# An Introduction to the Analysis of Ancient Sequencing Data

## Pre-Processing the Unmapped Reads of the Altai Neanderthal Sequencing
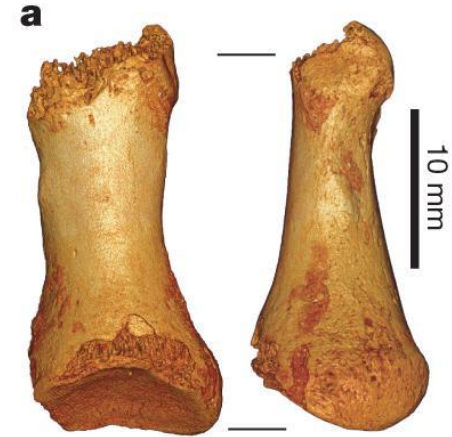
Maximilian Arlt
Master Student
Chair of RNA Bioinformatics and High-Throughput Sequencing,
Prof. Dr. Manja Marz
FSU Jena

# What is ancient DNA (aDNA) ?

-   ancient = before the fall of Roman Empire
-   extracted from 'old' remains (bones, tooth, soft tissue)



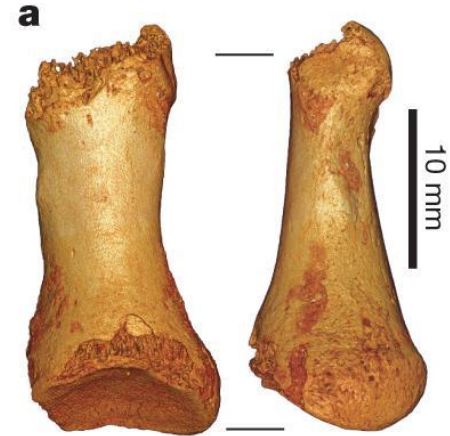adapted from  K Prüfer et al. Nature 000, 1-7 (2013)
doi:10.1038/nature12886

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# What is ancient DNA (aDNA) ?

- ancient = before the fall of Roman Empire
- extracted from 'old' remains (bones, tooth, soft tissue)

**It is degraded.**

**It is modified.**

**It is contaminated.**

# What is ancient DNA (aDNA) ?

- pedal toe phalanx, Neanderthal Individual
- ~ 50 000 years old
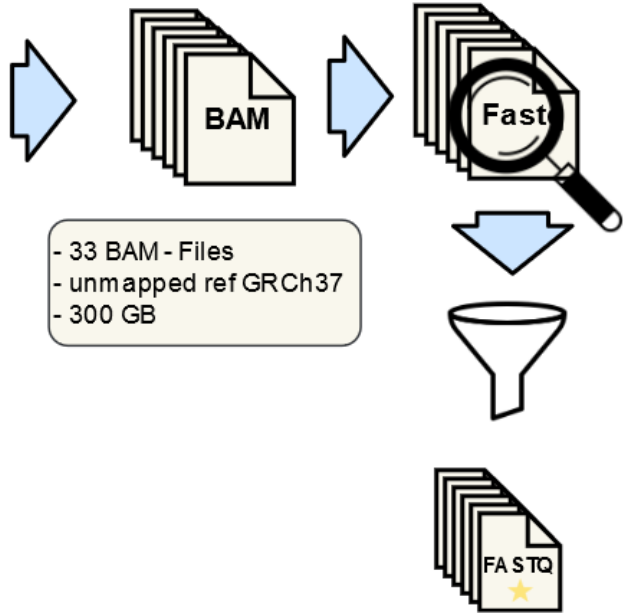- Extraction, Amplification, PE - Sequencing (Illumina)
- ~ 70 % endogenous DNA



adapted from K Prüfer et al. Nature 000, 1-7 (2013)
doi:10.1038/nature12886

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# What is ancient DNA (aDNA) ?

- pedal toe phalanx, Neanderthal Individual
- ~ 50 000 years old
- Extraction, Amplification, PE - Sequencing (Illumina)
- ~ 70 % endogenous DNA

