## 利用 Python 进行生信分析 | 基础 + 实战

——Python for Bioinformatics



赵华男

2022 一起加油

# 项目三 基因注释文件处理

## 基因注释与 GTF 和 GFF 文件的介绍

## 基因注释与 GTF 和 GFF 文件的介绍

- GFF 和 GTF 是两种最常用的数据库注释格式:
  - GFF 全称为 general feature format, 主要用来注释基因组
  - GTF 全称为 gene transfer format, 主要用来注释基因
- 注释文件的用途:
  - 在生物信息学分析中, 我们不但需要参考基因组信息 (FASTA) 和二代测序数据 (FASTQ) 来进行测序数据比对回贴
  - 还需要与之对应的注释信息 (GFF, GTF) 来进行下游分析, 比如常见的 RNA-seq (transcript) 和 ChIP (gene)
- GTF 是在 GFF 的基础上发展而来:
  - 本质上都是 TSV 文件 (以 TAB 制表符分隔的文本文件)
  - 都是 9 列文件 内容也比较接近
  - GFF 能够包含的信息更多更全, 可以包含染色体, 基因, 转录本的信息
  - 而 GTF 主要用来描述基因和转录本的信息
- 相互转化:
  - 如使用 Cufflinks 软件的的 gffread 命令
  - 我们自己写一个!
- 查看官方描述 (下面的链接)(GFF2 已经弃用, 特指 GFF3)

