

# Guide for trident v1.3.0.4

## Contents

<b>1</b>	<b>The trident CLI</b>	<b>1</b>
1.1	General notes	4
1.1.1	Logging and command line output	4
1.1.2	Duplicates	4
1.1.3	Group names in .fam files	4
1.1.4	Whitespaces in the .janno file	4
<b>2</b>	<b>Package creation and manipulation commands</b>	<b>5</b>
2.1	Init command	5
2.2	Fetch command	6
2.3	Forge command	7
2.3.1	The forge selection language	10
2.3.2	Treatment of the .janno file while merging	11
2.3.3	Treatment of the .ssf file while merging	12
2.3.4	Treatment of the .bib file while merging	12
2.3.5	Other options	12
2.4	Genoconvert command	13
2.5	Rectify command	15
<b>3</b>	<b>Inspection commands</b>	<b>16</b>
3.1	List command	16
3.2	Summarise command	17
3.3	Survey command	18
3.4	Validate command	18

## 1 The trident CLI

Trident is a command line software tool structured in multiple subcommands. If you installed it properly you can call it on the command line by typing `trident`. This will show an overview of the general options and all subcommands, which are explained in detail below.

```
Usage: trident [--version] [--logMode MODE | --debug] [--errLength INT]
           [--inPlinkPopName MODE] (COMMAND | COMMAND)
```

`trident` is a management and analysis tool for Poseidon packages. Report issues

34 here: <https://github.com/poseidon-framework/poseidon-hs/issues>

35  
36 Available options:

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

49  
50 Package creation and manipulation commands:

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

59  
60 Inspection commands:

61	<code>list</code>	List packages, groups or individuals from local or
62		remote Poseidon repositories
63	<code>summarise</code>	Get an overview over the content of one or multiple
64		Poseidon packages
65	<code>survey</code>	Survey the degree of context information completeness
66		for Poseidon packages
67	<code>validate</code>	Check Poseidon packages or package components for
68		structural correctness

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

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

76 `/path/to/poseidon/packages`  
77 `/modern`

```

78         /2019_poseidon_package1
79         /2019_poseidon_package2
80     /ancient
81         /...
82         /...
83     /Reference_Genomes
84         /...
85         /...

```

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

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

```

92 ~/my_project/my_project.geno
93 ~/my_project/my_project.snp
94 ~/my_project/my_project.ind

```

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

```

97 poseidonVersion: 2.7.1
98 title: My_awesome_project
99 description: Unpublished genetic data from my awesome project
100 contributor:
101   - name: Stephan Schiffels
102     email: schiffels@institute.org
103 packageVersion: 0.1.0
104 lastModified: 2020-10-07
105 genotypeData:
106   format: EIGENSTRAT
107   genoFile: my_project.geno
108   snpFile: my_project.snp
109   indFile: my_project.ind
110 jannoFile: my_project.janno
111 bibFile: sources.bib

```

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

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

```

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

```

120 -d /path/to/poseidon/packages/ReferenceGenomes  
121 -d ~/my\_project --packages

## 122 1.1 General notes

### 123 1.1.1 Logging and command line output

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

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

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

### 133 1.1.2 Duplicates

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

### 143 1.1.3 Group names in .fam files

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

### 149 1.1.4 Whitespaces in the .janno file

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

## 153 2 Package creation and manipulation commands

### 154 2.1 Init command

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

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

```
158 Usage: trident init ((-p|--genoOne FILE) | --inFormat FORMAT --genoFile FILE
159                  --snpFile FILE --indFile FILE) [--snpSet SET]
160                  (-o|--outPackagePath DIR) [-n|--outPackageName STRING]
161                  [--minimal]
```

