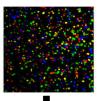# .bcl file
**binary base call format**

Raw output of sequencing machine; contains the **base call** and the **quality (confidence)** of the base call after each cycle (for example 2 x 100 chemistry has 100 forward and 100 reverse cycles) by cluster. Can be hundreds of GB worth of data.



## Demultiplexing

Demultiplexing is the process of sorting base calls into separate files by their unique indices

Performed with Illumina's bcl2fastq tool. Almost always done by the sequencing centre.

---

# .fastq file
**also shortened as .fq**

Information contained in .bcl file organized by sample and by sequence. Contains **every sequence read** (A,T,C, G, and N) and the **quality score (confidence) of each base call**. Quality scores represent the probability that the base was called in error. The quality score is reported in ASCII characters, 42 single digit numbers, characters, and letters. This makes fastq format easily human-readable, because each base (A, T, C, G, or N) is associated with a single digit quality score (!, *, 6, C, J, etc).

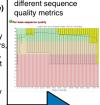In fastq format, each sequence is represented by four lines:

1. A header with seq info, including cluster coordinates from flow cell
2. The sequence
3. Qual score separater "+"
4. The quality scores

**Example of one read in .fastq format:**
@SIM:1:FCX:1:15:6329:1045 1:N:0:2
TCGCACTCAACGCCCTGCATATGACAAGACAGAATC
+
<>;##=>9=AAAAAAAAAA9#:<#<;<<<????#=

**Single end (SE)** sequencing runs produce one .fastq file (**R1**) per read, and **paired end (PE)** runs produce two files (**R1 - forward reads and R2 - reverse reads**)

## Quality control

Based on the quality metrics and the type of data you have (PE, SE, ancient, or modern), you will **trim the reads to remove adapter sequences, N base calls, and reads with low base quality** (low confidence in call). If you have ancient PE data, you may also want to merge the reads where they overlap.
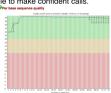
---

# quality filtered .fastq file
**also shortened as .fq**

Sequencing reads will be of various quality and need to be assessed. The program **fastqc** is used to visualize different sequence quality metrics



Common quality issues in aDNA that should be trimmed and filtered out of your fastq files:

-**Adapter readthrough** = your sequencing chemistry is longer than your fragments so the adapter sequence is read as part of the actual sequence.

-**Ns on the ends of reads**: complete adapter readthrough and poor quality data can result in Ns, meaning the sequencer could not call any base. These are not bases with typical aDNA "damage"

-**Low quality base calls**: low quality DNA, poor library construction, insufficient purificaiton can result in the sequencing machine not being able to make confident calls.



You may want to **merge** your forward and reverse reads. aDNA is short and reads are more likely to overlap. This helps remove contaminant DNA, which is generally longer.

---

Now that you are confident in your sequence read quality, you want to **map the reads to the reference genome of your target organism**.

There are different mapping algorithms, but essentially matching and mismatching bases between ref and seq reads are tabulated for an overall alignment score for each read. Common tools include BWA and bowtie.

A reference genome is a high quality genome assembly. There is generally one reference genome per organism that all scientists use. For example, there is one human reference genome that is regularly being updated as more sections of the genome are annotated (hg38).

## Mapping

---

# .sam/.bam files
**and their indexes (.sai and .bai)**

**SAMs are tsv** (tab-separated values) **files that contain a header and the alignment score for every read. A BAM file is the binary (compressed) version of a SAM file.** SAMs can be converted to BAMs, and vice versa. BAM files are not human readable (if you open it up it is a bunch of nonsense). Alignment files can be huge and take a lot of memory to parse. Compression into a BAM makes computation faster and more efficient. SAM/BAMs have associated index files (.sai/.bai) used for increased computatiion efficiency.

ATGCTGATGTAGTCGTAGCTG



SAM format includes header info beginnging with the "@" symbol; each row corresponds to a read, and each column stores information about the read alignment.

