

# Guide for trident v1.4.1.0

## Contents

<b>1</b>	<b>Installation</b>	<b>1</b>
<b>2</b>	<b>The trident CLI</b>	<b>2</b>
2.1	General notes . . . . .	4
2.1.1	Logging and command line output . . . . .	4
2.1.2	Package duplicates and versions . . . . .	4
2.1.3	Individual/sample duplicates . . . . .	4
2.1.4	Group names in .fam files . . . . .	5
2.1.5	Whitespaces in the .janno file . . . . .	5
<b>3</b>	<b>Package creation and manipulation commands</b>	<b>5</b>
3.1	Init command . . . . .	5
3.2	Fetch command . . . . .	6
3.3	Forge command . . . . .	8
3.3.1	The forge selection language . . . . .	11
3.3.2	Treatment of the genotype data while merging . . . . .	12
3.3.3	Treatment of the .janno file while merging . . . . .	13
3.3.4	Treatment of the .ssf file while merging . . . . .	14
3.3.5	Treatment of the .bib file while merging . . . . .	14
3.3.6	Other options . . . . .	14
3.4	Genoconvert command . . . . .	15
3.5	Jannocoalesce command . . . . .	16
3.6	Rectify command . . . . .	17
<b>4</b>	<b>Inspection commands</b>	<b>19</b>
4.1	List command . . . . .	19
4.2	Summarise command . . . . .	20
4.3	Survey command . . . . .	21
4.4	Validate command . . . . .	21

## 1 Installation

See the Poseidon website (<https://www.poseidon-adna.org/#/trident>) or the GitHub repository (<https://github.com/poseidon-framework/poseidon-hs>) for up-to-date installation instructions.

## 2 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>jannocoalesce</code>	Coalesce information from one or multiple janno files to another one
<code>rectify</code>	Adjust POSEIDON.yml files automatically to package changes

### Inspection commands:

<code>list</code>	List packages, groups or individuals from local or remote Poseidon repositories
<code>summarise</code>	Get an overview over the content of one or multiple Poseidon packages
<code>survey</code>	Survey the degree of context information completeness for Poseidon packages
<code>validate</code>	Check Poseidon packages or package components for

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

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

```
85 /path/to/poseidon/packages
86     /modern
87         /2019_poseidon_package1
88         /2019_poseidon_package2
89     /ancient
90         /...
91         /...
92     /Reference_Genomes
93         /...
94         /...
```

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

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

```
101 ~/my_project/my_project.geno
102 ~/my_project/my_project.snp
103 ~/my_project/my_project.ind
```

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

```
106 poseidonVersion: 2.7.1
107 title: My_awesome_project
108 description: Unpublished genetic data from my awesome project
109 contributor:
110     - name: Stephan Schiffels
111       email: schiffels@institute.org
112 packageVersion: 0.1.0
113 lastModified: 2020-10-07
114 genotypeData:
115     format: EIGENSTRAT
116     genoFile: my_project.geno
117     snpFile: my_project.snp
118     indFile: my_project.ind
```

119 jannoFile: my\_project.janno

120 bibFile: sources.bib

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

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

128 trident list -d /path/to/poseidon/packages/modern \

129 -d /path/to/poseidon/packages/ReferenceGenomes

130 -d ~/my\_project --packages

## 131 2.1 General notes

### 132 2.1.1 Logging and command line output

133 For all subcommands the general argument `--logMode` defines how trident reports messages (to stderr) on the  
134 command line:

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

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

### 142 2.1.2 Package duplicates and versions

- 143 • For **trident** multiple packages in a set of base directories can share the same **title**, if they have different  
144 **packageVersion** numbers. If the version numbers are identical or missing, then **trident** stops with an  
145 exception.
- 146 • The **trident** subcommands **genoconvert**, **list**, **rectify**, **survey** and **validate** by default consider all  
147 versions of each Poseidon package in the given base directories. The `--onlyLatest` flag causes them to  
148 instead only consider the latest versions.
- 149 • **fetch** and **forge** generally consider all package versions and their selection language (see below) allows  
150 for detailed version handling.
- 151 • **summarize** and **jannocoalesce** always only consider the latest package versions.

### 152 2.1.3 Individual/sample duplicates

- 153 • Individual/sample names (**Poseidon\_IDs**) within one package have to be unique, or trident will stop.
- 154 • We also discourage sample duplicates across packages in package repositories, but trident will generally  
155 continue with them. **validate** will fail though, if the `--ignoreDuplicates` flag is not set.
- 156 • **forge** offers a special mechanism to resolve sample duplicates within its selection language.

#### 157 2.1.4 Group names in .fam files

158 The .fam file of Plink-formatted genotype data is used inconsistently across different popular aDNA software  
159 tools to store group/population name information. The (global) option `--inPlinkPopName` with the arguments  
160 `asFamily` (default), `asPhenotype` and `asBoth` allows to control the reading of the population name from Plink  
161 .fam files. The subcommands that write genotype data (`forge`, `genoconvert`) have a corresponding option  
162 `--outPlinkPopName` to specify this for the output.

