

# Guide for trident v1.1.6.0

## Contents

<b>1</b>	<b>Poseidon package repositories</b>	<b>1</b>
<b>2</b>	<b>Analysing your own dataset outside of the main repository</b>	<b>2</b>
<b>3</b>	<b>Package creation and manipulation commands</b>	<b>3</b>
3.1	Init command . . . . .	3
3.2	Fetch command . . . . .	4
3.3	Forge command . . . . .	5
3.3.1	The forge selection language . . . . .	8
3.3.2	Other options . . . . .	9
3.3.3	Treatment of the .janno file while merging . . . . .	9
3.4	Genoconvert command . . . . .	10
3.5	Update command . . . . .	11
<b>4</b>	<b>Inspection commands</b>	<b>13</b>
4.1	List command . . . . .	13
4.2	Summarise command . . . . .	14
4.3	Survey command . . . . .	15
4.4	Validate command . . . . .	15

## 1 Poseidon package repositories

Trident generally requires Poseidon “packages” to work with (since version 0.28.0 it also supports direct interaction with “unpacked” genotype data – see `-p` below). 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
```

```

33     /ancient
34     /...
35     /...
36     /Reference_Genomes
37     /...
38     /...
39     /Archaic_Humans
40     /...
41     /...

```

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

44 Let's use the `list` command to list all packages in the `modern` and `Reference_Genomes` :

```

45 trident list -d /path/to/poseidon/packages/modern \
46 -d /path/to/poseidon/packages/ReferenceGenomes --packages

```

## 47 2 Analysing your own dataset outside of the main repository

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

```

52 ~/my_project/my_project.geno
53 ~/my_project/my_project.snp
54 ~/my_project/my_project.ind

```

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

```

57 poseidonVersion: 2.5.0
58 title: My_awesome_project
59 description: Unpublished genetic data from my awesome project
60 contributor:
61   - name: Stephan Schiffels
62     email: schiffels@institute.org
63 packageVersion: 0.1.0
64 lastModified: 2020-10-07
65 genotypeData:
66   format: EIGENSTRAT
67   genoFile: my_project.geno
68   snpFile: my_project.snp
69   indFile: my_project.ind
70 jannoFile: my_project.janno
71 bibFile: sources.bib

```

72 Two remarks: 1) all file paths are considered *relative* to the directory in which `POSEIDON.yml` resides. Here I

73 assume that you put this file into the same directory as the three genotype files. 2) Besides the genotype data  
 74 files there are two (technically optional) files referenced by this example `POSEIDON.yml` file: `sources.bib`  
 75 and `my_project.janno`. Of course you can add them manually - `init` automatically creates empty dummy  
 76 versions.

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

```
79 trident list -d /path/to/poseidon/packages/modern \  
80     -d /path/to/poseidon/packages/ReferenceGenomes  
81     -d ~/my_project --packages
```

## 82 3 Package creation and manipulation commands

### 83 3.1 Init command

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

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

```
87 Usage: trident init ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG  
88     --snpFile ARG --indFile ARG) [--snpSet ARG]  
89     (-o|--outPackagePath ARG) [-n|--outPackageName ARG]  
90     [--minimal]
```

91 Create a new Poseidon package from genotype data

92 Available options:

94	<code>-h,--help</code>	Show this help text
95	<code>-p,--genoOne ARG</code>	one of the input genotype data files. Expects <code>.bed</code> or
96		<code>.bim</code> or <code>.fam</code> for PLINK and <code>.geno</code> or <code>.snp</code> or <code>.ind</code> for
97		EIGENSTRAT. The other files must be in the same
98		directory and must have the same base name
99	<code>--inFormat ARG</code>	the format of the input genotype data: EIGENSTRAT or
100		PLINK (only necessary for data input with <code>--genoFile</code>
101		+ <code>--snpFile</code> + <code>--indFile</code> )
102	<code>--genoFile ARG</code>	the input geno file path
103	<code>--snpFile ARG</code>	the input snp file path
104	<code>--indFile ARG</code>	the input ind file path
105	<code>--snpSet ARG</code>	the snpSet of the new package: 1240K, HumanOrigins or
106		Other. Default: Other
107	<code>-o,--outPackagePath ARG</code>	the output package directory path
108	<code>-n,--outPackageName ARG</code>	the output package name - this is optional: If no
109		name is provided, then the package name defaults to
110		the basename of the (mandatory) <code>--outPackagePath</code>
111		argument
112	<code>--minimal</code>	should only a minimal output package be created?

113 The command

