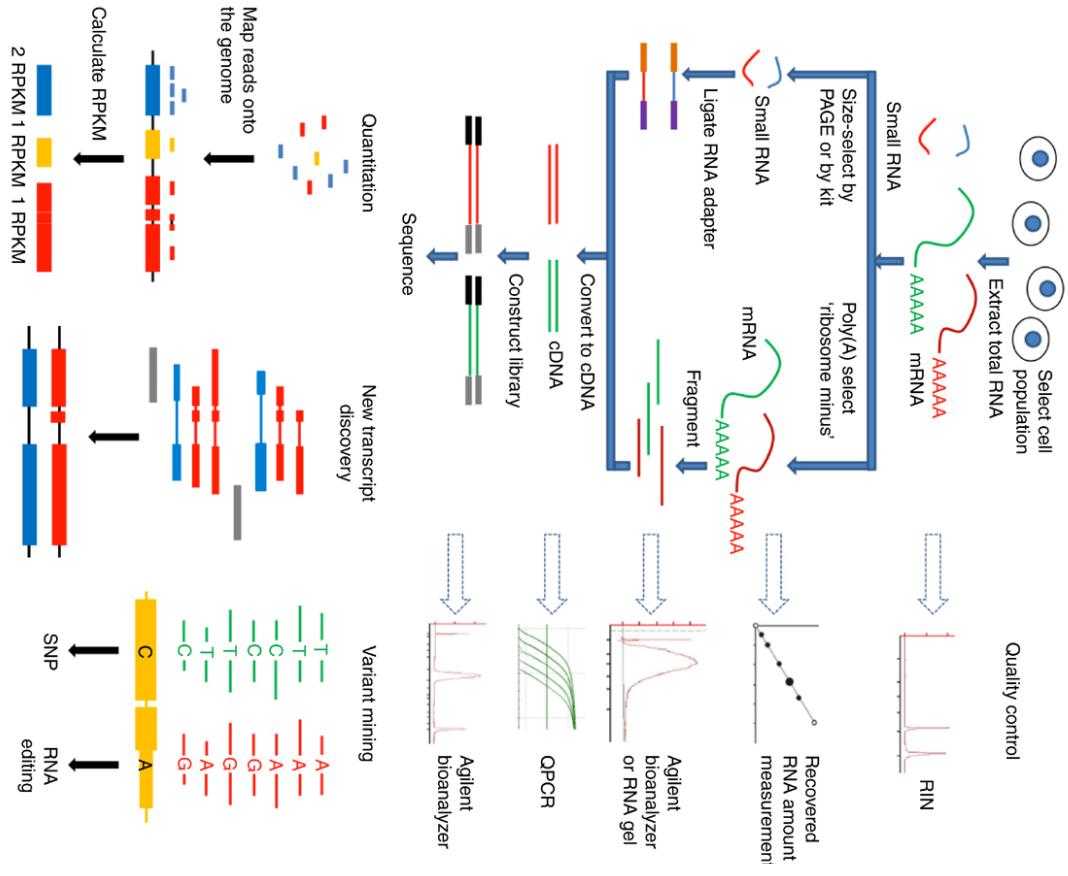


# Intro to Gene Annotations and RNA-seq

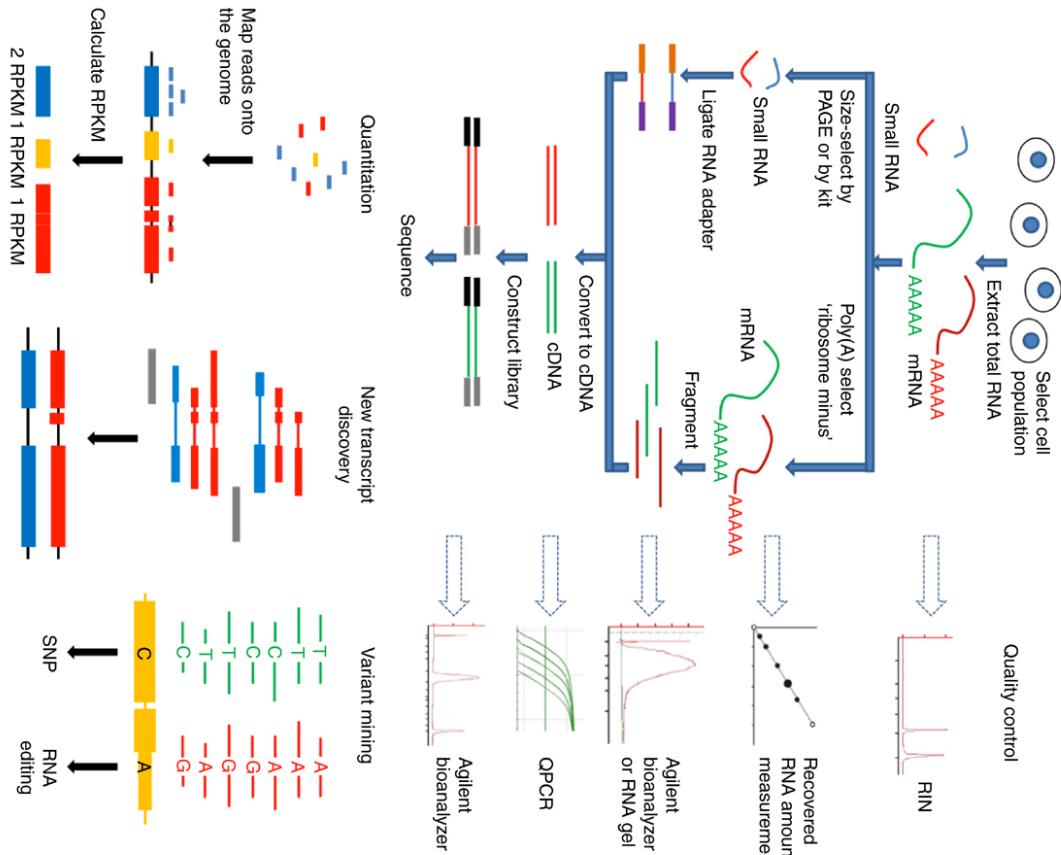
## Overview

- Overview of RNA-seq
  - Alignment
  - Visualization
- Genome Annotations
- Coordinate Systems
- Introduction to File Formats
- Explore GTF File
- RNA-seq analysis
  - Quantification
  - Differential Expression
  - Astrocyte Workflows

