

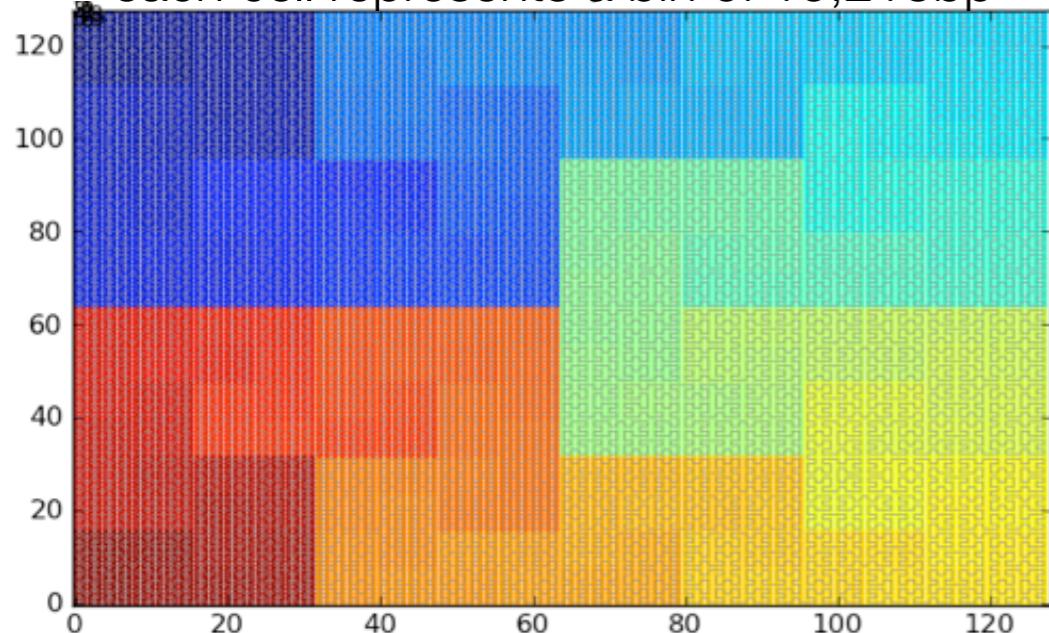


order=6;
 2^6 ($64 \times 64 = 4096$) cells
each cell represents a bin of 60,852bp



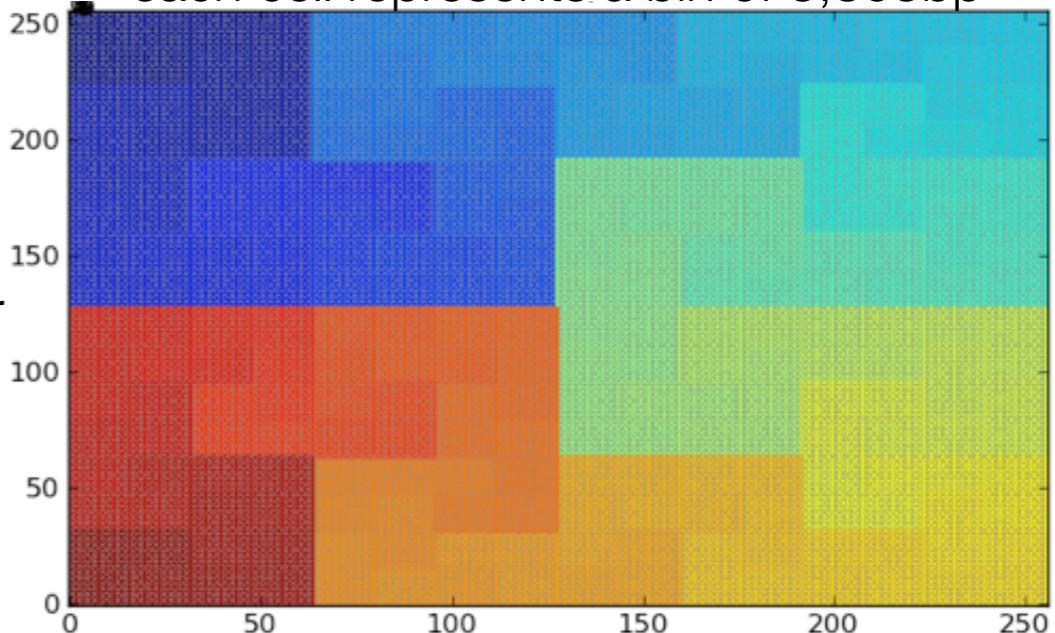
→
4x higher
res.

order=7;
 2^7 ($128 \times 128 = 16,384$) cells
each cell represents a bin of 15,213bp

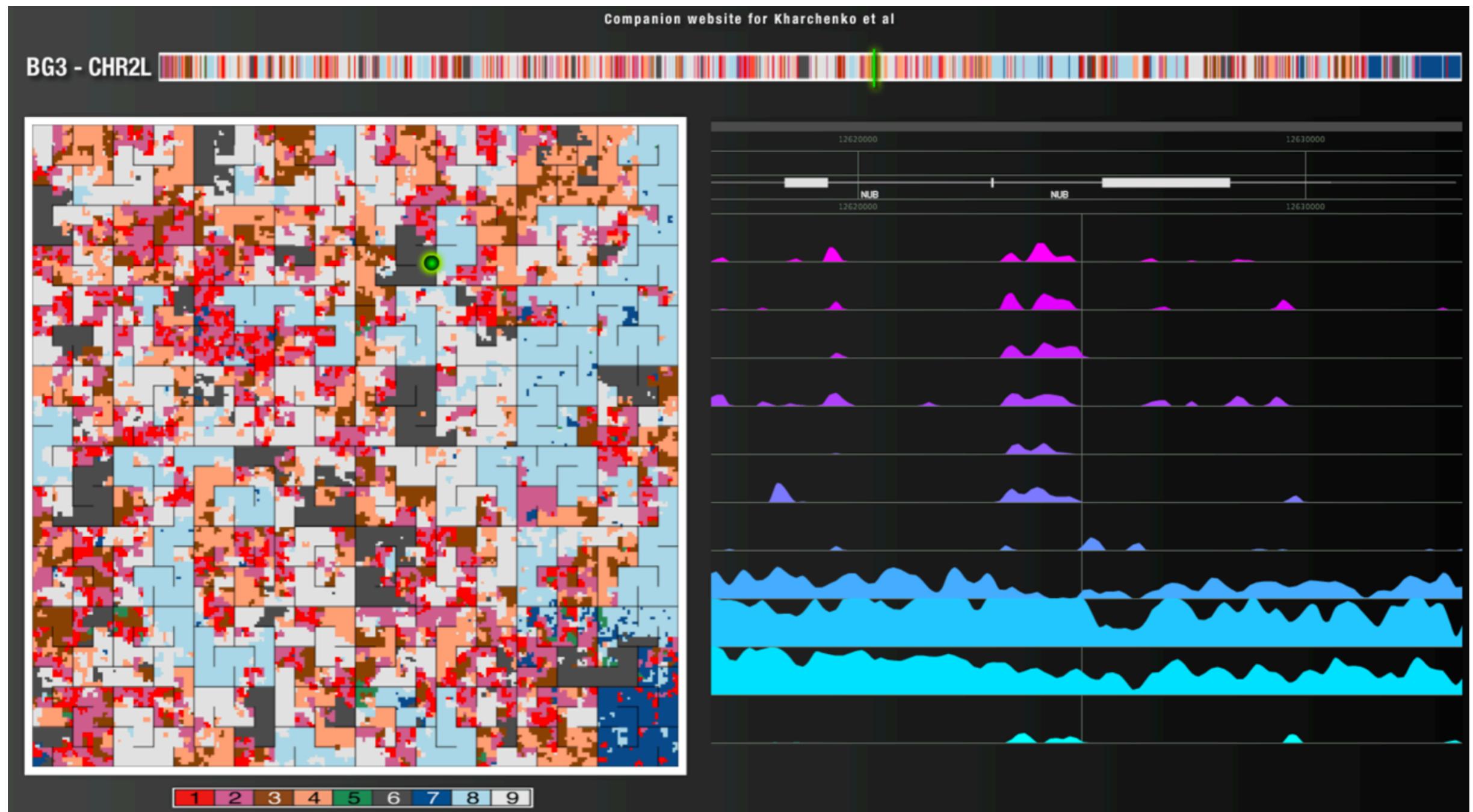


→
4x higher
res.

