Genomics Spring 2021

Units 3-4 Graded Homework

(For the following problems, submit just one master pdf file and submit it along with other necessary files via Blackboard. Grades will be given mainly based on the pdf file. Screen capture or copy & paste all the results into the pdf file.)

1. 1.25 pts. There are two attached files: znf214_mrna.txt and znf214_genomic.txt. Use Splign to find the mRNA and CDS coordinates in the genomic DNA.

    a.  0.625 pts: Report mRNA locations: **1..283, 284..430, 431..2668**



| Model 1 | Coverage 100.00% | CDS | 0.00% | Mismatches and indels 0 |
| | Overall 100.00% | In-frame | 0.00% | Exons (min/max/ave), bp 147 / 2238 / 889 |
| | Exon 100.00% | Primary transcript 2668 bp | | Introns (min/max/ave), bp 1154 / 17218 / 9186 |

Homo (+) sapiens zinc finger protein 214 (ZNF214), mRNA

1 ... 2668

chromosome:GRCh38:11:6998718:7020968:-1 (+)

614 ... 21651

Segments **Alignment**

**1** 2 3

```
   1 AAGACCTGGGACTTCCGTTTCGGTCCAGCCGGGCTGCGGCCATTGTTTGTGTGGACTGAGTGTTTTGGCC
     ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
 614 AAGACCTGGGACTTCCGTTTCGGTCCAGCCGGGCTGCGGCCATTGTTTGTGTGGACTGAGTGTTTTGGCC

  71 ATCCCGGTCCACTCTCACAGGCTCCGTTAAGTGACATGAACTCTCAGGAGGACTGAGCCAGAAGCGGACA
     ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
 684 ATCCCGGTCCACTCTCACAGGCTCCGTTAAGTGACATGAACTCTCAGGAGGACTGAGCCAGAAGCGGACA

 141 GGGCAAGACGAGCTGTGCTTGAAGGAAGAGGGGCAGAAAATCGAGGGTCAGGGACTGAGAAGCTACTCCG
     ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
 754 GGGCAAGACGAGCTGTGCTTGAAGGAAGAGGGGCAGAAAATCGAGGGTCAGGGACTGAGAAGCTACTCCG

 211 GTTTAGAAACCCCAGAGACACCCGTGTAGATGGGTACATCACGGCTTCTCTCCCACGTTAAGACTTAATA
     ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
 824 GTTTAGAAACCCCAGAGACACCCGTGTAGATGGGTACATCACGGCTTCTCTCCCACGTTAAGACTTAATA

 281 AAG.....
     |||
 894 AAGGCAAG
```

