

Contents

0.1	Guide for trident v1.4.0.2 to v1.4.0.3	1
0.1.1	The trident CLI	1
0.1.2	Package creation and manipulation commands	4
0.1.3	Inspection commands	15

0.1 Guide for trident v1.4.0.2 to v1.4.0.3

0.1.1 The trident CLI

Trident is a command line software tool structured in multiple subcommands. If you installed it properly you can call it on the command line by typing `trident`. This will show an overview of the general options and all subcommands, which are explained in detail below.

```
Usage: trident [--version] [--logMode MODE | --debug] [--errLength INT]
        [--inPlinkPopName MODE] (COMMAND | COMMAND)
```

`trident` is a management and analysis tool for Poseidon packages. Report issues here: <https://github.com/poseidon-framework/poseidon-hs/issues>

Available options:

<code>-h, --help</code>	Show this help text
<code>--version</code>	Show version number
<code>--logMode MODE</code>	How information should be reported: NoLog, SimpleLog, DefaultLog, ServerLog or VerboseLog. (default: DefaultLog)
<code>--debug</code>	Short for <code>--logMode VerboseLog</code> .
<code>--errLength INT</code>	After how many characters should a potential error message be truncated. "Inf" for no truncation. (default: CharCount 1500)
<code>--inPlinkPopName MODE</code>	Where to read the population/group name from the FAM file in Plink-format. Three options are possible: asFamily (default) asPhenotype asBoth.

Package creation and manipulation commands:

<code>init</code>	Create a new Poseidon package from genotype data
<code>fetch</code>	Download data from a remote Poseidon repository
<code>forge</code>	Select packages, groups or individuals and create a new Poseidon package from them
<code>genoconvert</code>	Convert the genotype data in a Poseidon package to a different file format
<code>rectify</code>	Adjust POSEIDON.yml files automatically to package changes

Inspection commands:

<code>list</code>	List packages, groups or individuals from local or
-------------------	----------------------------------------------------

```

43         remote Poseidon repositories
44     summarise      Get an overview over the content of one or multiple
45                    Poseidon packages
46     survey         Survey the degree of context information completeness
47                    for Poseidon packages
48     validate       Check Poseidon packages or package components for
49                    structural correctness

```

50 Trident allows to work directly with genotype data (see `-p` below), but its optimized for the interaction with
51 [Poseidon packages](#), which wrap and contextualize the data. Most trident subcommands therefore have a central
52 parameter, called `--baseDir` or simply `-d` to specify one or more base directories to look for packages. For example,
53 if all Poseidon packages live inside a repository at `/path/to/poseidon/packages` you would simply say `trident`
54 `<subcommand> -d /path/to/poseidon/dirs/` and `trident` would automatically search all subdirectories inside
55 of the repository for valid Poseidon packages (as identified by valid `POSEIDON.yml` files).

56 You can arrange a Poseidon repository in a hierarchical way. For example:

```

57 /path/to/poseidon/packages
58     /modern
59         /2019_poseidon_package1
60         /2019_poseidon_package2
61     /ancient
62         /...
63         /...
64     /Reference_Genomes
65         /...
66         /...

```

67 You can use this structure to select only the level of packages you're interested in, even individual ones, and you
68 can make use of the fact that `-d` can be given multiple times.

69 Being able to specify one or multiple repositories is often not enough, as you may have your own data to
70 co-analyse with the main repository. This is easy to do, as you simply need to provide your own genotype data as
71 yet another Poseidon package to be added to your `trident` command. For example, let's say you have genotype
72 data in `EIGENSTRAT` format (`trident` supports `EIGENSTRAT` and `PLINK` as formats.):

```

73 ~/my_project/my_project.geno
74 ~/my_project/my_project.snp
75 ~/my_project/my_project.ind

```

76 then you can make that to a skeleton Poseidon package with the `init` command. You can also do it manually by
77 simply adding a `POSEIDON.yml` file, with for example the following content:

```

78 poseidonVersion: 2.7.1
79 title: My_awesome_project
80 description: Unpublished genetic data from my awesome project
81 contributor:
82     - name: Stephan Schiffels
83       email: schiffels@institute.org
84 packageVersion: 0.1.0

```

```

85  lastModified: 2020-10-07
86  genotypeData:
87    format: EIGENSTRAT
88    genoFile: my_project.geno
89    snpFile: my_project.snp
90    indFile: my_project.ind
91    jannoFile: my_project.janno
92    bibFile: sources.bib

```

Two remarks: 1) all file paths are considered *relative* to the directory in which POSEIDON.yml resides. For this example we assume that this file is added into the same directory as the three genotype files. 2) Besides the genotype data files there are two (technically optional) files referenced by this example POSEIDON.yml file: **sources.bib** and **my_project.janno**. Of course you can add them manually - **init** automatically creates empty dummy versions.

Once you have set up your own Poseidon package (which is really only a skeleton so far), you can add it to your **trident** analysis, by simply adding your project directory to the command using **-d**, for example:

```

100 trident list -d /path/to/poseidon/packages/modern \
101   -d /path/to/poseidon/packages/ReferenceGenomes
102   -d ~/my_project --packages

```

103 0.1.1.1 General notes

104 **0.1.1.1.1 Logging and command line output** For all subcommands the general argument **--logMode**
 105 defines how trident reports messages (to stderr) on the command line:

- 106 • *NoLog*: Hides all messages.
- 107 • *SimpleLog*: Plain and simple output to stderr.
- 108 • *DefaultLog*: Adds severity indicators before each message. (default setting)
- 109 • *ServerLog*: Additionally adds timestamps before each message.
- 110 • *VerboseLog*: Shows not just messages on the log levels **Info**, **Warning** and **Error** like the other modes, but
 111 also on the more verbose level **Debug**. Use this for debugging.

