

# Concepts in genome coordinates, intervals, and arithmetic

Adapted from Aaron Quinlan's  
Applied Computational Genomics, Lecture 16 - 18

Jason Kunisaki  
Quinlan Lab  
University of Utah

# Learning objectives

## Part 1 (Mariam):

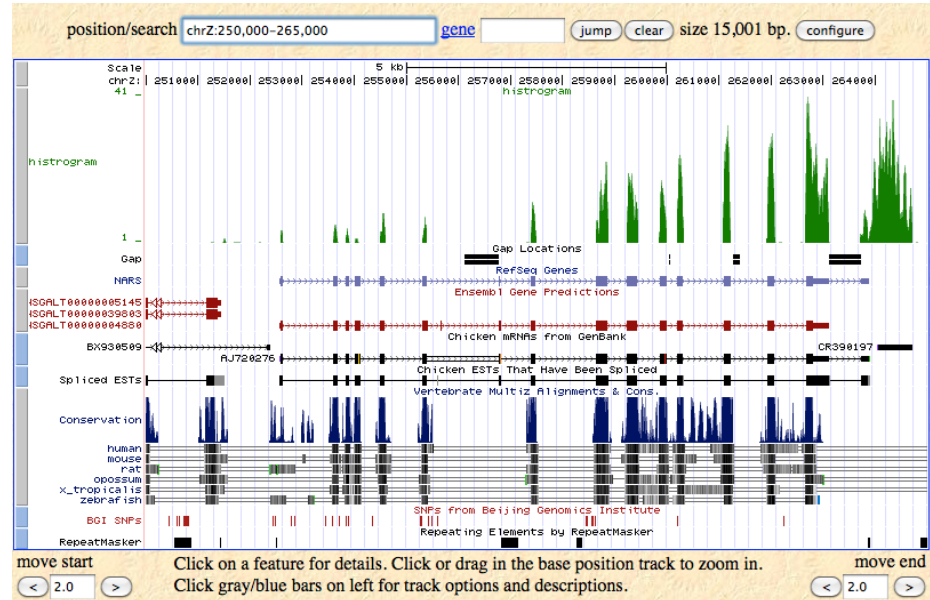
- Understand how “omic” analyses relies on features/annotations with specific genomic intervals or coordinates
- Introduce concepts in genome arithmetic
- Use the table browser to download data of interest

## Part 2 (Jason):

- Use **bedtools** to perform complex genome arithmetic analyses
- Understand the general utility of **bedtools** to answer many biological questions and problems

# Performing “omic” analyses requires an understanding of genome coordinates

- Genes: exons, introns, UTRs, promoters (BED, GFF, GTF)
- Genetic variation (VCF)
- Transcription factor binding sites (BED, BEDGRAPH)
- CpG islands (BED)
- Chromatin annotations (BED)
- Gene expression data (WIG, BIGWIG, BEDGRAPH)
- Your own observations: put them in context

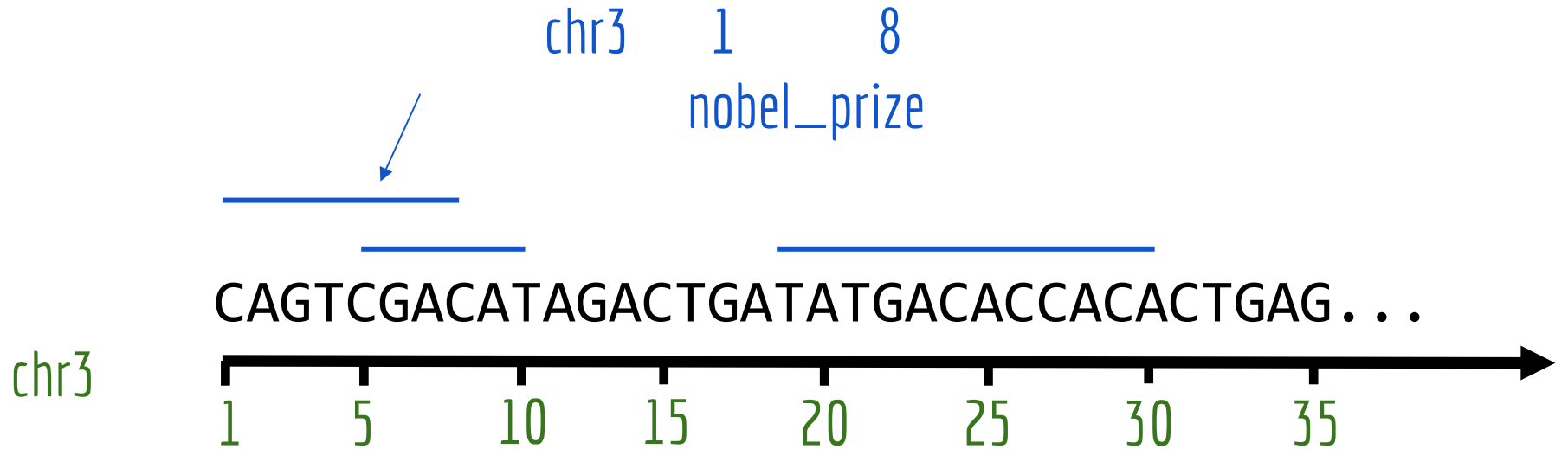




What is a coordinate system in genomics?



# The reference genome as a coordinate system



Great, genome formats using a common genome coordinate system! My life is going to be easy.

No.

# BED (browser extensible data) format

## BED format

[Index ▸](#)

BED format provides a flexible way to define the data lines that are displayed in an annotation track. BED lines have three required fields and nine additional optional fields. The number of fields per line must be consistent throughout any single set of data in an annotation track. The order of the optional fields is binding: lower-numbered fields must always be populated if higher-numbered fields are used.

If your data set is BED-like, but it is very large (over 50MB) and you would like to keep it on your own server, you should use the [bigBed](#) data format.

The first three required BED fields are:

1. **chrom** - The name of the chromosome (e.g. chr3, chrY, chr2\_random) or scaffold (e.g. scaffold10671).
2. **chromStart** - The starting position of the feature in the chromosome or scaffold. The first base in a chromosome is numbered 0.
3. **chromEnd** - The ending position of the feature in the chromosome or scaffold. The *chromEnd* base is not included in the display of the feature. For example, the first 100 bases of a chromosome are defined as *chromStart=0*, *chromEnd=100*, and span the bases numbered 0-99.

