



NGS data - Workflow, formats and programs

Thorfinn Sand Korneliussen

The GLOBE institute, University of Copenhagen

August 30, 2023



Outline

1 Background(GLOBE) and Learning objectives

2 Ancient DNA introduction

3 Low level formats

- data types
- Fasta and fastq files
- Computer exercise

4 Alignment and variant files

- SAMfiles
- Concepts and definitions
- VCF files
- Computer exercises

GLOBE Institute



- New institute
- Five sections
- Lundbeck Center for geogenetics
- Eight groups
- Korneliussen Group:
Rasmus Henriksen (tomorrow)
Abigail Ramsoe (Friday)
Isin Altinkya (computer exercises)



Learning objectives

Week1

Monday NGS intro

Tuesday Aligning data to a reference genome

Friday Modelling uncertainty of NGS data

After this week you will be able to navigate all commonly used files and software used in the context of NGS data. You will have a thorough understanding of how modern aligners work and will be able to apply this. NGS data is associated with high error rates and you will learn the theory of how this is dealt with.

Today

- Fileformats used in NGS data (fasta,fq,sam,vcf)
- Standard tools NGS projects(samtools,bcftools)
- Introduction to important concepts (depth,genotype, variable sites)



Standard workflow

- Lab stuff + sequencing
- Call bases and quality scores (FastQ)
- Alignment or *denovo* assembly. Alignment is topic for tomorrow

Genotypes

- Generate genotype likelihoods. Uncertainty is topic for friday.
- Estimate allele or haplotype frequencies
- Calculate genotype probabilities
- Call genotypes

Downstream analysis

- Variation detection
- Population genetic analysis

Ancient DNA introduction

After this sub-session you will:

- Be able to describe the characteristics of aDNA and the difficulties in studying aDNA.
- Understand important concepts (e.g. library, sequencing, complexity, clonality, DoC)
- Be able to participate in a meaningful discussion of aDNA projects and literature.

Sample types

- Bone
- Teeth, calculus(dental)
- Hair
- Fur, skin
- Shells
- Plants/seed/pollen
- Coprolites (Fossil feces)
- Sediments
- Ice cores
- Coral?

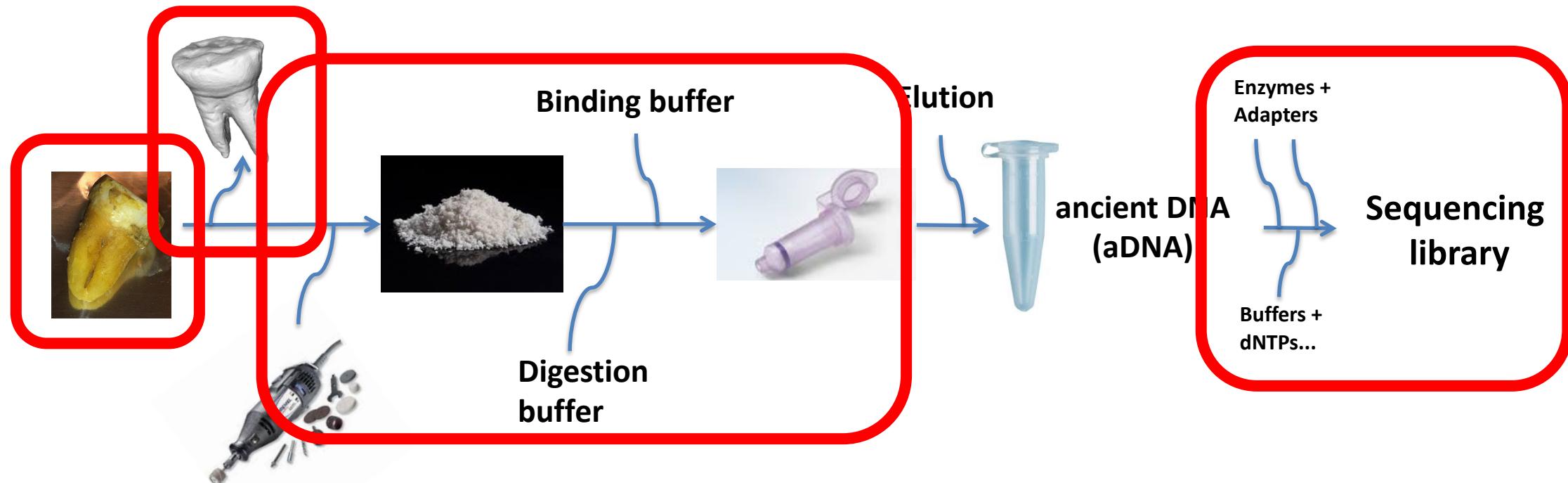


Sample collection

DNA wet lab work

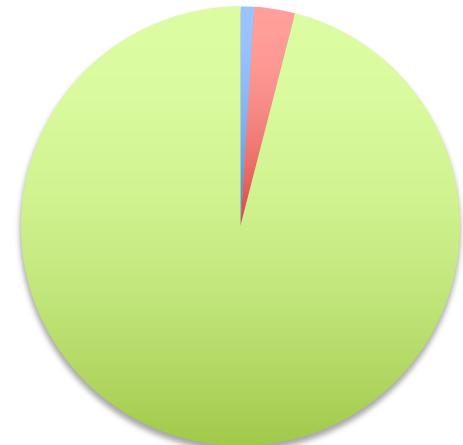


Data analysis



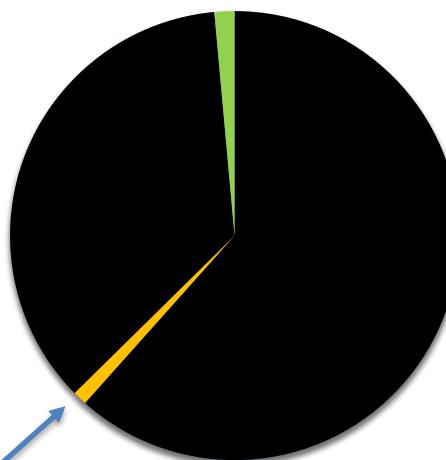
Low endogenous [DNA] in most ancient samples

Shotgun sequencing data (human sample)



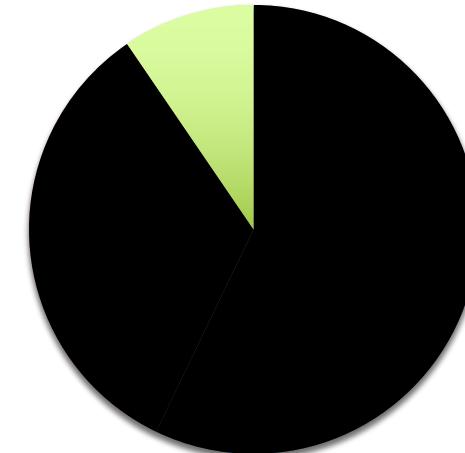
Modern

- Bacterial
- Unknown
- Human



e.g. *Y. pestis*

3- 14% Human



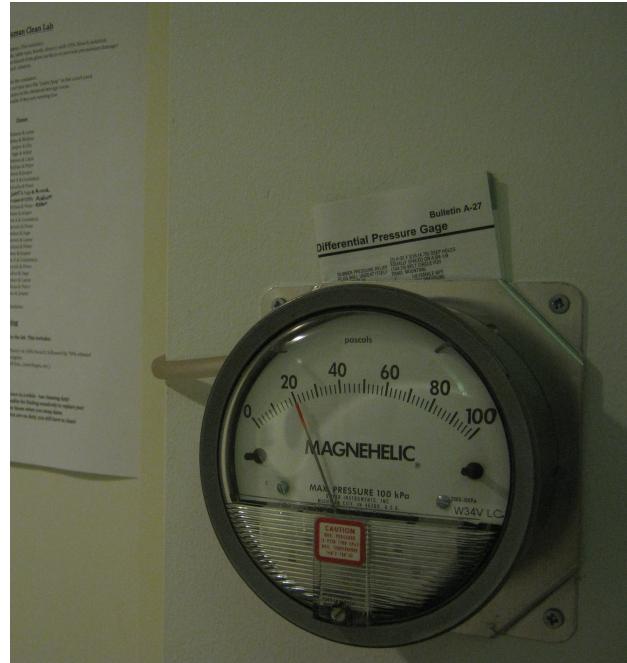
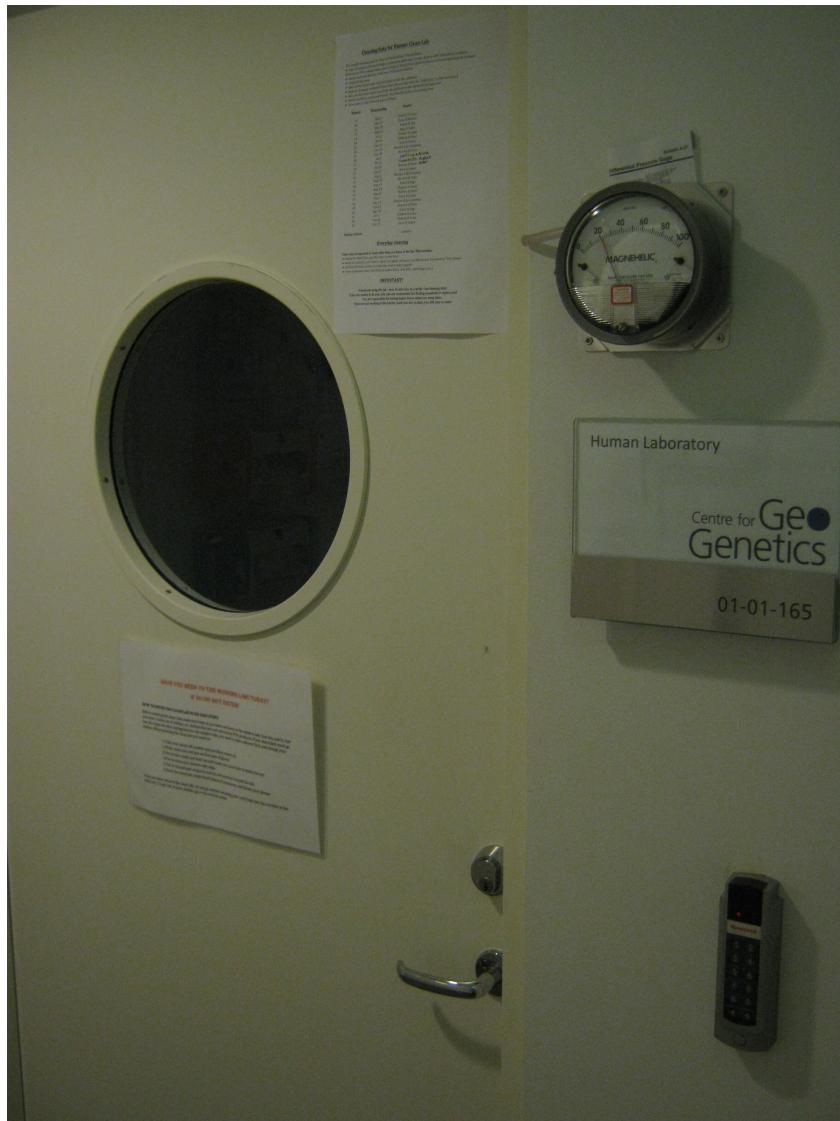
Ancient

Avoid contamination with modern human DNA by...

