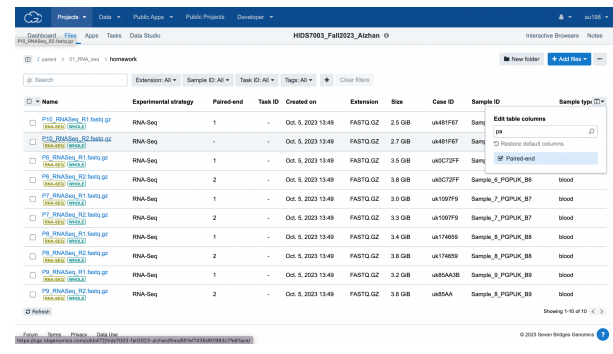
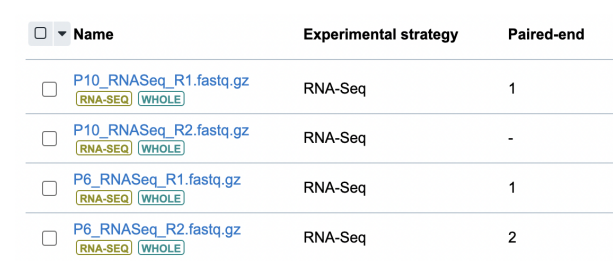
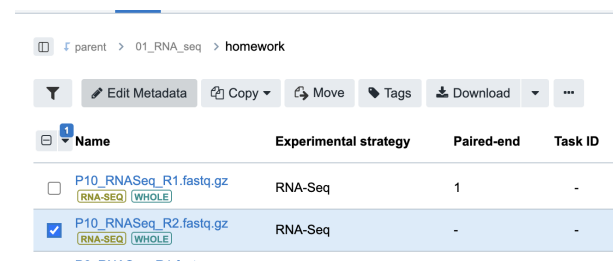

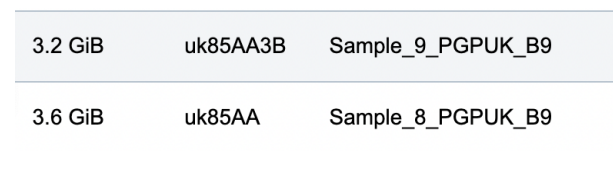
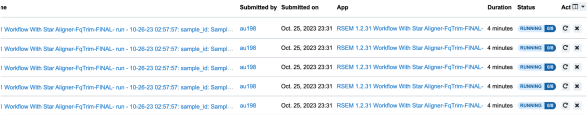

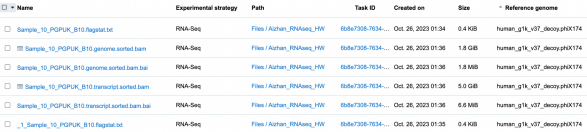










Homework 9: Data Processing of RNA-seq data HIDS-7003-01






I	<p>Edited the metadata</p> <p>Added the columns for:</p> <ul style="list-style-type: none"> - Pair-end - Case ID - Sample ID 	
II	<p>First error appeared to be in the pair-end column</p> <p>The P10_RNASeq_R2 does not have the label that its a reversed end</p>	
III	<p>Selecting the file to edit the metadata</p>	
IV	<p>Changed the paired-end to 2, as it's a reversed read</p>	
V	<p>The second error appears to be in the Case ID and Sample ID</p> <p>With the reversed sample's case ID - uk85AA and sample ID Sample_9</p>	