The 9 additional optional BED fields are:

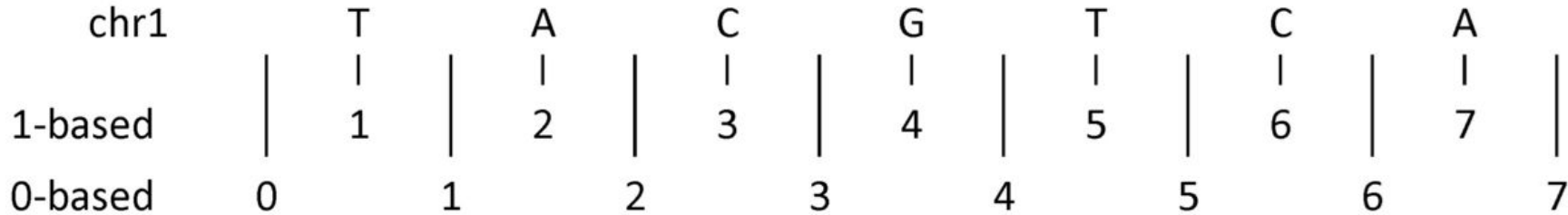
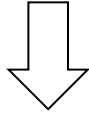
4. **name** - Defines the name of the BED line. This label is displayed to the left of the BED line in the Genome Browser window when the track is open to full display mode or directly to the left of the item in pack mode.
5. **score** - A score between 0 and 1000. If the track line *useScore* attribute is set to 1 for this annotation data set, the *score* value will determine the level of gray in which this feature is displayed (higher numbers = darker gray). This table shows the Genome Browser's translation of BED score values into shades of gray:

shade									
score in range	≤ 166	167-277	278-388	389-499	500-611	612-722	723-833	834-944	≥ 945

6. **strand** - Defines the strand - either '+' or '-'.
7. **thickStart** - The starting position at which the feature is drawn thickly (for example, the start codon in gene displays). When there is no thick part, thickStart and thickEnd are usually set to the chromStart position.
8. **thickEnd** - The ending position at which the feature is drawn thickly (for example, the stop codon in gene displays).
9. **itemRgb** - An RGB value of the form R,G,B (e.g. 255,0,0). If the track line *itemRgb* attribute is set to "On", this RGB value will determine the display color of the data contained in this BED line. NOTE: It is recommended that a simple color scheme (eight colors or less) be used with this attribute to avoid overwhelming the color resources of the Genome Browser and your Internet browser.
10. **blockCount** - The number of blocks (exons) in the BED line.
11. **blockSizes** - A comma-separated list of the block sizes. The number of items in this list should correspond to *blockCount*.
12. **blockStarts** - A comma-separated list of block starts. All of the *blockStart* positions should be calculated relative to *chromStart*. The number of items in this list should correspond to *blockCount*.



# BED files use 0 based, half open intervals. Say what?

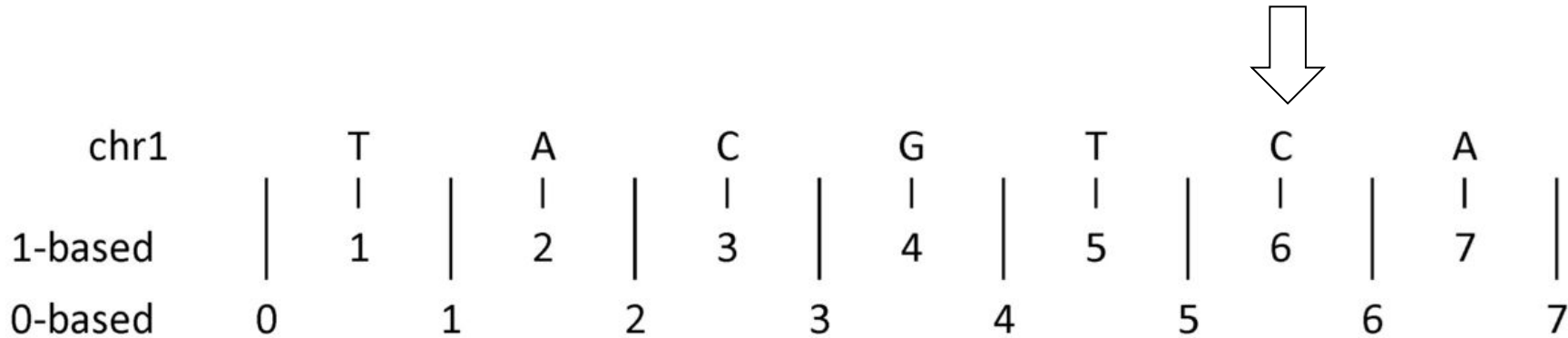


What interval defines the 1st nucleotide of chromosome 1?

1-based: chr1, start = 1, end = 1

0-based: chr1, start = 0, end = 1

# BED files use 0 based, half open intervals. Say what?



What interval defines the 1st nucleotide of chromosome 1?

1-based: chr1, start = 1, end = 1

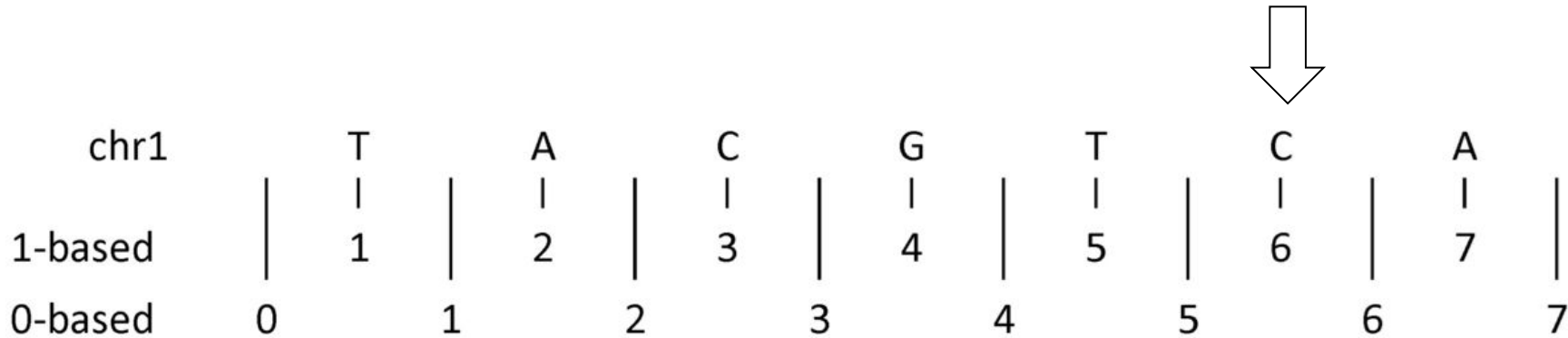
0-based: chr1, start = 0, end = 1

What interval defines the 6th nucleotide of chromosome 1?

1-based: chr1, start = ?, end = ?

0-based: chr1, start = ?, end = ?

# BED files use 0 based, half open intervals. Say what?



What interval defines the 1st nucleotide of chromosome 1?