order=8;
 2^8 ($256 \times 256 = 65,536$) cells
each cell represents a bin of 3,803bp



modENCODE browser

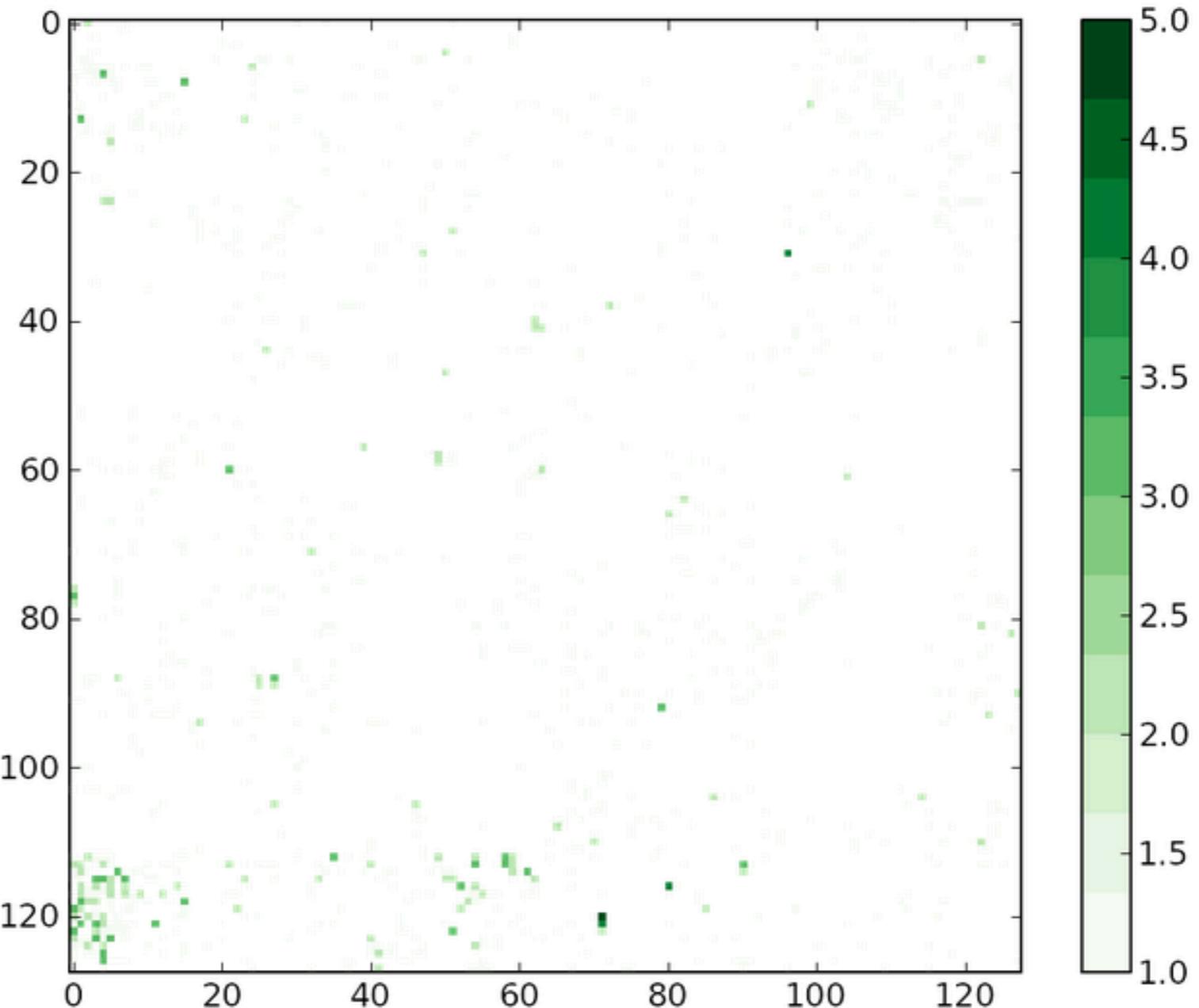


Very powerful, but an artistic toy.
Not reusable. Our motivation.

<http://compbio.med.harvard.edu/flychromatin/>

scurgen: Hilbert visualization toolkit for genomics

```
scurgen plot --chrom chr10 \  
    --cmap Greens \  
    --format png \  
    --dim 128 \  
    data/cpg-islands.hg19.chr10.bed
```



Ryan Layer

Graduate student

rl6sf @ virginia.edu

Research Interests: Sce
analysis; genome data mi
interpretation.



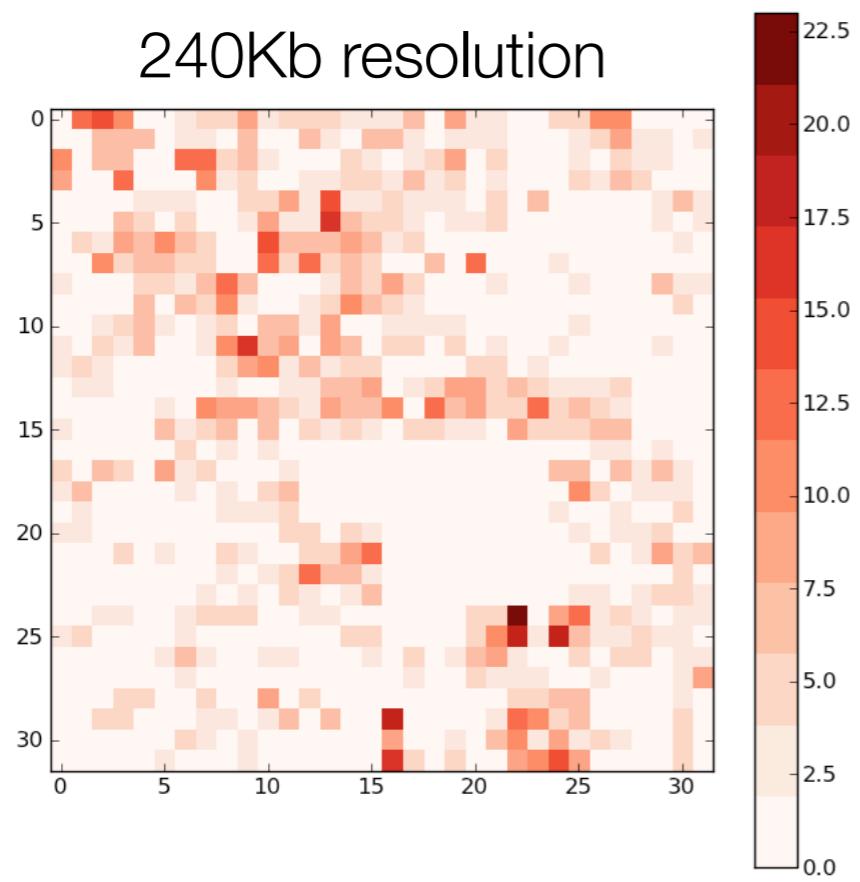
Ryan Dale

Algorithm for plotting genomic data

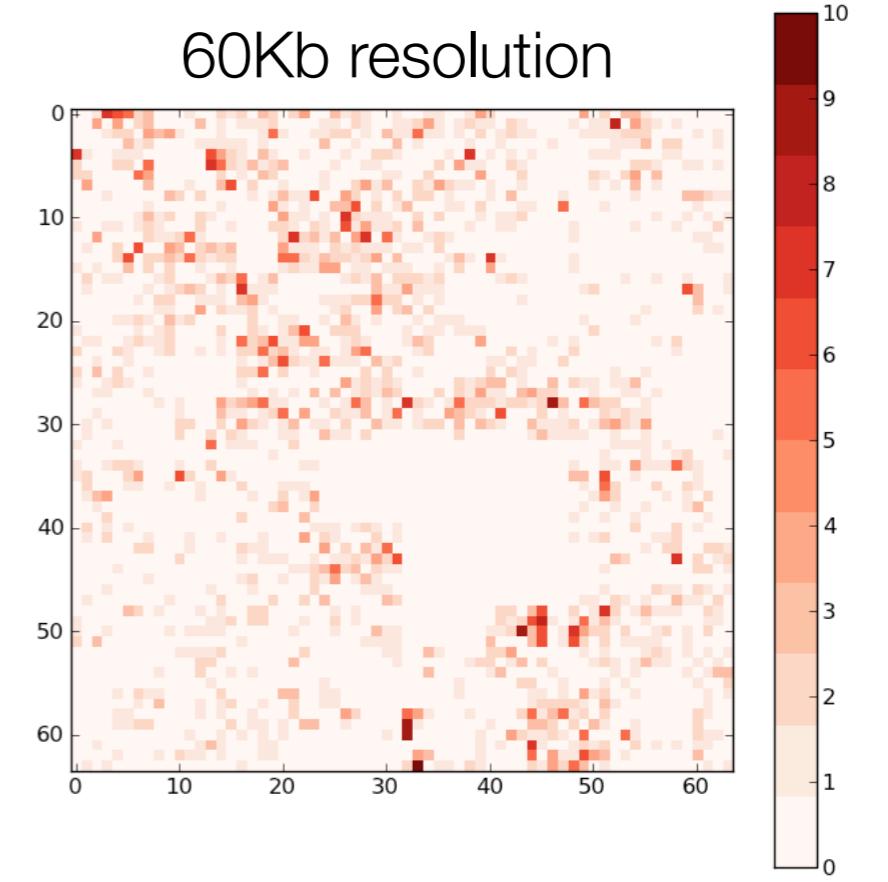
- Determine the size of the curve.
 - e.g. if 32x32 matrix, 1024 cells. curve length = 1024
- Determine the genomic size of each cell.
 - if 1024 cells and chrom length is 1048576 bp, then each cell represents 1024 bp.
- For each chrom interval in input file (BED, BAM, VCF, etc.)
 - figure out how at which point(s) along the curve (which cell) the interval belongs. **existing algorithms for this.**
 - update the cell. if BED, ++. if BEDGRAPH, += score.
- Plot matrix

