

# Introduction to Bioinformatics

Chris Miller, Ph.D.  
Washington University in St. Louis

# Bioinformatics Workshop 2023-2024

(aka bfx-workshop)



They're the same picture.

# Applied Bioinformatics for Genomics II

(aka BIOL.5625.01)

# Bioinformatics Workshop 2023-2024

---

## Supported by – ICTS Precision Health

- We aim to catalyze genomic research by providing grant review, development services, guidance and resources for genomic researchers and genomics education in the community.

Cite the **NIH CTSA Grant #UL1 TR002345** when research is supported by ICTS/CTSA funding or any ICTS Core Services

## BFX Workshop – contact Jenny if you haven't received the following

- Slack access, welcome email, Outlook bfx-workshop-2023 group invite



Register for BFX

<https://redcap.link/BFX2023>

[icts-precisionhealth.wustl.edu](http://icts-precisionhealth.wustl.edu)

[j.mckenzie@wustl.edu](mailto:j.mckenzie@wustl.edu)

# Applied Bioinformatics for Genomics II

Course: BIOL.5625.01

1 Credit Hour DBBS course

- **50% grade:** Attendance

- 75% (10 lectures) must be in person
- 3 can be viewed via recordings

- **50% grade:** Assignments

- choose 8 of the 10 assignments
- due by the end of the second Friday after the lecture

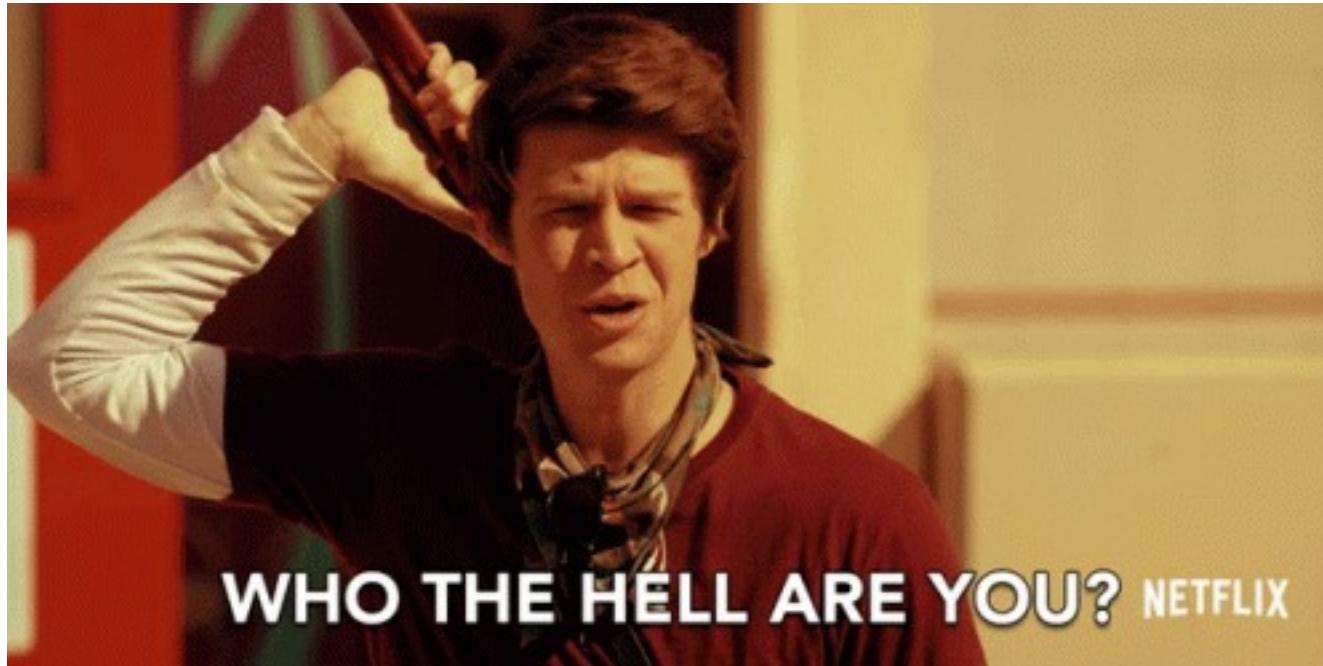


Register for BFX

<https://redcap.link/BFX2023>

[j.mckenzie@wustl.edu](mailto:j.mckenzie@wustl.edu)

[j.mckenzie@wustl.edu](mailto:j.mckenzie@wustl.edu)



**WHO THE HELL ARE YOU?**

NETFLIX

# Who we are, and why you should trust us



Chris Miller, Ph.D.

Course Director  
Associate Professor  
Division of Oncology

20 years of experience in  
Bioinformatics and Computational Biology



Jenny McKenzie, Ph.D.