1-based: chr1, start = 1, end = 1

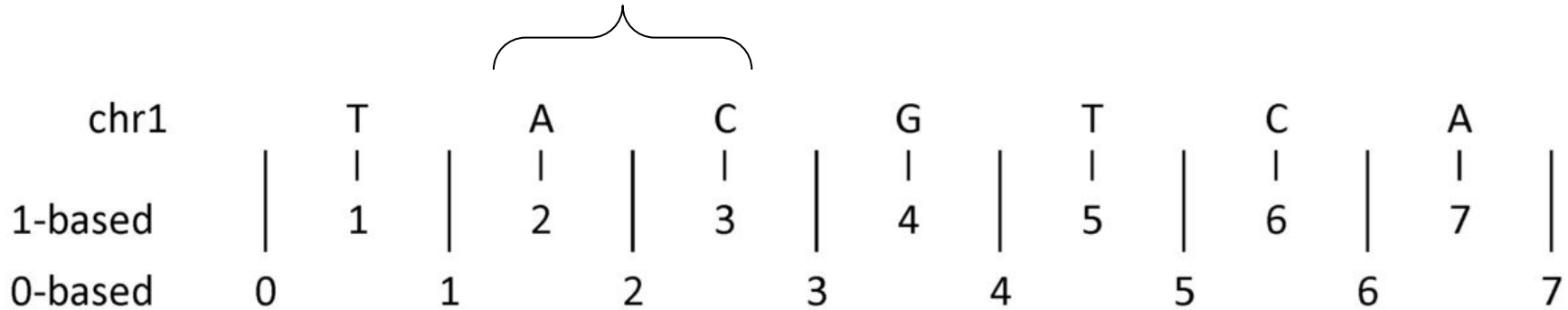
0-based: chr1, start = 0, end = 1

What interval defines the 6th nucleotide of chromosome 1?

1-based: chr1, start = 6, end = 6

0-based: chr1, start = 5, end = 6

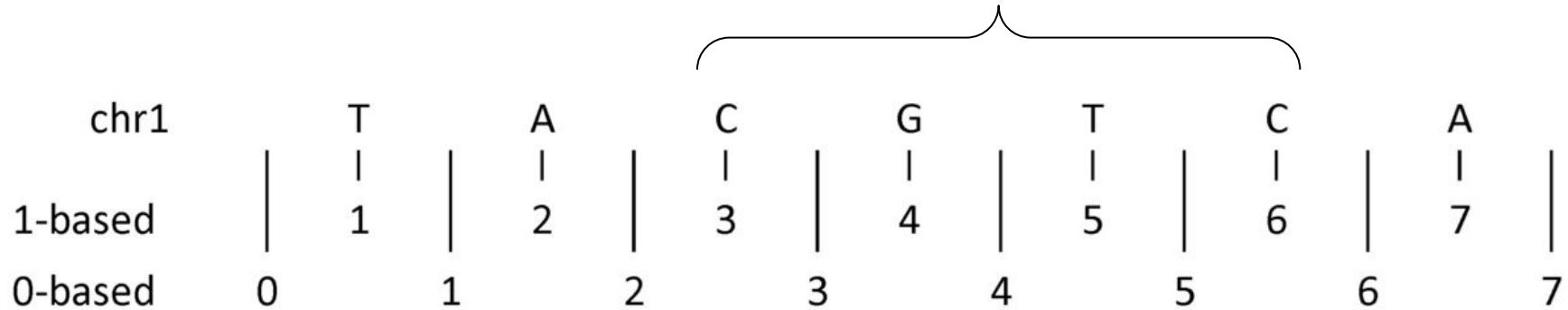
# BED files use 0 based, half open intervals. Say what?



What interval defines the 2nd and 3rd nucleotides of chromosome 1?

- 1-based? chr1, start = 2, end = 3
- 0-based? chr1, start = 1, end = 3

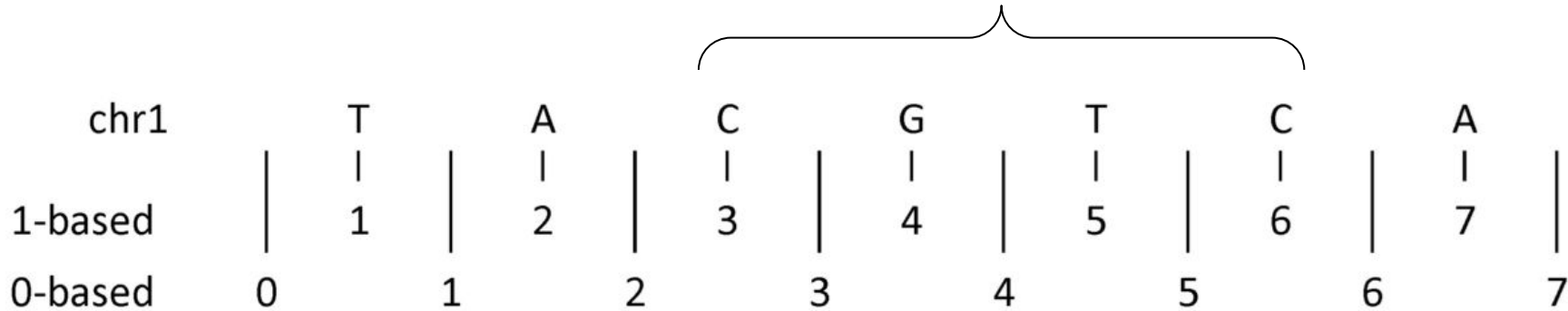
# BED files use 0 based, half open intervals. Say what?



What interval defines the 3rd and 6th nucleotides of chromosome 1?

- 1-based? chr1, start = ?, end = ?
- 0-based? chr1, start = ?, end = ?

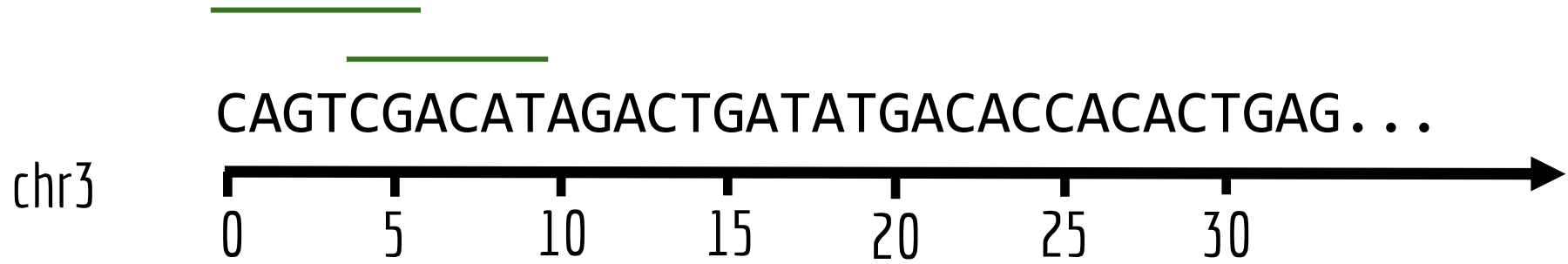
# BED files use 0 based, half open intervals. Say what?



What interval defines the 3rd and 6th nucleotides of chromosome 1?