112 **--debug** is short for **--logMode VerboseLog** to activate this important log level more easily.

113 0.1.1.1.2 Package duplicates and versions

- 114 • For **trident** multiple packages in a set of base directories can share the same **title**, if they have different
 115 **packageVersion** numbers. If the version numbers are identical or missing, then **trident** stops with an
 116 exception.
- 117 • The **trident** subcommands **genoconvert**, **list**, **rectify**, **survey** and **validate** by default consider all
 118 versions of each Poseidon package in the given base directories. The **--onlyLatest** flag causes them to
 119 instead only consider the latest versions.
- 120 • **fetch** and **forge** generally consider all package versions and their selection language (see below) allows
 121 for detailed version handling.
- 122 • **summarize** always only shows results for the latest package versions.

123 0.1.1.1.3 Individual/sample duplicates

- Individual/sample names (Poseidon_IDs) within one package have to be unique, or trident will stop.
- We also discourage sample duplicates across packages in package repositories, but trident will generally continue with them. `validate` will fail though, if the `--ignoreDuplicates` flag is not set.
- `forge` offers a special mechanism to resolve sample duplicates within its selection language.

0.1.1.1.4 Group names in .fam files The `.fam` file of Plink-formatted genotype data is used inconsistently across different popular aDNA software tools to store group/population name information. The (global) option `--inPlinkPopName` with the arguments `asFamily` (default), `asPhenotype` and `asBoth` allows to control the reading of the population name from Plink `.fam` files. The subcommands that write genotype data (`forge`, `genoconvert`) have a corresponding option `--outPlinkPopName` to specify this for the output.

0.1.1.1.5 Whitespaces in the .janno file While reading the `.janno` file `trident` trims all leading and trailing whitespaces around individual cells. Also all instances of the `No-Break Space` unicode character will be removed. This means these whitespaces will not be preserved when a package is `forged`.

0.1.2 Package creation and manipulation commands

0.1.2.1 Init command `init` creates a new, valid Poseidon package from genotype data files. It adds a valid `POSEIDON.yml` file, a dummy `.janno` file for context information and an empty `.bib` file for literature references.

[Click here for command line details](#)

```
Usage: trident init ((-p|--genoOne FILE) | --inFormat FORMAT --genoFile FILE
                  --snpFile FILE --indFile FILE) [--snpSet SET]
                  (-o|--outPackagePath DIR) [-n|--outPackageName STRING]
                  [--minimal]
```

Create a new Poseidon package from genotype data

Available options:

<code>-h,--help</code>	Show this help text
<code>-p,--genoOne FILE</code>	One of the input genotype data files. Expects <code>.bed</code> , <code>.bim</code> or <code>.fam</code> for PLINK and <code>.geno</code> , <code>.snv</code> or <code>.ind</code> for EIGENSTRAT. The other files must be in the same directory and must have the same base name.
<code>--inFormat FORMAT</code>	The format of the input genotype data: EIGENSTRAT or PLINK. Only necessary for data input with <code>--genoFile</code> + <code>--snpFile</code> + <code>--indFile</code> .
<code>--genoFile FILE</code>	Path to the input geno file.
<code>--snpFile FILE</code>	Path to the input snv file.
<code>--indFile FILE</code>	Path to the input ind file.
<code>--snpSet SET</code>	The <code>snpSet</code> of the package: 1240K, HumanOrigins or Other. Only relevant for data input with <code>-p --genoOne</code> or <code>--genoFile</code> + <code>--snpFile</code> + <code>--indFile</code> , because the packages in a <code>-d --baseDir</code> already have this information in their respective <code>POSEIDON.yml</code> files. (default: Other)

```

165 -o,--outPackagePath DIR Path to the output package directory.
166 -n,--outPackageName STRING
167                               The output package name. This is optional: If no name
168                               is provided, then the package name defaults to the
169                               basename of the (mandatory) --outPackagePath
170                               argument. (default: Nothing)
171 --minimal                     Should the output data be reduced to a necessary
172                               minimum and omit empty scaffolding?
173
174 The command
175
176 trident init \
177   --inFormat EIGENSTRAT/PLINK \
178   --genoFile path/to/geno_file \
179   --snpFile path/to/snp_file \
180   --indFile path/to/ind_file \
181   --snpSet 1240K|HumanOrigins|Other \
182   -o path/to/new_package_name
183
184 requires the format (--inFormat) of your input data (either EIGENSTRAT or PLINK), the paths to the respective
185 files (--genoFile, --snpFile, --indFile), and optionally the “shape” of these files (--snpSet), so if they cover
186 the 1240K, the HumanOrigins or an Other SNP set. A simpler interface is available with -p (+ --snpSet).

```

	EIGENSTRAT	PLINK
genoFile	.geno	.bed
snpFile	.snp	.bim
indFile	.ind	.fam

```

184 The output package of init is created as a new directory -o, which should not already exist, and gets the
185 package title corresponding to the basename of -o. You can also set the title explicitly with -n. The --minimal
186 flag causes init to create a minimal package with a very basic POSEIDON.yml and no .bib and .janno files.

```

```

187 0.1.2.2 Fetch command fetch allows to download Poseidon packages from a remote Poseidon server via a
188 Web API. Read more about the data available with it here.