```

114 trident init \
115     --inFormat EIGENSTRAT/PLINK \
116     --genoFile path/to/geno_file \
117     --snpFile path/to/snp_file \
118     --indFile path/to/ind_file \
119     --snpSet 1240K|HumanOrigins|Other \
120     -o path/to/new_package_name

```

121 requires the format ( `--inFormat` ) of your input data (either `EIGENSTRAT` or `PLINK` ), the paths to the  
122 respective files ( `--genoFile` , `--snpFile` , `--indFile` ), and optionally the “shape” of these files ( `--snpSet` ),  
123 so if they cover the `1240K` , the `HumanOrigins` or an `Other` SNP set. A simpler interface added in trident  
124 0.29.0 is available with `-p (+ --snpSet)` .

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

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

## 129 3.2 Fetch command

130 `fetch` allows to download poseidon packages from a remote poseidon server.

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

```

132 Usage: trident fetch (-d|--baseDir DIR)
133             (--downloadAll |
134             (--fetchFile ARG | (-f|--fetchString ARG)))
135             [--remoteURL ARG] [-u|--upgrade]

```

136 Download data from a remote Poseidon repository

137 Available options:

139 <code>-h,--help</code>	Show this help text
140 <code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (could be a Poseidon repository)
141 <code>--downloadAll</code>	download all packages the server is offering
142 <code>--fetchFile ARG</code>	A file with a list of packages. Works just as <code>-f</code> , but multiple values can also be separated by newline, not just by comma. <code>-f</code> and <code>--fetchFile</code> can be combined.
143 <code>-f,--fetchString ARG</code>	List of packages to be downloaded from the remote server. Package names should be wrapped in asterisks: <code>*package_title*</code> . You can combine multiple values with comma, so for example: <code>"*package_1*, *package_2*,</code>

```

150         *package_3*". fetchString uses the same parser as
151         forgeString, but does not allow excludes. If groups
152         or individuals are specified, then packages which
153         include these groups or individuals are included in
154         the download.
155     --remoteURL ARG        URL of the remote Poseidon server
156                           (default: "https://c107-224.cloud.gwdg.de")
157     -u,--upgrade           overwrite outdated local package versions

```

158 It works with

```

159 trident fetch -d ... -d ... \
160     -f "*package_title_1*,*package_title_2*,*package_title_3*,group_name,<Individual1>" \
161     --fetchFile path/to/forgFile

```

162 and the entities you want to download must be listed either in one or more simple strings with comma-separated  
163 values, which can be passed via one or multiple options `-f / --fetchString`, or in one or more text files  
164 (`--fetchFile`). Entities are then combined from these sources. Entities are specified using a special syntax:  
165 Package titles are wrapped in asterisks: *package\_title* (see also the documentation of `forge` below), group  
166 names are spelled as is, and individual names are wrapped in angular brackets, like `<Individual1>`. Fetch will  
167 figure out which packages need to be downloaded to include all specified entities. `--downloadAll`, which can be  
168 given instead of `-f` and `--fetchFile`, causes fetch to download all packages from the server. The downloaded  
169 packages are added in the first (!) `-d` directory (which gets created if it doesn't exist), but downloads are only  
170 performed if the respective packages are not already present in an up-to-date version in any of the `-d` dirs.

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

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

175 To overwrite outdated package versions with `fetch`, the `-u / --upgrade` flag has to be set. Note that many  
176 file systems do not offer a way to recover overwritten files. So be careful with this switch.

### 177 3.3 Forge command

178 `forge` creates new poseidon packages by extracting and merging packages, populations and individuals from  
179 your poseidon repositories.

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