- 1-based? chr1, start = 3, end = 6
- 0-based? chr1, start = 2, end = 6
- **0-based is convenient when measuring length of an interval**

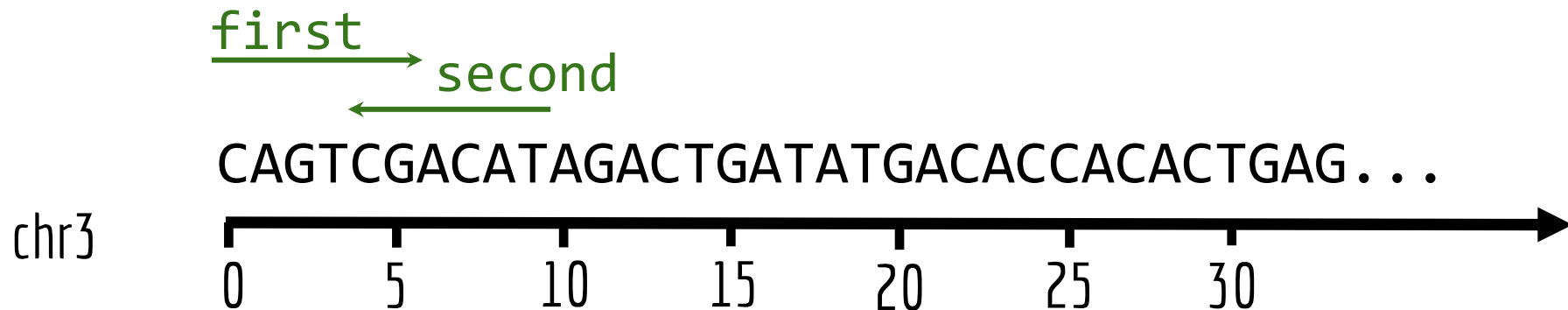
BED annotations also support "names" or "labels" (4th column)



chr3	0	7	exon_gene1
chr3	4	10	exon_gene2

Annotations

Start and end coordinates are agnostic to direction of the feature



chr3	0	7	+_strand_gene
chr3	4	10	-_strand_gene



# GFF/GTF format

**GFF format**Index ▾

GFF (General Feature Format) lines are based on the Sanger [GFF2 specification](#). GFF lines have nine required fields that *must* be tab-separated. If the fields are separated by spaces instead of tabs, the track will not display correctly. For more information on GFF format, refer to Sanger's [GFF page](#).

Note that there is also a GFF3 specification that is not currently supported by the Browser. All GFF tracks must be formatted according to Sanger's GFF2 specification.

If you would like to obtain browser data in GFF (GTF) format, please refer to [Genes in gtf or gff format](#) on the Wiki.

Here is a brief description of the GFF fields:

1. **seqname** - The name of the sequence. Must be a chromosome or scaffold.
2. **source** - The program that generated this feature.
3. **feature** - The name of this type of feature. Some examples of standard feature types are "CDS", "start\_codon", "stop\_codon", and "exon".
4. **start** - The starting position of the feature in the sequence. The first base is numbered 1.
5. **end** - The ending position of the feature (inclusive).
6. **score** - A score between 0 and 1000. If the track line *useScore* attribute is set to 1 for this annotation data set, the *score* value will determine the level of gray in which this feature is displayed (higher numbers = darker gray). If there is no score value, enter ".".
7. **strand** - Valid entries include '+', '-', or '.' (for don't know/don't care).
8. **frame** - If the feature is a coding exon, *frame* should be a number between 0-2 that represents the reading frame of the first base. If the feature is not a coding exon, the value should be '.'.
9. **group** - All lines with the same group are linked together into a single item.

**Example:**  
Here's an example of a GFF-based track. This [example](#) can be pasted into the browser without editing. NOTE: Paste operations on some operating systems will replace tabs with spaces, which will result in an error when the GFF track is uploaded. You can circumvent this problem by pasting the URL of the above example (<http://genome.ucsc.edu/goldenPath/help/regulatory.txt>) instead of the text itself into the custom annotation track text box. If you encounter an error when loading a GFF track, check that the data lines contain tabs rather than spaces.

```
browser position chr22:100000000-10025000
browser hide all
track name=regulatory description="TeleGene(tm) Regulatory Regions"
visibility=2
chr22 TeleGene enhancer 10000000 10001000 500 + . touch1
chr22 TeleGene promoter 10010000 10010100 900 + . touch1
chr22 TeleGene promoter 10020000 10025000 800 - . touch2
```

Click [here](#) to display this track in the Genome Browser.

chr22	TeleGene	enhancer	10000000	10001000	500	+	.	touch1
chr22	TeleGene	promoter	10010000	10010100	900	+	.	touch1
chr22	TeleGene	promoter	10020000	10025000	800	-	.	touch2



Note that the start and end coordinates are in different columns versus BED format

Formats use different coordinate systems. Because science.

**BED:** 0-based, half-open

**GFF:** 1-based, closed

**SAM:** 1-based, closed

**BAM:** 0-based, half-open.

**VCF:** 1-based, closed

...

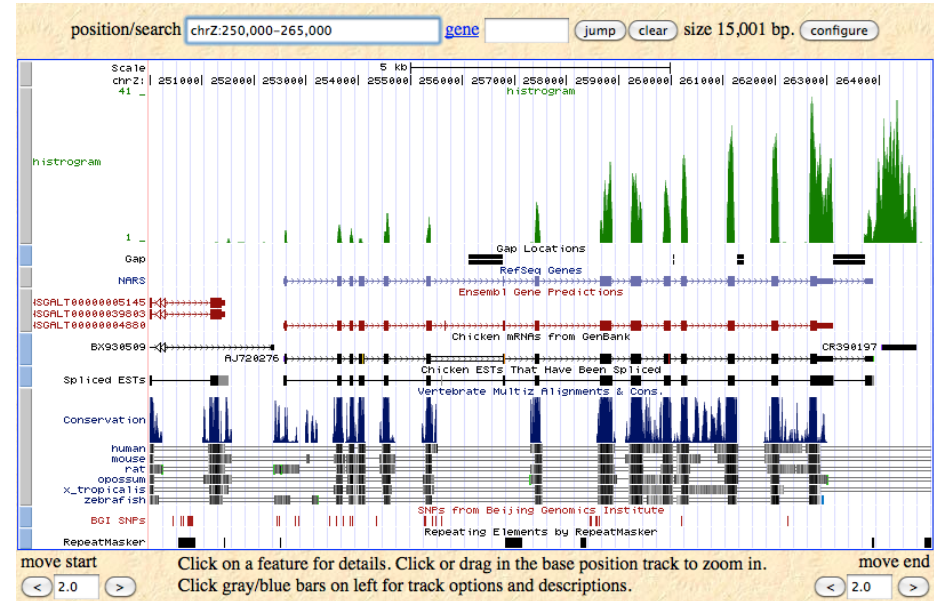
chr1	T	A	C	G	T	C	A
1-based	1	2	3	4	5	6	7
0-based	0	1	2	3	4	5	6



- inconsistent chromosome labels. (chr1 vs 1 vs c1)
- different sorting criteria. (lexicographic vs natural)
- mixed UNIX/Windows newlines. (transfer files across programs/systems)
- file violates spec with vigor.
- **file is gzipp'ed, not bgzipp'ed.**
- annotations use diff. genome builds.
- tool only works for one format.
- tool is hard-coded for specific build.
- **tool requires act of gods to compile.**