#### 163 2.1.5 Whitespaces in the .janno file

164 While reading the .janno file `trident` trims all leading and trailing whitespaces around individual cells. Also  
165 all instances of the `No-Break Space` unicode character will be removed. This means these whitespaces will not  
166 be preserved when a package is `forged`.

## 167 3 Package creation and manipulation commands

### 168 3.1 Init command

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

171 Command line details

```
172 Usage: trident init ((-p|--genoOne FILE) | --inFormat FORMAT --genoFile FILE
173                  --snpFile FILE --indFile FILE) [--snpSet SET]
174                  (-o|--outPackagePath DIR) [-n|--outPackageName STRING]
175                  [--minimal]
```

177 Create a new Poseidon package from genotype data

179 Available options:

180 <code>-h,--help</code>	Show this help text
181 <code>-p,--genoOne FILE</code>	One of the input genotype data files. Expects .bed, 182 .bim or .fam for PLINK and .geno, .snp or .ind for 183 EIGENSTRAT. The other files must be in the same 184 directory and must have the same base name.
185 <code>--inFormat FORMAT</code>	The format of the input genotype data: EIGENSTRAT or 186 PLINK. Only necessary for data input with <code>--genoFile</code> 187 + <code>--snpFile</code> + <code>--indFile</code> .
188 <code>--genoFile FILE</code>	Path to the input geno file.
189 <code>--snpFile FILE</code>	Path to the input snp file.
190 <code>--indFile FILE</code>	Path to the input ind file.
191 <code>--snpSet SET</code>	The snpSet of the package: 1240K, HumanOrigins or 192 Other. Only relevant for data input with <code>-p --genoOne</code> 193 or <code>--genoFile</code> + <code>--snpFile</code> + <code>--indFile</code> , because the 194 packages in a <code>-d --baseDir</code> already have this 195 information in their respective <code>POSEIDON.yml</code> files. 196 (default: Other)

```

197 -o,--outPackagePath DIR Path to the output package directory.
198 -n,--outPackageName STRING
199                               The output package name. This is optional: If no name
200                               is provided, then the package name defaults to the
201                               basename of the (mandatory) --outPackagePath
202                               argument. (default: Nothing)
203 --minimal                     Should the output data be reduced to a necessary
204                               minimum and omit empty scaffolding?
205
206 The command
207
208 trident init \
209   --inFormat EIGENSTRAT/PLINK \
210   --genoFile path/to/geno_file \
211   --snpFile path/to/snp_file \
212   --indFile path/to/ind_file \
213   --snpSet 1240K|HumanOrigins|Other \
214   -o path/to/new_package_name
215
216 requires the format (--inFormat) of your input data (either EIGENSTRAT or PLINK), the paths to the respective
217 files (--genoFile, --snpFile, --indFile), and optionally the “shape” of these files (--snpSet), so if they cover
218 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

```

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

```

## 219 3.2 Fetch command

```

220 fetch allows to download Poseidon packages from a remote Poseidon server via a Web API. This server provides
221 all packages in the Poseidon public archives.

```

```

222 Command line details

```

```

223 Usage: trident fetch (-d|--baseDir DIR)
224           (--downloadAll |
225           (--fetchFile FILE | (-f|--fetchString DSL)))
226           [--remoteURL URL] [--archive STRING]

```

```

227
228 Download data from a remote Poseidon repository

```

```

229
230 Available options:

```

```

231 -h,--help           Show this help text

```

232 `-d,--baseDir DIR` A base directory to search for Poseidon packages.  
 233 `--downloadAll` Download all packages the server is offering.  
 234 `--fetchFile FILE` A file with a list of packages. Works just as `-f`, but  
 235 multiple values can also be separated by newline, not  
 236 just by comma. `-f` and `--fetchFile` can be combined.  
 237 `-f,--fetchString DSL` List of packages to be downloaded from the remote  
 238 server. Package names should be wrapped in asterisks:  
 239 `*package_title*`. You can combine multiple values with  
 240 comma, so for example: `"*package_1*, *package_2*,`  
 241 `*package_3*"`. `fetchString` uses the same parser as  
 242 `forgeString`, but does not allow excludes. If groups  
 243 or individuals are specified, then packages which  
 244 include these groups or individuals are included in  
 245 the download.  
 246 `--remoteURL URL` URL of the remote Poseidon server.  
 247 (default: `"https://server.poseidon-adna.org"`)  
 248 `--archive STRING` The name of the Poseidon package archive that should  
 249 be queried. If not given, then the query falls back  
 250 to the default archive of the server selected with  
 251 `--remoteURL`. See the archive documentation at  
 252 `https://www.poseidon-adna.org/#/archive_overview` for  
 253 a list of archives currently available from the  
 254 official Poseidon Web API. (default: Nothing)

255 It works with