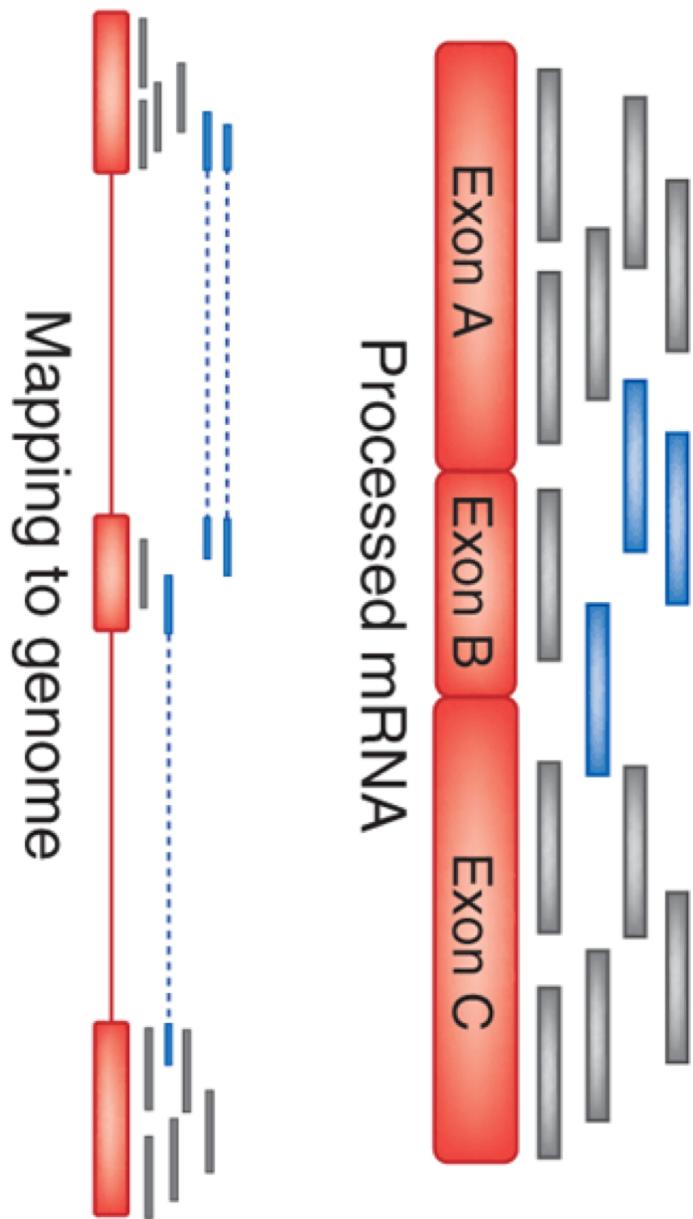


# Intro to RNA-seq and review

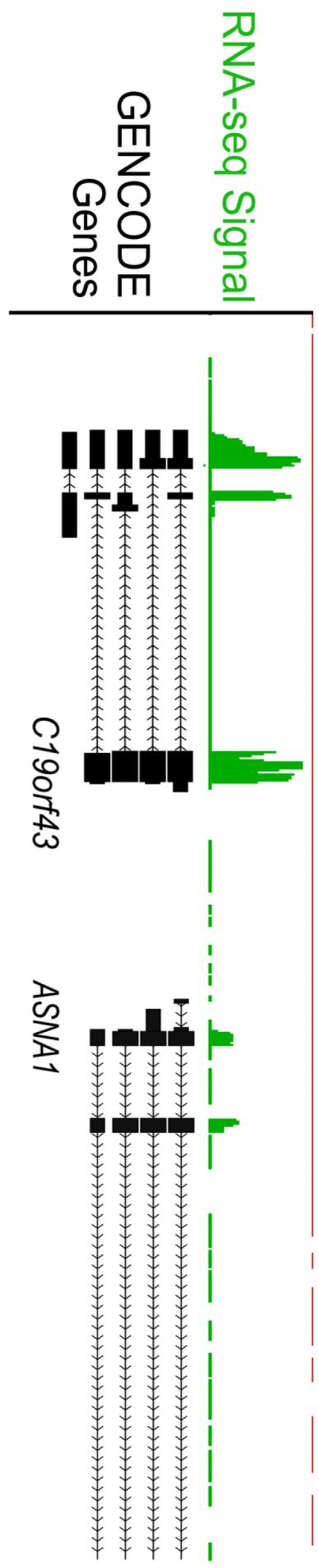
## Overview of RNA-seq



## Alignment to Genome and Reference Annotations



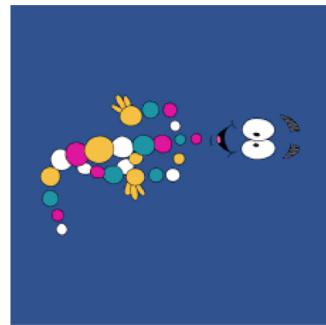
## Visualization of Alignment



# Review of Alignment and Intro to RNA-seq

## What is a genome annotation?

CTGCCGTCTGCTGCCATCGGAGCCAAAGCCGGCTGTGACTGCTCAGAC  
CAGCCGGCTGGAGGGGGCTCAGCAGGTCTGGCTTGGCCTGGAGA  
GCAGGGGAAGATCAGGCAGGCCATCGCTGCCACAGAACCCAGTGGATTG  
GCCTAGGGGGATCTCTGAGCTCAACAAGCCCTCTGGGGTAGGTG  
AGAGACGGGGCAGAGGCCAGGCACAGCCAAGAGGGCTGAAGAAT  
GGTAGAACGGAGCAGCTGGTGATGTGGGGCCACCGGCCCCAGGCTCT  
GTCTCCCCCAGGGTGTGATGCCAGGCATGCCCTCCCCAGCATCA  
GGTCTCAGAGCTGCAGAAAGACGACGCCGACTTGGATCACACTCTTG  
AGTGTCCCCAGTGTGCAGAGGTGAGAGGAGTAGACAGTGAAGGGAG  
TGGCGTCGCCCTAGGGCTCTACGGGCGGGCTCTGTCTCCTGGAG  
AGGCTTCGATGCCCTCACCCCTCTTGATCTTCCCCTGTGATGTCATCT  
GGAGCCTGCTGCTGGGGCTATAAGCCTCCTAGTGTGGCTCCAA  
GGCCTGGCAGAGCTTCCCAGGGAAGCTACAGCAGCAAACAGTCTGC  
ATGGGTCACTCCCTCACTCCAGCTCAGGCCAGGCCAGGGGGCTGGAG  
AGAAAGGCTCTGGAGAACCTGTGATGAAGGCTGTCAACCAGTCCAT  
AGGCAAGCCTGGCTGCCTCAGCTGGTCAGACAGCAGGGCTGGAGAAG  
GGGAGGAGAGGAAGTGAGGTTGCCTGCCCTGTCTACCTGAGGCTGA  
GGAAGGAGAAGGGGATGCACTGTTGGGAGGCAGCTGTAACCTAAAGCT  
TAGCCTCTGTCACAGGAAGGCCATCAGGCACCAAGGGATTCTG  
CCAGCATAGTGCTCCTGGACCTGATACACCCGGCACCCGTGCTGGAC  
ACGCTGTGGCTGGATCTGAGCCCTGGTGGAGGTCAAAGCACCTTGG  
TTCTGCCATTGCTGTGTTGAGTTCACTCTGCCTTTCTTCCT  
AGAGCCTCCACACCCCGAGATCACATTCTCACTGCCCTTGCTGCC  
AGTTTCACCAAGTAGGCCTCTTCTGACAGGCAGCTGCACCACTGCT  
GGCGCTGTGCCCTTCTGCTCTGCCGCTGGAGACGGTGTGTCATG

