

# Guide for trident v1.1.10.2

## Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>The trident CLI</b>                            | <b>1</b>  |
| 1.1      | General notes                                     | 4         |
| 1.1.1    | Logging and command line output                   | 4         |
| 1.1.2    | Duplicates  | 4         |
| 1.1.3    | Group names in .fam files                         | 4         |
| 1.1.4    | Whitespaces in the .janno file                    | 4         |
| <b>2</b> | <b>Package creation and manipulation commands</b> | <b>4</b>  |
| 2.1      | Init command                                      | 4         |
| 2.2      | Fetch command                                     | 6         |
| 2.3      | Forge command                                     | 7         |
| 2.3.1    | The forge selection language                      | 10        |
| 2.3.2    | Treatment of the .janno file while merging        | 11        |
| 2.3.3    | Other options                                     | 12        |
| 2.4      | Genoconvert command                               | 12        |
| 2.5      | Update command                                    | 14        |
| <b>3</b> | <b>Inspection commands</b>                        | <b>15</b> |
| 3.1      | List command                                      | 15        |
| 3.2      | Summarise command                                 | 17        |
| 3.3      | Survey command                                    | 17        |
| 3.4      | Validate command                                  | 18        |

## 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 ARG] [--errLength ARG]`

`[--inPlinkPopName ARG] (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:

```

34  -h,--help          Show this help text
35  --version          Show version number
36  --logMode ARG      How information should be reported: NoLog, SimpleLog,
37                    DefaultLog, ServerLog or VerboseLog
38                    (default: DefaultLog)
39  --errLength ARG    After how many characters should a potential error
40                    message be truncated. "Inf" for no truncation.
41                    (default: CharCount 1500)
42  --inPlinkPopName ARG Where to read the population/group name from the FAM
43                    file in Plink-format. Three options are possible:
44                    asFamily (default) | asPhenotype | asBoth.
45
46  Package creation and manipulation commands:
47  init              Create a new Poseidon package from genotype data
48  fetch             Download data from a remote Poseidon repository
49  forge             Select packages, groups or individuals and create a
50                    new Poseidon package from them
51  genoconvert       Convert the genotype data in a Poseidon package to a
52                    different file format
53  update            Update POSEIDON.yml files automatically
54
55  Inspection commands:
56  list              List packages, groups or individuals from local or
57                    remote Poseidon repositories
58  summarise         Get an overview over the content of one or multiple
59                    Poseidon packages
60  summarize         Synonym for summarise
61  survey            Survey the degree of context information completeness
62                    for Poseidon packages
63  validate          Check one or multiple Poseidon packages for
64                    structural correctness
65
66  Trident allows to work directly with genotype data (see -p below), but its optimized for the interaction with
67  Poseidon packages, which wrap and contextualize the data. Most trident subcommands therefore have a central
68  parameter, called --baseDir or simply -d to specify one or more base directories to look for packages. For example,
69  if all Poseidon packages live inside a repository at /path/to/poseidon/packages you would simply say trident
70  <subcommand> -d /path/to/poseidon/dirs/ and trident would automatically search all subdirectories inside
71  of the repository for valid Poseidon packages (as identified by valid POSEIDON.yml files).
72
73  You can arrange a poseidon repository in a hierarchical way. For example:
74
75  /path/to/poseidon/packages
76    /modern
77      /2019_poseidon_package1
78      /2019_poseidon_package2
79    /ancient
80    /...

```

```

78     /...
79     /Reference_Genomes
80     /...
81     /...