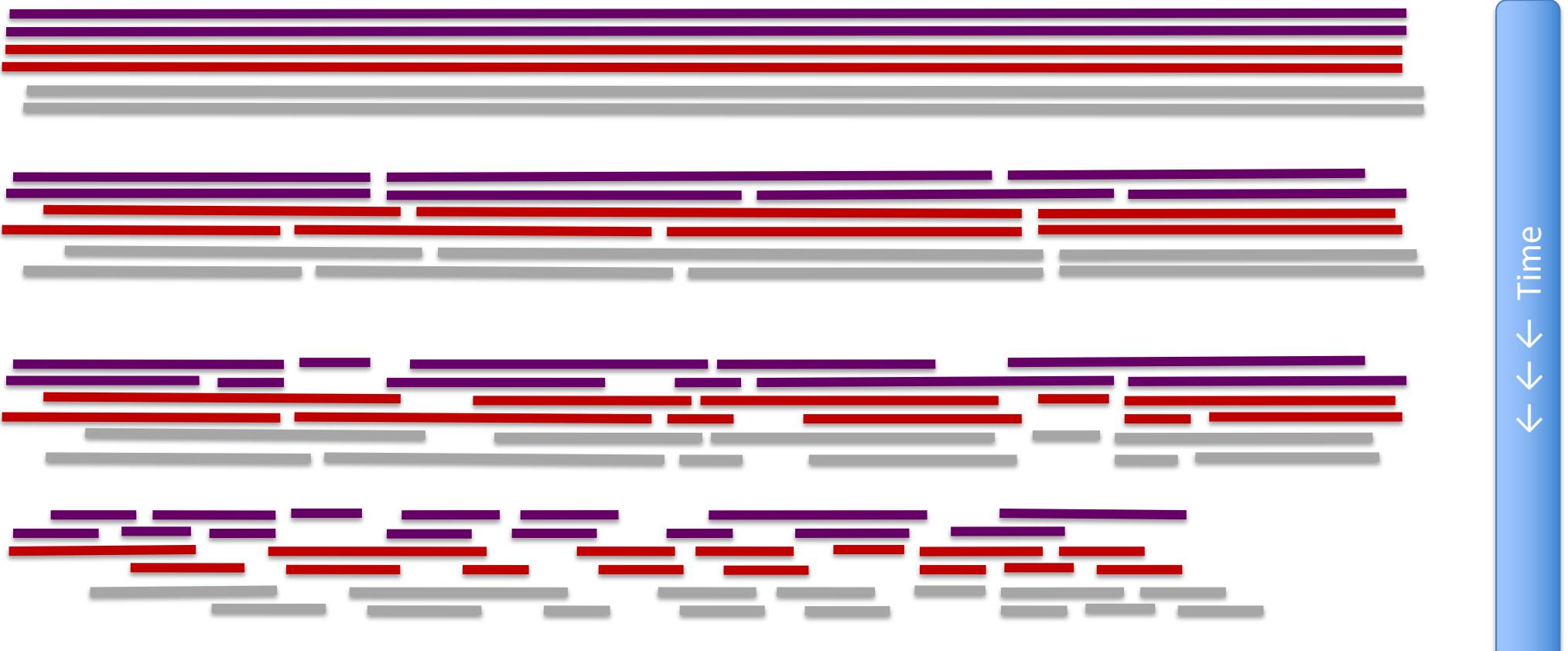
- Minimize contact at sample collection
 - Gloves,
 - Zip-lock bags,
 - (Face mask)
 - Minimal cleaning of specimens,
 - Refrigerated storage (-20°C).
- Apply aDNA lab procedures
- Track contamination
 - mtDNA
 - Y-chromosome
 - X-chromosome



Ancient DNA lab at Centre for GeoGenetics

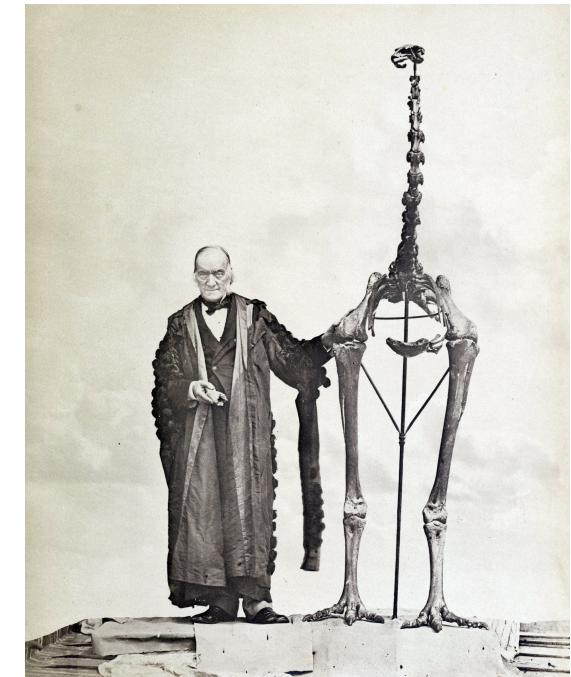
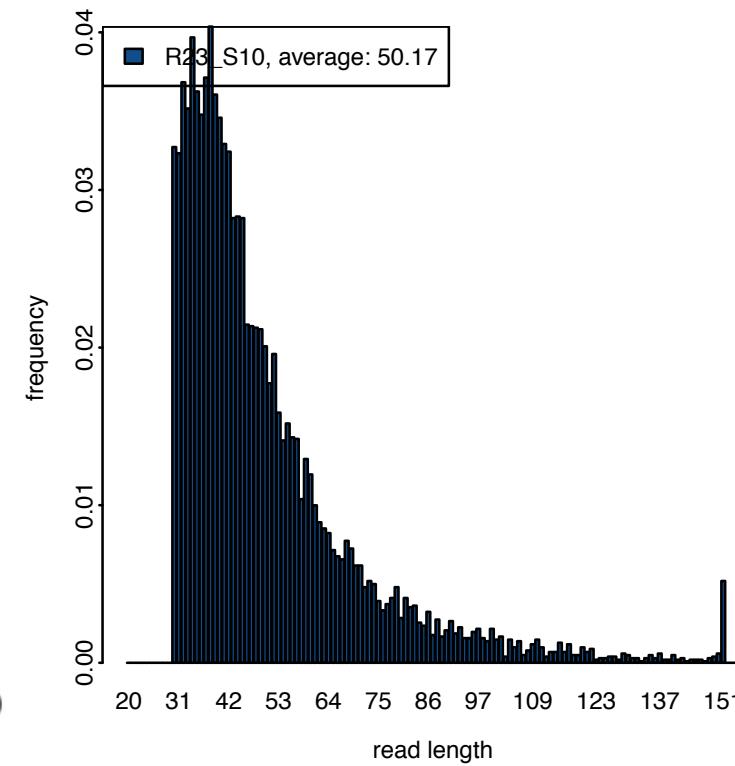
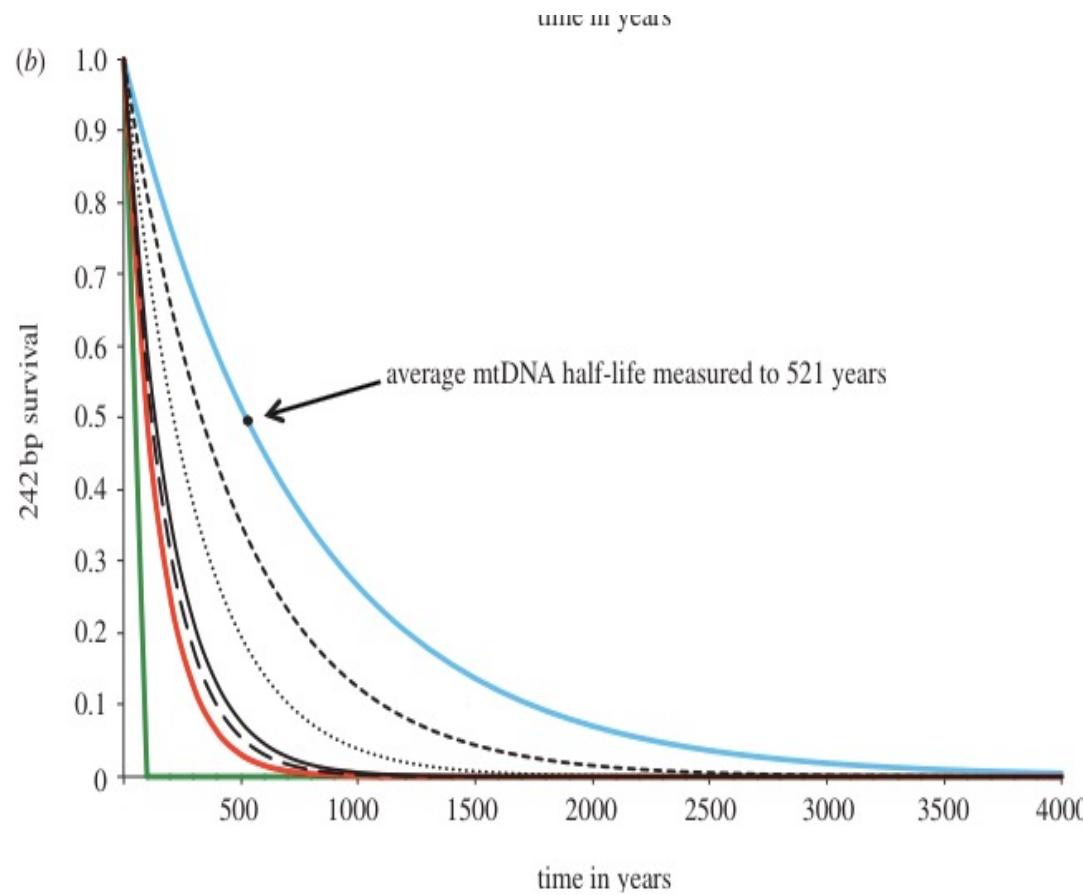