1 **2** 3

```
                                              M  A  V  T  F  E  D  V  T  I  I  F  T  W  E
   284 .....AAAGCCTGATCTTTGACCAGATGGCAGTAACATTTGAAGATGTGACTATTATTTTTACATGGGAG
             ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
 18109 CCTAGAAAGCCTGATCTTTGACCAGATGGCAGTAACATTTGAAGATGTGACTATTATTTTTACATGGGAG

         E  W  K  F  L  D  S  S  Q  K  R  L  Y  R  E  V  M  W  E  N  Y  T  N
   349 GAGTGGAAATTCCTGGATTCTTCTCAAAAAAGACTCTACAGGGAGGTCATGTGGGAGAACTACACAAATG
       ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
 18179 GAGTGGAAATTCCTGGATTCTTCTCAAAAAAGACTCTACAGGGAGGTCATGTGGGAGAACTACACAAATG

         V  M  S  V
   419 TCATGTCAGTAG.....
       |||||||||||||
 18249 TCATGTCAGTAGGTAAA
```

```
              E   N   W   N   E   S   Y   K   S   Q   E   E   K   F   R   Y   L   E   Y   E   N   F
   431 . . . . . A A A A C T G G A A T G A G A G C T A C A A A T C C C A A G A A G A A A A A T T C A G A T A C T T A G A A T A T G A A A A T T T T
             | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 19409 T C T A G A A A A C T G G A A T G A G A G C T A C A A A T C C C A A G A A G A A A A A T T C A G A T A C T T A G A A T A T G A A A A T T T T

              S   Y   W   Q   G   W   W   N   A   G   A   Q   M   Y   E   N   Q   N   Y   G   E   T   V
   496 T C C T A C T G G C A A G G C T G G T G G A A T G C T G G C G C C C A G A T G T A T G A G A A T C A G A A C T A T G G G G A A A C T G T T C
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 19479 T C C T A C T G G C A A G G C T G G T G G A A T G C T G G C G C C C A G A T G T A T G A G A A T C A G A A C T A T G G G G A A A C T G T T C

              Q   G   T   D   S   K   D   L   T   Q   Q   D   R   S   Q   C   Q   E   W   L   I   L   S   T
   566 A A G G G A C A G A T T C C A A A G A C C T C A C A C A G C A A G A T C G T T C C C A G T G T C A G G A A T G G T T A A T A C T C T C C A C
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 19549 A A G G G A C A G A T T C C A A A G A C C T C A C A C A G C A A G A T C G T T C C C A G T G T C A G G A A T G G T T A A T A C T C T C C A C

              Q   V   P   G   Y   G   N   Y   E   L   T   F   E   S   K   S   L   R   N   L   K   Y   K
   636 A C A A G T A C C A G G G T A T G G G A A C T A T G A A C T G A C T T T T G A A A G C A A A A G T C T C A G G A A C T T A A A A T A T A A A
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 19619 A C A A G T A C C A G G G T A T G G G A A C T A T G A A C T G A C T T T T G A A A G C A A A A G T C T C A G G A A C T T A A A A T A T A A A

              N   F   M   P   W   Q   S   L   E   T   K   T   T   Q   D   Y   G   R   E   I   Y   M   S
   706 A A T T T T A T G C C T T G G C A G T C C T T A G A A A C A A A A A C C A C T C A A G A C T A T G G T A G A G A A A T C T A C A T G A G T G
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 19689 A A T T T T A T G C C T T G G C A G T C C T T A G A A A C A A A A A C C A C T C A A G A C T A T G G T A G A G A A A T C T A C A T G A G T G

              G   S   H   G   F   Q   G   G   R   Y   R   L   G   I   S   R   K   N   L   S   M   E   K   E
   776 G T T C A C A T G G T T T T T C A A G G G G G C A G A T A C C G T C T T G G C A T A T C C A G G A A A A A C C T C T C C A T G G A A A A A G A
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 19759 G T T C A C A T G G T T T T T C A A G G G G G C A G A T A C C G T C T T G G C A T A T C C A G G A A A A A C C T C T C C A T G G A A A A A G A

              Q   K   L   I   V   Q   H   S   Y   I   P   V   E   E   A   L   P   Q   Y   V   G   V   I