```
256 trident fetch -d ... -d ... \  
257 -f "*package_title_1*,*package_title_2-1.0.1*,group_name,<individual1>"
```

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

261 Entities are specified using a special syntax (see also the documentation of `forge` below): packages are wrapped  
 262 in asterisks, with or without version appended after a dash (e.g. `*package_title*` or `*package_title-1.2.3*`),  
 263 group names are spelled as is, and individual names are wrapped in angular brackets (e.g. `<individual1>`).  
 264 Fetch will figure out which packages need to be downloaded to include all specified entities. `--downloadAll`,  
 265 which can be given instead of `-f` and `--fetchFile`, causes fetch to download all packages from the server. The  
 266 downloaded packages are added in the first (!) `-d` directory (which gets created if it doesn't exist), but downloads  
 267 are only performed if the respective packages are not already present in the latest version in any of the `-d` dirs.

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

270 `fetch` also has the optional arguments `--remote https://...` to name an alternative Poseidon server and  
 271 `--archive` to select a specific Poseidon public archive on the server.

### 272 3.3 Forge command

273 **forge** creates new Poseidon packages by extracting and merging packages, populations and individuals/samples  
274 from your Poseidon repositories.

275 Command line details

```
276 Usage: trident forge ((-d|--baseDir DIR) |  
277                      ((-p|--genoOne FILE) | --inFormat FORMAT --genoFile FILE  
278                      --snpFile FILE --indFile FILE) [--snpSet SET])  
279                      [--forgeFile FILE | (-f|--forgeString DSL)]  
280                      [--selectSnps FILE] [--intersect] [--outFormat FORMAT]  
281                      [--minimal] [--onlyGeno] (-o|--outPackagePath DIR)  
282                      [-n|--outPackageName STRING] [--packagewise]  
283                      [--outPlinkPopName MODE]  
284
```

285 Select packages, groups or individuals and create a new Poseidon package from  
286 them

287  
288 Available options:

289 -h,--help	Show this help text
290 -d,--baseDir DIR	A base directory to search for Poseidon packages.
291 -p,--genoOne FILE	One of the input genotype data files. Expects .bed, 292 .bim or .fam for PLINK and .geno, .snp or .ind for 293 EIGENSTRAT. The other files must be in the same 294 directory and must have the same base name.
295 --inFormat FORMAT	The format of the input genotype data: EIGENSTRAT or 296 PLINK. Only necessary for data input with --genoFile 297 + --snpFile + --indFile.
298 --genoFile FILE	Path to the input geno file.
299 --snpFile FILE	Path to the input snp file.
300 --indFile FILE	Path to the input ind file.
301 --snpSet SET	The snpSet of the package: 1240K, HumanOrigins or 302 Other. Only relevant for data input with -p --genoOne 303 or --genoFile + --snpFile + --indFile, because the 304 packages in a -d --baseDir already have this 305 information in their respective POSEIDON.yml files. 306 (default: Other)
307 --forgeFile FILE	A file with a list of packages, groups or individual 308 samples. Works just as -f, but multiple values can 309 also be separated by newline, not just by comma. 310 Empty lines are ignored and comments start with "#", 311 so everything after "#" is ignored in one line. 312 Multiple instances of -f and --forgeFile can be 313 given. They will be evaluated according to their 314 input order on the command line.
315 -f,--forgeString DSL	List of packages, groups or individual samples to be



combined in the output package. Packages follow the syntax `*package_title*`, populations/groups are simply `group_id` and individuals `<individual_id>`. You can combine multiple values with comma, so for example: `"*package_1*, <individual_1>, <individual_2>, group_1"`. Duplicates are treated as one entry. Negative selection is possible by prepending "-" to the entity you want to exclude (e.g. `"*package_1*, -<individual_1>, -group_1"`). `forge` will apply excludes and includes in order. If the first entity is negative, then `forge` will assume you want to merge all individuals in the packages found in the `baseDirs` (except the ones explicitly excluded) before the exclude entities are applied. An empty `forgeString` (and no `--forgeFile`) will therefore merge all available individuals. If there are individuals in your input packages with equal individual id, but different main group or source package, they can be specified with the special syntax `"<package:group:individual>"`.

`--selectSnps FILE` To extract specific SNPs during this `forge` operation, provide a Snp file. Can be either Eigenstrat (file ending must be `'.snp'`) or Plink (file ending must be `'.bim'`). 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 option `'--intersect'` is also set, only the SNPs overlapping between the SNP file and the forged packages are output. (default: Nothing)

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

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

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

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

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

`-n,--outPackageName STRING` The output package name. This is optional: If no name