```

```

189 Click here for command line details

```

```

190 Usage: trident fetch (-d|--baseDir DIR)
191           (--downloadAll |
192           (--fetchFile FILE | (-f|--fetchString DSL)))
193           [--remoteURL URL] [--archive STRING]
194

```

```

195 Download data from a remote Poseidon repository
196

```

```

197 Available options:

```

```

198 -h,--help           Show this help text
199 -d,--baseDir DIR    A base directory to search for Poseidon packages.
200 --downloadAll       Download all packages the server is offering.

```

201 --fetchFile FILE A file with a list of packages. Works just as -f, but
202 multiple values can also be separated by newline, not
203 just by comma. -f and --fetchFile can be combined.

204 -f,--fetchString DSL List of packages to be downloaded from the remote
205 server. Package names should be wrapped in asterisks:
206 *package_title*. You can combine multiple values with
207 comma, so for example: "*package_1*, *package_2*,
208 *package_3*". fetchString uses the same parser as
209 forgeString, but does not allow excludes. If groups
210 or individuals are specified, then packages which
211 include these groups or individuals are included in
212 the download.

213 --remoteURL URL URL of the remote Poseidon server.
214 (default: "https://server.poseidon-adna.org")

215 --archive STRING The name of the Poseidon package archive that should
216 be queried. If not given, then the query falls back
217 to the default archive of the server selected with
218 --remoteURL. See the archive documentation at
219 https://www.poseidon-adna.org/#/archive_overview for
220 a list of archives currently available from the
221 official Poseidon Web API. (default: Nothing)

222 It works with

223 trident fetch -d ... -d ... \
224 -f "*package_title_1*,*package_title_2-1.0.1*,group_name,<individual1>"

225 and the entities you want to download must be listed either in a simple string of comma-separated values, which
226 can be passed via -f/--fetchString, or in a text file (--fetchFile). Entities are then combined from these
227 sources.

228 Entities are specified using a special syntax (see also the documentation of `forge` below): packages are wrapped
229 in asterisks, with or without version appended after a dash (e.g. `*package_title*` or `*package_title-1.2.3*`),
230 group names are spelled as is, and individual names are wrapped in angular brackets (e.g. `<individual1>`).
231 Fetch will figure out which packages need to be downloaded to include all specified entities. `--downloadAll`,
232 which can be given instead of -f and --fetchFile, causes fetch to download all packages from the server. The
233 downloaded packages are added in the first (!) -d directory (which gets created if it doesn't exist), but downloads
234 are only performed if the respective packages are not already present in the latest version in any of the -d dirs.

235 Note that `trident fetch` makes most sense in combination with `trident list --remote`: First one can inspect
236 what is available on the server, then one can create a custom fetch command.

237 `fetch` also has the optional arguments `--remote https://...` to name an alternative Poseidon server and
238 `--archive` to select a Poseidon archive on the server. Here is a list of the [archives available on the official](#)
239 [Poseidon server](#).

240 **0.1.2.3 Forge command** `forge` creates new Poseidon packages by extracting and merging packages,
241 populations and individuals/samples from your Poseidon repositories.

242 Click here for command line details

```

243 Usage: trident forge ((-d|--baseDir DIR) |
244         ((-p|--genoOne FILE) | --inFormat FORMAT --genoFile FILE
245         --snpFile FILE --indFile FILE) [--snpSet SET])
246         [--forgeFile FILE | (-f|--forgeString DSL)]
247         [--selectSnps FILE] [--intersect] [--outFormat FORMAT]
248         [--minimal] [--onlyGeno] (-o|--outPackagePath DIR)
249         [-n|--outPackageName STRING] [--packagewise]
250         [--outPlinkPopName MODE]
251
252 Select packages, groups or individuals and create a new Poseidon package from
253 them
254
255 Available options:
256 -h,--help                Show this help text
257 -d,--baseDir DIR          A base directory to search for Poseidon packages.
258 -p,--genoOne FILE         One of the input genotype data files. Expects .bed,
259                             .bim or .fam for PLINK and .geno, .snp or .ind for
260                             EIGENSTRAT. The other files must be in the same
261                             directory and must have the same base name.
262 --inFormat FORMAT         The format of the input genotype data: EIGENSTRAT or
263                             PLINK. Only necessary for data input with --genoFile
264                             + --snpFile + --indFile.
265 --genoFile FILE           Path to the input geno file.
266 --snpFile FILE            Path to the input snp file.
267 --indFile FILE            Path to the input ind file.
268 --snpSet SET              The snpSet of the package: 1240K, HumanOrigins or
269                             Other. Only relevant for data input with -p|--genoOne
270                             or --genoFile + --snpFile + --indFile, because the
271                             packages in a -d|--baseDir already have this
272                             information in their respective POSEIDON.yml files.
273                             (default: Other)
274 --forgeFile FILE          A file with a list of packages, groups or individual
275                             samples. Works just as -f, but multiple values can
276                             also be separated by newline, not just by comma.
277                             Empty lines are ignored and comments start with "#",
278                             so everything after "#" is ignored in one line.
279                             Multiple instances of -f and --forgeFile can be
280                             given. They will be evaluated according to their
281                             input order on the command line.
282 -f,--forgeString DSL      List of packages, groups or individual samples to be
283                             combined in the output package. Packages follow the
284                             syntax *package_title*, populations/groups are simply
285                             group_id and individuals <individual_id>. You can
286                             combine multiple values with comma, so for example:
287                             "*package_1*, <individual_1>, <individual_2>,"

```

group_1". Duplicates are treated as one entry.