VI	Similarly to the above steps, I have changed the Case ID to math with the forward sample	<div> <div>Case</div> <div> <div>Case ID ?</div> <div>uk85AA3B</div> </div> <div> <div>Case UUID ?</div> <div></div> </div> </div>
VII	Similarly, I have changed the Sample ID to Sample_9	<div> <div>Sample</div> <div> <div>Sample ID ?</div> <div>Sample_9_PGPUK_B9</div> </div> <div> <div>Sample UUID ?</div> <div></div> </div> <div> <div>Sample type ?</div> <div>blood</div> </div> </div>
VIII	1. Selected all the input files 2. Set the batch ON	<div> <div>Input Read Files ?</div> <div>Change selection</div> <div> Batch by: File metadata </div> <div> This task will be batched by file metadata (Sample ID) and this will create 5 groups. </div> <div> <div>Sample_10_pgpu_k_b10 (2 items) x</div> <div>Sample_6_pgpu_k_b6 (2 items) x</div> <div>Sample_7_pgpu_k_b7 (2 items) x</div> <div>Sample_8_pgpu_k_b8 (2 items) x</div> <div>Sample_9_pgpu_k_b9 (2 items) x</div> </div> </div>
IX	3. Set the annotation 4. Referenced to the genome	<div> <div>Annotation GTF ?</div> <div>Change selection</div> <div> Batch by: None </div> <div> Homo_sapiens.GRCh37.75.gtf </div> <div> <div>Reference FASTA File or Index TAR Bundle * ?</div> <div>Change selection</div> <div> Batch by: None </div> <div> human_g1k_v37_decoy.phiX174_Homo_sapiens.GRCh37.... </div> </div> </div>

X	The output settings are set for us to obtain 12 categories of output/ or tasks	<div>Output Settings</div> <div><div>Genes results</div><div>No value</div></div> <div><div>Isoforms results</div><div>No value</div></div> <div><div>RSEM Plot Model PDF File ?</div><div>No value</div></div> <div><div>Transcript BAM ?</div><div>No value</div></div> <div><div>flagstat_metrics ?</div><div>No value</div></div> <div><div>flagstat_metrics_1 ?</div><div>No value</div></div> <div><div>report ?</div><div>No value</div></div> <div><div>report_zip ?</div><div>No value</div></div> <div><div>report_zip_1 ?</div><div>No value</div></div> <div><div>sample_name_genome_bam ?</div><div>No value</div></div> <div><div>star_log_files ?</div><div>No value</div></div> <div><div>trimmed_reads ?</div><div>No value</div></div>
XI	Screenshot of running the batch	
XII	Created a new folder for the output	
XIII	There are 110 files within the whole output showing in the output folder, with each patient having 22 (including the genome-referenced files)	<div>Filter by: Sample ID</div> <div>✕ Clear selected</div> <div>Showing 5</div> <div><div><input type="checkbox"/> Sample_10_PGPUK_B10 (22)</div><div><input type="checkbox"/> Sample_6_PGPUK_B6 (22)</div><div><input type="checkbox"/> Sample_7_PGPUK_B7 (22)</div><div><input type="checkbox"/> Sample_8_PGPUK_B8 (22)</div><div><input type="checkbox"/> Sample_9_PGPUK_B9 (22)</div></div>
XIV	For speeding up and execution purposes, the following parameters could have been turned off: <ul style="list-style-type: none">"Output genome BAM""Sort BAM by coordinate""Calculate credibility intervals"	

<p>XV</p>	<p>Looking specifically at the patient task of the sample</p> <p>There are 16 files generated as an output</p> <p>However, as the final output we would only use five output files:</p> <ol style="list-style-type: none"> 1) <i>Genes results</i> 2) <i>Isoforms results</i> 3) <i>RSEM Plot Model PDF file</i> 4) <i>Transcript BAM</i> 5) <i>Genome BAM</i> 	
<p>To summarize the output I will be looking at all files of the sample 10</p>		
<p>1.</p>	<p>Gene expression results</p> <p>This file contains data and information about the genes detected in the RNA sequencing experiment. It typically includes gene expression levels and related statistics.</p>	<p>▼ Genes results </p> <p>Sample_10_PGPUK_B10.genes.results</p>
<p>2.</p>	<p>Isoforms results</p> <p>This file stores results related to RNA isoforms. It contains data and information about the various alternative splicing variants of genes detected in the RNA sequencing analysis.</p>	<p>▼ Isoforms results </p> <p>Sample_10_PGPUK_B10.isoforms.results</p>