163 Create a new Poseidon package from genotype data

165 Available options:

166	<code>-h,--help</code>	Show this help text
167	<code>-p,--genoOne FILE</code>	One of the input genotype data files. Expects <code>.bed</code> , 168 <code>.bim</code> or <code>.fam</code> for PLINK and <code>.geno</code> , <code>.snp</code> or <code>.ind</code> for 169 EIGENSTRAT. The other files must be in the same 170 directory and must have the same base name.
171	<code>--inFormat FORMAT</code>	The format of the input genotype data: EIGENSTRAT or 172 PLINK. Only necessary for data input with <code>--genoFile</code> 173 + <code>--snpFile</code> + <code>--indFile</code> .
174	<code>--genoFile FILE</code>	Path to the input geno file.
175	<code>--snpFile FILE</code>	Path to the input snp file.
176	<code>--indFile FILE</code>	Path to the input ind file.
177	<code>--snpSet SET</code>	The snpSet of the package: 1240K, HumanOrigins or 178 Other. Only relevant for data input with <code>-p --genoOne</code> 179 or <code>--genoFile</code> + <code>--snpFile</code> + <code>--indFile</code> , because the 180 packages in a <code>-d --baseDir</code> already have this 181 information in their respective <code>POSEIDON.yml</code> files. 182 (default: Other)
183	<code>-o,--outPackagePath DIR</code>	Path to the output package directory.
184	<code>-n,--outPackageName STRING</code>	The output package name. This is optional: If no name 185 is provided, then the package name defaults to the 186 basename of the (mandatory) <code>--outPackagePath</code> 187 argument. (default: Nothing)
188	<code>--minimal</code>	Should the output data be reduced to a necessary 189 minimum and omit empty scaffolding? 190

191 The command

```
192 trident init \  
193   --inFormat EIGENSTRAT/PLINK \  
194   --genoFile path/to/geno_file \  

```

```

195 --snpFile path/to/snp_file \
196 --indFile path/to/ind_file \
197 --snpSet 1240K|HumanOrigins|Other \
198 -o path/to/new_package_name

```

199 requires the format (`--inFormat`) of your input data (either EIGENSTRAT or PLINK), the paths to the respective  
 200 files (`--genoFile`, `--snpFile`, `--indFile`), and optionally the “shape” of these files (`--snpSet`), so if they cover  
 201 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

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

## 205 2.2 Fetch command

206 `fetch` allows to download Poseidon packages from a remote Poseidon server via a [Web API](#). Read more about  
 207 the data available with it [here](#).

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

```

209 Usage: trident fetch (-d|--baseDir DIR)
210             (--downloadAll |
211             (--fetchFile FILE | (-f|--fetchString DSL)))
212             [--remoteURL URL] [--archive STRING]

```

214 Download data from a remote Poseidon repository

216 Available options:

217 <code>-h,--help</code>	Show this help text
218 <code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
219 <code>--downloadAll</code>	Download all packages the server is offering.
220 <code>--fetchFile FILE</code>	A file with a list of packages. Works just as <code>-f</code> , but
221	multiple values can also be separated by newline, not
222	just by comma. <code>-f</code> and <code>--fetchFile</code> can be combined.
223 <code>-f,--fetchString DSL</code>	List of packages to be downloaded from the remote
224	server. Package names should be wrapped in asterisks:
225	<code>*package_title*</code> . You can combine multiple values with
226	comma, so for example: <code>"*package_1*, *package_2*,</code>
227	<code>*package_3*"</code> . <code>fetchString</code> uses the same parser as
228	<code>forgeString</code> , but does not allow excludes. If groups
229	or individuals are specified, then packages which
230	include these groups or individuals are included in

231 the download.

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

233 `--archive STRING` The name of the Poseidon package archive that should  
 234 be queried. If not given, then the query falls back  
 235 to the default archive of the server selected with  
 236 `--remoteURL`. See the archive documentation at  
 237 [https://www.poseidon-adna.org/#/archive\\_overview](https://www.poseidon-adna.org/#/archive_overview) for  
 238 a list of archives currently available from the  
 239 official Poseidon Web API. (default: Nothing)

241 It works with

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

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

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

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

256 `fetch` also has the optional arguments `--remote https://...` to name an alternative Poseidon server and  
 257 `--archive` to select a Poseidon archive on the server. Here is a list of the [archives available on the official](#)  
 258 [Poseidon server](#).