Negative selection is possible by prepending "-" to the entity you want to exclude (e.g. "*package_1*, -<individual_1>, -group_1"). forge will apply excludes and includes in order. If the first entity is negative, then forge will assume you want to merge all individuals in the packages found in the baseDirs (except the ones explicitly excluded) before the exclude entities are applied. An empty forgeString (and no --forgeFile) will therefore merge all available individuals. If there are individuals in your input packages with equal individual id, but different main group or source package, they can be specified with the special syntax "<package:group:individual>".

--selectSnps FILE To extract specific SNPs during this forge operation, provide a Snp file. Can be either Eigenstrat (file ending must be '.snp') or Plink (file ending must be '.bim'). When this option is set, the output package will have exactly the SNPs listed in this file. Any SNP not listed in the file will be excluded. If option '--intersect' is also set, only the SNPs overlapping between the SNP file and the forged packages are output. (default: Nothing)

--intersect Whether to output the intersection of the genotype files to be forged. The default (if this option is not set) is to output the union of all SNPs, with genotypes defined as missing in those packages which do not have a SNP that is present in another package. With this option set, the forged dataset will typically have fewer SNPs, but less missingness.

--outFormat FORMAT The format of the output genotype data: EIGENSTRAT or PLINK. (default: PLINK)

--minimal Should the output data be reduced to a necessary minimum and omit empty scaffolding?

--onlyGeno Should only the resulting genotype data be returned? This means the output will not be a Poseidon package.

-o,--outPackagePath DIR Path to the output package directory.

-n,--outPackageName STRING The output package name. This is optional: If no name is provided, then the package name defaults to the basename of the (mandatory) --outPackagePath argument. (default: Nothing)

--packagewise Skip the within-package selection step in forge. This will result in outputting all individuals in the


```

333         relevant packages, and hence a superset of the
334         requested individuals/groups. It may result in better
335         performance in cases where one wants to forge entire
336         packages or almost entire packages. Details: Forge
337         conceptually performs two types of selection: First,
338         it identifies which packages in the supplied base
339         directories are relevant to the requested forge, i.e.
340         whether they are either explicitly listed using
341         *PackageName*, or because they contain selected
342         individuals or groups. Second, within each relevant
343         package, individuals which are not requested are
344         removed. This option skips only the second step, but
345         still performs the first.
346     --outPlinkPopName MODE   Where to write the population/group name into the FAM
347                             file in Plink-format. Three options are possible:
348                             asFamily (default) | asPhenotype | asBoth. See also
349                             --inPlinkPopName.
350
351     forge can be used with
352
353     trident forge -d ... -d ... \
354         -f "*package_name*, group_id, <individual_id>" \
355         -o path/to/new_package_name
356
357     where the entities (packages, groups/populations, individuals/samples) you want in the output package can be
358     denoted either as a string on the command line (-f/--forgeString), or in an input text file (--forgeFile).
359     See the section below for the syntax of this selection language. Do not forget to wrap the --forgeString query
360     in quotes.
361
362     Including one or multiple Poseidon packages with -d is not the only way to include data for a forge operation.
363     It is also possible to consider unpackaged genotype data directly with -p (+ --snpSet) or --inFormat +
364     --genoFile + --snpFile + --indFile (+ --snpSet). This makes the following example possible, where we
365     merge data from one Poseidon package and two genotype datasets to get a new EIGENSTRAT dataset.
366
367     trident forge \
368         -d 2017_GonzalesFortesCurrentBiology \
369         -p 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \
370         --inFormat PLINK \
371         --genoFile 2017_HaberAJHG/2017_HaberAJHG.bed \
372         --snpFile 2017_HaberAJHG/2017_HaberAJHG.bim \
373         --indFile 2017_HaberAJHG/2017_HaberAJHG.fam \
374         -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \
375         -o testpackage \
376         --outFormat EIGENSTRAT \
377         --onlyGeno

```

0.1.2.3.1 The forge selection language The text in --forgeString, --forgeFile (and with limited syntax also in --fetchString and --fetchFile) are parsed as a domain specific query language that describes

precisely which entities should be compiled in the output package of a given **forge** operation. The language has multiple syntactic elements and a specific evaluation logic.

In general a **--forgeString** query consists of multiple entities, separated by **,**. The main entities are Poseidon packages, groups/populations and individuals/samples:

- Each package title is surrounded by *****: ***package***. That means if you want all individuals of the Poseidon package **2019_Jeong_InnerEurasia** in the output package you would add ***2019_Jeong_InnerEurasia*** to the query.
- Groups/populations are not specially marked: **group**. So to get all individuals of the group **Swiss_Roman_period**, you would simply add **Swiss_Roman_period**.
- Individuals/samples are surrounded by **<** and **>**: **<individual>**. **ALA026** therefore becomes **<ALA026>**. A second way to denote individuals is with the more verbose and specific syntax **<package:group:individual>**. Such defined individuals take precedence over differently defined ones (so: directly with **<individual>** or as a subset of ***package*** or **group**). This allows to resolve duplication issues precisely – at least in cases where the duplicated individuals differ in source package or primary group.
- Package versions can be appended to package names, such as ***package-1.2.3***, or **<package-1.2.3:group:individual>**.

