

.janno file details

Contents

1	Background	1
1.1	The Poseidon_ID	1
1.1.1	What does the Poseidon_ID represent exactly?	2
1.2	Other identifiers	2
2	Relations among samples/individuals	3
3	Spatial position	3
4	Temporal position	4
4.1	General structure	4
4.2	The columns in detail	5
5	Genetic summary data	5
5.1	Individual properties	5
5.2	Library properties	6
5.3	Data yield	7
5.4	Data quality	7
6	Context information	8

1 Background

The .janno file columns are specified in the Poseidon package specification [here](#). The following documentation includes additional background information for many of the variables. This should make it more easy to compile the necessary information for both published and unpublished data. The .pdf version of the latest version of this document is available [here](#).

1.1 The Poseidon_ID

The Poseidon_ID column assigns each entity in a Poseidon package (so one row of the .janno file) a unique identifier string.

Often the Poseidon_ID can be readily taken from the respective accompanying publication introducing a given sample. If there are multiple samples from one ancient human individual, then they may share this identifier in the publication. For the Poseidon package they have to be clearly distinguished with relevant suffixes, though,

added to the `Poseidon_ID`. `Poseidon_ID`s are also employed in the genetic data files in a Poseidon package and therefore have to adhere to certain constraints.

1.1.1 What does the `Poseidon_ID` represent exactly?

Generally, archaeogenetics operates on burial contexts, e.g. graves, with one or multiple ancient human individuals. Usually, though not always, it is possible to attribute the skeletal remains within these graves to individuals based on the archaeological context and physical-anthropological analysis. Each individual can get sampled one or multiple times, either by directly probing their preserved tissue, mostly bones, or by sampling any reagent that contains their DNA (through whatever pathway or taphonomic process). From one such sample one or multiple extracts can be derived, which can be transformed into one or multiple libraries, which may or may not be subjected to a DNA capture protocol and then sequenced one or multiple times. The raw sequencing data can undergo various different forms of computational processing and eventually genotyping to produce the data relevant for most derived analyses and thus stored in Poseidon.

While the wetlab-processes can be understood as a relatively predictable tree of separate physical and digital products for any given ancient individual, the computational data-processing finally breaks the conceptual tree-ness by allowing for arbitrary conflation of sequencing data obtained through potentially separate means: Data from different libraries can very well be merged if they are from the same individual, even if they are not from the same sample.

A `Poseidon_ID`, and therefore the identifier for the main singular entity in a Poseidon package, could approximately be described as representing one end-point in the data preparation graph laid out above. Typically this end-point corresponds to an optimal result, consciously selected for a given individual, research question and publication. Unfortunately, in reality a `Poseidon_ID` is not suited to uniquely identify exactly one such end-point. The reality in the Poseidon ecosystem is rather that slightly different end-points can have the same `Poseidon_ID`, e.g. across package versions or public Poseidon archives. A single endpoint can only be uniquely identified from a combination of `Poseidon_ID`, Poseidon package and package version.

1.2 Other identifiers

The column `Alternative_IDs` provides a way to list other IDs used for the respective individual. These might for example be names used in different publications or popular names like “Iceman”, “Ötzi”, “Girl of the Uchter Moor”, “Tollund Man”, etc.. The `Relation_*` columns described below allow to more precisely express the relationship type “identical” among samples in a Poseidon package.

The `Collection_ID` column stores an additional, secondary identifier as it is often provided by collaboration partners (archaeologists, museums, collections) that provide the specimen for archaeogenetic research. These identifiers can have a very heterogenous structure and may not be unique across different projects or institutions. The `Collection_ID` column is therefore a free-form text field.

The `Group_Name` column contains one or multiple group or population names for each individual, separated by ;. The first entry must be identical to the one used in the genotype data for the respective sample in a Poseidon package, and whitespace is not allowed in any of the entries. Assigning group and population names is a hard problem in archeogenetics [1], so the `.janno` file allows for more than one identifier.