```

361         is provided, then the package name defaults to the
362         basename of the (mandatory) --outPackagePath
363         argument. (default: Nothing)
364     --packagewise      Skip the within-package selection step in forge. This
365                       will result in outputting all individuals in the
366                       relevant packages, and hence a superset of the
367                       requested individuals/groups. It may result in better
368                       performance in cases where one wants to forge entire
369                       packages or almost entire packages. Details: Forge
370                       conceptually performs two types of selection: First,
371                       it identifies which packages in the supplied base
372                       directories are relevant to the requested forge, i.e.
373                       whether they are either explicitly listed using
374                       *PackageName*, or because they contain selected
375                       individuals or groups. Second, within each relevant
376                       package, individuals which are not requested are
377                       removed. This option skips only the second step, but
378                       still performs the first.
379     --outPlinkPopName MODE Where to write the population/group name into the FAM
380                           file in Plink-format. Three options are possible:
381                           asFamily (default) | asPhenotype | asBoth. See also
382                           --inPlinkPopName.

```

383 `forge` can be used with

```

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

```

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

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

```

395 trident forge \
396     -d 2017_GonzalesFortesCurrentBiology \
397     -p 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \
398     --inFormat PLINK \
399     --genoFile 2017_HaberAJHG/2017_HaberAJHG.bed \
400     --snpFile 2017_HaberAJHG/2017_HaberAJHG.bim \
401     --indFile 2017_HaberAJHG/2017_HaberAJHG.fam \
402     -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \
403     -o testpackage \

```

```

404 --outFormat EIGENSTRAT \
405 --onlyGeno

```

### 406 3.3.1 The forge selection language

407 The text in `--forgeString`, `--forgeFile` (and with limited syntax also in `--fetchString` and `--fetchFile`)  
 408 are parsed as a domain specific query language that describes precisely which entities should be compiled in  
 409 the output package of a given **forge** operation. The language has multiple syntactic elements and a specific  
 410 evaluation logic.

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

- 413 • Each package title is surrounded by `*`: `*package*`. That means if you want all individuals of the Poseidon  
 414 package `2019_Jeong_InnerEurasia` in the output package you would add `*2019_Jeong_InnerEurasia*`  
 415 to the query.
- 416 • Groups/populations are not specially marked: `group`. So to get all individuals of the group  
 417 `Swiss_Roman_period`, you would simply add `Swiss_Roman_period`.
- 418 • Individuals/samples are surrounded by `<` and `>`: `<individual>`. `ALA026` therefore becomes `<ALA026>`. A sec-  
 419 ond way to denote individuals is with the more verbose and specific syntax `<package:group:individual>`.  
 420 Such defined individuals take precedence over differently defined ones (so: directly with `<individual>` or  
 421 as a subset of `*package*` or `group`). This allows to resolve duplication issues precisely – at least in cases  
 422 where the duplicated individuals differ in source package or primary group.
- 423 • Package versions can be appended to package names, such as `*package-1.2.3*`.
- 424 • This also works with the verbose individual syntax: `<package-1.2.3:group:individual>`.

425 In the `--forgeFile` each line is treated as a separate `forgeString`, empty lines are ignored and `#` symbols start  
 426 comments. So this is a valid example of a `forgeFile`:

```

427 # Packages
428 *package1*, *package2-1.2.3*
429
430 # Groups and individuals from other packages beyond package1 and package2
431 group1, <individual1>, group2, <individual2>, <pac1:group2:individual3>
432
433 # group2 has two outlier individuals that should be ignored
434 -<individual1> # This one has very low coverage
435 -<pac2:group3:individual4> # This one is from a different time period

```

436 By prepending `-` to entities, we can exclude them from the forged package (this feature is not avail-  
 437 able for `fetch`). **forge** figures out the final list of samples to include by executing all **forge**-entities in  
 438 order. So an entity list `*PackageA*,-<Individual1>,GroupA` may result in a different outcome than  
 439 `*PackageA*,GroupA,-<Individual1>`, depending on whether `<Individual1>` belongs to `GroupA` or not.

440 If the **forge** entity list starts with a negative entity, or if the entity list is empty, **forge** will implicitly assume  
 441 you want to include all individuals in all **latest** versions of packages found in the base directories (except the  
 442 ones explicitly excluded, of course).

443 The specific semantics of the various ways to include or exclude entities are:

### 444 3.3.1.1 Inclusion queries

- 445 • **\*Pac1\***: Select all individuals in the latest version of package “Pac1”
- 446 • **\*Pac1-1.0.1\***: Select all individuals in package “Pac1” with version “1.0.1”
- 447 • **Group1**: Select all individuals associated with “Group1” in all latest versions of all packages
- 448 • **<Ind1>**: Select the individual named “Ind1”, searching in all latest packages.
- 449 • **<Pac1:Group1:Ind1>**: Select the individual named “Ind1” associated with “Group1” in the latest version
- 450 of package “Pac1”
- 451 • **<Pac1-1.0.1:Group1:Ind1>**: Select the individual named “Ind1” associated with “Group1” in the package
- 452 “Pac1” with version “1.0.1”

### 453 3.3.1.2 Exclusion queries

- 454 • **-\*Pac1\***: Remove all individuals in all versions of package “Pac1”
- 455 • **-\*Pac1-1.0.1\***: Remove only individuals in package “Pac1” with version “1.0.1” (but leave other versions
- 456 in)
- 457 • **-Group1**: Remove all individuals associated with “Group1” in all versions of all packages (not just the
- 458 latest)
- 459 • **-<Ind1>**: Remove all individuals named “Ind1” in all versions of all packages (not just the latest).
- 460 • **-<Pac1:Group1:Ind1>**: Remove the individual named “Ind1” associated with “Group1”, searching in all
- 461 versions of package “Pac1”
- 462 • **-<Pac1-1.0.1:Group1:Ind1>**: Remove the individual named “Ind1” associated with “Group1”, but only
- 463 if they are in “Pac1” with version “1.0.1”

464 If a query results in multiple individuals with the same name, forge will throw an error.

### 465 3.3.2 Treatment of the genotype data while merging

466 Forge performs a series of steps to merge the genotype data of multiple source files:

- 467 1. Genotype data from each package is streamed in parallel. Because our packages may have different SNP
- 468 locations (specified by chromosome-position pairs) listed in their `.bim/.snp` file, we first perform a zipping-
- 469 operation, whose behaviour depends on whether `--intersect` is set or not. Without `--intersect`, any
- 470 SNP position listed in any package will be forwarded to the output, with missing values being filled in in
- 471 all packages that do not list that particular SNP. With `--intersect`, only SNP positions that are present
- 472 in all packages are considered. Note that relevant for this step is only whether a given SNP position is part
- 473 of the genotype data, not whether the actual genotypes are missing or not.
- 474 2. At each SNP, the consensus alleles are selected, by collecting all reference and alternative alleles from all
- 475 sources. If more than two non-dummy alleles (alleles different from N) are present in that collection, an
- 476 error is thrown. If exactly two non-dummy alleles are present (which should be the case for binary SNPs),
- 477 the two alleles are declared “reference” and “alternative” alleles for the output. If only one non-dummy
- 478 allele is present, it is set to be the reference allele, and “N” is set to be the alternative.
- 479 3. All source genotype data is then read and recoded in terms of the two chosen consensus alleles. This will
- 480 make sure that source data with flipped reference and alternative allele gets correctly merged in.
- 481 4. SNP IDs, as part of PLINK `.bim` files are checked across the source files. If all SNP IDs for a given SNP
- 482 are missing, then the result will also be missing. If there is only one SNP ID present in some or all source
- 483 packages, that ID gets forwarded to the output. In the (unusual) case that there are multiple different
- 484 non-missing SNP ids (of the form “rs” followed by a number), then a debug warning is output (which gets

printed to the screen when `--logMode DEBUG` is selected), and simply the first value is chosen to be output into the forged `.bim` file. We decided not to throw an error in that case, because we consider the physical position of the SNP (specified by Chromosome and position) to be definitive, and the SNP ID to be of secondary importance.

5. Genetic positions, as part of PLINK `.bim` files are checked in a similar manner, with “0.0” being interpreted as missing.

### 3.3.3 Treatment of the `.janno` file while merging

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

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

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

The following example illustrates the described behaviour:

#### **A.janno**

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn1	AdditionalColumn2
XXX011	POP1	M	A	D
XXX012	POP2	F	B	E
XXX013	POP1	M	C	F

#### **B.janno**

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn3	AdditionalColumn2
YYY022	POP5	F	G	J
YYY023	POP5	F	H	K
YYY024	POP5	M	I	L

#### **A.janno + B.janno**

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn1	AdditionalColumn2	AdditionalColumn3
XXX011	POP1	M	A	D	n/a
XXX012	POP2	F	B	E	n/a
XXX013	POP1	M	C	F	n/a
YYY022	POP5	F	n/a	J	G
YYY023	POP5	F	n/a	K	H

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn1	AdditionalColumn2	AdditionalColumn3
YYY024	POP5	M	n/a	L	I

### 3.3.4 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.

### 3.3.5 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).

### 3.3.6 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.

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

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

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

### 541 3.4 Genoconvert command

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

544 Command line details

```
545 Usage: trident genoconvert ((-d|--baseDir DIR) |  
546                             ((-p|--genoOne FILE) | --inFormat FORMAT  
547                             --genoFile FILE --snpFile FILE --indFile FILE)  
548                             [--snpSet SET]) --outFormat FORMAT [--onlyGeno]  
549                             [-o|--outPackagePath DIR] [--removeOld]  
550                             [--outPlinkPopName MODE] [--onlyLatest]
```

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

554 Available options:

555 -h,--help	Show this help text
556 -d,--baseDir DIR	A base directory to search for Poseidon packages.
557 -p,--genoOne FILE	One of the input genotype data files. Expects .bed, 558 .bim or .fam for PLINK and .geno, .snp or .ind for 559 EIGENSTRAT. The other files must be in the same 560 directory and must have the same base name.
561 --inFormat FORMAT	The format of the input genotype data: EIGENSTRAT or 562 PLINK. Only necessary for data input with --genoFile 563 + --snpFile + --indFile.
564 --genoFile FILE	Path to the input geno file.
565 --snpFile FILE	Path to the input snp file.
566 --indFile FILE	Path to the input ind file.
567 --snpSet SET	The snpSet of the package: 1240K, HumanOrigins or 568 Other. Only relevant for data input with -p --genoOne 569 or --genoFile + --snpFile + --indFile, because the 570 packages in a -d --baseDir already have this 571 information in their respective POSEIDON.yml files. 572 (default: Other)
573 --outFormat FORMAT	the format of the output genotype data: EIGENSTRAT or 574 PLINK.
575 --onlyGeno	Should only the resulting genotype data be returned? 576 This means the output will not be a Poseidon package.