DNA breaks down post mortem



at leas 20bp, otherwise the error just like a random mapping

DNA breaks down post mortem



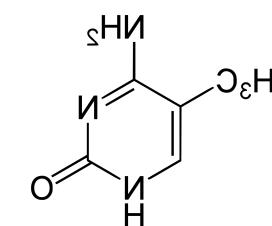
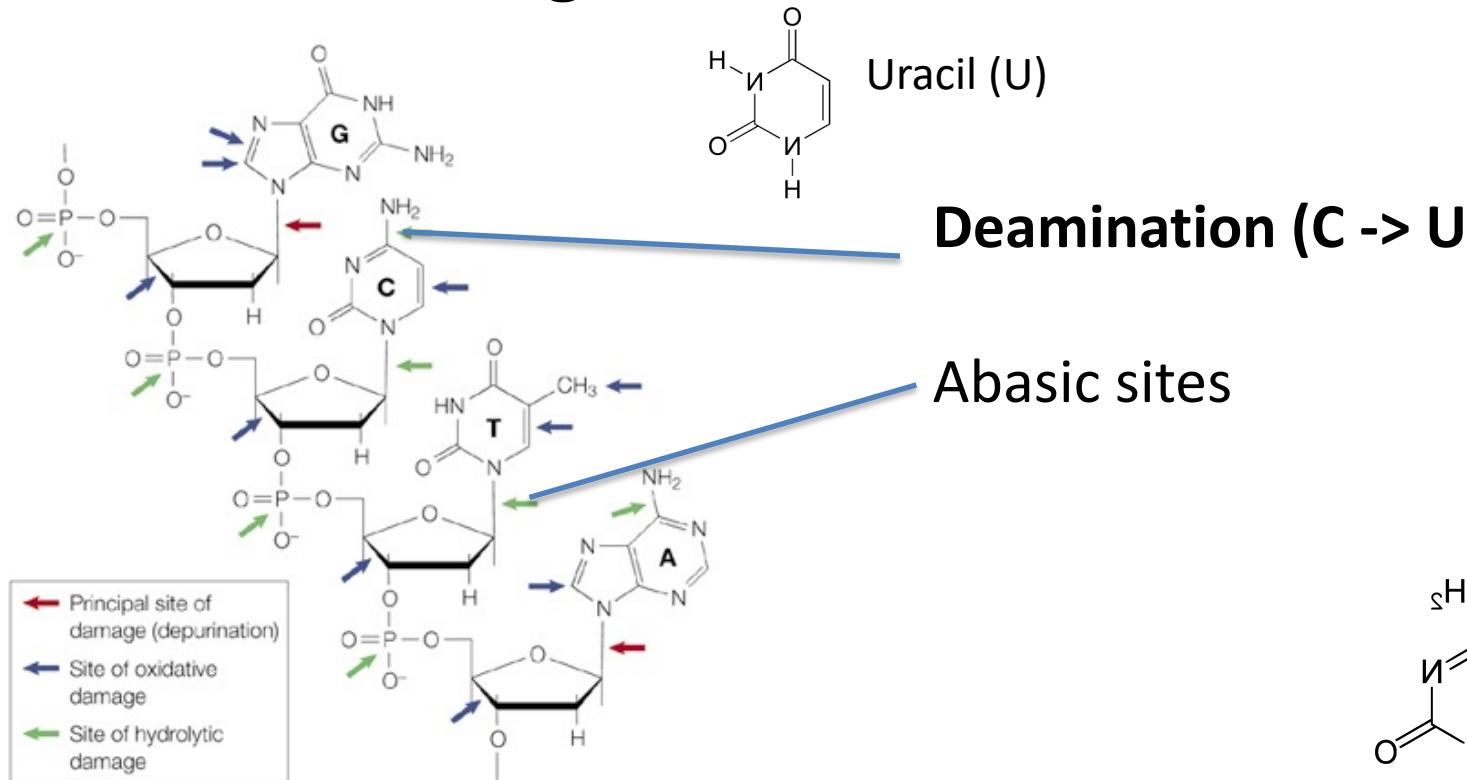
Allentoft 2012

Additional reading: Kistler NAR 2017

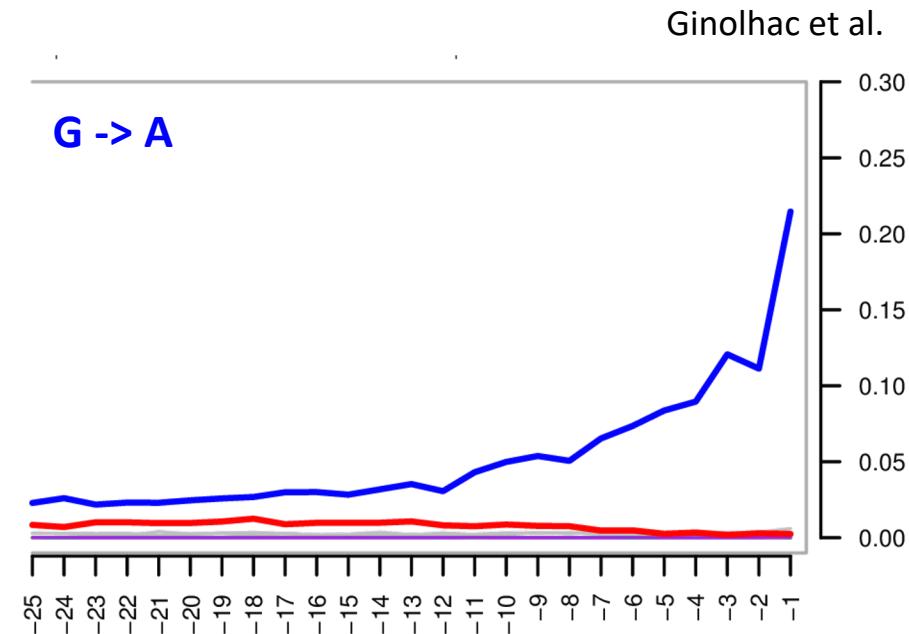
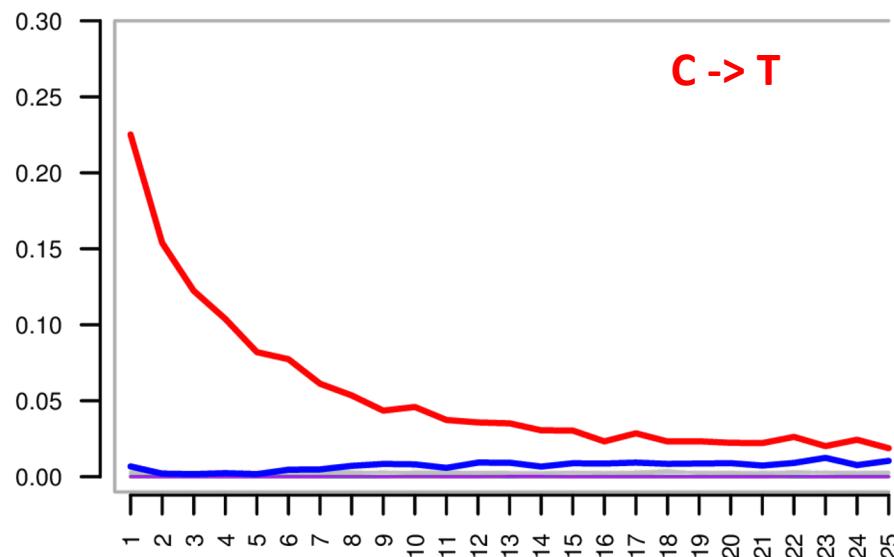
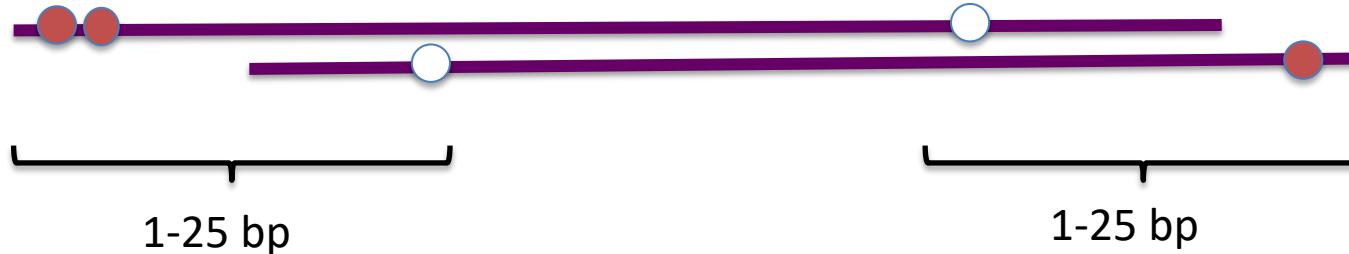
DNA decay and damage

DNA damage results in

- Short DNA lengths (break of backbone)
- Post mortem change or loss of bases.