```

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

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

```

88 ~/my_project/my_project.geno
89 ~/my_project/my_project.snp
90 ~/my_project/my_project.ind

```

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

```

93 poseidonVersion: 2.5.0
94 title: My_awesome_project
95 description: Unpublished genetic data from my awesome project
96 contributor:
97   - name: Stephan Schiffels
98     email: schiffels@institute.org
99 packageVersion: 0.1.0
100 lastModified: 2020-10-07
101 genotypeData:
102   format: EIGENSTRAT
103   genoFile: my_project.geno
104   snpFile: my_project.snp
105   indFile: my_project.ind
106   jannoFile: my_project.janno
107   bibFile: sources.bib

```

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

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

```

114 trident list -d /path/to/poseidon/packages/modern \
115   -d /path/to/poseidon/packages/ReferenceGenomes
116   -d ~/my_project --packages

```

## 1.1 General notes

### 1.1.1 Logging and command line output

For all subcommands the general argument `--logMode` defines how trident reports messages (to stderr) on the command line:

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

### 1.1.2 Duplicates

- If multiple packages in a package repository share the same **title**, then trident will try to select the one with the highest version number. If this is not sufficient to resolve the conflict, trident will stop.
- Individual/sample names (**Poseidon\_IDs**) within one package have to be unique, or trident will stop.
- We generally also discourage ID duplicates across packages in package repositories, but trident will generally continue with them after printing a warning. This does not apply for **validate**, by default (you can change this behaviour with `--ignoreDuplicates`), and **forge**. **forge** offers a special mechanism to resolve duplicates within its selection language (see below).

### 1.1.3 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.

### 1.1.4 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**.

## 2 Package creation and manipulation commands

### 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 ARG) | --inFormat ARG --genoFile ARG
                    --snpFile ARG --indFile ARG) [--snpSet ARG]
                    (-o|--outPackagePath ARG) [-n|--outPackageName ARG]
```

```

153             [--minimal]
154     Create a new Poseidon package from genotype data
155
156     Available options:
157     -h,--help                Show this help text
158     -p,--genoOne ARG         one of the input genotype data files. Expects .bed or
159                               .bim or .fam for PLINK and .geno or .snp or .ind for
160                               EIGENSTRAT. The other files must be in the same
161                               directory and must have the same base name
162     --inFormat ARG           the format of the input genotype data: EIGENSTRAT or
163                               PLINK (only necessary for data input with --genoFile
164                               + --snpFile + --indFile)
165     --genoFile ARG           the input geno file path
166     --snpFile ARG            the input snp file path
167     --indFile ARG            the input ind file path
168     --snpSet ARG             the snpSet of the package: 1240K, HumanOrigins or
169                               Other. (only relevant for data input with
170                               -p|--genoOne or --genoFile + --snpFile + --indFile,
171                               because the packages in a -d|--baseDir already have
172                               this information in their respective POSEIDON.yml
173                               files) Default: Other
174     -o,--outPackagePath ARG  the output package directory path
175     -n,--outPackageName ARG  the output package name - this is optional: If no
176                               name is provided, then the package name defaults to
177                               the basename of the (mandatory) --outPackagePath
178                               argument
179     --minimal                should only a minimal output package be created?
180
181     The command
182
183     trident init \
184     --inFormat EIGENSTRAT/PLINK \
185     --genoFile path/to/geno_file \
186     --snpFile path/to/snp_file \
187     --indFile path/to/ind_file \
188     --snpSet 1240K|HumanOrigins|Other \
189     -o path/to/new_package_name
190
191     requires the format (--inFormat) of your input data (either EIGENSTRAT or PLINK), the paths to the respective
192     files (--genoFile, --snpFile, --indFile), and optionally the “shape” of these files (--snpSet), so if they cover
193     the 1240K, the HumanOrigins or an Other SNP set. A simpler interface added in trident 0.29.0 is available with
194     -p (+ --snpSet).

```

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

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

## 195 2.2 Fetch command

196 `fetch` allows to download Poseidon packages from a remote Poseidon server. Read more about this repository  
197 [here](#).

198 Click here for command line details