**What about the other 30 % ?**



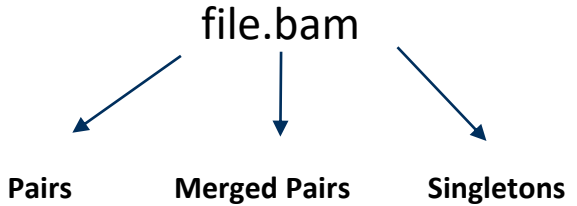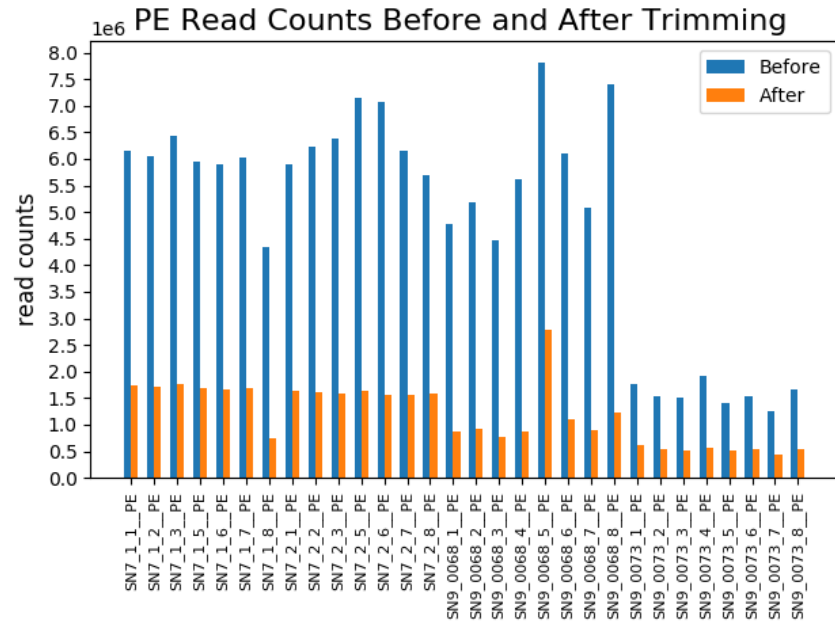adapted from  K Prüfer et al. Nature 000, 1-7 (2013)
doi:10.1038/nature12886

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

DATA SOURCE:
Prüfer, K., & Racimo, F. et al. The complete genome sequence of a Neandertal from the Altai Mountains. (2013). Nature. 505. 10.1038/nature12886.

BAM

- 33 BAM - Files
- unmapped ref GRCh37
- 300 GB

Fast

FASTQ

1. Quality assessment
- fastqc reports
- multiqc report

2. Quality enrichment
- fastp
- fastqc

file.bam

Pairs    Merged Pairs    Singletons

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

Merged Read Counts Before and After Trimming



PE Read Counts Before and After Trimming

**Before Filtering:**

| | |
|---|---|
| 1,309,313,311 | M |
| 144,357,489 | PE (2x) |
| **1,598,028,289** | **TOTAL** |

**Fastp parameters:**

Phred   >= 24
Length  >= 20 nt
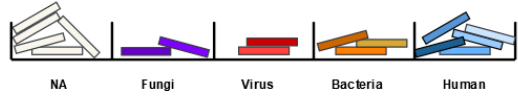Adapter Trimming
Low complexity Filtering
PolyX Trimming >= 10 nt

**After Filtering:**

| | |
|---|---|
| 1,293,551,377 | M |
| 35,876,296 | PE (2x) |
| **1,365,303,969** | **TOTAL** |

FASTQ

adapted from
Wood, Derrick E. and Steven L. Salzberg.
"Kraken: ultrafast metagenomics sequence classification using exact alignments."
Genome Biology (2013).

NA    Fungi    Virus    Bacteria    Human

**3. Kmer classification**
- kraken2
- clark

NA  +  Viral  =  Mix

**4. Binning**
- viral reads
- NA
- mixed

**4.1 Quality assess- and enrichment**
- viral reads
- NA
- mixed

**5. De-novo assembly**
- megahit
- spades

6.    6.    6.

**6. Remapping**
- HiSat2

**NCBI Viruses BLAST - Database**

**7. MEGABLAST on assembled contigs**

CSV    8. Data Assessment and Interpretation    CSV

FRIEDRICH-SCHILLER-
**UNIVERSITÄT
JENA**

Maximilian Arlt

8

# Thank you for your attention!

# References

1.  Prüfer, K., Racimo, F., Patterson, N. *et al.* The complete genome sequence of a Neanderthal from the Altai Mountains. *Nature* 505, 43–49 (2014).

2.  Prüfer, K., Racimo, F., Patterson, N. *et al.* The complete genome sequence of a Neanderthal from the Altai Mountains. *Nature* 505, 43–49 (2014). Supplementary information.

3.  Gansauge MT., Meyer M. (2019) A Method for Single-Stranded Ancient DNA Library Preparation. In: Shapiro B., Barlow A., Heintzman P., Hofreiter M., Paijmans J., Soares A. (eds) Ancient DNA. Methods in Molecular Biology, vol 1963. Humana Press, New York, NY.