Course Coordinator  
ICTS Precision Health  
Program Scientist



John Garza

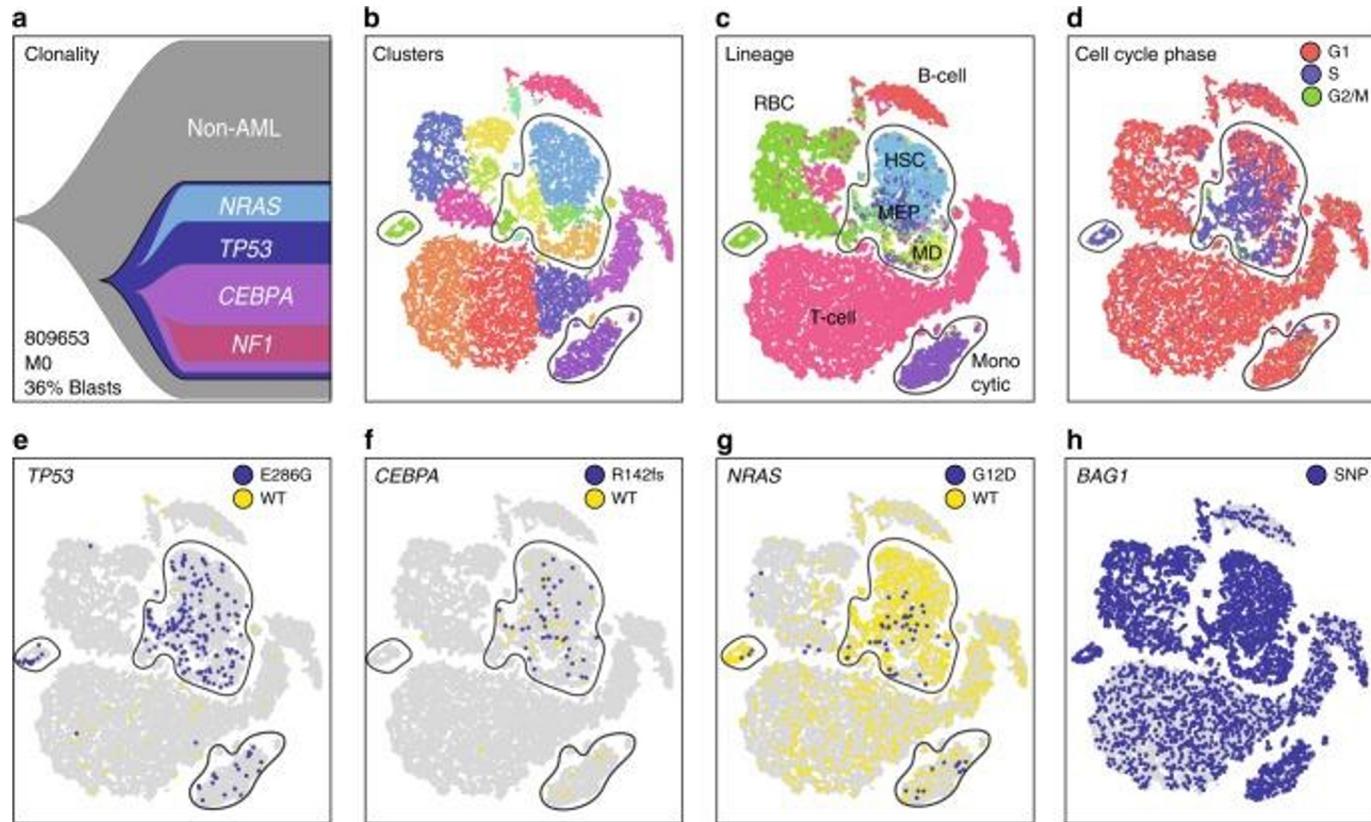
Teaching Assistant  
Bioinformatics/Genome Analytics  
Programmer

## Other Lecturers/Organizers include:

Jason Walker  
Malachi Griffith  
Jennifer Foltz  
Susanna Kiwala  
Juan Macias  
Brigida Rusconi

# Why learn bioinformatics?

- Biology is now a quantitative discipline - especially genomics



# Why learn bioinformatics?

- Biology is now a quantitative discipline - especially genomics
- Skills in programming, statistics, and visualization help you get the most out of your data



People who need complex data analysis



People who know how to do  
complex data analysis

# Why learn bioinformatics?

- Biology is now a quantitative discipline - especially genomics
- Skills in programming, statistics, and visualization help you get the most out of your data
- We're aiming to teach you the theory and practice of computational biology, with a focus on genomics but lessons that apply broadly