In the **--forgeFile** each line is treated as a separate **forgeString**, empty lines are ignored and **#**s start comments. So this is a valid example of a **forgeFile**:

```
# Packages
*package1*, *package2-1.2.3*

# Groups and individuals from other packages beyond package1 and package2
group1, <individual1>, group2, <individual2>, <pac1:group2:individual3>

# group2 has two outlier individuals that should be ignored
-<individual1> # This one has very low coverage
-<pac2:group3:individual4> # This one is from a different time period
```

By prepending **-** to entities, we can exclude them from the forged package (this feature is not available for **fetch**). **forge** figures out the final list of samples to include by executing all **forge**-entities in order. So an entity list ***PackageA*,-<Individual1>,GroupA** may result in a different outcome than ***PackageA*,GroupA,-<Individual1>**, depending on whether **<Individual1>** belongs to **GroupA** or not.

If the **forge** entity list starts with a negative entity, or if the entity list is empty, **forge** will implicitly assume you want to include all individuals in all **latest** versions of packages found in the base directories (except the ones explicitly excluded, of course).

The specific semantics of the various ways to include or exclude entities are:

Inclusion queries

- ***Pac1***: Select all individuals in the latest version of package “Pac1”
- ***Pac1-1.0.1***: Select all individuals in package “Pac1” with version “1.0.1”
- **Group1**: Select all individuals associated with “Group1” in all latest versions of all packages
- **<Ind1>**: Select the individual named “Ind1”, searching in all latest packages.
- **<Pac1:Group1:Ind1>**: Select the individual named “Ind1” associated with “Group1” in the latest version of package “Pac1”

- **<Pac1-1.0.1:Group1:Ind1>**: Select the individual named “Ind1” associated with “Group1” in the package “Pac1” with version “1.0.1”

Exclusion queries

- **-*Pac1***: Remove all individuals in all versions of package “Pac1”
- **-*Pac1-1.0.1***: Remove only individuals in package “Pac1” with version “1.0.1” (but leave other versions in)
- **-Group1**: Remove all individuals associated with “Group1” in all versions of all packages (not just the latest)
- **-<Ind1>**: Remove all individuals named “Ind1” in all versions of all packages (not just the latest).
- **-<Pac1:Group1:Ind1>**: Remove the individual named “Ind1” associated with “Group1”, searching in all versions of package “Pac1”
- **-<Pac1-1.0.1:Group1:Ind1>**: Remove the individual named “Ind1” associated with “Group1”, but only if they are in “Pac1” with version “1.0.1”

If a query results in multiple individuals with the same name, forge will throw an error.

0.1.2.3.2 Treatment of the .janno file while merging `forge` merges and subsets .janno files along with the genotype data. If a package lacks a .janno file, then a basic one will be created internally based on the information in the genotype data, and used for the output. Missing columns across packages will be filled with **n/a**.

For merging two .janno files **A** and **B** the following rules apply regarding undefined, arbitrary additional columns:

- If **A** has an additional column which is not in **B** then empty cells in the rows imported from **B** are filled with **n/a**.
- If **A** and **B** share additional columns with identical column name, then they are treated as semantically identical units and merged accordingly.
- In the resulting .janno file, all additional columns from both **A** and **B** are sorted alphabetically and appended after the normal, specified variables.

The following example illustrates the described behaviour:

A.janno

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn1	AdditionalColumn2
XXX011	POP1	M	A	D
XXX012	POP2	F	B	E
XXX013	POP1	M	C	F

B.janno

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn3	AdditionalColumn2
YYY022	POP5	F	G	J
YYY023	POP5	F	H	K
YYY024	POP5	M	I	L

A.janno + B.janno

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn1	AdditionalColumn2	AdditionalColumn3
XXX011	POP1	M	A	D	n/a
XXX012	POP2	F	B	E	n/a
XXX013	POP1	M	C	F	n/a
YYY022	POP5	F	n/a	J	G
YYY023	POP5	F	n/a	K	H
YYY024	POP5	M	n/a	L	I

0.1.2.3.3 Treatment of the .ssf file while merging The Sequencing Source File (short .ssf file) is forged in exactly the same way as the janno file. SSF files that are present are included in the forge product in the way that the user expects, following selection of those entities which are listed in the **poseidon_IDs** columns of the SSF files. Columns that are only present in some packages, including those not defined by our [Schema] are also included in the forged product in the same way as described for Janno above.

0.1.2.3.4 Treatment of the .bib file while merging In the forge process all relevant samples for the output package are determined. This includes their .janno entries and therefore the information on the publication keys documented for them in the .janno **Publication** column. The output .bib file compiles only the relevant references for the samples in the output package. It includes the references exactly once and is sorted alphabetically (by key).

0.1.2.3.5 Other options Just as for **init** the output package of **forge** is created as a new directory **-o**. The title can also be explicitly defined with **-n**.

--minimal allows for the creation of a minimal output package without **.bib** and **.janno**. This is especially useful for data analysis pipelines, where only the genotype data is required. Even more basic output comes with **--onlyGeno**, which means that only the genotype data is returned without any Poseidon package.

forge has a optional flag **--intersect**, that defines, if the genotype data from different packages should be merged with an **union** or an **intersect** operation. The default (if this option is not set) is to output the union of all SNPs, with genotypes defined as missing in samples from packages which do not have a SNP that is present in another package. With this option set, on the other hand, the forged dataset will typically have fewer SNPs, but less missingness.