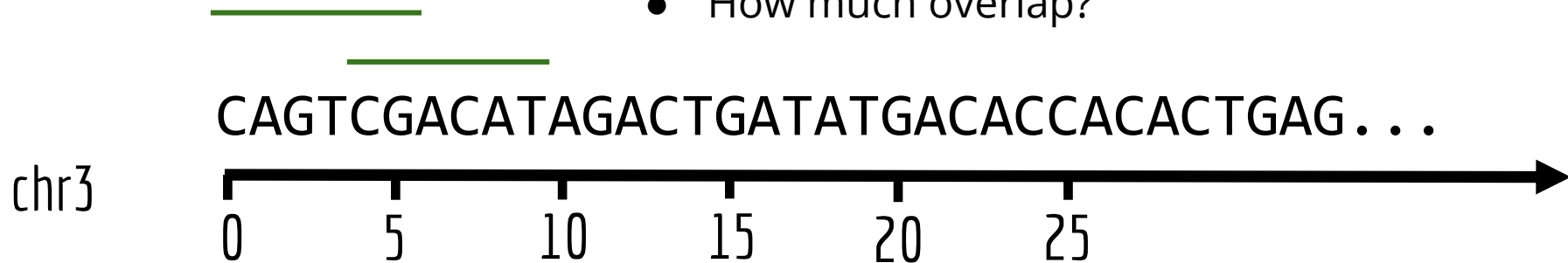
# Performing “omic” analyses requires an understanding of genome coordinates

- Genes: exons, introns, UTRs, promoters (BED, GFF, GTF)
- Genetic variation (VCF)
- Transcription factor binding sites (BED, BEDGRAPH)
- CpG islands (BED)
- Chromatin annotations (BED)
- Gene expression data (WIG, BIGWIG, BEDGRAPH)
- **Genome arithmetic:** the method of comparing, contrast and gain insight among multiple genome interval files



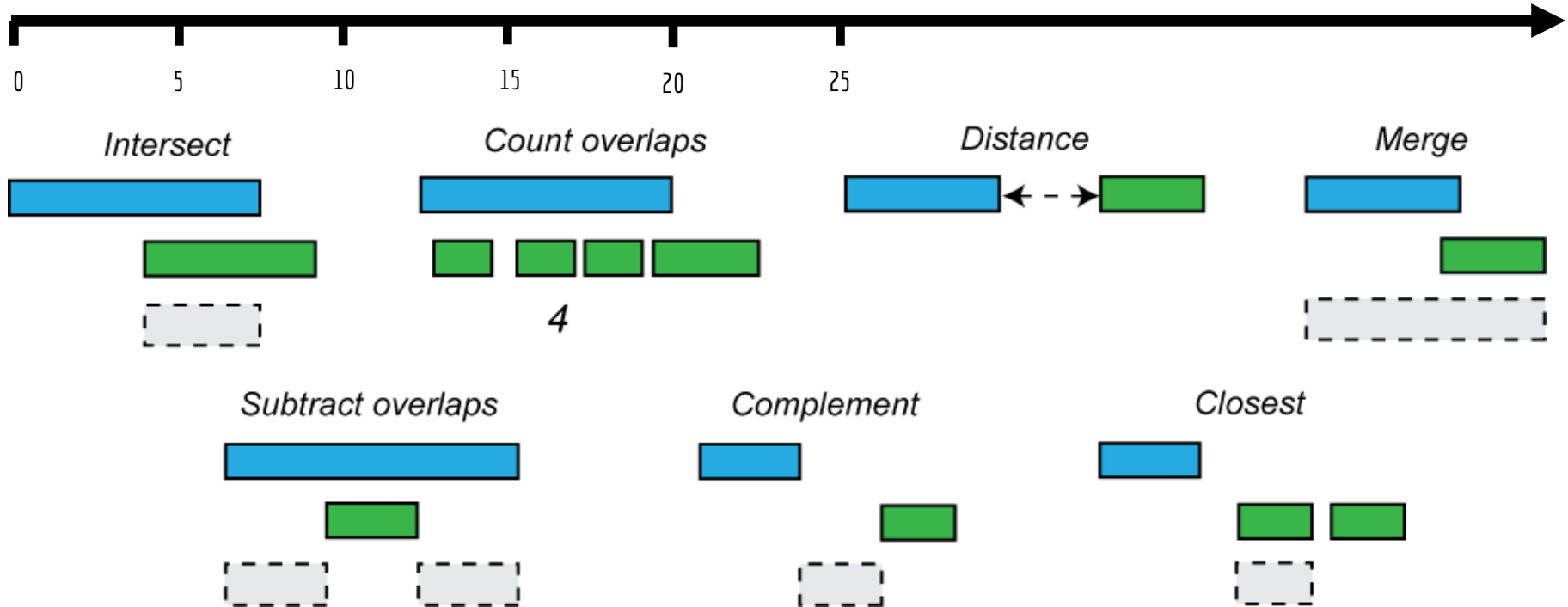
# Genome arithmetic depends upon the genome coordinate system

- Is there overlap?
- How much overlap?