```
199 Usage: trident fetch (-d|--baseDir DIR)
200             (--downloadAll |
201             (--fetchFile ARG | (-f|--fetchString ARG)))
202             [--remoteURL ARG] [-u|--upgrade]
203 Download data from a remote Poseidon repository
```

205 Available options:

|     |                                   |   |
|-----|-----------------------------------|---|
| 206 | <code>-h,--help</code>            | Show this help text   |
| 207 | <code>-d,--baseDir DIR</code>     | a base directory to search for Poseidon Packages<br>(could be a Poseidon repository)  |
| 209 | <code>--downloadAll</code>        | download all packages the server is offering  |
| 210 | <code>--fetchFile ARG</code>      | A file with a list of packages. Works just as <code>-f</code> , but<br>multiple values can also be separated by newline, not<br>just by comma. <code>-f</code> and <code>--fetchFile</code> can be combined.  |
| 213 | <code>-f,--fetchString ARG</code> | List of packages to be downloaded from the remote<br>server. Package names should be wrapped in asterisks:<br><code>*package_title*</code> . You can combine multiple values with<br>comma, so for example: <code>"*package_1*, *package_2*,<br/>*package_3*"</code> . <code>fetchString</code> uses the same parser as<br><code>forgeString</code> , but does not allow excludes. If groups<br>or individuals are specified, then packages which<br>include these groups or individuals are included in<br>the download. |
| 222 | <code>--remoteURL ARG</code>      | URL of the remote Poseidon server<br>(default: <code>"https://c107-224.cloud.gwdg.de"</code> )  |
| 224 | <code>-u,--upgrade</code>         | overwrite outdated local package versions   |

225 It works with

```
226 trident fetch -d ... -d ... \
227 -f "*package_title_1*,*package_title_2*,*package_title_3*,group_name,<Individual1>"
```

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

231 Entities are specified using a special syntax (see also the documentation of `forge` below): Package titles are  
232 wrapped in asterisks: `package_title`, group names are spelled as is, and individual names are wrapped in angular

brackets, like `<Individual1>`. Fetch will figure out which packages need to be downloaded to include all specified entities. `--downloadAll`, which can be given instead of `-f` and `--fetchFile`, causes fetch to download all packages from the server. The downloaded packages are added in the first (!) `-d` directory (which gets created if it doesn't exist), but downloads are only performed if the respective packages are not already present in an up-to-date version in any of the `-d` dirs.

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

`fetch` also has the optional arguments `--remote https://...` to name an alternative poseidon server. The default points to the **DAG server**.

To overwrite outdated package versions with `fetch`, the `-u/--upgrade` flag has to be set. Note that many file systems do not offer a way to recover overwritten files. So be careful with this switch.

## 2.3 Forge command

`forge` creates new Poseidon packages by extracting and merging packages, populations and individuals from your Poseidon repositories.

[Click here for command line details](#)