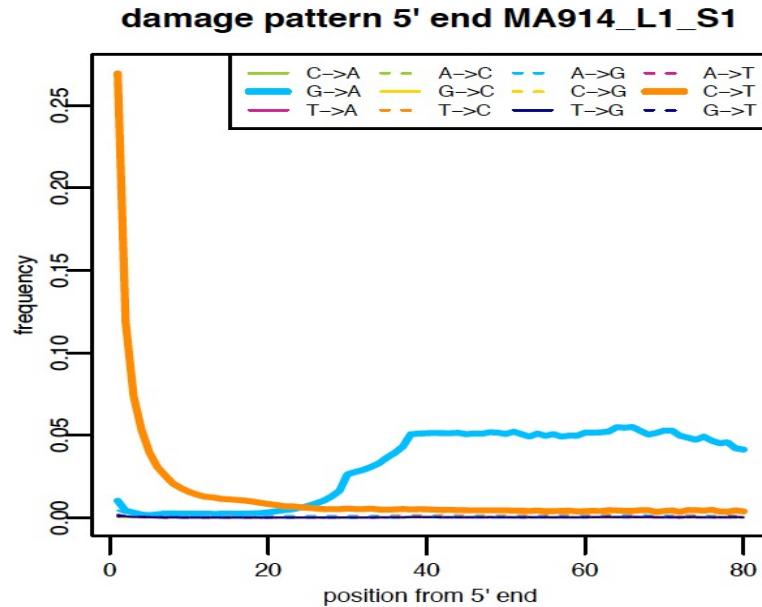
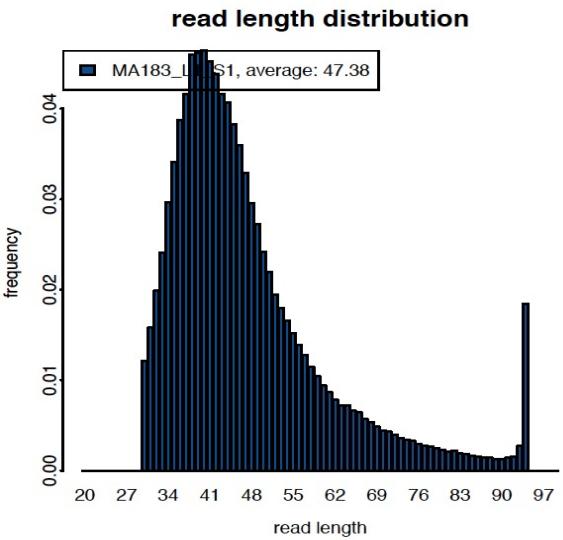


aDNA is damaged

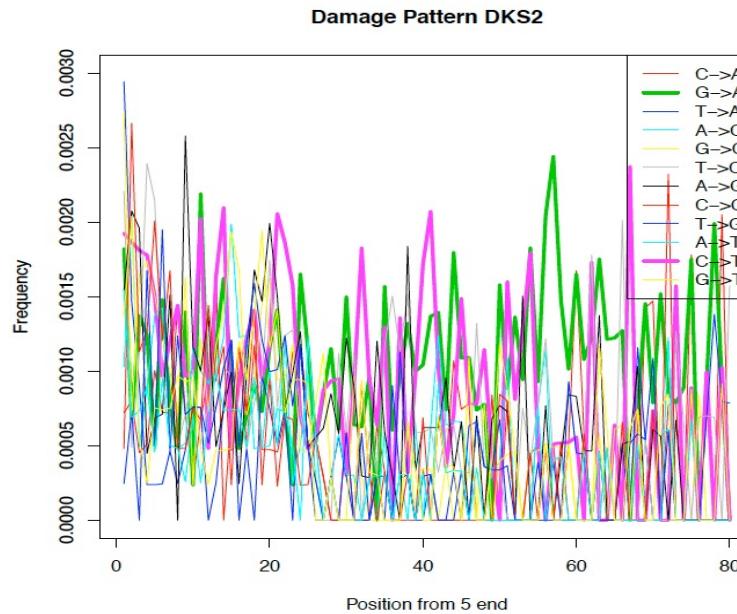
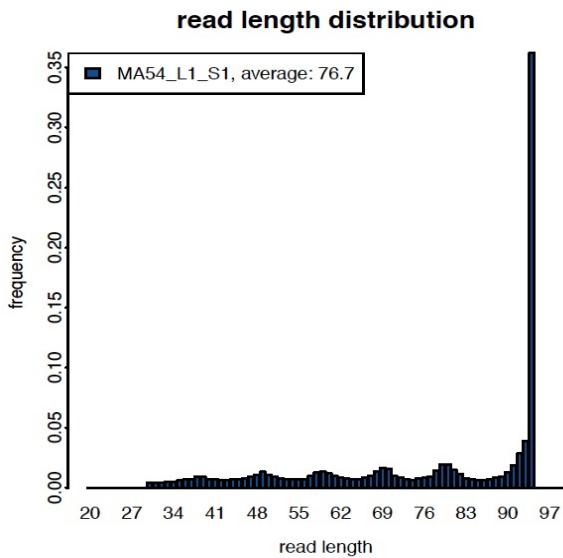


Which sample is ancient?

1



2

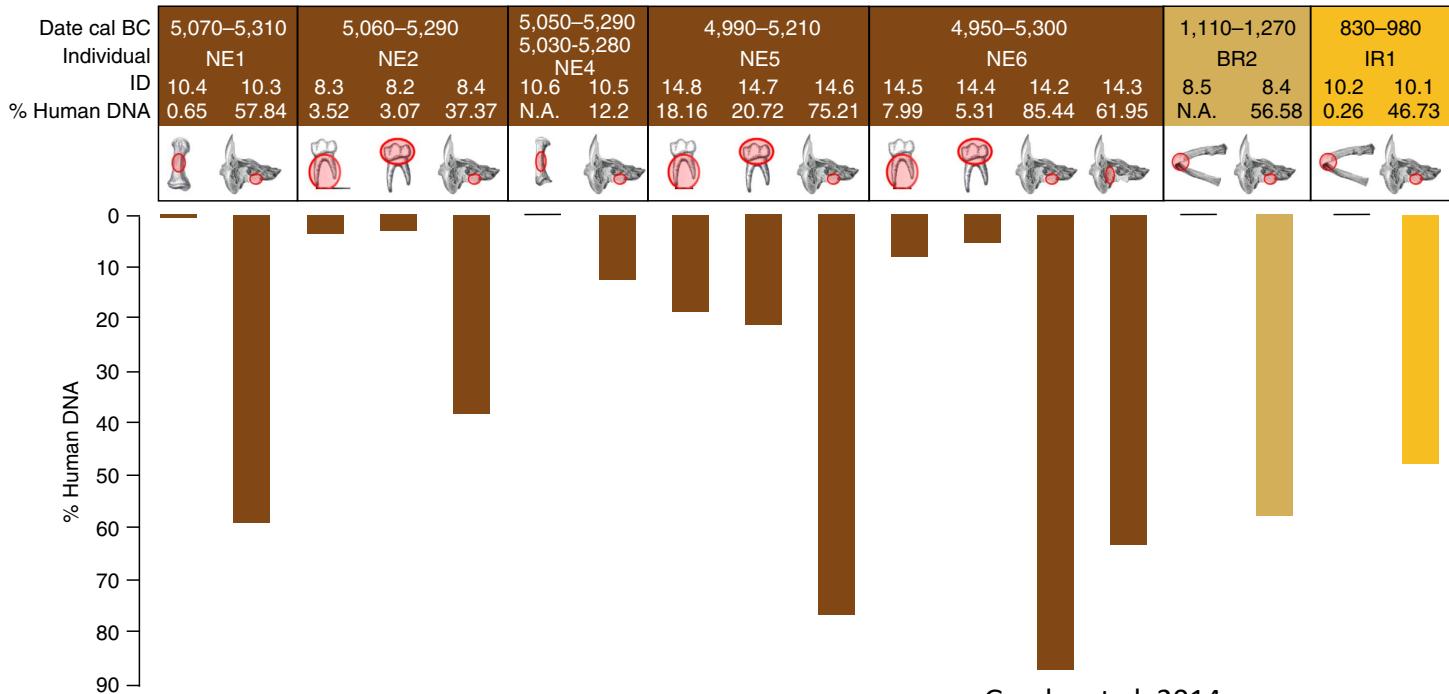
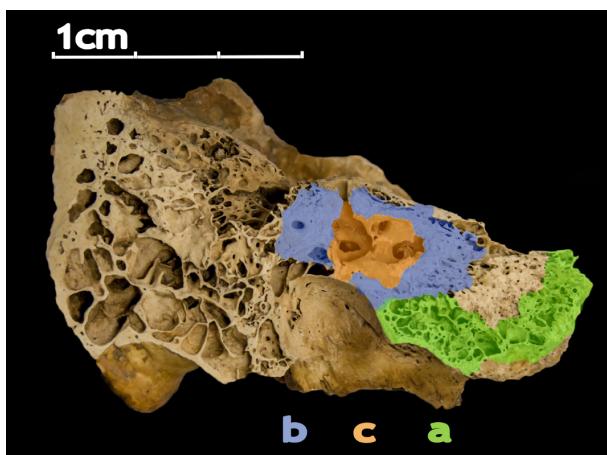
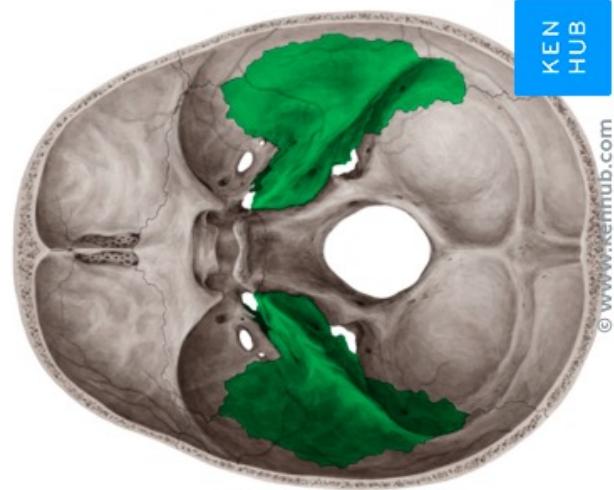


Degradation of DNA

1. What factors affects the rate of decay? **Exposure to sunlight, temperature, microbes**
2. How can the damage be useful? **Validate whether it is ancient or not e.g. C->T damage, or G->A damage.**
3. In what locations/regions/environments would you expect good samples? **Dry and cold**
4. Where would you predict problems?
5. How can we analyse the data in presence of damage?