2 Relations among samples/individuals

To systematically document biological relationships uncovered among samples/individuals in one or multiple Poseidon datasets (e.g. with software like READ [2] or BREADR [3]), the `.janno` file can be fit with a set of columns featuring the `Relation_*` prefix. Across these columns it should be possible to encode all kinds of pairwise, biological relationships an individual might have.

`Relation_To` is a string list column (so: multiple values are possible if separated by `;`) that stores the Poseidon_IDs of other samples/individuals to which the current individual has some relationship.

`Relation_Degree` stores a formal description of the closeness of this relationship as measured purely from aDNA data. It is therefore also a list column that can hold the following values for each relationship:

- `identical`: The two samples are from the same individual or from identical twins
- `first`: The two individuals are closely related – a first degree relationship (e.g. siblings, parent-offspring)
- `second`: A second degree relationship (e.g. cousins, grandparent to grandchild)
- `thirdToFifth`: A third to fifth degree relationship (e.g. great-grandparent to great-grandchild)
- `sixthToTenth`: A sixth to tenth degree relationship
- `unrelated`: Unrelated – this is the default state among all individuals, which does not have to be expressed explicitly. This category will therefore probably never be used
- `other`: Any other kind of relationship not covered by the aforementioned categories

For each entry in `Relation_To` there must be a corresponding entry in `Relation_Degree`.

`Relation_Type` allows to add more verbose details about the relationship type, if it was possible to reconstruct that from the archaeological or historical context. Because there are too many possible permutations, there is no pre-defined set of values for what can and cannot be entered here. It is advisable, though, to stick to a general scheme like the following, which describes a given relationship from the point of view of the current individual:

- `same_as`: This sample is from the same individual as another sample
- `identical_twin_of`: This individual is likely an identical twin of another individual
- `father_of`: This individual is likely the father of the partner individual
- `grandchild_of`: This individual is likely the grandchild of the partner individual
- `mother_or_daughter_of`: This individual is likely either the mother or daughter of the partner individual (which might be unclear, in case of imprecise archaeological dating)
- `unknown`: The relationship is unclear or not yet determined. This is the default state and does not have to be expressed, unless multiple relationships are present and some but not all are known.
- ...

Unlike `Relation_Degree`, `Relation_Type` can be left empty even if there are entries in `Relation_To`. But if it is filled, then the number of values must be equal to the number of entries in both `Relation_To` and `Relation_Degree`.

The `Relation_Note` column allows to add free-form text information about the relationships of this individual. This might also include information about the method used to infer the degree and type.

3 Spatial position

The `.janno` file contains six columns to describe the spatial origin of an individual sample: `Country`, `Country_ISO`, `Location`, `Site` and finally `Latitude` and `Longitude`.

106 The **Country** column should contain a present-day political country name following the **English short name** in
 107 [ISO 3166](#).

108 The **Country_ISO** column should contain the present-day political country of origin of the sample, expressed in
 109 codes using the standard [ISO 3166-1](#) alpha-2 code, i.e. “AR” for Argentina or “NO” for Norway.

110 The **Location** column allows for free-form text entry and can contain further, unspecified location information.
 111 This might be the name of an administrative or geographic region, or an arbitrary unit of reference like a
 112 mountain, lake or city close to the point of discovery of the respective sample.

113 The **Site** column should contain a site name, ideally in the latin alphabet and ideally the name that is commonly
 114 used in publications.

115 The **Latitude** and **Longitude** columns should contain geographic coordinates (WGS84) in decimal degrees (DD)
 116 with a precision of not more than five places after the decimal point. This yields a precision of about [1.1132m at](#)
 117 [the equator](#) which is sufficient to describe the position of an archaeological site. Coordinates in other formats
 118 like for example Degrees Minutes Seconds (DMS) or in completely different coordinate reference systems should
 119 be transformed. There exist many open source software solutions to do that, most based on the [PROJ library](#)
 120 e.g. the [The World Coordinate Converter](#).

