

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	16
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, but
67 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 with
68 Poseidon packages, which wrap and contextualize the data. Most trident subcommands therefore have a central
69 parameter, called --baseDir or simply -d to specify one or more base directories to look for packages. For example,
70 if all Poseidon packages live inside a repository at /path/to/poseidon/packages you would simply say trident
71 <subcommand> -d /path/to/poseidon/dirs/ and trident would automatically search all subdirectories inside
72 of the repository for valid Poseidon packages (as identified by valid POSEIDON.yml files).

```

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

```

74 /path/to/poseidon/packages
75     /modern

```

```

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

```

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

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

```

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

```

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

```

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97 description: Unpublished genetic data from my awesome project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

```

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

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

```

116 trident list -d /path/to/poseidon/packages/modern \
117   -d /path/to/poseidon/packages/ReferenceGenomes

```

118 `-d ~/my_project --packages`

119 1.2 Notes on duplicates

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

127 2 Package creation and manipulation commands

128 2.1 Init command

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

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

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

136 Create a new Poseidon package from genotype data

137 Available options:

139 <code>-h,--help</code>	Show this help text
140 <code>-p,--genoOne ARG</code>	one of the input genotype data files. Expects <code>.bed</code> or 141 <code>.bim</code> or <code>.fam</code> for PLINK and <code>.geno</code> or <code>.snp</code> or <code>.ind</code> for 142 EIGENSTRAT. The other files must be in the same 143 directory and must have the same base name
144 <code>--inFormat ARG</code>	the format of the input genotype data: EIGENSTRAT or 145 PLINK (only necessary for data input with <code>--genoFile</code> 146 + <code>--snpFile</code> + <code>--indFile</code>)
147 <code>--genoFile ARG</code>	the input geno file path
148 <code>--snpFile ARG</code>	the input snp file path
149 <code>--indFile ARG</code>	the input ind file path
150 <code>--snpSet ARG</code>	the snpSet of the package: 1240K, HumanOrigins or 151 Other. (only relevant for data input with 152 <code>-p --genoOne</code> or <code>--genoFile</code> + <code>--snpFile</code> + <code>--indFile</code> , 153 because the packages in a <code>-d --baseDir</code> already have 154 this information in their respective <code>POSEIDON.yml</code> 155 files) Default: Other
156 <code>-o,--outPackagePath ARG</code>	the output package directory path
157 <code>-n,--outPackageName ARG</code>	the output package name - this is optional: If no

158 name is provided, then the package name defaults to
 159 the basename of the (mandatory) `--outPackagePath`
 160 argument
 161 `--minimal` should only a minimal output package be created?

162 The command

```
163 trident init \  

164   --inFormat EIGENSTRAT/PLINK \  

165   --genoFile path/to/geno_file \  

166   --snpFile path/to/snp_file \  

167   --indFile path/to/ind_file \  

168   --snpSet 1240K|HumanOrigins|Other \  

169   -o path/to/new_package_name
```

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

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

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

177 2.2 Fetch command

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

180 Click here for command line details

```
181 Usage: trident fetch (-d|--baseDir DIR)  

182           (--downloadAll |  

183           (--fetchFile ARG | (-f|--fetchString ARG)))  

184           [--remoteURL ARG] [-u|--upgrade]
```

185 Download data from a remote Poseidon repository

186 Available options:

188 <code>-h,--help</code>	Show this help text
189 <code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (could be a Poseidon repository)
191 <code>--downloadAll</code>	download all packages the server is offering
192 <code>--fetchFile ARG</code>	A file with a list of packages. Works just as <code>-f</code> , but

```

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

```

207 It works with

```

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

```

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

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

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

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

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

226 2.3 Forge command

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

229 Click here for command line details

```

230 Usage: trident forge ((-d|--baseDir DIR) |
231                     ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG
232                     --snpFile ARG --indFile ARG) [--snpSet ARG])

```