Text

Bone types matter



Gamba et al. 2014

Figure 1 | Petrous bones versus non-petrous bones. Percentage of non-clonal endogenous DNA recovered after shotgun sequencing. The sampled bone tooth portion is circled in red. N.A. indicates that the library did not pass quality assessment for sequencing.

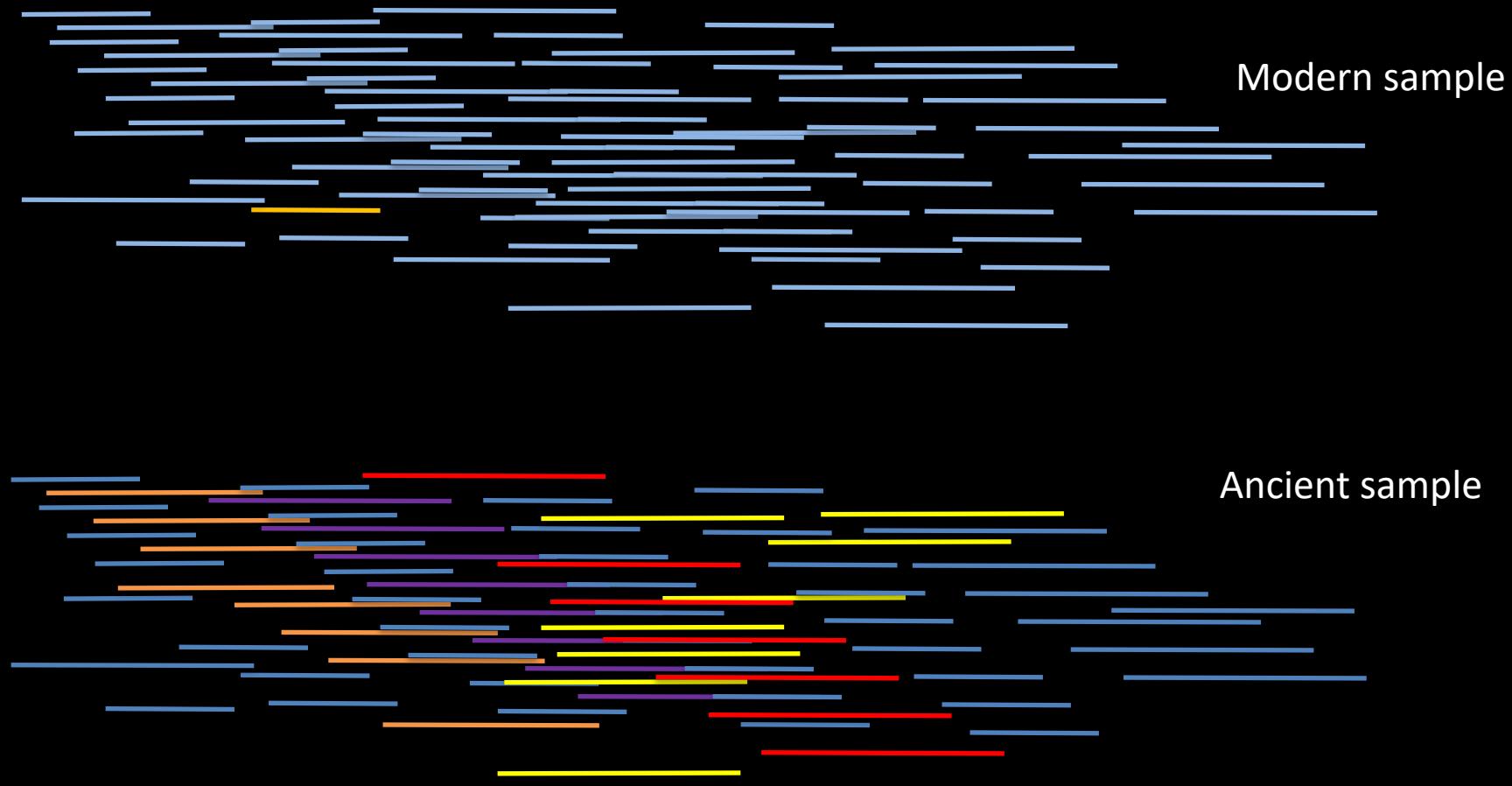


Libraries and complexity

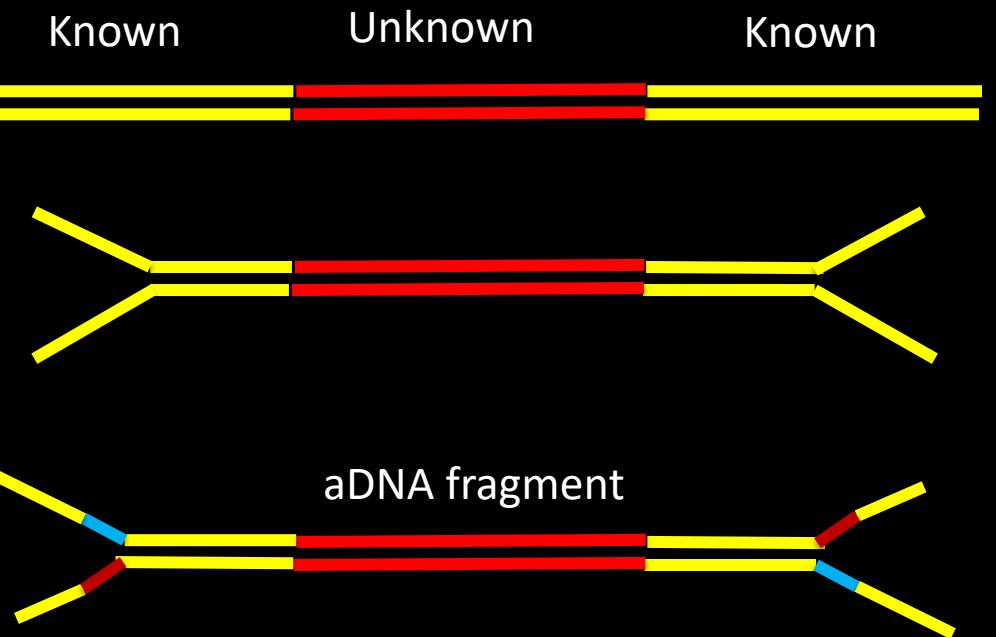


- **Molecular biology library:** A collection of DNA fragments that may be stored... propagated .. . representing a total swarm of... genome/ mutants/cDNA/.
- **Sequencing library:** A collection of DNA fragments prepared for sequencing, representing (or containing all?) the DNA content of the sample.

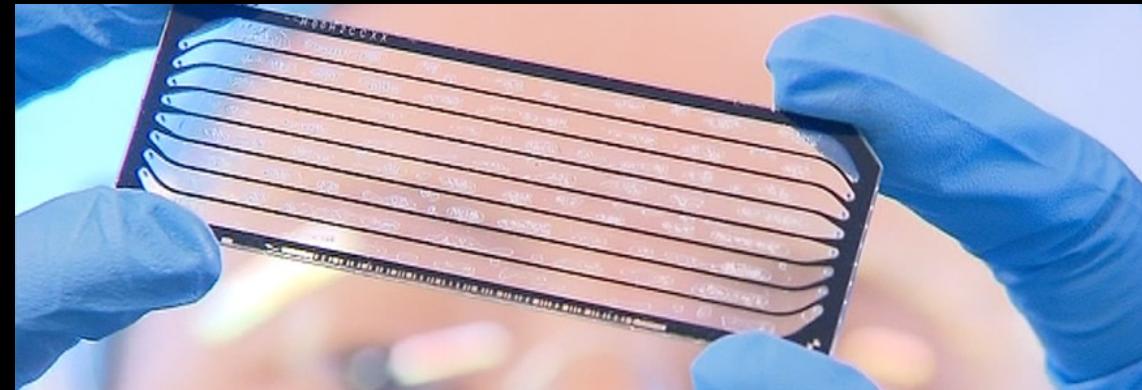
Sequencing library



Library fragment between adaptors

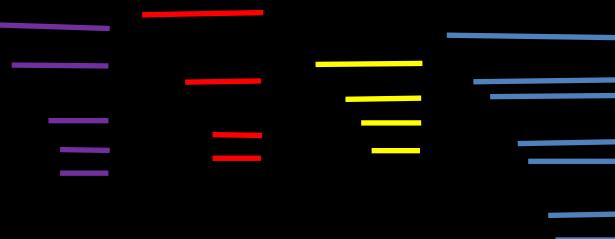
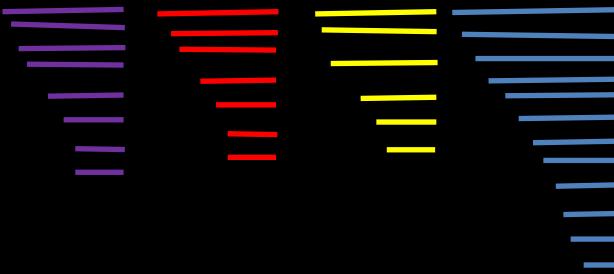
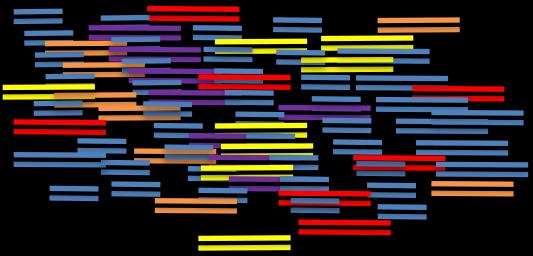
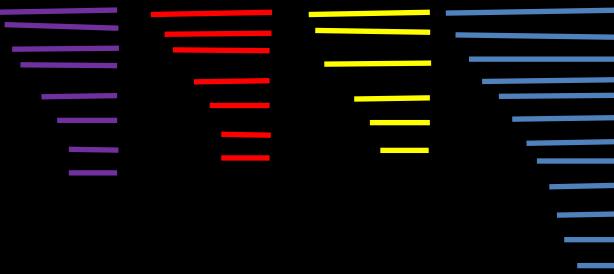