```
Usage: trident forge ((-d|--baseDir DIR) |
                    ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG
                    --snpFile ARG --indFile ARG) [--snpSet ARG])
                    [--forgeFile ARG | (-f|--forgeString ARG)]
                    [--selectSnps ARG] [--intersect] [--outFormat ARG]
                    [--minimal] [--onlyGeno] (-o|--outPackagePath ARG)
                    [-n|--outPackageName ARG] [--packagewise]
                    [--outPlinkPopName ARG]
```

Select packages, groups or individuals and create a new Poseidon package from them

Available options:

|                               |   |
|-------------------------------|---|
| <code>-h,--help</code>        | Show this help text   |
| <code>-d,--baseDir DIR</code> | a base directory to search for Poseidon Packages (could be a Poseidon repository)   |
| <code>-p,--genoOne ARG</code> | one of the input genotype data files. Expects .bed or .bim or .fam for PLINK and .geno or .snp or .ind for EIGENSTRAT. The other files must be in the same directory and must have the same base name |
| <code>--inFormat ARG</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 ARG</code>   | the input geno file path  |
| <code>--snpFile ARG</code>    | the input snp file path   |
| <code>--indFile ARG</code>    | the input ind file path   |
| <code>--snpSet ARG</code>     | the snpSet of the package: 1240K, HumanOrigins or Other. (only relevant for data input with   |

275 -p|--genoOne or --genoFile + --snpFile + --indFile,  
276 because the packages in a -d|--baseDir already have  
277 this information in their respective POSEIDON.yml  
278 files) Default: Other

279 --forgeFile ARG A file with a list of packages, groups or individual  
280 samples. Works just as -f, but multiple values can  
281 also be separated by newline, not just by comma.  
282 Empty lines are ignored and comments start with "#",  
283 so everything after "#" is ignored in one line.  
284 Multiple instances of -f and --forgeFile can be  
285 given. They will be evaluated according to their  
286 input order on the command line.

287 -f,--forgeString ARG List of packages, groups or individual samples to be  
288 combined in the output package. Packages follow the  
289 syntax \*package\_title\*, populations/groups are simply  
290 group\_id and individuals <individual\_id>. You can  
291 combine multiple values with comma, so for example:  
292 "\*package\_1\*, <individual\_1>, <individual\_2>,"  
293 group\_1". Duplicates are treated as one entry.  
294 Negative selection is possible by prepending "-" to  
295 the entity you want to exclude (e.g. "\*package\_1\*,"  
296 -<individual\_1>, -group\_1"). forge will apply  
297 excludes and includes in order. If the first entity  
298 is negative, then forge will assume you want to merge  
299 all individuals in the packages found in the baseDirs  
300 (except the ones explicitly excluded) before the  
301 exclude entities are applied. An empty forgeString  
302 (and no --forgeFile) will therefore merge all  
303 available individuals. If there are individuals in  
304 your input packages with equal individual id, but  
305 different main group or source package, they can be  
306 specified with the special syntax  
307 "<package:group:individual>".

308 --selectSnps ARG To extract specific SNPs during this forge operation,  
309 provide a Snp file. Can be either Eigenstrat (file  
310 ending must be '.snp') or Plink (file ending must be  
311 '.bim'). When this option is set, the output package  
312 will have exactly the SNPs listed in this file. Any  
313 SNP not listed in the file will be excluded. If  
314 option '--intersect' is also set, only the SNPs  
315 overlapping between the SNP file and the forged  
316 packages are output.

317 --intersect Whether to output the intersection of the genotype  
318 files to be forged. The default (if this option is  
319 not set) is to output the union of all SNPs, with



```

320         genotypes defined as missing in those packages which
321         do not have a SNP that is present in another package.
322         With this option set, the forged dataset will
323         typically have fewer SNPs, but less missingness.
324     --outFormat ARG         the format of the output genotype data: EIGENSTRAT or
325                             PLINK. Default: PLINK
326     --minimal               should only a minimal output package be created?
327     --onlyGeno              should only the resulting genotype data be returned?
328                             This means the output will not be a Poseidon package
329     -o,--outPackagePath ARG the output package directory path
330     -n,--outPackageName ARG the output package name - this is optional: If no
331                             name is provided, then the package name defaults to
332                             the basename of the (mandatory) --outPackagePath
333                             argument
334     --packagewise           Skip the within-package selection step in forge. This
335                             will result in outputting all individuals in the
336                             relevant packages, and hence a superset of the
337                             requested individuals/groups. It may result in better
338                             performance in cases where one wants to forge entire
339                             packages or almost entire packages. Details: Forge
340                             conceptually performs two types of selection: First,
341                             it identifies which packages in the supplied base
342                             directories are relevant to the requested forge, i.e.
343                             whether they are either explicitly listed using
344                             *PackageName*, or because they contain selected
345                             individuals or groups. Second, within each relevant
346                             package, individuals which are not requested are
347                             removed. This option skips only the second step, but
348                             still performs the first.
349     --outPlinkPopName ARG   Where to write the population/group name into the FAM
350                             file in Plink-format. Three options are possible:
351                             asFamily (default) | asPhenotype | asBoth. See also
352                             --inPlinkPopName.
353
354     forge can be used with
355
356     trident forge -d ... -d ... \
357         -f "*package_name*, group_id, <individual_id>" \
358         -o path/to/new_package_name
359
360     where the entities (packages, groups/populations, individuals/samples) you want in the output package can be
361     denoted either as a string on the command line (-f/--forgeString), or in an input text file (--forgeFile).
362     See the section below for the syntax of this selection language. Do not forget to wrap the --forgeString query
363     in quotes.
364
365     Including one or multiple Poseidon packages with -d is not the only way to include data for a forge operation.
366     It is also possible to consider unpackaged genotype data directly with -p (+ --snpSet) or --inFormat +
367     --genoFile + --snpFile + --indFile (+ --snpSet). This makes the following example possible, where we

```

merge data from one Poseidon package and two genotype datasets to get a new EIGENSTRAT dataset.

