







Poseidon – A framework for archaeogenetic human genotype data management

Clemens Schmid ^{1,2}, Ayshin Ghalichi ¹, Theseas C. Lamnidis ¹, Dhananjaya B. A. Mudiyansele ^{1,3}, Wolfgang Haak ¹, and Stephan Schiffels ¹

¹*Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany*

²*International Max Planck Research School for the Science of Human History, Max Planck Institute for Geoanthropology, Jena, Germany*

³*Saarland University, Saarbrücken, Germany*

1 Abstract

The study of ancient human genomes, archaeo- or palaeogenetics, has accelerated in the last ten years, with now thousands of new ancient genomes being released each year. Operating at the interface of genetics, anthropology and archaeology, this data includes features from all three fields, including rich meta- and context-data, for example regarding spatiotemporal provenience. While archives and standards for genetic sequencing data already exist, no such infrastructure exists for combined genetic and meta-data that could ensure FAIR principles across the field. Here, we present Poseidon, a framework for open and FAIR data handling in archaeogenetics, including a specified package format, software tools, and public, community-maintained online archives. Poseidon emphasises human- and machine-readable data storage, the development of convenient and interoperable command line software, and a high degree of source granularity to elevate the original data publication to the main unit of long-term curation.

Keywords: ancient DNA, data management, FAIR data

2 Introduction

The technology to extract and sequence DNA from human remains thousands of years old has revolutionised the study of the human past. This is documented by groundbreaking new insights, from our evolutionary relationships to distant relatives like Neanderthals [1, 2] to prehistoric [3, 4] and historic migrations [5, 6]. Since the sequencing of the first ancient modern human genome in 2010 [7], hundreds of studies have been published, accompanied by massive datasets of ancient human DNA sequences. A drop in sequencing costs and new technologies like hybridisation capture [8, 4, 9] have in fact lead to an acceleration of new published ancient genomes, with data now coming out faster than individual researchers typically can keep track of and co-analyse. Recently, the threshold of genome-wide data for 10,000 ancient human individuals has been surpassed [10].

To make all this new data publicly available, researchers can partly rely on existing infrastructure for the archival and distribution of modern genetic data, such as the Sequence Read Archive (SRA) [11], the European Nucleotide Archive (ENA) [12] or other INSDC databases (<https://www.insdc.org>). However, this infrastructure has not been prepared to also capture the rich context-data ranging from archaeological field

35 observations to radiocarbon dating that accompanies ancient samples. Nor is there a standardised archive yet for
36 derived genotype data that is routinely used to substantiate most if not all of the conclusions in archaeogenetic
37 papers. This raises multiple concrete issues:

- 38 • Ancient individuals only constitute meaningful observations if their spatiotemporal provenience is known.
39 Current practice renders it difficult to maintain the connection between archaeological context and sampled
40 genomic data, as this information is generally kept separately.
- 41 • Specific results of typical archaeogenetic analyses (e.g. PCA, F-Statistics, kinship estimation) can only
42 be fully reproduced with genotype-level data. But current practice is not to include this data with a
43 publication, be it for its unwieldy size or the lack of a central repository to easily share it.
- 44 • Meta-analyses involving large amounts of data require tedious curation. Despite the fact that archaeoge-
45 netic genotype data is largely standardized, and common practices for reporting processing steps, data
46 quality indicators, as well as spatial and temporal sample origin have emerged, combining information
47 from different papers is still hindered by severe structural variation.

48 A major project addressing some of these problems in human archaeogenetics is the Allen Ancient DNA
49 Resource (AADR), which is a curated dataset of public ancient DNA data assembled by the Reich Lab at
50 Harvard University [13]. While the AADR clearly fills a gap in the field, there continues to be a need for open
51 standardization and the creation of a community-maintained archive. Both standard and archive should be
52 well-specified to be human and machine-readable, fully open and transparent, and permanently downloadable
53 with their entire version history. They should be geared towards scientific practice and include well-documented
54 interfaces for research software. To ensure fairness, i.e. equal access for users and contributors from different
55 backgrounds, and long-term reliability, they should be as independent as possible from specific institutions, key
56 persons [14] or infrastructure providers, and ideally controlled by the community-of-practice as a whole.

57 Here, we present '*Poseidon*' (<https://www.poseidon-adna.org>), a computational framework including
58 an open data format, software, and community-maintained archives, to enable this standardised and FAIR
59 [15] handling of archaeogenetic data. The name *Poseidon* is inspired by the notion of a *sea of data* benefiting
60 from structure and governance. We imagine Poseidon to simplify a variety of concrete workflows and mitigate
61 technical challenges practitioners in the field of human archaeogenetics regularly face:

- 62 • **Data storage:** Archaeogenetic samples/ancient individuals (we often use these terms synonymously be-
63 low, see Supplementary Text 8 for a definition) can only be effectively analysed with context data. The
64 Poseidon package format allows one to store archaeogenetic genotype data with arbitrary spatiotemporal
65 or archaeological information on a per-sample level. The package format enforces human- and machine-
66 readability for the context data to enable its computational co-analysis with the genomic data, while also
67 maintaining a high level of flexibility and extensibility to accommodate arbitrary additional variables.
- 68 • **Data acquisition:** Research in archaeogenetics is strongly dependent on incorporating published reference
69 data. Poseidon features public archives with per-article packages that can be downloaded through an open
70 web API. The packages include genotypes, context data and machine-readable citation information. To
71 ensure computational reproducibility, Poseidon supports version tracking for these packages, with old
72 versions directly available through the web interface. Beyond hosting data, the web infrastructure also
73 provides options to report issues and suggest transparent data updates.
- 74 • **Data analysis:** Working with common software tools in human archaeogenetics (e.g. smartpca [16, 17],
75 ADMIXTURE [18], qpAdm [4, 19]) requires frequent merging, sub-setting and file format conversion for the
76 samples of interest. The Poseidon core team develops the software tool *trident* to simplify these operations

77 both with Poseidon packages and unpackaged genotype data. The unified package format allows for new
78 analysis tools with flexible data detection, low-memory stream processing and on-the-fly aggregation of
79 widely used genome-wide statistics as demonstrated in the *xerxes* software tool also developed by the core
80 team.

81 • **Data publication** For each sample, human archaeogenetic papers should include genotypes, context
82 information and quality-control data, so e.g. estimates for the aDNA damage or contamination. Poseidon
83 offers a standardized and reusable way to share this data directly with the publication and/or through
84 our public archives.

85 3 Overview

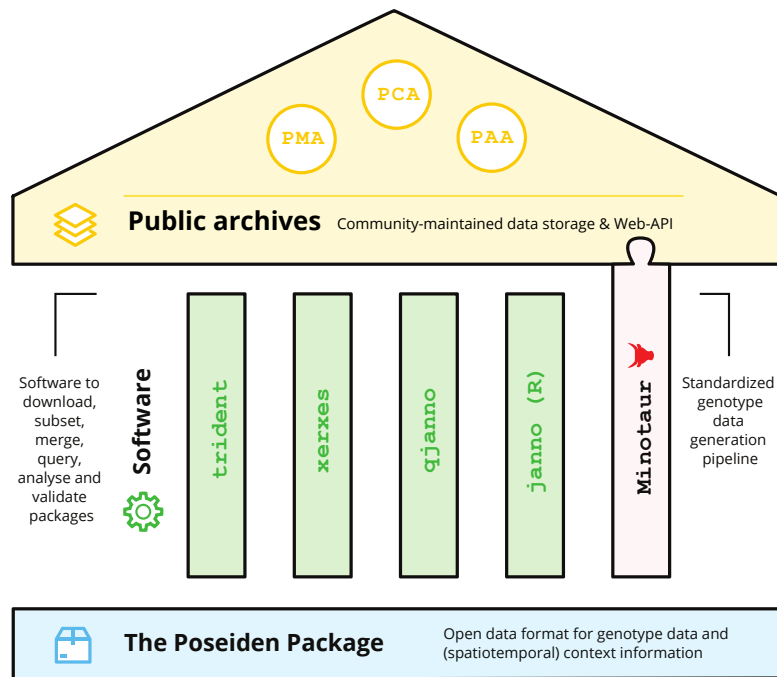


Figure 1: Schematic overview of the core components of the Poseidon framework. The Poseidon package specification forms the foundation on which various open source software tools and the Minotaur workflow are based. They, in turn, underpin and enable the public Poseidon data archives.

