

Guide for trident v1.1.7.0

Contents

| | | |
|----------|--|-----------|
| 1 | The trident CLI | 1 |
| 1.1 | Handling data with trident | 2 |
| 1.2 | Notes on duplicates | 4 |
| 2 | Package creation and manipulation commands | 4 |
| 2.1 | Init command | 4 |
| 2.2 | Fetch command | 5 |
| 2.3 | Forge command | 6 |
| 2.3.1 | The forge selection language | 9 |
| 2.3.2 | Treatment of the .janno file while merging | 10 |
| 2.3.3 | Other options | 11 |
| 2.4 | Genoconvert command | 12 |
| 2.5 | Update command | 13 |
| 3 | Inspection commands | 15 |
| 3.1 | List command | 15 |
| 3.2 | Summarise command | 16 |
| 3.3 | Survey command | 17 |
| 3.4 | Validate command | 17 |

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] (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 ARG</code> | How information should be reported: NoLog, SimpleLog, DefaultLog, ServerLog or VerboseLog |

```

34                                     (default: DefaultLog)
35  --errLength ARG                    After how many characters should a potential error
36                                     message be truncated. "Inf" for no truncation.
37                                     (default: CharCount 1500)
38
39 Package creation and manipulation commands:
40  init                               Create a new Poseidon package from genotype data
41  fetch                              Download data from a remote Poseidon repository
42  forge                              Select packages, groups or individuals and create a
43                                     new Poseidon package from them
44  genoconvert                        Convert the genotype data in a Poseidon package to a
45                                     different file format
46  update                             Update POSEIDON.yml files automatically
47
48 Inspection commands:
49  list                               List packages, groups or individuals from local or
50                                     remote Poseidon repositories
51  summarise                          Get an overview over the content of one or multiple
52                                     Poseidon packages
53  summarize                          Synonym for summarise
54  survey                             Survey the degree of context information completeness
55                                     for Poseidon packages
56  validate                           Check one or multiple Poseidon packages for
57                                     structural correctness
58
59 For all subcommands the general argument --logMode defines how trident reports messages (to stderr) on the
60 command line:
61
62 • NoLog: Hides all messages.
63 • SimpleLog: Plain and simple output to stderr.
64 • DefaultLog: Adds severity indicators before each message. (default setting)
65 • ServerLog: Additionally adds timestamps before each message.
66 • VerboseLog: Shows not just messages on the log levels Info, Warning and Error like the other modes,
67   but also on the more verbose level Debug. Use this for debugging.

```

66 1.1 Handling data with trident

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

74 You can arrange a poseidon repository in a hierarchical way. For example:

```
75 /path/to/poseidon/packages
```

```

76     /modern
77         /2019_poseidon_package1
78         /2019_poseidon_package2
79     /ancient
80         /...
81         /...
82     /Reference_Genomes
83         /...
84         /...

```

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

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

```

91 ~/my_project/my_project.geno
92 ~/my_project/my_project.snp
93 ~/my_project/my_project.ind

```

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

```

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

```

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

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

```

118 trident list -d /path/to/poseidon/packages/modern \
119     -d /path/to/poseidon/packages/ReferenceGenomes
120     -d ~/my_project --packages

```

1.2 Notes on 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_ID`s) 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).

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](#)

```

134 Usage: trident init ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG
135                   --snpFile ARG --indFile ARG) [--snpSet ARG]
136                   (-o|--outPackagePath ARG) [-n|--outPackageName ARG]
137                   [--minimal]

```

Create a new Poseidon package from genotype data

Available options:

| | |
|-------------------------------|--|
| <code>-h,--help</code> | Show this help text |
| <code>-p,--genoOne ARG</code> | one of the input genotype data files. Expects <code>.bed</code> or <code>.bim</code> or <code>.fam</code> for PLINK and <code>.geno</code> or <code>.snp</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 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 <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 |

```

158 -o,--outPackagePath ARG the output package directory path
159 -n,--outPackageName ARG the output package name - this is optional: If no
160 name is provided, then the package name defaults to
161 the basename of the (mandatory) --outPackagePath
162 argument
163 --minimal should only a minimal output package be created?