chr3	0	7	enhancer
chr3	4	10	TF_binding_site

# Genome arithmetic operations



# Do two intervals intersect (overlap)?

10 20



17 27



17 27



10 20



10 20



13 16



13 16

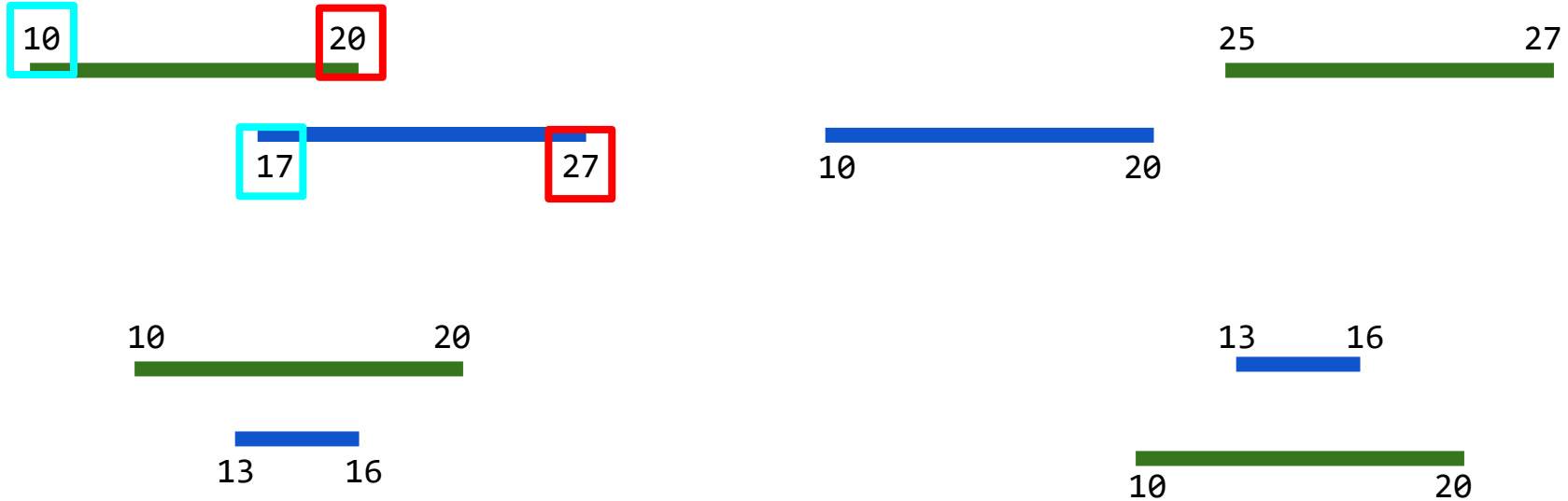


10 20



```
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!!!
```

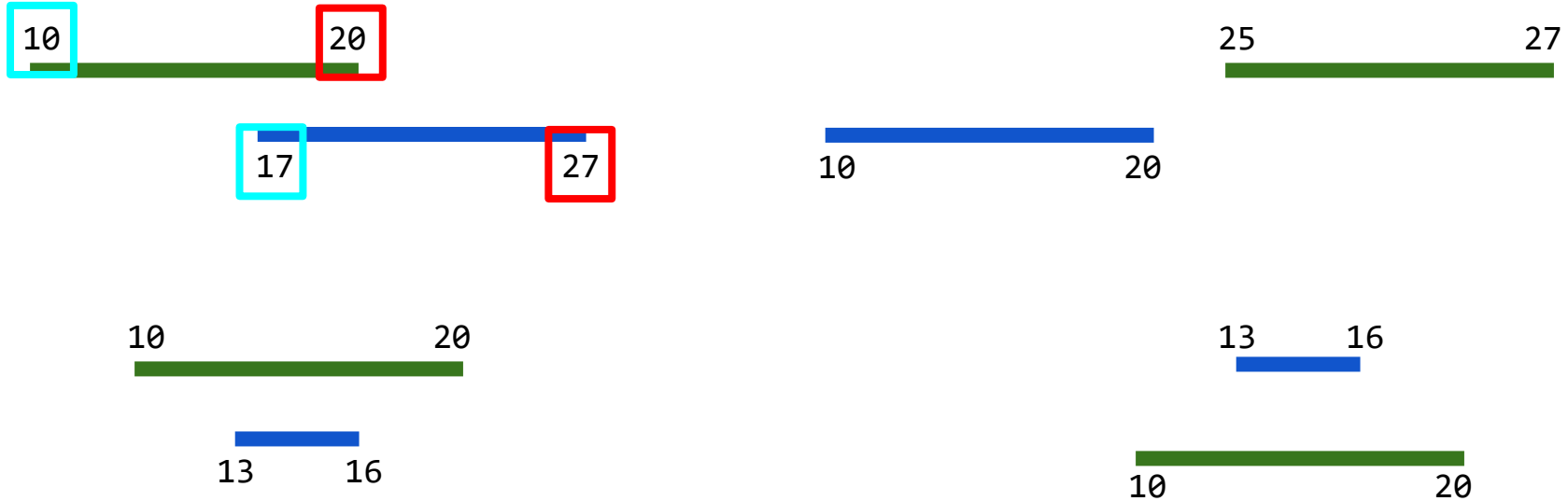
# Do two intervals intersect (overlap)? A simpler way.



$$\begin{aligned} I &= \min(a.\text{end}, b.\text{end}) - \max(a.\text{start}, b.\text{start}) \\ &= \min(20, 27) - \max(10, 17) \\ &= 20 - 17 = 3 \end{aligned}$$



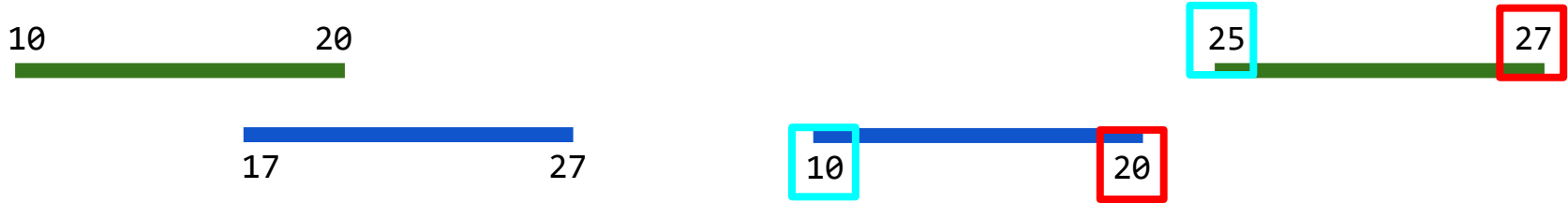
# Do two intervals intersect (overlap)? A simpler way.



$$\begin{aligned} I &= \min(a.\text{end}, b.\text{end}) - \max(a.\text{start}, b.\text{start}) \\ &= \min(20, 27) - \max(10, 17) \\ &= 20 - 17 = 3 \end{aligned}$$

If  $I > 0$  --> there is intersection  
If  $I \leq 0$  --> there is no intersection

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



$$\begin{aligned} I &= \min(a.\text{end}, b.\text{end}) - \max(a.\text{start}, b.\text{start}) \\ &= \min(20, 27) - \max(10, 25) \\ &= 20 - 25 = -5 \end{aligned}$$

If  $I > 0$  --> there is intersection  
If  $I \leq 0$  --> there is no intersection