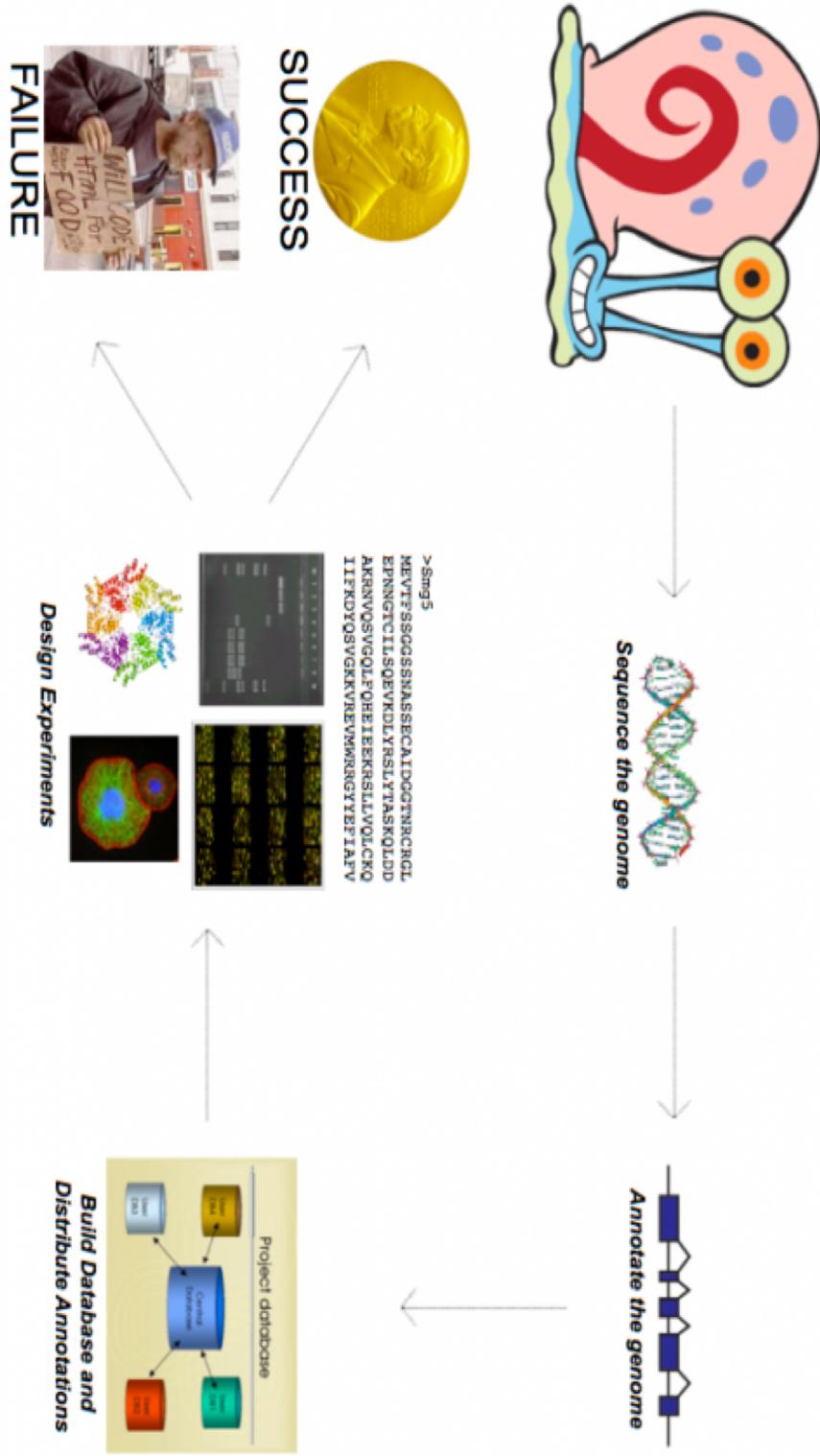


## What is an annotation?

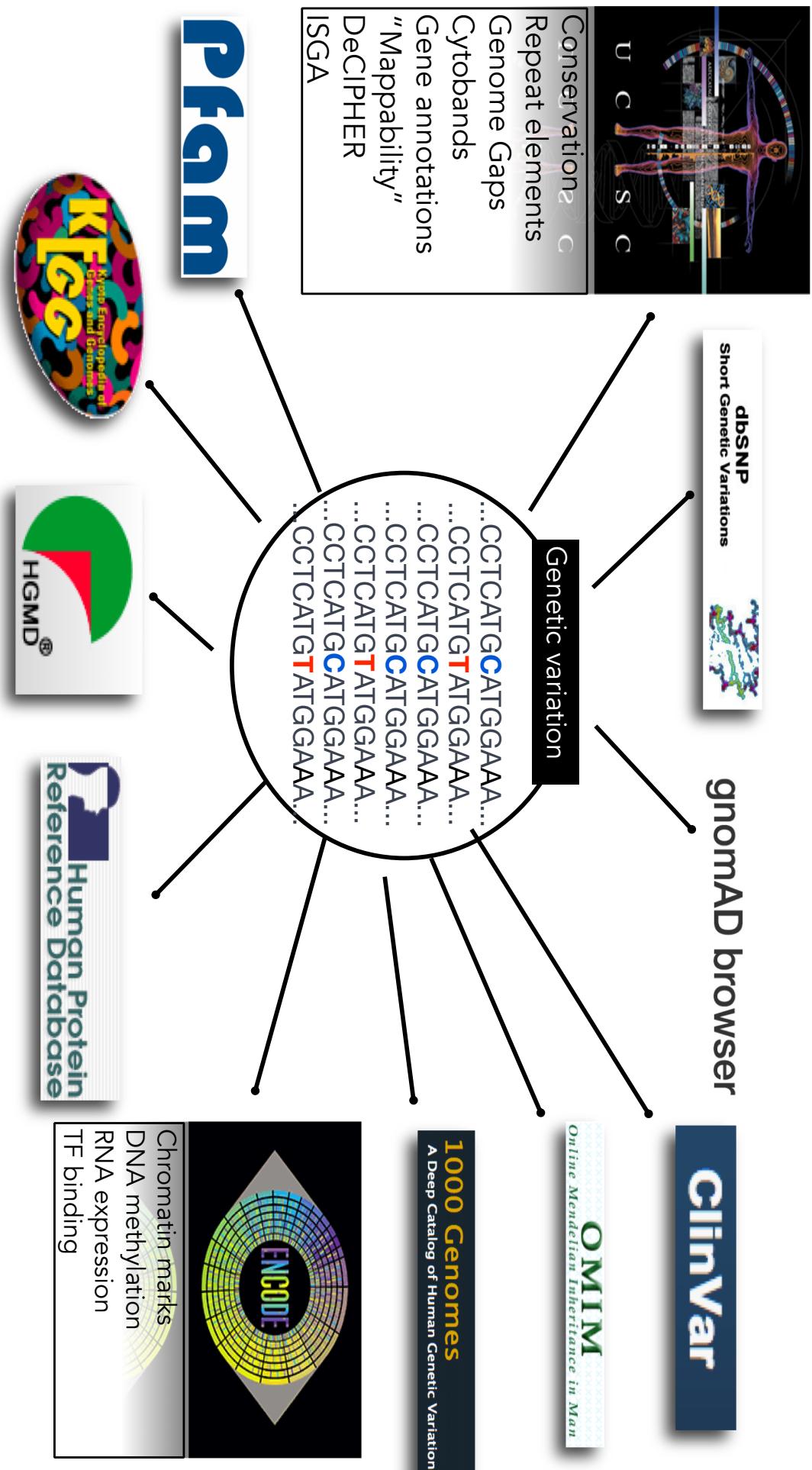
- Gene
- Repeat
- Disease Locus
- Something interesting

```
CTGCCGCTTGCTGCCATGGAGCCAAAGCCGGCTGACTGCTCAGAC  
CAGCCGGCTGGAGGGGGCTCAGGCTGGCTTGGCCTGGAGA  
GCAGGTGGAAGATCAGGCAGGCCATCGCTGCCACAGAACCCAGTGGATTG  
GCCTAGGTTGGATCTCTGAGCTAACAAAGCCCTCTGGGGTAGGTGC  
AGAGACGGGGAGCAGCTGGTGTGAGAACAGGGCTGAAGAAT  
GGTAGAAACGGAGCAGCTGGTGTGAGAACAGGGCTGAAGAAT  
GTCTCCCCAGGTGGTGTGAGAACAGGGCTGAAGAAT  
GGTCTCCAGAGCTGAGAACAGAACAGGGGACTTGGACACACTCTTG  
AGTGTCCCAGTGTGAGAGGTGAGAGGAGAGTAGACAGTGAGTGGAG  
TGGCGTCGCCCCTAGGGCTCTACGGGGGGCGTCTCTGCTCCTGGAG  
AGGCTTCGATGCCCTCACACCCCTTGTATCTCCCTGTATGTCATCT  
GGAGGCCCTGCTGCTGGGGGGCTATAAGCCTCTAGTCTGGCTCAA  