```

577 -o,--outPackagePath DIR Path to the output package directory. This is
578 optional: If no path is provided, then the output is
579 written to the directories where the input genotype
580 data file (.bed/.geno) is stored. (default: Nothing)
581 --removeOld Remove the old genotype files when creating the new
582 ones.
583 --outPlinkPopName MODE Where to write the population/group name into the FAM
584 file in Plink-format. Three options are possible:
585 asFamily (default) | asPhenotype | asBoth. See also
586 --inPlinkPopName.
587 --onlyLatest Consider only the latest versions of packages, or the
588 groups and individuals within the latest versions of
589 packages, respectively.

590 With the default setting

591 trident genoconvert -d ... -d ... --outFormat EIGENSTRAT|PLINK

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

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

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

600 trident genoconvert \
601 -p 2018_Mittnik_Baltic/Mittnik_Baltic.bed \
602 --outFormat EIGENSTRAT
603 -o my_directory

```

### 604 3.5 Jannocoalesce command

605 jannocoalesce merges information from one or multiple source .janno files into a target .janno file.

606 Command line details

```

607 Usage: trident jannocoalesce ((-s|--sourceFile FILE) | (-d|--baseDir DIR))
608 (-t|--targetFile FILE) [-o|--outFile FILE]
609 [--includeColumns ARG | --excludeColumns ARG]
610 [-f|--force] [--sourceKey ARG] [--targetKey ARG]
611 [--stripIdRegex ARG]
612

```

613 Coalesce information from one or multiple janno files to another one

614 Available options:

```

615 -h,--help Show this help text
616 -s,--sourceFile FILE The source .janno file.

```



```

618 -d,--baseDir DIR      A base directory to search for Poseidon packages.
619 -t,--targetFile FILE  The target .janno file to fill.
620 -o,--outFile FILE     An optional file to write the results to. If not
621                       specified, change the target file in place.
622                       (default: Nothing)
623 --includeColumns ARG  A comma-separated list of .janno column names to
624                       coalesce. If not specified, all columns that can be
625                       found in the source and target will get filled.
626 --excludeColumns ARG  A comma-separated list of .janno column names NOT to
627                       coalesce. All columns that can be found in the source
628                       and target will get filled, except the ones listed
629                       here.
630 -f,--force            With this option, potential non-missing content in
631                       target columns gets overridden with non-missing
632                       content in source columns. By default, only missing
633                       data gets filled-in.
634 --sourceKey ARG       The .janno column to use as the source key.
635                       (default: "Poseidon_ID")
636 --targetKey ARG       The .janno column to use as the target key.
637                       (default: "Poseidon_ID")
638 --stripIdRegex ARG    An optional regular expression to identify parts of
639                       the IDs to strip before matching between source and
640                       target. Uses POSIX Extended regular expressions.

```

641 A most basic run may just include two arguments:

```

642 trident jannocoalesce \
643   --sourceFile path/to/source.janno \
644   --targetFile path/to/target.janno

```

645 `jannocoalesce` generally works by reading a source .janno file with `-s|--sourceFile` (or all .janno files in a  
646 `-d|--baseDir`) and a target .janno file with `-t|--targetFile`.

647 It then merges these files by a key column, which can be selected with `--sourceKey` and `--targetKey`. The  
648 default for both of these key columns is the `Poseidon_ID`. In case the entries in the key columns slightly and  
649 systematically differ, e.g. because the `Poseidon_ID`s in either have a special suffix (for example `_SG`), then the  
650 `--stripIdRegex` option allows to strip these with a regular expression to thus match the keys.

651 `jannocoalesce` generally attempts to fill **all** empty cells in the target .janno file with information from the  
652 source. `--includeColumns` and `--excludeColumns` allow to select specific columns for which this should be  
653 done. In some cases it may be desirable to not just fill empty fields in the target, but overwrite the information  
654 already there with the `-f|--force` option. If the target file should be preserved, then the output can be directed  
655 to a new output .janno file with `-o|--outFile`.

## 656 3.6 Rectify command

```

657 rectify automatically harmonizes POSEIDON.yml files of one or multiple packages. This is not an automatic
658 update from one Poseidon version to the next, but rather a clean-up wizard after manual modifications.

```

## 659 Command line details