## 259 2.3 Forge command

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

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

```
263 Usage: trident forge ((-d|--baseDir DIR) |
264   ((-p|--genoOne FILE) | --inFormat FORMAT --genoFile FILE
265   --snpFile FILE --indFile FILE) [--snpSet SET])
266   [--forgeFile FILE | (-f|--forgeString DSL)]
267   [--selectSnps FILE] [--intersect] [--outFormat FORMAT]
268   [--minimal] [--onlyGeno] (-o|--outPackagePath DIR)
269   [-n|--outPackageName STRING] [--packagewise]
270   [--outPlinkPopName MODE]
```

```

272 Select packages, groups or individuals and create a new Poseidon package from
273 them
274
275 Available options:
276 -h,--help Show this help text
277 -d,--baseDir DIR A base directory to search for Poseidon packages.
278 -p,--genoOne FILE One of the input genotype data files. Expects .bed,
279 .bim or .fam for PLINK and .geno, .snp or .ind for
280 EIGENSTRAT. The other files must be in the same
281 directory and must have the same base name.
282 --inFormat FORMAT The format of the input genotype data: EIGENSTRAT or
283 PLINK. Only necessary for data input with --genoFile
284 + --snpFile + --indFile.
285 --genoFile FILE Path to the input geno file.
286 --snpFile FILE Path to the input snp file.
287 --indFile FILE Path to the input ind file.
288 --snpSet SET The snpSet of the package: 1240K, HumanOrigins or
289 Other. Only relevant for data input with -p|--genoOne
290 or --genoFile + --snpFile + --indFile, because the
291 packages in a -d|--baseDir already have this
292 information in their respective POSEIDON.yml files.
293 (default: Other)
294 --forgeFile FILE A file with a list of packages, groups or individual
295 samples. Works just as -f, but multiple values can
296 also be separated by newline, not just by comma.
297 Empty lines are ignored and comments start with "#",
298 so everything after "#" is ignored in one line.
299 Multiple instances of -f and --forgeFile can be
300 given. They will be evaluated according to their
301 input order on the command line.
302 -f,--forgeString DSL List of packages, groups or individual samples to be
303 combined in the output package. Packages follow the
304 syntax *package_title*, populations/groups are simply
305 group_id and individuals <individual_id>. You can
306 combine multiple values with comma, so for example:
307 "*package_1*, <individual_1>, <individual_2>,
308 group_1". Duplicates are treated as one entry.
309 Negative selection is possible by prepending "-" to
310 the entity you want to exclude (e.g. "*package_1*,
311 -<individual_1>, -group_1"). forge will apply
312 excludes and includes in order. If the first entity
313 is negative, then forge will assume you want to merge
314 all individuals in the packages found in the baseDirs
315 (except the ones explicitly excluded) before the
316 exclude entities are applied. An empty forgeString

```



317 (and no `--forgeFile`) will therefore merge all  
 318 available individuals. If there are individuals in  
 319 your input packages with equal individual id, but  
 320 different main group or source package, they can be  
 321 specified with the special syntax  
 322 "`<package:group:individual>`".

323 `--selectSnps FILE` To extract specific SNPs during this forge operation,  
 324 provide a Snp file. Can be either Eigenstrat (file  
 325 ending must be `'.snp'`) or Plink (file ending must be  
 326 `'.bim'`). When this option is set, the output package  
 327 will have exactly the SNPs listed in this file. Any  
 328 SNP not listed in the file will be excluded. If  
 329 option `'--intersect'` is also set, only the SNPs  
 330 overlapping between the SNP file and the forged  
 331 packages are output. (default: Nothing)

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

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

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

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

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

346 `-n,--outPackageName STRING`  
 347 The output package name. This is optional: If no name  
 348 is provided, then the package name defaults to the  
 349 basename of the (mandatory) `--outPackagePath`  
 350 argument. (default: Nothing)