86 Poseidon consists of three major components: A data format, software, and public archives (see Figure 1).
87 At the foundation stands a data format to store genotype data together with context information. The software
88 implements this specification, relies on and validates its promises, and builds convenient, useful functionality
89 on top of it, both for individual users and for our cloud infrastructure. Finally, the public archives store data
90 using the data format and employ the software for validation, modification and book-keeping.

91 The Poseidon schema (Supplementary Text 1) defines the structure and format of a Poseidon package (Figure
92 2). This includes the general layout and purpose of the main files, the POSEIDON.yml, the .janno and the .ssf
93 file, and detailed definitions of the variables within them, both regarding semantics and syntax (Supplementary
94 Text 2). The short definitions in the schema are explained in more detail on the Poseidon website, which also
95 serves as a central hub for all components of the framework.

96 The Poseidon software suite enables users to create, download, inspect, subset, merge and analyse Poseidon
97 packages (Supplementary Text 3-7). It is mostly implemented in Haskell [20], a purely functional programming
98 language, and split over multiple command line tools openly available as statically compiled executables for

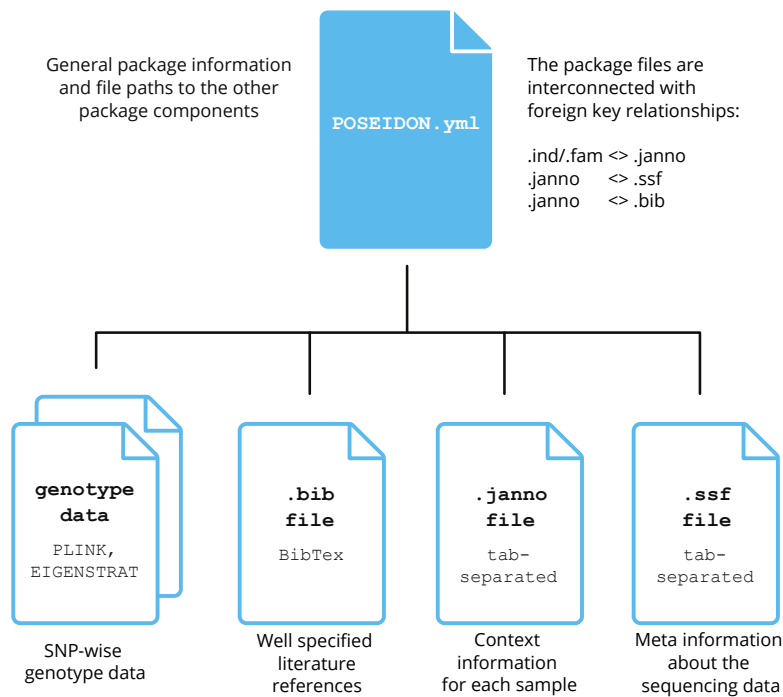


Figure 2: Schematic overview of the Poseidon package structure. The POSEIDON.yml file defines the package and interlinks the additional data files for genotype, context- and bibliography information in a relational structure.

99 the major desktop and server operating systems. The main tool, *trident*, provides data handling functionality,
100 including the code necessary for the client-server infrastructure of the public archives. The *xerxes* software
101 implements commonly used genome-wide allele-frequency statistics, *qjanno* allows for SQL-like queries on .janno
102 files and the *janno* R package offers a .janno file interface for the R programming language.

103 Using these tools, we conceived and implemented three public data archives for sharing and maintaining
104 published (ancient) human DNA data packaged in the Poseidon format: The Poseidon Community Archive
105 (PCA), the Poseidon Minotaur Archive (PMA) and the Poseidon AADR Archive (PAA). They already store
106 considerable amounts of public genotype data (Figure 4, Figure 5 and Supplementary Text 8). Each of these
107 archives is at their base represented by a public repository on GitHub, versioned with Git and capable of holding
108 large genotype data files through an integration with Git's Large File Storage (LFS) functionality. This setup
109 allows for community driven maintenance through Git pull requests (Figure 3). A custom webserver mirrors the
110 data from GitHub with an explicit version history, and allows to query the archive data and download packages
111 in the browser, from the command line and through web-frontends. While the PCA focuses on author-submitted
112 genotype data, the PMA only accepts data that went through a specific standardized genotype data processing
113 pipeline, the Minotaur workflow, with a semi-automatic interface on GitHub and a processing queue currently
114 hosted on computational infrastructure at the Max Planck Institute for Evolutionary Anthropology (MPI-EVA).
115 Finally, the PAA stores releases of the AADR dataset in Poseidon format, after a light-weight curation step to
116 ensure format compatibility.

117 4 Structure

118 The following sections explain the different components of Poseidon in detail: The Poseidon package, the software
119 tools, and the public archives including the Minotaur processing workflow. The underlying repositories with
120 specifications and code are available on GitHub (<https://github.com/poseidon-framework>) and, each in a
121 current version, at a long-term archive: <https://doi.org/10.17605/OSF.IO/ZUQGB>

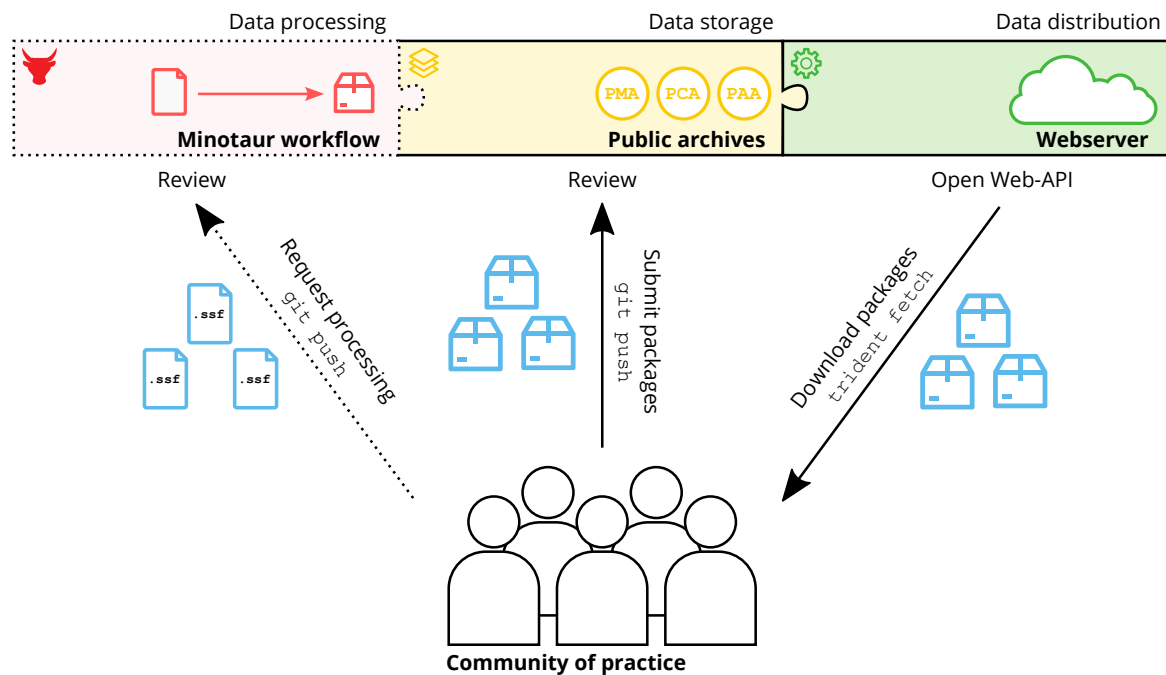


Figure 3: Schematic overview of the most common interaction pathways between the archaeogenetic community of practice and the Poseidon infrastructure. Community members share data either by preparing .ssf files as build-instructions for the Minotaur workflow or by directly submitting packages to the Community Archive. The data in the archives can then in turn be downloaded from the archives through the Poseidon webservice and its API.

