

Module 1.1: Learn - Biology

Overview

By the end of this module, you should be able to explain the general biology of the placenta, why we are interested in placental genes that have higher expression in one sex versus the other, the difficulties involved in studying genes on the sex chromosomes which are likely differentially expressed between the sexes, and have a thorough knowledge the data we will be working with. By the end of this course, you will have conducted differential gene expression analysis of real human placenta samples.

If you have questions about any of the content, please post in the course slack under the channel for module 1.

In this section we will work on:

1. Course description and research question
2. Placenta biology
3. Sex differences and sex chromosomes
4. Genomics/transcriptomics refresher
5. Gene expression and RNAseq preprocessing

Research Aims

What is this course?

This is a course-based research experience specifically designed to give ASUOnline students the opportunity to attain research experience, gain technical skills, and give back to the scientific community. The Genomics Research Experience is a 7-week research experience that answers *real* research questions with *real* research techniques. Notice the emphasis on the words *real* and *research*. This means we do not know the outcome of our research question and will be working as a team to identify best practices, evaluate, and

troubleshoot. The results we produce will be used in a real publication, so please be thorough in your understanding. Further, because this is an authentic research experience – meaning we have not conducted these analyses before – things could go wrong. In fact, as in all authentic research experiences, we will get both expected and unexpected results. You will analyze those results together.

Each module contains three sections that are aimed to help you learn (1) the **biology** behind our research, (2) the skills needed to perform our research, and (3) the **professional development** skills used in real genomics research. Let's take a look at the question we are seeking to answer throughout the duration of this course.

The research question

In this research experience, we are asking a very specific question:

What are the effects of trimming next-generation sequencing reads on the measurement of sex differences in gene expression in the human placenta?

For any computational analysis, researchers usually choose data processing parameters based on what they find in other groups' publications or what was used successfully in their own past studies. What was chosen for those data processing steps may (or may not) have big effects on findings and conclusions drawn from that analysis. In this project we are exploring one common choice of parameters for trimming, a very common early data processing step for next-generation sequencing analysis. This choice will be widely relevant since many groups use trimming as the first step in their sequencing data analysis workflow. You will learn about the basic steps of data processing and differential gene expression analysis in this course.

The Wilson lab at ASU is interested in studying the biology of the placenta, a crucial organ in human development. They analyzed human placenta samples and determined a set of genes that had significantly higher expression in one of the sexes.

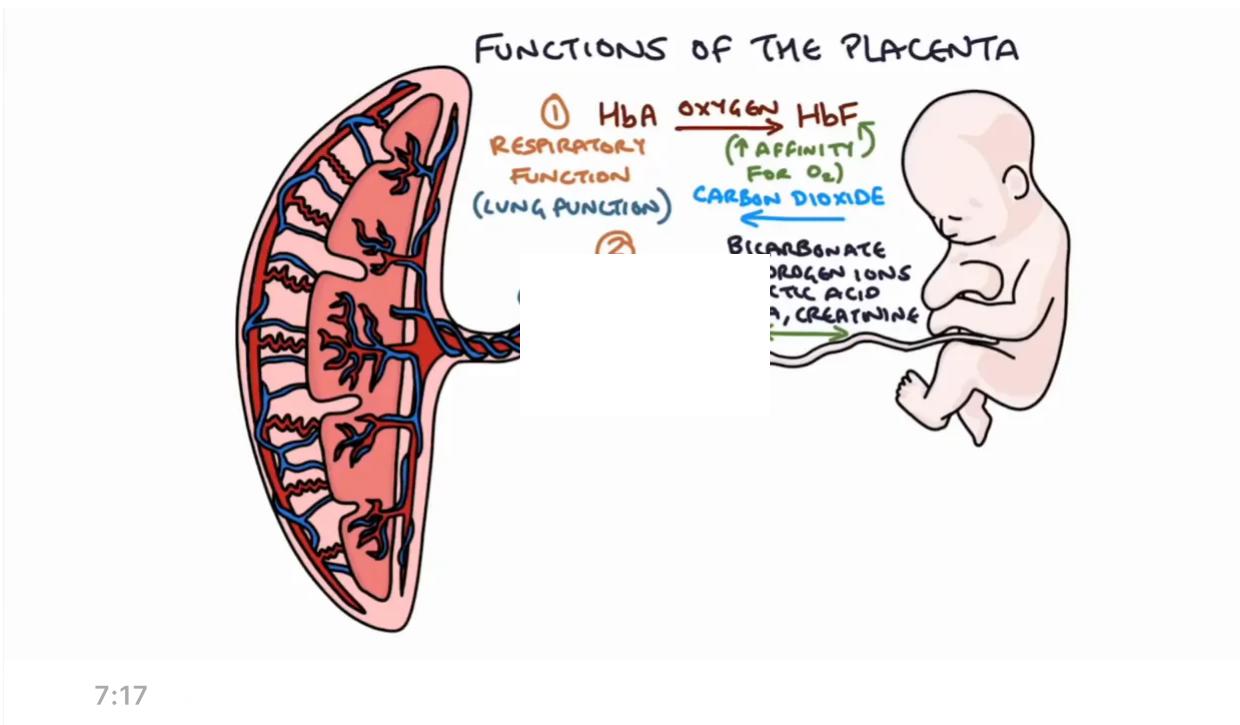
males) using specific data processing parameters, but we wanted to know how robust our results are. Since the differentially expressed between the sexes in the human placenta will serve as the test for effects of varying trimr Having worked with this data and these samples in the lab, we have already established that this data set is biolc high quality, and does in fact show a signal for gene expression differences; this establishes a baseline for this re

Placenta biology and sex differences

Placenta Biology

To give you some background on why we chose to use genes differentially expressed by sex in the placenta as a effects of trimming, we will start by describing the function of the placenta and talk about what is known about gei tissue.

The placenta is a highly specialized organ that develops in the uterus during pregnancy serving as an interface b and the growing fetus. The [placenta ↗ \(https://www.mayoclinic.org/healthy-lifestyle/pregnancy-week-by-week/in-depth-placenta/art-20044425#:~:text=The%20placenta%20is%20an%20organ,umbilical%20cord%20arises%20from%20it.\)](https://www.mayoclinic.org/healthy-lifestyle/pregnancy-week-by-week/in-depth-placenta/art-20044425#:~:text=The%20placenta%20is%20an%20organ,umbilical%20cord%20arises%20from%20it.) nourishes, pr removes waste from a developing fetus. The video below gives an overview of the structures of the placenta and



7:17

Video. Understanding the Placenta.

This video succinctly reviews the placental structures and the five key functions of the placenta.

[View Transcript. \(https://canvas.asu.edu/courses/122165/files/54792578?wrap=1\)](https://canvas.asu.edu/courses/122165/files/54792578?wrap=1) [\(https://canvas.asu.edu/courses/122165/files/54792578/download?download_frd=1\)](https://canvas.asu.edu/courses/122165/files/54792578/download?download_frd=1)

The placenta has two sides – one is connected to the fetus by the umbilical cord and the other is connected to the maternal tissue called “decidua”. Contained within the umbilical cord are two arteries that carry deoxygenated blood to the fetus and a single umbilical vein that carries oxygenated blood to the fetus.

Within the placenta itself, there is a network of umbilical and maternal veins and arteries. The umbilical arteries and veins form tree-like structures surrounded by the placental membrane; in between is the intervillous space where exchange between maternal and fetal blood occurs. Maternal blood forms pools outside the placental membrane and moves

between the mother and fetus and across the membrane without mixing the blood. This process of diffusion allows for the exchange of substances across the placenta.

Figure. The structures of the placenta.