Sequencing adapters
Incl. index



Sequencing library complexity vs clonality

Number of unique fragments in library



More library
Same complexity

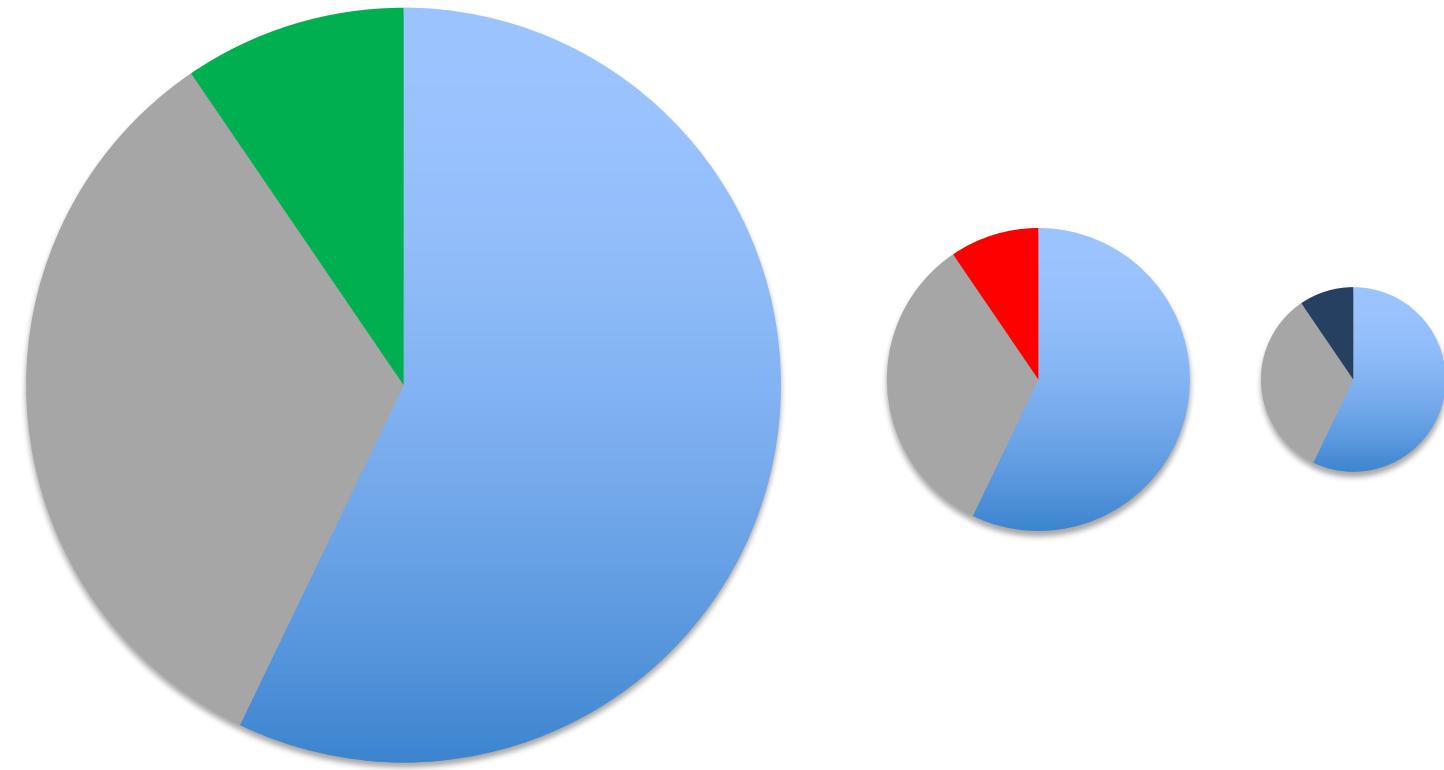
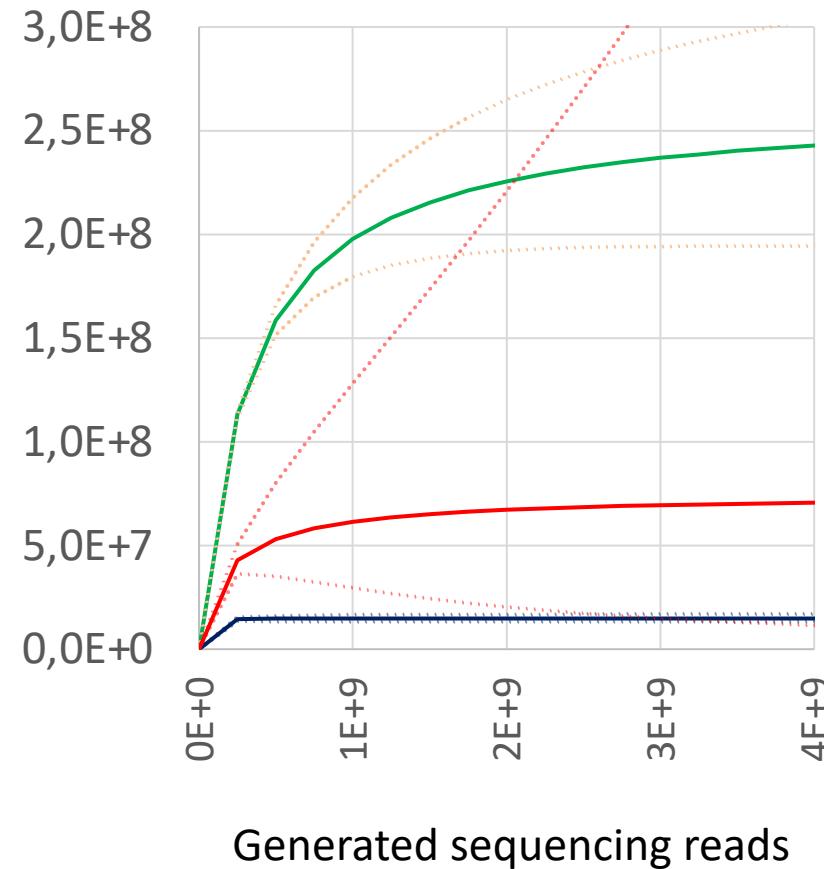
More library,
lower complexity
Higher clonality

DNA complexity

How much data can be squeezed out of the sample?



NEW sequencing reads

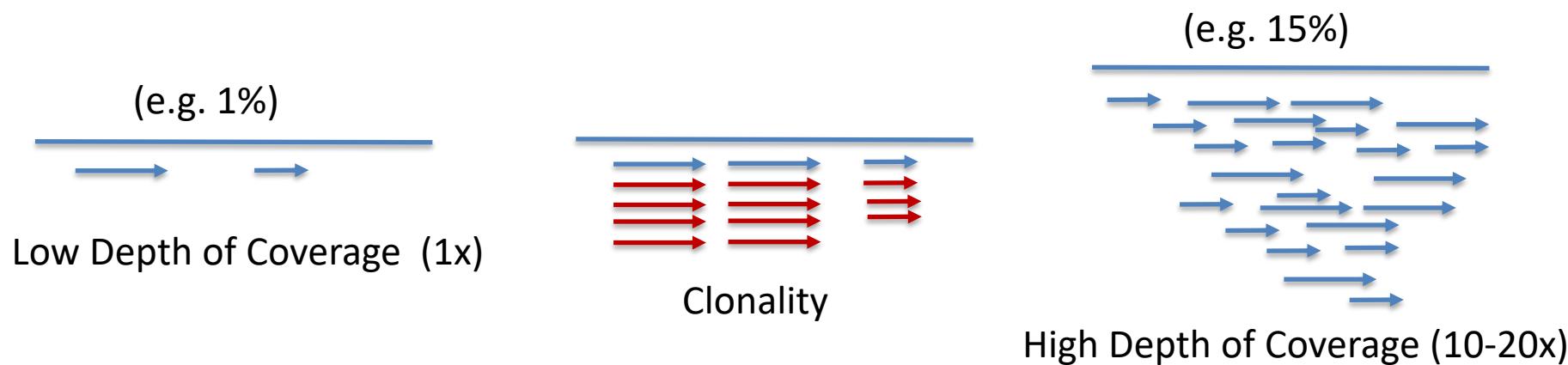


Illumina aDNA sequencing



agtccctggatagtcttagtccgatt

Reads/sample (trimmed)	Endo %	Mapping reads	Average length	Coverage (DoC \geq 1)	DoC	Possible DoC
20,000,000	1	200,000	50 bp	0.3 %	0.003x	1x
20,000,000	15	3,000,000	50 bp	4.5 %	0.05x	10-20x
20,000,000	100	20,000,000	125 bp	76 %	0.8x	>100x





Based on slides from Lasse Vinner, head of the sequencing center (on picture)

Summary of aDNA part

- Authentic aDNA is short and damaged
- Sample types and sample quality matter
- Library complexity determines the max seq output
- Risk of contamination
- Clean-lab facilities are required



Standard workflow

- Lab stuff + sequencing
- Call bases and quality scores (FastQ)
- Alignment or *denovo* assembly. Alignment is topic for tomorrow

Genotypes

- Generate genotype likelihoods. Uncertainty is topic for friday.
- Estimate allele or haplotype frequencies
- Calculate genotype probabilities
- Call genotypes

Downstream analysis

- Variation detection
- Population genetic analysis



Usual files



- Unmapped reads (`fastq`, `fasta` files with qualities).
- Mapped reads (`bam`, binary `sam` files).
- Variant information file (`.vcf` variant call format).