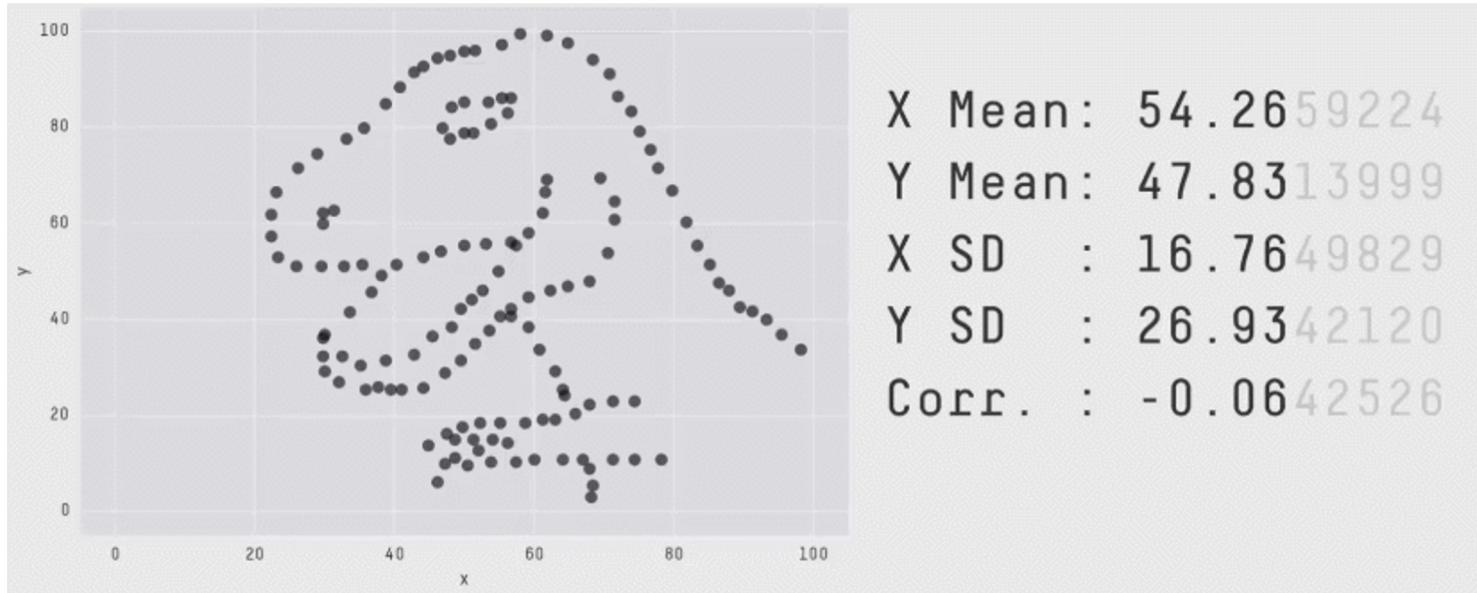
# Goals:

- To empower you to improve and expedite your research
- To expose you to new ideas and techniques that may advance your research program