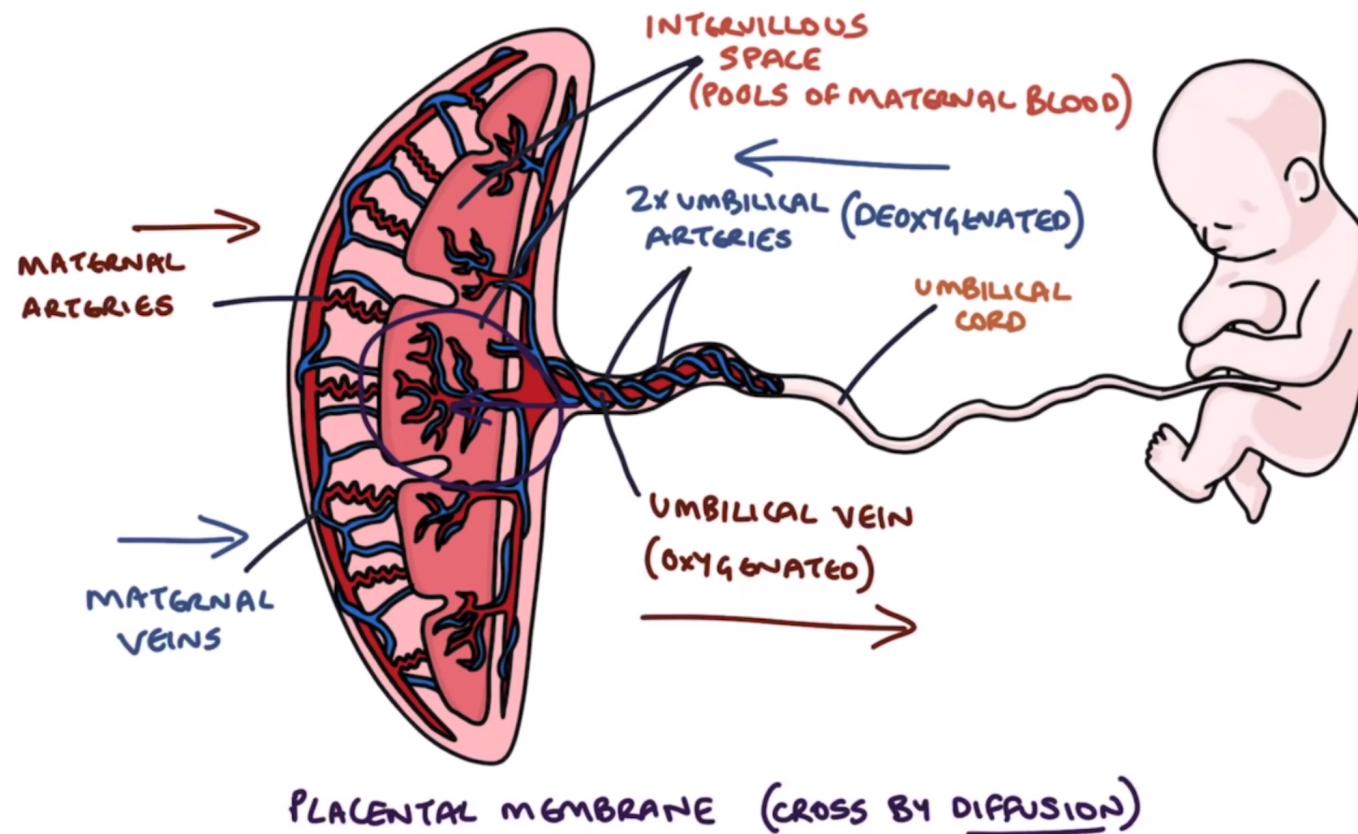


Figure. The structures of the placenta.

The main structures of the placenta include the arteries and veins on the maternal side, the intervillous space as the site of exchange, the two umbilical arteries for deoxygenated fetal blood, and the umbilical vein for oxygenated blood to get sent to the fetus.

The five key functions of the placenta

1. **Respiratory function:** allows the maternal blood to exchange oxygen with the fetal blood and for CO₂ waste the fetal blood
2. **Excretion function:** rid fetus of waste and provide osmoregulation
3. **Nutrition function:** carbohydrates, vitamins, and nutrients from maternal blood diffuse across the placental n blood stream
4. **Immunity function:** maternal antibodies from previous infections can cross the placental membrane to protect pregnancy and temporarily after birth
5. **Endocrine function:** the placental tissue creates hormones that help maintain the pregnancy, such as HCG (gonadotropin), estrogen, and progesterone

Sex differences in the placenta and pregnancy

In this research, we are comparing gene expression differences between placentas from individuals with two X chromosomes (XX) and placentas from individuals with an X and a Y chromosome (XY). Recall that the XX individual receives one X chromosome from each parent, one from the maternal genome and the other from their paternal genome, while the XY individual receives the X from the mother and the Y from the paternal genome. As the placenta and fetus are derived from the fertilized egg, they share the same genome. Therefore, an XX fetus will develop an XX placenta and an XY fetus develops an XY placenta.

In the past, the placenta had been viewed as an asexual organ and that resulted in biased studies that didn't consider sex differences in placental function. This is especially important in complicated pregnancies where the maternal condition combined with sex differences of the placenta impact the development of the fetus. For example, sex differences have been observed in various characteristics like birth weight (XY males tend to be larger) or pregnancy complications (XY males more often experience complications) and these outcomes can be intensified when the mother has a medical complication such as hypertension.

One study on pregnancies complicated by asthma found that if the mother went untreated, the XX fetus was smaller than the XY fetus. If the mother had the asthma under control, the XX fetus appeared to be unaffected; conversely, the XY fetus was unaffected by either condition. Another study observed that sex differences may even play a role in the first 48 hours of a premature birth, as XY fetuses have a higher survival rate as they are born more vasodilated than the premature XX individuals. Furthermore, sex differences are known to impact the immunity function of the human placenta and sex-specific fetal responses to complications may contribute to these adverse outcomes (Clifton, 2009).

In our data set, we aimed to have our sample set have as close to a 50:50 ratio of XX:XY as possible. We were able to analyze ten (10) XX and twelve (12) XY samples (22 total). Our research focuses on sex differences in placental health so that we can better understand how changes can lead to sex differences in poor outcomes for the fetus,

Sex chromosomes

To help you learn about and communicate results about sex differences, it is important to use the appropriate language when discussing them. We are interested in genetic differences due to genetic sex (sex chromosome genotype) – not gender, which are based on how a person self-identifies. The sex chromosome complement is the set of sex chromosomes contained in the nucleus. Although genetic males typically have one X chromosome and one Y chromosome, and females typically have two X chromosomes, the number of X and Y chromosomes can vary. Common sex chromosome variations include Klinefelter syndrome (47,XXY) and Turner syndrome (45,X). Further, sex chromosome content can vary over time as individuals can lose their Y chromosome within a proportion of their cells over time as they age; loss of Y chromosomes has been observed in cancer cells.

As part of the data processing for this research, we take one extra step - compared to typical genomics pipelines - to assess gene expression from sex chromosomes. This talk was featured in a meeting from Bioconductor which is a source software project that uses the R statistical programming language; you will be using Bioconductor packages throughout the remainder of this course. The video is included below as supplementary viewing.

Keynote: Sex biased Genomics and Methodology



Video. Bioconductor Conference 2021, Keynote Speaker: Melissa Wilson, PhD.

In this video, Dr. Melissa Wilson, Computational Evolutionary Biologist at ASU, presents *Sex-biased genome evolution* to discuss the structure, function, and evolution of the sex chromosomes.

[View Transcript. \(https://canvas.asu.edu/courses/122165/files/54912778?wrap=1\)](https://canvas.asu.edu/courses/122165/files/54912778?wrap=1) [Download \(https://canvas.asu.edu/courses/122165/files/54912778/download?download_frd=1\)](https://canvas.asu.edu/courses/122165/files/54912778/download?download_frd=1)

In this video, Dr. Wilson explains the structure, function, and evolution of the sex chromosomes. Dr. Wilson explains that sequence conservation between the X and Y chromosomes makes it difficult to align sequencing reads to the sex chromosomes. Take-home points that are important for this research experience are to know:

X and Y chromosomes share evolutionary history and were once homologous autosomes.

- The Y chromosome is much shorter than the X chromosome now, but still retains regions of shared ancestry.
- Because of this shared ancestry the modern X and Y chromosomes still share sequence similarity that can be used by alignment algorithms. Because of this, we need an approach to correctly represent the sex chromosome complement that matches the samples we are looking at.
- There are regions of the sex chromosomes that are 98 - 100% identical.

Sex chromosome complement is not always XX or XY, and may not match gonads or gender identity.

- Common sex chromosome variations include Klinefelter syndrome (47,XXY, which occurs in about 1/500 live individuals) and Turner syndrome (45,X), which occurs in about 1/2500 live assigned female at birth individuals.
- Swyer syndrome is uncommon (1/20,000 live births), and results when a person has XY sex chromosomes but only one testis (gonad).

It is possible for the sex chromosomes in a particular tissue to not match the overall karyotype of the person.

- For example, it is a common phenomenon that an individual who is born XY will lose their Y chromosome with other cells over time as they age.
- Further, it is common for cancer cells to lose an X chromosome (if starting as XX) or a Y chromosome (if starting as XY) as they rapidly divide.

The Wilson lab developed a technique called XYAlign which masks parts of the human reference genome to implement the sex of the sample. The XX samples in this study were aligned to the reference genome with the Y chromosome masked out. The XY samples were aligned to the reference genome with the pseudoautosomal regions which are identical. The Y chromosomes masked out, because you can not determine which sex chromosomes those reads came from.

Our Placenta Samples

In this video, Dr. Plaisier discusses how our samples of interest from the placenta were collected following deliveries of pregnancies (> 36 weeks). To minimize confounding factors that could skew results, the patient enrollment criteria included women with uncomplicated pregnancies, no delivery complications (e.g. premature delivery), and no medical conditions such as hypertension (high blood pressure). The tissue samples were collected from the fetal side of the placenta (indicated by the umbilical cord) and snap-frozen and stored by Yale Biobank. The Wilson Lab purchased these samples from the Yale Biobank and sent them to the Yale Genome Sequencing Center to extract DNA and RNA and perform whole exome (DNA) and RNA sequencing. We will use the .fastq.gz files we will use for this research. We will discuss the fastq files at the end of this module. The next section will show all members of this class, to process and analyze the data.