GCCCTGGCAGAGTCTTCCAGGGAAAGCTACAAGCAGCAAACAGTCTGC  
ATGGGTCACTCCCTCACTCCCAGCTCAGGCCAGGGCAGGGGCGCA  
AGAAAGGGCTGGTGAGAACCTGTGCATGAAGGCTGTCACCAGTCCAT  
AGGCAAGCCTGGCTGCCAGCTGGTCAGACAGAGGGCTGGAGAAG  
GGGAGAAGAGGAAGTGGAGGTTGCCCTGTCCTACCTGAGGCTGA  
GGAAGGGAGGGATGCACTGTTGGGAGGCAGTGTAACTCAAAGCCT  
TAGCCTCTGTTCCACGAAGCAGGGCATCAGGCACCAAGGGATTCTG  
CCAGCATAAGTGTCTCTGGACCAAGTGTGATACACCCGGCACCCCTGTCTGGAC  
ACGCTGTTGGCTGGATCTGAGCCCTGGAGGTCAAGCCACCTTGG  
TCTGCCATTGCTGCTGGAGATTCACTCTGCCCTTCTTCCCT  
AGAGCCTCACCACCCGGAGATCACATTCCTCACTGCTTGTCTGCC  
AGTTCACCAAGTAGGCCTCTTGTGACAGGCAGCTGCACCACTGCT  
GGCGCTGTGCCCTCTTGCTCTGCCCTGGAGACGGTGTGTCATG
```

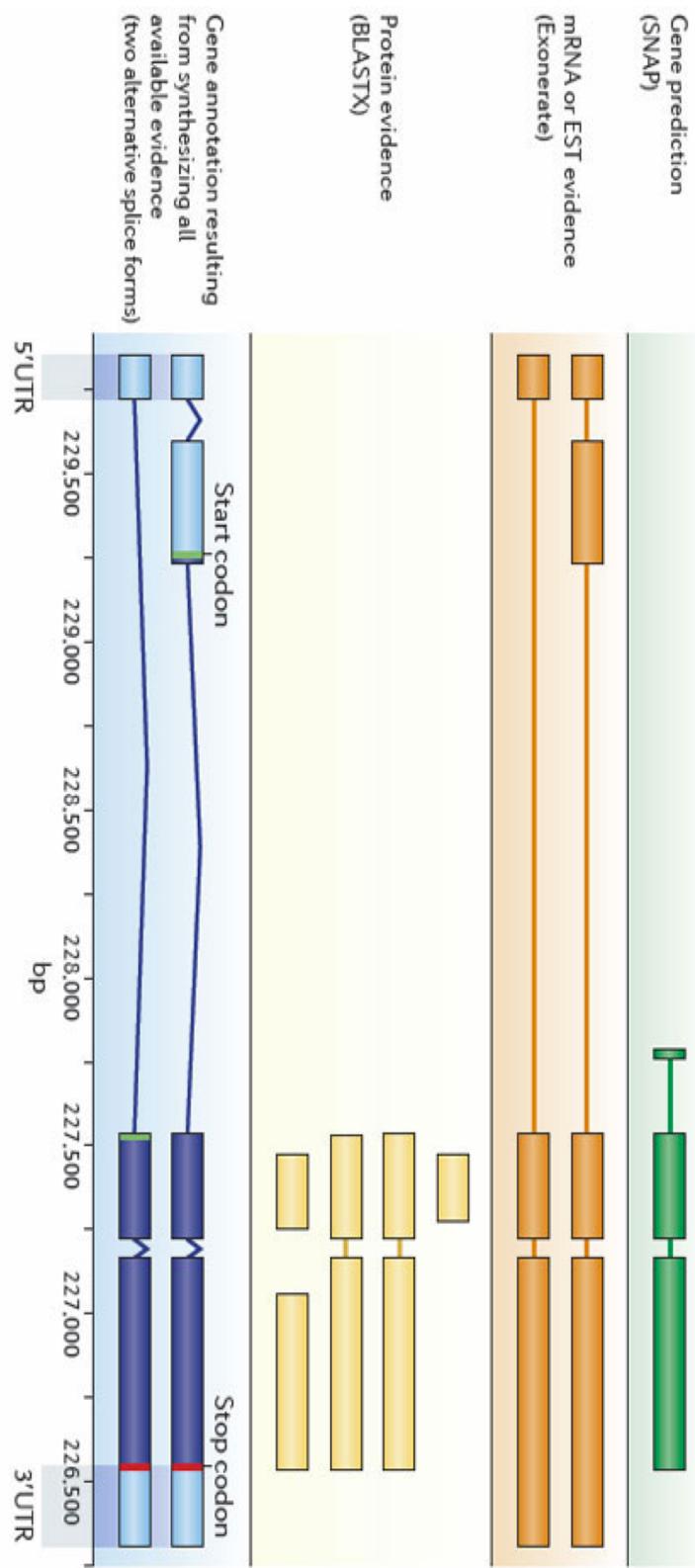
## Why are genome annotations important?