121 4 Temporal position

122 The temporal position of a sample is encoded with seven different columns in the `.janno` file: `Date_C14_Labnr`,
 123 `Date_C14_Uncal_BP`, `Date_C14_Uncal_BP_Err`, `Date_BC_AD_Median`, `Date_BC_AD_Start`, `Date_BC_AD_Stop`,
 124 `Date_Type`.

125 4.1 General structure

126 The `Date_Type` column handles the general distinction between the most common forms of age information:

- 127 • **modern**: Applies to present-day reference samples, so not ancient DNA.
- 128 • **C14**: Applies if there is a set of radiocarbon dates explicitly listed in the columns `Date_C14_Labnr`,
 129 `Date_C14_Uncal_BP` and `Date_C14_Uncal_BP_Err` whose post-calibration probability distribution is a
 130 meaningful prior for the individual’s year of death. The dates do not always have to be directly from the
 131 individual’s tissue, but they should be immediately relevant for their year of death (e.g. a date from a
 132 grain kernel recovered from the individual’s grave).
- 133 • **contextual**: Applies in all other cases if the columns `Date_BC_AD_Median`, `Date_BC_AD_Start`,
 134 `Date_BC_AD_Stop` can be filled. This includes age attribution based on the archaeologically determined
 135 stratigraphy or typological information. **contextual** should also be chosen if the sample is dated very
 136 indirectly with radiocarbon dating (e.g. radiocarbon dates from other, unrelated features of the same site)
 137 or dated with other physical or chemical dating methods (e.g. dendrochronology or optically stimulated
 138 luminescence).

139 So `Date_C14_Labnr`, `Date_C14_Uncal_BP` and `Date_C14_Uncal_BP_Err` only go along with `Date_Type = C14`,
 140 whereas `Date_BC_AD_Median`, `Date_BC_AD_Start`, `Date_BC_AD_Stop` complement both `Date_Type = C14` and
 141 `Date_Type = contextual`. Radiocarbon dates that only serve as secondary evidence for a contextual dating
 142 should NOT be reported in `Date_C14_Labnr`, `Date_C14_Uncal_BP` and `Date_C14_Uncal_BP_Err`.

143 4.2 The columns in detail

144 Each radiocarbon date has a unique identifier: the “lab number”. It consists of a lab code issued by the journal
145 [Radiocarbon](#) for each laboratory and a serial number. This lab number makes the date well identifiable and
146 should be reported in `Date_C14_Labnr` with the lab code separated from the serial number with a minus symbol.

147 The uncalibrated radiocarbon measurement can be described by a Gaussian distribution with mean and standard
148 deviation. So the column `Date_C14_Uncal_BP` holds the mean of that distribution in years before present (BP)
149 as usually reported by radiocarbon laboratories. The age is always a positive integer value starting from a zero
150 that corresponds to 1950 AD. The column `Date_C14_Uncal_BP_Err` holds the respective standard deviation for
151 each date in years. This should be the 1-sigma distance, so that the probability that the actual uncalibrated age
152 of the measured sample is within the $\text{Date_C14_Uncal_BP} \pm \text{Date_C14_Uncal_BP_Err}$ range is about 68%.

153 `Date_C14_Labnr`, `Date_C14_Uncal_BP` and `Date_C14_Uncal_BP_Err` each can hold multiple values separated
154 by ; to allow for multiple radiocarbon dates for each aDNA sample. With multiple values the number and order
155 of values in the columns must be consistent.

156 In the columns `Date_BC_AD_Median`, `Date_BC_AD_Start`, `Date_BC_AD_Stop` ages are reported in years BC and
157 AD, so in relation to the zero point of the Gregorian calendar. BC dates are represented with negative, AD with
158 positive integer values.