http://www.gmod.org/wiki/GFF3

### **Demo: GTF**

```
File: hg38 refseg from ucsc.rm XM XR.fix name.gtf
chr1
        hg38 ncbiRefSeq exon
                                 67092176
                                             67093604
                                                         0.000000
                                                                              gene id "Clorf141"; transcript id "NR 075077.1";
                                 67096252
                                             67096321
                                                         0.00000
chr1
        hq38 ncbiRefSeq exon
                                                                              gene id "Clorfoil"; transcript id "NR 075077.1";
cur1
        hq38 nefSeq emn
                                 671/1238
                                             67103382
                                                         0.000000
                                                                              gene id "Clorf141"; transcript id "NR 075077.1";
                                 67111577
                                                         0.000000
chr1
        hq38 ncbiRefSeq exon
                                             67111644
                                                                              gene id "Clorf141"; transcript id "NR 075077.1";
chr1
        hg38 ncbiRefSeg exon
                                 67113614
                                             67113756
                                                         0.000000
                                                                              gene id "Clorf141": transcript id "NR 075077.1":
        hg38 ncbiRefSeg exon
                                 67115352
                                                         0.000000
                                                                              gene id "Clorf141": transcript id "NR 075077.1":
chr1
                                             67115464
                                 67125752
                                                         0.000000
chr1
        hq38 ncbiRefSeq exon
                                             67125909
                                                                              gene id "Clorf141"; transcript id "NR 075077.1";
chr1
        hg38 ncbiRefSeg exon
                                67127166
                                             67127257
                                                         0.000000
                                                                              gene id "Clorf141": transcript id "NR 075077.1":
chr1
        hg38 ncbiRefSeg exon
                                 67131142
                                             67131227
                                                         0.000000
                                                                              gene id "Clorf141": transcript id "NR 075077.1":
chr1
        hg38 ncbiRefSeg exon
                                 67134930
                                             67134971
                                                         0.000000
                                                                              gene id "Clorf141": transcript id "NR 075077.1":
chr1
        hg38 ncbiRefSeg stop codon 67093580
                                                 67093582
                                                             0.000000
                                                                                  gene id "Clorf141": transcript id "NM 001276352.1":
chr1
        hg38 ncbiRefSeg CDS 67093583
                                         67093604
                                                     0.000000
                                                                          gene id "Clorf141"; transcript id "NM 001276352.1";
chr1
        hg38 ncbiRefSeg exon
                                 67092176
                                             67093604
                                                         0.000000
                                                                              gene id "Clorf141"; transcript id "NM 001276352.1";
chr1
        hg38 ncbiRefSeg CDS 67096252
                                         67096321
                                                     0.000000
                                                                          gene id "Clorf141"; transcript id "NM 001276352.1";
                                                         0.000000
chr1
        hq38 ncbiRefSeq exon
                                67096252
                                             67096321
                                                                              gene id "Clorf141"; transcript id "NM 001276352.1";
chr1
        hq38 ncbiRefSeq CDS 67103238
                                         67103382
                                                     0.000000
                                                                          gene id "Clorf141"; transcript id "NM 001276352.1";
chr1
        hq38 ncbiRefSeq exon
                                67103238
                                             67103382
                                                         0.000000
                                                                              gene id "Clorf141"; transcript id "NM 001276352.1";
                                                                          gene id "Clorf141"; transcript id "NM 001276352.1";
chr1
        hq38 ncbiRefSeq CDS 67111577
                                         67111644
                                                     0.000000
chr1
                                 67111577
                                             67111644
                                                         0.000000
                                                                              gene id "Clorf141"; transcript id "NM 001276352.1";
        hg38 ncbiRefSeg exon
        hg38 ncbiRefSeg CDS 67115352
                                         67115464
                                                                          gene id "Clorf141": transcript id "NM 001276352.1":
chr1
                                                     0.000000
chr1
        hg38 ncbiRefSeg exon
                                67115352
                                             67115464
                                                         0.000000
                                                                              gene id "Clorf141": transcript id "NM 001276352.1":
        hg38 ncbiRefSeg CDS 67125752
                                         67125909
                                                     0.000000
                                                                          gene id "Clorf141": transcript id "NM 001276352.1":
chr1
        hg38 ncbiRefSeg exon
                                67125752
                                             67125909
                                                         0.000000
                                                                              gene id "Clorf141": transcript id "NM 001276352.1":
chr1
chr1
        hg38 ncbiRefSeg CDS 67127166
                                         67127240
                                                     0.000000
                                                                          gene id "Clorf141": transcript id "NM 001276352.1":
chr1
        hg38 ncbiRefSeg start codon 67127238
                                                 67127240
                                                             0.000000
                                                                                  gene id "Clorf141": transcript id "NM 001276352.1":
chr1
        hg38 ncbiRefSeg exon
                                 67127166
                                             67127257
                                                         0.000000
                                                                              gene id "Clorf141": transcript id "NM 001276352.1":
chr1
        hq38 ncbiRefSeq exon
                                 67131142
                                             67131227
                                                         0.000000
                                                                              gene id "Clorf141"; transcript id "NM 001276352.1";
                                67134930
chr1
        hg38 ncbiRefSeg exon
                                             67134971
                                                         0.000000
                                                                              gene id "Clorf141"; transcript id "NM 001276352.1";
chr1
        hq38 ncbiRefSeq stop codon 67093005
                                                 67093007
                                                             0.000000
                                                                                  gene id "Clorf141"; transcript id "NM 001276351.1";
                                                                          gene id "Clorf141"; transcript id "NM 001276351.1";
chr1
        hq38 ncbiRefSeq CDS 67093008
                                         67093604
                                                     0.000000
                                 67092176
                                             67093604
                                                         0.000000
chr1
        hg38 ncbiRefSeq exon
                                                                              gene id "Clorf141"; transcript id "NM 001276351.1";
        hq38 ncbiRefSeq CDS 67095235
                                         67095421
                                                     0.000000
                                                                          gene id "Clorf141"; transcript id "NM 001276351.1";
chr1
                                 67095235
                                             67095421
                                                         0.000000
                                                                              gene id "Clorf141"; transcript id "NM 001276351.1";
chr1
        hq38 ncbiRefSeg exon
chr1
        hg38 ncbiRefSeg CDS 67096252
                                         67096321
                                                     0.000000
                                                                          gene id "Clorf141": transcript id "NM 001276351.1":
chr1
        hg38 ncbiRefSeg exon
                                 67096252
                                             67096321
                                                         0.000000
                                                                              gene id "Clorf141": transcript id "NM 001276351.1":
```