```

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

```

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

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

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

305 --outFormat ARG the format of the output genotype data: EIGENSTRAT or
 306 PLINK. Default: PLINK

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

308 --onlyGeno should only the resulting genotype data be returned?
 309 This means the output will not be a Poseidon package

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

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

315 --no-extract Skip the selection step in forge. This will result in
 316 outputting all individuals in the relevant packages,
 317 and hence a superset of the requested
 318 individuals/groups. It may result in better
 319 performance in cases where one wants to forge entire
 320 packages or almost entire packages. Note that this
 321 will also ignore any ordering in the output
 322 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
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`.

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

419 2.4 Genoconvert command

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

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

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

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

429
430 Available options:

431 <code>-h,--help</code>	Show this help text
432 <code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages 433 (could be a Poseidon repository)
434 <code>-p,--genoOne ARG</code>	one of the input genotype data files. Expects <code>.bed</code> or 435 <code>.bim</code> or <code>.fam</code> for PLINK and <code>.geno</code> or <code>.snp</code> or <code>.ind</code> for 436 EIGENSTRAT. The other files must be in the same 437 directory and must have the same base name
438 <code>--inFormat ARG</code>	the format of the input genotype data: EIGENSTRAT or 439 PLINK (only necessary for data input with <code>--genoFile</code> 440 + <code>--snpFile</code> + <code>--indFile</code>)
441 <code>--genoFile ARG</code>	the input geno file path
442 <code>--snpFile ARG</code>	the input snp file path
443 <code>--indFile ARG</code>	the input ind file path
444 <code>--snpSet ARG</code>	the snpSet of the package: 1240K, HumanOrigins or 445 Other. (only relevant for data input with 446 <code>-p --genoOne</code> or <code>--genoFile</code> + <code>--snpFile</code> + <code>--indFile</code> , 447 because the packages in a <code>-d --baseDir</code> already have 448 this information in their respective <code>POSEIDON.yml</code> 449 files) Default: Other
450 <code>--outFormat ARG</code>	the format of the output genotype data: EIGENSTRAT or 451 PLINK.
452 <code>--onlyGeno</code>	should only the resulting genotype data be returned? 453 This means the output will not be a Poseidon package
454 <code>-o,--outPackagePath ARG</code>	the output package directory path - this is optional: 455 If no path is provided, then the output is written to 456 the directories where the input genotype data file 457 (<code>.bed/.geno</code>) is stored
458 <code>--removeOld</code>	Remove the old genotype files when creating the new 459 ones

460 With the default setting

```

461 trident genoconvert -d ... -d ... --outFormat EIGENSTRAT|PLINK
462 all packages in -d will be converted to the desired --outFormat (either EIGENSTRAT or PLINK), if the data is
463 not already in this format. This includes updating the respective POSEIDON.yml files.
464 The “old” data is not deleted, but kept around. That means conversion can result in a package with both PLINK
465 and EIGENSTRAT data, but only one is linked in the POSEIDON.yml file, and that is what will be used by
466 trident. To delete the old data in the conversion you can add the --removeOld flag.
467 Instead of -d to change Poseidon packages, the -p (+ --snpSet) or --inFormat + --genoFile + --snpFile
468 + --indFile (+ --snpSet) allow to directly convert genotype data that is not wrapped in a Poseidon package
469 and store it to a directory given in -o. See this example:
470 trident genoconvert \
471   -p 2018_Mittnik_Baltic/Mittnik_Baltic.bed \
472   --outFormat EIGENSTRAT
473   -o my_directory

```

474 2.5 Update command

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

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

```

478 Usage: trident update (-d|--baseDir DIR) [--poseidonVersion ARG]
479                [--ignorePoseidonVersion] [--versionComponent ARG]
480                [--noChecksumUpdate] [--newContributors ARG]
481                [--logText ARG] [--force]
482 Update POSEIDON.yml files automatically

```

484 Available options:

485 -h,--help	Show this help text
486 -d,--baseDir DIR	a base directory to search for Poseidon Packages (could be a Poseidon repository)
488 --poseidonVersion ARG	Poseidon version the packages should be updated to: e.g. "2.5.3" (default: Nothing)
490 --ignorePoseidonVersion	Read packages even if their poseidonVersion is not compatible with the trident version. The assumption is, that the package is already structurally adjusted to the trident version and only the version number is lagging behind.
495 --versionComponent ARG	Part of the package version number in the POSEIDON.yml file that should be updated: Major, Minor or Patch (see https://semver.org) (default: Patch)
499 --noChecksumUpdate	Should update of checksums in the POSEIDON.yml file be skipped
501 --ignoreGeno	ignore SNP and GenoFile