# 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

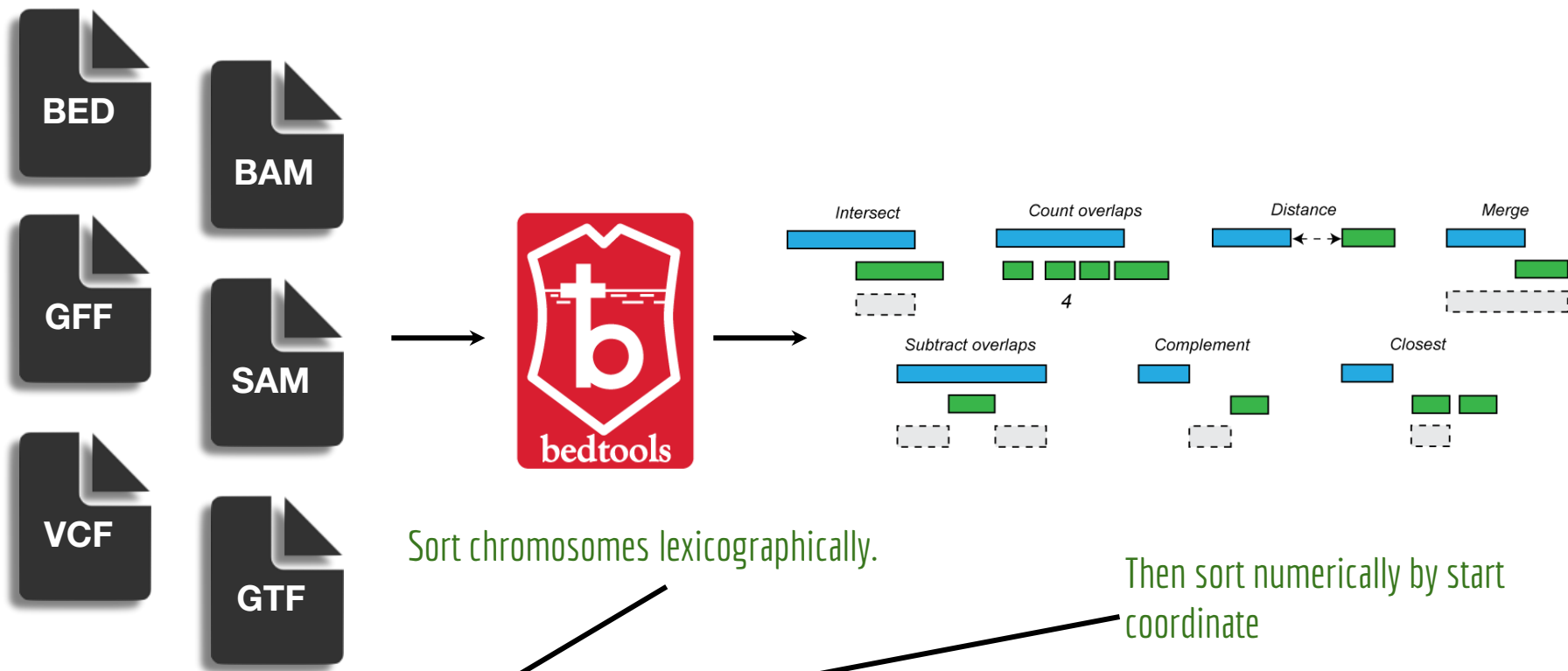
## Documentation:

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

## Code:

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

# Supports most interval formats & handles diff. coordinate systems




Sort chromosomes lexicographically.

Then sort numerically by start coordinate

```
sort -k1,1 -k2,2n interval_1.bed > interval_1.sorted.bed
```

# Quick overview of the UCSC Genome Browser

# The UCSC Table Browser

 Genomes Genome Browser Tools Mirrors Downloads My Data Help About Us

**Table Browser**

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). To examine the biological function of your set through annotation enrichments, send the data to [GREAT](#). Send data to [GenomeSpace](#) for use with diverse computational tools. Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data. All tables can be downloaded in their entirety from the [Sequence and Annotation Downloads](#) page.

**clade:** Mammal **genome:** Human **assembly:** Feb. 2009 (GRCh37/hg19)

**group:** Genes and Gene Predictions **track:** RefSeq Genes [add custom tracks](#) [track hubs](#)

**table:** refGene [describe table schema](#)

**region:** ☒ genome ☐ ENCODE Pilot regions ☐ position chr21:33031597-33041570 [lookup](#) [define regions](#)

**identifiers (names/accessions):** [paste list](#) [upload list](#)

**filter:** [create](#)

**intersection:** [create](#)

**correlation:** [create](#)

**output format:** ☒ all fields from selected table ☐ selected fields from primary and related tables

**output file:** ☐ sequence

**file type return:** ☐ GTF - gene transfer format ☐ CDS FASTA alignment from multiple alignment ☒ BED - browser extensible data ☐ custom track ☐ hyperlinks to Genome Browser

[get output](#) [summary](#)

Send output to ☐ [Galaxy](#) ☐ [GREAT](#) ☐ [GenomeSpace](#) ☒ keep output in browser)

To reset all user preferences (including custom tracks), [click here](#).



Slide from Jim Havrilla

# Let's save the annotations as a file called "refseq.bed"

Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). To examine the biological function of your set through annotation enrichments, send the data to [GREAT](#). Send data to [GenomeSpace](#) for use with diverse computational tools. Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data. All tables can be downloaded in their entirety from the [Sequence and Annotation Downloads](#) page.

clade:  genome:  assembly:

group:  track:

table:

region: ☒ genome ☐ ENCODE Pilot regions ☐ position

identifiers (names/accessions):

filter:

intersection:

correlation:

output format:  Send output to ☐ Galaxy ☐ GREAT ☐ GenomeSpace

output file:  (leave blank to keep output in browser)

file type returned: ☒ plain text ☐ gzip compressed

To reset all user cart settings (including custom tracks), [click here](#).

Then click "get output"

# We want the “Whole Gene” in this case. Other options...

[Home](#) [Genomes](#) [Genome Browser](#) [Tools](#) [Mirrors](#) [Downloads](#) [My Data](#) [Help](#) [About Us](#)

## Output refGene as BED

☐ Include [custom track](#) header:

name=

description=

visibility=

url=

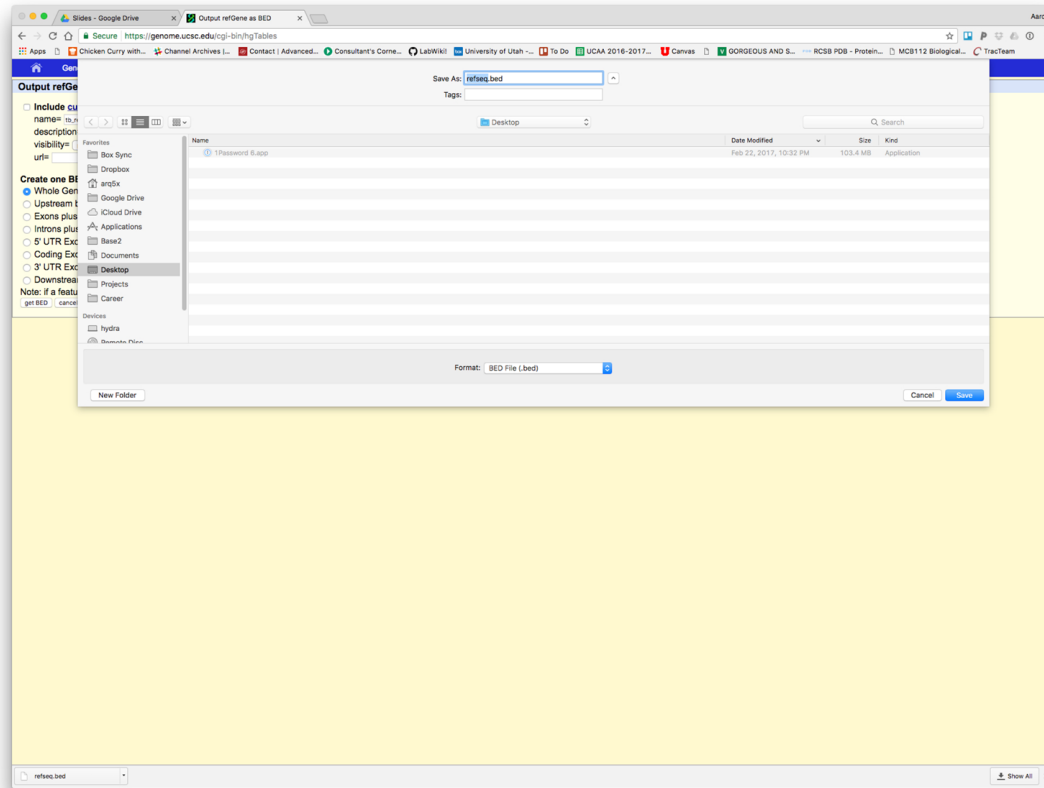
**Create one BED record per:**

- ☒ Whole Gene
- ☐ Upstream by  bases
- ☐ Exons plus  bases at each end
- ☐ Introns plus  bases at each end
- ☐ 5' UTR Exons
- ☐ Coding Exons
- ☐ 3' UTR Exons
- ☐ Downstream by  bases