122 4.1 The Poseidon Package specification (v2.7.1)

123 The core idea of Poseidon is to organize genotype data together with relevant meta- and context data in a
124 structured yet flexible, human- and machine-readable format. This format is the Poseidon package, defined in
125 a semantically versioned [21] specification, openly available online (<https://github.com/poseidon-framework/poseidon-schema>), and also part of this publication as Supplementary Text 1. A Poseidon package must
126 contain a POSEIDON.yml file and genotype data in PLINK or EIGENSTRAT format. It should additionally
127 contain a .janno file to store sample-wise context information and a .bib file for literature references. Option-
128 ally, it can also contain an .ssf file with information on the underlying raw sequencing data, an unstructured
129 README.md file for arbitrary meta information and a CHANGELOG.md file to document changes to the
130 package.
131

```
132 1 poseidonVersion: 2.7.1  
133 2 title: Switzerland_LNBA  
134 3 description: Genetic data from Late Neolithic and Bronze Age Switzerland  
135 4 contributor:  
136 5   - name: Roswita Malone  
137 6     email: roswita.malone@example.org  
138 7 packageVersion: 1.1.2  
139 8 lastModified: 2024-03-15  
140 9 genotypeData:  
141 10   format: PLINK  
142 11   genoFile: Switzerland_LNBA.bed  
143 12   snpFile: Switzerland_LNBA.bim  
144 13   indFile: Switzerland_LNBA.fam  
145 14   snpSet: 1240K  
146 15 jannoFile: Switzerland_LNBA.janno
```

147.6 BibFile: Switzerland_LNBA.bib

Example 1: A typical, neither minimal nor maximal, POSEIDON.yml file.

148 A Poseidon package is defined by a POSEIDON.yml file, using the YAML markup language ([https://ya](https://yaml.org)
149 [ml.org](https://yaml.org)) with a set of predefined fields, which stores some general information and, most importantly, relative
150 paths to the other files in the package (see Example 1). A full list of the specified fields is provided as part of
151 the schema, but only `poseidonVersion` (the schema version the package adheres to), `title` (a short package
152 name), and `genotypeData` (context information for and relative paths to the genotype data files) are mandatory.
153 Besides its function to define the package, the POSEIDON.yml also enables software to easily crawl for packages
154 in a file system.

155 At the time of writing, Poseidon supports genotype data for single nucleotide polymorphisms (SNPs) in two
156 common file formats: a binary-encoded format as introduced by the plink software package (PLINK) [22] and a
157 plain-text, human-readable format defined for the eigensoft software tools (EIGENSTRAT) [17]. Both formats
158 are structurally similar, and split into three files: A genotype table (`.bed/.geno`); a SNP file (`.bim/.snp`), defining
159 the SNPs in the genotype table; and an individual file (`.fam/.ind`) for the respective samples. In the future,
160 Poseidon may add support for other formats, e.g. the more flexible vcf file format [23]. Arbitrary SNP sets are
161 allowed and supported, but special keywords are reserved for the commonly used Affymetrix *HumanOrigins*
162 array [24], and the so-called *1240k* in-solution hybridization capture reagent [25].

163 The `.janno` file accompanies the genotype data to provide context for each sample. It is designed as a tabular,
164 tab-separated text file with a set of predefined columns. Each row corresponds to one entry in the individual file
165 (`.fam/.ind`), featuring at least the following mandatory columns: `Poseidon_ID` (a sample identifier), `Genetic_Sex`
166 (Female, Male, Unknown) and `Group_Name` (one or multiple group identifiers). Other, optional columns include
167 standard information on the spatiotemporal provenience of a given sample, its genetic data quality metrics, and
168 which major laboratory procedures it was subjected to. Supplementary Text 2 includes documentation for all
169 specified columns of the `.janno` file. One particular column documents the scientific publication(s) within which
170 the genetic or contextual data for a given sample was originally reported. This is a list column with BibTeX
171 keys, which, in turn, must be specified in the `.bib` file of the Poseidon package, thus ensuring that each sample
172 can be properly cited.

173 The `.ssf` file is similarly structured as the `.janno` file and stores sequencing source data, i.e. meta-information
174 about the raw sequencing data behind the genotypes in a Poseidon package. The rows in this table correspond
175 to sequencing entities, typically the set of unprocessed reads sequenced from DNA libraries or even multiple
176 runs/lanes of the same library. At the time of writing, the specified columns mostly focus on how to download
177 a given dataset from the INSDC databases such as ENA or SRA (<https://www.insdc.org>). The `.ssf` file does
178 not have any mandatory variables, but entries can be linked to the package with the list column `poseidon_IDs`
179 in a many-to-many relationship.

180 4.2 Software tools

181 The following software tools were developed to work with Poseidon packages and other infrastructure of the
182 Poseidon ecosystem. They can be considered reference implementations of the Poseidon schema, but are by no
183 means exclusive nor comprehensive. The Poseidon schema is defined as an independent, versioned entity, and
184 software developers can easily implement other tools in addition to what we as the Poseidon core team currently
185 offer.

186 All Poseidon software is open-source (MIT-License) and available on GitHub ([https://github.com/posei](https://github.com/poseidon-framework)
187 [don-framework](https://github.com/poseidon-framework)). That is also where users should report issues.

188 4.2.1 trident (v1.4.1.0)

189 *trident* is a command line software tool to create, download, inspect, merge, and subset Poseidon packages –
190 and therefore the central tool of the Poseidon framework. It is implemented in Haskell as one executable for
191 the *poseidon-hs* library and can be installed from source with the build systems cabal [26] or stack [27]. To
192 ease installation we provide pre-compiled, static binary executables for Linux, Windows and macOS. We also
193 maintain a package recipe under *poseidon-trident* on the bioconda channel [28], enabling installation of *trident*
194 through the conda (<https://docs.anaconda.com>) package management system.

195 *trident* includes multiple sub-commands for different operations on and with Poseidon packages – Supple-
196 mentary Text 3 gives a detailed overview including all specific arguments. It supports two mechanisms to obtain
197 Poseidon packages: A user can create them from genotype data with `trident init` or download them from our
198 community-maintained archives with `trident fetch`. The most involved and technically complex sub-command
199 in *trident* is `trident forge`, which allows users to both subset and merge Poseidon packages. To simplify the
200 process of package maintenance, `trident genoconvert` converts the genotype data between the formats Po-
201 seidon supports and `trident rectify` automates some common tasks such as checksum and version updates.
202 For package inspection, *trident* includes `trident list`, which returns lists of packages, groups, or samples for
203 a given package selection (either locally, or remotely by accessing the web API for the community-maintained
204 online archives). `trident summarise` compiles some basic summary statistics about a package collection and
205 `trident survey` gives an overview to which degree .janno files in a package collection are filled with data, i.e.
206 the level of completeness for context information. The last and arguably the most important inspection sub-
207 command is `trident validate`, which parses all components of a package to report violations of the Poseidon
208 schema. This ensures the structural integrity of Poseidon packages and maintains machine-readability.

```
209 1 trident list --remote --individuals --raw | grep "Finland\\|Volga0ka"  
210 2 trident fetch -d . -f "*2018_Lamnidis_Fennoscandia*,*2023_Peltola_Volga0ka*"  
211 3 trident forge -d . -d ../myData/ -f "Finland_Levanluhta,Volga0ka_IA,<mySample>" -o newPac
```

Example 2: A basic *trident* command line workflow to explore the community archive, download relevant packages and create a new package from the downloaded data, as well as a local/private data collection. `trident list` queries the webserver (`--remote`), and returns a tab separated (`--raw`) table of individuals/samples (`--individuals`) available in the public Community Archive (PCA). This list can then be filtered with standard command line tools (here `grep`). With `trident fetch` two packages are selected for download (`-f ...`) into the current working directory (`-d .`). These two packages as well as an additional, local package collection (`-d ../myData/`) can then be read into `forge` to create a new package `newPac`, specifying (`-f ...`) the groups `Finland_Levanluhta`, `Volga0ka_IA` and the local sample `mySample`.