```

181 Usage: trident forge ((-d|--baseDir DIR) |
182                     ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG
183                     --snpFile ARG --indFile ARG) [--snpSet ARG])
184                     [--forgeFile ARG | (-f|--forgeString ARG)]
185                     [--selectSnps ARG] [--intersect] [--outFormat ARG]
186                     [--minimal] [--onlyGeno] (-o|--outPackagePath ARG)
187                     [-n|--outPackageName ARG] [--no-extract]
188     Select packages, groups or individuals and create a new Poseidon package from
189     them
190

```

```

191 Available options:
192   -h,--help           Show this help text
193   -d,--baseDir DIR    a base directory to search for Poseidon Packages
194                       (could be a Poseidon repository)
195   -p,--genoOne ARG    one of the input genotype data files. Expects .bed or
196                       .bim or .fam for PLINK and .geno or .snp or .ind for
197                       EIGENSTRAT. The other files must be in the same
198                       directory and must have the same base name
199   --inFormat ARG      the format of the input genotype data: EIGENSTRAT or
200                       PLINK (only necessary for data input with --genoFile
201                       + --snpFile + --indFile)
202   --genoFile ARG      the input geno file path
203   --snpFile ARG       the input snp file path
204   --indFile ARG       the input ind file path
205   --snpSet ARG        the snpSet of the new package: 1240K, HumanOrigins or
206                       Other. Default: Other
207   --forgeFile ARG     A file with a list of packages, groups or individual
208                       samples. Works just as -f, but multiple values can
209                       also be separated by newline, not just by comma.
210                       Empty lines are ignored and comments start with "#",
211                       so everything after "#" is ignored in one line.
212                       Multiple instances of -f and --forgeFile can be
213                       given. They will be evaluated according to their
214                       input order on the command line.
215   -f,--forgeString ARG List of packages, groups or individual samples to be
216                       combined in the output package. Packages follow the
217                       syntax *package_title*, populations/groups are simply
218                       group_id and individuals <individual_id>. You can
219                       combine multiple values with comma, so for example:
220                       "*package_1*, <individual_1>, <individual_2>,
221                       group_1". Duplicates are treated as one entry.
222                       Negative selection is possible by prepending "-" to
223                       the entity you want to exclude (e.g. "*package_1*,
224                       -<individual_1>, -group_1"). forge will apply
225                       excludes and includes in order. If the first entity
226                       is negative, then forge will assume you want to merge
227                       all individuals in the packages found in the baseDirs
228                       (except the ones explicitly excluded) before the
229                       exclude entities are applied. An empty forgeString
230                       (and no --forgeFile) will therefore merge all
231                       available individuals.
232   --selectSnps ARG    To extract specific SNPs during this forge operation,
233                       provide a Snp file. Can be either Eigenstrat (file
234                       ending must be '.snp') or Plink (file ending must be
235                       '.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.

`--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 ARG` the format of the output genotype data: EIGENSTRAT or PLINK. Default: PLINK

`--minimal` should only a minimal output package be created?

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

`-o,--outPackagePath ARG` the output package directory path

`-n,--outPackageName ARG` the output package name - this is optional: If no name is provided, then the package name defaults to the basename of the (mandatory) `--outPackagePath` argument

`--no-extract` Skip the selection step in forge. This will result in outputting all individuals in the relevant packages, and hence a superset of the requested individuals/groups. It may result in better performance in cases where one wants to forge entire packages or almost entire packages. Note that this will also ignore any ordering in the output groups/individuals. With this option active, individuals from the relevant packages will just be written in the order that they appear in the original packages.

`forge` can be used with

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

where the entities (packages, groups/populations, individuals/samples) you want in the output package can be denoted either as one or more simple strings with comma-separated values via one or more (`-f` / `--forgeString`) options, or in one or more text files (`--forgeFile`). Because the order in which inclusions and exclusions are given, the order strictly follows the order as these strings are given via options `-f` / `--forgeString` and `--forgeFile`.