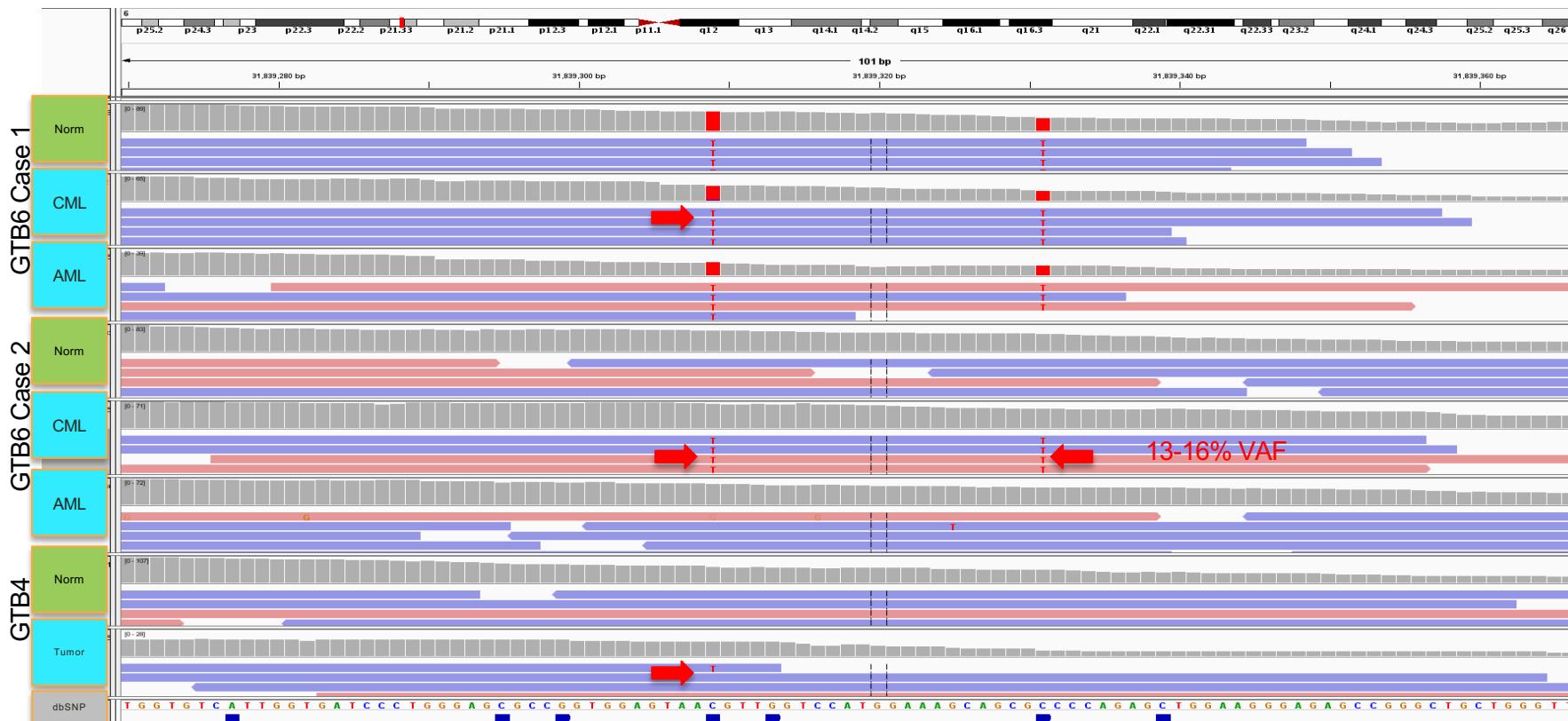
Don't trust your data

# Summary statistics are dangerous

- Visualize your data!
- A picture is worth a thousand p-values



# Contamination of CML samples



# Watch out!

- Computational analyses require controls too!
- Look at the data and understand it's limitations!
- Don't assume that the data is clean – prove to yourself that it is!

# Expectations:

- Check the prerequisites from fall week 01. Install the software, be familiar with the unix command line, know how to use docker to launch analyses
  - [https://github.com/genome/bfx-workshop/tree/master/lectures/week\\_01](https://github.com/genome/bfx-workshop/tree/master/lectures/week_01)
- Most of you are new to computational analysis – *ask questions!*
- Work hard, follow along, and get your money's worth from this course
- The folks teaching and the TAs all know their stuff, *ask questions!*

# Course Structure:

- Weekly lecture introducing topic
- Practical exercise allowing you to apply that knowledge
- <https://github.com/genome/bfx-workshop>
- ICTS Slack instance: #bfx-workshop channel
- Email list: bfx-workshop-2023@gowustl.onmicrosoft.com

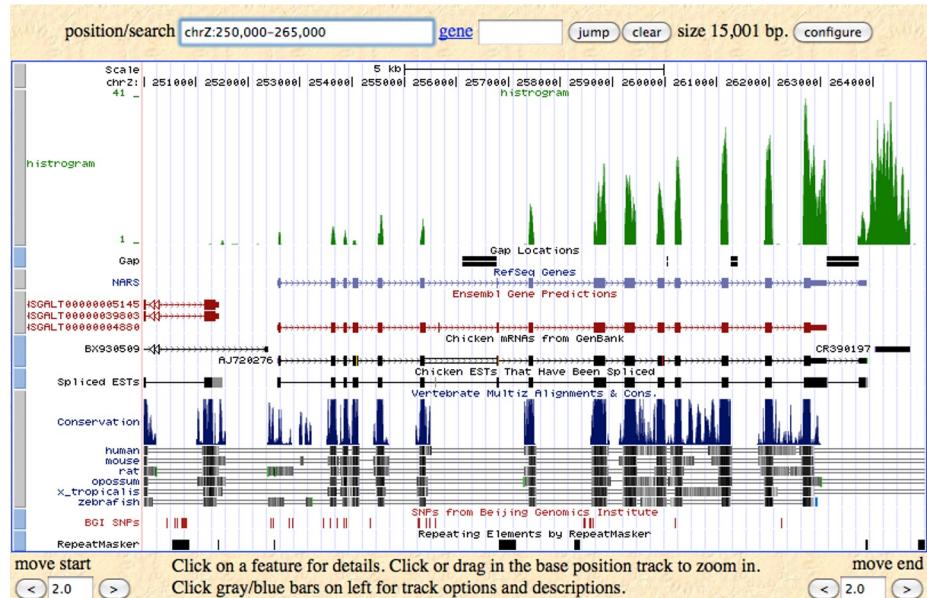
# Genome arithmetic and bedtools

some slides adapted from :  
Aaron Quinlan's Applied Computational Genomics course <https://github.com/quinlan-lab/applied-computational-genomics>  
Griffith Lab's RNAbio course: <https://rnabio.org/>

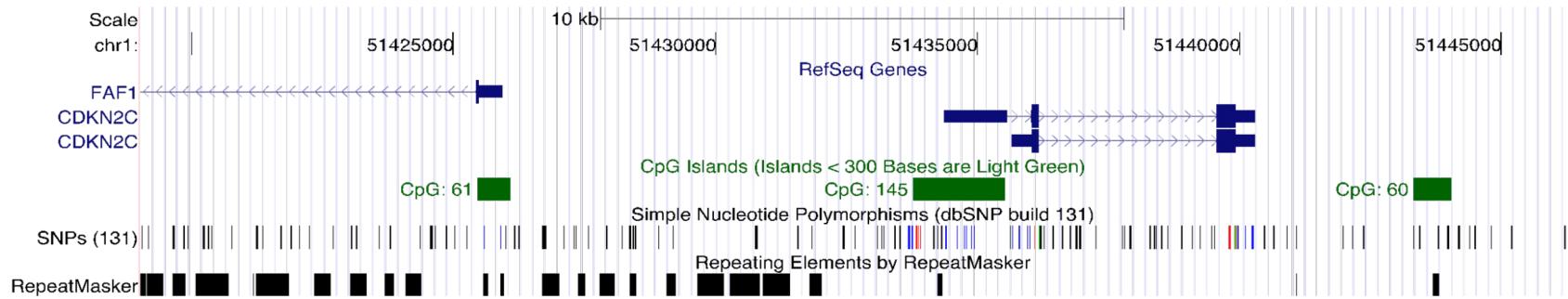
# What is a genome interval?

