

Guide for trident v1.4.1.0

Contents

1	Installation	1
2	The trident CLI	2
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
3	Package creation and manipulation commands	5
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
4	Inspection commands	19
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

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

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

```
/path/to/poseidon/packages
  /modern
    /2019_poseidon_package1
    /2019_poseidon_package2
  /ancient
    /...
    /...
  /Reference_Genomes
    /...
    /...
```

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

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

```
~/my_project/my_project.geno
~/my_project/my_project.snp
~/my_project/my_project.ind
```

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

```
poseidonVersion: 2.7.1
title: My_awesome_project
description: Unpublished genetic data from my awesome project
contributor:
  - name: Stephan Schiffels
    email: schiffels@institute.org
packageVersion: 0.1.0
lastModified: 2020-10-07
genotypeData:
  format: EIGENSTRAT
  genoFile: my_project.geno
  snpFile: my_project.snp
```

```

119 indFile: my_project.ind
120 jannoFile: my_project.janno
121 bibFile: sources.bib

```

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

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

```

129 trident list -d /path/to/poseidon/packages/modern \
130     -d /path/to/poseidon/packages/ReferenceGenomes
131     -d ~/my_project --packages

```

132 2.1 General notes

133 2.1.1 Logging and command line output

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

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

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

143 2.1.2 Package duplicates and versions

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

153 2.1.3 Individual/sample duplicates

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

158 2.1.4 Group names in .fam files

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

164 2.1.5 Whitespaces in the .janno file

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

168 3 Package creation and manipulation commands

169 3.1 Init command

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

172 Command line details

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

177
178 Create a new Poseidon package from genotype data

179
180 Available options:

181 <code>-h,--help</code>	Show this help text
182 <code>-p,--genoOne FILE</code>	One of the input genotype data files. Expects <code>.bed</code> , 183 <code>.bim</code> or <code>.fam</code> for PLINK and <code>.geno</code> , <code>.snp</code> or <code>.ind</code> for 184 EIGENSTRAT. The other files must be in the same 185 directory and must have the same base name.
186 <code>--inFormat FORMAT</code>	The format of the input genotype data: EIGENSTRAT or 187 PLINK. Only necessary for data input with <code>--genoFile</code> 188 + <code>--snpFile</code> + <code>--indFile</code> .
189 <code>--genoFile FILE</code>	Path to the input geno file.
190 <code>--snpFile FILE</code>	Path to the input snp file.
191 <code>--indFile FILE</code>	Path to the input ind file.
192 <code>--snpSet SET</code>	The snpSet of the package: 1240K, HumanOrigins or 193 Other. Only relevant for data input with <code>-p --genoOne</code> 194 or <code>--genoFile</code> + <code>--snpFile</code> + <code>--indFile</code> , because the 195 packages in a <code>-d --baseDir</code> already have this 196 information in their respective <code>POSEIDON.yml</code> files. 197 (default: Other)

```

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

```

requires the format (`--inFormat`) of your input data (either `EIGENSTRAT` or `PLINK`), the paths to the respective files (`--genoFile` , `--snpFile` , `--indFile`), and optionally the “shape” of these files (`--snpSet`), so if they cover 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

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

222 3.2 Fetch command

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

225 Command line details

```

226 Usage: trident fetch (-d|--baseDir DIR)
227
228 (--downloadAll |
229
230 (--fetchFile FILE | (-f|--fetchString DSL)))
231
232 [--remoteURL URL] [--archive STRING]

