

# Contents

0.1	Guide for trident v1.3.0.4	1
0.1.1	The trident CLI	1
0.1.2	Package creation and manipulation commands	4
0.1.3	Inspection commands	15

## 0.1 Guide for trident v1.3.0.4

### 0.1.1 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 MODE | --debug] [--errLength INT]
        [--inPlinkPopName MODE] (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 MODE</code>	How information should be reported: NoLog, SimpleLog, DefaultLog, ServerLog or VerboseLog. (default: DefaultLog)
<code>--debug</code>	Short for <code>--logMode VerboseLog</code> .
<code>--errLength INT</code>	After how many characters should a potential error message be truncated. "Inf" for no truncation. (default: CharCount 1500)
<code>--inPlinkPopName MODE</code>	Where to read the population/group name from the FAM file in Plink-format. Three options are possible: asFamily (default)   asPhenotype   asBoth.

#### Package creation and manipulation commands:

<code>init</code>	Create a new Poseidon package from genotype data
<code>fetch</code>	Download data from a remote Poseidon repository
<code>forge</code>	Select packages, groups or individuals and create a new Poseidon package from them
<code>genoconvert</code>	Convert the genotype data in a Poseidon package to a different file format
<code>rectify</code>	Adjust POSEIDON.yml files automatically to package changes

#### Inspection commands:

<code>list</code>	List packages, groups or individuals from local or
-------------------	--

```

43         remote Poseidon repositories
44     summarise      Get an overview over the content of one or multiple
45                    Poseidon packages
46     survey         Survey the degree of context information completeness
47                    for Poseidon packages
48     validate       Check Poseidon packages or package components for
49                    structural correctness

```

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

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

```

57 /path/to/poseidon/packages
58     /modern
59         /2019_poseidon_package1
60         /2019_poseidon_package2
61     /ancient
62         /...
63         /...
64     /Reference_Genomes
65         /...
66         /...

```

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

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

```

73 ~/my_project/my_project.geno
74 ~/my_project/my_project.snp
75 ~/my_project/my_project.ind

```

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

```

78 poseidonVersion: 2.7.1
79 title: My_awesome_project
80 description: Unpublished genetic data from my awesome project
81 contributor:
82     - name: Stephan Schiffels
83       email: schiffels@institute.org
84 packageVersion: 0.1.0

```

```

85  lastModified: 2020-10-07
86  genotypeData:
87    format: EIGENSTRAT
88    genoFile: my_project.geno
89    snpFile: my_project.snp
90    indFile: my_project.ind
91    jannoFile: my_project.janno
92    bibFile: sources.bib

```

93 Two remarks: 1) all file paths are considered *relative* to the directory in which POSEIDON.yml resides. For this  
 94 example we assume that this file is added into the same directory as the three genotype files. 2) Besides the  
 95 genotype data files there are two (technically optional) files referenced by this example POSEIDON.yml file:  
 96 **sources.bib** and **my\_project.janno**. Of course you can add them manually - **init** automatically creates empty  
 97 dummy versions.

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

```

100 trident list -d /path/to/poseidon/packages/modern \
101   -d /path/to/poseidon/packages/ReferenceGenomes
102   -d ~/my_project --packages

```

### 103 0.1.1.1 General notes

104 **0.1.1.1.1 Logging and command line output** For all subcommands the general argument **--logMode**  
 105 defines how trident reports messages (to stderr) on the command line:

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

112 **--debug** is short for **--logMode VerboseLog** to activate this important log level more easily.

### 113 0.1.1.1.2 Duplicates

- 114 • If multiple packages in a package repository share the same **title**, then trident will try to select the  
 115 one with the highest version number. If this is not sufficient to resolve the conflict, trident will stop. An  
 116 exception for that is the **list** subcommand, which will read and report all packages/groups/individuals in  
 117 all versions.
- 118 • Individual/sample names (**Poseidon\_IDs**) within one package have to be unique, or trident will stop.
- 119 • We generally also discourage ID duplicates across packages in package repositories, but trident will generally  
 120 continue with them after printing a warning. This does not apply for **validate**, by default (you can  
 121 change this behaviour with **--ignoreDuplicates**), and **forge**. **forge** offers a special mechanism to resolve  
 122 duplicates within its selection language (see below).

