

Guide for trident v1.2.0.0 to v1.2.1.0

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	13
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.
```

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
```

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. An exception for that is the **list** subcommand, which will read and report all packages/groups/individuals in all versions.
- 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`

```

154         --snpFile ARG --indFile ARG) [--snpSet ARG]
155         (-o|--outPackagePath ARG) [-n|--outPackageName ARG]
156         [--minimal]
157     Create a new Poseidon package from genotype data
158
159     Available options:
160     -h,--help                Show this help text
161     -p,--genoOne ARG         one of the input genotype data files. Expects .bed or
162                             .bim or .fam for PLINK and .geno or .snp or .ind for
163                             EIGENSTRAT. The other files must be in the same
164                             directory and must have the same base name
165     --inFormat ARG           the format of the input genotype data: EIGENSTRAT or
166                             PLINK (only necessary for data input with --genoFile
167                             + --snpFile + --indFile)
168     --genoFile ARG           the input geno file path
169     --snpFile ARG            the input snp file path
170     --indFile ARG            the input ind file path
171     --snpSet ARG             the snpSet of the package: 1240K, HumanOrigins or
172                             Other. (only relevant for data input with
173                             -p|--genoOne or --genoFile + --snpFile + --indFile,
174                             because the packages in a -d|--baseDir already have
175                             this information in their respective POSEIDON.yml
176                             files) Default: Other
177     -o,--outPackagePath ARG  the output package directory path
178     -n,--outPackageName ARG  the output package name - this is optional: If no
179                             name is provided, then the package name defaults to
180                             the basename of the (mandatory) --outPackagePath
181                             argument
182     --minimal                should only a minimal output package be created?

```

183 The command

```

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

	EIGENSTRAT	PLINK
genoFile	.geno	.bed

	EIGENSTRAT	PLINK
snpFile	.snp	.bim
indFile	.ind	.fam

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

2.2 Fetch command

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

Click here for command line details

```
Usage: trident fetch (-d|--baseDir DIR)
           (--downloadAll |
           (--fetchFile ARG | (-f|--fetchString ARG)))
           [--remoteURL ARG]
```

Download data from a remote Poseidon repository

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>--downloadAll</code>	download all packages the server is offering
<code>--fetchFile ARG</code>	A file with a list of packages. Works just as <code>-f</code> , but multiple values can also be separated by newline, not just by comma. <code>-f</code> and <code>--fetchFile</code> can be combined.
<code>-f,--fetchString ARG</code>	List of packages to be downloaded from the remote server. Package names should be wrapped in asterisks: <code>*package_title*</code> . You can combine multiple values with comma, so for example: <code>"*package_1*, *package_2*, *package_3*"</code> . <code>fetchString</code> uses the same parser as <code>forgeString</code> , but does not allow excludes. If groups or individuals are specified, then packages which include these groups or individuals are included in the download.
<code>--remoteURL ARG</code>	URL of the remote Poseidon server (default: <code>"https://server.poseidon-adna.org"</code>)

It works with