```

502 --newContributors ARG    Contributors to add to the POSEIDON.yml file in the
503                          form "[Firstname Lastname](Email address);..."
504 --logText ARG            Log text for this version jump in the CHANGELOG file
505                          (default: "not specified")
506 --force                  Normally the POSEIDON.yml files are only changed if
507                          the poseidonVersion is adjusted or any of the
508                          checksums change. With --force a package version
509                          update can be triggered even if this is not the case.

```

510 It can be called with a lot of optional arguments

```

511 trident update -d ... -d ... \
512   --poseidonVersion "X.X.X" \
513   --versionComponent Major/Minor/Patch \
514   --noChecksumUpdate
515   --ignoreGeno
516   --newContributors "[Firstname Lastname](Email address);..."
517   --logText "short description of the update"
518   --force

```

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

- 523 • any of the files with checksums (e.g. the genotype data) in it were modified,
- 524 • the **--poseidonVersion** argument differs from the **poseidonVersion** in the package's POSEIDON.yml
525 file
- 526 • or the **--force** flag was set in **update**.

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

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

542 :heavy_exclamation_mark: As **update** reads and rewrites POSEIDON.yml files, it may change their inner order,
543 layout or even content (e.g. if they have fields which are not in the [Poseidon package definition](#)). Create a backup
544 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:

```
-h,--help          Show this help text
-d,--baseDir DIR   a base directory to search for Poseidon Packages
                   (could be a Poseidon repository)
--remote           list packages from a remote server instead the local
                   file system
--remoteURL ARG    URL of the remote Poseidon server
                   (default: "https://c107-224.cloud.gwdg.de")
--packages         list all packages
--groups           list all groups, ignoring any group names after the
                   first as specified in the Janno-file
--individuals      list individuals
-j,--jannoColumn JANNO_HEADER
                   list additional fields from the janno files, using
                   the Janno column heading name, such as Country, Site,
                   Date_C14_Uncal_BP, Endogenous, ...
--raw             output table as tsv without header. Useful for piping
                   into grep or awk
--ignoreGeno      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

```
.------.------.------.
|           Title           |    Date    | Nr Individuals |
:=====:=====:=====:
| 2015_1000Genomes_1240K_haploid_pulldown | 2020-08-10 | 2535          |
| 2016_Mallick_SGDP1240K_diploid_pulldown | 2020-08-10 | 280           |
| 2018_BostonDatashare_modern_published   | 2020-08-10 | 2772          |
| ...                                     | ...         |               |
'------'------'-----'
```

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

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

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

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

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

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

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

599 3.2 Summarise command

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

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

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

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

```
604
605 Available options:
606     -h,--help                Show this help text
607     -d,--baseDir DIR         a base directory to search for Poseidon Packages
608                             (could be a Poseidon repository)
609     --raw                    output table as tsv without header. Useful for piping
610                             into grep or awk
```

611 You can run it with

```
612 trident summarise -d ... -d ...
```

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

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

618 3.3 Survey command

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

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

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

```
622     Survey the degree of context information completeness for Poseidon packages
```


623

624 Available options:

625	<code>-h,--help</code>	Show this help text
626	<code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages
627		(could be a Poseidon repository)
628	<code>--raw</code>	output table as tsv without header. Useful for piping
629		into grep or awk

630 Running

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

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

633 means what.

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

635 3.4 Validate command

636 `validate` checks poseidon datasets for structural correctness.

637 Click here for command line details

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

639 Check one or multiple Poseidon packages for structural correctness

640

641 Available options:

642	<code>-h,--help</code>	Show this help text
643	<code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages
644		(could be a Poseidon repository)
645	<code>--ignoreGeno</code>	ignore SNP and GenoFile
646	<code>--noExitCode</code>	do not produce an explicit exit code
647	<code>--ignoreDuplicates</code>	do not stop on duplicated individual names in the
648		package collection

649 You can run it with

650 `trident validate -d ... -d ...`

651 and it will either report a success (**Validation passed**) or failure with specific error messages to simplify fixing

652 the issues.

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

654 what is checked:

- 655 • Presence of the necessary files
- 656 • Full structural correctness of `.bib` and `.janno` file
- 657 • Superficial correctness of genotype data files. A full check would be too computationally expensive
- 658 • Correspondence of BibTeX keys in `.bib` and `.janno`
- 659 • Correspondence of individual and group IDs in `.janno` and genotype data files

660 In fact much of this validation already runs as part of the general package reading pipeline invoked for many

661 `trident` subcommands (e.g. `forge`). `validate` is meant to be more thorough, though, and will explicitly fail if

662 even a single package is broken.

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