123 **0.1.1.1.3 Group names in .fam files** The .fam file of Plink-formatted genotype data is used inconsistently  
 124 across different popular aDNA software tools to store group/population name information. The (global) option  
 125 `--inPlinkPopName` with the arguments `asFamily` (default), `asPhenotype` and `asBoth` allows to control the  
 126 reading of the population name from Plink .fam files. The subcommands that write genotype data (`forge`,  
 127 `genoconvert`) have a corresponding option `--outPlinkPopName` to specify this for the output.

128 **0.1.1.1.4 Whitespaces in the .janno file** While reading the .janno file `trident` trims all leading and  
 129 trailing whitespaces around individual cells. Also all instances of the `No-Break Space` unicode character will be  
 130 removed. This means these whitespaces will not be preserved when a package is `forged`.

## 131 0.1.2 Package creation and manipulation commands

132 **0.1.2.1 Init command** `init` creates a new, valid Poseidon package from genotype data files. It adds a valid  
 133 `POSEIDON.yml` file, a dummy .janno file for context information and an empty .bib file for literature references.

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

```
135 Usage: trident init ((-p|--genoOne FILE) | --inFormat FORMAT --genoFile FILE
136                  --snpFile FILE --indFile FILE) [--snpSet SET]
137                  (-o|--outPackagePath DIR) [-n|--outPackageName STRING]
138                  [--minimal]
```

139  
 140 Create a new Poseidon package from genotype data

141  
 142 Available options:

143	<code>-h,--help</code>	Show this help text
144	<code>-p,--genoOne FILE</code>	One of the input genotype data files. Expects .bed, 145 .bim or .fam for PLINK and .geno, .snp or .ind for 146 EIGENSTRAT. The other files must be in the same 147 directory and must have the same base name.
148	<code>--inFormat FORMAT</code>	The format of the input genotype data: EIGENSTRAT or 149 PLINK. Only necessary for data input with <code>--genoFile</code> 150 + <code>--snpFile</code> + <code>--indFile</code> .
151	<code>--genoFile FILE</code>	Path to the input geno file.
152	<code>--snpFile FILE</code>	Path to the input snp file.
153	<code>--indFile FILE</code>	Path to the input ind file.
154	<code>--snpSet SET</code>	The snpSet of the package: 1240K, HumanOrigins or 155 Other. Only relevant for data input with <code>-p --genoOne</code> 156 or <code>--genoFile</code> + <code>--snpFile</code> + <code>--indFile</code> , because the 157 packages in a <code>-d --baseDir</code> already have this 158 information in their respective <code>POSEIDON.yml</code> files. 159 (default: Other)
160	<code>-o,--outPackagePath DIR</code>	Path to the output package directory.
161	<code>-n,--outPackageName STRING</code>	
162		The output package name. This is optional: If no name 163 is provided, then the package name defaults to the 164 basename of the (mandatory) <code>--outPackagePath</code>

```

165         argument. (default: Nothing)
166     --minimal           Should the output data be reduced to a necessary
167                         minimum and omit empty scaffolding?

```

168 The command

```

169 trident init \
170     --inFormat EIGENSTRAT/PLINK \
171     --genoFile path/to/geno_file \
172     --snpFile path/to/snp_file \
173     --indFile path/to/ind_file \
174     --snpSet 1240K|HumanOrigins|Other \
175     -o path/to/new_package_name

```

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

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

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

182 **0.1.2.2 Fetch command** `fetch` allows to download Poseidon packages from a remote Poseidon server via a  
183 [Web API](#). Read more about the data available with it [here](#).

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

```

185 Usage: trident fetch (-d|--baseDir DIR)
186         (--downloadAll |
187         (--fetchFile FILE | (-f|--fetchString DSL)))
188         [--remoteURL URL] [--archive STRING]

```

190 Download data from a remote Poseidon repository

192 Available options:

```

193     -h,--help           Show this help text
194     -d,--baseDir DIR    A base directory to search for Poseidon packages.
195     --downloadAll       Download all packages the server is offering.
196     --fetchFile FILE    A file with a list of packages. Works just as -f, but
197                         multiple values can also be separated by newline, not
198                         just by comma. -f and --fetchFile can be combined.
199     -f,--fetchString DSL List of packages to be downloaded from the remote
200                         server. Package names should be wrapped in asterisks:

```

201           \*package\_title\*. You can combine multiple values with  
202           comma, so for example: "\*package\_1\*, \*package\_2\*,  
203           \*package\_3\*". fetchString uses the same parser as  
204           forgeString, but does not allow excludes. If groups  
205           or individuals are specified, then packages which  
206           include these groups or individuals are included in  
207           the download.