Human chromosome 1: ENCODE ChIPseq of Nrf1

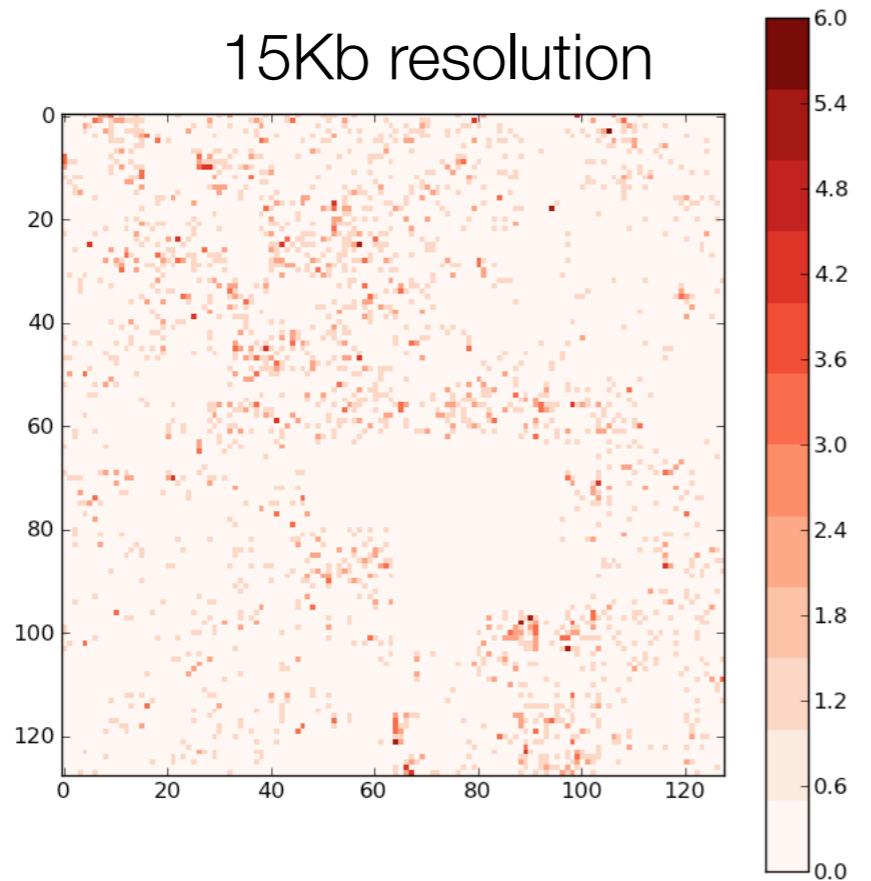
240Kb resolution



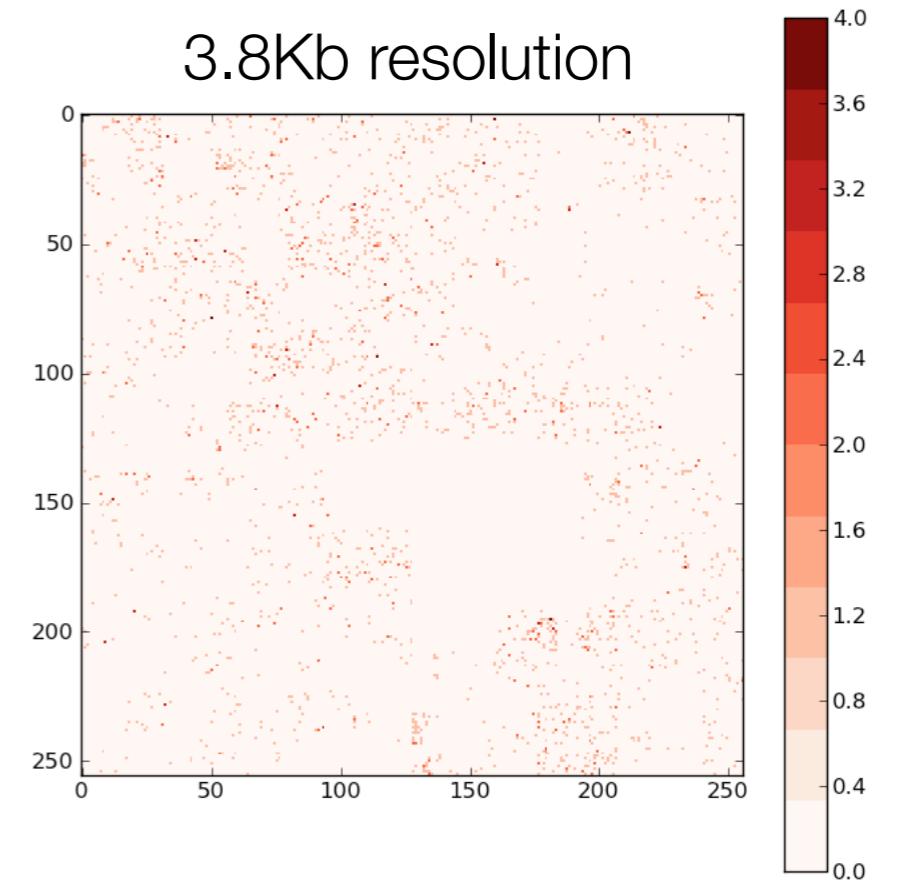
60Kb resolution



15Kb resolution



3.8Kb resolution

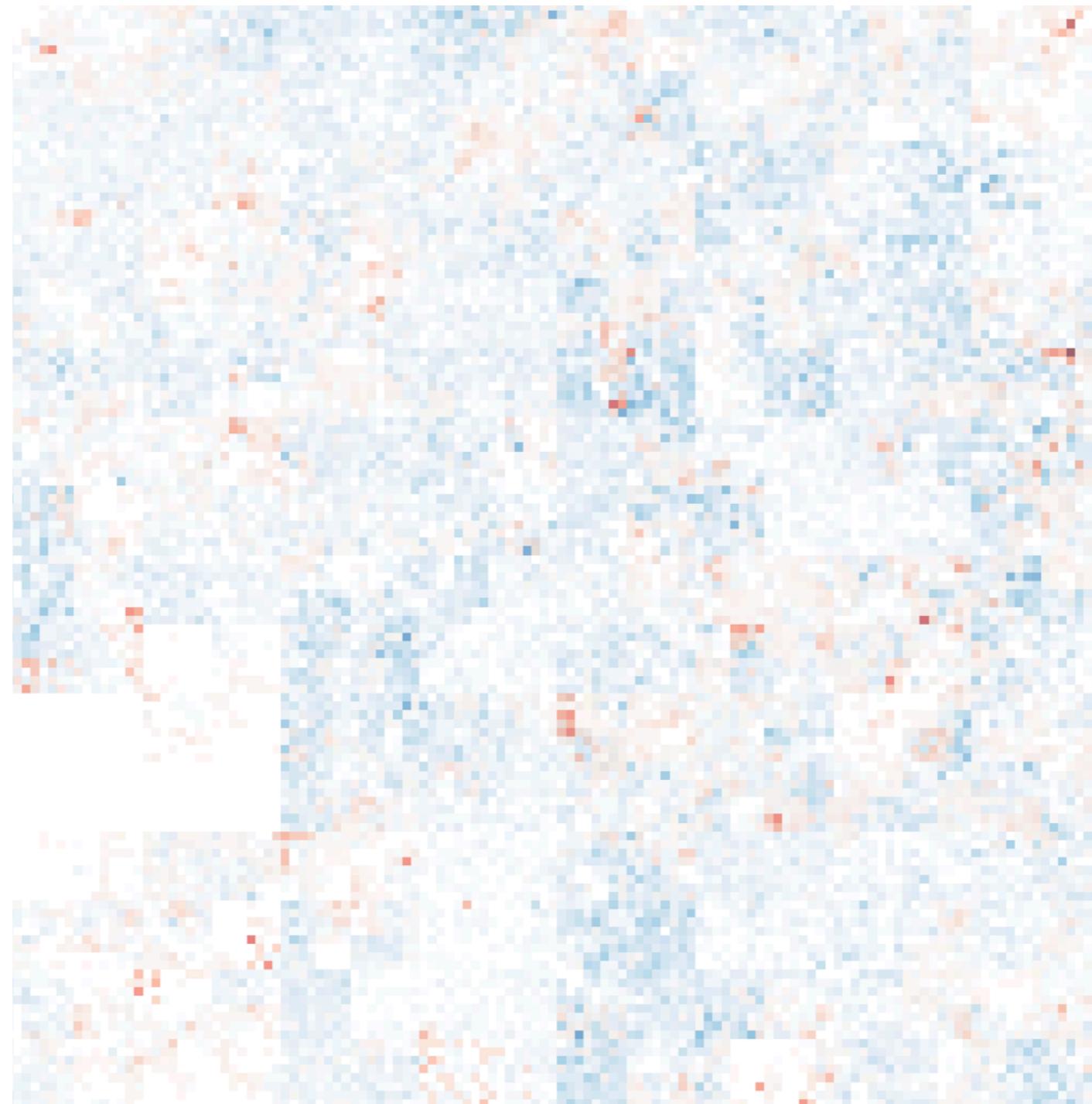


One dataset is okay. Ideal for comparisons

chr10

conserved elements: blue

coding exons: red



each cell is 8kb

centromere →

What does this tell us (that we already know)?