- 159 • If radiocarbon dates are available (`Date_Type = C14`): `Date_BC_AD_Median` should report the median age
160 after calibration. With multiple dates this can be determined either with sum calibration or more complex
161 (e.g. bayesian) age modelling. `Date_BC_AD_Start` and `Date_BC_AD_Stop` should report the starting/ending
162 age of a 95% probability window around the age median.
- 163 • If only contextual (e.g. from archaeological typology) age information is available (`Date_Type =`
164 `contextual`): `Date_BC_AD_Start` and `Date_BC_AD_Stop` should simply report the approximate start and
165 end date determined by the respective source of scientific authority (e.g. an archaeologist knowledgeable
166 about the relevant typological sequences). In this case `Date_BC_AD_Median` should be calculated as the
167 mean of `Date_BC_AD_Start` and `Date_BC_AD_Stop` rounded to an integer value.
- 168 • If the sample is a modern reference sample (`Date_Type = modern`): `Date_BC_AD_Median`, `Date_BC_AD_Start`,
169 `Date_BC_AD_Stop` should all be set to the value 2000, for 2000 AD.

170 The column `Date_Note` stores arbitrary free-form text information about the dating of a sample.

171 5 Genetic summary data

172 5.1 Individual properties

173 The `Genetic_Sex` column should encode the biological sex as determined from the DNA read distribution on
174 the X and Y chromosome. It only allows for the entries

- 175 • F: female
- 176 • M: male
- 177 • U: unknown

178 This limitation stems from the genotype data formats by Plink and the Eigensoft software package. Edge cases
179 (e.g. XXY, XYY, X0, ...) can not be expressed with this format and should be reported as U with an additional
180 comment in the free text `Note` field. Genetic sex determination for ancient DNA can be performed for example
181 with `Sex.DetERRmine` [4].

182 The **MT_Haplogroup** column is meant to store the human mitochondrial DNA haplogroup for the respective
183 individual in a simple string. The entry can be arbitrarily precise. A software tool to determine the MT haplogroup
184 is for example Haplogrep [5].

185 The **Y_Haplogroup** column holds the respective human Y-chromosome DNA haplogroup in a simple string. To
186 avoid confusion from using different haplotype naming systems, the notation should follow a syntax with the
187 main branch + the most terminal derived Y-SNP separated with a minus symbol (e.g. R1b-P312), similar to
188 that used by Yfull.

189 5.2 Library properties

190 The **Source_Tissue** column documents the skeletal, soft tissue or other elements from which source material
191 for DNA library preparation was extracted. If multiple samples have been taken from different elements, these
192 can be listed separated by ;. Specific bone names should be reported with an underscore (e.g. bone_phalanx,
193 tooth_molar).

194 The **Nr_Libraries** column holds a simple integer value of the number of libraries that have been prepared for
195 an individual.

196 The **Library_Names** column should list the names for the libraries as used in the publication, separated by ;.

197 The **Capture_Type** column specifies the general pre-sequencing preparation methods that have been applied to
198 the library. See [6] for a review of the different techniques (not including newer developments). This field can
199 hold one of multiple different values, but also multiple of these separated by ; if different methods have been
200 applied for different libraries.

- 201 • **Shotgun**: Sequencing without any enrichment (whole genome sequencing, screening etc.).
- 202 • **1240K**: Target enrichment with hybridization capture optimised for sequences covering the 1240k SNP
203 array [7], [8], [9].
- 204 • **ArborComplete**, **ArborPrimePlus**, **ArborAncestralPlus**: Target enrichment with hybridization capture
205 as provided by Arbor Biosciences in three different kits branded myBaits Expert Human Affinities.
- 206 • **TwistAncientDNA**: Target enrichment with hybridization capture as provided by Twist Bioscience [10].
- 207 • **OtherCapture**: Target enrichment with hybridization capture for any other set of sequences.
- 208 • **ReferenceGenome**: Modern reference genomes where aDNA fragmentation is not an issue and other sample
209 preparation techniques apply.

210 The **UDG** column documents if the libraries for the respective individual went through UDG (or USER enzyme)
211 treatment. This wet lab protocol step removes molecular damage in the form of deaminated cytosines characteristic
212 of ancient DNA.

- 213 • **minus**: A protocol without UDG treatment (e.g. [11]).
- 214 • **half**: A protocol with UDG-half treatment (e.g. [12]).
- 215 • **plus**: A protocol with UDG-full treatment (e.g. [13]).
- 216 • **mixed**: Multiple libraries that went through different UDG treatment approaches, and whose data were
217 later merged.