```
trident forge \  
  -d 2017_GonzalesFortesCurrentBiology \  
  -p 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \  
  --inFormat PLINK \  
  --genoFile 2017_HaberAJHG/2017_HaberAJHG.bed \  
  --snpFile 2017_HaberAJHG/2017_HaberAJHG.bim \  
  --indFile 2017_HaberAJHG/2017_HaberAJHG.fam \  
  -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \  
  -o testpackage \  
  --outFormat EIGENSTRAT \  
  --onlyGeno
```

### 2.3.1 The forge selection language

The text in `--forgeString` and `--forgeFile` 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.

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

```
# Packages  
*package1*, *package2*  
  
# Groups and individuals from other packages beyond package1 and package2  
group1, <individual1>, group2, <individual2>, <individual3>  
  
# group2 has two outlier individuals that should be ignored  
-<bad_individual1> # This one has very low coverage  
-<bad_individual2> # This one is from a different time period
```

By prepending `-` to the bad individuals, we can exclude them from the forged package. `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 packages found in the baseDirs (except the ones explicitly excluded, of course).

An empty forgeString will therefore merge all available individuals.

### 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                 |

### 2.3.3 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 an 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 | <code>--intersect</code> | 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.

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

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

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

## 2.4 Genoconvert command

`genoconvert` converts the genotype data in a Poseidon package to a different file format. The respective entries in the `POSEIDON.yml` file are changed accordingly.

[Click here for command line details](#)

```
Usage: trident genoconvert ((-d|--baseDir DIR) |
                           ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG
                           --snpFile ARG --indFile ARG) [--snpSet ARG])
                           --outFormat ARG [--onlyGeno]
                           [-o|--outPackagePath ARG] [--removeOld]
```

```

458             [--outPlinkPopName ARG]
459 Convert the genotype data in a Poseidon package to a different file format
460
461 Available options:
462 -h,--help                Show this help text
463 -d,--baseDir DIR          a base directory to search for Poseidon Packages
464                           (could be a Poseidon repository)
465 -p,--genoOne ARG          one of the input genotype data files. Expects .bed or
466                           .bim or .fam for PLINK and .geno or .snp or .ind for
467                           EIGENSTRAT. The other files must be in the same
468                           directory and must have the same base name
469 --inFormat ARG            the format of the input genotype data: EIGENSTRAT or
470                           PLINK (only necessary for data input with --genoFile
471                           + --snpFile + --indFile)
472 --genoFile ARG            the input geno file path
473 --snpFile ARG             the input snp file path
474 --indFile ARG             the input ind file path
475 --snpSet ARG              the snpSet of the package: 1240K, HumanOrigins or
476                           Other. (only relevant for data input with
477                           -p|--genoOne or --genoFile + --snpFile + --indFile,
478                           because the packages in a -d|--baseDir already have
479                           this information in their respective POSEIDON.yml
480                           files) Default: Other
481 --outFormat ARG           the format of the output genotype data: EIGENSTRAT or
482                           PLINK.
483 --onlyGeno                should only the resulting genotype data be returned?
484                           This means the output will not be a Poseidon package
485 -o,--outPackagePath ARG   the output package directory path - this is optional:
486                           If no path is provided, then the output is written to
487                           the directories where the input genotype data file
488                           (.bed/.geno) is stored
489 --removeOld               Remove the old genotype files when creating the new
490                           ones
491 --outPlinkPopName ARG     Where to write the population/group name into the FAM
492                           file in Plink-format. Three options are possible:
493                           asFamily (default) | asPhenotype | asBoth. See also
494                           --inPlinkPopName.
495
496 With the default setting
497
498 trident genoconvert -d ... -d ... --outFormat EIGENSTRAT|PLINK
499
500 all packages in -d will be converted to the desired --outFormat (either EIGENSTRAT or PLINK), if the data is
501 not already in this format. This includes updating the respective POSEIDON.yml files.
502
503 The “old” data is not deleted, but kept around. That means conversion can result in a package with both PLINK
504 and EIGENSTRAT data, but only one is linked in the POSEIDON.yml file, and that is what will be used by
505 trident. To delete the old data in the conversion you can add the --removeOld flag.

```

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

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

## 509 2.5 Update command

510 `update` automatically harmonizes POSEIDON.yml files of one or multiple packages if the packages were changed.  
 511 This is not an automatic update from one Poseidon version to the next!

512 Click here for command line details