208   --remoteURL URL           URL of the remote Poseidon server.  
209                               (default: "https://server.poseidon-adna.org")

210   --archive STRING          The name of the Poseidon package archive that should  
211                               be queried. If not given, then the query falls back  
212                               to the default archive of the server selected with  
213                               --remoteURL. See the archive documentation at  
214                               [https://www.poseidon-adna.org/#/archive\\_overview](https://www.poseidon-adna.org/#/archive_overview) for  
215                               a list of archives currently available from the  
216                               official Poseidon Web API. (default: Nothing)

217   It works with

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

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

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

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

232 **fetch** also has the optional arguments --remote <https://...> to name an alternative Poseidon server and  
233 --archive to select a Poseidon archive on the server. Here is a list of the [archives available on the official](#)  
234 [Poseidon server](#).

235 **0.1.2.3 Forge command** **forge** creates new Poseidon packages by extracting and merging packages,  
236 populations and individuals from your Poseidon repositories.

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

238 Usage: trident forge ((-d|--baseDir DIR) |  
239                       ((-p|--genoOne FILE) | --inFormat FORMAT --genoFile FILE  
240                       --snpFile FILE --indFile FILE) [--snpSet SET])  
241                       [--forgeFile FILE | (-f|--forgeString DSL)]

```

242         [--selectSnps FILE] [--intersect] [--outFormat FORMAT]
243         [--minimal] [--onlyGeno] (-o|--outPackagePath DIR)
244         [-n|--outPackageName STRING] [--packagewise]
245         [--outPlinkPopName MODE]
246
247     Select packages, groups or individuals and create a new Poseidon package from
248     them
249
250     Available options:
251     -h,--help                Show this help text
252     -d,--baseDir DIR         A base directory to search for Poseidon packages.
253     -p,--genoOne FILE        One of the input genotype data files. Expects .bed,
254                             .bim or .fam for PLINK and .geno, .snp or .ind for
255                             EIGENSTRAT. The other files must be in the same
256                             directory and must have the same base name.
257     --inFormat FORMAT        The format of the input genotype data: EIGENSTRAT or
258                             PLINK. Only necessary for data input with --genoFile
259                             + --snpFile + --indFile.
260     --genoFile FILE          Path to the input geno file.
261     --snpFile FILE           Path to the input snp file.
262     --indFile FILE           Path to the input ind file.
263     --snpSet SET             The snpSet of the package: 1240K, HumanOrigins or
264                             Other. Only relevant for data input with -p|--genoOne
265                             or --genoFile + --snpFile + --indFile, because the
266                             packages in a -d|--baseDir already have this
267                             information in their respective POSEIDON.yml files.
268                             (default: Other)
269     --forgeFile FILE         A file with a list of packages, groups or individual
270                             samples. Works just as -f, but multiple values can
271                             also be separated by newline, not just by comma.
272                             Empty lines are ignored and comments start with "#",
273                             so everything after "#" is ignored in one line.
274                             Multiple instances of -f and --forgeFile can be
275                             given. They will be evaluated according to their
276                             input order on the command line.
277     -f,--forgeString DSL     List of packages, groups or individual samples to be
278                             combined in the output package. Packages follow the
279                             syntax *package_title*, populations/groups are simply
280                             group_id and individuals <individual_id>. You can
281                             combine multiple values with comma, so for example:
282                             "*package_1*, <individual_1>, <individual_2>,
283                             group_1". Duplicates are treated as one entry.
284                             Negative selection is possible by prepending "-" to
285                             the entity you want to exclude (e.g. "*package_1*,
286                             -<individual_1>, -group_1"). forge will apply

```

287 excludes and includes in order. If the first entity  
 288 is negative, then forge will assume you want to merge  
 289 all individuals in the packages found in the baseDirs  
 290 (except the ones explicitly excluded) before the  
 291 exclude entities are applied. An empty forgeString  
 292 (and no --forgeFile) will therefore merge all  
 293 available individuals. If there are individuals in  
 294 your input packages with equal individual id, but  
 295 different main group or source package, they can be  
 296 specified with the special syntax  
 297 "<package:group:individual>".