# Fastqc Reports

# Quality Reports indicate artificial Sequence Content

## FastQC - Before

# Quality Reports indicate artificial Sequence Content

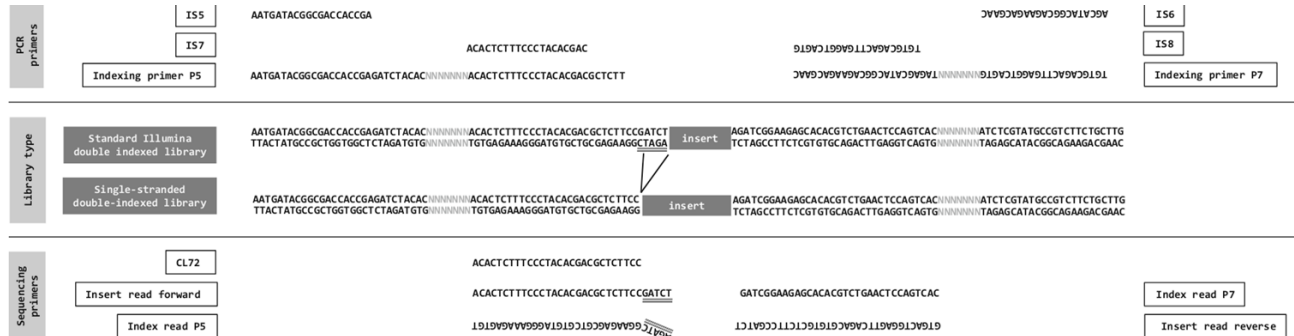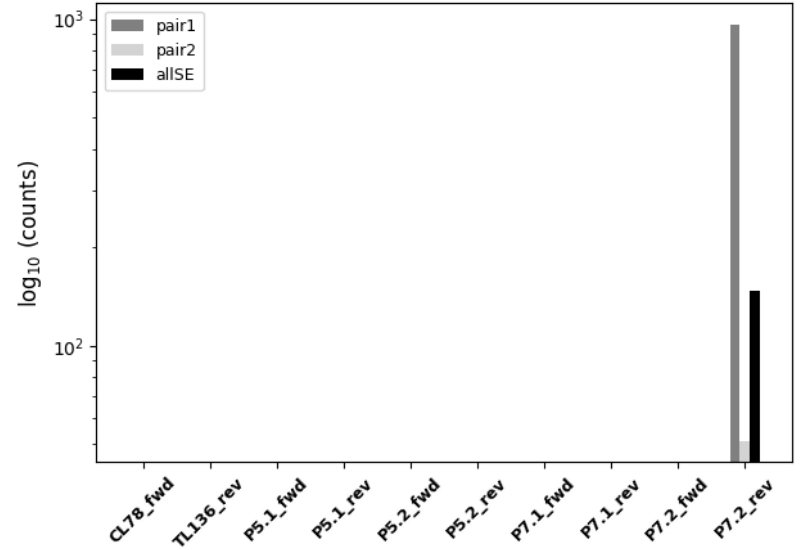## FastQC - After

# Ancient DNA is commonly:

**degraded**

**modified**

**contaminated**



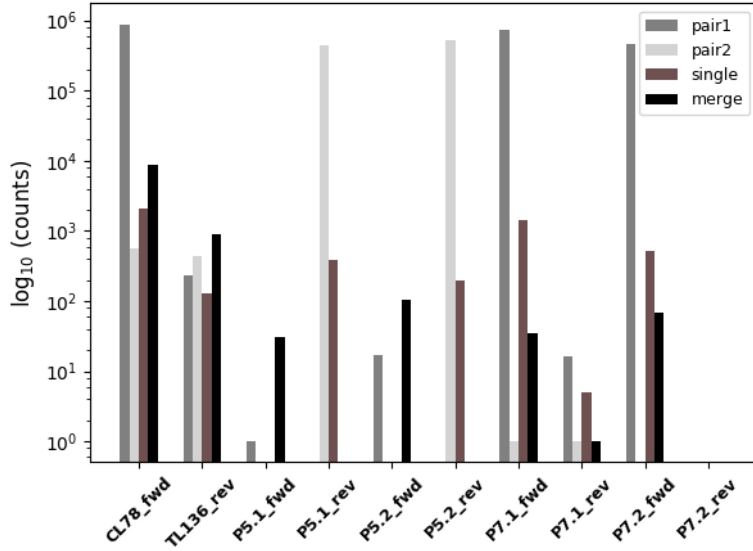adapted from: K Prüfer et al. Nature 000, 1-7 (2013); doi:10.1038/nature12886; supplements Figure S5a.3

adapted from MT. Gansauge , M. Meyer, Methods in Molecular Biology(2019), https://doi.org/10.1007/978-1-4939-9176-1_9

# Increasing resolution - coverage depth plots on MEGAHIT viral candidates