Digging deeper into genomic “relationships”

PREPRINT

Vol. 00 no. 00 2012
Pages 1–8

Binary Interval Search (BITS):

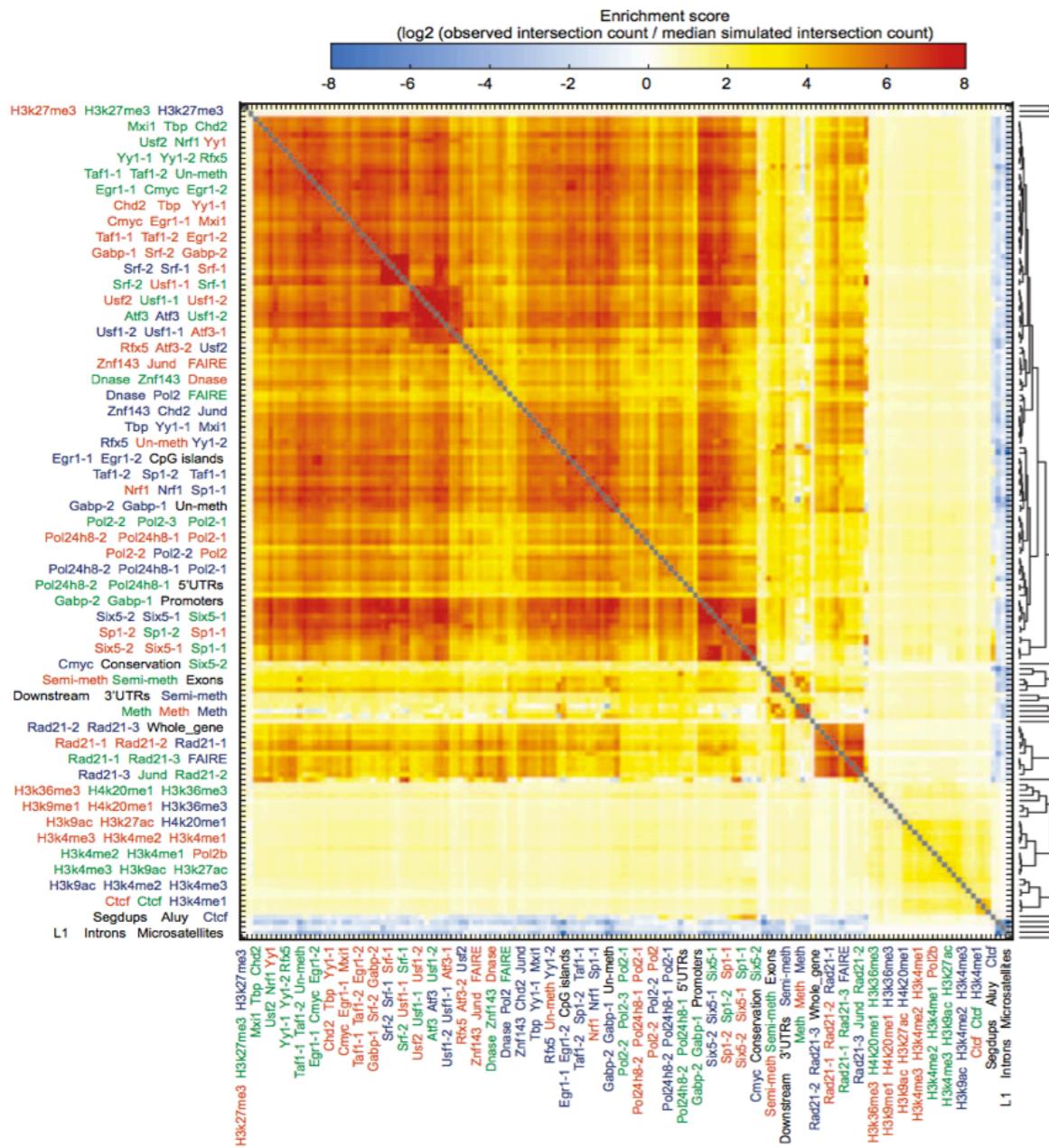
A Scalable Algorithm for Counting Interval Intersections

Ryan M. Layer¹, Kevin Skadron¹, Gabriel Robins¹, Ira M. Hall², and Aaron R. Quinlan^{3*}

¹Department of Computer Science, University of Virginia, Charlottesville, VA

²Department of Biochemistry and Molecular Genetics, University of Virginia, Charlottesville, VA

³Department of Public Health Sciences and Center for Public Health Genomics, University of Virginia, Charlottesville, VA