298 --selectSnps FILE To extract specific SNPs during this forge operation,  
 299 provide a Snp file. Can be either Eigenstrat (file  
 300 ending must be '.snp') or Plink (file ending must be  
 301 '.bim'). When this option is set, the output package  
 302 will have exactly the SNPs listed in this file. Any  
 303 SNP not listed in the file will be excluded. If  
 304 option '--intersect' is also set, only the SNPs  
 305 overlapping between the SNP file and the forged  
 306 packages are output. (default: Nothing)

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

314 --outFormat FORMAT The format of the output genotype data: EIGENSTRAT or  
 315 PLINK. (default: PLINK)

316 --minimal Should the output data be reduced to a necessary  
 317 minimum and omit empty scaffolding?

318 --onlyGeno Should only the resulting genotype data be returned?  
 319 This means the output will not be a Poseidon package.

320 -o,--outPackagePath DIR Path to the output package directory.

321 -n,--outPackageName STRING  
 322 The output package name. This is optional: If no name  
 323 is provided, then the package name defaults to the  
 324 basename of the (mandatory) --outPackagePath  
 325 argument. (default: Nothing)

326 --packagewise Skip the within-package selection step in forge. This  
 327 will result in outputting all individuals in the  
 328 relevant packages, and hence a superset of the  
 329 requested individuals/groups. It may result in better  
 330 performance in cases where one wants to forge entire  
 331 packages or almost entire packages. Details: Forge



```

332         conceptually performs two types of selection: First,
333         it identifies which packages in the supplied base
334         directories are relevant to the requested forge, i.e.
335         whether they are either explicitly listed using
336         *PackageName*, or because they contain selected
337         individuals or groups. Second, within each relevant
338         package, individuals which are not requested are
339         removed. This option skips only the second step, but
340         still performs the first.
341     --outPlinkPopName MODE Where to write the population/group name into the FAM
342         file in Plink-format. Three options are possible:
343         asFamily (default) | asPhenotype | asBoth. See also
344         --inPlinkPopName.
345
346     forge can be used with
347
348     trident forge -d ... -d ... \
349         -f "*package_name*, group_id, <individual_id>" \
350         -o path/to/new_package_name
351
352     where the entities (packages, groups/populations, individuals/samples) you want in the output package can be
353     denoted either as a string on the command line (-f/--forgeString), or in an input text file (--forgeFile).
354     See the section below for the syntax of this selection language. Do not forget to wrap the --forgeString query
355     in quotes.
356
357     Including one or multiple Poseidon packages with -d is not the only way to include data for a forge operation.
358     It is also possible to consider unpackaged genotype data directly with -p (+ --snpSet) or --inFormat +
359     --genoFile + --snpFile + --indFile (+ --snpSet). This makes the following example possible, where we
360     merge data from one Poseidon package and two genotype datasets to get a new EIGENSTRAT dataset.
361
362     trident forge \
363         -d 2017_GonzalesFortesCurrentBiology \
364         -p 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \
365         --inFormat PLINK \
366         --genoFile 2017_HaberAJHG/2017_HaberAJHG.bed \
367         --snpFile 2017_HaberAJHG/2017_HaberAJHG.bim \
368         --indFile 2017_HaberAJHG/2017_HaberAJHG.fam \
369         -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \
370         -o testpackage \
371         --outFormat EIGENSTRAT \
372         --onlyGeno

```

**0.1.2.3.1 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

374 package 2019\_Jeong\_InnerEurasia in the output package you would add `*2019_Jeong_InnerEurasia*`  
 375 to the query.

- 376 • Groups/populations are not specially marked: `group`. So to get all individuals of the group  
 377 `Swiss_Roman_period`, you would simply add `Swiss_Roman_period`.
- 378 • Individuals/samples are surrounded by `<` and `>`: `<individual>`. ALA026 therefore becomes `<ALA026>`. A sec-  
 379 ond way to denote individuals is with the more verbose and specific syntax `<package:group:individual>`.  
 380 Such defined individuals take precedence over differently defined ones (so: directly with `<individual>` or  
 381 as a subset of `*package*` or `group`). This allows to resolve duplication issues precisely – at least in cases  
 382 where the duplicated individuals differ in source package or primary group.