218 The **Library_Built** column describes the library preparation method regarding single- or double-stranded
219 protocols. See e.g. [14] for more information.

- 220 • **ds**: Double-stranded library preparation.
- 221 • **ss**: Single-stranded library preparation.

- **mixed**: If multiple libraries with different strandedness were combined. See also the Sequencing Source File in the Poseidon package as a way to provide details.

The **Genotype_Ploidy** column stores whether the genotype calls for this individual are originally haploid or diploid. Even for diploid organisms, it is often useful to represent genotypes by single haploid alleles (so-called pseudo-haploid genotypes), for example to generate relatively unbiased genotype calls from low coverage data. Because both the PLINK and EIGENSTRAT genotyping formats always *encode* genotype calls as diploid (by “doubling” the pseudo-haploid genotypes), the information on the original Ploidy of the call gets lost. This column is therefore used to record the underlying calling procedure. This becomes important, for example, when sample sizes are queried to compute bias-correction factors when computing F-Statistics or FST. The **Genotype_Ploidy** column can contain one of the following values:

- **diploid**: True diploid genotype calls were made.
- **haploid**: Haploid genotypes were called and then doubled.

The column **Data_Preparation_Pipeline_URL** should finally store an URL that links to a complete and human-readable description of the computational pipeline (for example a specific configuration for nf-core/eager [15]) by which the sample data was processed.

5.3 Data yield

The **Endogenous** column holds the percentage of mapped reads over the total amount of reads that went into the mapping pipeline. That boils down to the DNA percentage of the library that matches the (human) reference. It should be determined from Shotgun libraries (so before any hybridization capture), not on target (i.e. across the whole genome, not specific positions), and before any mapping quality filtering. In case of multiple libraries only the highest value should be reported. The % endogenous DNA can be calculated for example with the [endorS.py](#) script.

The **Nr_SNPs** column gives the number of SNPs reported in the genotype data files for this individual.

The **Coverage_on_Target_SNPs** column reports the mean fold coverage on the SNP set of the genotype dataset (e.g. 1240K) for the merged libraries of this sample. To calculate the coverage it is necessary to determine which SNPs are covered how many times by the mapped reads. Individual SNPs might be covered multiple times, whereas others may not be covered at all by the highly deteriorated ancient DNA. The coverage for each SNP is therefore a number between 0 and n. The statistic can be determined for example with the QualiMap [16] software package. In case of multiple libraries, the total coverage should be given across all libraries.

5.4 Data quality

The **Damage** column contains the % damage on the first position of the 5' end for the main Shotgun library used for sequencing or capture. This is an important statistic to verify the age of ancient DNA. In case of multiple libraries you should report a value from the merged read alignment.

Contamination of ancient DNA with foreign reads is a major challenge for archaeogenetics. There exist multiple competing ideas, algorithms and software tools to estimate the degree of contamination for individual samples (e.g. ANGSD [17], contamLD [18] or hapCon [19]), with some methods only applicable under certain circumstances (e.g. popular X-chromosome based approaches only work on male individuals). Also the results of different methods tend to differ both in the degree of contamination they estimate and in the way the output is usually encoded. To cover the multitude of methods in this domain, and to make the results representable in the .janno file, we offer the **Contamination_*** column family.

262 **Contamination** is a list column to represent the different contamination values estimated for a sample with one
263 or multiple software tools. As usual multiple values are separated by ;.

264 **Contamination_Err** is another list column to store the respective (standard) error term for the values in
265 **Contamination**.

266 Some tools for contamination estimation do not return a mean plus a standard error. ContamMix, for example,
267 yields a 95% confidence interval instead, to better represent assymetric output distributions. **Contamination**
268 and **Contamination_Err** can not represent this. We suggest to derive a mean and a standard error from these
269 alternative outputs. The latter can be calculated as the largest distance from the mean to the limits of the
270 confidence interval.