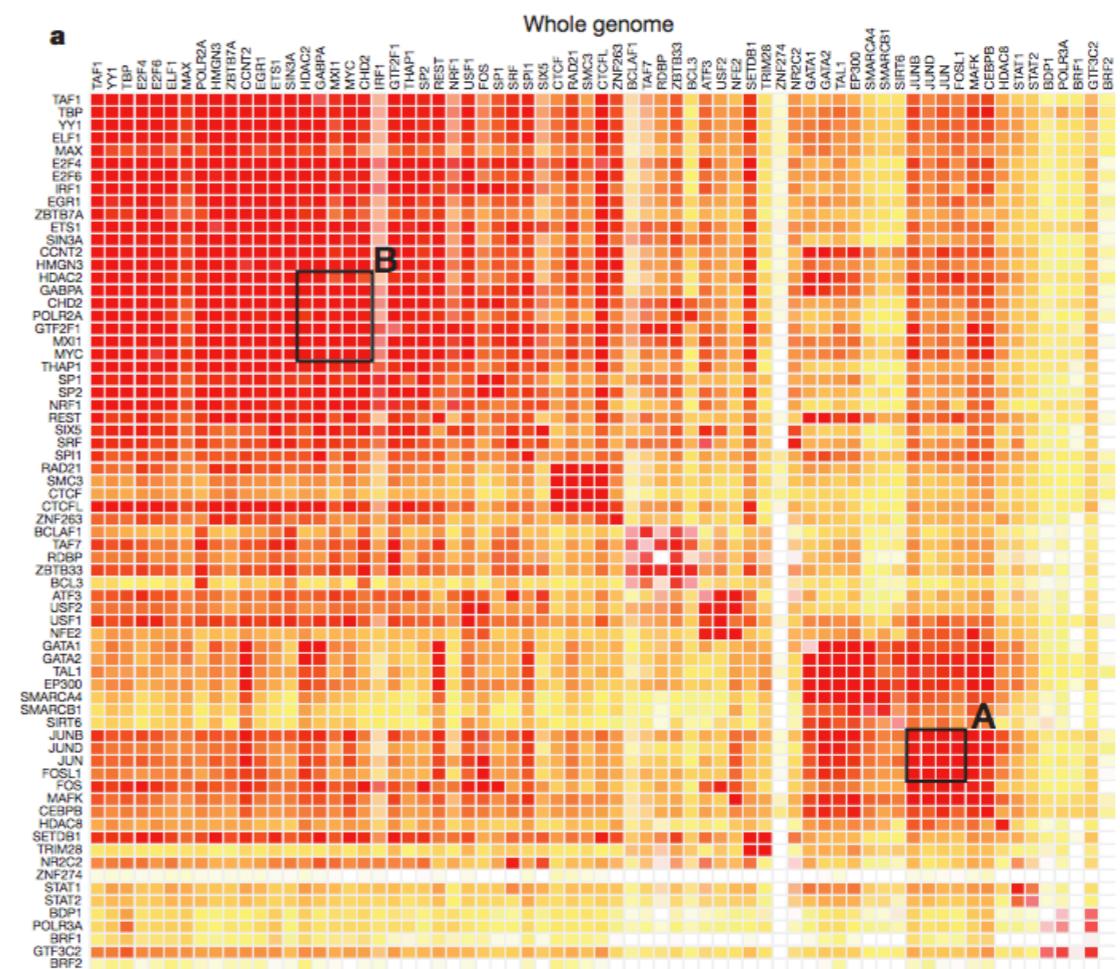


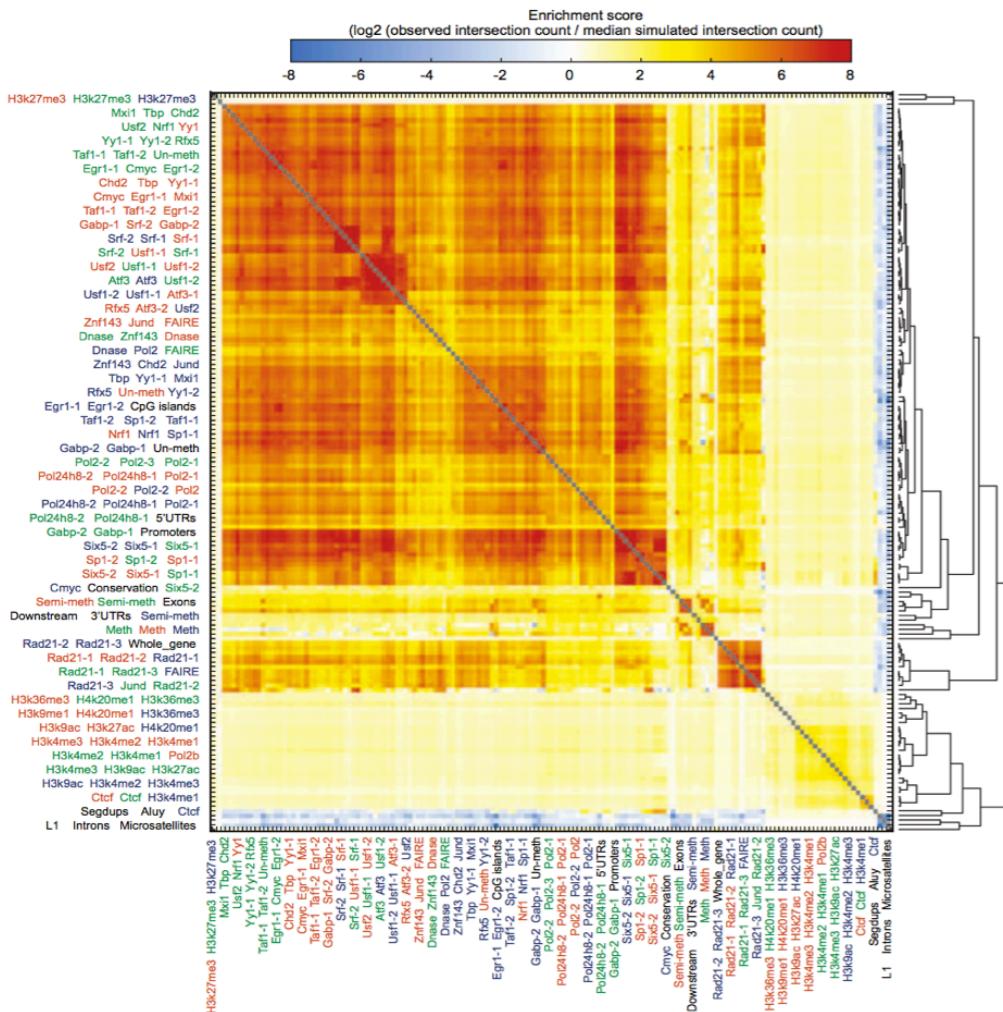
ARTICLE

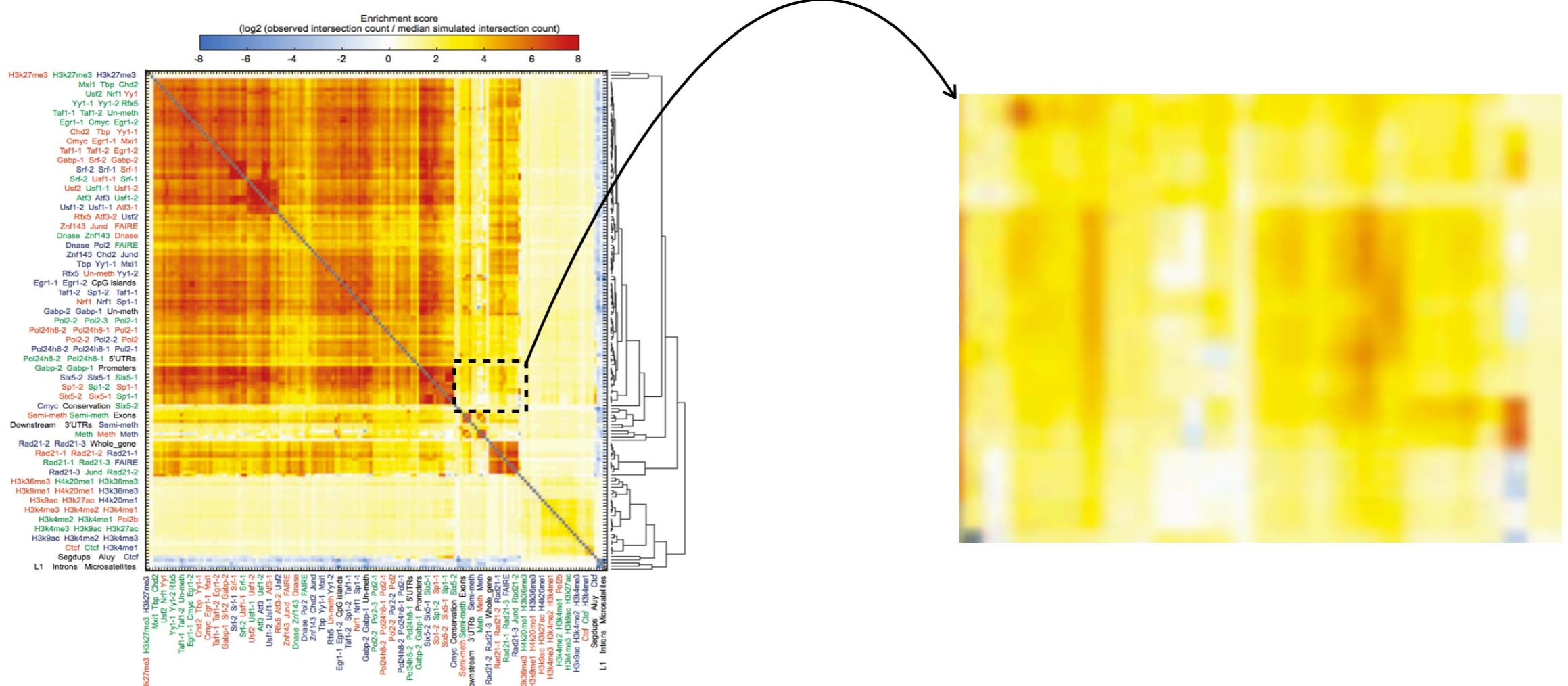
doi:10.1038/nature11247

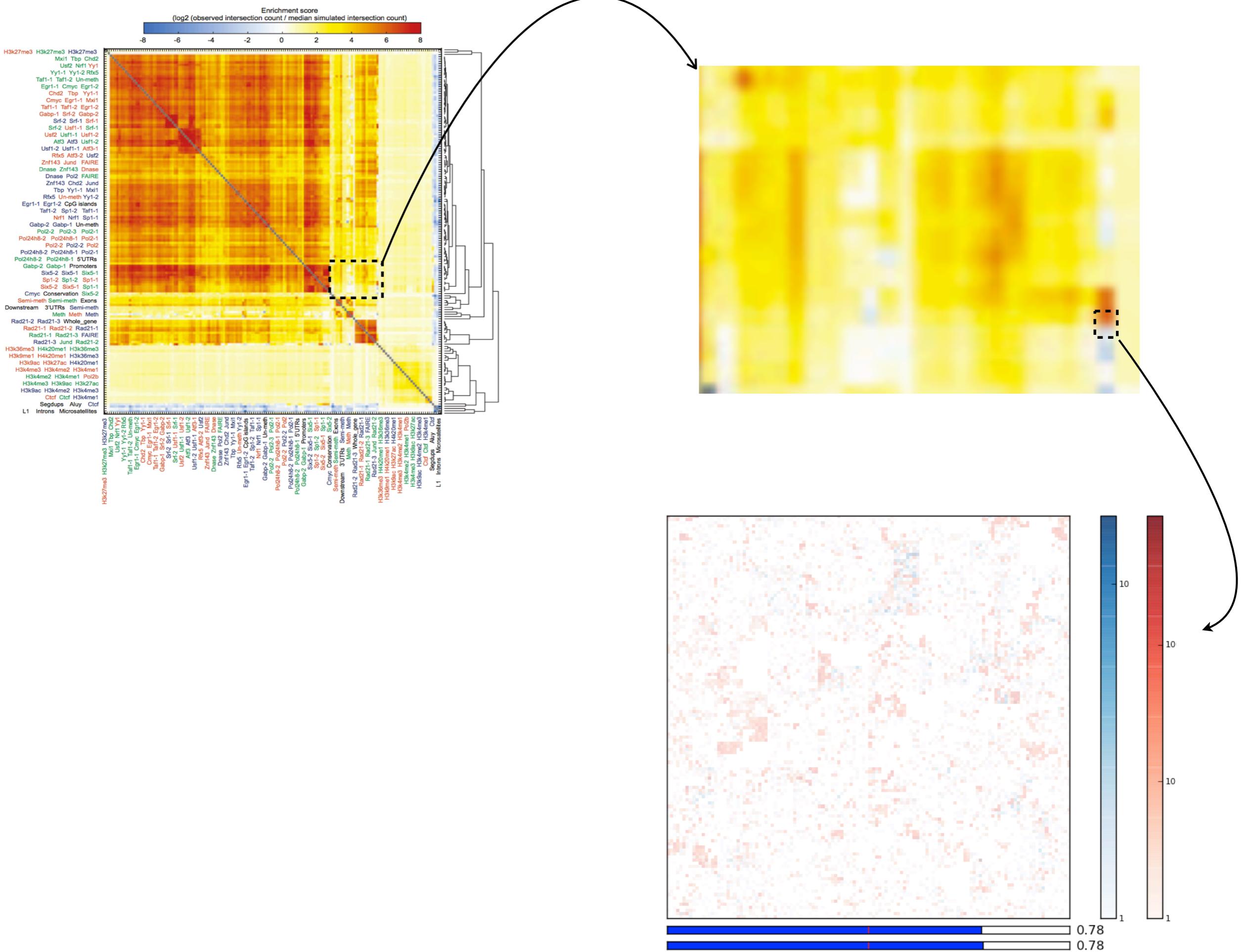
An integrated encyclopedia of DNA elements in the human genome