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

```
385 # Packages
386 *package1*, *package2*
387
388 # Groups and individuals from other packages beyond package1 and package2
389 group1, <individual1>, group2, <individual2>, <individual3>
390
391 # group2 has two outlier individuals that should be ignored
392 -<bad_individual1> # This one has very low coverage
393 -<bad_individual2> # This one is from a different time period
```

394 By prepending `-` to the bad individuals, we can exclude them from the forged package. `forge` fig-  
 395 ures out the final list of samples to include by executing all `forge`-entities in order. So an entity list  
 396 `*PackageA*, -<Individual1>, GroupA` may result in a different outcome than `*PackageA*, GroupA, -<Individual1>`,  
 397 depending on whether `<Individual1>` belongs to `GroupA` or not. If the `forge` entity list starts with a negative  
 398 entity, or if the entity list is empty, `forge` will implicitly assume you want to include all individuals in all  
 399 packages found in the `baseDirs` (except the ones explicitly excluded, of course).

400 An empty `forgeString` will therefore merge all available individuals.

401 **0.1.2.3.2 Treatment of the .janno file while merging** `forge` merges and subsets `.janno` files along with  
 402 the genotype data. If a package lacks a `.janno` file, then a basic one will be created internally based on the  
 403 information in the genotype data, and used for the output. Missing columns across packages will be filled with  
 404 `n/a`.

405 For merging two `.janno` files **A** and **B** the following rules apply regarding undefined, arbitrary additional columns:

- 406 • If **A** has an additional column which is not in **B** then empty cells in the rows imported from **B** are filled  
 407 with `n/a`.
- 408 • If **A** and **B** share additional columns with identical column name, then they are treated as semantically  
 409 identical units and merged accordingly.
- 410 • In the resulting `.janno` file, all additional columns from both **A** and **B** are sorted alphabetically and  
 411 appended after the normal, specified variables.

412 The following example illustrates the described behaviour:

413 **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

**0.1.2.3.3 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.

**0.1.2.3.4 Treatment of the .bib file while merging** In the forge process all relevant samples for the output package are determined. This includes their .janno entries and therefore the information on the publication keys documented for them in the .janno `Publication` column. The output .bib file compiles only the relevant references for the samples in the output package. It includes the references exactly once and is sorted alphabetically (by key).

**0.1.2.3.5 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.

**0.1.2.4 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 FILE) | --inFormat FORMAT
                           --genoFile FILE --snpFile FILE --indFile FILE)
                           [--snpSet SET]) --outFormat FORMAT [--onlyGeno]
                           [-o|--outPackagePath DIR] [--removeOld]
                           [--outPlinkPopName MODE]
```

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

Available options:

<code>-h,--help</code>	Show this help text
<code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
<code>-p,--genoOne FILE</code>	One of the input genotype data files. Expects <code>.bed</code> , <code>.bim</code> or <code>.fam</code> for PLINK and <code>.geno</code> , <code>.snp</code> or <code>.ind</code> for EIGENSTRAT. The other files must be in the same directory and must have the same base name.
<code>--inFormat FORMAT</code>	The format of the input genotype data: EIGENSTRAT or