```

164 The command

```

165 trident init \
166 --inFormat EIGENSTRAT/PLINK \
167 --genoFile path/to/geno_file \
168 --snpFile path/to/snp_file \
169 --indFile path/to/ind_file \
170 --snpSet 1240K|HumanOrigins|Other \
171 -o path/to/new_package_name

```

172 requires the format (`--inFormat`) of your input data (either `EIGENSTRAT` or `PLINK`), the paths to the
173 respective files (`--genoFile` , `--snpFile` , `--indFile`), and optionally the “shape” of these files (`--snpSet`),
174 so if they cover the `1240K` , the `HumanOrigins` or an `Other` SNP set. A simpler interface added in trident
175 0.29.0 is available with `-p` (+ `--snpSet`) .

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

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

180 2.2 Fetch command

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

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

```

184 Usage: trident fetch (-d|--baseDir DIR)
185         (--downloadAll |
186         (--fetchFile ARG | (-f|--fetchString ARG)))
187         [--remoteURL ARG] [-u|--upgrade]

```

188 Download data from a remote Poseidon repository

189 Available options:

```

191 -h,--help          Show this help text
192 -d,--baseDir DIR   a base directory to search for Poseidon Packages

```

```

193         (could be a Poseidon repository)
194     --downloadAll      download all packages the server is offering
195     --fetchFile ARG    A file with a list of packages. Works just as -f, but
196                        multiple values can also be separated by newline, not
197                        just by comma. -f and --fetchFile can be combined.
198     -f,--fetchString ARG List of packages to be downloaded from the remote
199                        server. Package names should be wrapped in asterisks:
200                        *package_title*. You can combine multiple values with
201                        comma, so for example: "*package_1*, *package_2*,
202                        *package_3*". fetchString uses the same parser as
203                        forgeString, but does not allow excludes. If groups
204                        or individuals are specified, then packages which
205                        include these groups or individuals are included in
206                        the download.
207     --remoteURL ARG    URL of the remote Poseidon server
208                        (default: "https://c107-224.cloud.gwdg.de")
209     -u,--upgrade       overwrite outdated local package versions

```

210 It works with

```

211 trident fetch -d ... -d ... \
212     -f "*package_title_1*,*package_title_2*,*package_title_3*,group_name,<Individual1>"

```

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

216 Entities are specified using a special syntax (see also the documentation of `forge` below): Package titles
217 are wrapped in asterisks: *package_title*, group names are spelled as is, and individual names are wrapped in
218 angular brackets, like `<Individual1>`. Fetch will figure out which packages need to be downloaded to include
219 all specified entities. `--downloadAll`, which can be given instead of `-f` and `--fetchFile`, causes fetch to
220 download all packages from the server. The downloaded packages are added in the first (!) `-d` directory (which
221 gets created if it doesn't exist), but downloads are only performed if the respective packages are not already
222 present in an up-to-date version in any of the `-d` dirs.

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

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

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

229 2.3 Forge command

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

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

```

233 Usage: trident forge ((-d|--baseDir DIR) |
234                      ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG
235                      --snpFile ARG --indFile ARG) [--snpSet ARG])
236                      [--forgeFile ARG | (-f|--forgeString ARG)]
237                      [--selectSnps ARG] [--intersect] [--outFormat ARG]
238                      [--minimal] [--onlyGeno] (-o|--outPackagePath ARG)
239                      [-n|--outPackageName ARG] [--no-extract]
240 Select packages, groups or individuals and create a new Poseidon package from
241 them
242
243 Available options:
244 -h,--help                Show this help text
245 -d,--baseDir DIR         a base directory to search for Poseidon Packages
246                          (could be a Poseidon repository)
247 -p,--genoOne ARG         one of the input genotype data files. Expects .bed or
248                          .bim or .fam for PLINK and .geno or .snp or .ind for
249                          EIGENSTRAT. The other files must be in the same
250                          directory and must have the same base name
251 --inFormat ARG           the format of the input genotype data: EIGENSTRAT or
252                          PLINK (only necessary for data input with --genoFile
253                          + --snpFile + --indFile)
254 --genoFile ARG           the input geno file path
255 --snpFile ARG            the input snp file path
256 --indFile ARG            the input ind file path
257 --snpSet ARG             the snpSet of the package: 1240K, HumanOrigins or
258                          Other. (only relevant for data input with
259                          -p|--genoOne or --genoFile + --snpFile + --indFile,
260                          because the packages in a -d|--baseDir already have
261                          this information in their respective POSEIDON.yml
262                          files) Default: Other
263 --forgeFile ARG          A file with a list of packages, groups or individual
264                          samples. Works just as -f, but multiple values can
265                          also be separated by newline, not just by comma.
266                          Empty lines are ignored and comments start with "#",
267                          so everything after "#" is ignored in one line.
268                          Multiple instances of -f and --forgeFile can be
269                          given. They will be evaluated according to their
270                          input order on the command line.
271 -f,--forgeString ARG     List of packages, groups or individual samples to be
272                          combined in the output package. Packages follow the
273                          syntax *package_title*, populations/groups are simply
274                          group_id and individuals <individual_id>. You can
275                          combine multiple values with comma, so for example:
276                          "*package_1*, <individual_1>, <individual_2>,
277                          group_1". Duplicates are treated as one entry.