```
513 Usage: trident update (-d|--baseDir DIR) [--poseidonVersion ARG]  
514               [--ignorePoseidonVersion] [--versionComponent ARG]  
515               [--noChecksumUpdate] [--newContributors ARG]  
516               [--logText ARG] [--force]
```

517 Update POSEIDON.yml files automatically

518  
 519 Available options:

|     |                                      |   |
|-----|--------------------------------------|---|
| 520 | <code>-h,--help</code>               | Show this help text   |
| 521 | <code>-d,--baseDir DIR</code>        | a base directory to search for Poseidon Packages                          |
| 522 |                                      | (could be a Poseidon repository)  |
| 523 | <code>--poseidonVersion ARG</code>   | Poseidon version the packages should be updated to:                       |
| 524 |                                      | e.g. "2.5.3" (default: Nothing)   |
| 525 | <code>--ignorePoseidonVersion</code> | Read packages even if their poseidonVersion is not                        |
| 526 |                                      | compatible with the trident version. The assumption                       |
| 527 |                                      | is, that the package is already structurally adjusted                     |
| 528 |                                      | to the trident version and only the version number is                     |
| 529 |                                      | lagging behind.   |
| 530 | <code>--versionComponent ARG</code>  | Part of the package version number in the                                 |
| 531 |                                      | POSEIDON.yml file that should be updated: Major,                          |
| 532 |                                      | Minor or Patch (see <a href="https://semver.org">https://semver.org</a> ) |
| 533 |                                      | (default: Patch)  |
| 534 | <code>--noChecksumUpdate</code>      | Should update of checksums in the POSEIDON.yml file                       |
| 535 |                                      | be skipped  |
| 536 | <code>--ignoreGeno</code>            | ignore SNP and GenoFile   |
| 537 | <code>--newContributors ARG</code>   | Contributors to add to the POSEIDON.yml file in the                       |
| 538 |                                      | form "[Firstname Lastname](Email address);..."                            |
| 539 | <code>--logText ARG</code>           | Log text for this version jump in the CHANGELOG file                      |
| 540 |                                      | (default: "not specified")  |
| 541 | <code>--force</code>                 | Normally the POSEIDON.yml files are only changed if                       |
| 542 |                                      | the poseidonVersion is adjusted or any of the                             |
| 543 |                                      | checksums change. With <code>--force</code> a package version             |
| 544 |                                      | update can be triggered even if this is not the case.                     |

545 It can be called with a lot of optional arguments