```

469         PLINK. Only necessary for data input with --genoFile
470         + --snpFile + --indFile.
471     --genoFile FILE      Path to the input geno file.
472     --snpFile FILE       Path to the input snp file.
473     --indFile FILE       Path to the input ind file.
474     --snpSet SET         The snpSet of the package: 1240K, HumanOrigins or
475                         Other. Only relevant for data input with -p|--genoOne
476                         or --genoFile + --snpFile + --indFile, because the
477                         packages in a -d|--baseDir already have this
478                         information in their respective POSEIDON.yml files.
479                         (default: Other)
480     --outFormat FORMAT   the format of the output genotype data: EIGENSTRAT or
481                         PLINK.
482     --onlyGeno           Should only the resulting genotype data be returned?
483                         This means the output will not be a Poseidon package.
484     -o,--outPackagePath DIR Path to the output package directory. This is
485                         optional: If no path is provided, then the output is
486                         written to the directories where the input genotype
487                         data file (.bed/.geno) is stored. (default: Nothing)
488     --removeOld          Remove the old genotype files when creating the new
489                         ones.
490     --outPlinkPopName MODE Where to write the population/group name into the FAM
491                         file in Plink-format. Three options are possible:
492                         asFamily (default) | asPhenotype | asBoth. See also
493                         --inPlinkPopName.
494
495     With the default setting
496
497     trident genoconvert -d ... -d ... --outFormat EIGENSTRAT|PLINK
498
499     all packages in -d will be converted to the desired --outFormat (either EIGENSTRAT or PLINK), if the data is
500     not already in this format. This includes updating the respective POSEIDON.yml files.
501
502     The “old” data is not deleted, but kept around. That means conversion can result in a package with both PLINK
503     and EIGENSTRAT data, but only one is linked in the POSEIDON.yml file, and that is what will be used by
504     trident. To delete the old data in the conversion you can add the --removeOld flag.
505
506     Instead of -d to change Poseidon packages, the -p (+ --snpSet) or --inFormat + --genoFile + --snpFile
507     + --indFile (+ --snpSet) allow to directly convert genotype data that is not wrapped in a Poseidon package
508     and store it to a directory given in -o. See this example:
509
510     trident genoconvert \
511         -p 2018_Mittnik_Baltic/Mittnik_Baltic.bed \
512         --outFormat EIGENSTRAT
513         -o my_directory
514
515 0.1.2.5 Rectify command rectify automatically harmonizes POSEIDON.yml files of one or multiple
516 packages. This is not an automatic update from one Poseidon version to the next, but rather a clean-up wizard
517 after manual modifications.

```

511 Click here for command line details

```
512 Usage: trident rectify (-d|--baseDir DIR) [--ignorePoseidonVersion]
513         [--poseidonVersion ?.??.?]
514         [--packageVersion VPART [--logText STRING]]
515         [--checksumAll | [--checksumGeno] [--checksumJanno]
516         [--checksumSSF] [--checksumBib]]
517         [--newContributors DSL]
```

518  
519 Adjust POSEIDON.yml files automatically to package changes

520  
521 Available options:

```
522  -h,--help           Show this help text
523  -d,--baseDir DIR    A base directory to search for Poseidon packages.
524  --ignorePoseidonVersion Read packages even if their poseidonVersion is not
525                      compatible with trident.
526  --poseidonVersion ?.??.? Poseidon version the packages should be updated to:
527                      e.g. "2.5.3".
528  --packageVersion VPART Part of the package version number in the
529                      POSEIDON.yml file that should be updated: Major,
530                      Minor or Patch (see https://semver.org).
531  --logText STRING    Log text for this version in the CHANGELOG file.
532  --checksumAll       Update all checksums.
533  --checksumGeno      Update genotype data checksums.
534  --checksumJanno     Update .janno file checksum.
535  --checksumSSF       Update .ssf file checksum
536  --checksumBib       Update .bib file checksum.
537  --newContributors DSL Contributors to add to the POSEIDON.yml file in the
538                      form "[Firstname Lastname](Email address);..."
```

539 It can be called with a lot of optional arguments:

```
540 trident rectify -d ... -d ... \
541   --poseidonVersion "X.X.X" \
542   --packageVersion Major|Minor|Patch \
543   --logText "short description of the update"
544   --checksumAll
545   --newContributors "[Firstname Lastname](Email address);..."
```

546 These arguments determine which fields of the POSEIDON.yml file should be modified.

- 547 • `--poseidonVersion` allows a simple change of the `poseidonVersion` field in the POSEIDON.yml file.
- 548 • `--packageVersion` increments the package version number in the first, the second or the third position.  
549 It can optionally be called with `--logText`, which appends an entry to the CHANGELOG file for the  
550 respective package version update. `--logText` also creates a new CHANGELOG file if it does not exist  
551 yet.
- 552 • `--checksumGeno`, `--checksumJanno`, `--checksumSSF` and `--checksumBib` add or modify the respective  
553 checksum fields in the POSEIDON.yml file. `--checksumAll` is a wrapper to call all of them at once.

554 • `--newContributors` adds new contributors.

555 :warning: As `rectify` reads and rewrites `POSEIDON.yml` files, it may change their inner order, layout or  
556 even content (e.g. if they have fields which are not in the **POSEIDON.yml definition**). Create a backup of the  
557 `POSEIDON.yml` file before running `rectify` if you are uncertain if this might affect you negatively.

### 558 0.1.3 Inspection commands

559 **0.1.3.1 List command** `list` lists packages, groups and individuals of the datasets you use, or of the  
560 packages available on the server.

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