Protocol for sample ext

- Isolate fetal component of
placenta from viable tissue
dead tissue

Video. Sample Collection.

In this video, Wilson Lab Research Scientist, Dr. Seema Plaisier, discusses the origin of the placenta samples we will be using in this class, from patient criteria to sample extraction and sequencing.

[View Transcript](https://canvas.asu.edu/courses/122165/files/54792269?wrap=1) (https://canvas.asu.edu/courses/122165/files/54792269/download?download_frd=1)

Genomics and gene expression

What is transcriptomics?

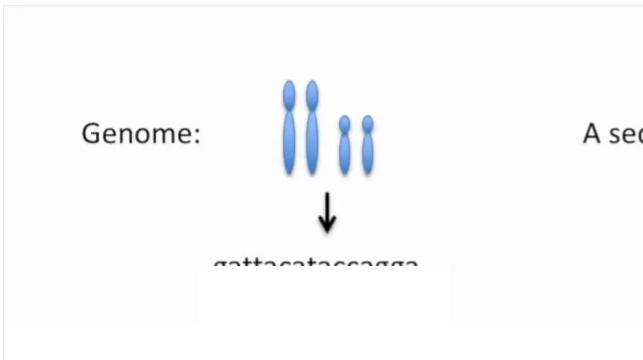
In this course, we will be using next generation sequencing to measure global gene expression, which is one type (sometimes referred to as transcriptomics). Genomics is the study of the genome (the complete set of genes for those genes interact with one another and the environment. Some of the ways genomics can be used are to predict disease, evaluate drug effectiveness, or to understand biological processes. Transcriptomics is when we assay mRNA in a cell. Using high-throughput sequencing technologies, we can lyse cells from specific tissues, stabilize

sequencing to detect and quantify the set of mRNA transcripts being transcribed, giving us a window into genes that are specifically involved in specific cellular processes and conditions. We can do this in different tissues in different conditions to tease out what genes are specifically involved in which contexts, thereby allowing us to sleuth out mechanisms underlying the biology of a gene. A gene is considered expressed when it is actively transcribed; some genes are constantly expressed while others are not expressed, and vary with tissue, age, sex, and environmental response (e.g., stress). If you need more information on the biology used in bioinformatics, refer to your [Biostar Handbook: Biology for Bioinformaticians](https://www.biostarhandbook.com/biology-for-bioinformaticians.html) (<https://www.biostarhandbook.com/biology-for-bioinformaticians.html>).

RNAseq data pre-processing approaches

In this course we will focus on the RNA that is transcribed in the cells of the human placenta tissue. We extracted RNA from XX and XY placentas and sequenced that RNA using RNAseq (RNA sequencing).

We have already done the RNAseq data processing for you but it is important as a researcher to understand where the data came from and how it was processed. In the video below, we will explore what RNAseq is, the processes behind RNA sequencing, and how to analyze RNAseq data.



Video. StatQuest: A Gentle Guide to RNAseq.

In this video, you will learn what RNAseq is, how biological samples go from *in vitro* (test tube) to *in silico* (virtual), and what that means for downstream analysis.

[View Transcript \(<https://canvas.asu.edu/courses/122165/files/54792276?wrap=1>\)](https://canvas.asu.edu/courses/122165/files/54792276?wrap=1) [\(https://canvas.asu.edu/courses/122165/files/54792276/download?download_frd=1\)](https://canvas.asu.edu/courses/122165/files/54792276/download?download_frd=1)

(<https://canvas.asu.edu/courses/122165/files/54792276?wrap=1>) RNAseq is a type of high throughput or next-generation sequencing method that tells us which genes are active by measuring the amount they are transcribed. In the last section, we learned that genes that are active produce mRNA transcripts while those that are inactive do not, but there could also be differences in the amounts of mRNA transcripts being produced between two conditions. If we take normal cells as our control and compare the expression levels of genes between mutated cells, we will see that some genes are expressed in both conditions at the same level, while other genes have increased expression in one condition producing more transcripts.

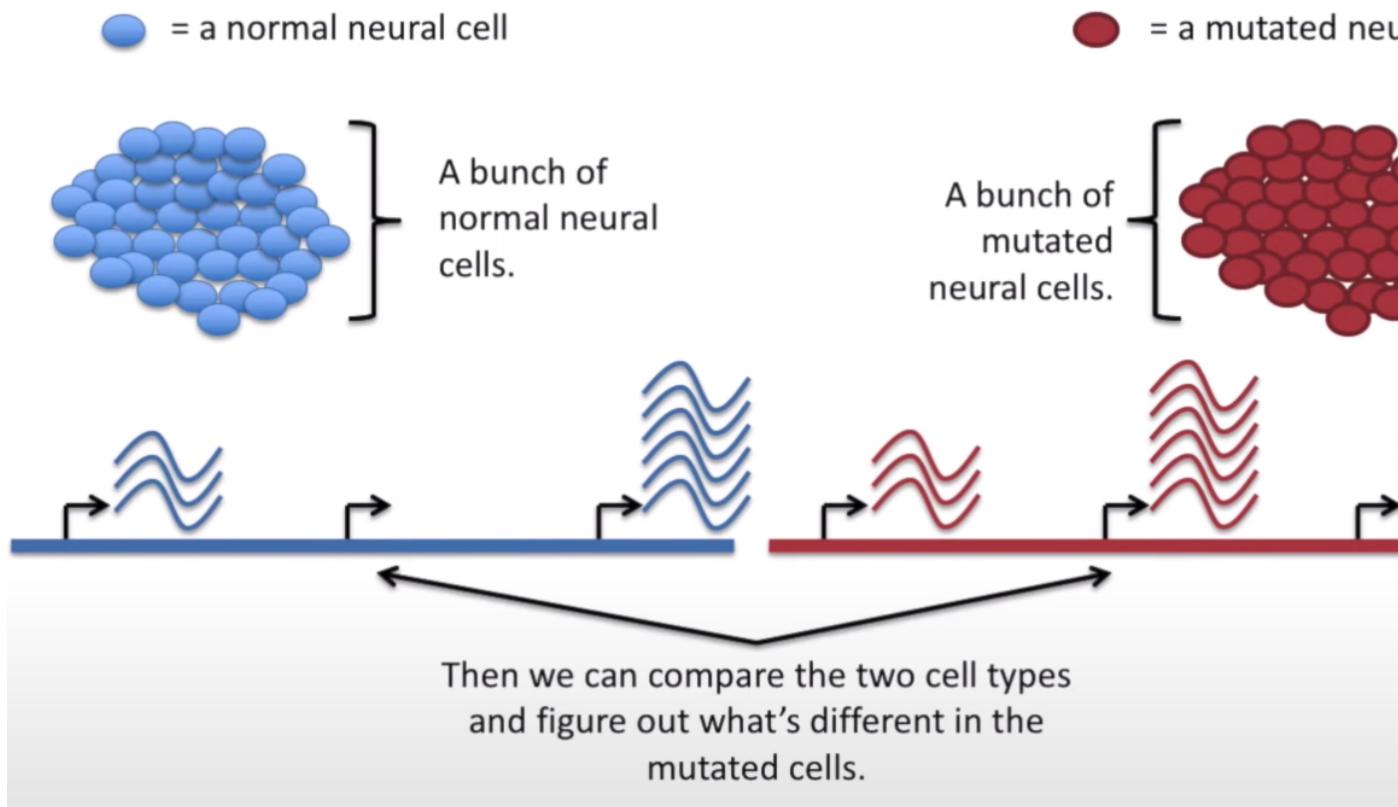


Figure. RNA sequencing overview.

We can use RNAseq from normal cells as our control and compare the results to mutated cells to find differences in gene expression. In the example, each right arrow indicates a transcription start site, or beginning of a gene, and the wavy lines beside them depict the amount of mRNA transcripts being produced by each gene.

There are three main steps for RNAseq: preparing a sequencing library, sequencing, and data analysis. Let's take each step in turn.