```
546 trident update -d ... -d ... \  
547 --poseidonVersion "X.X.X" \  
548 --versionComponent Major/Minor/Patch \  
549 --noChecksumUpdate  
550 --ignoreGeno  
551 --newContributors "[Firstname Lastname](Email address);..."  
552 --logText "short description of the update"  
553 --force
```

554 By default `update` will not edit a package's POSEIDON.yml file, even when arguments like `--versionComponent`,  
555 `--newContributors` or `--logText` are explicitly set. This default exists to run the function on a large set of  
556 packages where only few of them were edited and need an active update. A package will only be modified by  
557 `update` if either

- 558 • any of the files with checksums (e.g. the genotype data) in it were modified,
- 559 • the `--poseidonVersion` argument differs from the `poseidonVersion` in the package's POSEIDON.yml  
560 file
- 561 • or the `--force` flag was set in `update`.

562 If any of these applies to a package in the search directory (`--baseDir/-d`), it will be updated. This includes  
563 the following steps:

- 564 • If `--poseidonVersion` is different from the `poseidonVersion` field in the package, then that will be  
565 updated.
- 566 • The `packageVersion` will be incremented. If `--versionComponent` is not set, then it falls back to `Patch`,  
567 so a change in the last position of the three digit version number. `Minor` increments the middle, and `Major`  
568 the first position (see [semantic versioning](#)).
- 569 • The `lastModified` field will be updated to the current day (based on your computer's system time).
- 570 • The contributors in `--newContributors` will be added to the `contributor` field if they're not there already.
- 571 • If any checksums changed, then they will be updated. If certain checksums are not set yet, then they will  
572 be added. The checksum update can be skipped with `--noChecksumUpdate` or partially skipped for the  
573 genotype data with `--ignoreGeno`.
- 574 • The CHANGELOG.md file will be updated with a new row for the new version and the text in `--logText`  
575 (default: "not specified"), which will be appended as the first line of the file. If no CHANGELOG.md file  
576 exists, then it will be created and referenced in the POSEIDON.yml file.

577 :heavy\_exclamation\_mark: As `update` reads and rewrites POSEIDON.yml files, it may change their inner order,  
578 layout or even content (e.g. if they have fields which are not in the [Poseidon package definition](#)). Create a backup  
579 of the POSEIDON.yml file before running `update` if you are uncertain.

## 580 3 Inspection commands

### 581 3.1 List command

582 `list` lists packages, groups and individuals of the datasets you use, or of the packages available on the server.

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

```

584 Usage: trident list ((-d|--baseDir DIR) | --remote [--remoteURL ARG])
585                (--packages | --groups | --individuals
586                [-j|--jannoColumn JANNO_HEADER])) [--raw]
587 List packages, groups or individuals from local or remote Poseidon
588 repositories
589
590 Available options:
591  -h,--help                Show this help text
592  -d,--baseDir DIR         a base directory to search for Poseidon Packages
593                           (could be a Poseidon repository)
594  --remote                 list packages from a remote server instead the local
595                           file system
596  --remoteURL ARG          URL of the remote Poseidon server
597                           (default: "https://c107-224.cloud.gwdg.de")
598  --packages               list all packages
599  --groups                 list all groups, ignoring any group names after the
600                           first as specified in the Janno-file
601  --individuals            list individuals
602  -j,--jannoColumn JANNO_HEADER
603                           list additional fields from the janno files, using
604                           the Janno column heading name, such as Country, Site,
605                           Date_C14_Uncal_BP, Endogenous, ...
606  --raw                    output table as tsv without header. Useful for piping
607                           into grep or awk
608  --ignoreGeno             ignore SNP and GenoFile
609
610 To list packages from your local repositories, as seen above you can run
611
612 trident list -d ... -d ... --packages
613
614 This will yield a table like this
615
616 .------.------.------.
617 |                Title                |    Date    | Nr Individuals |
618 :=====:=====:=====:
619 | 2015_1000Genomes_1240K_haploid_pull | 2020-08-10 | 2535           |
620 | 2016_Mallick_SGDP1240K_diploid_pull | 2020-08-10 | 280            |
621 | 2018_BostonDatashare_modern_published | 2020-08-10 | 2772           |
622 | ...                                | ...        |                |
623 '-----'-----'-----'

```

so a nicely formatted table of all packages, their last update and the number of individuals in it.

To view packages on the remote server, instead of using directories to specify the locations of repositories on your system, you can use `--remote` to show packages on the remote server. For example