271 **Contamination_Meas** finally is the third necessary list column, which contextualizes the values in **Contamination**
272 and **Contamination_Err**. Each measure in these columns has to be accompanied by the software and software
273 version used to calculate it. The individual entries might e.g. look like this:

- 274 • ANGSD v0.935
- 275 • hapCon v0.4a1
- 276 • custom script

277 This setup has the consequence that the columns **Contamination**, **Contamination_Err**, **Contamination_Meas**
278 always have to have the same number of ;-separated values.

279 The **Contamination_Note** column is a free text field to add additional information about the contamination
280 estimates, e.g. which parameters where used with the respective software tools.

281 6 Context information

282 The **Genetic_Source_Accession_IDs** column was introduced to link the derived genotype data in Poseidon
283 with the raw sequencing data typically uploaded to archives like the ENA [20] or SRA [21]. There, projects and
284 individual samples are given clear unique identifiers: Accession IDs. This janno column is supposed to store one
285 or multiple of these Accessions IDs for each individual/sample in Poseidon. If multiple are entered, then they
286 should be arranged by descending specificity from left to right (e.g. project id > sample id > sequencing run id).

287 The **Primary_Contact** column is a free-form text field that stores the name of the main or the corresponding
288 author of the respective paper for published data.

289 The **Publication** column holds either the value **unpublished** for (yet) unpublished samples or – for published
290 data – one or multiple citation-keys of the form **AuthorJournalYear** without any spaces or special characters.
291 These keys have to be identical to the **BibTeX** citation-keys identifying the respective entries in the **.bib** file of
292 the package. BibTeX is a file format to store bibliographic information, where each entry (article, book, website,
293 ...) is defined by a series of parameters (authors, year of publication, journal, ...). Here's an example **.bib** file
294 with two entries for [22] and [23]:

```
@article{CassidyPNAS2015,  
  doi = {10.1073/pnas.1518445113},  
  url = {https://doi.org/10.1073%2Fpnas.1518445113},  
  year = 2015,  
  month = {dec},  
  publisher = {Proceedings of the National Academy of Sciences},  
  volume = {113},
```



```

    number = {2},
    pages = {368--373},
    author = {Lara M. Cassidy and Rui Martiniano and Eileen M. Murphy and
              Matthew D. Teasdale and James Mallory and Barrie Hartwell
              and Daniel G. Bradley},
    title = {Neolithic and Bronze Age migration to Ireland and establishment
              of the insular Atlantic genome},
    journal = {Proceedings of the National Academy of Sciences}
}

@article{FeldmanScienceAdvances2019,
    doi = {10.1126/sciadv.aax0061},
    url = {https://doi.org/10.1126%2Fsciadv.aax0061},
    year = 2019,
    month = {jul},
    publisher = {American Association for the Advancement of Science ({AAAS})},
    volume = {5},
    number = {7},
    pages = {eaax0061},
    author = {Michal Feldman and Daniel M. Master and Raffaella A. Bianco and
              Marta Burri and Philipp W. Stockhammer and Alissa Mittnik and
              Adam J. Aja and Choongwon Jeong and Johannes Krause},
    title = {Ancient {DNA} sheds light on the genetic origins of early Iron Age
              Philistines},
    journal = {Science Advances}
}

```

295 The string `CassidyPNAS2015` is the citation-key of the first entry. To cite both publications in the `Publication`
 296 column, one would enter `CassidyPNAS2015;FeldmanScienceAdvances2019`.

297 When creating a new Poseidon package the `.bib` file should be filled together with the `Publication` column.
 298 One of the most simple ways to obtain the BibTeX entries may be to request them with the doi from the [doi2bib](#)
 299 web app. It could be necessary to adjust the result manually, though. The citation-key, for example, has to be
 300 replaced by the one used in the `Publication` column.

301 The `Note` column is a free-form text field that can contain small amounts of additional information that is not
 302 yet expressed in a more systematic form in the the other `.janno` file columns.

303 The `Keywords` column was introduced to allow for tagging individuals with arbitrary keywords. This should
 304 simplify sorting and filtering in personal Poseidon package repositories. Each keyword is a string and multiple
 305 keywords can be separated with `;`.