Step 1: Preparing a sequencing library

Prepare a sequencing library (based on the protocol from Illumina, the company who makes the library prep kit w

Preparing an RNA-seq library

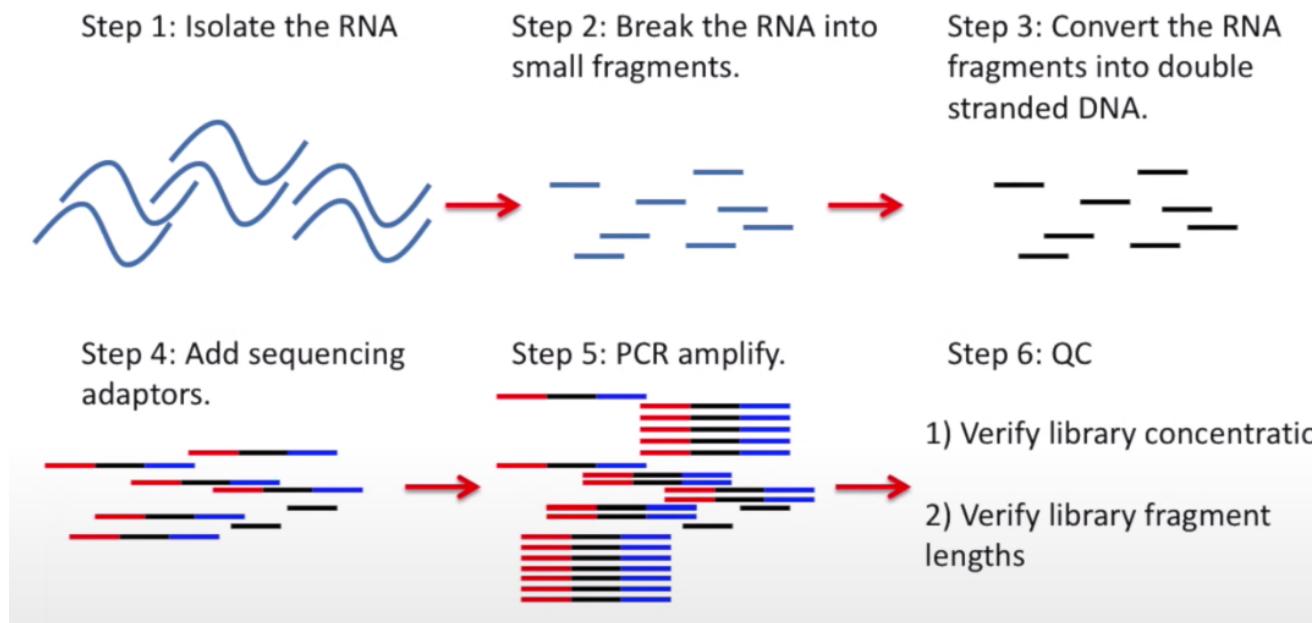


Figure. Preparing an RNaseq library.

The steps in preparing an RNA-seq library are to isolate the RNA, break it into small fragments, convert it into cDNA (complementary DNA), add sequencing adaptors, PCR amplify, and then perform QC to verify library concentration and fragment length.

1. Isolate the RNA from cells that were treated to prevent RNA degradation (stored in cold temperatures with RNase inhibitor). Break the RNA into fragments because the sequencing machine performs best with 200-300bp fragments and can be thousands of bases long.
2. Convert RNA into double stranded cDNA because DNA is more stable than RNA and can be amplified using a polymerase chain reaction (PCR).
3. Add sequencing primers which will allow the sequencer to read the fragments and identify the nucleotide sequence.
4. PCR amplify the library of fragments from the adapter to the end of the fragment.

5. Perform quality control to verify library concentration and library fragment lengths to ensure the sequencing results are accurate and precise (quality control to make sure there were any unexpected technical problems).

Step 2: Sequencing

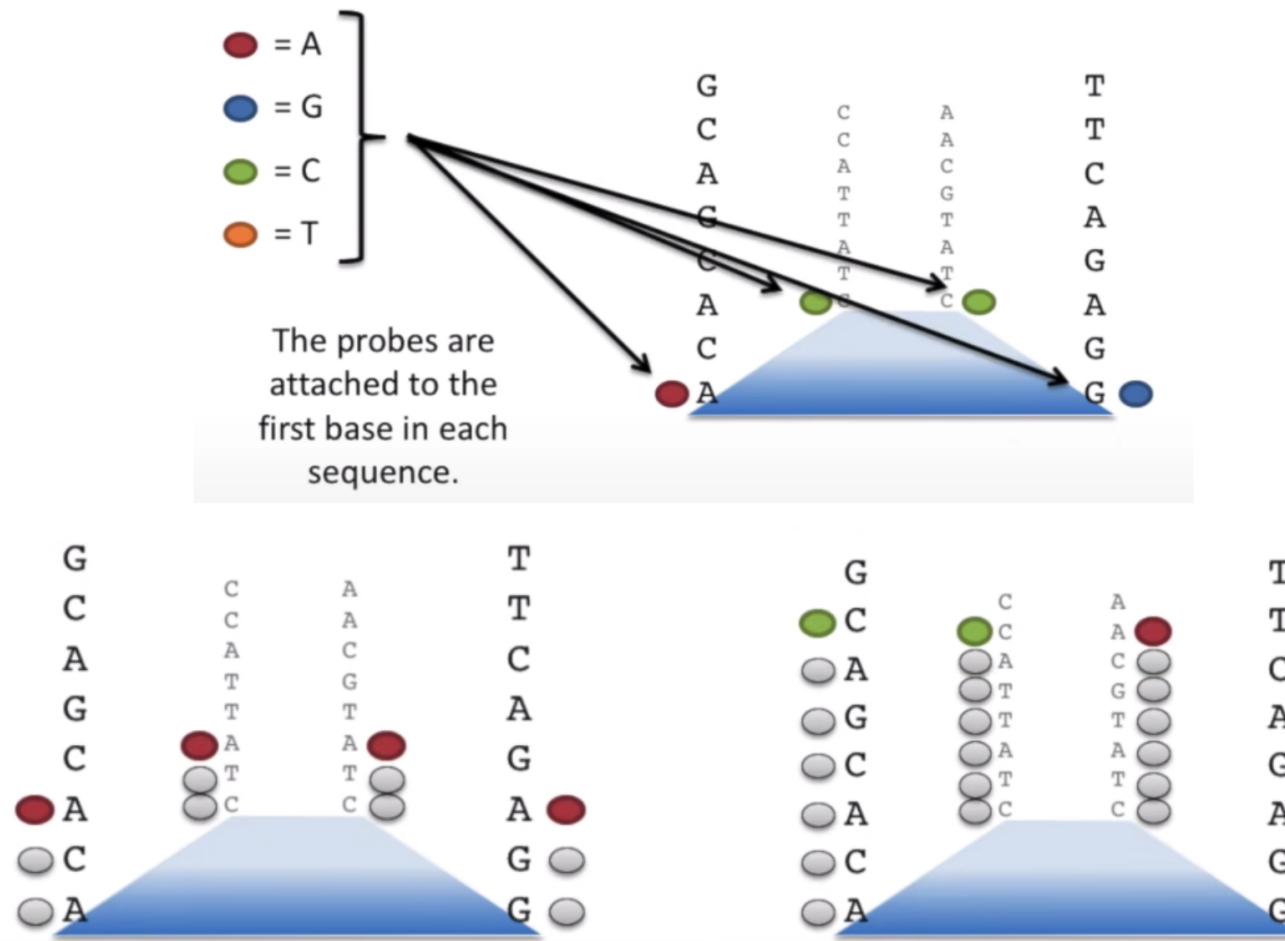


Figure. The process of sequencing using a flo cell and fluorescent probes.

The probes attach to the first base in each sequence, are photographed, and then washed away; this process is repeated until the full sequence of each fragment is determined.

After the library is prepared with fragmented cDNA that has been PCR amplified, the fragments are laid out in a grid. The machine attaches fluorescent probes to the bottom base of each sequence with a different color for each base. The machine takes an aerial photo from above where it can detect each base by color; after it registers the lit base in the image, it washes off the color from the probe. Now the machine can attach the next row of probes and repeat the process until it has sequenced the entire sequence of each fragment.

Since the average run can be around 400,000,000 reads, probe colors can be too light to detect or have “low diversity” if there is a probe of one color in a region saturated with probes of another color, making it difficult to identify with confidence what base is at that position. Either of these would be considered low quality bases and, if many are in the same sequence, can be flagged for trimming or we trim.

Step 3: Data Analysis

After reading the sequences from these fragments, we are given a file that contains the nucleotide sequences for a huge number of software packages designed to help you make sense of these large files of read sequences. The first step in this analysis is to use software to match these reads to genes, which is called **data preprocessing**. Using software to compare the amount of gene expression between samples in one set of samples versus another is called differential gene expression analysis.

Data Preprocessing

This diagram shows the typical steps of sequencing data preprocessing along with the specific software we used. Starting from the read sequence files from RNA sequencing (.fastq), we first used tools like FastQC and MultiQC for sequence quality checking and the trimming tool bbduk to trim low quality sequence and adapter sequence. Next, the reads are aligned to the reference genome so we know what genes those reads map to. Once aligned, we do some post-processing analysis to make the alignment output files easier to work with (sorting, indexing, etc). Finally, the reads mapped to each gene are tallied and used as a measure of gene expression. We will go into more detail with these steps, specifically as they relate to our research aims.