- Genes: exons, introns, UTRs, promoters (BED, GFF, GTF)
- Conservation (BEDGRAPH)
- Genetic variation (VCF)
- Sequence alignments (BAM)
- Transcription factor binding sites (BED, BEDGRAPH)
- CpG islands (BED)
- Segmental duplications (BED)
- Chromatin annotations (BED)
- Gene expression data (WIG, BIGWIG, BEDGRAPH)

**Your own observations: put them in context**

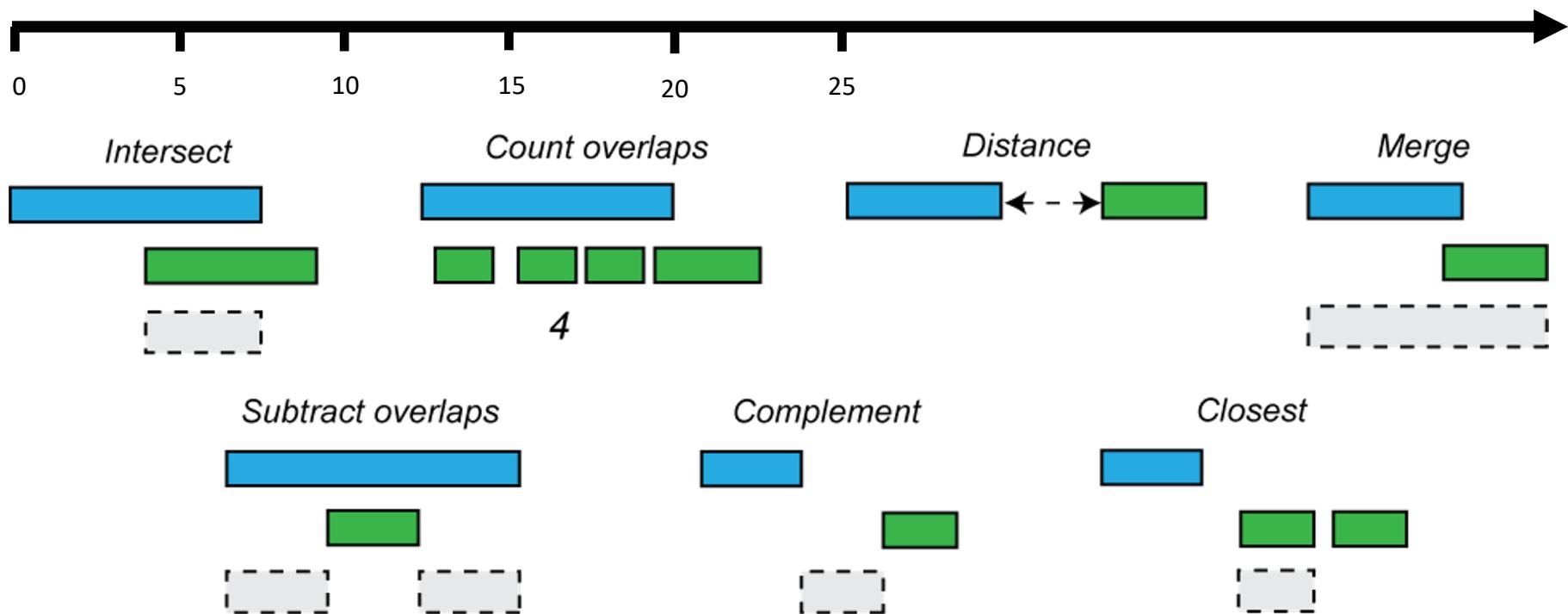


# Genome intervals

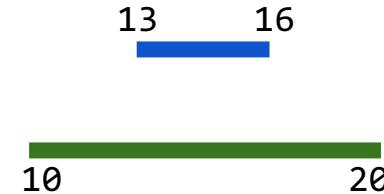
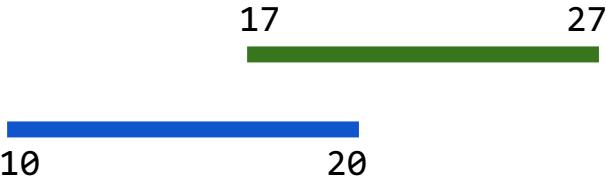


**Genome arithmetic:** the method of comparing, contrasting, and gaining insight using multiple genome interval files

# Genome arithmetic operations



# Do two intervals intersect (overlap)?

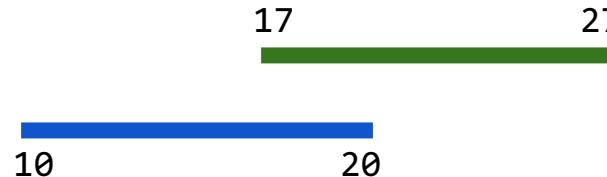


```
if ((a.start <= b.start and a.end >= b.start) or
    (b.start <= a.start and b.end >= a.start) or
    (a.start <= b.start and a.end >= b.end)    or
    (b.start <= a.start and b.end >= a.end))
{
    INTERSECTION!!!
}
else NADA!!!
```