   846 A C A G A A G C T C A T A G T T C A G C A T T C T T A T A T C C C A G T G G A G G A A G C C C T T C C A C A G T A T G T T G G G G T G A T A
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 19829 A C A G A A G C T C A T A G T T C A G C A T T C T T A T A T C C C A G T G G A G G A A G C C C T T C C A C A G T A T G T T G G G G T G A T A

              C   Q   E   D   L   L   R   D   S   M   E   E   K   Y   C   G   C   N   K   C   K   G   I
   916 T G T C A A G A A G A C C T A C T G A G A G A T T C A A T G G A A G A A A A G T A C T G T G G A T G T A A T A A A T G T A A A G G A A T T T
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 19899 T G T C A A G A A G A C C T A C T G A G A G A T T C A A T G G A A G A A A A G T A C T G T G G A T G T A A T A A A T G T A A A G G A A T T T

              Y   Y   W   N   S   R   C   V   F   H   K   R   N   Q   P   G   E   N   L   C   Q   C   S   I
   986 A T T A T T G G A A C T C A C G G T G T G T T T T C C A C A A G A G A A A T C A A C C T G G A G A A A A C C T C T G T C A A T G C T C C A T
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 19969 A T T A T T G G A A C T C A C G G T G T G T T T T C C A C A A G A G A A A T C A A C C T G G A G A A A A C C T C T G T C A A T G C T C C A T

              C   K   A   C   F   S   Q   R   S   D   L   Y   R   H   P   R   N   H   I   G   K   K   L
  1056 C T G T A A A G C A T G C T T C T C T C A G A G A T C A G A C T T G T A T A G A C A T C C A A G A A A C C A C A T A G G T A A G A A G C T G
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 20039 C T G T A A A G C A T G C T T C T C T C A G A G A T C A G A C T T G T A T A G A C A T C C A A G A A A C C A C A T A G G T A A G A A G C T G

              Y   G   C   D   E   V   D   G   N   F   H   Q   S   S   G   V   H   F   H   Q   R   V   H
  1126 T A C G G A T G T G A T G A A G T T G A C G G T A A C T T T C A T C A G A G C T C C G G A G T T C A C T T T C A T C A G A G A G T T C A C A
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 20109 T A C G G A T G T G A T G A A G T T G A C G G T A A C T T T C A T C A G A G C T C C G G A G T T C A C T T T C A T C A G A G A G T T C A C A

              I   G   E   V   P   Y   S   C   N   A   C   G   K   S   F   S   Q   I   S   S   L   H   N   H
  1196 T A G G G G A G G T A C C T T A T A G C T G T A A T G C A T G T G G T A A G A G C T T C A G C C A G A T C T C T A G T C T T C A C A A T C A
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 20179 T A G G G G A G G T A C C T T A T A G C T G T A A T G C A T G T G G T A A G A G C T T C A G C C A G A T C T C T A G T C T T C A C A A T C A

              Q   R   V   H   T   E   E   K   F   Y   K   I   E   C   D   K   D   L   S   R   N   S   L
  1266 T C A A A G A G T C C A C A C A G A A G A G A A A T T C T A T A A A A T T G A G T G T G A T A A A G A C C T C A G T A G A A A T T C A T T A
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 20249 T C A A A G A G T C C A C A C A G A A G A G A A A T T C T A T A A A A T T G A G T G T G A T A A A G A C C T C A G T A G A A A T T C A T T A

              L   H   I   H   Q   R   L   H   I   G   E   K   P   F   K   C   N   Q   C   G   K   S   F
  1336 C T T C A C A T T C A C C A G A G A C T T C A C A T A G G A G A G A A G C C T T T T A A A T G T A A T C A G T G T G G T A A G A G T T T T A
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 20319 C T T C A C A T T C A C C A G A G A C T T C A C A T A G G A G A G A A G C C T T T T A A A T G T A A T C A G T G T G G T A A G A G T T T T A