3.	<p>RSEM Plot (PDF)</p> <p>This is the output file that contains a PDF representation of the plot model generated by the RNA-Seq by Expectation-Maximization (RSEM) analysis. It visualizes gene expression data.</p>	<p>▼ RSEM Plot Model PDF File ? </p> <p> Sample_10_PGPUK_B10_plot_model.pdf</p>
4.	<p>Transcript BAM file after alignment (transcript.sorted.bam)</p> <p>This is a BAM file containing transcript-level data, specifically aligned RNA-Seq reads that correspond to the transcribed regions of the genome. It's a critical intermediate file in RNA-Seq analysis.</p>	<p>▼ Transcript BAM ? </p> <p>Sample_10_PGPUK_B10.transcript.sorted.bam</p>
5.	<p>Quality check of BAM files after alignment (flagstat.txt)</p> <p>This TXT file contains metrics for the genome BAM file, providing alignment statistics for the entire genome</p>	<p>▼ flagstat_metrics ? </p> <p>Sample_10_PGPUK_B10.flagstat.txt</p>
6.	<p>Quality check of BAM files after alignment (flagstat.txt)</p> <p>This TXT file contains metrics generated by the 'flagstat' command for the transcript BAM file. It provides statistics on the alignment quality and characteristics of the transcript-level data.</p>	<p>▼ flagstat_metrics_1 ? </p> <p>_1_Sample_10_PGPUK_B10.flagstat.txt</p>
7.	<p>Report (TXT)</p> <p>The report summarizing the results and findings of the RNA sequencing analysis, which may include quality control and statistics was NOT GENERATED</p>	<p>report ? No value</p>

8.	<p>Quality check of input FASTQ files (<i>fastqc.zip</i>) for forward (R1) and reversed (R2) ends</p> <p>FastQC is used for quality control of RNA-Seq data, helping to assess the quality and characteristics of the sequencing reads.</p>	<p>▼ report_zip ? </p> <p>P10_RNASeq_R1_fastqc.zip</p> <p>P10_RNASeq_R2_fastqc.zip</p>
9.	<p>Quality check after trimming (<i>trimmed_fastqc.zip</i>)</p> <p>Similar to the previous item, this ZIP file array contains FastQC reports specifically generated after trimming the sequencing reads, helping to assess the quality of the post-processed data</p>	<p>▼ report_zip_1 ? </p> <p>P10_RNASeq_R2.trimmed_fastqc.zip</p> <p>P10_RNASeq_R1.trimmed_fastqc.zip</p>
10.	<p>Genome BAM file after alignment (<i>genome.sorted.bam</i>)</p> <p>This BAM file contains RNA-Seq reads mapped to the entire genome. It represents the alignment of reads to the genomic regions.</p>	<p>▼ sample_name_genome_bam ? </p> <p>Sample_10_PGPUK_B10.genome.sorted.bam</p>
11.	<p>Log files (x3)</p> <p>These are log files generated during the alignment process using the STAR aligner. They provide information about the alignment and mapping of sequencing reads to the reference genome for RNA-Seq analysis.</p>	<p>▼ star_log_files ? </p> <p>Sample_10_PGPUK_B10Log.progress.out</p> <p>Sample_10_PGPUK_B10Log.out</p> <p>Sample_10_PGPUK_B10Log.final.out</p>
12.	<p>Trimmed FASTQ files (<i>trimmed.fastq</i>) for forward (R1) and reversed (R2) ends</p> <p>This is an array of file types, such as FASTQ and related formats, containing RNA-Seq reads after quality trimming and adapter removal. It represents the cleaned and processed sequencing data.</p>	<p>▼ trimmed_reads ? </p> <p>P10_RNASeq_R1.trimmed.fastq</p> <p>P10_RNASeq_R2.trimmed.fastq</p>