```
660 Usage: trident rectify (-d|--baseDir DIR) [--ignorePoseidonVersion]
661                [--poseidonVersion ?.??.?]
662                [--packageVersion VPART [--logText STRING]]
663                [--checksumAll | [--checksumGeno] [--checksumJanno]
664                [--checksumSSF] [--checksumBib]]
665                [--newContributors DSL] [--onlyLatest]
```

667 Adjust POSEIDON.yml files automatically to package changes

### 669 Available options:

```
670 -h,--help                Show this help text
671 -d,--baseDir DIR         A base directory to search for Poseidon packages.
672 --ignorePoseidonVersion  Read packages even if their poseidonVersion is not
673                          compatible with trident.
674 --poseidonVersion ?.??.? Poseidon version the packages should be updated to:
675                          e.g. "2.5.3".
676 --packageVersion VPART   Part of the package version number in the
677                          POSEIDON.yml file that should be updated: Major,
678                          Minor or Patch (see https://semver.org).
679 --logText STRING         Log text for this version in the CHANGELOG file.
680 --checksumAll            Update all checksums.
681 --checksumGeno           Update genotype data checksums.
682 --checksumJanno          Update .janno file checksum.
683 --checksumSSF            Update .ssf file checksum
684 --checksumBib            Update .bib file checksum.
685 --newContributors DSL    Contributors to add to the POSEIDON.yml file in the
686                          form "[Firstname Lastname](Email address);...".
687 --onlyLatest             Consider only the latest versions of packages, or the
688                          groups and individuals within the latest versions of
689                          packages, respectively.
```

690 It can be called with a lot of optional arguments. Note that `rectify` by default does **not** apply any changes if  
691 none of these arguments are set.

```
692 trident rectify -d ... -d ... \
693   --poseidonVersion "X.X.X" \
694   --packageVersion Major|Minor|Patch \
695   --logText "short description of the update" \
696   --checksumAll \
697   --newContributors "[Firstname Lastname](Email address);..."
```

698 The following arguments determine which fields of the POSEIDON.yml file should be modified:

- 699 • `--poseidonVersion` allows a simple change of the `poseidonVersion` field in the POSEIDON.yml file.
- 700 • `--packageVersion` increments the package version number in the first, the second or the third position.
- 701 It can optionally be called with `--logText`, which appends an entry to the CHANGELOG file for the

702       respective package version update. `--logText` also creates a new CHANGELOG file if it does not exist  
703       yet.

- 704       • `--checksumGeno`, `--checksumJanno`, `--checksumSSF` and `--checksumBib` add or modify the respective  
705       checksum fields in the POSEIDON.yml file. `--checksumAll` is a wrapper to call all of them at once.
- 706       • `--newContributors` adds new contributors.

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

## 710 4 Inspection commands

### 711 4.1 List command

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

713       Command line details