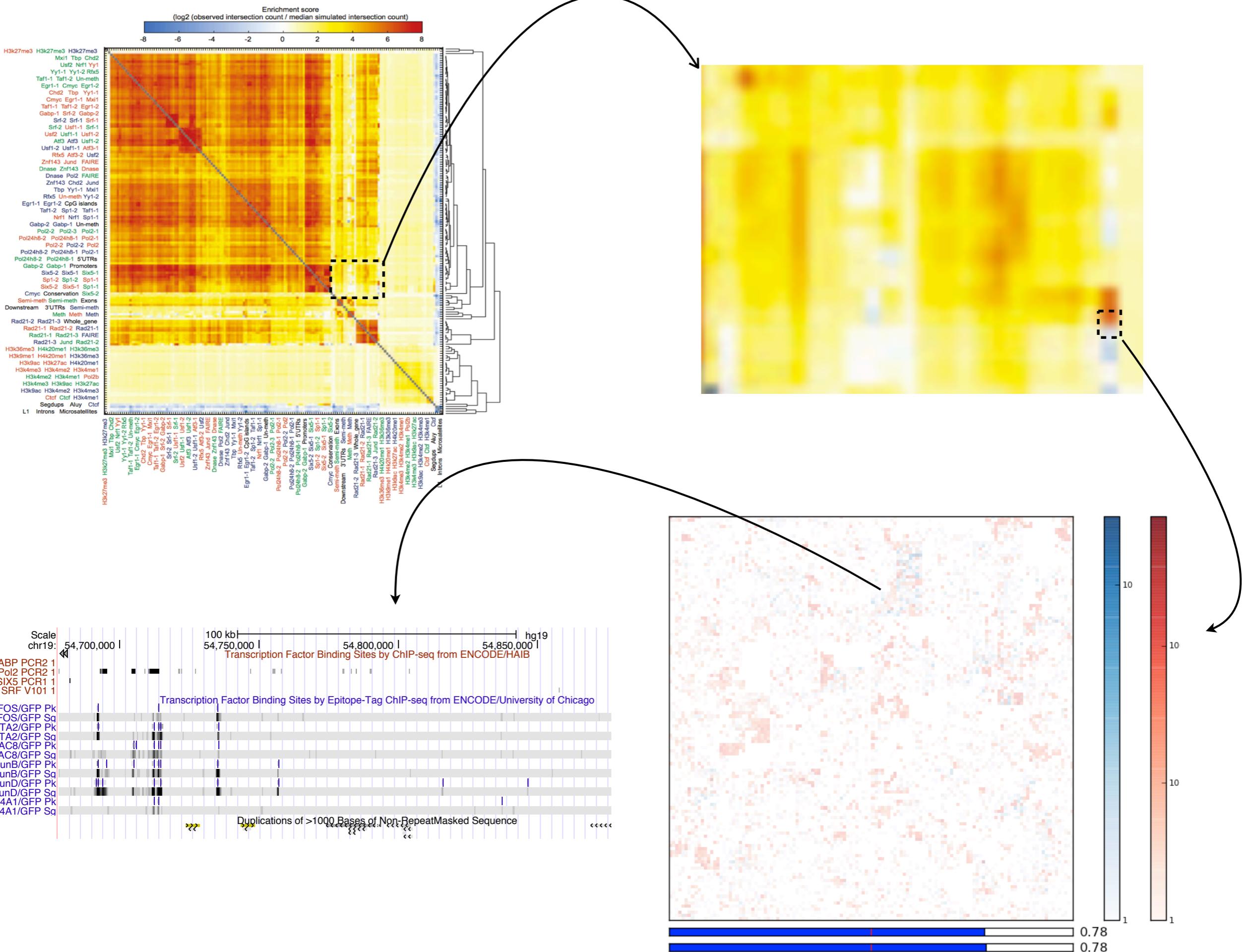
The ENCODE Project Consortium*



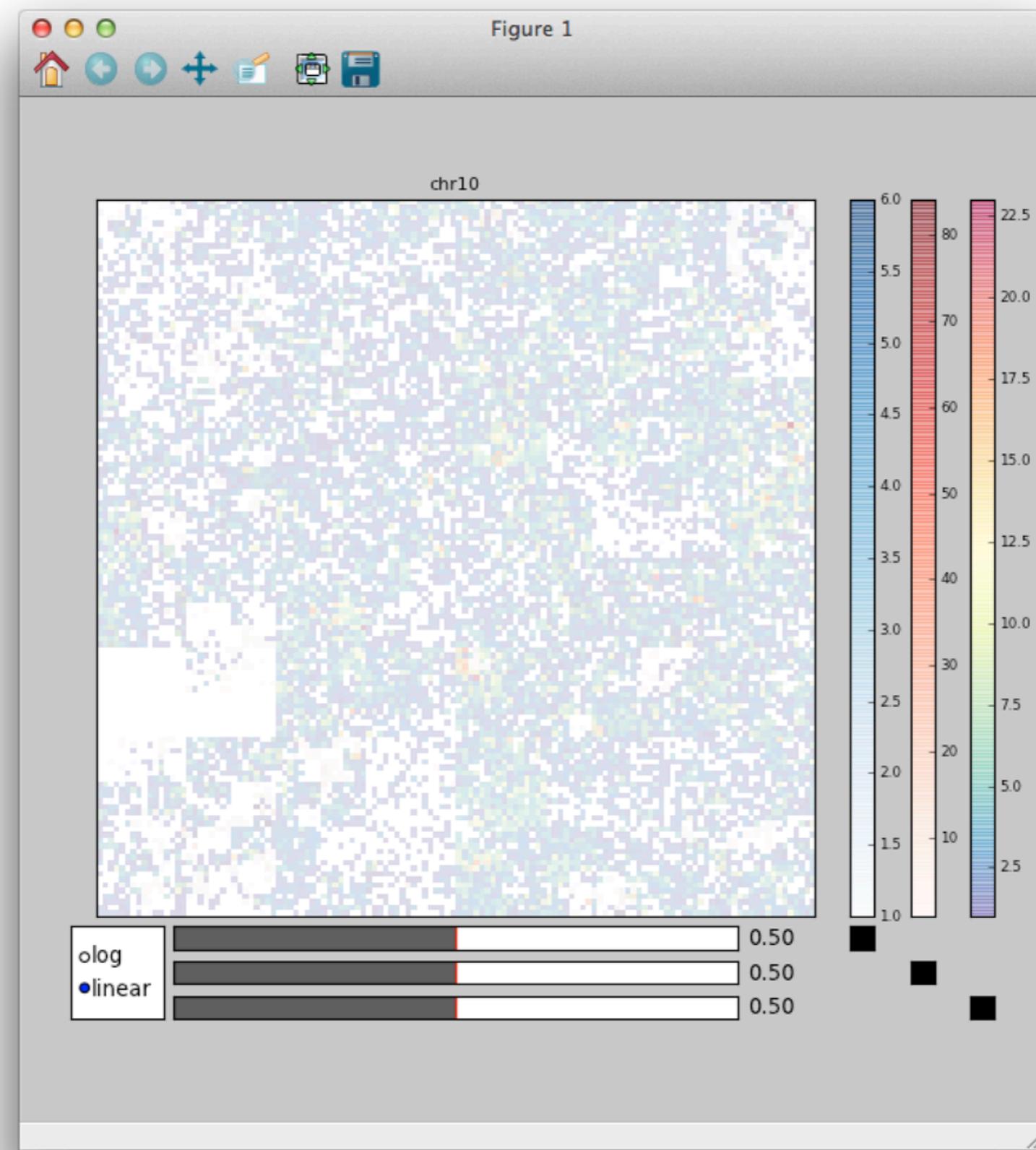








scurgen GUI



Ideal for comparative analyses:

- Spatial genomic relationships among multiple experiments: ChIP-seq, RNA-seq
- Spatial and quantitative
- Before and after comparisons
- Comparing replicates
- Unbiased data mining