```

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

292 --selectSnps ARG To extract specific SNPs during this forge operation,
293 provide a Snp file. Can be either Eigenstrat (file
294 ending must be '.snp') or Plink (file ending must be
295 '.bim'). When this option is set, the output package
296 will have exactly the SNPs listed in this file. Any
297 SNP not listed in the file will be excluded. If
298 option '--intersect' is also set, only the SNPs
299 overlapping between the SNP file and the forged
300 packages are output.

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

308 --outFormat ARG the format of the output genotype data: EIGENSTRAT or
309 PLINK. Default: PLINK

310 --minimal should only a minimal output package be created?

311 --onlyGeno should only the resulting genotype data be returned?

312 This means the output will not be a Poseidon package

313 -o,--outPackagePath ARG the output package directory path

314 -n,--outPackageName ARG the output package name - this is optional: If no
315 name is provided, then the package name defaults to
316 the basename of the (mandatory) --outPackagePath
317 argument

318 --no-extract Skip the selection step in forge. This will result in
319 outputting all individuals in the relevant packages,
320 and hence a superset of the requested
321 individuals/groups. It may result in better
322 performance in cases where one wants to forge entire

packages or almost entire packages. Note that this will also ignore any ordering in the output groups/individuals. With this option active, individuals from the relevant packages will just be written in the order that they appear in the original packages.

`forge` can be used with

```
trident forge -d ... -d ... \
  -f "*package_name*, group_id, <individual_id>" \
  -o path/to/new_package_name
```

where the entities (packages, groups/populations, individuals/samples) you want in the output package can be denoted either as a string on the command line (`-f / --forgeString`), or in an input text file (`--forgeFile`). See the section below for the syntax of this selection language. Do not forget to wrap the `--forgeString` query in quotes.

Including one or multiple Poseidon packages with `-d` is not the only way to include data for a forge operation. It is also possible to consider unpackaged genotype data directly with `-p (+ --snpSet)` or `--inFormat + --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 |

| Poseidon_ID | Group_Name | Genetic_Sex | AdditionalColumn1 | AdditionalColumn2 |
|-------------|------------|-------------|-------------------|-------------------|
| 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 `snpSet` s 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 |

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

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

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

424 2.4 Genoconvert command

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

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

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

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

435 Available options:

| | |
|-----------------------------------|--|
| 436 <code>-h,--help</code> | Show this help text |
| 437 <code>-d,--baseDir DIR</code> | a base directory to search for Poseidon Packages (could be a Poseidon repository) |
| 439 <code>-p,--genoOne ARG</code> | one of the input genotype data files. Expects <code>.bed</code> or 440 <code>.bim</code> or <code>.fam</code> for PLINK and <code>.geno</code> or <code>.snp</code> or <code>.ind</code> for 441 EIGENSTRAT. The other files must be in the same 442 directory and must have the same base name |
| 443 <code>--inFormat ARG</code> | the format of the input genotype data: EIGENSTRAT or 444 PLINK (only necessary for data input with <code>--genoFile</code> 445 <code>+ --snpFile + --indFile</code>) |
| 446 <code>--genoFile ARG</code> | the input geno file path |
| 447 <code>--snpFile ARG</code> | the input snp file path |
| 448 <code>--indFile ARG</code> | the input ind file path |
| 449 <code>--snpSet ARG</code> | the snpSet of the package: 1240K, HumanOrigins or 450 Other. (only relevant for data input with 451 <code>-p --genoOne</code> or <code>--genoFile + --snpFile + --indFile</code> , 452 because the packages in a <code>-d --baseDir</code> already have 453 this information in their respective POSEIDON.yml 454 files) Default: Other |
| 455 <code>--outFormat ARG</code> | the format of the output genotype data: EIGENSTRAT or 456 PLINK. |
| 457 <code>--onlyGeno</code> | should only the resulting genotype data be returned? 458 This means the output will not be a Poseidon package |