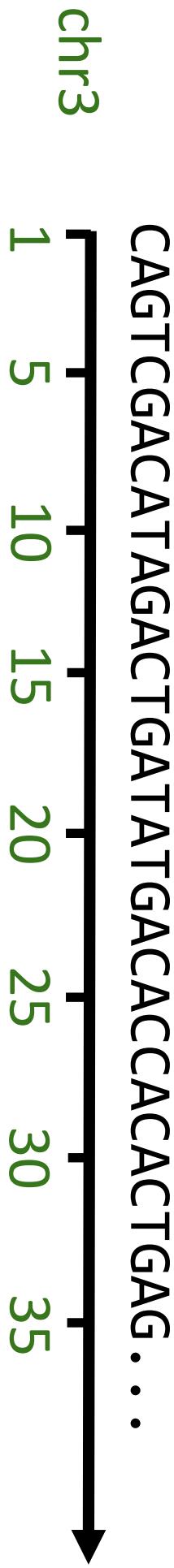
## Annotations Provide Context



## How do we annotate the genome?



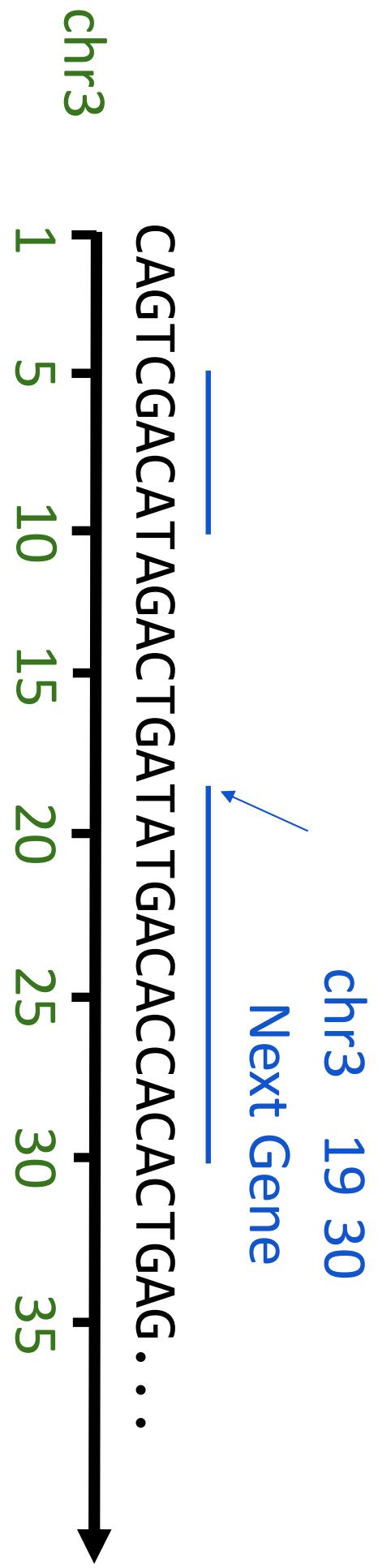
## The genome as coordinates



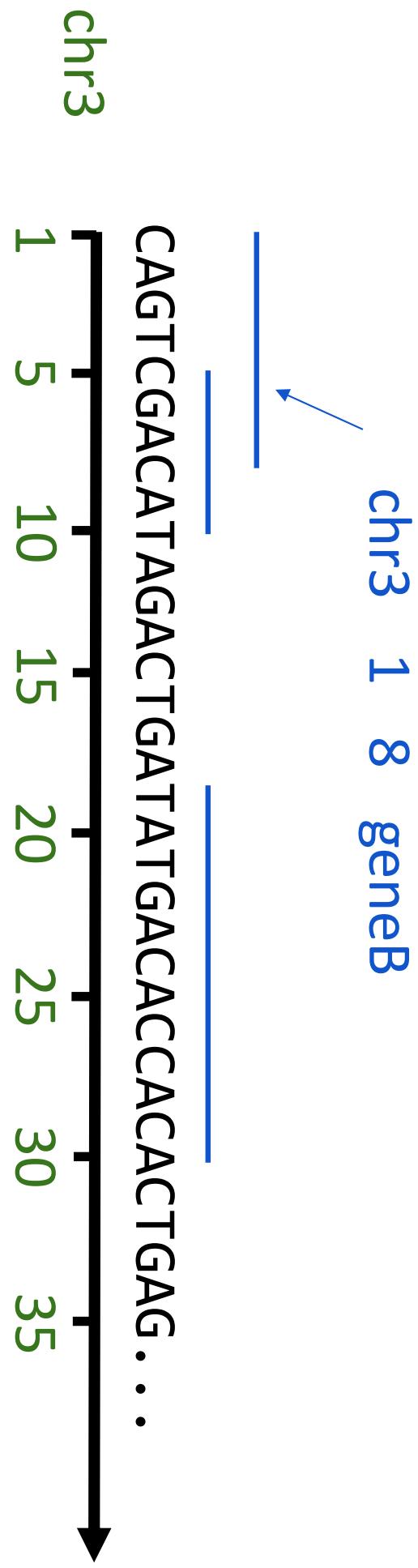
How do you describe a gene?



## What about 2 consecutive genes?



## What about overlapping genes?



## File Formats

- BED (Browser Extensible Format)
- GFF (General Feature Format)
- GTF (Gene Transfer Format, GTF2.2)

# BED File 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:

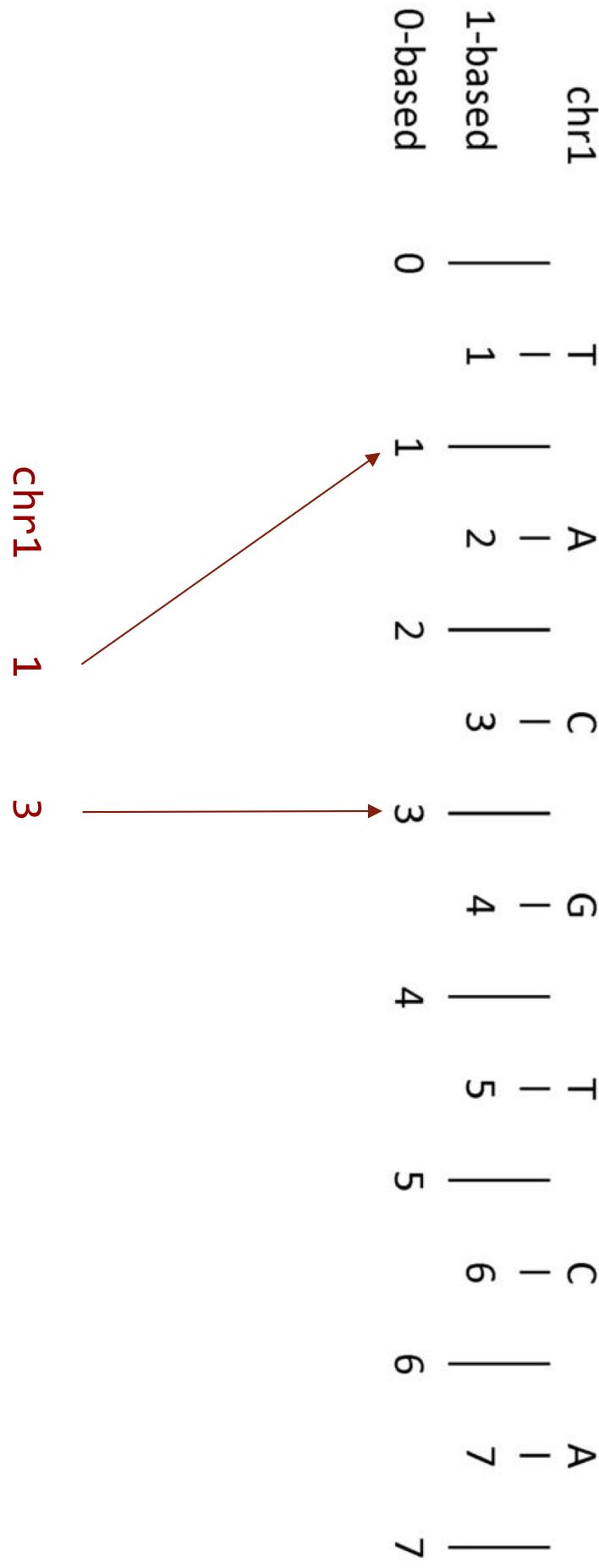
1. **chrom** - The name of the chromosome (e.g. chr3, chrY, chr2\_random) or scaffold (e.g. scaffold10671).
2. **chromStart** - The starting position of the feature in the chromosome or scaffold. The first base in a chromosome is numbered 0.
3. **chromEnd** - The ending position of the feature in the 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:

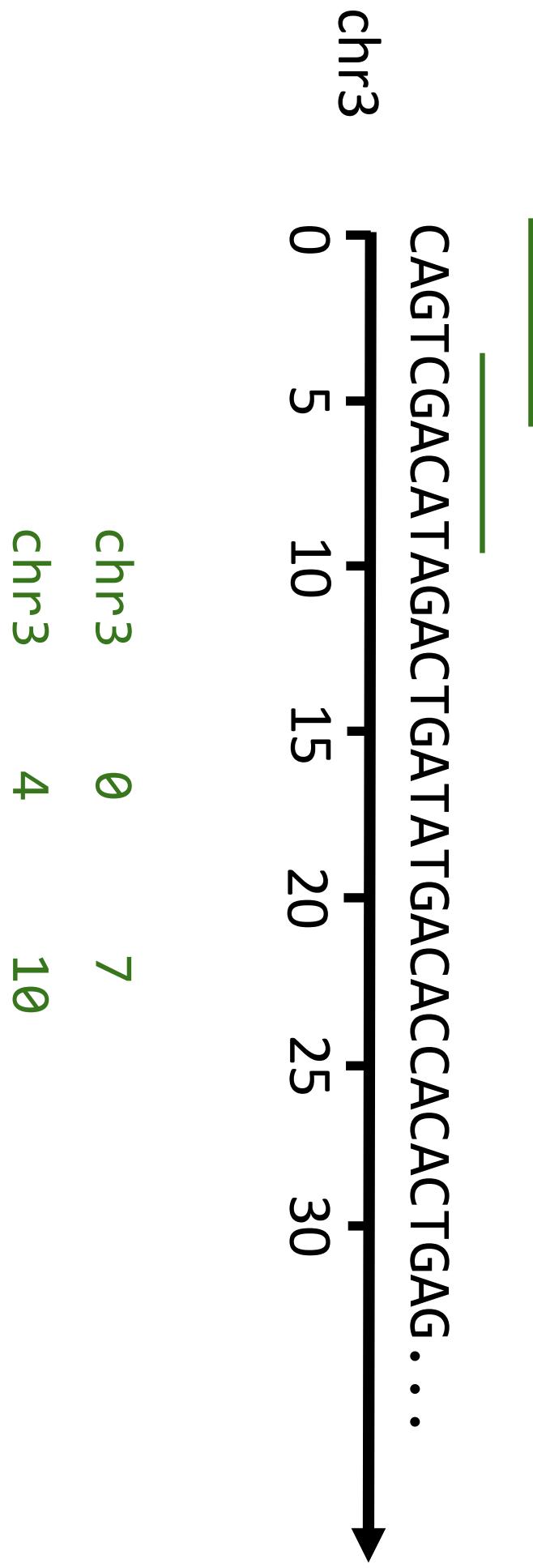
4. **name** - Defines the name of the BED line. This label is displayed to the left of the BED line in the Genome Browser window when the track is open to full display mode or directly to the left of the item in pack mode.
5. **score** - A score between 0 and 1000. If the track line *useScore* attribute 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*.

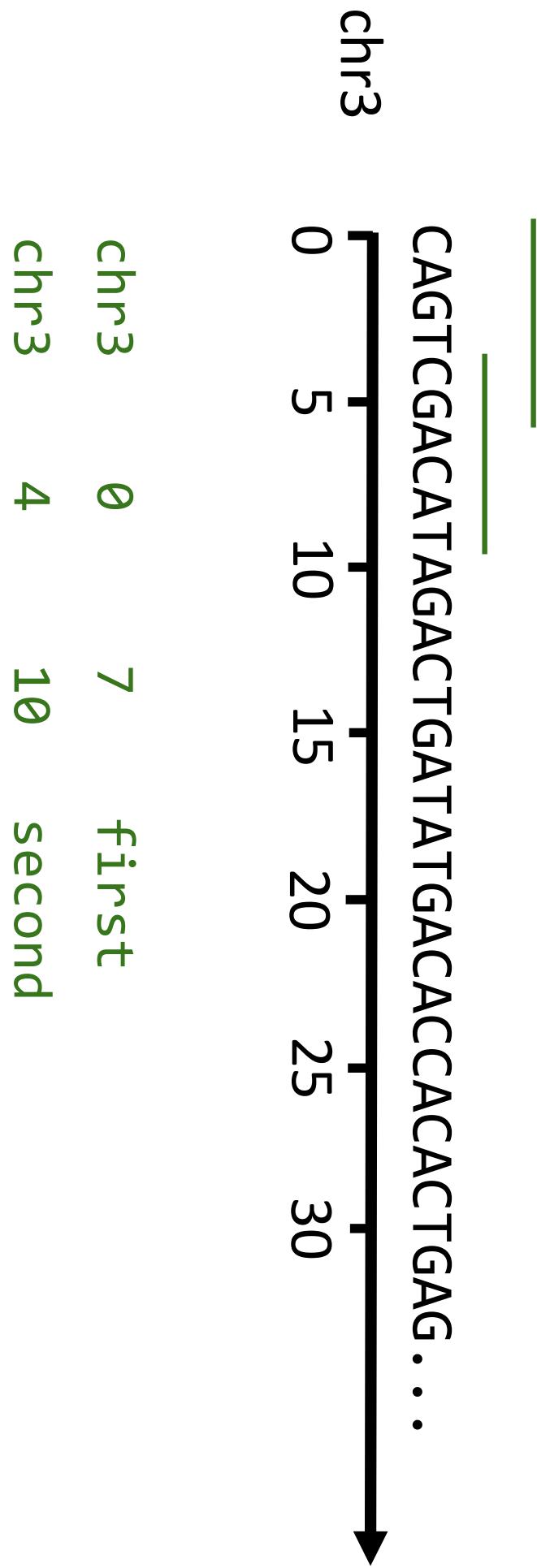
## BED File use 0-based, half open intervals