@HD VN:1.6 SO:coordinate
@SQ SN:ref LN:45
r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *

## Filtering bams

---

Unfiltered SAM/BAM files contain a lot of information we don't need (unmapped reads) and a lot of poor quality alignments (reads with bad alignments scores) that we don't want to include in final analyses.

Common to filter out:
-**Unmapped reads**
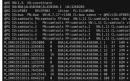-**Reads with low mapping quality scores**: low confidence that the read actually goes where placed.

-**Duplicate reads:** PCR and optical duplicates can make your data seem higher coverage. Generally we only want unique reads in our analyses.

-**Reads with more than one alignment:** one read can be aligned equally well to multiple segments of the genome. Since we can't be sure where they belong, they are often removed.

---

# quality filtered .sam/.bam files
**and their indexes (.sai and .bai)**

You now have a quality filtered (QF) aligned file (.sam/.bam)!



**samtools** is the primary tool for sam/bam file manipulation and is used for all NGS data, modern and ancient.

**A cleaned .sam/.bam file is generally analysis-ready,** meaning that now you can calculate mapping statistics (% reference coverage, depth of coverage, read length distribution, etc) and begin genome analysis of point mutations (SNPs aka SNVs) or structural vairants.

Ancient DNA that has not been fully treated with UDG to remove C->T and G->A transitions caused deamination usually requires one more step to rescale these positions so they are not called as variants.
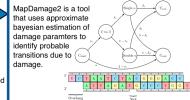
## Rescaling with MapDamage

Following .sam/.bam filtering, it is typical to generate mapping statistics using a tool called Qualimap.

The aDNA bioinformatic workflow is pretty standardized until this point. Contamination assessment, however, is organism-dependent and can be done many ways.
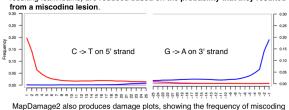
## Decontamination with pmdtools

---

# rescaled quality filtered .bam files **and their indexes (.bai)**

C->T and G->A transitions accumulate near the ends of aDNA fragments. Miscoding lesion may not be flagged as poor quality from sequencing or mapping, so they remainin in the quality filtered reads.

MapDamage2 is a tool that uses approximate bayesian estimation of damage paramters to identify probable transitions due to damage.

MapDamage2 modifies the .bam file such that **base quality scores** (data from the sequencing run that come from the fastq files that are also stored in the resulting .sam/.bams) **are reduced based on the probability that they resulted from a miscoding lesion.**



C -> T on 5' strand    G -> A on 3' strand

MapDamage2 also produces damage plots, showing the frequency of miscoding lesions on the read ends.

---

# decontaminated .bam files **and their indexes (.sai/.bai)**

Miscoding lesions due to degradation can also be used directly for read authentication.

**PMDtools (post-mortem damage tools) removes contaminant reads from .sam/.bam files** using a likelihood framework that models three processes:

**1.** C -> T changes resulting from post mortem damage
**2.** True C -> T biological polymorphisms (SNPs/SNVs)
**3.** C -> T changes caused by sequencing errors



**PMDtools assigns each read a PMD score, and reads with a score below the threshold are discarded**.

This approach has been shown to reduce contamination to aDNA datasets to negligible amounts while maintain important biological data that could be lost under more conseravtive approaches.

**PMDtools modifies the .sam/.bam by removing low-scoring reads.**

## Variant calling

---

# .vcf files

**Variant call format (VCF) files are tsv files with a standardized format used to store information about variants** observed in reads after they have been aligned to a reference genome. VCF files can store information about point mutations (SNPs/SNVs), insertions, deletions, and large structural variants.

Like .sam/.bam files, .vcf files have headers with mandatory and optional information. After the header, each row corresponds to a variant and provides information on the genomic position of the variant, the reference allele, and data used to assess the quality and confidence of the variant call, such as read depth and genotype likelihood. There are MANY variant calling tools. Some are better for aDNA.



.vcf files are commonly used as input for analytical tools, such as HaploGrep2, PLINK, and SNPeff.