



Chapter 04. Visualization

Python Programming for Bioinformatics

Robert C. Chi

Agenda

- **Introduction**
- **Visualize Genome**
- **Visualize Chromosome**
- **Summary**

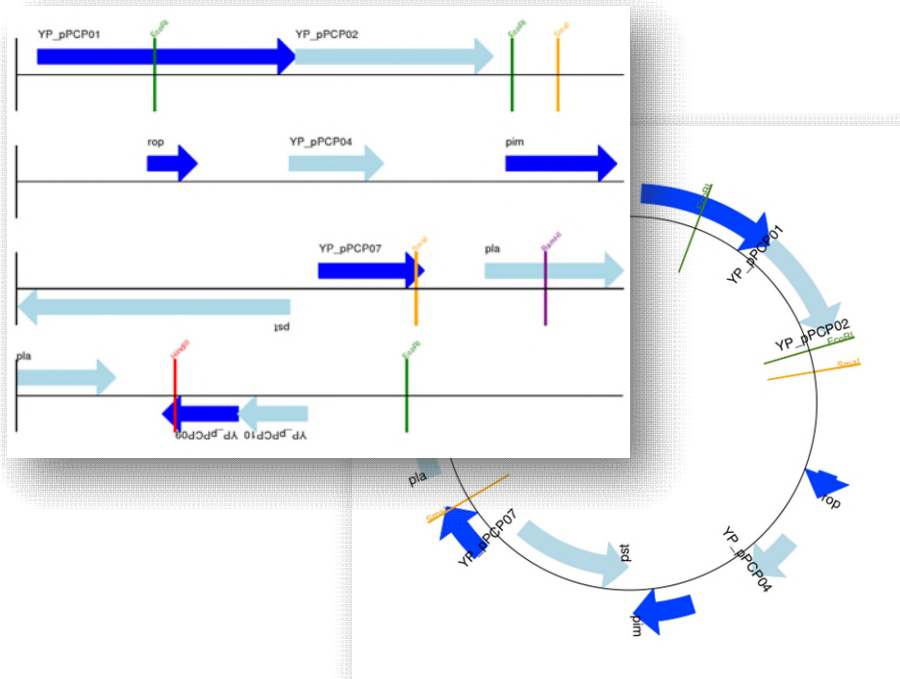




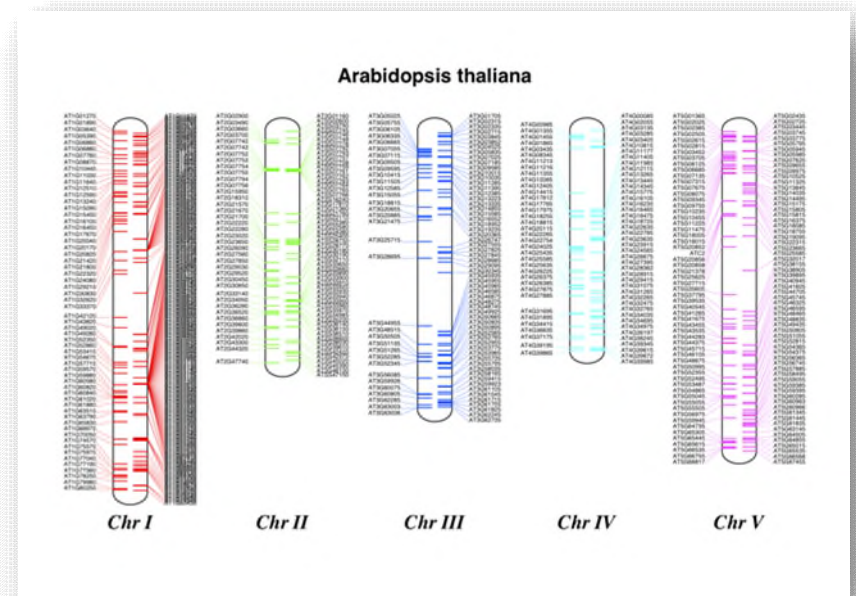
INTRODUCTION

What can it do?

- Visualize **Genome**



- Visualize **Chromosome**



What Packages are Required?

