

Suggested strategy:

- First, split all the alignments in the SAM file based on chromosomes. Create a dictionary where the key is the chromosome name and the value is a list of start coordinates. Note that some lines do not have a value for chromosomes which means that sequence did not align to the genome, so we can ignore them.
- Next go through each exon in the GTF file and count the number of reads that are between the start and end coordinates. Since we know which chromosome the gene is in you should have to go through the start positions one at a time to see how many are within the start and end coordinate.

- 1) **Implement without sorting the sam file:** Use the original sam file to determine the number of reads that match a gene, count the number of start positions from the SAM file that fall between the start and stop coordinates of an exon that belongs to a gene.
- 2) **Implement after sorting the sam file:** For each chromosome sort the alignment coordinates. You can use any of the functions discussed in class. Now again count the number of alignments that match a gene. *****Note***** Since the alignment coordinate is sorted, you don't have to check every single value in the list. Once your alignment coordinate is greater than the end coordinate of the exon, you can stop. Additionally, you can remove the alignment coordinates from the list because the genes are sorted so you will never have to look before the position you have already visited.

In both the above cases calculate the total time taken with sorting and without sorting the coordinates from the sam file and clearly explain why sorting reduces/increases the total time for execution. You can use the below-mentioned code to calculate the total time taken to execute a function.

```
# in Python
import time
start = time.process_time()
# your code here
print(time.process_time() - start)
```