              N   R   S   S   V   L   H   V   H   Q   R   V   H   T   G   E   K   P   Y   K   C   D   E   C
  1406 A T C G G A G T T C A G T A C T T C A T G T T C A T C A G A G A G T C C A C A C A G G A G A A A A A C C A T A T A A G T G T G A T G A G T G
         | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
 20389 A T C G G A G T T C A G T A C T T C A T G T T C A T C A G A G A G T C C A C A C A G G A G A A A A A C C A T A T A A G T G T G A T G A G T G
```

```
              G   K   G   F   S   Q   S   S   N   L   R   I   H   Q   L   V   H   T   G   E   K   S   Y
 1476  T G G T A A G G G T T T C A G C C A G A G C T C A A A T C T T C G A A T T C A T C A G T T A G T A C A C A C A G G A G A G A A G T C T T A T
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |
20459  T G G T A A G G G T T T C A G C C A G A G C T C A A A T C T T C G A A T T C A T C A G T T A G T A C A C A C A G G A G A G A A G T C T T A T

              K   C   E   D   C   G   K   G   F   T   Q   R   S   N   L   Q   I   H   Q   R   V   H   T
 1546  A A A T G T G A A G A C T G T G G T A A A G G C T T T A C C C A G C G C T C A A A T C T T C A A A T T C A T C A G A G A G T G C A T A C A G
       | | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
20529  A A A T G T G A A G A C T G T G G T A A A G G C T T T A C C C A G C G C T C A A A T C T T C A A A T T C A T C A G A G A G T G C A T A C A G

          G   E   K   P   Y   K   C   D   D   C   G   K   D   F   S   H   S   S   D   L   R   I   H   Q
 1616  G A G A G A A A C C T T A T A A A T G T G A T G A C T G T G G A A A G G A C T T T A G T C A C A G C T C A G A T C T T C G C A T T C A T C A
       | | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
20599  G A G A G A A A C C T T A T A A A T G T G A T G A C T G T G G A A A G G A C T T T A G T C A C A G C T C A G A T C T T C G C A T T C A T C A

              R   V   H   T   G   E   K   P   Y   T   C   P   E   C   G   K   G   F   S   K   S   S   K
 1686  G A G A G T C C A T A C A G G G G A G A A A C C C T A T A C T T G T C C T G A A T G T G G G A A G G G C T T C A G T A A G A G T T C A A A G
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
20669  G A G A G T C C A T A C A G G G G A G A A A C C C T A T A C T T G T C C T G A A T G T G G G A A G G G C T T C A G T A A G A G T T C A A A G

              L   H   T   H   Q   R   V   H   T   G   E   K   P   Y   K   C   E   E   C   G   K   G   F
 1756  C T T C A C A C T C A T C A A A G A G T A C A T A C T G G A G A G A A A C C C T A C A A A T G T G A A G A G T G T G G C A A G G G A T T C A
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
20739  C T T C A C A C T C A T C A A A G A G T A C A T A C T G G A G A G A A A C C C T A C A A A T G T G A A G A G T G T G G C A A G G G A T T C A

          S   Q   R   S   H   L   L   I   H   Q   R   V   H   T   G   E   K   P   Y   K   C   H   D   C
 1826  G T C A G C G T T C A C A T C T T C T C A T T C A T C A G A G A G T C C A T A C A G G A G A A A A G C C C T A T A A A T G T C A T G A T T G
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
20809  G T C A G C G T T C A C A T C T T C T C A T T C A T C A G A G A G T C C A T A C A G G A G A A A A G C C C T A T A A A T G T C A T G A T T G

              G   K   G   F   S   H   S   S   N   L   H   I   H   Q   R   V   H   T   G   E   K   P   Y
 1896  T G G A A A G G G T T T T A G T C A C A G T T C T A A T C T T C A C A T T C A T C A G A G G G T C C A T A C A G G A G A G A A G C C T T A T
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