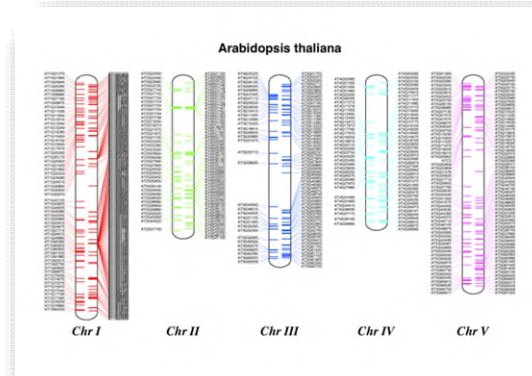
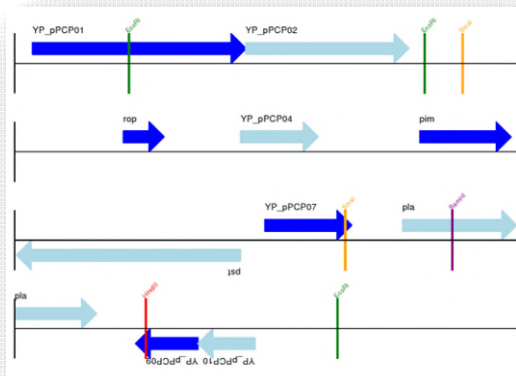
ReportLab



```
pip install reportlab  
pip install biopython
```

