

Guide for trident v1.1.11.0 to v1.1.12.0

Contents

1	The trident CLI	1
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
2	Package creation and manipulation commands	4
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	Treatment of the .ssf file while merging	12
2.3.4	Other options	12
2.4	Genoconvert command	12
2.5	Update command	14
3	Inspection commands	16
3.1	List command	16
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>

34 Available options:

```
35  -h,--help          Show this help text
36  --version          Show version number
37  --logMode ARG      How information should be reported: NoLog, SimpleLog,
38                     DefaultLog, ServerLog or VerboseLog
39                     (default: DefaultLog)
40  --errLength ARG     After how many characters should a potential error
41                     message be truncated. "Inf" for no truncation.
42                     (default: CharCount 1500)
43  --inPlinkPopName ARG Where to read the population/group name from the FAM
44                     file in Plink-format. Three options are possible:
45                     asFamily (default) | asPhenotype | asBoth.
```

46
47 Package creation and manipulation commands:

```
48  init              Create a new Poseidon package from genotype data
49  fetch             Download data from a remote Poseidon repository
50  forge             Select packages, groups or individuals and create a
51                     new Poseidon package from them
52  genoconvert        Convert the genotype data in a Poseidon package to a
53                     different file format
54  update             Update POSEIDON.yml files automatically
```

55
56 Inspection commands:

```
57  list              List packages, groups or individuals from local or
58                     remote Poseidon repositories
59  summarise          Get an overview over the content of one or multiple
60                     Poseidon packages
61  summarize          Synonym for summarise
62  survey             Survey the degree of context information completeness
63                     for Poseidon packages
64  validate           Check one or multiple Poseidon packages for
65                     structural correctness
```

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

```
73 /path/to/poseidon/packages
74   /modern
75     /2019_poseidon_package1
76     /2019_poseidon_package2
77   /ancient
```

```

78         /...
79         /...
80     /Reference_Genomes
81         /...
82         /...

```

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

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

```

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

```

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

```

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

```

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

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

```

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

```

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

```

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

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

196 2.2 Fetch command

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

199 Click here for command line details

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

206 Available options:

207	<code>-h,--help</code>	Show this help text
208	<code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (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 212 multiple values can also be separated by newline, not 213 just by comma. <code>-f</code> and <code>--fetchFile</code> can be combined.
214	<code>-f,--fetchString ARG</code>	List of packages to be downloaded from the remote 215 server. Package names should be wrapped in asterisks: 216 <code>*package_title*</code> . You can combine multiple values with 217 comma, so for example: <code>"*package_1*, *package_2*,</code> 218 <code>*package_3*"</code> . <code>fetchString</code> uses the same parser as 219 <code>forgeString</code> , but does not allow excludes. If groups 220 or individuals are specified, then packages which 221 include these groups or individuals are included in 222 the download.
223	<code>--remoteURL ARG</code>	URL of the remote Poseidon server (default: "https://c107-224.cloud.gwdg.de")
224	<code>-u,--upgrade</code>	overwrite outdated local package versions

226 It works with

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

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

232 Entities are specified using a special syntax (see also the documentation of `forge` below): Package titles are
233 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

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

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

288 -f,--forgeString ARG List of packages, groups or individual samples to be
289 combined in the output package. Packages follow the
290 syntax *package_title*, populations/groups are simply
291 group_id and individuals <individual_id>. You can
292 combine multiple values with comma, so for example:
293 "*package_1*, <individual_1>, <individual_2>,"
294 group_1". Duplicates are treated as one entry.
295 Negative selection is possible by prepending "-" to
296 the entity you want to exclude (e.g. "*package_1*,"
297 -<individual_1>, -group_1"). forge will apply
298 excludes and includes in order. If the first entity
299 is negative, then forge will assume you want to merge
300 all individuals in the packages found in the baseDirs
301 (except the ones explicitly excluded) before the
302 exclude entities are applied. An empty forgeString
303 (and no --forgeFile) will therefore merge all
304 available individuals. If there are individuals in
305 your input packages with equal individual id, but
306 different main group or source package, they can be
307 specified with the special syntax
308 "<package:group:individual>".

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

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


```

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

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

2.3.4 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 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>	Output 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.

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

```
460 Usage: trident genoconvert ((-d|--baseDir DIR) |
461      ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG
462      --snpFile ARG --indFile ARG) [--snpSet ARG])
463      --outFormat ARG [--onlyGeno]
464      [-o|--outPackagePath ARG] [--removeOld]
465      [--outPlinkPopName ARG]
```