20879  T G G A A A G G G T T T T A G T C A C A G T T C T A A T C T T C A C A T T C A T C A G A G G G T C C A T A C A G G A G A G A A G C C T T A T

          Q   C   A   K   C   G   K   G   F   S   H   S   S   A   L   R   I   H   Q   R   V   H   A
 1966  C A A T G T G C T A A G T G T G G T A A A G G T T T C A G T C A T A G C T C A G C T C T T C G A A T T C A T C A A A G A G T C C A T G C A G
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
20949  C A A T G T G C T A A G T G T G G T A A A G G T T T C A G T C A T A G C T C A G C T C T T C G A A T T C A T C A A A G A G T C C A T G C A G

          G   E   K   P   Y   K   C   R   E   Y   Y   K   G   F   D   H   N   S   H   L   H   N   N   H
 2036  G A G A G A A A C C T T A C A A A T G C C G T G A A T A T T A T A A G G G A T T T G A T C A T A A T T C A C A T C T T C A C A A T A A T C A
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
21019  G A G A G A A A C C T T A C A A A T G C C G T G A A T A T T A T A A G G G A T T T G A T C A T A A T T C A C A T C T T C A C A A T A A T C A

          R   R   G   N   L   *
 2106  T A G A A G A G G A A A C T T A T A A A T A T T G T T C A T T T A G T T A A C A G C T T T A A T C A A A G T T T A C C T A A C C T T T A A A
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
21089  T A G A A G A G G A A A C T T A T A A A T A T T G T T C A T T T A G T T A A C A G C T T T A A T C A A A G T T T A C C T A A C C T T T A A A

 2176  C C C T A T A A A T C C T G C T G T T A A G G A A A T C T T A T A A A T A A C A C A A G T A A T C C C A A G C A A C A T T T A T A G T T T C
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
21159  C C C T A T A A A T C C T G C T G T T A A G G A A A T C T T A T A A A T A A C A C A A G T A A T C C C A A G C A A C A T T T A T A G T T T C

 2246  C C C T A T C T C C C A C T A A G A A T T A T T T G C T T C A A A A G G A G A T C T T T A G A A A A A A C C C T A T A T A T T T A A A A T T
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
21229  C C C T A T C T C C C A C T A A G A A T T A T T T G C T T C A A A A G G A G A T C T T T A G A A A A A A C C C T A T A T A T T T A A A A T T

 2316  A T A G T G T A T T T T T T C T T T A C C T A C T A T A A A T A T A A T A C A G T C A T A A A T A T A T T A A A C A T T T A A G G A G A A A A
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
21299  A T A G T G T A T T T T T T C T T T A C C T A C T A T A A A T A T A A T A C A G T C A T A A A T A T A T T A A A C A T T T A A G G A G A A A A

 2386  C T C T T C A T T C T A T T T C A T T C T A G T C T T T T T T T C T G T G C A T T T T A A T G T G C A T G A A A T T G T G T G T T C A A T T
       | | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |   | | | | | |   | | | | |   | | | | |
21369  C T C T T C A T T C T A T T T C A T T C T A G T C T T T T T T T C T G T G C A T T T T A A T G T G C A T G A A A T T G T G T G T T C A A T T
```

```
 2456 TTGTATTTTCACTGTCTTCAAAATATTTACTTAAATTTTGGTTTGAATTGAAACTGGTTGGCCAACTGTT
      |||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
21439 TTGTATTTTCACTGTCTTCAAAATATTTACTTAAATTTTGGTTTGAATTGAAACTGGTTGGCCAACTGTT

 2526 AAACGACATCTCTTAACTCCCCTAAAAACTCCCTAGGAAGTAACAGAAAGATGGAAACACATAGAACTT
      |||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
21509 AAACGACATCTCTTAACTCCCCTAAAAACTCCCTAGGAAGTAACAGAAAGATGGAAACACATAGAACTT

 2596 AAAACTCAGTTTTGGCCGGTAGAATTCAATTGTTTATGGACAAAAGCCACCTAATAAAAGATAGGAAAGC
      |||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
21579 AAAACTCAGTTTTGGCCGGTAGAATTCAATTGTTTATGGACAAAAGCCACCTAATAAAAGATAGGAAAGC

 2666 ATT
      |||
21649 ATT
```