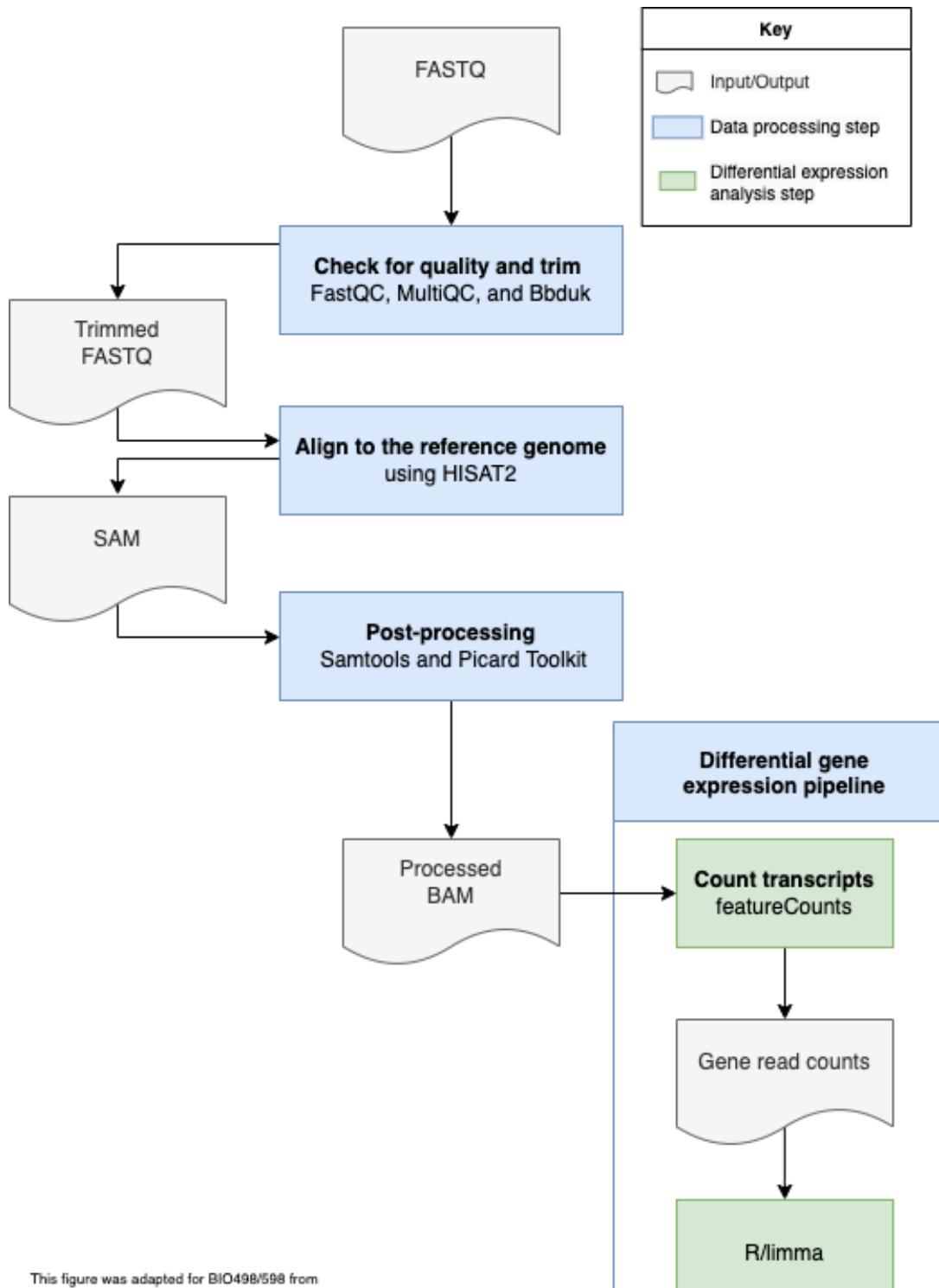


Figure. RNAseq data preprocessing and differential expression workflow.

This workflow was adapted for this course from <https://github.com/SexChrLab/Nasonia> (https://github.com/)

The

lab followed the data generation protocol (left) of checking the sequencing FASTQs for quality, trimming, aligning to a reference genome, and post-processing to generate gene counts. You will conduct the differential gene expression analysis (boxed in blue on the right) of taking the counts and processing them through the R/limma-voom pipeline.

Step 1: FastQC and Trimming

FastQC is a tool that provides data visualizations that can help you to identify samples that had technical issues or of your overall data quality. This step is done routinely and can be referred to at any point in the study if we suspect problems. Trimming software is often used to trim or filter sequencing reads with poor quality. We will use the trim tool (<https://jgi.doe.gov/data-and-tools/software-tools/bbtools/bb-tools-user-guide/bbduk-guide/>) which we have been able to use to improve data from studies with many different types of sequencing data.

The first figure FastQC gives you shows the quality scores across your reads. This is an example of this figure for a set. You can see as you go towards the end of the reads, the quality scores fall to zero, going from a warning level to a poor level marked red. Trimming software can be used to filter poor quality sequence reads, reads where the sequence of which base was at a specific position is not that confident. Results as bad as what is shown in this example would completely be removed with trimming software, but most data with a smaller proportion of poor quality reads can be filtered successfully. Trimming software can be used in conjunction with other steps such as alignment to the reference genome proceed more smoothly. The parameter ‘trimq’ in bbduk is used to set a minimum quality score. If bases in a read are less than the minimum threshold, they are trimmed off in the output fastq sequence. The parameter of bbduk called ‘minlen’ sets a minimum length of read that has to be left after quality trimming. In the figure, the thumb rule is that if trimming removes over half of the read, then we should filter out that read as a low quality read. So, for the data we will be using, the reads were 150 base pairs in length, so the original value of minlen we chose for data processing was 75. We will be investigating the effect of changing those values in this research project.

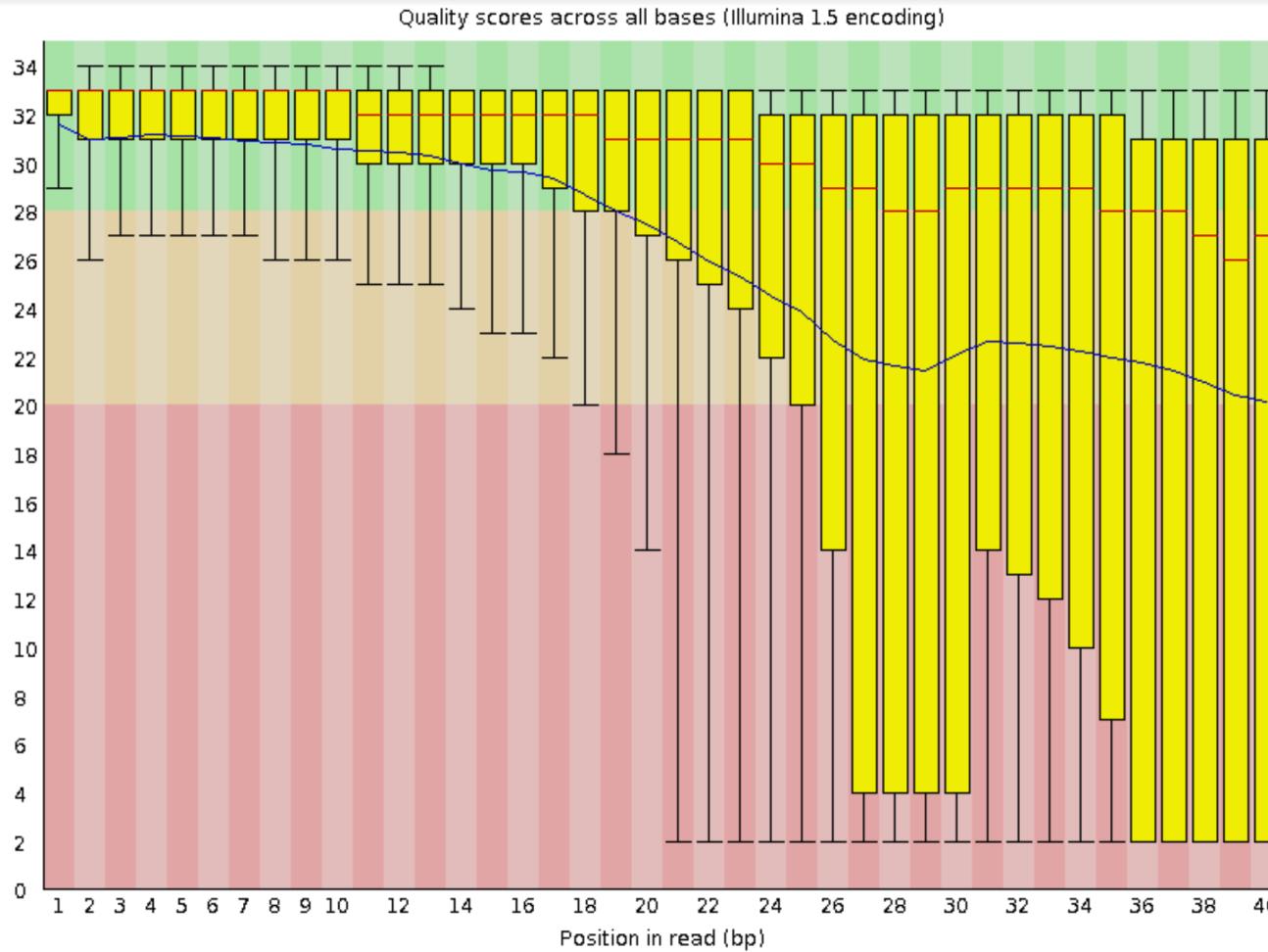


Figure. Quality scores from poor quality RNAseq data.