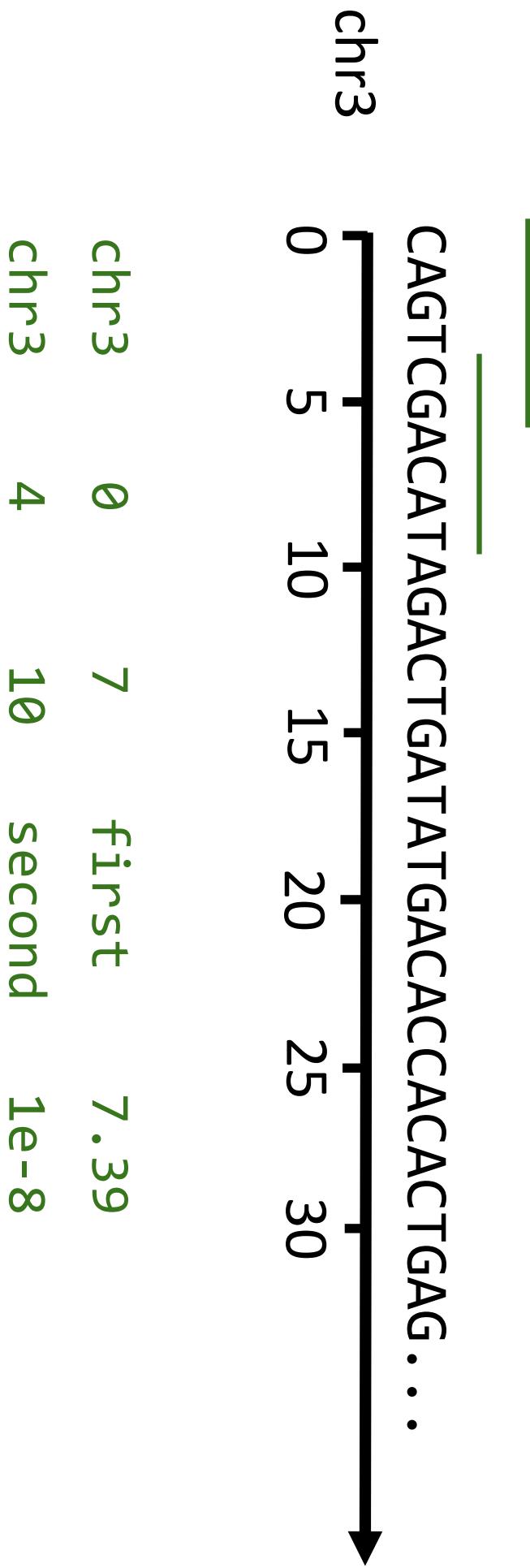
### Example BED3



## Add Name in 4<sup>th</sup> column to BED format



## Add Intensity in 5<sup>th</sup> column to BED format



## Add orientation/strand to the 6<sup>th</sup> field

first  
→ second

CAGTCGACATAGACTGATATGACACCACTGAG...

chr3      !      !      !      !      !  
      0      5      10     15     20     25     30

chr3      0      7      first      7.39     +

chr3      4      10     second     1e-8    -

## BED can also track Genome Annotations

### 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:

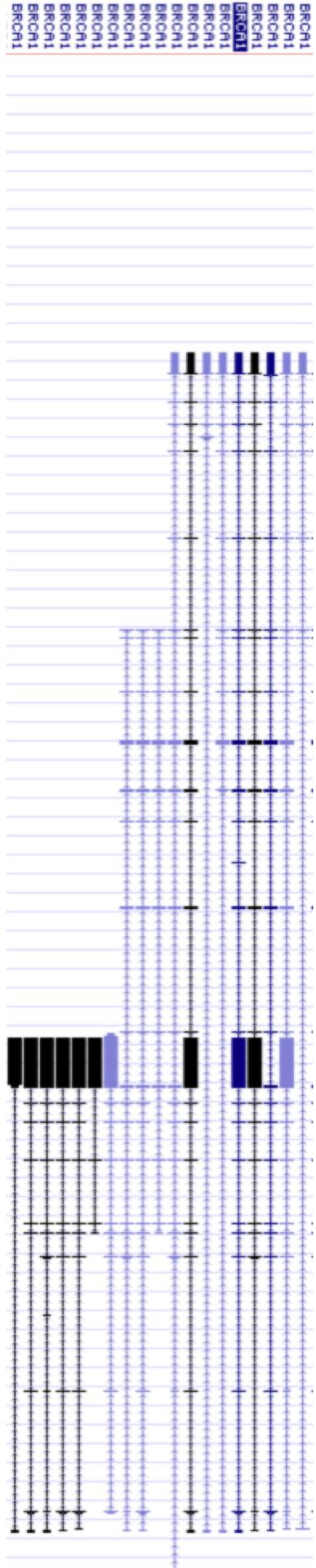
shade



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*.

BED12 Example



## GTF/GFF

**GFF (General Feature Format)**

**GTF (General Transfer Format)**

**Note: GFFv2 == GTF**

Fields:

1. Seqname – name of chrom
2. Source – name of data source
3. Feature - type
4. Start – Start position
5. End – End position
6. Score
7. Strand
8. Frame – Indicates that base of codon
9. Attribute – Semicolon list

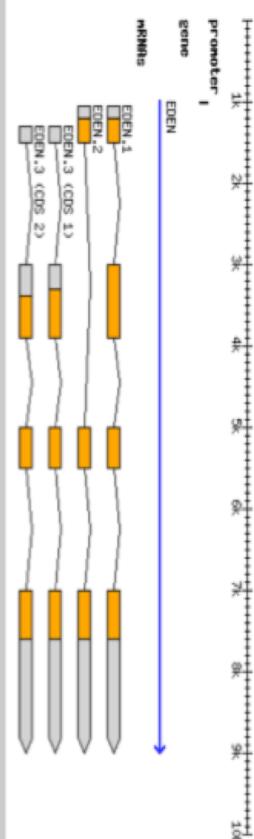
chr22	TeleGene	enhancer	10000000	10001000	500	+	.	touch1
chr22	TeleGene	promoter	10010000	10010100	900	+	.	touch1
chr22	TeleGene	promoter	10020000	10025000	800	-	.	touch2