#### Demo: another GTF

```
##description: evidence-based annotation of the human genome (GRCh38), version 29 (Ensembl 94)
##provider: GENCODE
##contact: gencode-help#ebi.ac.uk
##format: gtf
##date: 2018-08-30
       HAVANA gene
                                                                gene id "ENSC00000223972.5"; gene type "transcribed unprocessed pseudogene"; gene name "DDX11L1"; level 2; bayana gen
                                                                        gene id "ENSG00000223972.5": transcript id "ENST00000456328.2": gene type "transcribed unprocessed pseudogene
       HAVANA
                                                                gene id "ENSG00000223972.5"; transcript id "ENST00000456328.2"; gene type "transcribed unprocessed pseudogene"; gene
                                                                gene id "ENSG00000223972.5": transcript id "ENST00000456328.2"; gene type "transcribed unprocessed pseudogene"; gene
                                                                gene id "ENSG00000223972.5": transcript id "ENST00000456328.2": gene type "transcribed unprocessed pseudogene": gene
       HAVANA
                                                                        gene id "ENSG00000223972.5"; transcript id "ENST00000450305.2"; gene type "transcribed unprocessed pseudogene
       HAVANA
                                                                gene id "ENSG00000223972.5"; transcript id "ENST00000450305.2"; gene type "transcribed unprocessed pseudogene"; gene
       HAVANA exon
                                                                gene id "ENSG00000223972.5"; transcript id "ENST00000450305.2"; gene type "transcribed unprocessed pseudogene"; gene
       HAVANA
               exon
                                                                gene id "ENSG00000223972.5": transcript id "ENST00000450305.2": gene type "transcribed unprocessed pseudogene": gene
       HAVANA
                                                                gene id "ENSG00000223972.5": transcript id "ENST00000450305.2": gene type "transcribed unprocessed pseudogene": gene
                                                                gene id "ENSG00000223972.5"; transcript id "ENST00000450305.2"; gene type "transcriped unprocessed pseudogene"; gene
       HAVANA
                                                                gene id "ENSG00000223972.5"; transcript id "ENST00000450305.2"; gene type "transcribed unprocessed pseudogene"; gene
       HAVANA
               gene
                                                                gene id "ENSC00000227232.5"; gene type "upprocessed pseudogene"; gene name "WASH7P"; level 2; hayana gene "OTTHUMG000
                                                                        gene id "ENSG00000227232.5"; transcript id "ENST00000488147.1"; gene type "unprocessed pseudogene"; gene pame
       HAVANA exon
                                                                gene id "ENSG00000227232.5"; transcript id "ENST00000488147.1"; gene type "unprocessed pseudogene"; gene name "WASH7I
       HAVANA
                                                                gene id "ENSG00000227232.5": transcript id "ENST00000488147.1": gene type "unprocessed pseudogene": gene name "WASH71
                                                                gene id "ENSG00000227232.5"; transcript id "ENST00000488147.1"; gene type "unprocessed pseudogene"; gene name
       MAVANA
               evon
                                                                gene id "ENSG00000227232.5"; transcript id "ENST00000488147.1"; gene type "unprocessed pseudogene"; gene name "WASH71
       HAVANA
               exon
                                                                gene id "ENSG00000227232.5": transcript id "ENST00000488147.1": gene type "upprocessed pseudogene": gene name
               exon
                                                                gene id "ENSG00000227232.5"; transcript id "ENST00000488147.1"; gene type "unprocessed pseudogene"; gene name
               exon
                                                                gene id "ENSG00000227232.5": transcript id "ENST00000488147.1": gene type "unprocessed pseudogene": gene name
       HAVANA evon
                                                                gene id "ENSG00000227232.5": transcript id "ENST00000488147.1": gene type "unprocessed pseudogene": gene name "WASH71
                                                                gene id "ENSG00000227232.5": transcript id "ENST00000488147.1": gene type "unprocessed pseudogene": gene name
       MAVANA
                                                                gene id "ENSG00000227232.5"; transcript id "ENST00000488147.1"; gene type "unprocessed pseudogene"; gene name "WASH7E
       HAVANA
               exon
                                                                gene id "ENSG00000227232.5"; transcript id "ENST00000488147.1"; gene type "unprocessed pseudogene"; gene name "WASH71
       ENSEMBL gene
                                                                gene id "ENSG00000278267.1"; gene type "miRNA"; gene name "MIR6859-1"; level 3;
       ENSEMBL transcript
                                                                        gene id "ENSG00000278267.1"; transcript id "ENST00000619216.1"; gene type "miRNA"; gene name "MIR6859-1"; tra
       ENSEMBL exon
                                                                gene id "ENSG00000278267.1": transcript id "ENST00000619216.1": gene type "miRNA": gene name "MIR6859-1": transcript
                       29554
                                                                gene id "ENSG00000243485.5"; gene type "lingRNA"; gene name "MIR1302-2HG"; level 2; tag "ngRNA host"; hayana gene "OT
       HAVANA gene
       HAVANA transcript
                                                                        gene id "ENSG00000243485.5"; transcript id "ENST00000473358.1"; gene type "lincRNA"; gene name "MIR1302-2HG";
       HAVANA evon
                                                                gene id "ENSC00000243485.5": transcript id "ENST00000473358.1": gene type "linchal": gene name "MIR1302-2HG": transcr
                                                                        "FNSC00000243485 5", transcript id "FNST00000473358 1", gane type "lincDNA", gane name "MTD1302-2MG", transcript
```