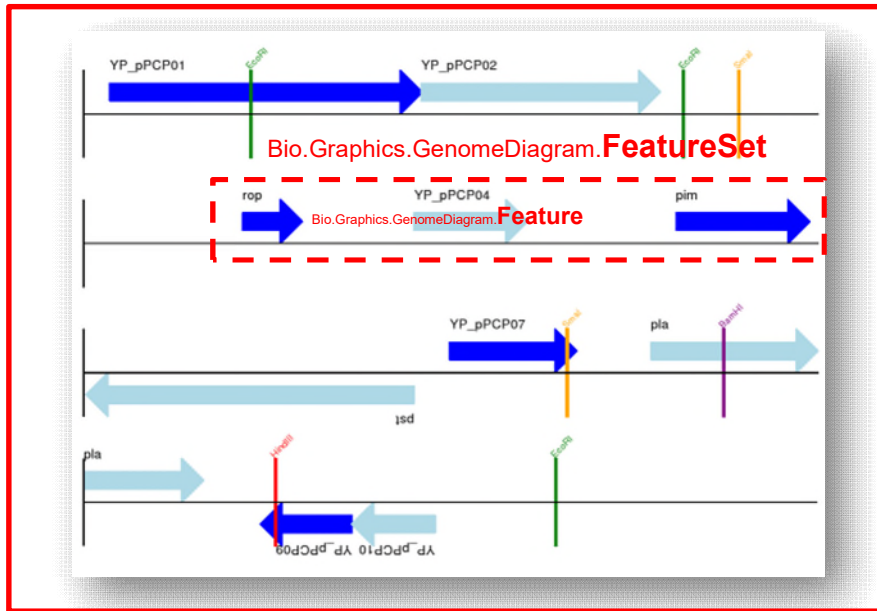
Bio.Graphics.GenomeDiagram

Bio.Graphics.BasicChromosome



Data Structure of GenomeDiagram

Bio.Graphics.GenomeDiagram.**Diagram**



Bio.Graphics.GenomeDiagram.**Track**

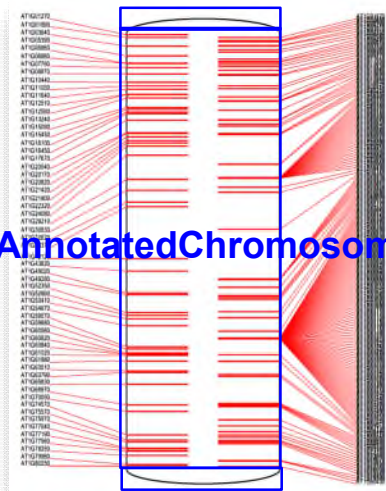
Bio.Graphics.GenomeDiagram.**Track**

Bio.Graphics.GenomeDiagram.**Track**

Bio.Graphics.GenomeDiagram.**Track**

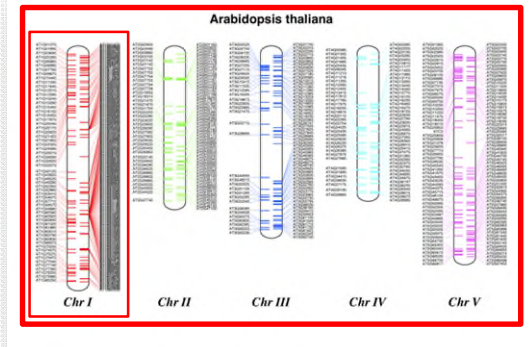
Data Structure of BasicChromosome

Bio.Graphics.BasicChromosome. **TelomereSegment**



Bio.Graphics.BasicChromosome. **AnnotatedChromosomeSegment**

Bio.Graphics.BasicChromosome. **Organism**



Bio.Graphics.BasicChromosome. **TelomereSegment**

Bio.Graphics.BasicChromosome. **Chromosome**



VISUALIZE GENOME

Environment Preparation

- **Install Packages**

```
1 !pip install reportlab
2 !pip install biopython
```

- **Download the File**

```
1 # GenBank file for the pPCP1 plasmid from Yersinia pestis biovar Microtus
2 import os
3 if not os.path.isfile("NC_005816.gb"):
4     os.system("wget https://raw.githubusercontent.com/biopython/biopython/master/Tests/GenBank/NC\_005816.gb")
```

Practice

- **Visualize Genome: Environment Preparation**

- **Write** and **Run** the following codes on a Colab page called “[VisualizeGenome.ipynb](#)”:

Install Biopython

```
[9] 1 !pip install reportlab
    2 !pip install biopython
```

Download the Related File

```
[6] 1 # GenBank file for the pPCP1 plasmid from Yersinia pestis biovar Microtus
    2 import os
    3 if not os.path.isfile("NC_005816.gb"):
    4     os.system("wget https://raw.githubusercontent.com/biopython/biopython/master/Tests/GenBank/NC_005816.gb")
```

(Solution [URL](#) of this Practice)



Get Data from the File

```
1 # Create a list to store all parsed SeqRecords
2 seq_records = []
3
4 # Parse one Record a time and add into seq_records
5 from Bio import SeqIO
6 for rec in SeqIO.parse("NC_005816.gb", "genbank"):
7     seq_records.append(rec)
8
9 # Show how many records have been parsed
10 print("Total Number of Records:", len(seq_records))
11 print()
12
13 # Show the first record to prove the parsing was successful
14 print("--- First Record ---")
15 print(seq_records[0])
```



```
Total Number of Records: 1
--- First Record ---
ID: NC_005816.1
Name: NC_005816
Description: Yersinia pestis biovar Microtus str. 91001 plasmid pPCP1, complete sequence
Database cross-references: Project:58037
Number of features: 41
/molecule_type=DNA
/topology=circular
/data_file_division=BCT
/date=21-JUL-2008
/accessions=['NC_005816']
/sequence_version=1
/gi=45478711
/keywords=[]
/source=Yersinia pestis biovar Microtus str. 91001
/organism=Yersinia pestis biovar Microtus str. 91001
/taxonomy=['Bacteria', 'Proteobacteria', 'Gammaproteobacteria', ... 'Yersinia']
/references=[Reference(title='Genetics of metabolic variations ...'),
             Reference(title='Complete genome sequence ...'),
             Reference(title='Direct Submission', ...),
             Reference(title='Direct Submission', ...)]
/comment=PROVISIONAL REFSEQ: This record has not yet been subject to final
NCBI review. The reference sequence was derived from AE017046.
COMPLETENESS: full length.
Seq('TGTAACGAACGGTGCAATAGTGATCCACCCCAACGCCTGAAATCAGATCCAGG...CTG')
```

Practice

- **Visualize Genome: Get Data**

- **Write** and **Run** the following codes on a Colab page called “**VisualizeGenome.ipynb**”:

```
1 # Create a list to store all parsed SeqRecords
2 seq_records = []
3
4 # Parse one Record a time and add into seq_records
5 from Bio import SeqIO
6 for rec in SeqIO.parse("NC_005816.gb", "genbank"):
7     seq_records.append(rec)
8
9 # Show how many records have been parsed
10 print("Total Number of Records:", len(seq_records))
11 print()
12
13 # Show the first record to prove the parsing was successful
14 print("--- First Record ---")
15 print(seq_records[0])
```

(Solution [URL](#) of this Practice)



Prepare Data Structures

```
1 # Get the first record as example
2 rec = seq_records[0]
3
4 # Create a Diagram
5 from Bio.Graphics import GenomeDiagram
6 diag = GenomeDiagram.Diagram(name=rec.id)
7
8 # Create a Track
9 trac = diag.new_track(1, name="Annotated Features")
10
11 # Create a Feature Set
12 feat_set = trac.new_set()
```

(1) Get the first record

(2) Create the Diagram

- name= String. Identifier for the diagram.