459 `-o,--outPackagePath ARG` the output package directory path - this is optional:
 460 If no path is provided, then the output is written to
 461 the directories where the input genotype data file
 462 (`.bed/.geno`) is stored
 463 `--removeOld` Remove the old genotype files when creating the new
 464 ones

465 With the default setting

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

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

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

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

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

479 2.5 Update command

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

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

483 Usage: `trident update (-d|--baseDir DIR) [--poseidonVersion ARG]`
 484 `[--ignorePoseidonVersion] [--versionComponent ARG]`
 485 `[--noChecksumUpdate] [--newContributors ARG]`
 486 `[--logText ARG] [--force]`

487 Update `POSEIDON.yml` files automatically

488 Available options:

490 `-h,--help` Show this help text
 491 `-d,--baseDir DIR` a base directory to search for Poseidon Packages
 492 (could be a Poseidon repository)
 493 `--poseidonVersion ARG` Poseidon version the packages should be updated to:
 494 e.g. "2.5.3" (default: Nothing)
 495 `--ignorePoseidonVersion` Read packages even if their `poseidonVersion` is not
 496 compatible with the `trident` version. The assumption
 497 is, that the package is already structurally adjusted
 498 to the `trident` version and only the version number is
 499 lagging behind.

```

500 --versionComponent ARG    Part of the package version number in the
501                             POSEIDON.yml file that should be updated: Major,
502                             Minor or Patch (see https://semver.org)
503                             (default: Patch)
504 --noChecksumUpdate         Should update of checksums in the POSEIDON.yml file
505                             be skipped
506 --ignoreGeno               ignore SNP and GenoFile
507 --newContributors ARG      Contributors to add to the POSEIDON.yml file in the
508                             form "[Firstname Lastname](Email address);..."
509 --logText ARG              Log text for this version jump in the CHANGELOG file
510                             (default: "not specified")
511 --force                    Normally the POSEIDON.yml files are only changed if
512                             the poseidonVersion is adjusted or any of the
513                             checksums change. With --force a package version
514                             update can be triggered even if this is not the case.

```

515 It can be called with a lot of optional arguments

```

516 trident update -d ... -d ... \
517   --poseidonVersion "X.X.X" \
518   --versionComponent Major/Minor/Patch \
519   --noChecksumUpdate
520   --ignoreGeno
521   --newContributors "[Firstname Lastname](Email address);..."
522   --logText "short description of the update"
523   --force

```

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

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

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

- 534 • If `--poseidonVersion` is different from the `poseidonVersion` field in the package, then that will be
535 updated.
- 536 • The `packageVersion` will be incremented. If `--versionComponent` is not set, then it falls back to
537 `Patch`, so a change in the last position of the three digit version number. `Minor` increments the middle,
538 and `Major` the first position (see [semantic versioning](#)).
- 539 • The `lastModified` field will be updated to the current day (based on your computer's system time).
- 540 • The contributors in `--newContributors` will be added to the `contributor` field if they're not there
541 already.
- 542 • If any checksums changed, then they will be updated. If certain checksums are not set yet, then they will

be added. The checksum update can be skipped with `--noChecksumUpdate` or partially skipped for the genotype data with `--ignoreGeno`.

- The CHANGELOG.md file will be updated with a new row for the new version and the text in `--logText` (default: “not specified”), which will be appended as the first line of the file. If no CHANGELOG.md file exists, then it will be created and referenced in the POSEIDON.yml file.