Low level format: Fasta files

```

1 >HN988947 |Severe acute respiratory syndrome coronavirus 2 isolate Wuhan-Hu-1
2 ATTAAGGTTTACCTTCCCAGTAAACCAACCAACTTTCGATCTCTGTAGATCT
3 GTTCTCAAACGAACTTAAATCTGTGGCTGTCACTCGGCTGCATCTAGTGACT
4 CACGAGTATAATTAAACTAATTACTGTGTTGACAGGACACGAGTAACCTCGTCTAC
5 TTCTCGAGGCTGCTTACGGTTCTGGTCCGGTGGACGGCATCAGCACATCTAGGTTT
6 CGTCCCCGGTGTGACGAAAGGTAAGATGGAGGACCTTGTCCCCCTGGTTAACGAGAAAAAC
7 ACACGTCAACTCAGTTGGCTGTTCAGGGTTCGACGGCTGCTGCAGCTGGCTTTGG
8 AGACTCGGTGGAGGAGGCTTATCAGAGGCACGTCAACATCTAAAGATGGACTTGTGG
9 CTTAGTAGAAGTTGAAAAGGGCGTTTGGCTCAACTCTGAACGCCCTATGTGTTCATCAA
10 ACAGTCGGATGCTGAACTGCACCTCATGGTCATGTTATGGTGGCTGAGCAGACT
11 CGAAGGCAATTCTAGTACGGTCTGGTAGTGGTGGAGACACTTGGTGGCTGGTCCATGTTGG
12 CGAAAGGCAATTCTAGTACGGTCTGGTAGTGGTGGAGACACTTGGTGGCTGGTCCATGTTGG
13 TGGCCATAGTTACGGCGCCGACTTCTGGCTTAAGAACGGTAAATAAAGGAGCTGG
14 TCCCTATGAAGATTTCAAGAAAACGGTGGACACTAAACATAGCTGGCTGGTGGTGGCTGG
15 ACTCATCGTGTGGCTACCTCTTGTGGCTGGCTGGACGGGGCATACCTCGCTATGGCTGATAACAACTCTGG
16 CCCCTGATGGCTACCTCTTGTGGCTGGCTGGCTGGACGGGGTGTGGTGGTGGTGGCTGG
17 ATGCACTTGTGGCTGGCAACACTGGGACTTATGGACATAAGAGGGGTGTGACTCTGG
18 TGAACATGAGCATGAAATTCTGGTGGACACGGAACTTGTGAAAGAGGCTATGAAATTGCA
19 GACACATTGAAATTAAATTGGCAAAGAAATTGGACACCTTCATGGGGAAATGCTAAA
20 TTTGGTATTTCTTAAATTCTTCTTAACTCAGACTATTCAACCAAGGGTTGAAAGAAAAAA
21 CCTGGATGGCTTTATGGTGGAAATTGACTGTCTCATCTGGCTCATCTGGCTGAAATG
22 CAACCAAATGTGCTTCAACTCTCATGGGTGATCATTTGGTGGTGAACCTCATGGCA
23 GACGGGGATTTGTGAAAGCACTTGTGAAATTGGTGGACTGTGAAATTGGTGGTGG
24 AGGGTGGCAACTACTGTGGTACTTACCCCCAAAATGCTGTGTTAAAATTGGTGGCAGC
25 ATGTCACATTCAAGAGTGGGGCTGGACATAGTCTGGCGAAATACCATATGAATCTGG
26 CTTGAACAAACCATCTGGTGGGGTGGCTGCACTATTGCTTGGAGGGCTGTGTTCTC
27 TTATGTTGGTGGCCATAACAACTGTGGCTATTGGGTGCTTCAAGCTGCTAGGGCTAACATAGG
28 TTGTAACCATACAGGTGTTGGAGAAGGTTGGAGGGTCAAGGCTTAAATGACAACTCTTGTG
29 AATACTCCAAAAGAGAAGTCACATCAATATTGGTGGTGACTTAAACTTAATGAGA
30 GATCGCCCATTTTGGCATTTCTGCTTCCACAAGTGTGCTTTGGAAACTGTGAA

```

- Text file representation of a nucleotide sequence
- Consists of an identifier starting with > followed by actual sequence which might span multiple lines
- A single fasta file might contain multiple entries. E.g chr1 chr2 etc.

Figure: First sequence of the SARS-CoV-2 virus



Low level format: Fastq files

FASTQ

```
@FC42BF1AAXX:6:1:5:732#0/1      <-- read ID
TGATTCTCTCGATATCCAGTCCTTAGTGNCATAGN <-- read (bases)
+
a^_aaaa'aa'_aaa_aaa'_--'`-'VBBBBBBBBB <-- quality score
@FC42BF1AAXX:6:1:5:492#0/1
AACAGTGGGAGGCTGCAGCAGGAGGATTNCTGAAN
+
ababb_abbbZbabaab^`aaTaabbaBBBBBBBBB
@FC42BF1AAXX:6:1:5:480#0/1
ACCTCCTCAGAGTTCTCGAGCTCGAGAANTCTGGN
```



Quality scores/Phred scores

Qscore

- The ASCII values can be interpreted as a probability
- A Q20 (ASCII '5') score is probability of 1%
- The score is the probability, ϵ , that the base is incorrect
-

$$Q_{score} = -10 \log_{10}(\epsilon)$$

-

$$\epsilon = 10^{\frac{-Q}{10}}$$

offset of 33 (sometimes 64)



Second round of computer exercises

These exercises revolve around being (still) able to login to the ricco server. And playing around with these lowlevel formats previously described. How long?

https://github.com/ANGSD/adv_binf_2023_week1



Alignment

Commonly used aligners

- Bowtie
- BWA
- MAQ
- SOAP2
- Mosaik
- Eland
- DRAGEN



Alignment

Commonly used aligners

- Bowtie
- BWA
- SOAP2
- DRAGEN

Burrows-Wheeler Alignment: are extremely fast. but not commonly used ones are not quality sensitive!



Alignment file

Alignments of the reads to a reference genome

an alignment file includes

reads TTTGTTCTTCTTTCTCTCTAGTCTTCTT ...

Qscore NVFVN] ^] ' ^ _] ^ ^ U]] '] [_ vs [_ ^ Z] _ ...

start position chr4 53351385

multiple best hits 1

Number of mismatch 2

sequence strand -

mapping quality V

SAM/BAM format 1/3

```

@HD      VN:1.0  SO:coordinate
@SQ      SN:chr20      LN:64444167
@PG      ID:TopHat    VN:2.0.14      CL:/srv/dna_tools/tophat/tophat -N 3 --read-edit-dist 5
        lign-edit-dist 2 -i 50 -I 5000 --max-coverage-intron 5000 -M -o out /data/user446/mapping_tophat
        20 /data/user446/mapping_tophat/L6_18_GTGAAA_L007_R1_001.fastq
HWI-ST1145:74:C101DACXX:7:1102:4284:73714    16    chr20    190930  3      100M   *
        CCGTGTAAAGGTGGATCGGGTCACCTCCAGCTAGGCTTAGGGATTCTTAGTTGGCCTAGGAAATCCAGCTAGTCTGTCTCTCA
C      BBDCDDCCDDDDCDDDDDDCDDCCDBC?DDDDDDDDDDDDDDDDCCDCDDDDDDDDDDCCCEDDC?DDDDDDDDDDDDDDDDDDDDDE
        AS:i:-15     XM:i:3  X0:i:0  XG:i:0  MD:Z:55C20C13A9 NM:i:3  NH:i:2  CC:Z:= CP:i:553527
HWI-ST1145:74:C101DACXX:7:1114:2759:41961    16    chr20    193953  50     100M   *
        TGCTGGATCATCTGGTTAGTGGCTCTGACTCAGAGGACCTCGTCCCCCTGGGCAGTGGACCTTCAGTGATTCCCTGACATAAGC
G      DCDDDDDEDDDDDDCDDDDDDCDDCCDDCDDDDDEEC>DFFEJJJJJIGJJJIHGBHHGJIJJJJJJJJJJJJJJJJJJJJJJJJJJH
        AS:i:-16     XM:i:3  X0:i:0  XG:i:0  MD:Z:60G16T18T3 NM:i:3  NH:i:1
HWI-ST1145:74:C101DACXX:7:1204:14760:4030    16    chr20    270877  50     100M   *
        GGCTTATTGGTAAAAAAGGAATAGCAGATTAATCAGAAATTCCACCTGGGCCAGCAGCACCAACCAGAAAGAAGGGAAAGAACACA
C      DDDDDDDDDCDDDDDDDDDDDEEEEEEEFFFEFFEGHHHFGDJJIHJJIIJJIIIIIGGFJJIHIIIIJJJJJJIGHHFAGFHJHF
        AS:i:-11     XM:i:2  X0:i:0  XG:i:0  MD:Z:0A85G13  NM:i:2  NH:i:1
HWI-ST1145:74:C101DACXX:7:1210:11167:8699    0      chr20    271218  50     50M4700N50M
        0      GTGGCTTCCACAGGAATGTTGAGGATGACATCCATGTCGGGTGCACTTGGGTCTCCGAAGCAGAACATCCTCAA1
accepted_hits.sam

```

Figure: Example of a small sam file. Header part starts with @. The actual alignment information for the reads follows after the header. One read correspond to one line