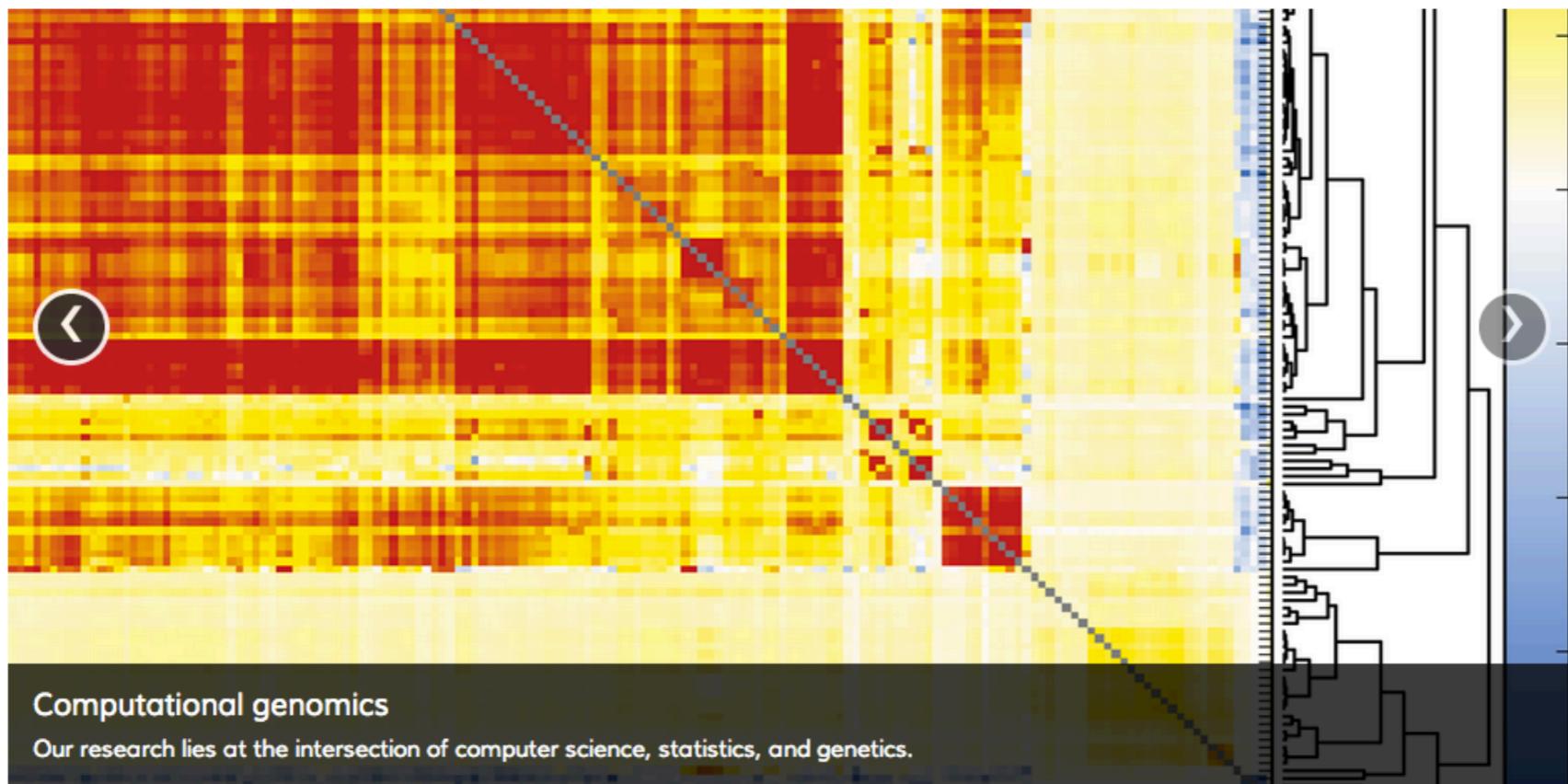
Acknowledgements

Ira Hall	<i>Univ. of Virginia</i>	Pat Concannon	<i>Univ. of Virginia</i>
Ankit Malhotra	<i>Univ. of Virginia</i>	Steve Rich	<i>Univ. of Virginia</i>
Michael Lindberg	<i>Univ. of Virginia</i>	Suna Onengut-Gumuscu	<i>Univ. of Virginia</i>
Royden Clark	<i>Univ. of Virginia</i>	Chris Moskaluk	<i>Univ. of Virginia</i>
Svetlana Sokolova	<i>Univ. of Virginia</i>	Shu-Man Fu	<i>Univ. of Virginia</i>
Mitchell Leibowitz	<i>Univ. of Virginia</i>	Gabor Marth	<i>Boston College</i>
		Jim Robinson	<i>Broad Institute</i>
		James Taylor	<i>Emory</i>
Nik Krumm	<i>Univ. of Washington</i>	Nick Navin	<i>MD Anderson</i>
Evan Eichler	<i>Univ. of Washington</i>	Kristin Baldwin	<i>Scripps</i>
Debbie Nickerson	<i>Univ. of Washington</i>		
Chris Carlson	<i>Fred Hutchinson CRC</i>		
Mark Rieder	<i>Univ. of Washington</i>		
Josh Smith	<i>Univ. of Washington</i>		
Peter Sudmant	<i>Univ. of Washington</i>		

quinlanlab.org

Quinlan Lab @UVA

Home Research Publications Software Teaching People Contact Blog



The Quinlan Lab at UVa.

We are a new computational genomics group in the [Center for Public Health Genomics](#) at the [University of Virginia](#). Our research marries genetics and genomics techniques with computer science, machine learning, and engineering to develop new ways of gaining insight into genome biology and the genetic basis of traits. We try to tackle problems with practical importance to understanding genome variation, chromosome evolution and mining genetic variation for improved understanding of human disease. Understanding the genome is a hard problem; we try to develop new approaches to make genomic research easier.

[Learn more...](#)

We are grateful to be funded by:



Mining the genome.

Aaron R. Quinlan
quinlanlab.org



University of Virginia, Charlottesville VA
Center for Public Health Genomics
Biochemistry and Molecular Genetics

Spark.



BC Cancer Agency

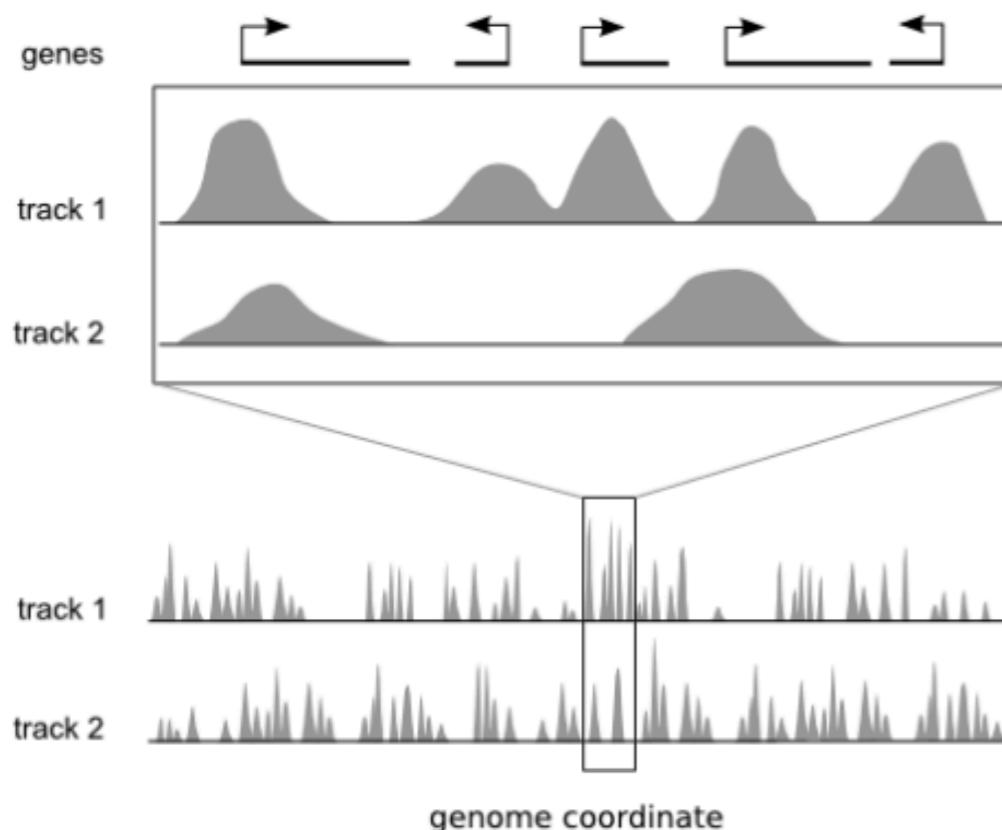
CARE + RESEARCH

An agency of the Provincial Health Services Authority



CANADA'S MICHAEL SMITH
GENOME
SCIENCES
C E N T R E

Motivation

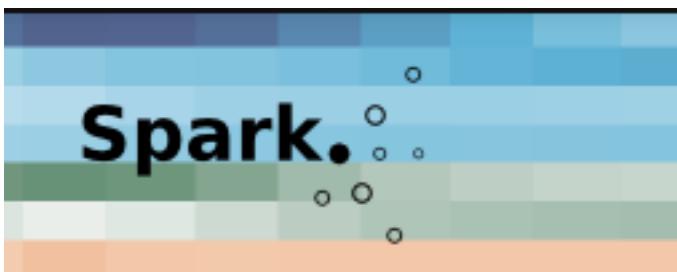


Genome browsers are ideal for viewing local regions of interest.

Many genomics techniques produce measurements that have both a value and a position on a reference genome, for example [ChIP-sequencing](#). Popular genome browsers arrange the linear genome coordinate along the x axis and express the data value range on the y axis. This approach enables integration of diverse data sets by plotting them as vertically stacked tracks across a common genomic x axis. Genome browsers are designed for viewing local regions of interest (e.g. an individual gene) and are frequently used during the initial data inspection and exploration phases.

But they do not provide a global overview of these regions.

Most genome browsers support zooming along the genome coordinate. This type of overview is not always useful because it produces a summary across a continuous genomic range (e.g. chromosome 1) and not across the subset of regions that are of interest (e.g. genes on chromosome 1). There is a need for tools that help answer questions like: "What are the common data patterns across genes start sites in my data set?"

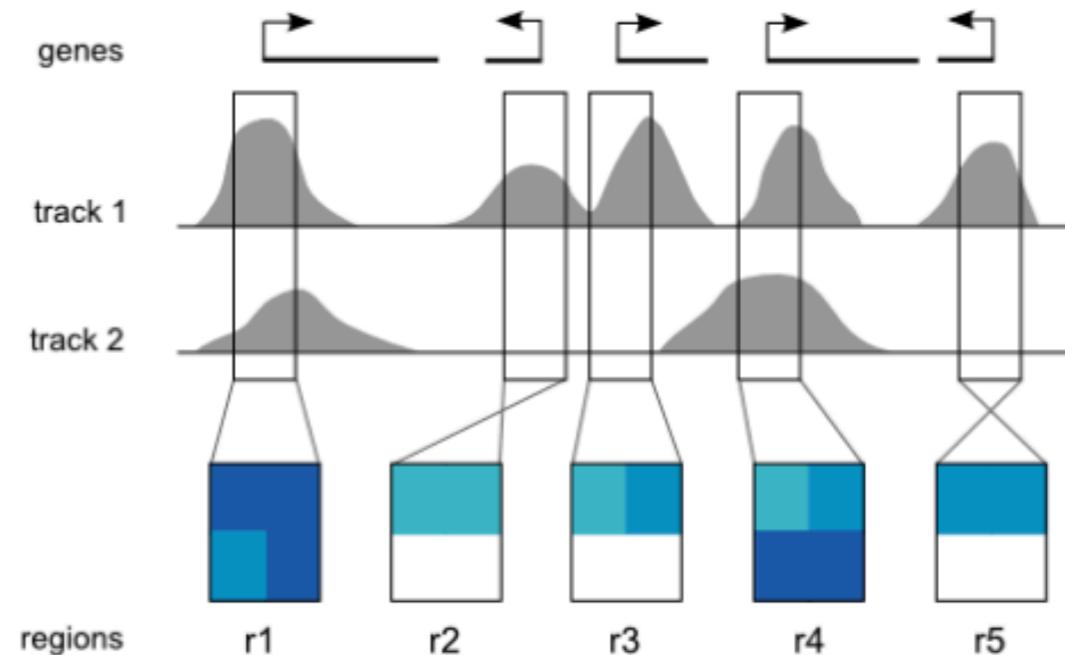


Spark.