212 A simple *trident* workflow could look like Example 2. The first two commands here require interaction with
213 the Poseidon webserver. This server software is, in fact, also implemented in Haskell as a part of *poseidon-hs*
214 and can be started with a hidden (and for most end users irrelevant) *trident* sub-command `trident serve`.
215 On the client side `trident list` and `trident fetch` use this API and interact, by default, with the endpoints
216 at <https://server.poseidon-adna.org>. A different server can be set with `--remoteURL`. More about the
217 available endpoints below.

218 The `forge` sub-command is the technically most involved operation in *trident*. It discovers all Poseidon
219 packages under a list of base directories (`-d`), reads them and their components into dedicated in-memory
220 data structures, parses the query language in `-f` (or from a file with `--forgeFile`) to decide which entities
221 (individuals, groups, packages) should be selected, and ultimately generates the new package including genotypes
222 and context data in a single, low-memory stream processing run. All components of the Poseidon package format
223 (i.e. the POSEIDON.yml file, the genotype data, the .janno, the .ssf, and the .bib file) are modelled in *poseidon-*
224 *hs* as algebraic data types with dedicated parsers. Corrupted files are rejected upon reading, including cross-file

225 compatibility checks within each package (e.g. bibtex keys mentioned in the .janno file must be specified in the
226 .bib file). The query language supported by `trident forge` in `-f` is a flexible domain-specific language (DSL)
227 to describe arbitrary positive and negative entity selection scenarios. For the stream processing of genotype
228 data in binary PLINK and EIGENSTRAT format, *poseidon-hs* relies on the `sequence-formats` library
229 (<https://github.com/stschiff/sequence-formats>), which allows *trident* to read, filter and write data from
230 a large amount of files at once, while also unifying SNPs to a consensus by recoding alleles on-the-fly. Varying
231 SNP sets in multiple packages can be merged using either an intersection or a union operation.

232 4.2.2 xerxes (v1.0.1.0)

233 *xerxes* is another command line software tool based on the *poseidon-hs* Haskell library. While *trident* is meant
234 for data management, *xerxes* is intended for basic, every-day data analysis operations. Its most advanced and
235 stable sub-command `xerxes fstats` makes extensive use of the genotype data stream processing introduced
236 above to implement commonly used genome-wide statistics. These are F-Statistics (F_2 , F_3 and F_4 , see Example
237 3, in several variants differing in bias-correction and normalisation) [24], F_{ST} [29], heterozygosity, and pairwise
238 nucleotide mismatch rates for assessing pairwise relatedness. All statistics are evaluated with error-estimation
239 using weighted block-Jackknife [30]. See Supplementary Text 4 for the *xerxes* user guide and Supplementary
240 Text 5 for a whitepaper documenting the mathematical underpinnings of the implemented algorithms.

```
241 1 xerxes fstats -d . \  
242 2     --stat "F4(Mbuti,Nganasan,Lithuanian,Finland_Levanluhta)" \  
243 3     --stat "F4(Mbuti,Nganasan,Lithuanian,Russia_Bolshoy)"
```

Example 3: Calculating two F_4 statistics with *xerxes*, reproducing an analysis of Lamnidis et al. 2018 [31]. We assume the base directory (`-d .`) to include Poseidon packages with the groups `Mbuti`, `Nganasan`, `Lithuanian`, `Finland_Levanluhta` and `Russia_Bolshoy`. These happen to be described in the Community Archive packages `2018_Lamnidis_Fennoscandia` [31], `2014_LazaridisNature` [3] and `2012_PattersonGenetics` [24], but *xerxes* will find those automatically if they are stored below the base directory. The individual statistics are specified with `--stats` and a dedicated domain specific language. For complex analyses *xerxes* offers a more powerful configuration file interface with the `--statConfig` argument.

```
244 1 groupDefs:  
245 2   CEU2: ["CEU.SG", "-<NA12889.SG>", "-<NA12890.SG>"]  
246 3   FIN2: ["FIN.SG", "-<HG00383.SG>", "-<HG00384.SG>"]  
247 4 fstats:  
248 5 - type: FST # this will create 2x2 = 4 statistics  
249 6   a: ["French", "Spanish"]  
250 7   b: ["Han", "CEU2"]  
251 8 - type: F4 # this will create 2x2x2x1 = 8 statistics  
252 9   a: ["French", "Spanish"]  
253 0   b: ["Han", "CEU2"]  
254 1   c: ["Nganasan", "Lithuanian"]  
255 2   d: ["Mbuti"]
```

Example 4: An example of a configuration file for *xerxes*, specifying a total of 12 statistics of F_{ST} and F_4 to be computed for all combinations of listed populations. In addition, the example features adhoc group definitions in the `groupDefs` section, which shows how to exclude individuals from a group using the minus-sign. More complex features of these configuration files include frequency-ascertainment, and selection of entire packages to act as groups.

256 We highlight three *xerxes* features that stand out: First, the user interface provides for a powerful way to
257 specify families of statistics using combinatorial expansions. For example, users can specify in a configuration file
258 (Example 4) that F_{ST} should be computed between all combinations of two lists of populations, and multiple

259 such families can be requested simultaneously. Second, statistics can be easily adapted with custom group-
260 definitions, in which individuals, groups or entire packages (similar to the `trident forge` selection language)
261 can be excluded or included. If a specific individual turns out to be an outlier that should not be included in
262 the computation of allele frequencies of its group, one can define a custom group and exclude that individual on
263 the fly, without the need to change the underlying genotype definition files. Third, users do not have to specify
264 the exact packages that their populations or individuals reside in, but can pass a base directory with dozens or
265 hundreds of packages, and *xerxes* will automatically select only the relevant packages, involving the groups or
266 individuals needed, thereby reducing memory- and run-time overhead: Similar to `trident forge`, *xerxes* uses
267 stream-processing to merge the relevant genotype-files without the need to load large genotype matrices into
268 memory. This is particularly important for dense genotyping data sets, for example involving tens of millions
269 of SNPs as with the 1000 Genomes dataset [32].

270 4.2.3 `qjanno` (v1.0.0.0)

271 *qjanno* is a command line tool implemented in Haskell to run SQL queries on `.janno` files (or arbitrary delimiter-
272 separated text files) and therefore to conveniently subject the context data provided with Poseidon packages
273 to a wide range of data transformations. It was built on top of a hard fork of v0.3.3 of the `qsh` package
274 (<https://github.com/itchyny/qhs>, MIT-License). See Supplementary Text 6 for the *qjanno* user guide.

275 On startup, *qjanno* creates an SQLite [33] database in memory. It then reads the requested, well-structured
276 text files, attributes each column a type and writes the contents of the files to tables in the in-memory database.
277 It finally sends the user-provided SQL query to the database, waits for the result, parses it and returns it on the
278 command line. The query gets pre-parsed to extract file names and then forwarded to an SQLite database server
279 via the Haskell library `sqlite-simple` (<https://github.com/nurpax/sqlite-simple>). That means *qjanno* can
280 parse and understand most SQLite3 syntax (see Example 5 and 6) with some minor exceptions (e.g. PRAGMA
281 functions).

```
282 1 qjanno "SELECT Poseidon_ID FROM 2018_Lamnidis_Fennoscandia.janno WHERE Country <> 'Finland'"
```

Example 5: A simple *qjanno* query on the command line. It extracts the columns `Poseidon_ID` and `Country` of the file
`2018_Lamnidis_Fennoscandia.janno` and returns all rows where `Country` is not `'Finland'`.