306 Arbitrary additional columns can be included in a `.janno` file, but they should be named in a way that they do
 307 not conflict with the Poseidon package specification. These columns will not be validated (assumed free-form
 308 text), but they will be preserved in the Poseidon package, and propagated during operations with `trident`
 309 `forge`.

310

- 311 [1] S. Eisenmann *et al.*, “Reconciling material cultures in archaeology with genetic data: The nomenclature of clusters emerging from archaeogenomic analysis,” *Scientific Reports*, vol. 8, no. 1, Aug. 2018, doi: [10.1038/s41598-018-31123-z](https://doi.org/10.1038/s41598-018-31123-z).
- 312 [2] J. M. Monroy Kuhn, M. Jakobsson, and T. Günther, “Estimating genetic kin relationships in prehistoric populations,” *PLOS ONE*, vol. 13, no. 4, p. e0195491, Apr. 2018, doi: [10.1371/journal.pone.0195491](https://doi.org/10.1371/journal.pone.0195491).
- 313 [3] A. B. Rohrlach, J. Tuke, D. Popli, and W. Haak, “BREADR: An R package for the bayesian estimation of genetic relatedness from low-coverage genotype data,” Apr. 2023, doi: [10.1101/2023.04.17.537144](https://doi.org/10.1101/2023.04.17.537144).
- 314 [4] T. C. Lamnidis *et al.*, “Ancient Fennoscandian genomes reveal origin and spread of Siberian ancestry in Europe,” *Nature Communications*, vol. 9, no. 1, Nov. 2018, doi: [10.1038/s41467-018-07483-5](https://doi.org/10.1038/s41467-018-07483-5).
- 315 [5] S. Schönherr, H. Weissensteiner, F. Kronenberg, and L. Forer, “Haplogrep 3 - an interactive haplogroup classification and analysis platform,” *Nucleic Acids Research*, vol. 51, no. W1, pp. W263–W268, Apr. 2023, doi: [10.1093/nar/gkad284](https://doi.org/10.1093/nar/gkad284).
- 316 [6] M. Knapp and M. Hofreiter, “Next generation sequencing of ancient DNA: Requirements, strategies and perspectives,” *Genes*, vol. 1, no. 2, pp. 227–243, Jul. 2010, doi: [10.3390/genes1020227](https://doi.org/10.3390/genes1020227).
- 317 [7] Q. Fu *et al.*, “An early modern human from Romania with a recent Neanderthal ancestor,” *Nature*, vol. 524, no. 7564, pp. 216–219, Jun. 2015, doi: [10.1038/nature14558](https://doi.org/10.1038/nature14558).
- 318 [8] W. Haak *et al.*, “Massive migration from the steppe was a source for Indo-European languages in Europe,” *Nature*, vol. 522, no. 7555, pp. 207–211, Mar. 2015, doi: [10.1038/nature14317](https://doi.org/10.1038/nature14317).
- 319 [9] I. Mathieson *et al.*, “Genome-wide patterns of selection in 230 ancient Eurasians,” *Nature*, vol. 528, no. 7583, pp. 499–503, Nov. 2015, doi: [10.1038/nature16152](https://doi.org/10.1038/nature16152).
- 320 [10] N. Rohland, S. Mallick, M. Mah, R. Maier, N. Patterson, and D. Reich, “Three assays for in-solution enrichment of ancient human DNA at more than a million SNPs,” *Genome Research*, vol. 32, no. 11–12, pp. 2068–2078, Nov. 2022, doi: [10.1101/gr.276728.122](https://doi.org/10.1101/gr.276728.122).
- 321 [11] F. Aron, G. U Neumann, and G. Brandt, “Non-UDG treated double-stranded ancient DNA library preparation for Illumina sequencing v1,” Dec. 2019, doi: [10.17504/protocols.io.bakricv6](https://doi.org/10.17504/protocols.io.bakricv6).
- 322 [12] F. Aron, G. U Neumann, and G. Brandt, “Half-UDG treated double-stranded ancient DNA library preparation for Illumina sequencing v1,” Sep. 2020, doi: [10.17504/protocols.io.bmh6k39e](https://doi.org/10.17504/protocols.io.bmh6k39e).
- 323 [13] F. Aron, G. U Neumann, and G. Brandt, “Full-UDG treated double-stranded ancient DNA library preparation for Illumina sequencing v1,” Dec. 2020, doi: [10.17504/protocols.io.bqbpmsmn](https://doi.org/10.17504/protocols.io.bqbpmsmn).
- 324 [14] M.-T. Gansauge and M. Meyer, “Single-stranded DNA library preparation for the sequencing of ancient or damaged DNA,” *Nature Protocols*, vol. 8, no. 4, pp. 737–748, Mar. 2013, doi: [10.1038/nprot.2013.038](https://doi.org/10.1038/nprot.2013.038).
- 325 [15] J. A. Fellows Yates *et al.*, “Reproducible, portable, and efficient ancient genome reconstruction with nf-core/eager,” *PeerJ*, vol. 9, p. e10947, Mar. 2021, doi: [10.7717/peerj.10947](https://doi.org/10.7717/peerj.10947).
- 326 [16] K. Okonechnikov, A. Conesa, and F. García-Alcalde, “Qualimap 2: Advanced multi-sample quality control for high-throughput sequencing data,” *Bioinformatics*, vol. 32, no. 2, pp. 292–294, Oct. 2015, doi: [10.1093/bioinformatics/btv566](https://doi.org/10.1093/bioinformatics/btv566).
- 327 [17] T. S. Korneliussen, A. Albrechtsen, and R. Nielsen, “ANGSD: Analysis of next generation sequencing data,” *BMC Bioinformatics*, vol. 15, no. 1, Nov. 2014, doi: [10.1186/s12859-014-0356-4](https://doi.org/10.1186/s12859-014-0356-4).
- 328 [18] N. Nakatsuka, É. Harney, S. Mallick, M. Mah, N. Patterson, and D. Reich, “ContamLD: Estimation of ancient nuclear DNA contamination using breakdown of linkage disequilibrium,” *Genome Biology*, vol. 21, no. 1, Aug. 2020, doi: [10.1186/s13059-020-02111-2](https://doi.org/10.1186/s13059-020-02111-2).
- 329 [19] Y. Huang and H. Ringbauer, “hapCon: Estimating contamination of ancient genomes by copying from reference haplotypes,” *Bioinformatics*, vol. 38, no. 15, pp. 3768–3777, Jun. 2022, doi: [10.1093/bioinformatics/btac390](https://doi.org/10.1093/bioinformatics/btac390).