```

231 Download data from a remote Poseidon repository

232

233 Available options:

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

258 It works with

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

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

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

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

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

276 3.3 Forge command

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

279 Command line details

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

289 Select packages, groups or individuals and create a new Poseidon package from
 290 them

292 Available options:

293 -h,--help	Show this help text
294 -d,--baseDir DIR	A base directory to search for Poseidon packages.
295 -p,--genoOne FILE	One of the input genotype data files. Expects .bed, 296 .bim or .fam for PLINK and .geno, .snp or .ind for 297 EIGENSTRAT. The other files must be in the same 298 directory and must have the same base name.
299 --inFormat FORMAT	The format of the input genotype data: EIGENSTRAT or 300 PLINK. Only necessary for data input with --genoFile 301 + --snpFile + --indFile.
302 --genoFile FILE	Path to the input geno file.
303 --snpFile FILE	Path to the input snp file.
304 --indFile FILE	Path to the input ind file.
305 --snpSet SET	The snpSet of the package: 1240K, HumanOrigins or 306 Other. Only relevant for data input with -p --genoOne 307 or --genoFile + --snpFile + --indFile, because the 308 packages in a -d --baseDir already have this 309 information in their respective POSEIDON.yml files. 310 (default: Other)
311 --forgeFile FILE	A file with a list of packages, groups or individual 312 samples. Works just as -f, but multiple values can 313 also be separated by newline, not just by comma. 314 Empty lines are ignored and comments start with "#", 315 so everything after "#" is ignored in one line. 316 Multiple instances of -f and --forgeFile can be 317 given. They will be evaluated according to their 318 input order on the command line.
319 -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

```

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

```

387 `forge` can be used with

```

388 trident forge -d ... -d ... \
389     -f "*package_name*, group_id, <individual_id>" \
390     -o path/to/new_package_name

```

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

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

```

400 trident forge \
401     -d 2017_GonzalesFortesCurrentBiology \
402     -p 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \
403     --inFormat PLINK \
404     --genoFile 2017_HaberAJHG/2017_HaberAJHG.bed \
405     --snpFile 2017_HaberAJHG/2017_HaberAJHG.bim \
406     --indFile 2017_HaberAJHG/2017_HaberAJHG.fam \
407     -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \

```

```

408 -o testpackage \
409 --outFormat EIGENSTRAT \
410 --onlyGeno

```

411 3.3.1 The forge selection language

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

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

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

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

```

433 # Packages
434 *package1*, *package2-1.2.3*
435
436 # Groups and individuals from other packages beyond package1 and package2
437 group1, <individual1>, group2, <individual2>, <pac1:group2:individual3>
438
439 # group2 has two outlier individuals that should be ignored
440 -<individual1> # This one has very low coverage
441 -<pac2:group3:individual4> # This one is from a different time period

```

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

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

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

450 3.3.1.1 Inclusion queries

- 451 • ***Pac1***: Select all individuals in the latest version of package “Pac1”
- 452 • ***Pac1-1.0.1***: Select all individuals in package “Pac1” with version “1.0.1”
- 453 • **Group1**: Select all individuals associated with “Group1” in all latest versions of all packages
- 454 • **<Ind1>**: Select the individual named “Ind1”, searching in all latest packages.
- 455 • **<Pac1:Group1:Ind1>**: Select the individual named “Ind1” associated with “Group1” in the latest version
- 456 of package “Pac1”
- 457 • **<Pac1-1.0.1:Group1:Ind1>**: Select the individual named “Ind1” associated with “Group1” in the package
- 458 “Pac1” with version “1.0.1”

459 3.3.1.2 Exclusion queries

- 460 • **-*Pac1***: Remove all individuals in all versions of package “Pac1”
- 461 • **-*Pac1-1.0.1***: Remove only individuals in package “Pac1” with version “1.0.1” (but leave other versions
- 462 in)
- 463 • **-Group1**: Remove all individuals associated with “Group1” in all versions of all packages (not just the
- 464 latest)
- 465 • **-<Ind1>**: Remove all individuals named “Ind1” in all versions of all packages (not just the latest).
- 466 • **-<Pac1:Group1:Ind1>**: Remove the individual named “Ind1” associated with “Group1”, searching in all
- 467 versions of package “Pac1”
- 468 • **-<Pac1-1.0.1:Group1:Ind1>**: Remove the individual named “Ind1” associated with “Group1”, but only
- 469 if they are in “Pac1” with version “1.0.1”

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

471 3.3.2 Treatment of the genotype data while merging

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

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

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

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn1	AdditionalColumn2	AdditionalColumn3
YYY023	POP5	F	n/a	K	H
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 `snpSet` s 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.