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

3 Inspection commands

3.1 List command

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

[Click here for command line details](#)

```
Usage: trident list ((-d|--baseDir DIR) | --remote [--remoteURL ARG])
        (--packages | --groups | --individuals
        [-j|--jannoColumn JANNO_HEADER]) [--raw]
```

List packages, groups or individuals from local or remote Poseidon repositories

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>--remote</code> | list packages from a remote server instead the local file system |
| <code>--remoteURL ARG</code> | URL of the remote Poseidon server (default: "https://c107-224.cloud.gwdg.de") |
| <code>--packages</code> | list all packages |
| <code>--groups</code> | list all groups, ignoring any group names after the first as specified in the Janno-file |
| <code>--individuals</code> | list individuals |
| <code>-j,--jannoColumn JANNO_HEADER</code> | list additional fields from the janno files, using the Janno column heading name, such as Country, Site, Date_C14_Uncal_BP, Endogenous, ... |
| <code>--raw</code> | output table as tsv without header. Useful for piping into grep or awk |
| <code>--ignoreGeno</code> | ignore SNP and GenoFile |

To list packages from your local repositories, as seen above you can run

```
trident list -d ... -d ... --packages
```

This will yield a table like this

```

583 .------.------.------.
584 |           Title           |   Date   | Nr Individuals |
585 :=====:=====:=====:
586 | 2015_1000Genomes_1240K_haploid_pulldown | 2020-08-10 | 2535          |
587 | 2016_Mallick_SGDP1240K_diploid_pulldown | 2020-08-10 | 280           |
588 | 2018_BostonDatashare_modern_published   | 2020-08-10 | 2772          |
589 | ...                                     | ...       |               |
590 '-----'-----'-----'

```

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

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

594 `trident list --packages --remote`

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

596 You can also list groups, as defined in the third column of EIGENSTRAT `.ind` files (or the first column of a
597 PLINK `.fam` file), and individuals with `--groups` and `--individuals` instead of `--packages`.

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

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

605 3.2 Summarise command

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

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

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

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

611 Available options:

```

612  -h,--help           Show this help text
613  -d,--baseDir DIR    a base directory to search for Poseidon Packages
614                      (could be a Poseidon repository)
615  --raw               output table as tsv without header. Useful for piping
616                      into grep or awk

```

617 You can run it with

618 `trident summarise -d ... -d ...`

619 which will show you context information like – among others – the number of individuals in the dataset, their
620 sex distribution, the mean age of the samples (for ancient data) or the mean coverage on the 1240K SNP array

621 in a table. `summarise` depends on complete .janno files and will silently ignore missing information for some
622 statistics.

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

624 3.3 Survey command

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

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

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

628 Survey the degree of context information completeness for Poseidon packages

629

630 Available options:

| | | |
|-----|-------------------------------|---|
| 631 | <code>-h,--help</code> | Show this help text |
| 632 | <code>-d,--baseDir DIR</code> | a base directory to search for Poseidon Packages (could be a Poseidon repository) |
| 633 | | |
| 634 | <code>--raw</code> | output table as tsv without header. Useful for piping 635 into <code>grep</code> or <code>awk</code> |

636 Running

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

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

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

641 3.4 Validate command

642 `validate` checks poseidon datasets for structural correctness.

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

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

645 Check one or multiple Poseidon packages for structural correctness

646

647 Available options:

| | | |
|-----|---------------------------------|--|
| 648 | <code>-h,--help</code> | Show this help text |
| 649 | <code>-d,--baseDir DIR</code> | a base directory to search for Poseidon Packages (could be a Poseidon repository) |
| 650 | | |
| 651 | <code>--ignoreGeno</code> | ignore SNP and GenoFile |
| 652 | <code>--noExitCode</code> | do not produce an explicit exit code |
| 653 | <code>--ignoreDuplicates</code> | do not stop on duplicated individual names in the 654 package collection |

655 You can run it with

656 `trident validate -d ... -d ...`

657 and it will either report a success (`Validation passed`) or failure with specific error messages to simplify
658 fixing the issues.

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

- 661 • Presence of the necessary files
- 662 • Full structural correctness of .bib and .janno file
- 663 • Superficial correctness of genotype data files. A full check would be too computationally expensive
- 664 • Correspondence of BibTeX keys in .bib and .janno
- 665 • Correspondence of individual and group IDs in .janno and genotype data files

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

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