# Do two intervals intersect (overlap)? A simpler way.



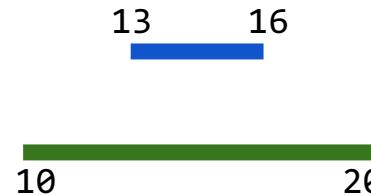
25                    27



10                    20



13                    16



10                    20

```
I = min(a.end, b.end) - max(a.start, b.start)
```

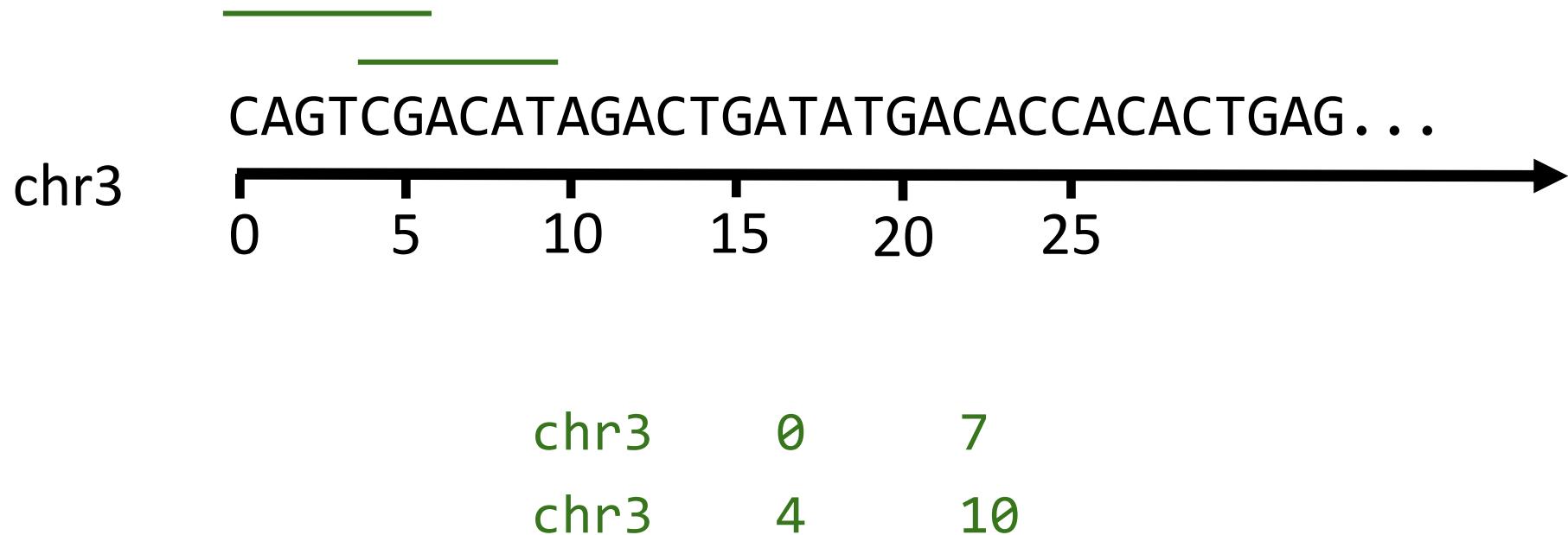
```
    if I > 0, intersection,
```

```
    if I <= 0, distance between the intervals
```

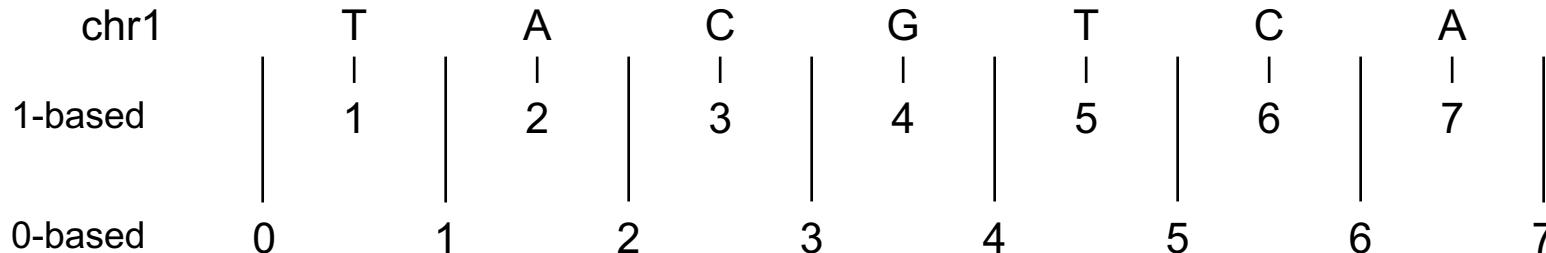
$$= \min(20, 27) - \max(10, 17)$$

$$= 20 - 17 = 3$$

# Genome arithmetic depends upon the genome coordinate system



# Genomic coordinates – 1 vs 0 based



	1-based	0-based
Indicate a single nucleotide	chr1:4-4 G	chr1:3-4 G
Indicate a range of nucleotides	chr1:2-4 ACG	chr1:1-4 ACG

- 1-based : Single nucleotides, variant positions, or ranges are specified directly by their corresponding nucleotide numbers
  - GFF, SAM, VCF, Ensembl browser, ...
- 0-based: Single nucleotides, variant positions, or ranges are specified by the coordinates that flank them
  - BED, BAM, UCSC browser, ...