```
562 Usage: trident list ((-d|--baseDir DIR) | --remote [--remoteURL URL]
563                  [--archive STRING])
564                  (--packages | --groups | --individuals
565                  [-j|--jannoColumn COLNAME]) [--raw]
```

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

570 Available options:

571 <code>-h,--help</code>	Show this help text
572 <code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
573 <code>--remote</code>	List packages from a remote server instead the local 574 file system.
575 <code>--remoteURL URL</code>	URL of the remote Poseidon server. 576 (default: "https://server.poseidon-adna.org")
577 <code>--archive STRING</code>	The name of the Poseidon package archive that should 578 be queried. If not given, then the query falls back 579 to the default archive of the server selected with 580 <code>--remoteURL</code> . See the archive documentation at 581 <a href="https://www.poseidon-adna.org/#/archive_overview">https://www.poseidon-adna.org/#/archive_overview</a> for 582 a list of archives currently available from the 583 official Poseidon Web API. (default: Nothing)
584 <code>--packages</code>	List all packages.
585 <code>--groups</code>	List all groups, ignoring any group names after the 586 first as specified in the <code>.janno-file</code> .
587 <code>--individuals</code>	List all individuals/samples.
588 <code>-j,--jannoColumn COLNAME</code>	List additional fields from the janno files, using 589 the <code>.janno</code> column heading name, such as "Country", 590 "Site", "Date_C14_Uncal_BP", etc..
591 <code>--raw</code>	Return the output table as tab-separated values 592 without header. This is useful for piping into <code>grep</code> 593 or <code>awk</code> .

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

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

596 This will yield a nicely formatted table of all packages, their version and the number of individuals in them.

597 You can use `--remote` to show packages on the remote server. For example

598 `trident list --packages --remote --archive "community-archive"`

599 will result in a view of all packages available in one of the [public online archives](#). Just as for `fetch`, the `--archive`

600 flag allows to choose which public archive to query.

601 Independent of whether you query a local or an online archive, you can not just list packages, but also groups,

602 as defined in the third column of EIGENSTRAT `.ind` files (or the first/last column of a PLINK `.fam` file), and

603 individuals with the flags `--groups` and `--individuals` (instead of `--packages`).

604 The `--individuals` flag additionally provides a way to immediately access information from `.janno` files

605 on the command line. This works with the `-j/--jannoColumn` option. For example adding `-j Country -j`

606 `Date_C14_Uncal_BP` to the commands above will add the `Country` and the `Date_C14_Uncal_BP` columns to the

607 respective output tables.

608 Note that if you want a less fancy table, for example because you want to load this into Excel, or pipe into

609 another command that cannot deal with the table layout, you can use the `--raw` option to output that table as

610 a simple tab-delimited stream.

611 **0.1.3.2 Summarise command** `summarise` prints some general summary statistics for a given poseidon

612 dataset taken from the `.janno` files.

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

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

615

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

617

618 Available options:

619 <code>-h,--help</code>	Show this help text
620 <code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
621 <code>--raw</code>	Return the output table as tab-separated values
622	without header. This is useful for piping into <code>grep</code>
623	or <code>awk</code> .

624 You can run it with

625 `trident summarise -d ... -d ...`

626 which will show you context information like – among others – the number of individuals in the dataset, their

627 sex distribution, the mean age of the samples (for ancient data) or the mean coverage on the 1240K SNP array

628 in a table. `summarise` depends on complete `.janno` files and will silently ignore missing information.

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

630 **0.1.3.3 Survey command** `survey` tries to indicate package completeness (mostly focused on `.janno` files)

631 for poseidon datasets.

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



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

634