Note: if a feature is close to the beginning or end of a chromosome and upstream/downstream bases are added, they may be truncated in order to avoid extending past the edge of the chromosome.



# Choose the location to which to download refseq.bed



# What is in the file?

```
1. head -30 ~/Desktop/refseq.bed | column -t | less -S (less)

chr1 66999251 67216822 NM_001308203 0 + 67000041 67208778 0 22 104,123,64,25,57,55,176,25,52,86,93,75,128,127,66,112,156,133,203,65,165,8067, 0,677,92278,99501,106208,109241,109975,137426,138375,139712,143435,146109,155579,156621,160870,18
chr1 66999638 67216822 NM_032291 0 + 67000041 67208778 0 25 413,64,25,72,57,55,176,12,12,25,52,86,93,75,501,128,127,60,112,156,133,203,65,165,8067, 0,91891,99114,101988,105821,108854,109588,126557,133574,137039,137988,139325,143048,145722,147913
chr1 16767166 16786584 NM_001145277 0 + 16767256 16785491 0 7 182,101,105,82,109,178,1248, 0,2960,7198,7388,8421,11166,18170,
chr1 16767166 16786584 NM_001145278 0 + 16767256 16785385 0 8 104,101,105,82,109,178,76,1248, 0,2960,7198,7388,8421,11166,15146,18170,
chr1 16767166 16786584 NM_018090 0 + 16767256 16785385 0 8 182,101,105,82,109,178,76,1248, 0,2960,7198,7388,8421,11166,15146,18170,
chr1 33547778 33567493 NR_126031 0 + 33567493 33567493 0 8 177,174,173,172,166,163,113,60, 0,1776,9872,11067,12370,14529,15889,19655,
chr1 48998526 50489626 NM_001323575 0 - 48999844 50489468 0 13 1439,97,163,153,112,115,90,40,217,95,125,123,192, 0,6787,54149,57978,101638,120482,130297,334336,512729,712915,1164458,1318541,1490908,
chr1 48998526 50489626 NM_001323574 0 - 48999844 50489468 0 14 1439,27,97,163,153,112,115,90,40,217,95,161,123,192, 0,2035,6787,54149,57978,101638,120482,130297,334336,512729,712915,1164422,1318541,1490908,
chr1 48998526 50489626 NR_136623 0 - 50489626 50489626 0 13 1439,97,163,153,112,115,90,124,40,217,95,123,192, 0,6787,54149,57978,101638,120482,130297,156925,334336,512729,712915,1318541,1490908,
chr1 25071759 25178815 NM_013943 0 + 25072044 25167428 0 6 357,110,126,107,182,3552, 0,52473,68825,81741,94591,95504,
chr1 48998526 50489626 NM_032785 0 - 48999844 50489468 0 14 1439,27,97,163,153,112,115,90,40,217,95,125,123,192, 0,2035,6787,54149,57978,101638,120482,130297,334336,512729,712915,1164458,1318541,1490908,
chr1 33546713 33586132 NM_001293562 0 + 33547850 33585783 0 11 182,118,177,174,173,135,166,163,113,215,488, 0,278,1065,2841,10937,12169,13435,15594,16954,36789,38931,
chr1 48998526 50489626 NM_001323573 0 - 48999844 50489468 0 13 1439,97,163,153,112,115,90,40,217,95,161,123,192, 0,6787,54149,57978,101638,120482,130297,334336,512729,712915,1164422,1318541,1490908,
chr1 33546713 33586132 NM_052998 0 + 33547850 33585783 0 12 182,121,212,177,174,173,135,166,163,113,215,488, 0,275,488,1065,2841,10937,12169,13435,15594,16954,36789,38931,
chr1 8378144 8404227 NM_001080397 0 + 8378168 8404073 0 9 102,421,93,225,728,154,177,206,421, 0,6221,7213,7733,12124,17352,19731,21408,25662,
chr1 33547778 33567493 NM_001301826 0 + 33547850 33567493 0 8 177,174,173,135,166,163,113,60, 0,1776,9872,11104,12370,14529,15889,19655,
chr1 33547778 33586132 NM_001301825 0 + 33547850 33585783 0 9 177,174,173,135,166,163,173,215,488, 0,1776,9872,11104,12370,14529,15829,35724,37866,
chr1 33546729 33586132 NM_001301824 0 + 33557656 33585783 0 8 380,173,135,166,163,113,215,488, 0,10921,12153,13419,15578,16938,36773,38915,
chr1 33546729 33586132 NM_001301823 0 + 33557656 33585783 0 9 380,174,173,135,166,163,113,215,488, 0,2825,10921,12153,13419,15578,16938,36773,38915,
chr1 92145899 92371559 NM_001195684 0 - 92149295 92327088 0 18 3515,108,42,121,300,159,141,153,335,190,148,169,184,138,185,174,61,177, 0,15329,17746,28320,31900,35893,36225,38969,39550,41612,47316,49462,54433,78270,116944,181128,219
chr1 92145899 92351836 NM_001195683 0 - 92149295 92327088 0 17 3515,108,42,121,300,159,141,153,335,190,148,169,184,138,185,174,402, 0,15329,17746,28320,31900,35893,36225,38969,39550,41612,47316,49462,54433,78270,116944,181128,205
chr1 100652477 100715409 NM_001918 0 - 100661810 100715376 0 11 9501,72,192,78,167,217,122,182,76,124,84, 0,19308,19523,23772,27895,29061,31704,43811,48514,53839,62848,
chr1 92145899 92351836 NM_003243 0 - 92149295 92327088 0 17 3515,108,42,121,300,159,141,153,338,190,148,169,184,138,185,174,402, 0,15329,17746,28320,31900,35893,36225,38969,39550,41612,47316,49462,54433,78270,116944,181128,205
chr1 92145899 92351836 NR_036634 0 - 92351836 92351836 0 18 3515,108,42,121,300,159,141,153,338,190,148,169,184,138,185,97,174,402, 0,15329,17746,28320,31900,35893,36225,38969,39550,41612,47316,49462,54433,78270,116944,120616,181
```