Step 1: Pre-processing

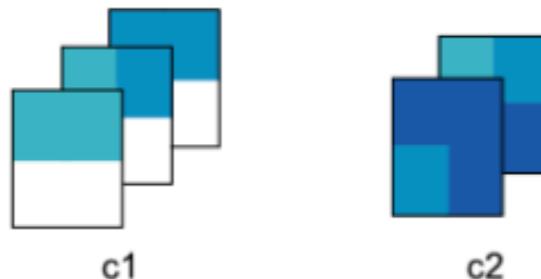
A Spark analysis begins with two user inputs: (i) one or more data tracks, and (ii) a set of regions of interest. Spark extracts a data matrix for each input region and orients it according to strand. Rows in these matrices correspond to data tracks and columns represent data bins along the genomic x-axis (two bins per region are used the diagram). The values are then normalized between 0 and 1 that can be visualized as heatmaps.

Approach



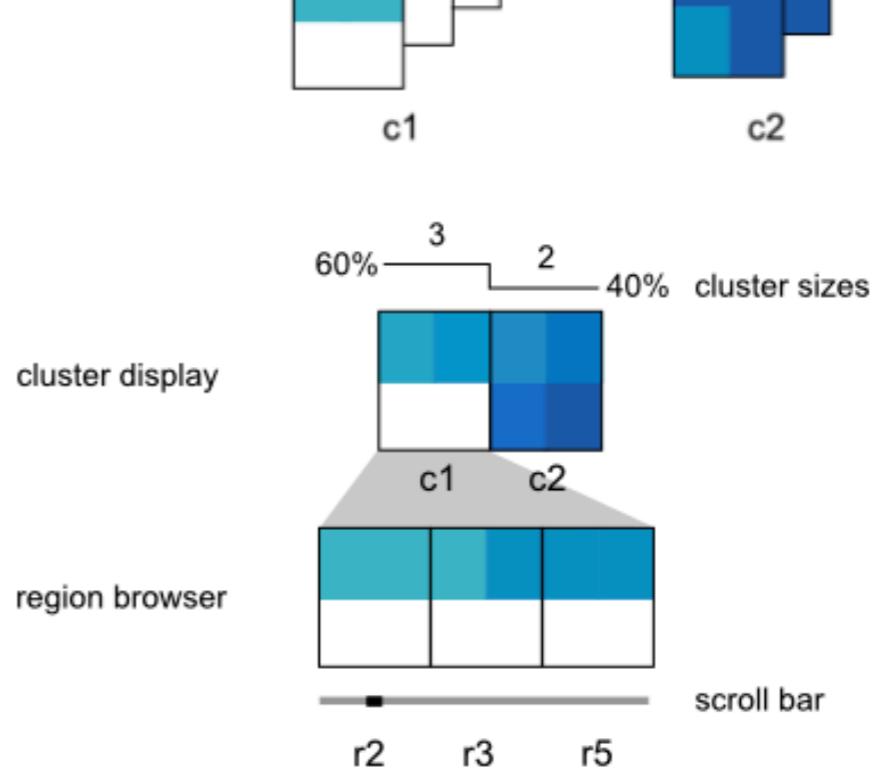
Step 2: Clustering

The preprocessed data are then clustered using k-means clustering using a user specified number of clusters (k). This technique was chosen for its effectiveness, relative simplicity, and runtime speed.



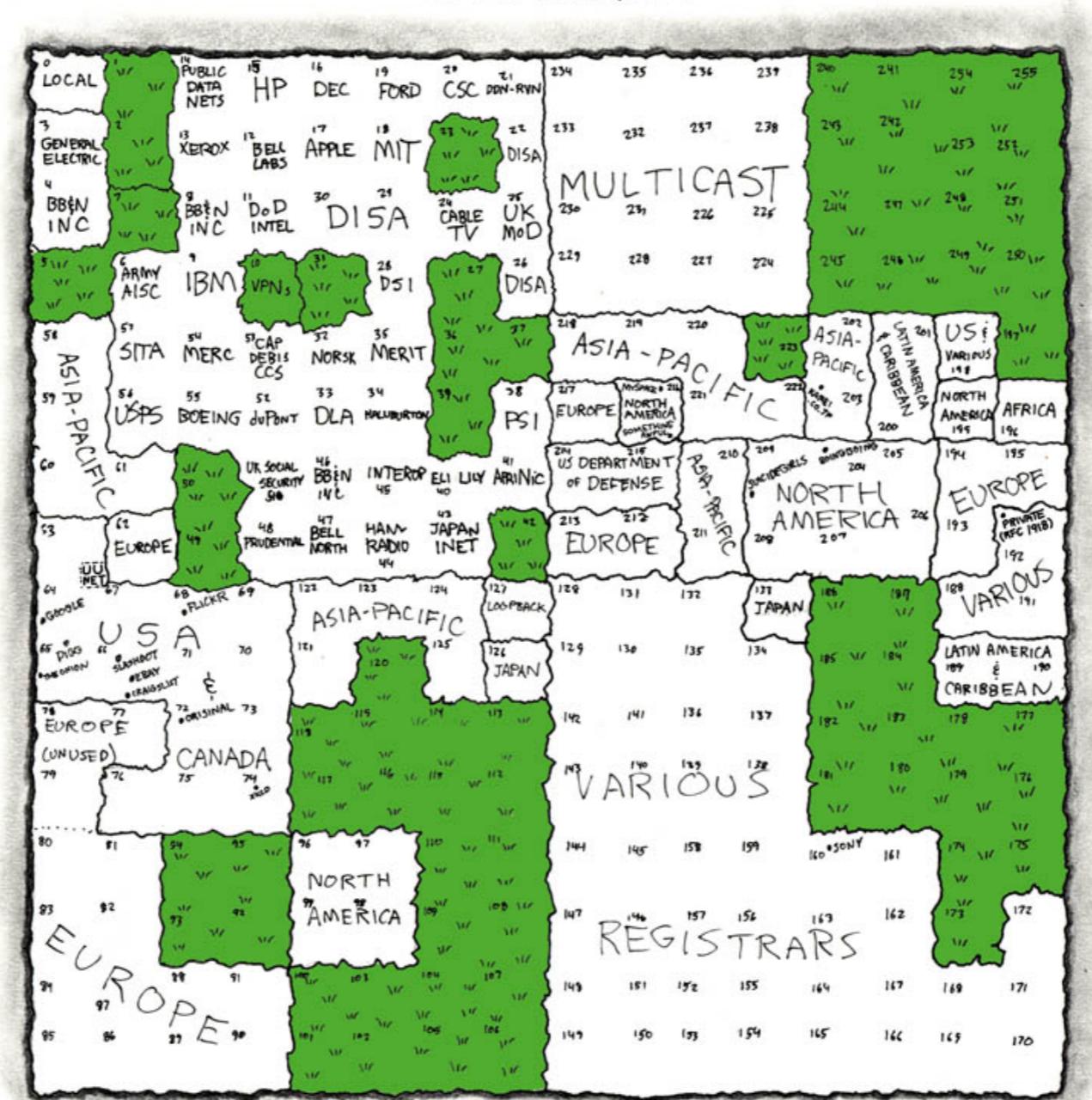
Step 3: Interactive visualization

The visualization is composed of two components: the Cluster Display, which provides a summary of each cluster, and the Region Browser, which displays individual cluster members. The user can therefore see both a high-level picture of the patterns in their data, while also being able to drill-down to individual regions of interest. The low-level region view is supported by links out to the [UCSC Genome Browser](#), and the high-level cluster view supports interactive analyses such as manual cluster refinement and links out to the [DAVID](#) gene ontology tool set.



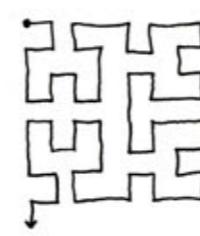
MAP OF THE INTERNET

THE IPv4 SPACE, 2006



THIS CHART SHOWS THE IP ADDRESS SPACE ON A PLANE USING A FRACTAL MAPPING WHICH PRESERVES GROUPING -- ANY CONSECUTIVE STRING OF IPs WILL TRANSLATE TO A SINGLE COMPACT, CONTIGUOUS REGION ON THE MAP. EACH OF THE 256 NUMBERED BLOCKS REPRESENTS ONE /8 SUBNET (CONTAINING ALL IPs THAT START WITH THAT NUMBER). THE UPPER LEFT SECTION SHOWS THE BLOCKS SOLD DIRECTLY TO CORPORATIONS AND GOVERNMENTS IN THE 1990's BEFORE THE RIRs TOOK OVER ALLOCATION.

0	1	14	15	16	19	→
3	2	13	12	17	18	
4	7	8	11			
5	6	9	10			



= UNALLOCATED
BLOCK