b.  0.625 pts: Report CDS locations: **614..891, 18109..18265, 19409..21651**

2. 1.25 pts. Create a BED6 file with 2 lines based on the attached paper (Takenaka_et_al-2015-FEBS_Journal.pdf). Figure 3 shows the location of transcription factor DdlR binding to the promoter region of the ddlR-ddl operon in Brevibacillus brevis. The chromosomal location of the ddlR CDS is 2968133..2969623. The zero-based BED6 file should contain the location information of two genomic regions:

a.  0.625 pts: The region bound by the DdlR transcription factor, which we will call the promoter. It is 170 bp in length, begins 140 nucleotides upstream from the start codon, and ends 29 nucleotides downstream from the start codon.

$$2968133-140 = 2967993(\text{start}) \qquad 2968133 + 29 = 2968162(\text{end}) \text{ *need 0base*}$$

$$2968162 - 2967993 = 169(\text{…reading graded homework collaboration in slack to attempt to determine why missing one})$$

chr1 2967992 2968162 promoter 0 +

b.  0.625 pts: The 5' UTR, noting that the transcription start site, as predicted by BPROM, begins 38 nucleotides upstream from the start codon. The 5' UTR is defined as the region from the transcription start site through the nucleotide that immediately precedes the start codon.

$$2968133 - 38 = 2968095(\text{start}) \qquad 2968132(\text{end})$$

$$2968132-2968095= 37+1 = 38\text{bp length of 5' UTR}$$

chr1 2968095 2968132 5UTR 0 +

takenaka bed6 file.txt - Notepad

File  Edit  Format  View  Help

```
chr1 2967992 2968162 promoter 0 +
chr1 2968095 2968132 5UTR 0 +
```

takenaka bed6 file.bed                                    HDD 1TB  (F:)\Individualized Genomics and ...\Bioinformatics Tools for Genome Analysis

saved these two lines in this .bed file for question 3

3. 1.25 pts. Submit a screenshot of the BED6 from Problem 2. Using the NCBI Genome Browser for Brevibacillus brevis NBRC 100599, load your BED6 file. Take a screenshot showing the entire promoter, 5' UTR region, and CDS of ddlR. Be sure to zoom in so that these regions take up a majority of the shot.

***I forgot to set the spaces between each individual item in the .bed file as a 'tab', gave me quite the headache until I figured this out!***



https://www.ncbi.nlm.nih.gov/nuccore/226092535?report=graph&tracks=[key:sequence_track,name:Sequence,display_name:Sequence,id:STD649220238,annots:Sequence,ShowLabel:false,ColorGaps:false,shown:true,order:1][key:gene_model_track,name:Genes,display_name:Genes,id:STD3194982005,annots:Unnamed,Options:MergeAll,CDSProductFeats:false,NtRuler:true,AaRuler:true,HighlightMode:2,ShowLabel:true,shown:true,order:5][key:feature_track,name:U1CtIwAlPC90B,display_name:(U) takenaka bed6 file.bed,id:U1CtIwAlPC90B,data_key:LZ63RzGc7kVCsqBCYaOWvMYJ7eGykLyVsJOYhYyBnq8Po1Mp6ROZzlp6b0ZN1OOpsrHvpfGBqoTtnvmT_5XzjsGn_KnQlfo,subkey:region,dbname:NetCache,annots:takenaka bed6 file.bed_UUD1614012630DUU_region,Layout:Adaptive,LinkedFeat:Packed,shown:true,order:6]&assm_context=GCA_000010165.1&v=2967880:2968251&c=FF0000&select=gi|226092535-002d49b8-002d4a61-0153-df2ec903-2e039e4c;&slim=0



Brevibacillus brevis NBRC 100599 DNA, complete genome

4. 1.25 pts. Use the web-based Biomart in Ensembl to create a dataset and save it as a TSV, CSV, or XLS file. Use the following parameters to make the dataset:

Dataset:
Ensembl Genes 100 (or the latest version)
Mouse genes (GRCm38.p6) (or the latest version)
Filters:
Chromosome 11
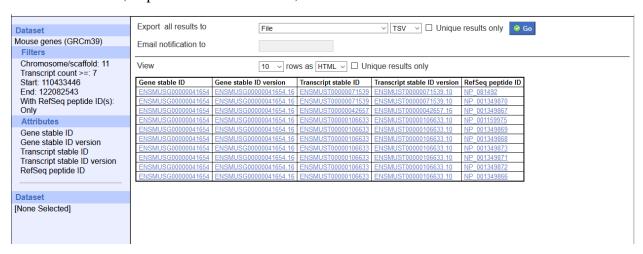Band E2 only (found via ensembl)
Transcript count >=7
Limit to genes with RefSeq protein (peptide) IDs only
Attributes:
Default attributes
Add "RefSeq Protein (peptide) ID"

Get all the results, export the results to a file, and submit the file.



**File will be attached to the submission of this assignment!**

(Hint) Band E2 is not available directly from Ensembl anymore. But, you may search for other databases, and filter via coordinates. Or you can use Ensembl archive with the specific archive version number.

(Optional) If you are doing this for human genes, you can still filter by karyoband. In this case, try with p12 band.