```
trident fetch -d ... -d ... \
-f "*package_title_1*,*package_title_2*,*package_title_3*,group_name,<individual1>"
```

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

234 Entities are specified using a special syntax (see also the documentation of `forge` below): Package titles are
235 wrapped in asterisks: `*package_title*`, group names are spelled as is, and individual names are wrapped in
236 angular brackets, so `<individual1>`. Fetch will figure out which packages need to be downloaded to include all
237 specified entities. `--downloadAll`, which can be given instead of `-f` and `--fetchFile`, causes fetch to download
238 all packages from the server. The downloaded packages are added in the first (!) `-d` directory (which gets created
239 if it doesn't exist), but downloads are only performed if the respective packages are not already present in the
240 latest version in any of the `-d` dirs.

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

243 `fetch` also has the optional arguments `--remote https://..."` to name an alternative poseidon server. The
244 default points to the [DAG server](#).

245 2.3 Forge command

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

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

249 Usage: `trident forge ((-d|--baseDir DIR) |`
250 `((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG`
251 `--snpFile ARG --indFile ARG) [--snpSet ARG])`
252 `[--forgeFile ARG | (-f|--forgeString ARG)]`
253 `[--selectSnps ARG] [--intersect] [--outFormat ARG]`
254 `[--minimal] [--onlyGeno] (-o|--outPackagePath ARG)`
255 `[-n|--outPackageName ARG] [--packagewise]`
256 `[--outPlinkPopName ARG]`
257 Select packages, groups or individuals and create a new Poseidon package from
258 them

260 Available options:

261 <code>-h,--help</code>	Show this help text
262 <code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages 263 (could be a Poseidon repository)
264 <code>-p,--genoOne ARG</code>	one of the input genotype data files. Expects .bed or 265 .bim or .fam for PLINK and .geno or .snp or .ind for 266 EIGENSTRAT. The other files must be in the same 267 directory and must have the same base name
268 <code>--inFormat ARG</code>	the format of the input genotype data: EIGENSTRAT or 269 PLINK (only necessary for data input with <code>--genoFile</code> 270 + <code>--snpFile</code> + <code>--indFile</code>)
271 <code>--genoFile ARG</code>	the input geno file path
272 <code>--snpFile ARG</code>	the input snp file path

273 --indFile ARG the input ind file path

274 --snpSet ARG the snpSet of the package: 1240K, HumanOrigins or

275 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  forge can be used with
355  trident forge -d ... -d ... \
356    -f "*package_name*, group_id, <individual_id>" \
357    -o path/to/new_package_name
358  where the entities (packages, groups/populations, individuals/samples) you want in the output package can be
359  denoted either as a string on the command line (-f/--forgeString), or in an input text file (--forgeFile).
360  See the section below for the syntax of this selection language. Do not forget to wrap the --forgeString query
361  in quotes.

```

362 Including one or multiple Poseidon packages with `-d` is not the only way to include data for a forge operation.
 363 It is also possible to consider unpackaged genotype data directly with `-p (+ --snpSet)` or `--inFormat +`
 364 `--genoFile + --snpFile + --indFile (+ --snpSet)`. This makes the following example possible, where we
 365 merge data from one Poseidon package and two genotype datasets to get a new EIGENSTRAT dataset.

```
366 trident forge \
367   -d 2017_GonzalesFortesCurrentBiology \
368   -p 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \
369   --inFormat PLINK \
370   --genoFile 2017_HaberAJHG/2017_HaberAJHG.bed \
371   --snpFile 2017_HaberAJHG/2017_HaberAJHG.bim \
372   --indFile 2017_HaberAJHG/2017_HaberAJHG.fam \
373   -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \
374   -o testpackage \
375   --outFormat EIGENSTRAT \
376   --onlyGeno
```

377 2.3.1 The forge selection language

378 The text in `--forgeString` and `--forgeFile` are parsed as a domain specific query language that describes
 379 precisely which entities should be compiled in the output package of a given `forge` operation. The language has
 380 multiple syntactic elements and a specific evaluation logic.

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

- 383 • Each package title is surrounded by `*`: `*package*`. That means if you want all individuals of the Poseidon
 384 package `2019_Jeong_InnerEurasia` in the output package you would add `*2019_Jeong_InnerEurasia*`
 385 to the query.
- 386 • Groups/populations are not specially marked: `group`. So to get all individuals of the group
 387 `Swiss_Roman_period`, you would simply add `Swiss_Roman_period`.
- 388 • Individuals/samples are surrounded by `<` and `>`: `<individual>`. `ALA026` therefore becomes `<ALA026>`. A sec-
 389 ond way to denote individuals is with the more verbose and specific syntax `<package:group:individual>`.
 390 Such defined individuals take precedence over differently defined ones (so: directly with `<individual>` or
 391 as a subset of `*package*` or `group`). This allows to resolve duplication issues precisely – at least in cases
 392 where the duplicated individuals differ in source package or primary group.