# Genome builds

## Reference Genome builds

Current human: GRCh38, hg38, b38

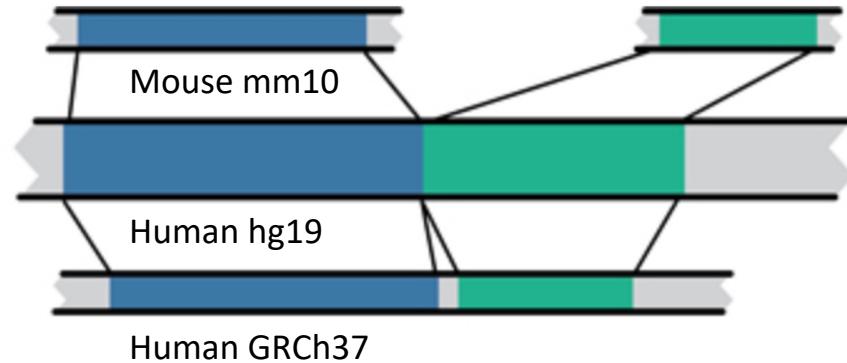
alternates: GRCh38v2\_ccdg,  
GRCh38\_full\_analysis\_set\_plus\_decoy\_hla

Previous human: GRCh37, hg19, b37

Current mouse: GRCm39

Previous mouse: GRCm38, mm10

## Lift-over



New human assembly: T2T-CHM13, pan-genomes

# Variant representation

- Does my SNP affect a \_\_\_\_\_?
- Single nucleotide events are generally easy

in pseudo-bed:

chr1	1935820	1935821	A/T
------	---------	---------	-----

or in VCF:

chr1	1935821	.	A	T
------	---------	---	---	---

# Variant shifting (alignment) and parsimony/trimming

Reference and alternative alleles of a CA short tandem repeat (STR)		REF	GGGCACACACAGGG		
		ALT	GGGCACACAGGG		
					
Genome Reference		Variant Call Format			
GGGCACACACAGGG		POS	REF	ALT	
REF	CA	8	CA	.	Not left aligned and alternate allele is empty
ALT	.				
REF	CAC	6	CAC	C	Not left aligned but parsimonious
ALT	C				
REF	GCACA	3	GCACA	GCA	Not right trimmed
ALT	GCA				
REF	GGCA	2	GGCA	GG	Not left trimmed
ALT	GG				
REF	GCA	3	GCA	G	Normalized (left aligned & parsimonious)
ALT	G				
Alleles represented against the human genome reference. Allele pairs are colored the same, all are representations of the same variant.			Alleles represented in Variant Call Format, all are representations of the same variant.		

**Parsimony:** representing variant in as few nucleotides as possible without reducing the length of any allele to 0

**Left (right) aligning =**  
shifting the start position of a variant as far to the left (right) as possible

There's a tool for doing this in VCFs!  
<https://gatk.broadinstitute.org/hc/en-us/articles/360037225872-LeftAlignAndTrimVariants>

# Intervals are often represented in the BED format

- There are several flavors of BED format: BED3, BED4, BED6, BED8, etc
- First 3 fields always required: **chr, start, stop**
- Followed by up to 9 additional optional fields: name, score, strand, thickStart, thickEnd, itemRGB, blockCount, blockSizes, blockStarts

chr7	127471196	127472363	Pos1	0	+
chr7	127472363	127473530	Pos2	0	+
chr7	127473530	127474697	Pos3	0	+
chr7	127474697	127475864	Pos4	0	+
chr7	127475864	127477031	Neg1	0	-
chr7	127477031	127478198	Neg2	0	-
chr7	127478198	127479365	Neg3	0	-
chr7	127479365	127480532	Pos5	0	+
chr7	127480532	127481699	Neg4	0	-

# Manipulation of SAM/BAM and BED files

- Several tools are used ubiquitously in sequence analysis to manipulate these files
- SAM/BAM files
  - samtools
  - bamtools
  - Picard
- BED files
  - bedtools
  - bedops



# Bedtools: a swiss army knife for genome analysis



## BEDTools: a flexible suite of utilities for comparing genomic features

Aaron R. Quinlan ; Ira M. Hall 

Bioinformatics (2010) 26 (6): 841-842.

DOI: <https://doi.org/10.1093/bioinformatics/btq033>

Published: 28 January 2010 Article history ▾

### Abstract

**Motivation:** Testing for correlations between different sets of genomic features is a fundamental task in genomics research. However, searching for overlaps between features with existing web-based methods is complicated by the massive datasets that are routinely produced with current sequencing technologies. Fast and flexible tools are therefore required to ask complex questions of these data in an efficient manner.