**Note that the start and end coordinates are 1-based versus 0-based BED format**

## GFF Example

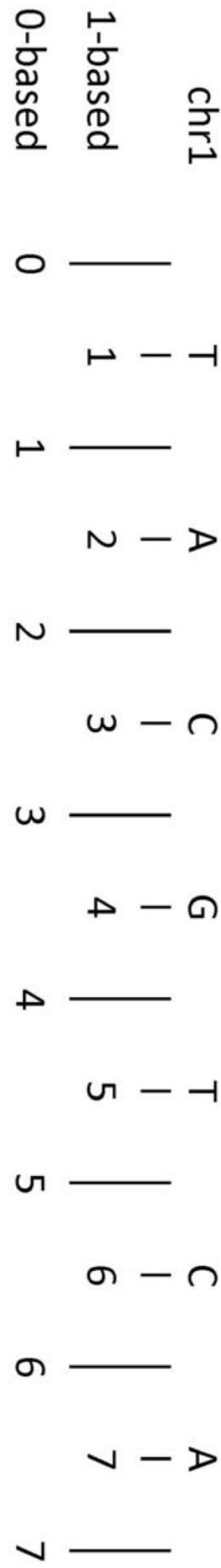
### GFF Example

Gene “EDEN” with 3 alternatively spliced transcripts, isoform 3 has two alternative translation start sites



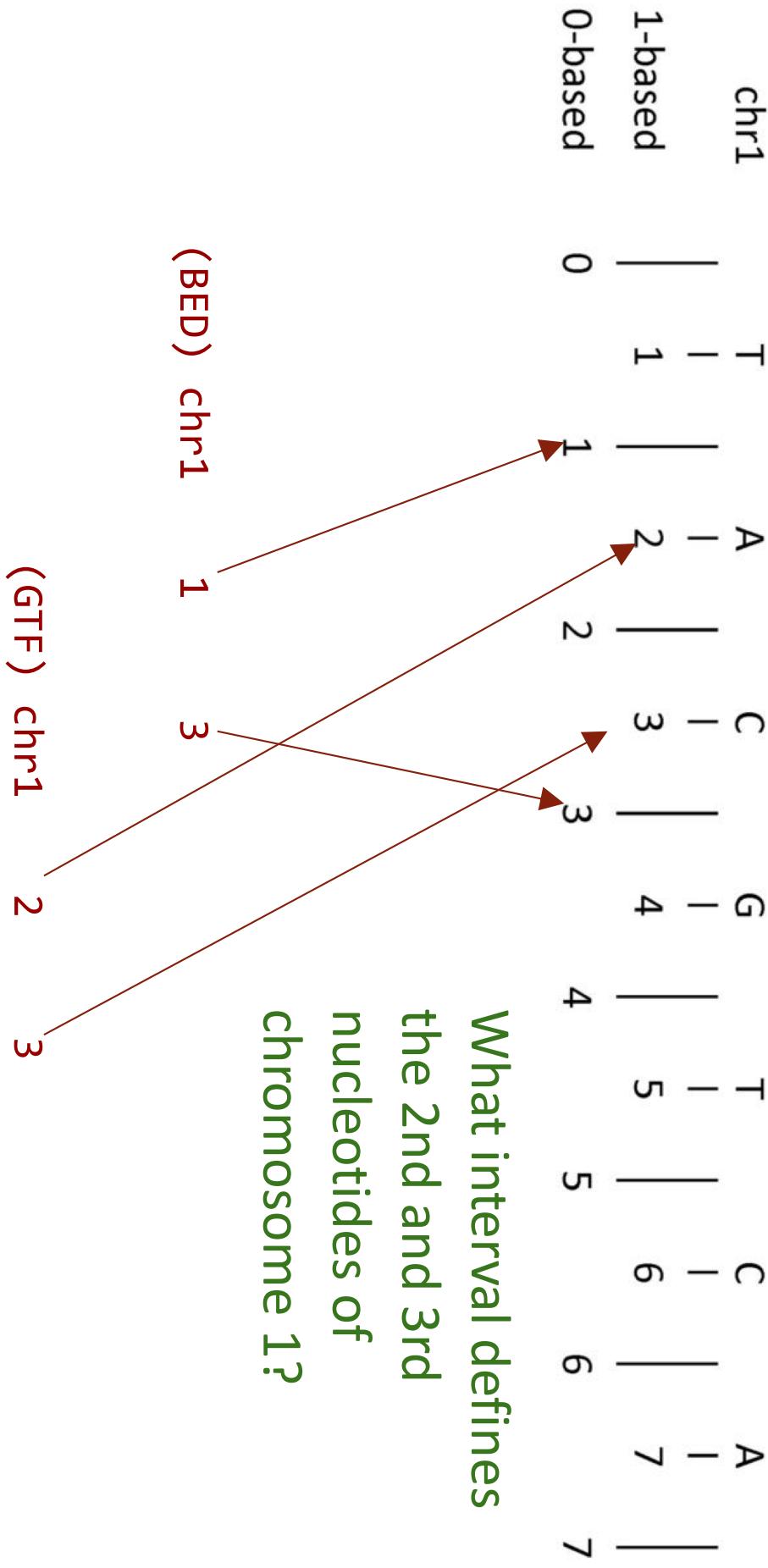
```
##gff-version 3
##sequence-region ctg123 1 1497228
ctg123 . gene 1000 9000 - + .
ID=tfb00001;Name=EDEN
ctg123 . TF_binding_site 1000 1012 - + .
ID=tfb00001;Parent=gene00001
ctg123 . mRNA 1050 9000 - + .
ID=mRNA00001;Parent=gene00001;Name=EDEN.1
ctg123 . mRNA 1050 9000 - + .
ID=mRNA00002;Parent=gene00001;Name=EDEN.2
ctg123 . mRNA 1300 9000 - + .
ID=mRNA00003;Parent=gene00001;Name=EDEN.3
ctg123 . exon 1300 1500 - + .
ID=exon00001;Parent=mRNA00003
ctg123 . exon 1050 1500 - + .
ID=exon00002;Parent=mRNA00001;mRNA00002
ctg123 . exon 3000 3902 - + .
ID=exon00003;Parent=mRNA00001;mRNA00003
ctg123 . exon 5000 - + .
ID=exon00004;Parent=mRNA00002;mRNA00003
ctg123 . exon 7000 9000 - + .
ID=exon00005;Parent=mRNA00001;mRNA00002;mRNA00003
ctg123 . CDS 1201 1500 - + 0
ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
ctg123 . CDS 3000 3902 - + 0
ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
ctg123 . CDS 5000 5500 - + 0
ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
ctg123 . CDS 7000 7600 - + 0
ID=cds00001;Parent=mRNA00002;Name=edenprotein.2
ctg123 . CDS 1201 1500 - + 0
ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
ctg123 . CDS 5000 5500 - + 0
ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
ctg123 . CDS 7000 7600 - + 0
ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
ctg123 . CDS 3301 3902 - + 0
ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
ctg123 . CDS 5000 5500 - + 1
ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
ctg123 . CDS 7000 7600 - + 1
ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
ctg123 . CDS 3391 3902 - + 0
ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
ctg123 . CDS 5000 5500 - + 1
ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
ctg123 . CDS 7000 7600 - + 1
ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
```

## How does BED and GTF define a genomic range?



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

## How does BED and GTF define a genomic range?



## Not all file formats are created Equal

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

**GFF:** 1-based, closed

**SAM:** 1-based, closed

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

**VCF:** 1-based, closed

...