283 *qjanno* does not have a complete understanding of the `.janno`-file structure, and mostly treats it like a normal
284 `.tsv` file. But `.janno` files are still given special consideration with a number of pseudo-functions, which allow
285 to search Poseidon packages and `.janno` files recursively to load them together into one database table (see
286 Example 6).

```
287 1 qjanno "SELECT Date_Type ,count(*) AS n FROM d(.) GROUP BY Date_Type"
```

Example 6: Another *qjanno* query. With `'FROM d(.)'` *qjanno* searches all latest versions of Poseidon packages under
the current working directory, reads their `.janno` files and appends them to a common database table. This
table is then grouped by the `Date_Type` column (*C14*, *contextual*, *modern*). The output features a row-count
for each group in a new summary column `n`.

288 4.2.4 `janno` R package (v1.0.0)

289 The *janno* R package simplifies loading and handling `.janno` files in R and the popular tidyverse [34] R package
290 ecosystem. It provides a dedicated R S3 class `janno` that inherits from the `tibble` class to allow tidy reading and
291 manipulating the context information in a Poseidon package (see Example 7). Supplementary Text 7 features
292 its user guide.

```
293 1 janno::read_janno(".", ".") |>
```

```
294 2 dplyr::group_by(Date_Type) |>
295 3 dplyr::summarise(n = dplyr::n())
```

Example 7: A small R workflow enabled by the *janno* R package. The code is equivalent to the *qjanno* Example 6. `read_janno()` discovers and reads all *.janno* files under the current working directory into an object of class *janno*. The dplyr verbs `group_by` and `summarise` can then be applied to this object to count the rows per *Date_Type* group.

296 As *trident* and *qjanno* (with `-d .../d(...)`), the package's `read_janno()` function searches *.janno* files re-
297 cursively under a given directory and loads them into a single data frame. The reading process includes validation
298 according to the Poseidon schema. String list columns in *.janno* files are translated to true list columns in R.
299 Beyond basic functionality (`print()`, `write_janno()`), the package features one more major function for *janno*
300 objects: `process_age()`. This function processes the age information in the *Date_** columns of the *.janno* file to
301 derive a set of new columns useful for further chronological analysis via radiocarbon calibration: *Date_BC_AD_Prob*
302 is a list column with the complete post-calibration probability distribution, *Date_BC_AD_Median_Derived* is the
303 median of this distribution and *Date_BC_AD_Sample* stores *n* random samples drawn from it. The radiocarbon
304 calibration is implemented with the *Bchron* R package [35].

305 4.3 Public archives

306 With a standardized, versioned data format, Poseidon does not require central infrastructure. Its software tools
307 work independently and users can apply them locally on their own data. But Poseidon was also developed to
308 mitigate the issue of increasingly tedious data preparation for future research projects. To that end, the Poseidon
309 ecosystem includes three public, openly curated archives, which share implementation and infrastructure, but
310 differ in their goals, mode of maintenance, and data content. See the Figures 4 and 5 for an overview of the
311 current data content in the archives. Supplementary Text 8 explains this comparison in more detail.

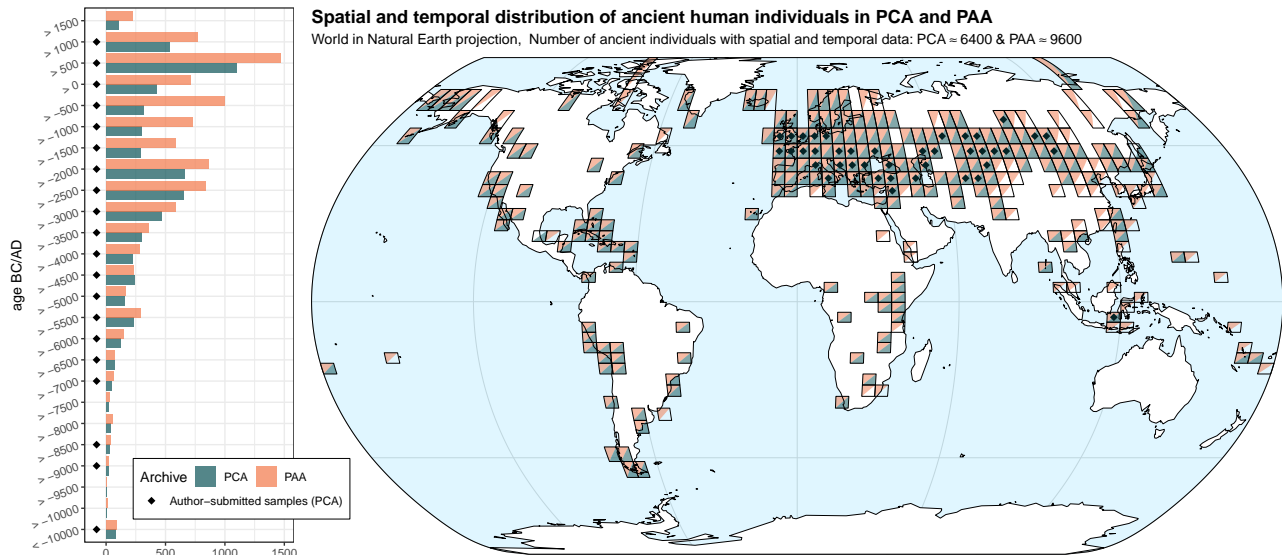


Figure 4: Spatiotemporal distribution of ancient individuals in the PCA and PAA (AADR v54.1.p1) public Poseidon archives. See Supplementary Text 8 for an explanation of how individuals were counted. The map shows the qualitative presence and absence of samples from both archives in a 5°-resolution grid. Especially highlighted are areas and time periods for which the PCA includes samples from currently 13 author-submitted Poseidon packages.

Comparison of PCA and PAA by different metrics

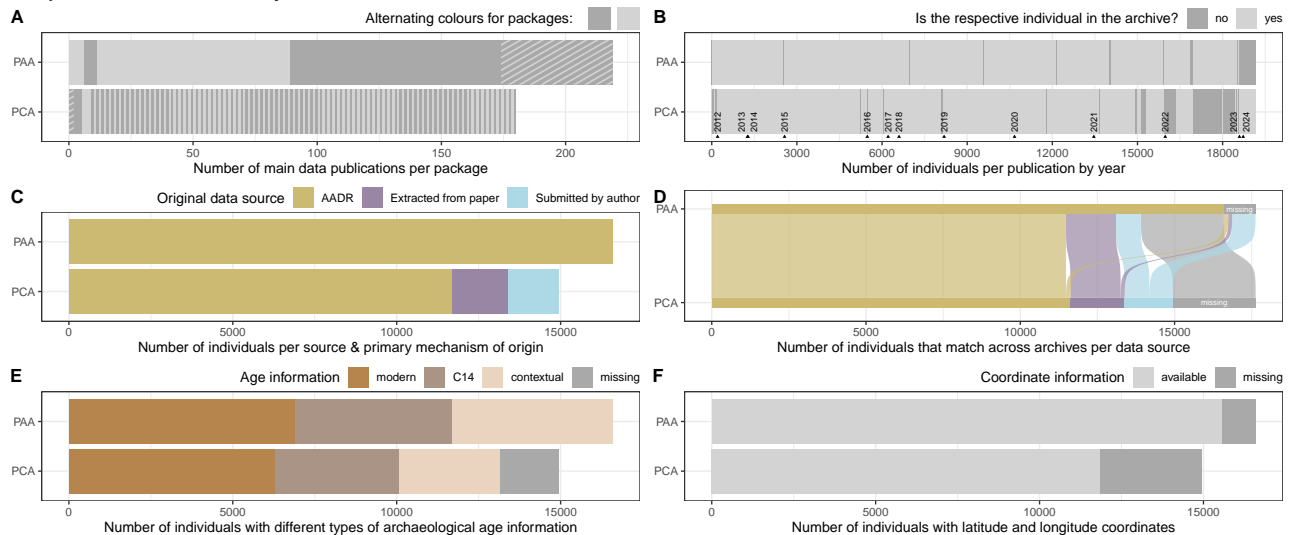


Figure 5: Multiple charts to compare the current content of the PCA and PAA (AADR v54.1.p1) public Poseidon archives. See Supplementary Text 8 for an explanation of how individuals were counted and more in-depth descriptions of each chart. A) Stacked barchart of publications and how they are distributed across packages in PCA and PAA. Publications represented in multiple packages are counted towards the shaded area to get a correct total, B) Barcode plot of individuals available in each archive per publication through time, C) Stacked barchart of individuals and how they were added to the archives, D) Sankey diagram of individuals matching across PCA and PAA, highlighting individuals unique to each archive on the right, E) Stacked barchart of dating information per individual available in the archives, F) Stacked barchart of spatial coordinate coverage in the archives.