**Results:** This article introduces a new software suite for the comparison, manipulation and annotation of genomic features in Browser Extensible Data (BED) and General Feature Format (GFF) format. BEDTools also supports the comparison of sequence alignments in BAM format to both BED and GFF features. The tools are extremely efficient and allow the user to compare large datasets (e.g. next-generation sequencing data) with both public and custom genome annotation tracks. BEDTools can be combined with one another as well as with standard UNIX commands, thus facilitating routine genomics tasks as well as pipelines that can quickly answer intricate questions of large genomic datasets.

## Papers:

<https://doi.org/10.1093/bioinformatics/btq033>

DOI: 10.1002/0471250953.bi1112s47

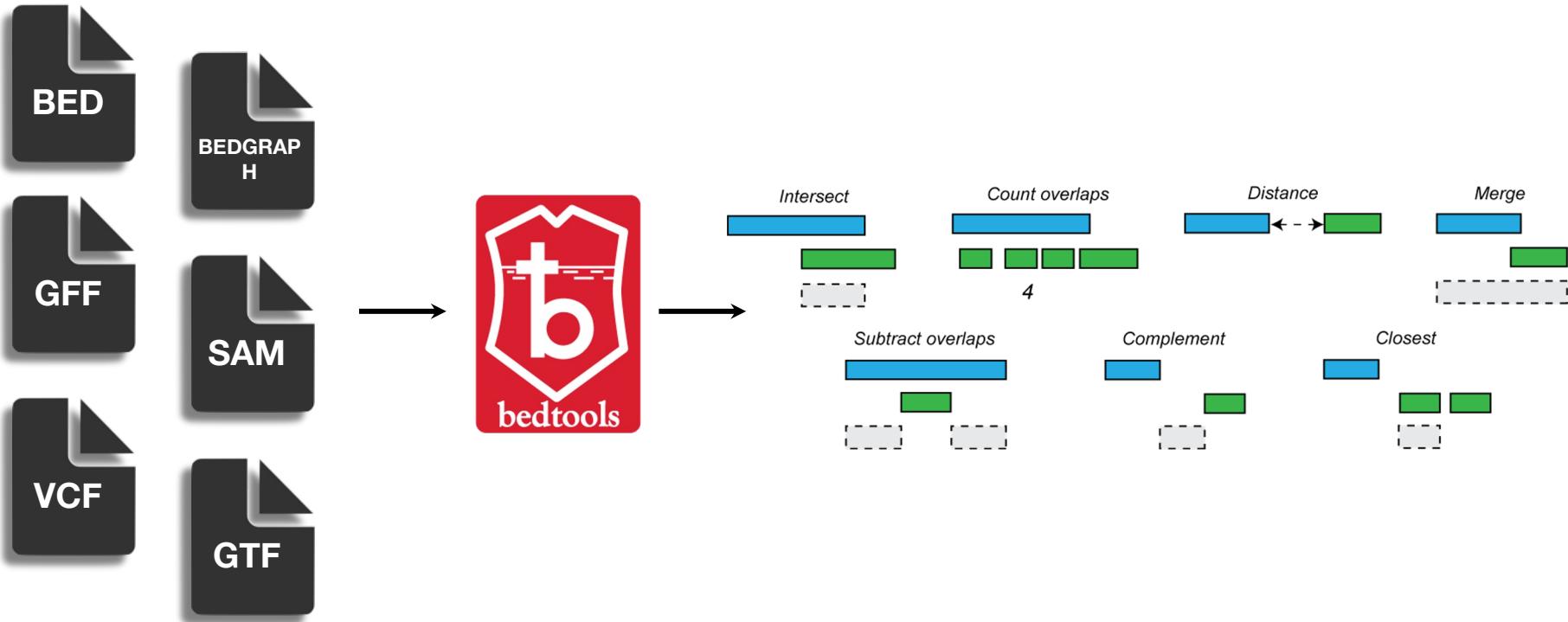
## Documentation:

<http://bedtools.readthedocs.io/en/latest/>

## Code:

<https://github.com/arq5x/bedtools2>

# Supports most interval formats & handles diff. coordinate systems



# Bedtools: example analyses

- Closest gene to a ChIP-seq peak.
- Is my latest discovery novel?
- 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?
- What is the density of variants in "windows" along the genome?

# Assignment: work through the bedtools tutorial

<https://sandbox.bio/tutorials/bedtools-intro>

For-credit students:

Pages 15-19 contain 5 exercises – submit a screenshot of the validated exercise on slide 19 after completing it – like this one:

Send it to [j.mckenzie@wustl.edu](mailto:j.mckenzie@wustl.edu) with "bfx exercise week 13" as the subject

## Exercises

Find non-exons

Create a BED file called `notexons.bed` that contains all of the intervals in the genome that are NOT exonic. Use the files `exons.bed` and `genome.txt` as input.

Exercise Criteria:

2 / 2

File `notexons.bed` exists

File `notexons.bed` contains non-exonic regions

[Check my work](#)