(3) Create the First Track

- name= String describing the track.

(4) Create the Feature Set

Practice

- **Visualize Genome: Create Data Structures**
 - **Write** and **Run** the following codes on a Colab page called “**VisualizeGenome.ipynb**”:

```
1 # Get the first record as example
2 rec = seq_records[0]
3
4 # Create a Diagram
5 from Bio.Graphics import GenomeDiagram
6 diag = GenomeDiagram.Diagram(name=rec.id)
7
8 # Create a Track
9 trac = diag.new_track(1, name="Annotated Features")
10
11 # Create a Feature Set
12 feat_set = trac.new_set()
```

(Solution [URL](#) of this Practice)



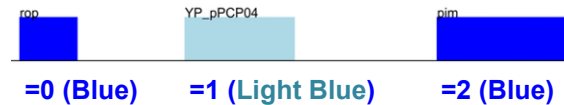
Create each Features

```
1 # Find the Features == "gene" in the first record
2 from reportlab.lib import colors
3
4 for feat in rec.features:
5     # Skip all the types not equal to "gene"
6     if feat.type != "gene":
7         continue
8
9     # If the current index of Feature Set is even, set the color to Blue
10    # If the current index of Feature Set is odd, set the color to Light Blue
11    if len(feats_set) % 2 == 0:
12        color = colors.blue
13    else:
14        color = colors.lightblue
15
16    # Add this feature to the Feature Set
17    feat_set.add_feature(feat, color=color,
18                        label=True, label_size=14, label_angle=0)
```

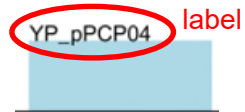
(1) Only fetch those features as "gene"

```
FEATURES   Location/Qualifiers
gene       87..1109
           /locus_tag="YP_pPCP01"
           /db_xref="GeneID:2767718"
```

(2) Even as "Blue", Odd as "Light Blue"



(3) Add Feature into Feature Set



Practice

- **Visualize Genome: Create each Feature**

- **Write** and **Run** the following codes on a Colab page called “**VisualizeGenome.ipynb**”:

```
1 # Find the Features == "gene" in the first record
2 from reportlab.lib import colors
3
4 for feat in rec.features:
5     # Skip all the types not equal to "gene"
6     if feat.type != "gene":
7         continue
8
9     # If the current index of Feature Set is even, set the color to Blue
10    # If the current index of Feature Set is odd, set the color to Light Blue
11    if len(feet_set) % 2 == 0:
12        color = colors.blue
13    else:
14        color = colors.lightblue
15
16    # Add this feature to the Feature Set
17    feat_set.add_feature(feat, color=color,
18                        label=True, label_size=14, label_angle=0)
```

(Solution [URL](#) of this Practice)

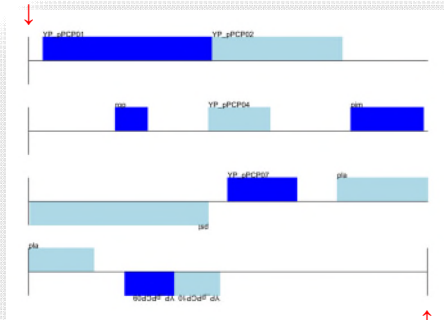


Output the Result

```
1 # Draw the Diagram with Linear Format as PNG File
2 diag.draw(format="linear", orientation="landscape", pagesize='A4',
3           fragments=4, start=0, end=len(rec))
4 diag.write("plasmid_linear.png", "PNG")
```

fragments=4

start=0

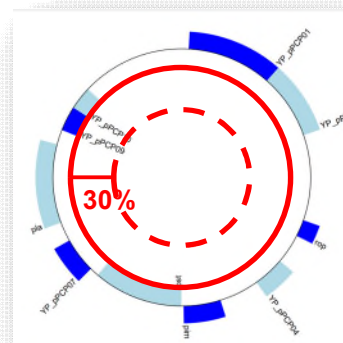


end=len(rec)