312 4.3.1 Technical infrastructure

313 All archives are hosted in dedicated Git repositories on GitHub (e.g. <https://github.com/poseidon-frame>
 314 `work/community-archive` for the PCA), where each individual package is stored in a directory named after
 315 the package. This setup has a number of advantages: Git provides version control down to the individual lines
 316 of each meta- and context data file, Git and GitHub together include co-working features that allow users to
 317 submit new packages and suggest concrete changes to existing ones, and GitHub is a comparatively affordable
 318 host with advanced automation features. Git is by default not suitable for large, binary files, but GitHub offers
 319 large file storage with the Git LFS extension. The archives make use of this for all `.bed` and `.bim` files. Each
 320 change in the public archives is automatically validated using GitHub Actions on their cloud infrastructure,
 321 running a number of scripts on the new state, including `trident validate`, to ensure continuous structural
 322 integrity.

323 As already introduced above for `trident list --remote` and `trident fetch`, the data in the public
 324 archives is not just available on GitHub, but also and more conveniently via a web API provided by an open web
 325 server running `trident serve`. This service is hosted by the scientific IT service provider of the Max Planck
 326 Society (Gesellschaft für wissenschaftliche Datenverarbeitung mbH Göttingen, <https://gwdg.de>). Once per
 327 day the server fetches the latest changes to the archives on GitHub and incorporates new packages and package
 328 versions. It does so through a Git-integrated bookkeeping mechanism using the hidden `trident` subcommands
 329 `chronicle` and `timetravel`. Note that the server also provides outdated package versions (for the community
 330 archive starting from 2023-06-12) to maintain computational reproducibility. Here are the endpoints the server
 331 supports:

- 332 • <https://server.poseidon-adna.org/packages> returns a JSON list of all packages

- 333 • <https://server.poseidon-adna.org/groups> returns a JSON list of all groups
- 334 • <https://server.poseidon-adna.org/individuals> returns a JSON list of all samples/individuals
- 335 • https://server.poseidon-adna.org/zip_file/<package_name> returns a complete zip file of the
336 package with the given name

337 The most important arguments for these are `?archive` to select the archive that should be
338 queried, `?additionalJannoColumns` to add more detailed information to the `/individuals` response and
339 `?package_version` to select a specific package version with `/zip_file`.

340 4.3.2 The Community Archive

341 The Poseidon Community Archive (PCA) stores author-submitted, article-wise Poseidon packages. It focuses
342 on packages prepared by the authors of the respective publication, containing the exact genotype data used for
343 the paper, to ensure a maximum of computational reproducibility. Author submissions are also ideal for the
344 context data in the `.janno` file, because the respective domain-experts are generally most knowledgeable on data
345 quality and the spatiotemporal origin of their samples.

346 For historical reasons the PCA does not only contain author submissions, though. To kickstart the public
347 archive development in 2020, we prepopulated it with packages derived from in-house data and previous versions
348 of the AADR, which have since then been further modified and edited, as is transparent in the version history of
349 these packages. This legacy data will remain in the PCA to maintain established workflows. Authors and other
350 community members can take ownership, update entries if need be, and thus have the possibility to (further)
351 improve the quality of these datasets. A contributing guide on the Poseidon webpage explain the details of
352 how to submit a new paper-associated dataset or suggest changes to an existing one. Each submission passes
353 through a checklist-based review process and is eventually confirmed by the Poseidon core team. Contributors
354 and original and intermediate authors are credited via publication keys and corresponding `.bib` entries, as well
355 as a dedicated `Contributor` list in the package-defining `POSEIDON.yml` file.

356 4.3.3 The AADR Archive

357 The Poseidon AADR Archive (PAA) stores releases of the AADR dataset [13] reworked into Poseidon packages.
358 It thus deviates from the PCA and the PMA in multiple important ways: It is not organized by individual
359 publications, includes the versioning of the original provider on top of our own versioning, and relies to a lesser
360 degree on community contributions. The cleaning and repackaging process is documented in an extra repository
361 (<https://github.com/poseidon-framework/aadr2poseidon>) and mostly has the following goals: i) Creation
362 of a version of the AADR that follows the Poseidon package standard and is thus directly compatible with
363 `trident` and other Poseidon tooling, ii) increasing the machine-readability of the AADR, especially regarding
364 the sample age information, and iii) providing clean `.bib` files with all references of publications in the AADR
365 for convenient citation management.

366 4.3.4 The Minotaur Archive

367 The Poseidon Minotaur Archive (PMA) mirrors the PCA in that it stores publication-wise packages, often the
368 very same as the PCA. However, Packages in the PMA do not rely on author-submitted genotype data, but
369 instead include genotypes consistently reprocessed from raw sequencing data, run through the Minotaur workflow
370 (see below). The motivation for this bioinformatic reprocessing is to generate an internally consistent dataset,
371 which optimises cross-package comparability, rather than per-author reproducibility of individual packages –
372 like the AADR.

373 The submission of packages to the PMA is less direct as for the PCA and involves the preparation of a per-
374 package recipe to parameterize the relevant processing run. Package recipes are archived in a dedicated GitHub
375 repository, where targeted GitHub Actions guide users through the necessary steps to create and submit a new
376 recipe.

377 4.4 The Minotaur workflow

378 The reproducibility of our processing behind the PMA is achieved with a semi-automatic computational workflow
379 to generate Poseidon packages from raw sequencing data: the Minotaur workflow. The entry point for this
380 processing pipeline is a *package recipe*, a collection of files containing all the information required to download,
381 validate, and process the raw reads into a Poseidon package. Each recipe must contain an .ssf file, a .tsv file
382 formatted like a valid input .tsv for nf-core/eager [36], a .config file outlining the nf-core/eager configuration
383 parameters for processing, a .sh script that adapts the .tsv file to the local cluster at the time of processing,
384 and finally a .txt file listing all the versions of scripts used when creating the recipe, to ensure reproducibility.

385 Contributors are able to request packages via GitHub issues on the dedicated minotaur-recipes repository
386 (<https://github.com/poseidon-framework/minotaur-recipes>), and actively prepare package recipes by
387 providing a configured .ssf file. This .ssf file gets validated through GitHub Actions, and then complemented to
388 a full recipe with all the required files. Each recipe submission passes through a checklist-based review process
389 and is finally confirmed by the Poseidon core team. This approach is designed to standardise and streamline the
390 processing through the Minotaur workflow, while still allowing enough flexibility in its configuration to account
391 for the heterogeneity present in raw sequencing data.

392 The minotaur-recipes repository is mirrored in the computational cluster of MPI-EVA, where the actual
393 processing takes place. In the future, we may outsource this processing to a publicly accessible cloud service
394 to make it fully independent from a particular institution. The raw data is downloaded there, and processed
395 through nf-core/eager with the parameters specified in the .config file of the package recipe. By default, this
396 processing includes adapter trimming and read-pair collapsing, aligning to the human reference, removal or
397 PCR duplicates, masking of the ends of reads to mitigate aDNA damage artefacts, and genotyping. Changes
398 to processing parameters are permitted, and can be specified, explained and recorded in the package recipe.
399 The genotypes generated from this processing are then turned into a Poseidon package, whose .janno file is
400 populated with descriptive statistics generated during the processing.

401 All the code responsible for this pipeline can be found at [https://github.com/poseidon-framework/pos](https://github.com/poseidon-framework/poseidon-eager)
402 [eidon-eager](https://github.com/poseidon-framework/poseidon-eager). The resulting Poseidon package is finally uploaded to the PMA, where it again undergoes review
403 before being added to the archive. During the review process, missing information for the package is filled in
404 either manually, or pulled from the Community Archive, if appropriate.