### Demo: GFF3

```
##qff-version 3.1.26
##sequence-region ctg123 1 1497228
cta123 . gene
                          1000
                                9000
                                               ID=gene00001:Name=EDEN
ctg123 . TF_binding_site 1000
                                               ID=tfbs00001;Parent=gene00001
                                1012
                          1050
                                9000
                                               ID=mRNA00001; Parent=gene00001; Name=EDEN.1
ctq123 . mRNA
ctq123 . mRNA
                          1050
                                9000
                                               ID=mRNA00002:Parent=gene00001:Name=EDEN.2
                                         +
cta123 . mRNA
                          1300
                                9000
                                               ID=mRNA00003: Parent=gene00001: Name=EDEN.3
ctg123 . exon
                          1300
                                1500
                                               ID=exon00001:Parent=mRNA00003
cta123 . exon
                          1050
                                1500
                                               ID=exon00002:Parent=mRNA00001.mRNA00002
ctq123 . exon
                          3000
                                3902
                                               ID=exon00003; Parent=mRNA00001, mRNA00003
ctq123 . exon
                          5000
                                5500
                                               ID=exon00004: Parent=mRNA00001.mRNA00002.mRNA00003
cta123 . exon
                          7000
                                9000
                                               ID=exon00005: Parent=mRNA00001.mRNA00002.mRNA00003
cta123 . CDS
                                500 6
                                               ID=cds000001:Parent=mRNA00001:Name=edenprotein.1
                          171
                                               ID=2ds00001:Parent=mRNA00001;Name=edenprotein.1
ctg132 (3)
                          3000
                                3902
cta123 . CDS
                          5000
                                5500
                                               ID=cds00001:Parent=mRNA00001:Name=edenprotein.1
ctq123 . CDS
                                7600
                                               ID=cds00001:Parent=mRNA00001:Name=edenprotein.1
                          7000
cta123 . CDS
                          1201
                                1500
                                               ID=cds00002:Parent=mRNA00002:Name=edenprotein.2
ctq123 . CDS
                          5000
                                5500
                                               ID=cds00002; Parent=mRNA00002; Name=edenprotein.2
                                         +
                                               ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
cta123 . CDS
                          7000
                                7600
                                3902
                                               ID=cds00003:Parent=mRNA00003:Name=edenprotein.3
cta123 . CDS
                          3301
```

### • 1. seq\_id

● 序列的编号,一般为 chr 或者 scanfold 编号,每条染色体拥有一个唯一的 ID。

#### 2. source

● 注释的来源,代表基因结构的来源,可以是数据库的名称,比如来自 RefSeq 数据库,也可以是软件的名称,比如用 GeneScan 软件预测得到, 当然。也可以为空,用""点号填充。

### • 3. type

- 代表区间对应的特征类型, 在 GTF 中, 常见的类型如下:
  - Gene, cDNA, mRNA, 5UTR, 3UTR, exon, CDS, start\_codon, stop\_codon

#### • 4. start

● 该基因或转录本在参考序列上的起始位置。

#### • 5. end

● 该基因或转录本在参考序列上的终止位置。

#### 6. score

- 得分, 软件提供了统计值, 是注释信息可能性的说明, 可以是序列相似性比对时的 E-values 值或者基因预测是的 P-values 值
- ""表示为空。

#### • 7. strand

- 代表正负链的信息,+表示正链、-表示负链、?表示不清楚正负链的信息
- 当正负链信息没有意义时,可以用""填充。

### Info: 以 GTF 为例

### • 8. phase

- 仅对注释类型为 "CDS" 有效,表示起始编码的位置,有效值为 0、1、2
  - 对于编码蛋白质的 CDS 来说,本列指定下一个密码子开始的位置。每 3 个核苷酸翻译一个氨基酸,从 0 开始, CDS 的起始位置,除以 3.余数就是这个值,表示到达下一个密码子需要跳过的碱基个数。
  - 0表示该编码框的第一个密码子第一个碱基位于其 5' 末端
  - 1表示该编码框的第一个密码子的第一个碱基位于该编码区外
  - 2表示该编码框的第一个密码子的第一、二个碱基位于该编码区外
- 如果Feature 为 CDS 时, 必须指明具体值!