SAM/BAM format 2/3

Tag	Description
QHD	The header line. The first line if present.
VN*	Format version. Accepted format: /{[0-9]+\.[0-9]+\$/}.
SO	Sorting order of alignments. Valid values: unknown (default), unsorted, coordinate. For coordinate sort, the major sort key is the RNAME field, by the order of QSQ lines in the header. The minor sort key is the POS field with equal RNAME and POS, order is arbitrary. All alignments with '*' in QNAME alignments with some other value but otherwise are in arbitrary order.
GT	Grouping of alignments, indicating that similar alignment records are group file is not necessarily sorted overall. Valid values: none (default), query (align by QNAME), and reference (alignments are grouped by RNAME/POS).
GSQ	Reference sequence dictionary. The order of GSQ lines defines the alignment
SN*	Reference sequence name. Each GSQ line must have a unique SN tag. The value in the alignment records in RNAME and RNEXT fields. Regular expression: [A-Z]{1,255}
LN*	Reference sequence length. Range: [1,2 ³¹ -1]
AS	Genome assembly identifier.
MS	MD5 checksum of the sequence in the uppercase, excluding spaces but including punctuation.
SP	Species.
UR	URI of the sequence. This value may start with one of the standard protocols. If it does not start with one of these protocols, it is assumed to be a file-system path.
RG	Read group. Unordered multiple RG lines are allowed.
ID*	Read group identifier. Each RG line must have a unique ID. The value of ID tags of alignment records. Must be unique among all read groups in header. IDs may be modified when merging SAM files in order to handle collisions.
CN	Name of sequencing center producing the read.
DS	Description.
DT	Date the run was produced (ISO8601 date or date/time).
FO	Flow order. The array of nucleotide bases that correspond to the nucleic acid flow of each read. Multi-base flows are encoded in IUPAC format, and non-various other characters. Format: /*\ [ACGTRSVWYHKDBN]+\ /
KS	The array of nucleotide bases that correspond to the key sequence of each read.
LB	Library.
PG	Programs used for processing the read group.
PI	Predicted median insert size.
PL	Platform/technology used to produce the reads. Valid values: CAPILLARY, SOLID, HELICOS, IONTORRENT, ONT, and PACBIO.
PM	Platform model. Free-form text providing further details of the platform/technology.
PU	Platform unit (e.g. flowcell-barcode.lane for Illumina or slide for SOLiD). Unit of the platform.
SM	Sample. Use pool name where a pool is being sequenced.
QPG	Program.
ID*	Program record identifier. Each QPG line must have a unique ID. The value of PG tag and PP tags of other QPG lines. PG IDs may be modified when merging SAM files in order to handle collisions.

SAM/BAM format 3/3

SAM format

QNAME ID of the read

FLAGS flag of the alignment

RNAME Reference name

POS 1-indexed position of leftmost first matching position
in sequence

MAPQ mapping quality (phred scaled)

CIGAR Alignment info

RNEXT mate reference name (paired end)

RPOS mate position (paired end)

ISIZE inferred insert size

SEQ the sequence

QUAL qscore of the sequence



Mapped reads

My definitions (The literature is not consistent)

Depth The number of reads that maps to a position

Counts The number of different alleles mapped to a position

Coverage The fraction of the genome (region) with data

Q20



Distribution of depths

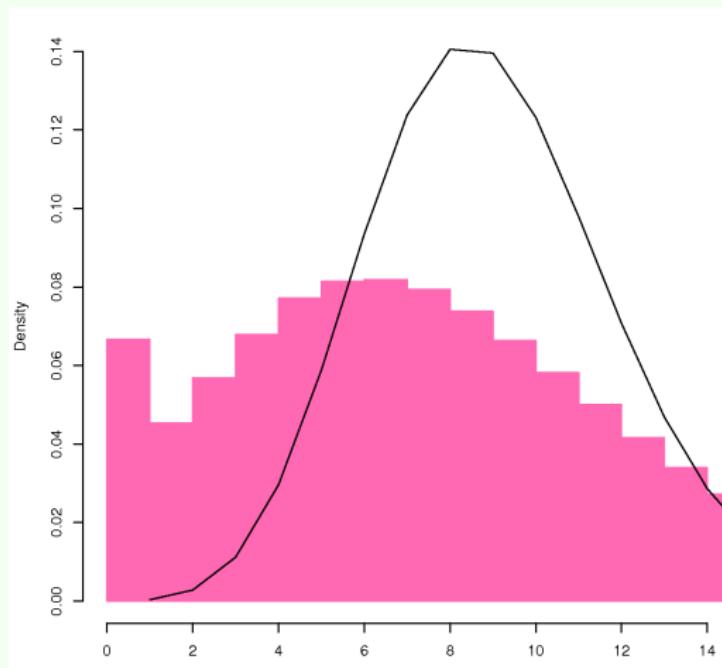


Figure: black line:expected. pink: observed

- NGS data with mean depth 8
- If the reads are random then we expect the depth to follow a Poisson distribution
- you need a very large depth to cover all sites



Definitions

- * An *allele* is a variant form of a gene.
- * A *genotype* are the alleles for a specific individual.
- * *genotype calling* is the process of determining the underlying 'true' genotype that has generated the data we observe.

SNP-chip and high throughput sequencing

SNP-chip Genotypes are readily available for a predetermined list of positions.

High-throughput sequencing Data are short sequences that are sampled from either chromosome with replacement.

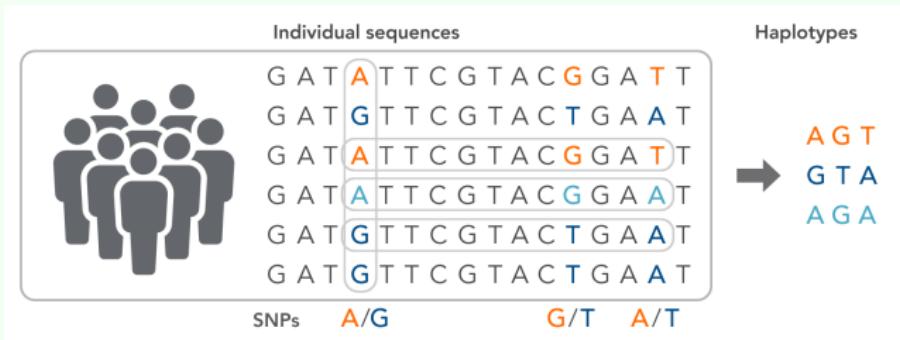
AGCCACATCAGGCCATTGCTGCAGCAGCACGGTCACCAGACAGAAATCTTGC
CAGGCCACACCCAGCCAATTGCTGCAGCAGCACGGTCACCAGACAGAAATCTTGC
CAGGCCACACACAGGCCATTGCTGCAGCAGCACGGTCACCAGACAGAAATCTTGC
TGACAGGCCACATCAGGCCATTGCTGCAGCAGCACGGTCACCAGACAGAAATCTTGC
CTGACAGGCCACATCAGGCCATTGCTGCAGCAGCACGGTCACCAGACAGAAATCTTGC
GTCTGACAGGCCACATCAGGCCATTGCTGCAGCAGCACGGTCAAC
TGCCAGTCTGACAGGCCACATCAGGCCATTGCTGCAGCAGCACGGTCACCAGACGGAAATCTTCT
CATTGCCAGTCTGACAGGCCACATCAGGCCATTGCTGCAGCAGCACGGTCACCAGACAGAAAT
ACCCATTGCCAGTCTGACAGGCCACATCAGGCCATTGCTGCAGCAGCACGGTCACCAGACAGAA
AGAGATGAAAACCCATTGCCAGTCTGACAGGCCACATCAGGCCATTGCTGCAGCAGCACGGTC
AGACCAGAGATGAAAACCCATTGCCAGTCTGACAGGCCACATCAGGCCATTGCTGCAGCAGCA



Question

Questions?

- Identify some of the alleles from the figure?
- Identify some of the genotypes.
- Calculate allele frequencies for the variable positions.



○○



VCF format

(a) VCF example

```

Header { ##fileformat=VCFv4.1
##fileDate=20110413
##source=VCFtools
##reference=file:///refs/human_NCBI36.fasta
##contig=<ID=L,length=249250621,md5=1b22b98cdeba49304cb5d48026a85128,species="Homo Sapiens">
##contig=<ID=X,length=155270560,md5=7e0e2e580297b7764e31dbc80c2540dd,species="Homo Sapiens">
##INFO<=ID=AA,Number=1,Type=String,Description="Ancestral Allele"
##INFO<=ID=H2,Number=0,Type=Flag,Description="HapMap2 membership"
##FORMAT<=ID=GT,Number=1,Type=String,Description="Genotype"
##FORMAT<=ID=GQ,Number=1,Type=Integer,Description="Genotype Quality"
##FORMAT<=ID=DP,Number=1,Type=Integer,Description="Read Depth"
##ALT<=ID=DEL,Description="Deletion"
##INFO<=ID=SVTYPE,Number=1,Type=String,Description="Type of structural variant"
##INFO<=ID=END,Number=1,Type=Integer,Description="End position of the variant"
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMPLE1 SAMPLE2
1 1 . ACG A,AT 40 PASS .
1 2 . C T,CT . PASS H2;AA=T GT 0|1 2/2
1 5 rs12 A G 67 PASS .
X 100 . <DEL> . PASS SVTYPE=DEL:END=299 GT:DP:G 1|1:.- 0|0:29:30
}

```

(b) SNP

```
(b) SNV
Alignment      VCF representation
1234          POS REF ALT
ACGT          2   C   T
ATGT
```

(c) Insertion

12345	POS	REF	ALT
AC-GT	2	C	CT
ACTGT			

(d) Deletion

1234 POS REF ALT
ACGT 1 ACG A
A--T

(e) Replacement

1234 POS REF ALT
ACGT 1 ACG AT
A-TT

(f) Large structural variant

Alignment	100	110	120	290	300
ACGTACGTACGTACGTACGTACGT[...]	ACGTACGTACGTACGTACGT				
ACGT[...]	[...]				GTAC

```
VCF representation  
POS REF ALT      INFO  
100  T    <DEL>  SVTYPE=DEL;END=299
```

(q) Resolving ambiguity

Alignment	Possible representation			Possible representation			Recommended VCF representation		
	POS	REF	ALT	POS	REF	ALT	POS	REF	ALT
1234567890									
TTTCCCTCTA	1	TTTCCCTCT	CTTACCTA	1	T	C	1	T	C
CTTACCT--A				4	C	A	4	C	A
^ ^ ^				7	TCT	T	5	CCT	C



Third and last round of computer exercises

These exercises revolve around alignment files and variant information files previously described. How long?

https://github.com/ANGSD/adv_binf_2023_week1