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

467 Available options:

```
468 -h,--help          Show this help text
469 -d,--baseDir DIR    a base directory to search for Poseidon Packages
470                    (could be a Poseidon repository)
471 -p,--genoOne ARG    one of the input genotype data files. Expects .bed or
472                    .bim or .fam for PLINK and .geno or .snp or .ind for
473                    EIGENSTRAT. The other files must be in the same
474                    directory and must have the same base name
475 --inFormat ARG       the format of the input genotype data: EIGENSTRAT or
476                    PLINK (only necessary for data input with --genoFile
477                    + --snpFile + --indFile)
478 --genoFile ARG       the input geno file path
479 --snpFile ARG        the input snp file path
480 --indFile ARG        the input ind file path
481 --snpSet ARG         the snpSet of the package: 1240K, HumanOrigins or
482                    Other. (only relevant for data input with
483                    -p|--genoOne or --genoFile + --snpFile + --indFile,
484                    because the packages in a -d|--baseDir already have
485                    this information in their respective POSEIDON.yml
486                    files) Default: Other
487 --outFormat ARG      the format of the output genotype data: EIGENSTRAT or
488                    PLINK.
489 --onlyGeno           should only the resulting genotype data be returned?
490                    This means the output will not be a Poseidon package
491 -o,--outPackagePath ARG the output package directory path - this is optional:
492                    If no path is provided, then the output is written to
493                    the directories where the input genotype data file
494                    (.bed/.geno) is stored
495 --removeOld          Remove the old genotype files when creating the new
496                    ones
497 --outPlinkPopName ARG Where to write the population/group name into the FAM
498                    file in Plink-format. Three options are possible:
499                    asFamily (default) | asPhenotype | asBoth. See also
500                    --inPlinkPopName.
501
```

502 With the default setting

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

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

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

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

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

516 2.5 Update command

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

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

```
520 Usage: trident update (-d|--baseDir DIR) [--poseidonVersion ARG]  
521                   [--ignorePoseidonVersion] [--versionComponent ARG]  
522                   [--noChecksumUpdate] [--newContributors ARG]  
523                   [--logText ARG] [--force]
```

524 Update POSEIDON.yml files automatically

525
526 Available options:

527	<code>-h,--help</code>	Show this help text
528	<code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages 529 (could be a Poseidon repository)
530	<code>--poseidonVersion ARG</code>	Poseidon version the packages should be updated to: 531 e.g. "2.5.3" (default: Nothing)
532	<code>--ignorePoseidonVersion</code>	Read packages even if their poseidonVersion is not 533 compatible with the trident version. The assumption 534 is, that the package is already structurally adjusted 535 to the trident version and only the version number is 536 lagging behind.
537	<code>--versionComponent ARG</code>	Part of the package version number in the 538 POSEIDON.yml file that should be updated: Major, 539 Minor or Patch (see https://semver.org) 540 (default: Patch)
541	<code>--noChecksumUpdate</code>	Should update of checksums in the POSEIDON.yml file 542 be skipped
543	<code>--ignoreGeno</code>	ignore SNP and GenoFile

```

544 --newContributors ARG    Contributors to add to the POSEIDON.yml file in the
545                          form "[Firstname Lastname](Email address);..."
546 --logText ARG            Log text for this version jump in the CHANGELOG file
547                          (default: "not specified")
548 --force                  Normally the POSEIDON.yml files are only changed if
549                          the poseidonVersion is adjusted or any of the
550                          checksums change. With --force a package version
551                          update can be triggered even if this is not the case.

```

552 It can be called with a lot of optional arguments

```

553 trident update -d ... -d ... \
554   --poseidonVersion "X.X.X" \
555   --versionComponent Major/Minor/Patch \
556   --noChecksumUpdate
557   --ignoreGeno
558   --newContributors "[Firstname Lastname](Email address);..."
559   --logText "short description of the update"
560   --force

```

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

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

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

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

584 :heavy_exclamation_mark: As `update` reads and rewrites POSEIDON.yml files, it may change their inner order,
585 layout or even content (e.g. if they have fields which are not in the [Poseidon package definition](#)). Create a backup
586 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.

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

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

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

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

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

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

641 3.2 Summarise command

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

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

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

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

```
646
647 Available options:
648     -h,--help                Show this help text
649     -d,--baseDir DIR         a base directory to search for Poseidon Packages
650                             (could be a Poseidon repository)
651     --raw                    output table as tsv without header. Useful for piping
652                             into grep or awk
```

653 You can run it with

```
654 trident summarise -d ... -d ...
```

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

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

660 3.3 Survey command

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

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

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

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

665

666 Available options:

667 -h,--help Show this help text

668 -d,--baseDir DIR a base directory to search for Poseidon Packages

669 (could be a Poseidon repository)

670 --raw output table as tsv without header. Useful for piping

671 into grep or awk

672 Running

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

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

675 means what.

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

677 **3.4 Validate command**

678 `validate` checks poseidon datasets for structural correctness.

679 Click here for command line details

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

681 Check one or multiple Poseidon packages for structural correctness

682

683 Available options:

684 -h,--help Show this help text

685 -d,--baseDir DIR a base directory to search for Poseidon Packages

686 (could be a Poseidon repository)

687 --ignoreGeno ignore SNP and GenoFile

688 --fullGeno test parsing of all SNPs (by default only the first

689 100 SNPs are probed)

690 --noExitCode do not produce an explicit exit code

691 --ignoreDuplicates do not stop on duplicated individual names in the

692 package collection

693 You can run it with

694 `trident validate -d ... -d ...`

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

696 the issues.

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

698 what is checked:

699 • Presence of the necessary files

700 • Full structural correctness of .bib and .janno file

701 • Superficial correctness of genotype data files by parsing the first 100 SNPs. A full check that parses all

702 SNPs can be run with the `--fullGeno` option

703 • Correspondence of BibTeX keys in .bib and .janno

704 • Correspondence of individual and group IDs in .janno and genotype data files

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

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