If the genome is
circular in reality

pixel x pixel

```
6 # Draw the Diagram with Circular Format as PDF File
7 from reportlab.lib.units import cm
8 diag.draw(format="circular", circular=True, pagesize=(20*cm,20*cm),
9           start=0, end=len(rec), circle_core=0.7)
10 diag.write("plasmid_circular.pdf", "PDF")
```



Practice

- **Visualize Genome: Output the Result**

- **Write** and **Run** the following codes on a Colab page called “[VisualizeGenome.ipynb](#)”:

```
1 # Draw the Diagram with Linear Format as PNG File
2 diag.draw(format="linear", orientation="landscape", pagesize='A4',
3           fragments=4, start=0, end=len(rec))
4 diag.write("plasmid_linear.png", "PNG")
5
6 # Draw the Diagram with Circular Format as PDF File
7 from reportlab.lib.units import cm
8 diag.draw(format="circular", circular=True, pagesize=(20*cm,20*cm),
9           start=0, end=len(rec), circle_core=0.7)
10 diag.write("plasmid_circular.pdf", "PDF")
```



(Solution [URL](#) of this Practice)



VISUALIZE CHROMOSOME

Environment Settings

Install Biopython

```
▶ 1 |pip install reportlab
   2 |pip install biopython
```

Download the Related File

```
[3] 1 # bio is a tool to download GenBank file
     2 |pip install bio
     3
     4 # Since the files are big, we only download the first chromosome
     5 |bio fetch NC_003070 > NC_003070.gbk
     6 #bio fetch NC_003071 > NC_003071.gbk
     7 #bio fetch NC_003074 > NC_003074.gbk
     8 #bio fetch NC_003075 > NC_003075.gbk
     9 #bio fetch NC_003076 > NC_003076.gbk
```


Practice

- **Visualize Chromosome: Environment Settings**
 - **Write** and **Run** the following codes on a Colab page called “**VisualizeChromosome.ipynb**”:

Install Biopython

```
1 !pip install reportlab
2 !pip install biopython
```

Download the Related File

```
[3] 1 # bio is a tool to download GenBank file
    2 !pip install bio
    3
    4 # Since the files are big, we only download the first chromosome
    5 !bio fetch NC_003070 > NC_003070.gb
    6 #bio fetch NC_003071 > NC_003071.gb
    7 #bio fetch NC_003074 > NC_003074.gb
    8 #bio fetch NC_003075 > NC_003075.gb
    9 #bio fetch NC_003076 > NC_003076.gb
```

(Solution [URL](#) of this Practice)



Settings of Parameters

```
1 from reportlab.lib.units import cm
2 from Bio import SeqIO
3 from Bio.Graphics import BasicChromosome
4
5 # The chromosome name vs. its file
6 entries = [("Chr I", "NC_003070.gbk"),
7 #           ("Chr II", "NC_003071.gbk"),
8 #           ("Chr III", "NC_003074.gbk"),
9 #           ("Chr IV", "NC_003075.gbk"),
10 #          ("Chr V", "NC_003076.gbk")
11          ]
12 # Set the scale factor of each chromosome
13 # Otherwise each chromosome will have same size visually
14 max_len = 30432563 #Could compute this
15 telomere_length = 1000000 #For illustration
16
17 chr_diagram = BasicChromosome.Organism()
18 chr_diagram.page_size = (29.7*cm, 21*cm) #A4 landscape
```

Practice

- **Visualize Chromosome: Settings of Parameters**
 - **Write** and **Run** the following codes on a Colab page called “[VisualizeChromosome.ipynb](#)”:

```
1 from reportlab.lib.units import cm
2 from Bio import SeqIO
3 from Bio.Graphics import BasicChromosome
4
5 # The chromosome name vs. its file
6 entries = [("Chr I", "NC_003070.gb"),
7 #           ("Chr II", "NC_003071.gb"),
8 #           ("Chr III", "NC_003074.gb"),
9 #           ("Chr IV", "NC_003075.gb"),
10 #          ("Chr V", "NC_003076.gb")
11          ]
12 # Set the scale factor of each chromosome
13 # Otherwise each chromosome will have same size visually
14 max_len = 30432563 #Could compute this
15 telomere_length = 1000000 #For illustration
16
17 chr_diagram = BasicChromosome.Organism()
18 chr_diagram.page_size = (29.7*cm, 21*cm) #A4 landscape
```