```
714 Usage: trident list ((-d|--baseDir DIR) | --remote [--remoteURL URL]
715                   [--archive STRING])
716                   (--packages | --groups | --individuals
717                   [-j|--jannoColumn COLNAME]) [--raw] [--onlyLatest]
```

718  
719       List packages, groups or individuals from local or remote Poseidon  
720       repositories

721  
722       Available options:

723       -h,--help	Show this help text
724       -d,--baseDir DIR	A base directory to search for Poseidon packages.
725       --remote	List packages from a remote server instead the local 726       file system.
727       --remoteURL URL	URL of the remote Poseidon server. 728       (default: "https://server.poseidon-adna.org")
729       --archive STRING	The name of the Poseidon package archive that should 730       be queried. If not given, then the query falls back 731       to the default archive of the server selected with 732       --remoteURL. See the archive documentation at 733 <a href="https://www.poseidon-adna.org/#/archive_overview">https://www.poseidon-adna.org/#/archive_overview</a> for 734       a list of archives currently available from the 735       official Poseidon Web API. (default: Nothing)
736       --packages	List all packages.
737       --groups	List all groups, ignoring any group names after the 738       first as specified in the .janno-file.
739       --individuals	List all individuals/samples.
740       -j,--jannoColumn COLNAME	List additional fields from the janno files, using 741       the .janno column heading name, such as "Country", 742       "Site", "Date_C14_Uncal_BP", etc..

743     **--raw**                     Return the output table as tab-separated values  
744                                 without header. This is useful for piping into grep  
745                                 or awk.

746     **--onlyLatest**             Consider only the latest versions of packages, or the  
747                                 groups and individuals within the latest versions of  
748                                 packages, respectively.

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

750   **trident list -d ... -d ... --packages**

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

752   You can use **--remote** to show packages on the remote server. For example

753   **trident list --packages --remote --archive "community-archive"**

754   will result in a view of all packages available in one of the Poseidon public archives. Just as for **fetch**, the  
755   **--archive** flag allows to choose which public archive to query.

756   Independent of whether you query a local or an online archive, you can not just list packages, but also groups,  
757   as defined in the third column of EIGENSTRAT **.ind** files (or the first/last column of a PLINK **.fam** file), and  
758   individuals with the flags **--groups** and **--individuals** (instead of **--packages**).

759   The **--individuals** flag additionally provides a way to immediately access information from **.janno** files  
760   on the command line. This works with the **-j/--jannoColumn** option. For example adding **-j Country -j**  
761   **Date\_C14\_Uncal\_BP** to the commands above will add the **Country** and the **Date\_C14\_Uncal\_BP** columns to the  
762   respective output tables.

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

## 766   4.2   Summarise command

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

768   Command line details

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

770

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

772

773   Available options:

774     **-h,--help**                 Show this help text

775     **-d,--baseDir DIR**         A base directory to search for Poseidon packages.

776     **--raw**                     Return the output table as tab-separated values  
777                                 without header. This is useful for piping into grep  
778                                 or awk.

779   You can run it with

780   **trident summarise -d ... -d ...**

781 which will show you context information like – among others – the number of individuals in the dataset, their  
782 sex distribution, the mean age of the samples (for ancient data) or the mean coverage on the 1240K SNP array  
783 in a table. `summarise` depends on complete `.janno` files and will silently ignore missing information.  
784 You can use the `--raw` option to output the summary table in a simple, tab-delimited layout.

## 785 4.3 Survey command

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

787 Command line details

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

789 Survey the degree of context information completeness for Poseidon packages

792 Available options:

793 <code>-h,--help</code>	Show this help text
794 <code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
795 <code>--raw</code>	Return the output table as tab-separated values 796 without header. This is useful for piping into <code>grep</code> 797 or <code>awk</code> .
798 <code>--onlyLatest</code>	Consider only the latest versions of packages, or the 799 groups and individuals within the latest versions of 800 packages, respectively.

801 Running

802 `trident survey -d ... -d ...`

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

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

## 806 4.4 Validate command

807 `validate` checks Poseidon packages and individual package components for structural correctness.

808 Command line details

809 Usage: `trident validate ((-d|--baseDir DIR) [--ignoreGeno] [--fullGeno]  
810 [--ignoreDuplicates] [-c|--ignoreChecksums]  
811 [--ignorePoseidonVersion] |  
812 --pyml FILE | (-p|--genoOne FILE) | --inFormat FORMAT  
813 --genoFile FILE --snpFile FILE --indFile FILE |  
814 --janno FILE | --ssf FILE | --bib FILE) [--noExitCode]  
815 [--onlyLatest]`

816 Check Poseidon packages or package components for structural correctness

817 Available options:

```

820 -h,--help                Show this help text
821 -d,--baseDir DIR        A base directory to search for Poseidon packages.
822 --ignoreGeno            Ignore snp and geno file.
823 --fullGeno              Test parsing of all SNPs (by default only the first
824                        100 SNPs are probed).
825 --ignoreDuplicates      Do not stop on duplicated individual names in the
826                        package collection.
827 -c,--ignoreChecksums    Whether to ignore checksums. Useful for speedup in
828                        debugging.
829 --ignorePoseidonVersion Read packages even if their poseidonVersion is not
830                        compatible with trident.
831 --pym1 FILE             Path to a POSEIDON.yml file.
832 -p,--genoOne FILE       One of the input genotype data files. Expects .bed,
833                        .bim or .fam for PLINK and .geno, .snp or .ind for
834                        EIGENSTRAT. The other files must be in the same
835                        directory and must have the same base name.
836 --inFormat FORMAT       The format of the input genotype data: EIGENSTRAT or
837                        PLINK. Only necessary for data input with --genoFile
838                        + --snpFile + --indFile.
839 --genoFile FILE         Path to the input geno file.
840 --snpFile FILE          Path to the input snp file.
841 --indFile FILE          Path to the input ind file.
842 --janno FILE            Path to a .janno file.
843 --ssf FILE              Path to a .ssf file.
844 --bib FILE              Path to a .bib file.
845 --noExitCode            Do not produce an explicit exit code.
846 --onlyLatest            Consider only the latest versions of packages, or the
847                        groups and individuals within the latest versions of
848                        packages, respectively.

```

849 You can run it with

```

850 trident validate -d ... -d ...

```

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

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

857 When applied to packages, `validate` tries to ensure that each package adheres to the Poseidon package specification. Here is a list of what is checked:

- 859 • Structural correctness of the POSEIDON.yml file.
- 860 • Presence of all files references in the POSEIDON.yml file.
- 861 • Full structural correctness of .janno, .ssf and .bib file.
- 862 • Superficial correctness of genotype data files by parsing the first 100 SNPs. A full check that parses all

863       SNPs can be triggered with the `--fullGeno` option. `--ignoreGeno`, on the other hand, causes `validate`  
864       to ignore the genotype data entirely, which speeds up the validation significantly.

- 865       • Correspondence of BibTeX keys in `.bib` and `.janno`
- 866       • Correspondence of sample IDs in `.janno` and `.ssf`.
- 867       • Correspondence of sample and group IDs in `.janno` and genotype data files.

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

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