The boxplot was generated using FastQC and shows the quality scores across all bases using Illumina 1.5 encoding. The x-axis is the position of the read in basepairs (bp) and the y-axis shows the quality score (PHRED).

Another FastQC figure that is relevant to our research shows the percentage of Illumina sequencing adapters. The difference in adapter content before and after trimming with bbduk. Each line on the figure represents one sample. The x-axis of the figures show the position along the 100-base pair reads and the y-axis shows the proportion of reads to an Illumina adapter used during library preparation at that position. We should not see Illumina adapter sequence because they are supposed to be removed as part of the Illumina library generation protocol before the sequencing.

(.fastq). This is reflected in the max value of the y-axis is 7%. Since the library prep kit adapter removal step is not in the untrimmed data figure on the left that the overall percent adapter sequence is low, but some samples have been expected for high quality data (ranges marked in yellow for warning, marked in red for poor data quality). Removal of adapter sequence is one of the key functions of trimming software and improves alignment of these reads to the reference genome. We can see that after computationally trimming with bbduk, the percent adapter sequence falls below 5 percent for all samples. For analysis, any data set for which trimming was done at all included removal of Illumina adapter sequences, regardless of the trimq and minlen.

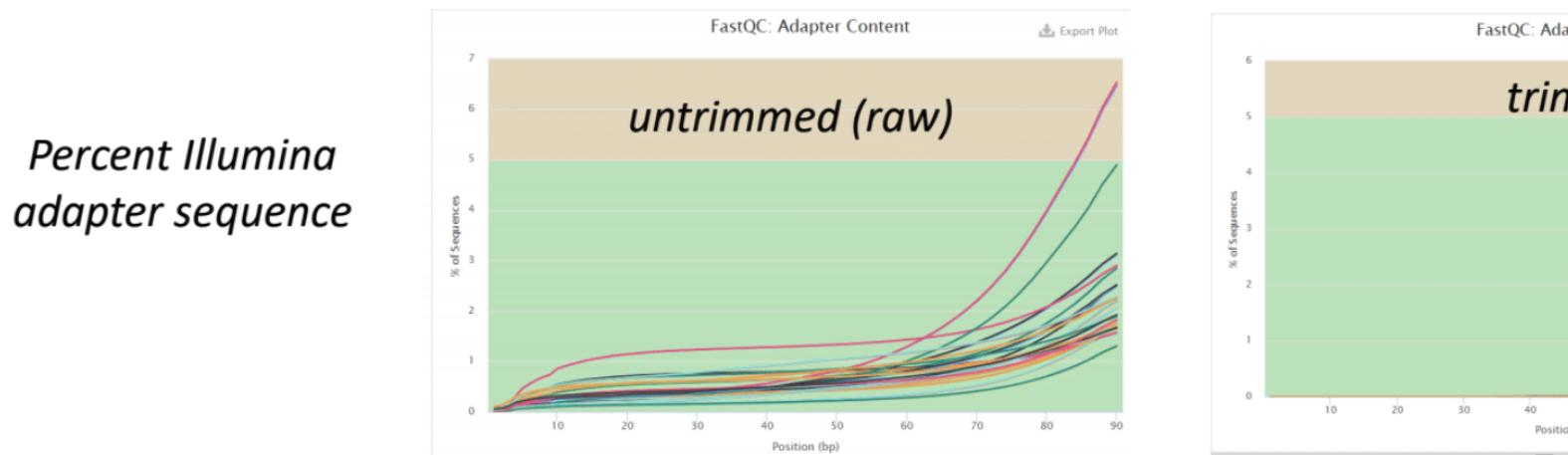


Figure. FastQC figures of percent Illumina adapter sequences.

We can use this to compare adapter content from our untrimmed batch to the trimmed batch.

Step 2: Alignment

To perform alignment, we can download the sequence of the entire genome for the species of the organism whose genome we are analyzing, human in our case. We used the GRCh38 build of the human genome, an established human genome assembly with many publications. Since we are using sequencing at the transcriptome level, we downloaded the entire set of transcripts from the genome build to use as an alignment reference. Software to align short read sequences to the reference genome, such as HISAT2, was used to do this step.

Aligning to autosomes (non sex chromosomes) is generally quite straightforward as genes can be expressed from multiple copies of each autosome which are nearly identical. Alignment to the sex chromosomes is more complicated because

chromosomes are quite different from one another in size and sequence, but contain specific regions of identity c evolutionary history. We need to make additional considerations in order to accurately align reads to the sex chro

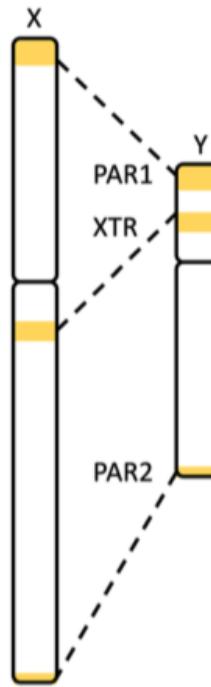
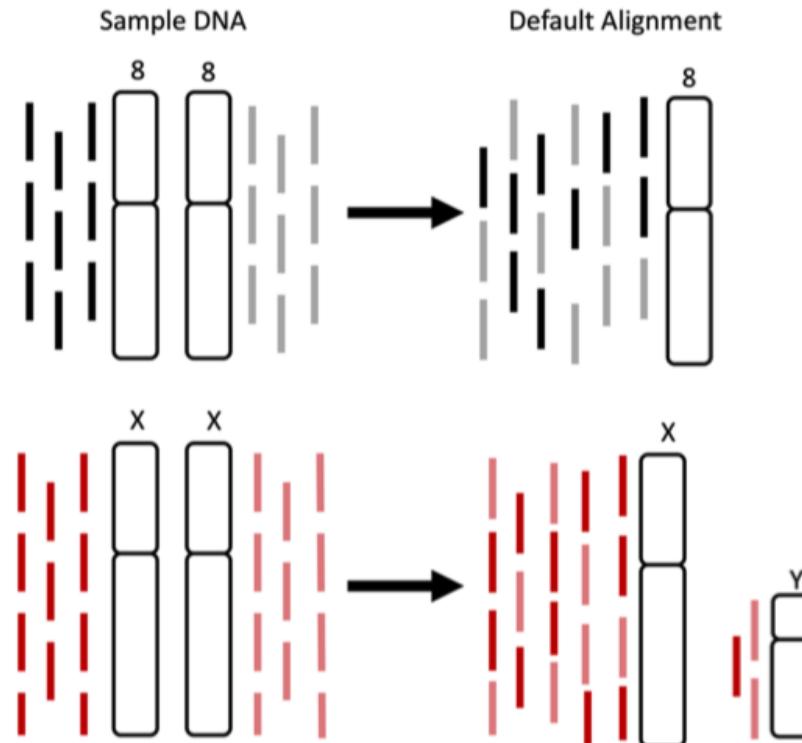
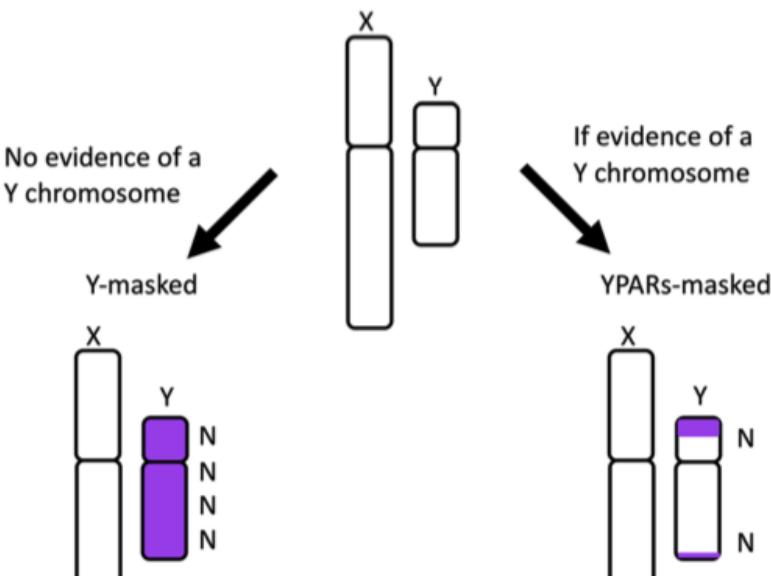
A X and Y sequence homology**B** RNA-seq alignment to a default reference genome**C** Sex chromosome complement informed alignment



Figure. Sex-informed alignments to a reference genome.

- (A) X and Y sequence homology, (B) RNAseq default alignment to reference genome, (C) sex chromosome complement informed alignment.

This figure depicts how we do alignments taking sex into account. Figure (A) illustrates the sequence homology chromosomes, XTR and PAR regions. In (B), the standard RNAseq alignment shows that the sample DNA from chromosome 8 align to the sequence of chromosome 8 in the reference genome. This works for autosomes, but the sex chromosomes? Since there is sequence for both the X and Y chromosomes in the reference genome sec download, so reads from the homologous XTR or PAR regions from an XX sample could randomly align to chrom the sample has no chromosome Y. Similarly for an XY individual, reads that are mapping to the regions found in t chromosomes will be randomly aligned to either the X or Y chromosome sequence in the reference genome. Thi samples that have a Y chromosome differently from samples that do not have a Y chromosome.