### ● 9. attributes (更正!)

- 一个包含众多属性的列表
  - GFF3: 格式为 "标签 = 值" | "(tag=value)", 如果一个序列元件没有 Parent 属性, 说明他的父元件就是 scaffold 或者 chromosome
  - GTF: 格式为 "标签值;" | "tag value;", 每个特征之后都要有分号 (包括最后一个特征)
  - GTF: 其内容必须包括 gene id 和 transcript id. 其 value 可以为空 ""

## GTF 和 GFF 文件操作

### ● ensembl 示例文件下载

- https://www.ensembl.org/
- All genomes (Search 下方)
- 选择 Human
- Gene annotation
- Download GTF or GFF3 files for genes, cDNAs, ncRNA, proteins

#### Index of /pub/release-107/gtf/homo\_sapiens

Name	Last modified	Size I	Description
Parent Directory		-	
? CHECKSUMS	2022-06-05 20:13	225	
Homo sapiens.GRCh38.107.abinitio.gtf.gz	2022-06-04 11:52	3.3M	
Homo_sapiens.GRCh38.107.chr,gtf.gz	2022-06-04 11:45	50M	
Homo sapiens.GRCh38.107.chr patch hapl scaff.gtf.gz	2022-06-04 11:49	54M	
Homo_sapiens.GRCh38.107.gtf.gz	2022-06-04 11:45	50M	_
README	2022-06-04 11:50	9.9K	_

# GTF 和 GFF 文件操作

### Demo: GFF3 feature

```
import pandas as pd
    import numpy as np
    df = pd.read csv(
             'data/Homo_sapiens.GRCh38.107.gff3.gz',
             sep="\t",
             comment="#",
             names = Annotation . REQUIRED_COLUMNS .
             skipinitialspace=True,
             skip_blank_lines=True,
10
             on bad lines='skip'.
11
12
             # chunksize=chunksize,
13
             engine="c".
             dtype={
14
15
                 "start": np.int64.
16
                 "end": np.int64.
                 "score": np.float32,
17
18
                 "segname": str.
19
             }.
             na values=".".
20
21
```

GTF 和 GFF 文件操作

# GTF 和 GFF 文件操作

### **Demo: GFF3 feature**

```
df.feature.unique()
    array(['chromosome', 'biological_region', 'pseudogene', 'lnc_RNA', 'exon',
           'pseudogenic_transcript', 'ncRNA_gene', 'miRNA', 'gene', 'mRNA',
           'five_prime_UTR', 'CDS', 'three_prime_UTR', 'snRNA', 'ncRNA',
           'unconfirmed_transcript', 'snoRNA', 'scRNA', 'rRNA',
           'V gene segment', 'D gene segment', 'J gene segment',
           'C gene segment', 'transcript', 'scaffold', 'tRNA'], dtype=object)
 8
    # feature required:
9
10
    # - "CDS".
    # - "start codon" # 简化, 先不考虑, 其实是CDS前三个碱基
11
    # - "stop codon" # 简化, 先不考虑
12
13
14
    # optional:
    # - "5UTR".
15
16
    # - "3UTR",
17
    # - "*RNA"
```

## 练习

```
|- 仅在必要时才对 self 或 cls 注释(实战课三当中进行讲解)
  - cls + Oclassmethod
   - GFF与GTF文件进行转换
   - 统计各条染色体的基因密度
   - 获得基因的的TSS. TES及启动子坐标 (GFF3)
   - 计算全基因组可转录区域长度及所占基因组比例 (GFF3, GTF)
   - 计算基因的转录长度、外显子数目及翻译区长度 (GFF3)
     - 计算基因的转录长度
        - 统计平均长度, 中位数长度
10
     - 计算基因的外显子(exon)个数
11
        - 统计平均个数,中位数个数
12
13
        - 统计最多exon的基因,最少exon的gene
14
   - 程序 CLI 化
15
     - sys.argv
16
      - argparse
```

程序 CLI 化

# 程序 CLI 化

# 程序 CLI 化

项目实战三基因注释文件处理

程序 CLI 化

Q & A!

98 / 98

程序 CLI 化

# Thank you!