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

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

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

547 3.4 Genoconvert command

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

550 Command line details

```
551 Usage: trident genoconvert ((-d|--baseDir DIR) |
552                             ((-p|--genoOne FILE) | --inFormat FORMAT
553                             --genoFile FILE --snpFile FILE --indFile FILE)
554                             [--snpSet SET]) --outFormat FORMAT [--onlyGeno]
555                             [-o|--outPackagePath DIR] [--removeOld]
556                             [--outPlinkPopName MODE] [--onlyLatest]
```

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

560 Available options:

561 <code>-h,--help</code>	Show this help text
562 <code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
563 <code>-p,--genoOne FILE</code>	One of the input genotype data files. Expects .bed, 564 .bim or .fam for PLINK and .geno, .snp or .ind for 565 EIGENSTRAT. The other files must be in the same 566 directory and must have the same base name.
567 <code>--inFormat FORMAT</code>	The format of the input genotype data: EIGENSTRAT or 568 PLINK. Only necessary for data input with <code>--genoFile</code> 569 <code>+ --snpFile + --indFile</code> .
570 <code>--genoFile FILE</code>	Path to the input geno file.
571 <code>--snpFile FILE</code>	Path to the input snp file.
572 <code>--indFile FILE</code>	Path to the input ind file.
573 <code>--snpSet SET</code>	The snpSet of the package: 1240K, HumanOrigins or 574 Other. Only relevant for data input with <code>-p --genoOne</code> 575 or <code>--genoFile + --snpFile + --indFile</code> , because the 576 packages in a <code>-d --baseDir</code> already have this 577 information in their respective POSEIDON.yml files. 578 (default: Other)
579 <code>--outFormat FORMAT</code>	the format of the output genotype data: EIGENSTRAT or 580 PLINK.
581 <code>--onlyGeno</code>	Should only the resulting genotype data be returned? 582 This means the output will not be a Poseidon package.

```

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

596  With the default setting

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

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

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

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

606  trident genoconvert \
607    -p 2018_Mittnik_Baltic/Mittnik_Baltic.bed \
608    --outFormat EIGENSTRAT
609    -o my_directory

```

610 3.5 Jannocoalesce command

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

612 Command line details

```

613 Usage: trident jannocoalesce ((-s|--sourceFile FILE) | (-d|--baseDir DIR))
614                               (-t|--targetFile FILE) [-o|--outFile FILE]
615                               [--includeColumns ARG | --excludeColumns ARG]
616                               [-f|--force] [--sourceKey ARG] [--targetKey ARG]
617                               [--stripIdRegex ARG]
618

```

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

621 Available options:

```

622  -h,--help             Show this help text
623  -s,--sourceFile FILE  The source .janno file.

```



```

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

```

647 A most basic run may just include two arguments:

```

648 trident jannocoalesce \
649   --sourceFile path/to/source.janno \
650   --targetFile path/to/target.janno

```

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

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

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

662 3.6 Rectify command

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

665 Command line details

```
666 Usage: trident rectify (-d|--baseDir DIR) [--ignorePoseidonVersion]
667         [--poseidonVersion ?.??.?]
668         [--packageVersion VPART [--logText STRING]]
669         [--checksumAll | [--checksumGeno] [--checksumJanno]
670         [--checksumSSF] [--checksumBib]]
671         [--newContributors DSL] [--onlyLatest]
```

672
673 Adjust POSEIDON.yml files automatically to package changes

675 Available options:

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

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

```
698 trident rectify -d ... -d ... \
699     --poseidonVersion "X.X.X" \
700     --packageVersion Major|Minor|Patch \
701     --logText "short description of the update" \
702     --checksumAll \
703     --newContributors "[Firstname Lastname](Email address);..."
```

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

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

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

- 710 • `--checksumGeno`, `--checksumJanno`, `--checksumSSF` and `--checksumBib` add or modify the respec-
 711 tive checksum fields in the POSEIDON.yml file. `--checksumAll` is a wrapper to call all of them at
 712 once.
- 713 • `--newContributors` adds new contributors.