The solution (and way we processed our dataset) is depicted in figure (C), using a sex complement informed alig technique called masking, which replaces specific sequences in the reference genome with “N”s, such that no se those masked areas and those reads are forced to align somewhere else in the genome. To prevent shared seq misaligning, we prepare a reference genome with the Y chromosome sequence masked to use for alignment of r samples. We prepare a second reference genome sequence file with the PARs masked for alignment of XY sampl not really know where those reads are supposed to map. Using this alignment method, we were able to see diffe of many genes on the sex chromosomes and characterize those genes in a manuscript.

Step 3: Gene Expression Quantification

After alignment, the number of reads that map to a specific gene is used as a way to quantify gene expression, s align our reads, the better our estimation of gene expression and the more likely we are to pick up true difference for genes on the sex chromosomes. We use a program called featureCounts to take the coordinates of genes giv file we downloaded from the same database as the reference genome sequence, featureCounts uses the output to count the number of reads that map to each gene. These are given as a single file per samples. We use a sir merge all of those files into one table that contains the counts for each gene in all samples in our study.

Differential Expression (DE) Analysis

You will learn a lot more about the differential gene expression pipeline we will be using later in the course, but we will provide a broad overview here. Our aim after data processing is to find genes that are truly expressed at a different level in XX female and XY male placentas. To do this, we first filter out genes with no reads detected in our samples. We next normalize the gene expression levels to account for technical issues such as the fact that longer genes are going to have more reads coming from them than shorter genes. Once the data has been normalized, we can use data visualization techniques to group genes together based on their gene expression profiles to see if they cluster by sex as well as use molecular modeling to find genes that show significant differences in expression between XX female and XY male placentas.

How we came up with the research question for this course

Now that you have a good idea of the biological system we are studying and the way that we are measuring gene expression in the placenta, we wanted to tell you how the idea for this project came to be. As with all scientific research, it starts with an observation. In the video below, Dr. Seema Plaisier explains how applying an established data processing workflow to a new set of samples led to the research question we are interested in. She discusses how our lab applied trimming parameters from the original study to the placenta RNAseq we will be analyzing in this course to a new set of placenta RNAseq. In the original research, the parameter values we used (trimq=30, minlen=75) gave us data that passed all quality control requirements and contained very little loss in overall read count after trimming. However, in the sequencing data for the new placenta samples, trimming with these parameters resulted in overfiltering and erroneous results. This got us thinking about how trimming affects downstream analysis. In this class, we will be processing the samples from the original study using a range of values for trimming parameters and investigating what changes in the lists of genes differentially expressed in XX vs XY placentas.

Troubleshooting

- Step 1 of data processing is **trimming**: removes nucleotide identification or corresponds up in the results
- **Fc** filters and data quality

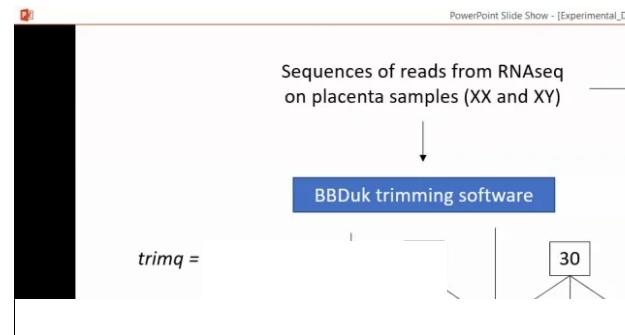
Video. Origin of the research question.

In this video, Wilson Lab Research Scientist, Dr. Seema Plaisier, discusses the origin of the research question and how it was organically derived from replicating research workflows.

[View transcript. \(https://canvas.asu.edu/courses/122165/files/54792240?wrap=1\)](https://canvas.asu.edu/courses/122165/files/54792240?wrap=1) ↓
[\(https://canvas.asu.edu/courses/122165/files/54792240/download?download_frd=1\)](https://canvas.asu.edu/courses/122165/files/54792240/download?download_frd=1)

As explained in the video above, the main underlying question we are asking in this course is: ***How does trimming affect gene expression analysis?***

Since we were questioning the trimming parameters of the original study, we wanted to get an idea of what other researchers did a quick search for published research papers that used the same trimming software (BBduk) to trim RNAseq samples. We found that of all the trimming parameters, trimq (trimming quality) and minlen (minimum length of reads after quality trimming) were the most frequently reported in the methods and the most highly variable. We used this to determine the range of values that spanned the range observed in literature and seemed appropriate given our experience with using BBduk in projects.



Video. Experimental design.

In this video, Wilson Lab Research Scientist, Dr. Seema Plaisier, discusses the experimental design of our research, the trimming parameters we are investigating, and why.

[View transcript. \(https://canvas.asu.edu/courses/122165/files/54792262?wrap=1\)](https://canvas.asu.edu/courses/122165/files/54792262?wrap=1) [\(https://canvas.asu.edu/courses/122165/files/54792262/download?download_frd=1\)](https://canvas.asu.edu/courses/122165/files/54792262/download?download_frd=1)

For this course we chose to investigate a range of 0, 10, and 30 for trimq and 10, 30, and 75 for minlen to test for differential gene expression analysis. We used a workflow management tool called Snakemake to generate the files for this study in an automated, reproducible way. The datasets we will be working with consist of one set of untrimmed reads and nine sets of trimmed datasets. For the untrimmed reads, we did not apply any trimming, therefore, do not use minlen parameters. The untrimmed dataset will serve as a control and the entire class will be using the same controls. We generated nine trimmed datasets by running the untrimmed data with bbduk using all combinations of trimq (0, 10, 30, and 75) and minlen (10, 30, and 75). The untrimmed and trimmed data sets were independently run through data preprocessing steps, putting them through differentially gene expression analysis so that the results can be examined to find trends. Each student will analyze one trimmed dataset then we will compare results in groups and as a class to analyze the downstream effects of different trimming parameters.

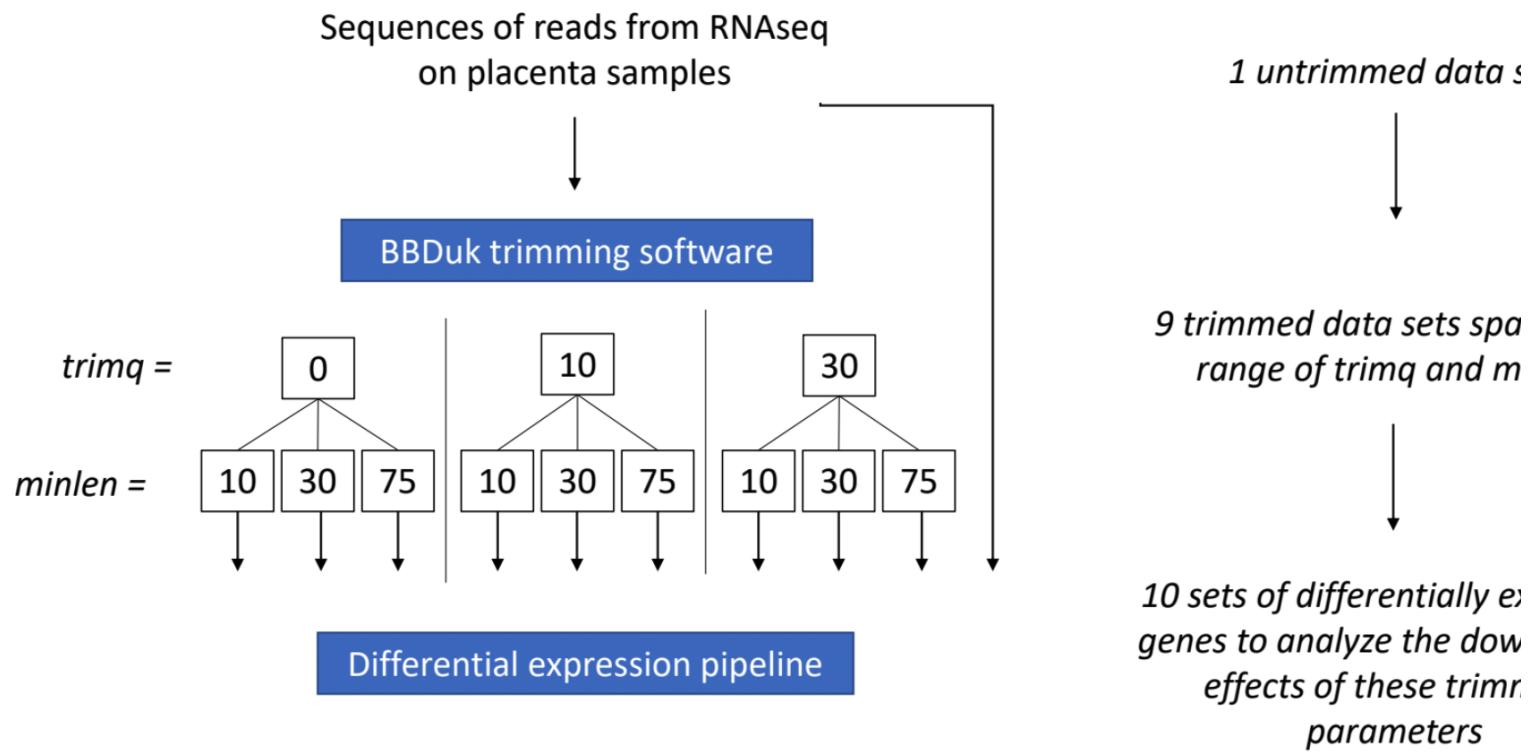


Figure. Experimental design of trimming parameters.