--intersect also influences the automatic determination of the **snpSet** field in the POSEIDON.yml file for the resulting package. If the **snpSets** of all input packages are identical, then the resulting package will just inherit this configuration. Otherwise **forge** applies the following pairwise merging logic:

Input snpSet A	Input snpSet B	--intersect	Ouput snpSet
Other	*	*	Other
1240K	HumanOrigins	True	HumanOrigins
1240K	HumanOrigins	False	1240K

--selectSnps allows to provide **forge** with a SNP file in EIGENSTRAT (**.snp**) or PLINK (**.bim**) format to create a package with a specific selection. When this option is set, the output package will have exactly the SNPs listed in this file. Any SNP not listed in the file will be excluded. If **--intersect** is also set, only the SNPs overlapping between the SNP file and the forged packages are output.

472 Merging genotype data across different data sources and file formats is tricky. **forge** is more verbose about
473 potential issues, if the **--logMode** flag is set to **VerboseLog**.

474 The **--onlyGeno** command specifies that only genotype data should be output, not an entire Poseidon package.

475 With **--packagewise** the within-package selection step in **forge** can be skipped. This will result in outputting
476 all individuals in the relevant packages, and hence a superset of the requested individuals/groups. It may result
477 in better performance in cases where one wants to forge entire packages.

478 **0.1.2.4 Genoconvert command** **genoconvert** converts the genotype data in a Poseidon package to a
479 different file format. The respective entries in the **POSEIDON.yml** file are changed accordingly.

480 [Click here for command line details](#)

```
481 Usage: trident genoconvert ((-d|--baseDir DIR) |  
482                             ((-p|--genoOne FILE) | --inFormat FORMAT  
483                             --genoFile FILE --snpFile FILE --indFile FILE)  
484                             [--snpSet SET]) --outFormat FORMAT [--onlyGeno]  
485                             [-o|--outPackagePath DIR] [--removeOld]  
486                             [--outPlinkPopName MODE] [--onlyLatest]
```

488 Convert the genotype data in a Poseidon package to a different file format

490 Available options:

491	-h,--help	Show this help text
492	-d,--baseDir DIR	A base directory to search for Poseidon packages.
493	-p,--genoOne FILE	One of the input genotype data files. Expects .bed, 494 .bim or .fam for PLINK and .geno, .snp or .ind for 495 EIGENSTRAT. The other files must be in the same 496 directory and must have the same base name.
497	--inFormat FORMAT	The format of the input genotype data: EIGENSTRAT or 498 PLINK. Only necessary for data input with --genoFile 499 + --snpFile + --indFile.
500	--genoFile FILE	Path to the input geno file.
501	--snpFile FILE	Path to the input snp file.
502	--indFile FILE	Path to the input ind file.
503	--snpSet SET	The snpSet of the package: 1240K, HumanOrigins or 504 Other. Only relevant for data input with -p --genoOne 505 or --genoFile + --snpFile + --indFile, because the 506 packages in a -d --baseDir already have this 507 information in their respective POSEIDON.yml files. 508 (default: Other)
509	--outFormat FORMAT	the format of the output genotype data: EIGENSTRAT or 510 PLINK.
511	--onlyGeno	Should only the resulting genotype data be returned? 512 This means the output will not be a Poseidon package.
513	-o,--outPackagePath DIR	Path to the output package directory. This is 514 optional: If no path is provided, then the output is

written to the directories where the input genotype data file (.bed/.geno) is stored. (default: Nothing)

`--removeOld` Remove the old genotype files when creating the new ones.

`--outPlinkPopName MODE` Where to write the population/group name into the FAM file in Plink-format. Three options are possible: asFamily (default) | asPhenotype | asBoth. See also `--inPlinkPopName`.

`--onlyLatest` Consider only the latest versions of packages, or the groups and individuals within the latest versions of packages, respectively.

With the default setting

```
trident genoconvert -d ... -d ... --outFormat EIGENSTRAT|PLINK
```

all packages in `-d` will be converted to the desired `--outFormat` (either EIGENSTRAT or PLINK), if the data is not already in this format. This includes updating the respective POSEIDON.yml files.

The “old” data is not deleted, but kept around. That means conversion can result in a package with both PLINK and EIGENSTRAT data, but only one is linked in the POSEIDON.yml file, and that is what will be used by trident. To delete the old data in the conversion you can add the `--removeOld` flag.

Instead of `-d` to change Poseidon packages, the `-p` (+ `--snpSet`) or `--inFormat` + `--genoFile` + `--snpFile` + `--indFile` (+ `--snpSet`) allow to directly convert genotype data that is not wrapped in a Poseidon package and store it to a directory given in `-o`. See this example:

```
trident genoconvert \
  -p 2018_Mittnik_Baltic/Mittnik_Baltic.bed \
  --outFormat EIGENSTRAT
  -o my_directory
```

0.1.2.5 Rectify command

`rectify` automatically harmonizes POSEIDON.yml files of one or multiple packages. This is not an automatic update from one Poseidon version to the next, but rather a clean-up wizard after manual modifications.