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

```
395 # Packages
396 *package1*, *package2*
397
398 # Groups and individuals from other packages beyond package1 and package2
399 group1, <individual1>, group2, <individual2>, <individual3>
400
401 # group2 has two outlier individuals that should be ignored
402 -<bad_individual1> # This one has very low coverage
403 -<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

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn1	AdditionalColumn2	AdditionalColumn3
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 a an optional flag `--intersect`, that defines, if the genotype data from different packages should be merged with an **union** or an **intersect** operation. The default (if this option is not set) is to output the union of all SNPs, with genotypes defined as missing in samples from packages which do not have a SNP that is present in another package. With this option set, on the other hand, the forged dataset will typically have fewer SNPs, but less missingness.

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

Input snpSet A	Input snpSet B	<code>--intersect</code>	Ouput snpSet
Other	*	*	Other
1240K	HumanOrigins	True	HumanOrigins
1240K	HumanOrigins	False	1240K

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

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

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

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

2.4 Genoconvert command

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

[Click here for command line details](#)

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

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

Available options:

-h,--help	Show this help text
-d,--baseDir DIR	a base directory to search for Poseidon Packages (could be a Poseidon repository)
-p,--genoOne ARG	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
--inFormat ARG	the format of the input genotype data: EIGENSTRAT or PLINK (only necessary for data input with --genoFile + --snpFile + --indFile)
--genoFile ARG	the input geno file path
--snpFile ARG	the input snp file path
--indFile ARG	the input ind file path
--snpSet ARG	the snpSet of the package: 1240K, HumanOrigins or Other. (only relevant for data input with -p --genoOne or --genoFile + --snpFile + --indFile, because the packages in a -d --baseDir already have this information in their respective POSEIDON.yml files) Default: Other
--outFormat ARG	the format of the output genotype data: EIGENSTRAT or PLINK.
--onlyGeno	should only the resulting genotype data be returned? This means the output will not be a Poseidon package
-o,--outPackagePath ARG	the output package directory path - this is optional: If no path is provided, then the output is written to the directories where the input genotype data file (.bed/.geno) is stored
--removeOld	Remove the old genotype files when creating the new ones
--outPlinkPopName ARG	Where to write the population/group name into the FAM file in Plink-format. Three options are possible:

```

500             asFamily (default) | asPhenotype | asBoth. See also
501             --inPlinkPopName.

```

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 Available options:

527	-h,--help	Show this help text
528	-d,--baseDir DIR	a base directory to search for Poseidon Packages
529		(could be a Poseidon repository)
530	--poseidonVersion ARG	Poseidon version the packages should be updated to:
531		e.g. "2.5.3" (default: Nothing)
532	--ignorePoseidonVersion	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	--versionComponent ARG	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 --noChecksumUpdate      Should update of checksums in the POSEIDON.yml file
542                          be skipped
543 --ignoreGeno            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.

587 3 Inspection commands

588 3.1 List command

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

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

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

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

597 Available options:

599 -h,--help	Show this help text
600 -d,--baseDir DIR	a base directory to search for Poseidon Packages 601 (could be a Poseidon repository)
602 --remote	list packages from a remote server instead the local 603 file system
604 --remoteURL ARG	URL of the remote Poseidon server 605 (default: "https://server.poseidon-adna.org")
606 --packages	list all packages
607 --groups	list all groups, ignoring any group names after the 608 first as specified in the Janno-file
609 --individuals	list individuals
610 -j,--jannoColumn JANNO_HEADER	list additional fields from the janno files, using 611 the Janno column heading name, such as Country, Site, 612 Date_C14_Uncal_BP, Endogenous, ...
613 --raw	output table as tsv without header. Useful for piping 614 into grep or awk

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

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

618 This will yield a nicely formatted table of all packages, their last update and the number of individuals in it.

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

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

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

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

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

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

632 3.2 Summarise command

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

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

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

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

637
638 Available options:
639 -h,--help Show this help text
640 -d,--baseDir DIR a base directory to search for Poseidon Packages
641 (could be a Poseidon repository)
642 --raw output table as tsv without header. Useful for piping
643 into grep or awk

644 You can run it with

645 trident summarise -d ... -d ...

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

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

651 3.3 Survey command

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

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

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

655 Survey the degree of context information completeness for Poseidon packages

656
657 Available options:
658 -h,--help Show this help text
659 -d,--baseDir DIR a base directory to search for Poseidon Packages
660 (could be a Poseidon repository)

661 --raw output table as tsv without header. Useful for piping
662 into grep or awk

663 Running

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

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

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

668 3.4 Validate command

669 `validate` checks poseidon datasets for structural correctness.

670 Click here for command line details

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

672 Check one or multiple Poseidon packages for structural correctness

673

674 Available options:

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

684 You can run it with

685 `trident validate -d ... -d ...`

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

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

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

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

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