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

717 4 Inspection commands

718 4.1 List command

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

720 Command line details

```
721 Usage: trident list ((-d|--baseDir DIR) | --remote [--remoteURL URL]
722                   [--archive STRING])
723                   (--packages | --groups | --individuals
724                   [-j|--jannoColumn COLNAME]) [--raw] [--onlyLatest]
```

725
 726 List packages, groups or individuals from local or remote Poseidon
 727 repositories

728
 729 Available options:

730 -h,--help	Show this help text
731 -d,--baseDir DIR	A base directory to search for Poseidon packages.
732 --remote	List packages from a remote server instead the local 733 file system.
734 --remoteURL URL	URL of the remote Poseidon server. 735 (default: "https://server.poseidon-adna.org")
736 --archive STRING	The name of the Poseidon package archive that should 737 be queried. If not given, then the query falls back 738 to the default archive of the server selected with 739 --remoteURL. See the archive documentation at 740 https://www.poseidon-adna.org/#/archive_overview for 741 a list of archives currently available from the 742 official Poseidon Web API. (default: Nothing)
743 --packages	List all packages.
744 --groups	List all groups, ignoring any group names after the 745 first as specified in the .janno-file.
746 --individuals	List all individuals/samples.
747 -j,--jannoColumn COLNAME	List additional fields from the janno files, using 748 the .janno column heading name, such as "Country",

749 "Site", "Date_C14_Uncal_BP", etc..
750 `--raw` Return the output table as tab-separated values
751 without header. This is useful for piping into grep
752 or awk.
753 `--onlyLatest` Consider only the latest versions of packages, or the
754 groups and individuals within the latest versions of
755 packages, respectively.

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

757 `trident list -d ... -d ... --packages`

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

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

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

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

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

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

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

773 4.2 Summarise command

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

775 Command line details

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

777

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

779

780 Available options:

781 <code>-h,--help</code>	Show this help text
782 <code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
783 <code>--raw</code>	Return the output table as tab-separated values 784 without header. This is useful for piping into grep 785 or awk.

786 You can run it with

787 `trident summarise -d ... -d ...`

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

792 4.3 Survey command

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

794 Command line details

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

796
797 Survey the degree of context information completeness for Poseidon packages

798
799 Available options:

800 <code>-h,--help</code>	Show this help text
801 <code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
802 <code>--raw</code>	Return the output table as tab-separated values 803 without header. This is useful for piping into grep 804 or awk.
805 <code>--onlyLatest</code>	Consider only the latest versions of packages, or the 806 groups and individuals within the latest versions of 807 packages, respectively.

808 Running

809 `trident survey -d ... -d ...`

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

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

813 4.4 Validate command

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

815 Command line details

816 Usage: trident validate ((-d|--baseDir DIR) [--ignoreGeno] [--fullGeno]
817 [--ignoreDuplicates] [-c|--ignoreChecksums]
818 [--ignorePoseidonVersion] |
819 --pym1 FILE | (-p|--genoOne FILE) | --inFormat FORMAT
820 --genoFile FILE --sn1File FILE --indFile FILE |
821 --janno FILE | --ssf FILE | --bib FILE) [--noExitCode]
822 [--onlyLatest]

823
824 Check Poseidon packages or package components for structural correctness

825

826 Available options:

```

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

```

856 You can run it with

```
857 trident validate -d ... -d ...
```

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

859 Instead of validating entire packages with `-d` you can also apply it to individual files and package components:

```

860 --pym1 (POSEIDON.yml), -p | --inFormat + --genoFile + --snpFile + --indFile (genotype data),
861 --janno (.janno file), --ssf (.ssf file) or --bib (.bib file). In this case validate attempts to read and
862 parse the respective files individually and reports any issues it encounters. Note that this considers the files in
863 isolation and does not include any cross-file consistency checks.

```

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

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

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

- 872 • Correspondence of BibTeX keys in .bib and .janno
- 873 • Correspondence of sample IDs in .janno and .ssf.
- 874 • Correspondence of sample and group IDs in .janno and genotype data files.

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

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