We will be using one (1) untrimmed dataset and nine (9) trimmed datasets with a range of parameters to generate ten (10) sets of differentially expressed genes to analyze.

In the figure above, you can see the nine trimmed datasets that have a different value for trimq and minlen parameter generated for you. Each person has been assigned the control plus one of the following:

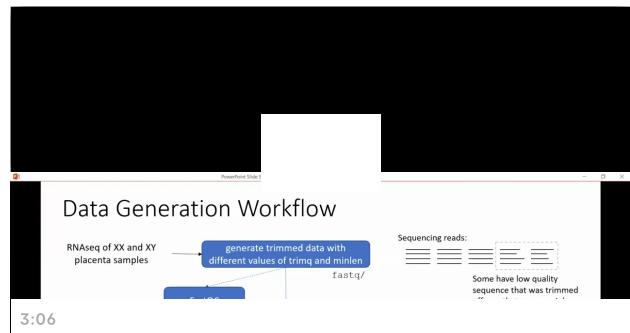
Table. List of the trimming parameter combinations selected for this experiment.

Directory name	Trimq	Minlen

trimq_NONE_minlen_NONE	NONE (control)	NONE (control)
trimq_0_minlen_10	0	10
trimq_0_minlen_30	0	30
trimq_0_minlen_75	0	75
trimq_10_minlen_10	10	10
trimq_10_minlen_30	10	30
trimq_10_minlen_75	10	75
trimq_30_minlen_10	30	10
trimq_30_minlen_30	30	30
trimq_30_minlen_75	30	75

Data Generation Workflow

This video will provide a quick recap of the data preprocessing steps to create the gene count matrices you will be using in class.

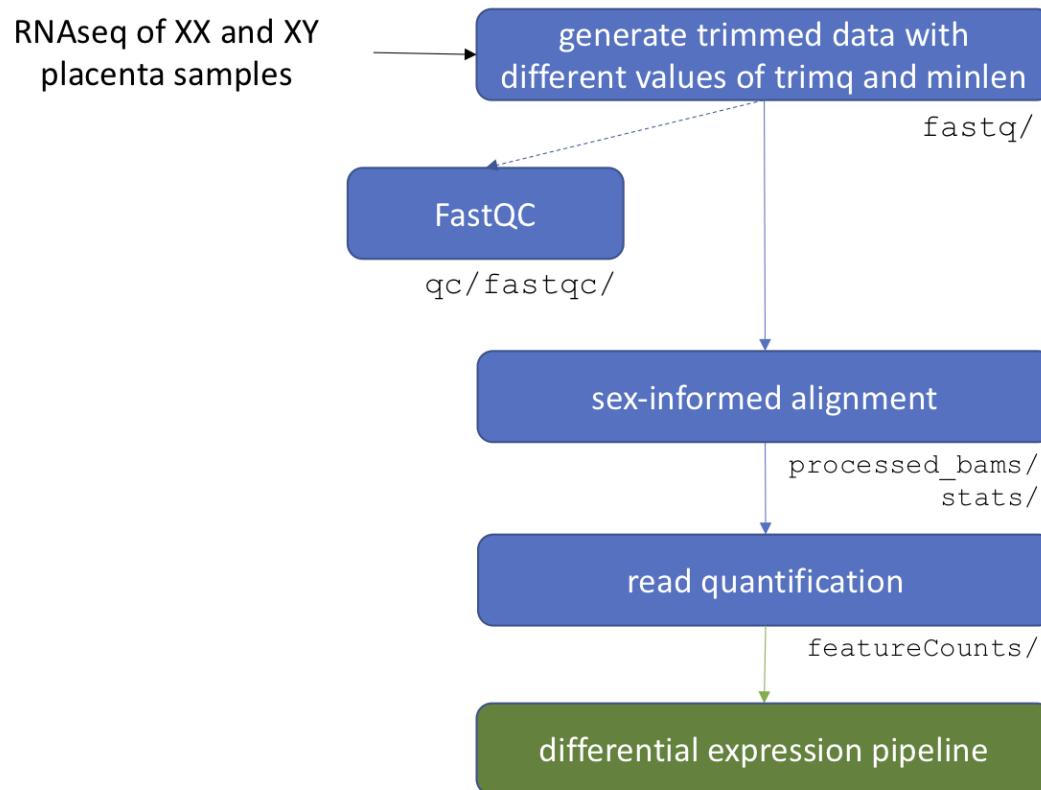


Video. Data generation workflow.

In this video, Wilson Lab Research Scientist, Dr. Seema Plaisier, discusses the data generation workflow that has already been conducted on the dataset you will be using throughout the course. RNAseq samples were trimmed with the different parameters, evaluated for quality, aligned to a sex-informed reference genome, and then quantified into gene counts.

[View transcript. \(https://canvas.asu.edu/courses/122165/files/54792240?wrap=1\)](https://canvas.asu.edu/courses/122165/files/54792240?wrap=1) ↴
[\(https://canvas.asu.edu/courses/122165/files/54792240/download?download_frd=1\)](https://canvas.asu.edu/courses/122165/files/54792240/download?download_frd=1)

Data Generation Workflow



Sequencing reads:

A sequence of horizontal lines representing sequencing reads. Some lines are solid black, while others are dashed grey. A dashed box highlights a group of three lines, with the text "Some have sequences off, any too long were" explaining the filtering process.

XX reads to ref Y-masked: XY reads

Two sets of horizontal lines representing reads. The top set, labeled 'XX reads to ref Y-masked:', has several lines crossed out with a red bar. The bottom set, labeled 'XY reads', has several lines crossed out with a blue bar. Below these, two genes are shown: geneA with a value of 3 and geneB with a value of 5, indicating expression levels.

geneA higher expression in X
geneB higher expression in X

https://github.com/SexChrLab/Genomics_CURE/tree/main

Figure. Data Generation Workflow.

The diagram on the left shows the steps that have already been conducted for you (blue) and the analysis you will be conducting (green). The figures on the right are examples of the corresponding step to the left.

Module 1.1 Additional Resources

- [Genomics ↗ \(https://www.genome.gov/about-genomics/fact-sheets/A-Brief-Guide-to-Genomics\)](https://www.genome.gov/about-genomics/fact-sheets/A-Brief-Guide-to-Genomics)
- ↗ (https://www.genome.gov/about-genomics/fact-sheets/A-Brief-Guide-to-Genomics) Placenta (detailed article)
Burton GJ, Fowden AL. The placenta: a multifaceted, transient organ. Philos Trans R Soc Lond B Biol Sci. 2015;368(1663):20140066. doi: 10.1098/rstb.2014.0066. PMID: 25602070; PMCID: PMC4305167.
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4305167/> ↗ (https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4305167/)
- [Central Dogma ↗ \(https://www.nature.com/scitable/topicpage/translation-dna-to-mrna-to-protein-393/\)](https://www.nature.com/scitable/topicpage/translation-dna-to-mrna-to-protein-393/):
- [Biostar: Biology for Bioinformaticians ↗ \(https://www.biostarhandbook.com/biology-for-bioinformaticians.html\)](https://www.biostarhandbook.com/biology-for-bioinformaticians.html)
- [Snakemake ↗ \(https://snakemake.readthedocs.io/en/stable/tutorial/basics.html\)](https://snakemake.readthedocs.io/en/stable/tutorial/basics.html)
- The X Chromosome
[TED-ed: secrets of the x chr \[5:05\] ↗ \(https://youtu.be/veB31XmUQm8\)](#)
- Clifton. (2010). Review: Sex and the Human Placenta: Mediating Differential Strategies of Fetal Growth and Sex (Eastbourne), 31(3), S33–S39. <https://doi.org/10.1016/j.placenta.2009.11.010> ↗ (https://doi.org/10.1016/j.placenta.2009.11.010)
- Olney, Brotman, S. M., Andrews, J. P., Valverde-Vesling, V. A., & Wilson, M. A. (2020). Reference genome analysis informed by the sex chromosome complement of the sample increase ability to detect sex differences in gene RNA-Seq data. *Biology of Sex Differences*, 11(1), 1–42. <https://doi.org/10.1186/s13293-020-00312-9> ↗ (https://doi.org/10.1186/s13293-020-00312-9)