```
trident list --packages --remote
```

will result in a view of all published packages in our [public online repository](#).

You can also list groups, as defined in the third column of EIGENSTRAT .ind files (or the first column of a



626 PLINK .fam file), and individuals with `--groups` and `--individuals` instead of `--packages`.

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

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

## 634 3.2 Summarise command

635 `summarise` prints some general summary statistics for a given poseidon dataset taken from the .janno files.

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

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

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

639

640 Available options:

|     |                               |  |
|-----|-------------------------------|--|
| 641 | <code>-h,--help</code>        | Show this help text  |
| 642 | <code>-d,--baseDir DIR</code> | a base directory to search for Poseidon Packages<br>(could be a Poseidon repository) |
| 643 |                               |  |
| 644 | <code>--raw</code>            | output table as tsv without header. Useful for piping<br>into grep or awk            |
| 645 |                               |  |

646 You can run it with

647 `trident summarise -d ... -d ...`

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

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

## 653 3.3 Survey command

654 `survey` tries to indicate package completeness (mostly focused on .janno files) for poseidon datasets.

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

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

657 Survey the degree of context information completeness for Poseidon packages

658

659 Available options:

|     |                               |  |
|-----|-------------------------------|--|
| 660 | <code>-h,--help</code>        | Show this help text  |
| 661 | <code>-d,--baseDir DIR</code> | a base directory to search for Poseidon Packages<br>(could be a Poseidon repository) |
| 662 |                               |  |
| 663 | <code>--raw</code>            | output table as tsv without header. Useful for piping                                |

664                                   into grep or awk

665 Running

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

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

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

### 670 3.4 Validate command

671 `validate` checks poseidon datasets for structural correctness.

672 Click here for command line details

673 Usage: `trident validate (-d|--baseDir DIR)`

674     Check one or multiple Poseidon packages for structural correctness

675  
676 Available options:

|                                     |  |
|-------------------------------------|--|
| 677 <code>-h,--help</code>          | Show this help text  |
| 678 <code>-d,--baseDir DIR</code>   | a base directory to search for Poseidon Packages<br>679     (could be a Poseidon repository) |
| 680 <code>--ignoreGeno</code>       | ignore SNP and GenoFile  |
| 681 <code>--fullGeno</code>         | test parsing of all SNPs (by default only the first<br>682     100 SNPs are probed)          |
| 683 <code>--noExitCode</code>       | do not produce an explicit exit code   |
| 684 <code>--ignoreDuplicates</code> | do not stop on duplicated individual names in the<br>685     package collection              |

686 You can run it with

687 `trident validate -d ... -d ...`

688 and it will either report a success (**Validation passed**) or failure with specific error messages to simplify fixing  
689 the issues.

690 `validate` tries to ensure that each package in the dataset adheres to the [schema definition](#). Here is a list of  
691 what is checked:

- 692     • Presence of the necessary files
- 693     • Full structural correctness of `.bib` and `.janno` file
- 694     • Superficial correctness of genotype data files by parsing the first 100 SNPs. A full check that parses all  
695     SNPs can be run with the `--fullGeno` option
- 696     • Correspondence of BibTeX keys in `.bib` and `.janno`
- 697     • Correspondence of individual and group IDs in `.janno` and genotype data files

698 In fact much of this validation already runs as part of the general package reading pipeline invoked for many  
699 trident subcommands (e.g. `forge`). `validate` is meant to be more thorough, though, and will explicitly fail if  
700 even a single package is broken.

701 Remember to run it with `--logMode VerboseLog` to get more information if the output is not sufficient to debug  
702 an issue.