[Click here for command line details](#)

```
Usage: trident rectify (-d|--baseDir DIR) [--ignorePoseidonVersion]
      [--poseidonVersion ??.?]
      [--packageVersion VPART [--logText STRING]]
      [--checksumAll | [--checksumGeno] [--checksumJanno]
      [--checksumSSF] [--checksumBib]]
      [--newContributors DSL] [--onlyLatest]
```

Adjust POSEIDON.yml files automatically to package changes

Available options:

<code>-h,--help</code>	Show this help text
<code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.

```

556 --ignorePoseidonVersion Read packages even if their poseidonVersion is not
557 compatible with trident.
558 --poseidonVersion ?.??.? Poseidon version the packages should be updated to:
559 e.g. "2.5.3".
560 --packageVersion VPART Part of the package version number in the
561 POSEIDON.yml file that should be updated: Major,
562 Minor or Patch (see https://semver.org).
563 --logText STRING Log text for this version in the CHANGELOG file.
564 --checksumAll Update all checksums.
565 --checksumGeno Update genotype data checksums.
566 --checksumJanno Update .janno file checksum.
567 --checksumSSF Update .ssf file checksum
568 --checksumBib Update .bib file checksum.
569 --newContributors DSL Contributors to add to the POSEIDON.yml file in the
570 form "[Firstname Lastname](Email address);...".
571 --onlyLatest Consider only the latest versions of packages, or the
572 groups and individuals within the latest versions of
573 packages, respectively.

```

574 It can be called with a lot of optional arguments:

```

575 trident rectify -d ... -d ... \
576 --poseidonVersion "X.X.X" \
577 --packageVersion Major|Minor|Patch \
578 --logText "short description of the update"
579 --checksumAll
580 --newContributors "[Firstname Lastname](Email address);..."

```

581 These arguments determine which fields of the POSEIDON.yml file should be modified.

- 582 • `--poseidonVersion` allows a simple change of the `poseidonVersion` field in the POSEIDON.yml file.
- 583 • `--packageVersion` increments the package version number in the first, the second or the third position.
584 It can optionally be called with `--logText`, which appends an entry to the CHANGELOG file for the
585 respective package version update. `--logText` also creates a new CHANGELOG file if it does not exist
586 yet.
- 587 • `--checksumGeno`, `--checksumJanno`, `--checksumSSF` and `--checksumBib` add or modify the respective
588 checksum fields in the POSEIDON.yml file. `--checksumAll` is a wrapper to call all of them at once.
- 589 • `--newContributors` adds new contributors.

590 :warning: As `rectify` reads and rewrites POSEIDON.yml files, it may change their inner order, layout or
591 even content (e.g. if they have fields which are not in the **POSEIDON.yml definition**). Create a backup of the
592 POSEIDON.yml file before running `rectify` if you are uncertain if this might affect you negatively.

593 0.1.3 Inspection commands

594 **0.1.3.1 List command** `list` lists packages, groups and individuals of the datasets you use, or of the
595 packages available on the server.

596 [Click here for command line details](#)

```

597 Usage: trident list ((-d|--baseDir DIR) | --remote [--remoteURL URL]
598                 [--archive STRING])
599                 (--packages | --groups | --individuals
600                 [-j|--jannoColumn COLNAME]) [--raw] [--onlyLatest]
601
602 List packages, groups or individuals from local or remote Poseidon
603 repositories
604
605 Available options:
606  -h,--help                Show this help text
607  -d,--baseDir DIR         A base directory to search for Poseidon packages.
608  --remote                 List packages from a remote server instead the local
609                           file system.
610  --remoteURL URL          URL of the remote Poseidon server.
611                           (default: "https://server.poseidon-adna.org")
612  --archive STRING         The name of the Poseidon package archive that should
613                           be queried. If not given, then the query falls back
614                           to the default archive of the server selected with
615                           --remoteURL. See the archive documentation at
616                           https://www.poseidon-adna.org/#/archive_overview for
617                           a list of archives currently available from the
618                           official Poseidon Web API. (default: Nothing)
619  --packages               List all packages.
620  --groups                 List all groups, ignoring any group names after the
621                           first as specified in the .janno-file.
622  --individuals            List all individuals/samples.
623  -j,--jannoColumn COLNAME List additional fields from the janno files, using
624                           the .janno column heading name, such as "Country",
625                           "Site", "Date_C14_Uncal_BP", etc..
626  --raw                    Return the output table as tab-separated values
627                           without header. This is useful for piping into grep
628                           or awk.
629  --onlyLatest             Consider only the latest versions of packages, or the
630                           groups and individuals within the latest versions of
631                           packages, respectively.
632
633 To list packages from your local repositories, as seen above you can run
634
635 trident list -d ... -d ... --packages
636
637 This will yield a nicely formatted table of all packages, their version and the number of individuals in them.
638
639 You can use --remote to show packages on the remote server. For example
640
641 trident list --packages --remote --archive "community-archive"
642
643 will result in a view of all packages available in one of the public online archives. Just as for fetch, the --archive
644 flag allows to choose which public archive to query.
645
646 Independent of whether you query a local or an online archive, you can not just list packages, but also groups,

```


640 as defined in the third column of EIGENSTRAT .ind files (or the first/last column of a PLINK .fam file), and
641 individuals with the flags --groups and --individuals (instead of --packages).

642 The --individuals flag additionally provides a way to immediately access information from .janno files
643 on the command line. This works with the -j/--jannoColumn option. For example adding -j Country -j
644 Date_C14_Uncal_BP to the commands above will add the Country and the Date_C14_Uncal_BP columns to the
645 respective output tables.

646 Note that if you want a less fancy table, for example because you want to load this into Excel, or pipe into
647 another command that cannot deal with the table layout, you can use the --raw option to output that table as
648 a simple tab-delimited stream.

649 **0.1.3.2 Summarise command** summarise prints some general summary statistics for a given poseidon
650 dataset taken from the .janno files.