351 `--packagewise` Skip the within-package selection step in forge. This  
 352 will result in outputting all individuals in the  
 353 relevant packages, and hence a superset of the  
 354 requested individuals/groups. It may result in better  
 355 performance in cases where one wants to forge entire  
 356 packages or almost entire packages. Details: Forge  
 357 conceptually performs two types of selection: First,  
 358 it identifies which packages in the supplied base  
 359 directories are relevant to the requested forge, i.e.  
 360 whether they are either explicitly listed using  
 361 `*PackageName*`, or because they contain selected

```

362         individuals or groups. Second, within each relevant
363         package, individuals which are not requested are
364         removed. This option skips only the second step, but
365         still performs the first.
366     --outPlinkPopName MODE   Where to write the population/group name into the FAM
367                             file in Plink-format. Three options are possible:
368                             asFamily (default) | asPhenotype | asBoth. See also
369                             --inPlinkPopName.
370
371     forge can be used with
372
373     trident forge -d ... -d ... \
374         -f "*package_name*, group_id, <individual_id>" \
375         -o path/to/new_package_name
376
377     where the entities (packages, groups/populations, individuals/samples) you want in the output package can be
378     denoted either as a string on the command line (-f/--forgeString), or in an input text file (--forgeFile).
379     See the section below for the syntax of this selection language. Do not forget to wrap the --forgeString query
380     in quotes.
381
382     Including one or multiple Poseidon packages with -d is not the only way to include data for a forge operation.
383     It is also possible to consider unpackaged genotype data directly with -p (+ --snpSet) or --inFormat +
384     --genoFile + --snpFile + --indFile (+ --snpSet). This makes the following example possible, where we
385     merge data from one Poseidon package and two genotype datasets to get a new EIGENSTRAT dataset.
386
387     trident forge \
388         -d 2017_GonzalesFortesCurrentBiology \
389         -p 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \
390         --inFormat PLINK \
391         --genoFile 2017_HaberAJHG/2017_HaberAJHG.bed \
392         --snpFile 2017_HaberAJHG/2017_HaberAJHG.bim \
393         --indFile 2017_HaberAJHG/2017_HaberAJHG.fam \
394         -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \
395         -o testpackage \
396         --outFormat EIGENSTRAT \
397         --onlyGeno

```

### 393 2.3.1 The forge selection language

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

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

- 399 • Each package title is surrounded by \*: **\*package\***. That means if you want all individuals of the Poseidon  
400 package 2019\_Jeong\_InnerEurasia in the output package you would add **\*2019\_Jeong\_InnerEurasia\***  
401 to the query.

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

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

```
# Packages
*package1*, *package2*

# Groups and individuals from other packages beyond package1 and package2
group1, <individual1>, group2, <individual2>, <individual3>

# group2 has two outlier individuals that should be ignored
-<bad_individual1> # This one has very low coverage
-<bad_individual2> # This one is from a different time period
```

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

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

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

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

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

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

The following example illustrates the described behaviour:

**A.janno**

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

#### B.janno

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

#### A.janno + B.janno

Poseidon_ID	Group_Name	Genetic_Sex	AdditionalColumn1	AdditionalColumn2	AdditionalColumn3
XXX011	POP1	M	A	D	n/a
XXX012	POP2	F	B	E	n/a
XXX013	POP1	M	C	F	n/a
YYY022	POP5	F	n/a	J	G
YYY023	POP5	F	n/a	K	H
YYY024	POP5	M	n/a	L	I

#### 2.3.3 Treatment of the .ssf file while merging

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

#### 2.3.4 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).

#### 2.3.5 Other options

Just as for `init` the output package of `forge` is created as a new directory `-o`. The title can also be explicitly defined with `-n`.

`--minimal` allows for the creation of a minimal output package without `.bib` and `.janno`. This is especially useful for data analysis pipelines, where only the genotype data is required. Even more basic output comes with `--onlyGeno`, which means that only the genotype data is returned without any Poseidon package.

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

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

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

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

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

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

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

## 477 2.4 Genoconvert command

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

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

```
481 Usage: trident genoconvert ((-d|--baseDir DIR) |
482                             ((-p|--genoOne FILE) | --inFormat FORMAT
483                             --genoFile FILE --snpFile FILE --indFile FILE)
484                             [--snpSet SET]) --outFormat FORMAT [--onlyGeno]
485                             [-o|--outPackagePath DIR] [--removeOld]
486                             [--outPlinkPopName MODE]
```