- 330 [20] J. Burgin *et al.*, “The European Nucleotide Archive in 2022,” *Nucleic Acids Research*, vol. 51, no. D1, pp. D121–D125, Nov. 2022, doi: [10.1093/nar/gkac1051](https://doi.org/10.1093/nar/gkac1051).
- 331 [21] K. Katz, O. Shutov, R. Lapoint, M. Kimelman, J. R. Brister, and C. O’Sullivan, “The Sequence Read Archive: A decade more of explosive growth,” *Nucleic Acids Research*, vol. 50, no. D1, pp. D387–D390, Nov. 2021, doi: [10.1093/nar/gkab1053](https://doi.org/10.1093/nar/gkab1053).
- 332 [22] L. M. Cassidy *et al.*, “Neolithic and Bronze Age migration to Ireland and establishment of the insular Atlantic genome,” *Proceedings of the National Academy of Sciences*, vol. 113, no. 2, pp. 368–373, Dec. 2015, doi: [10.1073/pnas.1518445113](https://doi.org/10.1073/pnas.1518445113).
- 333 [23] M. Feldman *et al.*, “Ancient DNA sheds light on the genetic origins of early Iron Age Philistines,” *Science Advances*, vol. 5, no. 7, p. eaax0061, Jul. 2019, doi: [10.1126/sciadv.aax0061](https://doi.org/10.1126/sciadv.aax0061).