651 [Click here for command line details](#)

652 Usage: trident summarise (-d|--baseDir DIR) [--raw]

653

654 Get an overview over the content of one or multiple Poseidon packages

655

656 Available options:

657 -h,--help	Show this help text
658 -d,--baseDir DIR	A base directory to search for Poseidon packages.
659 --raw	Return the output table as tab-separated values 660 without header. This is useful for piping into grep 661 or awk.

662 You can run it with

663 trident summarise -d ... -d ...

664 which will show you context information like – among others – the number of individuals in the dataset, their
665 sex distribution, the mean age of the samples (for ancient data) or the mean coverage on the 1240K SNP array
666 in a table. summarise depends on complete .janno files and will silently ignore missing information.

667 You can use the --raw option to output the summary table in a simple, tab-delimited layout.

668 **0.1.3.3 Survey command** survey tries to indicate package completeness (mostly focused on .janno files)
669 for poseidon datasets.

670 [Click here for command line details](#)

671 Usage: trident survey (-d|--baseDir DIR) [--raw] [--onlyLatest]

672

673 Survey the degree of context information completeness for Poseidon packages

674

675 Available options:

676 -h,--help	Show this help text
677 -d,--baseDir DIR	A base directory to search for Poseidon packages.
678 --raw	Return the output table as tab-separated values

without header. This is useful for piping into grep or awk.

`--onlyLatest` Consider only the latest versions of packages, or the groups and individuals within the latest versions of packages, respectively.

Running

```
trident survey -d ... -d ...
```

will yield a table with one row for each package. See `trident survey -h` for a legend which cell of this table means what.

Again you can use the `--raw` option to output the survey table in a tab-delimited format.

0.1.3.4 Validate command `validate` checks Poseidon packages and individual package components for structural correctness.

[Click here for command line details](#)

Usage: `trident validate ((-d|--baseDir DIR) [--ignoreGeno] [--fullGeno] [--ignoreDuplicates] [-c|--ignoreChecksums] [--ignorePoseidonVersion] | --pym1 FILE | (-p|--genoOne FILE) | --inFormat FORMAT --genoFile FILE --snpFile FILE --indFile FILE | --janno FILE | --ssf FILE | --bib FILE) [--noExitCode] [--onlyLatest]`

Check Poseidon packages or package components for structural correctness

Available options:

<code>-h,--help</code>	Show this help text
<code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
<code>--ignoreGeno</code>	Ignore snp and geno file.
<code>--fullGeno</code>	Test parsing of all SNPs (by default only the first 100 SNPs are probed).
<code>--ignoreDuplicates</code>	Do not stop on duplicated individual names in the package collection.
<code>-c,--ignoreChecksums</code>	Whether to ignore checksums. Useful for speedup in debugging.
<code>--ignorePoseidonVersion</code>	Read packages even if their poseidonVersion is not compatible with trident.
<code>--pym1 FILE</code>	Path to a POSEIDON.yml file.
<code>-p,--genoOne FILE</code>	One of the input genotype data files. Expects .bed, .bim or .fam for PLINK and .geno, .snp or .ind for EIGENSTRAT. The other files must be in the same directory and must have the same base name.
<code>--inFormat FORMAT</code>	The format of the input genotype data: EIGENSTRAT or PLINK. Only necessary for data input with --genoFile

```

721         + --snpFile + --indFile.
722 --genoFile FILE      Path to the input geno file.
723 --snpFile FILE       Path to the input snp file.
724 --indFile FILE       Path to the input ind file.
725 --janno FILE         Path to a .janno file.
726 --ssf FILE           Path to a .ssf file.
727 --bib FILE           Path to a .bib file.
728 --noExitCode         Do not produce an explicit exit code.
729 --onlyLatest         Consider only the latest versions of packages, or the
730                       groups and individuals within the latest versions of
731                       packages, respectively.

```

732 You can run it with

```

733 trident validate -d ... -d ...

```

734 to check packages and it will either report a success (**Validation passed**) or failure with specific error messages.

735 Instead of validating entire packages with **-d** you can also apply it to individual files and package components: **--pyml** (POSEIDON.yml), **-p | --inFormat + --genoFile + --snpFile + --indFile** (genotype data), **--janno** (.janno file), **--ssf** (.ssf file) or **--bib** (.bib file). In this case **validate** attempts to read and parse the respective files individually and reports any issues it encounters. Note that this considers the files in isolation and does not include any cross-file consistency checks.

740 When applied to packages, **validate** tries to ensure that each package adheres to the **schema definition**. Here is a list of what is checked:

- 742 • Structural correctness of the POSEIDON.yml file.
- 743 • Presence of all files references in the POSEIDON.yml file.
- 744 • Full structural correctness of .janno, .ssf and .bib file.
- 745 • Superficial correctness of genotype data files by parsing the first 100 SNPs. A full check that parses all SNPs can be triggered with the **--fullGeno** option. **--ignoreGeno**, on the other hand, causes **validate** to ignore the genotype data entirely, which speeds up the validation significantly.
- 748 • Correspondence of BibTeX keys in .bib and .janno
- 749 • Correspondence of sample IDs in .janno and .ssf.
- 750 • Correspondence of sample and group IDs in .janno and genotype data files.

751 In fact much of this validation already runs as part of the general package reading pipeline invoked for other trident subcommands (e.g. **forge**). **validate** is meant to be more thorough/brittle, though, and will explicitly fail if even a single package is broken. For special cases more flexibility can be enabled with the options **--ignoreDuplicates**, **--ignoreChecksums** and **--ignorePoseidonVersion**.

755 Remember to run **validate** it with **--debug** to get more information in case the default output is not sufficient to analyse an issue.