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

489 Available options:

491 -h,--help	Show this help text
492 -d,--baseDir DIR	A base directory to search for Poseidon packages.
493 -p,--genoOne FILE	One of the input genotype data files. Expects .bed,

```

494         .bim or .fam for PLINK and .geno, .snp or .ind for
495         EIGENSTRAT. The other files must be in the same
496         directory and must have the same base name.
497     --inFormat FORMAT      The format of the input genotype data: EIGENSTRAT or
498                             PLINK. Only necessary for data input with --genoFile
499                             + --snpFile + --indFile.
500     --genoFile FILE        Path to the input geno file.
501     --snpFile FILE         Path to the input snp file.
502     --indFile FILE         Path to the input ind file.
503     --snpSet SET           The snpSet of the package: 1240K, HumanOrigins or
504                             Other. Only relevant for data input with -p|--genoOne
505                             or --genoFile + --snpFile + --indFile, because the
506                             packages in a -d|--baseDir already have this
507                             information in their respective POSEIDON.yml files.
508                             (default: Other)
509     --outFormat FORMAT     the format of the output genotype data: EIGENSTRAT or
510                             PLINK.
511     --onlyGeno             Should only the resulting genotype data be returned?
512                             This means the output will not be a Poseidon package.
513     -o,--outPackagePath DIR Path to the output package directory. This is
514                             optional: If no path is provided, then the output is
515                             written to the directories where the input genotype
516                             data file (.bed/.geno) is stored. (default: Nothing)
517     --removeOld            Remove the old genotype files when creating the new
518                             ones.
519     --outPlinkPopName MODE Where to write the population/group name into the FAM
520                             file in Plink-format. Three options are possible:
521                             asFamily (default) | asPhenotype | asBoth. See also
522                             --inPlinkPopName.
523
524     With the default setting
525
526     trident genoconvert -d ... -d ... --outFormat EIGENSTRAT|PLINK
527
528     all packages in -d will be converted to the desired --outFormat (either EIGENSTRAT or PLINK), if the data is
529     not already in this format. This includes updating the respective POSEIDON.yml files.
530
531     The “old” data is not deleted, but kept around. That means conversion can result in a package with both PLINK
532     and EIGENSTRAT data, but only one is linked in the POSEIDON.yml file, and that is what will be used by
533     trident. To delete the old data in the conversion you can add the --removeOld flag.
534
535     Instead of -d to change Poseidon packages, the -p (+ --snpSet) or --inFormat + --genoFile + --snpFile
536     + --indFile (+ --snpSet) allow to directly convert genotype data that is not wrapped in a Poseidon package
537     and store it to a directory given in -o. See this example:
538
539     trident genoconvert \
540         -p 2018_Mittnik_Baltic/Mittnik_Baltic.bed \
541         --outFormat EIGENSTRAT
542         -o my_directory

```

## 2.5 Rectify command

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

Click here for command line details

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

Adjust POSEIDON.yml files automatically to package changes

Available options:

<code>-h,--help</code>	Show this help text
<code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
<code>--ignorePoseidonVersion</code>	Read packages even if their poseidonVersion is not compatible with trident.
<code>--poseidonVersion ?.??.?</code>	Poseidon version the packages should be updated to: e.g. "2.5.3".
<code>--packageVersion VPART</code>	Part of the package version number in the POSEIDON.yml file that should be updated: Major, Minor or Patch (see <a href="https://semver.org">https://semver.org</a> ).
<code>--logText STRING</code>	Log text for this version in the CHANGELOG file.
<code>--checksumAll</code>	Update all checksums.
<code>--checksumGeno</code>	Update genotype data checksums.
<code>--checksumJanno</code>	Update .janno file checksum.
<code>--checksumSSF</code>	Update .ssf file checksum
<code>--checksumBib</code>	Update .bib file checksum.
<code>--newContributors DSL</code>	Contributors to add to the POSEIDON.yml file in the form "[Firstname Lastname](Email address);..."

It can be called with a lot of optional arguments:

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

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

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

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

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

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

## 587 3 Inspection commands

### 588 3.1 List command

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

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

```
591 Usage: trident list ((-d|--baseDir DIR) | --remote [--remoteURL URL]
592                  [--archive STRING])
593                  (--packages | --groups | --individuals
594                  [-j|--jannoColumn COLNAME]) [--raw]
```

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

598  
599 Available options:

600 <h>-h,--help</h>	Show this help text
601 <h>-d,--baseDir DIR</h>	A base directory to search for Poseidon packages.
602 <h>--remote</h>	List packages from a remote server instead the local 603 file system.
604 <h>--remoteURL URL</h>	URL of the remote Poseidon server. 605 (default: "https://server.poseidon-adna.org")
606 <h>--archive STRING</h>	The name of the Poseidon package archive that should 607 be queried. If not given, then the query falls back 608 to the default archive of the server selected with 609 --remoteURL. See the archive documentation at 610 https://www.poseidon-adna.org/#/archive_overview for 611 a list of archives currently available from the 612 official Poseidon Web API. (default: Nothing)
613 <h>--packages</h>	List all packages.
614 <h>--groups</h>	List all groups, ignoring any group names after the 615 first as specified in the .janno-file.
616 <h>--individuals</h>	List all individuals/samples.
617 <h>-j,--jannoColumn COLNAME</h>	List additional fields from the janno files, using 618 the .janno column heading name, such as "Country", 619 "Site", "Date_C14_Uncal_BP", etc..