(Solution [URL](#) of this Practice)

Visualize Chromosomes

```
1 for index, (name, filename) in enumerate(entries):
2     # Fetch all records & features of this chromosome
3     record = SeqIO.read(filename, "genbank")
4     length = len(record)
5     features = [f for f in record.features if f.type=="tRNA"]
6
7     # Set the color of each chromosome
8     #Record an Artemis style integer color in the feature's qualifiers,
9     #1 = Black, 2 = Red, 3 = Green, 4 = blue, 5 =cyan, 6 = purple
10    for f in features: f.qualifiers["color"] = [index+2]
11
12    # Create a Chromosome
13    cur_chromosome = BasicChromosome.Chromosome(name)
14    # Set the scale to the MAXIMUM length plus the two telomeres in bp,
15    # Allows each chromosome to be of unequal length, reflecting their true length
16    cur_chromosome.scale_num = max_len + 2 * telomere_length
17
18    #Add an opening telomere
19    start = BasicChromosome.TelomereSegment()
20    start.scale = telomere_length
21    cur_chromosome.add(start)
22
23    #Add a body - again using bp as the scale length here.
24    body = BasicChromosome.AnnotatedChromosomeSegment(length, features)
25    body.scale = length
26    cur_chromosome.add(body)
27
28    #Add a closing telomere
29    end = BasicChromosome.TelomereSegment(inverted=True)
30    end.scale = telomere_length
31    cur_chromosome.add(end)
32
33    #This chromosome is done
34    chr_diagram.add(cur_chromosome)
35
36 chr_diagram.draw("tRNA_chrom.pdf", "Arabidopsis thaliana")
```


Practice

```
1 for index, (name, filename) in enumerate(entries):
2     # Fetch all records & features of this chromosome
3     record = SeqIO.read(filename, "genbank")
4     length = len(record)
5     features = [f for f in record.features if f.type=="tRNA"]
6
7     # Set the color of each chromosome
8     #Record an Artemis style integer color in the feature's qualifiers,
9     #1 = Black, 2 = Red, 3 = Green, 4 = blue, 5 =cyan, 6 = purple
10    for f in features: f.qualifiers["color"] = [index+2]
11
12    # Create a Chromosome
13    cur_chromosome = BasicChromosome.Chromosome(name)
14    # Set the scale to the MAXIMUM length plus the two telomeres in bp,
15    # Allows each chromosome to be of unequal length, reflecting their true length
16    cur_chromosome.scale_num = max_len + 2 * telomere_length
17
18    #Add an opening telomere
19    start = BasicChromosome.TelomereSegment()
20    start.scale = telomere_length
21    cur_chromosome.add(start)
22
23    #Add a body - again using bp as the scale length here.
24    body = BasicChromosome.AnnotatedChromosomeSegment(length, features)
25    body.scale = length
26    cur_chromosome.add(body)
27
28    #Add a closing telomere
29    end = BasicChromosome.TelomereSegment(inverted=True)
30    end.scale = telomere_length
31    cur_chromosome.add(end)
32
33    #This chromosome is done
34    chr_diagram.add(cur_chromosome)
35
36 chr_diagram.draw("tRNA_chrom.pdf", "Arabidopsis thaliana")
```

- **Visualize Chromosome**

- **Write** and **Run** the following codes on a Colab page called “**VisualizeChromosome.ipynb**”:

(Solution [URL](#) of this Practice)



Summary

- **Bio.Graphics.GenomeDiagram**
 - Bio.Graphics.GenomeDiagram.Diagram
 - Bio.Graphics.GenomeDiagram.Track
 - Bio.Graphics.GenomeDiagram.FeatureSet
 - Bio.Graphics.GenomeDiagram.Feature
- **Bio.Graphics.BasicChromosome**
 - Bio.Graphics.BasicChromosome.Organism
 - Bio.Graphics.BasicChromosome.Chromosome
 - Bio.Graphics.BasicChromosome.TelomereSegment
 - Bio.Graphics.BasicChromosome.AnnotatedChromosomeSegment