405 5 Discussion

406 Archaeogenetic research, as many other fast growing, data-driven fields, is challenged by data heterogeneity, a
407 lack of systematically applied standards, and from difficulties to discover published data. Poseidon addresses
408 these issues on multiple levels: First, Poseidon provides a standardised, yet flexible data format for day-to-day
409 scientific data analysis, large scale automation, and tidy data storage. Our design choice to build Poseidon around
410 the notion of packages, together with integration of bibliography information (the .bib file), emphasises good
411 citation practice, with any downstream merging necessarily resulting in complete bibliography files listing all
412 original source papers. Second, the Poseidon software, such as *trident*, to discover, validate, update, merge, subset
413 and analyse such packages, complements the standard to ease adoption. Third, perhaps the most ambitious

414 component of Poseidon, our public archives, make use of the standard and its versioning feature to host published
415 data and make it findable and transparently maintainable via GitHub community features.

416 This multi-layer architecture, with loosely coupled components, allows for a variety of adoption paths or
417 starting points for users and data analysts: The package format alone can serve as a useful storage format
418 for local work, even without using the software or archives. The software can help to create and work with
419 such packages locally, even without the cloud-features and server-access. Finally, our archives can be seen as a
420 transparent hub to download and discover data, even without our software or adoption of the package format.

421 In light of our developing research field and its position between disciplines, we designed the Poseidon package
422 definition to allow for flexibility, e.g. by allowing arbitrary non-schema columns, as for example used in the PAA
423 to incorporate some fields specific to that data source. Furthermore, we have placed the standard definition itself
424 on GitHub to enable its long-term use as a "living standard", which can be modified and developed further
425 given emerging use-cases and practice from the community, after open discussion and review.

426 While the Poseidon package format combines various standard formats (e.g. YAML, .tsv, EIGENSTRAT)
427 into its own package specification, opportunities exist to integrate it with larger systems in the Linked Open
428 Data (LOD) world [37]. For example, many of the concepts used in our meta-data definitions exist already in
429 public ontologies, which could be more tightly integrated into our format in the future. Specifically, for example,
430 our **Country** field in the .janno file definition could link to entities in Wikidata [38] or other comparable LOD
431 databases. Our decision for a light-weight flat-file setup has given us leeway to adjust the system exactly to our
432 preferences and the requirements of the field, but integration with the Web of Data thus remains an open tasks
433 for the future. To find partners to establish this uplink in future versions Poseidon is part of the NFDI4Objects
434 initiative (<https://www.nfdi4objects.net>).

435 Regarding infrastructure, Poseidon is currently very much dependent on GitHub. Following the example
436 of other research data standards [39], all code, data, and issue-tracking is stored there, relying extensively on
437 GitHub's CI system ('GitHub Actions') for automatic code compilation and data validation. The lock-in into
438 GitHub's proprietary platform is slightly mitigated by the open Git format used for code- and data storage,
439 but this is still a serious dependency, including factual drawbacks like a strict bandwidth limitation for large
440 file data downloads.

441 Beyond these technical questions, finally, Poseidon is also exposed to some of the broader social challenges
442 of scientific data management: Guiding the growth of a healthy community of developers, contributors and
443 maintainers is not trivial. Poseidon currently depends on a small core team consisting of the authors of this
444 paper. A growing number of active collaborators will require committing to a suitable governance system [40].
445 The work of the core team is currently funded by their employing institution, which also provides computational
446 infrastructure and computing hours. If this commitment gets reduced in the middle- to long-term future, funding
447 may become an increasingly pressing issue – a challenge shared by many research data management projects
448 [41].

449 Critically the long-term success of Poseidon depends on scientists and generally practitioners in the field
450 of archaeogenetics to reroute resources, not least time, into its development and maintenance, if they consider
451 it valuable for their research and publications. With emerging initiatives like the HAAM Community (<https://haam-community.github.io>) and various working groups recently embracing Poseidon-based workflows and
452 first papers referencing it explicitly ([42, 43] and other forthcoming work) we indeed see positive signals towards
453 wider adoption. Regardless of whether this development subsists and a community forms around Poseidon, the
454 Poseidon data format and the software developed for it will remain permanently and openly available for future
455 reference.
456

457 6 Acknowledgements

458 This project has received funding from the Department of Archaeogenetics at the Max Planck Institute for
459 Evolutionary Anthropology (MPI-EVA), the International Max Planck Research School for the Science of Hu-
460 man History at the Max Planck Institute for Geoanthropology (MPI-GEA), NFDI4Objects, so the the National
461 research data infrastructure initiative by the German National Science foundation (DFG), and the European
462 Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant
463 agreement number 851511). The data processing with the Minotaur workflow heavily relies on computational
464 facilities of MPI-EVA. We gratefully acknowledge insightful discussions with many current and former members
465 of our department at MPI-EVA, most notably Selina Carlhoff, Luca Traverso, and Harald Ringbauer. Special
466 thanks go to Michelle O’Reilly (MPI-GEA) who designed the Poseidon logo, specified the main colour palette
467 for the website and revised the schematic overview Figures 1 and 2.

468 Supplementary Texts

- 469 • Supplementary Text 1: Poseidon package specification v2.7.1
- 470 • Supplementary Text 2: .janno file details
- 471 • Supplementary Text 3: Guide for *trident* v1.4.1.0
- 472 • Supplementary Text 4: Guide for *xerxes* v1.0.1.0
- 473 • Supplementary Text 5: *xerxes* theoretical background
- 474 • Supplementary Text 6: Guide for *qjanno* v1.0.0.0
- 475 • Supplementary Text 7: Guide for the *janno* R package v1.0.0
- 476 • Supplementary Text 8: Comparison of the public archive content

477 References

- 478 [1] Matthias Meyer et al. “A high-coverage genome sequence from an archaic Denisovan individual”. In:
479 *Science* 338.6104 (2012), pp. 222–226. DOI: 10.1126/science.1224344.
- 480 [2] Kay Prüfer et al. “The complete genome sequence of a Neanderthal from the Altai Mountains”. In: *Nature*
481 505.7481 (2014), pp. 43–49. DOI: 10.1038/nature12886.
- 482 [3] Iosif Lazaridis et al. “Ancient human genomes suggest three ancestral populations for present-day Euro-
483 peans”. In: *Nature* 513.7518 (2014), pp. 409–413. DOI: 10.1038/nature13673.
- 484 [4] Wolfgang Haak et al. “Massive migration from the steppe was a source for Indo-European languages in
485 Europe”. In: *Nature* 522.7555 (2015), pp. 207–211. DOI: 10.1038/nature14317.
- 486 [5] Ashot Margaryan et al. “Population genomics of the Viking world”. In: *Nature* 585.7825 (2020), pp. 390–
487 396. DOI: 10.1038/s41586-020-2688-8.
- 488 [6] Joscha Gretzinger et al. “The Anglo-Saxon migration and the formation of the early English gene pool”.
489 In: *Nature* (2022), pp. 1–8. DOI: 10.1038/s41586-022-05247-2.
- 490 [7] Morten Rasmussen et al. “Ancient human genome sequence of an extinct Palaeo-Eskimo”. In: *Nature*
491 463.7282 (2010), pp. 757–762. DOI: 10.1038/nature08835.