620     --raw                         Return the output table as tab-separated values  
621                                     without header. This is useful for piping into grep  
622                                     or awk.

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

624 `trident list -d ... -d ... --packages`

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

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

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

628 will result in a view of all packages available in one of the [public online archives](#). Just as for `fetch`, the `--archive`  
629 flag allows to choose which public archive to query.

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

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

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

## 640 3.2 Summarise command

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

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

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

644

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

646

647 Available options:

648     -h,--help                     Show this help text

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

650     --raw                         Return the output table as tab-separated values  
651                                     without header. This is useful for piping into grep  
652                                     or awk.

653 You can run it with

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

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

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

### 659 3.3 Survey command

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

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

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

663  
664 Survey the degree of context information completeness for Poseidon packages

665  
666 Available options:

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

672 Running

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

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

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

### 677 3.4 Validate command

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

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

680 Usage: `trident validate ((-d|--baseDir DIR) [--ignoreGeno] [--fullGeno]`  
681  `[--ignoreDuplicates] [-c|--ignoreChecksums]`  
682  `[--ignorePoseidonVersion] |`  
683  `--pyml FILE | (-p|--genoOne FILE) | --inFormat FORMAT`  
684  `--genoFile FILE --snpFile FILE --indFile FILE |`  
685  `--janno FILE | --ssf FILE | --bib FILE) [--noExitCode]`

686  
687 Check Poseidon packages or package components for structural correctness

688  
689 Available options:

690 <code>-h,--help</code>	Show this help text
691 <code>-d,--baseDir DIR</code>	A base directory to search for Poseidon packages.
692 <code>--ignoreGeno</code>	Ignore snp and geno file.
693 <code>--fullGeno</code>	Test parsing of all SNPs (by default only the first 100 SNPs are probed).
695 <code>--ignoreDuplicates</code>	Do not stop on duplicated individual names in the package collection.

```

697 -c,--ignoreChecksums      Whether to ignore checksums. Useful for speedup in
698                          debugging.
699 --ignorePoseidonVersion  Read packages even if their poseidonVersion is not
700                          compatible with trident.
701 --pym1 FILE               Path to a POSEIDON.yml file.
702 -p,--genoOne FILE         One of the input genotype data files. Expects .bed,
703                          .bim or .fam for PLINK and .geno, .snp or .ind for
704                          EIGENSTRAT. The other files must be in the same
705                          directory and must have the same base name.
706 --inFormat FORMAT         The format of the input genotype data: EIGENSTRAT or
707                          PLINK. Only necessary for data input with --genoFile
708                          + --snpFile + --indFile.
709 --genoFile FILE           Path to the input geno file.
710 --snpFile FILE            Path to the input snp file.
711 --indFile FILE            Path to the input ind file.
712 --janno FILE              Path to a .janno file.
713 --ssf FILE                Path to a .ssf file.
714 --bib FILE                Path to a .bib file.
715 --noExitCode              Do not produce an explicit exit code.

```

716 You can run it with

```

717 trident validate -d ... -d ...

```

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

719 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.

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

- 726 • Structural correctness of the POSEIDON.yml file.
- 727 • Presence of all files references in the POSEIDON.yml file.
- 728 • Full structural correctness of .janno, .ssf and .bib file.
- 729 • Superficial correctness of genotype data files by parsing the first 100 SNPs. A full check that parses all SNPs can be triggered with the `--fullGeno` option. `--ignoreGeno`, on the other hand, causes `validate` to ignore the genotype data entirely, which speeds up the validation significantly.
- 732 • Correspondence of BibTeX keys in .bib and .janno
- 733 • Correspondence of sample IDs in .janno and .ssf.
- 734 • Correspondence of sample and group IDs in .janno and genotype data files.

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

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