635 Survey the degree of context information completeness for Poseidon packages

636

637 Available options:

638	-h,--help	Show this help text
639	-d,--baseDir DIR	A base directory to search for Poseidon packages.
640	--raw	Return the output table as tab-separated values
641		without header. This is useful for piping into grep
642		or awk.

643 Running

644 trident survey -d ... -d ...

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

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

648 **0.1.3.4 Validate command** validate checks Poseidon packages and individual package components for  
649 structural correctness.

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

651 Usage: trident validate ((-d|--baseDir DIR) [--ignoreGeno] [--fullGeno]  
652 [--ignoreDuplicates] [-c|--ignoreChecksums]  
653 [--ignorePoseidonVersion] |  
654 --pyml FILE | (-p|--genoOne FILE) | --inFormat FORMAT  
655 --genoFile FILE --snpFile FILE --indFile FILE |  
656 --janno FILE | --ssf FILE | --bib FILE) [--noExitCode]

657

658 Check Poseidon packages or package components for structural correctness

659

660 Available options:

661	-h,--help	Show this help text
662	-d,--baseDir DIR	A base directory to search for Poseidon packages.
663	--ignoreGeno	Ignore snp and geno file.
664	--fullGeno	Test parsing of all SNPs (by default only the first
665		100 SNPs are probed).
666	--ignoreDuplicates	Do not stop on duplicated individual names in the
667		package collection.
668	-c,--ignoreChecksums	Whether to ignore checksums. Useful for speedup in
669		debugging.
670	--ignorePoseidonVersion	Read packages even if their poseidonVersion is not
671		compatible with trident.
672	--pyml FILE	Path to a POSEIDON.yml file.
673	-p,--genoOne FILE	One of the input genotype data files. Expects .bed,
674		.bim or .fam for PLINK and .geno, .snp or .ind for

675 EIGENSTRAT. The other files must be in the same  
676 directory and must have the same base name.

677 `--inFormat FORMAT` The format of the input genotype data: EIGENSTRAT or  
678 PLINK. Only necessary for data input with `--genoFile`  
679 `+ --snpFile + --indFile`.

680 `--genoFile FILE` Path to the input geno file.  
681 `--snpFile FILE` Path to the input snp file.  
682 `--indFile FILE` Path to the input ind file.  
683 `--janno FILE` Path to a .janno file.  
684 `--ssf FILE` Path to a .ssf file.  
685 `--bib FILE` Path to a .bib file.  
686 `--noExitCode` Do not produce an explicit exit code.

687 You can run it with

688 `trident validate -d ... -d ...`

689 to check packages and it will either report a success (`Validation passed`) or failure with specific error messages.

690 Instead of validating entire packages with `-d` you can also apply it to individual files and package com-  
691 ponents: `--pym1` (POSEIDON.yml), `-p | --inFormat + --genoFile + --snpFile + --indFile` (genotype  
692 data), `--janno` (.janno file), `--ssf` (.ssf file) or `--bib` (.bib file). In this case `validate` attempts to read and  
693 parse the respective files individually and reports any issues it encounters. Note that this considers the files in  
694 isolation and does not include any cross-file consistency checks.

695 When applied to packages, `validate` tries to ensure that each package adheres to the `schema definition`. Here is  
696 a list of what is checked:

- 697 • Structural correctness of the POSEIDON.yml file.
- 698 • Presence of all files references in the POSEIDON.yml file.
- 699 • Full structural correctness of .janno, .ssf and .bib file.
- 700 • Superficial correctness of genotype data files by parsing the first 100 SNPs. A full check that parses all  
701 SNPs can be triggered with the `--fullGeno` option. `--ignoreGeno`, on the other hand, causes `validate`  
702 to ignore the genotype data entirely, which speeds up the validation significantly.
- 703 • Correspondence of BibTeX keys in .bib and .janno
- 704 • Correspondence of sample IDs in .janno and .ssf.
- 705 • Correspondence of sample and group IDs in .janno and genotype data files.

706 In fact much of this validation already runs as part of the general package reading pipeline invoked for other  
707 trident subcommands (e.g. `forge`). `validate` is meant to be more thorough/brittle, though, and will explicitly  
708 fail if even a single package is broken. For special cases more flexibility can be enabled with the options  
709 `--ignoreDuplicates`, `--ignoreChecksums` and `--ignorePoseidonVersion`.

710 Remember to run `validate` it with `--debug` to get more information in case the default output is not sufficient  
711 to analyse an issue.