- 492 [8] Qiaomei Fu et al. “DNA analysis of an early modern human from Tianyuan Cave, China”. In: *Proceedings*
493 *of the National Academy of Sciences of the United States of America* 110.6 (2013), pp. 2223–2227. DOI:
494 10.1073/pnas.1221359110.
- 495 [9] Nadin Rohland et al. “Three assays for in-solution enrichment of ancient human DNA at more than a
496 million SNPs”. In: *Genome research* 32.11-12 (2022), pp. 2068–2078. DOI: 10.1101/gr.276728.122.
- 497 [10] Ewen Callaway. “‘Truly gobsmacked’: Ancient-human genome count surpasses 10, 000”. In: *Nature* 617.7959
498 (2023), pp. 20–20. DOI: 10.1038/d41586-023-01403-4.
- 499 [11] Kenneth Katz et al. “The Sequence Read Archive: a decade more of explosive growth”. In: *Nucleic Acids*
500 *Research* 50.D1 (2021), pp. D387–D390. DOI: 10.1093/nar/gkab1053.
- 501 [12] Josephine Burgin et al. “The European Nucleotide Archive in 2022”. In: *Nucleic Acids Research* 51.D1
502 (2022), pp. D121–D125. DOI: 10.1093/nar/gkac1051.
- 503 [13] Swapan Mallick et al. “The Allen Ancient DNA Resource (AADR) a curated compendium of ancient
504 human genomes”. In: *Scientific Data* 11.1 (2024). DOI: 10.1038/s41597-024-03031-7.
- 505 [14] Elgun Jabrayilzade et al. “Bus factor in practice”. In: *Proceedings of the 44th International Conference on*
506 *Software Engineering: Software Engineering in Practice*. ACM, 2022. DOI: 10.1145/3510457.3513082.
- 507 [15] Mark D. Wilkinson et al. “The FAIR Guiding Principles for scientific data management and stewardship”.
508 In: *Scientific Data* 3.1 (2016). DOI: 10.1038/sdata.2016.18.
- 509 [16] Alkes L Price et al. “Principal components analysis corrects for stratification in genome-wide association
510 studies”. In: *Nature Genetics* 38.8 (2006), pp. 904–909. DOI: 10.1038/ng1847.
- 511 [17] Nick Patterson, Alkes L. Price, and David Reich. “Population Structure and Eigenanalysis”. In: *PLoS*
512 *Genetics* 2.12 (2006), e190. DOI: 10.1371/journal.pgen.0020190.
- 513 [18] David H. Alexander, John Novembre, and Kenneth Lange. “Fast model-based estimation of ancestry in
514 unrelated individuals”. In: *Genome Research* 19.9 (2009), pp. 1655–1664. DOI: 10.1101/gr.094052.109.
- 515 [19] Iosif Lazaridis et al. “Genetic origins of the Minoans and Mycenaeans”. In: *Nature* 548.7666 (2017),
516 pp. 214–218. DOI: 10.1038/nature23310.
- 517 [20] Simon Marlow, Simon Peyton Jones, et al. *Haskell 2010 Language Report*. Tech. rep. [Online; accessed
518 2024-03-05]. Haskell.org, 2010.
- 519 [21] Tom Preston-Werner. *Semantic Versioning*. web. [Online; accessed 2024-03-19]. 2013.
- 520 [22] Shaun Purcell et al. “PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage
521 Analyses”. In: *The American Journal of Human Genetics* 81.3 (2007), pp. 559–575. DOI: 10.1086/519795.
- 522 [23] Petr Danecek et al. “The variant call format and VCFtools”. In: *Bioinformatics* 27.15 (2011), pp. 2156–
523 2158. DOI: 10.1093/bioinformatics/btr330.
- 524 [24] Nick Patterson et al. “Ancient Admixture in Human History”. In: *Genetics* 192.3 (2012), pp. 1065–1093.
525 DOI: 10.1534/genetics.112.145037.
- 526 [25] Iain Mathieson et al. “Genome-wide patterns of selection in 230 ancient Eurasians”. In: *Nature* 528.7583
527 (2015), pp. 499–503. DOI: 10.1038/nature16152.
- 528 [26] Isaac Jones. “The Haskell cabal: A common architecture for building applications and libraries”. In: *6th*
529 *Symposium on Trends in Functional Programming*. 2005, pp. 340–354.
- 530 [27] Stefania Loredana Nita and Marius Mihailescu. “Haskell Stack”. In: *Haskell Quick Syntax Reference*.
531 Apress, 2019, pp. 165–171. DOI: 10.1007/978-1-4842-4507-1_23.

- 532 [28] Björn Grüning et al. “Bioconda: sustainable and comprehensive software distribution for the life sciences”.
533 In: *Nature Methods* 15.7 (2018), pp. 475–476. DOI: 10.1038/s41592-018-0046-7.
- 534 [29] G Bhatia et al. “Estimating and interpreting FST: The impact of rare variants”. In: *Genome research*
535 23.9 (2013), pp. 1514–1521. DOI: 10.1101/gr.154831.113.
- 536 [30] Frank M T A Busing, Erik Meijer, and Rien Van Der Leeden. “Delete-m Jackknife for Unequal m”. In:
537 *Statistics and computing* 9.1 (1999), pp. 3–8. DOI: 10.1023/A:1008800423698.
- 538 [31] Theseas C. Lamnidis et al. “Ancient Fennoscandian genomes reveal origin and spread of Siberian ancestry
539 in Europe”. In: *Nature Communications* 9.1 (2018). DOI: 10.1038/s41467-018-07483-5.
- 540 [32] 1000 Genomes Project Consortium et al. “A global reference for human genetic variation”. In: *Nature*
541 526.7571 (2015), pp. 68–74. DOI: 10.1038/nature15393.
- 542 [33] Kevin P. Gaffney et al. “SQLite: Past, Present, and Future”. In: *Proc. VLDB Endow.* 15.12 (2022),
543 pp. 3535–3547. DOI: 10.14778/3554821.3554842.
- 544 [34] Hadley Wickham et al. “Welcome to the Tidyverse”. In: *Journal of Open Source Software* 4.43 (2019),
545 p. 1686. DOI: 10.21105/joss.01686.
- 546 [35] John Haslett and Andrew Parnell. “A Simple Monotone Process with Application to Radiocarbon-Dated
547 Depth Chronologies”. In: *Journal of the Royal Statistical Society Series C: Applied Statistics* 57.4 (2008),
548 pp. 399–418. DOI: 10.1111/j.1467-9876.2008.00623.x.
- 549 [36] James A. Fellows Yates et al. “Reproducible, portable, and efficient ancient genome reconstruction with
550 nf-core/eager”. In: *PeerJ* 9 (2021), e10947. DOI: 10.7717/peerj.10947.
- 551 [37] Christian Bizer, Tom Heath, and Tim Berners-Lee. “Linked Data - The Story So Far”. In: *Linking the*
552 *World’s Information*. ACM, 2023, pp. 115–143. DOI: 10.1145/3591366.3591378.
- 553 [38] Denny Vrandečić and Markus Krötzsch. “Wikidata: a free collaborative knowledgebase”. In: *Communi-*
554 *cations of the ACM* 57.10 (2014), pp. 78–85. DOI: 10.1145/2629489.
- 555 [39] Robert Crystal-Ornelas et al. “A Guide to Using GitHub for Developing and Versioning Data Standards
556 and Reporting Formats”. In: *Earth and Space Science* 8.8 (2021). DOI: 10.1029/2021ea001797.
- 557 [40] Dany Di Tullio and D. Sandy Staples. “The Governance and Control of Open Source Software Projects”. In:
558 *Journal of Management Information Systems* 30.3 (2013), pp. 49–80. DOI: 10.2753/mis0742-1222300303.
- 559 [41] Margreet Bloemers and Annalisa Montesanti. “The FAIR Funding Model: Providing a Framework for
560 Research Funders to Drive the Transition toward FAIR Data Management and Stewardship Practices”.
561 In: *Data Intelligence* 2.1–2 (2020), pp. 171–180. DOI: 10.1162/dint_a_00039.
- 562 [42] Sanni Peltola et al. “Genetic admixture and language shift in the medieval Volga-Oka interfluvium”. In:
563 *Current Biology* 33.1 (2023), 174–182.e10. DOI: 10.1016/j.cub.2022.11.036.
- 564 [43] Selina Carlhoff et al. “Genomic portrait and relatedness patterns of the Iron Age Log Coffin culture in
565 northwestern Thailand”. In: *Nature Communications* 14.1 (2023). DOI: 10.1038/s41467-023-44328-2.