Including one or multiple Poseidon packages with `-d` is not the only way to include data for a forge

operation. It is also possible to include unpackaged genotype data directly with `-p (+ --snpSet)` or `--inFormat + --genoFile + --snpFile + --indFile (+ --snpSet)`. This makes the following example possible, where we merge data from one Poseidon package and two genotype datasets to get a new EIGENSTRAT dataset.

```
trident forge \
  -d 2017_GonzalesFortesCurrentBiology \
  -p 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \
  --inFormat PLINK \
  --genoFile 2017_HaberAJHG/2017_HaberAJHG.bed \
  --snpFile 2017_HaberAJHG/2017_HaberAJHG.bim \
  --indFile 2017_HaberAJHG/2017_HaberAJHG.fam \
  -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \
  -o testpackage \
  --outFormat EIGENSTRAT \
  --onlyGeno
```

### 3.3.1 The forge selection language

Entities in the `--forgeString` or the `--forgeFile` have to be marked in a certain way:

- Each package is surrounded by `*`, so if you want all individuals of `2019_Jeong_InnerEurasia` in the output package you would add `*2019_Jeong_InnerEurasia*` to the list.
- Groups/populations are not specially marked. So to get all individuals of the group `Swiss_Roman_period`, you would simply add `Swiss_Roman_period`.
- Individuals/samples are surrounded by `<` and `>`, so `ALA026` becomes `<ALA026>`.

Do not forget to wrap the `forgeString` in quotes.

You can use both `-f / --forgeString` and `--forgeFile` and even combine multiple of each. They are evaluated in order.

In the file 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



321 packages found in the baseDirs (except the ones explicitly excluded, of course). An empty forgeString will  
 322 therefore merge all available individuals.

### 323 3.3.2 Other options

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

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

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

334 `--intersect` also influences the automatic determination of the `snpSet` field in the POSEIDON.yml file for  
 335 the resulting package. If the `snpSet` s of all input packages are identical, then the resulting package will just  
 336 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

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

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

### 343 3.3.3 Treatment of the .janno file while merging

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

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

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

354 The following example illustrates the described behaviour:

#### 355 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

#### 356 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

#### 357 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

### 358 3.4 Genoconvert command

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

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

362 Usage: trident genoconvert ((-d|--baseDir DIR) |  
 363 ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG  
 364 --snpFile ARG --indFile ARG) [--snpSet ARG])  
 365 --outFormat ARG [--onlyGeno]  
 366 [-o|--outPackagePath ARG] [--removeOld]

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

368 Available options:

370 -h,--help Show this help text  
 371 -d,--baseDir DIR a base directory to search for Poseidon Packages  
 372 (could be a Poseidon repository)  
 373 -p,--genoOne ARG one of the input genotype data files. Expects .bed or

```

374         .bim or .fam for PLINK and .geno or .snp or .ind for
375         EIGENSTRAT. The other files must be in the same
376         directory and must have the same base name
377     --inFormat ARG         the format of the input genotype data: EIGENSTRAT or
378                             PLINK (only necessary for data input with --genoFile
379                             + --snpFile + --indFile)
380     --genoFile ARG         the input geno file path
381     --snpFile ARG          the input snp file path
382     --indFile ARG          the input ind file path
383     --snpSet ARG           the snpSet of the new package: 1240K, HumanOrigins or
384                             Other. Default: Other
385     --outFormat ARG        the format of the output genotype data: EIGENSTRAT or
386                             PLINK.
387     --onlyGeno             should only the resulting genotype data be returned?
388                             This means the output will not be a Poseidon package
389     -o,--outPackagePath ARG the output package directory path - this is optional:
390                             If no path is provided, then the output is written to
391                             the directories where the input genotype data file
392                             (.bed/.geno) is stored
393     --removeOld            Remove the old genotype files when creating the new
394                             ones

```

395 With the default setting

```

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

```

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

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

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

```

405 trident genoconvert \
406     -p 2018_Mittnik_Baltic/Mittnik_Baltic.bed \
407     --outFormat EIGENSTRAT
408     -o my_directory

```

### 409 3.5 Update command

410 `update` automatically harmonizes POSEIDON.yml files of one or multiple packages if the packages were  
411 changed. This is not an automatic update from one Poseidon version to the next!

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

```

413 Usage: trident update (-d|--baseDir DIR) [--poseidonVersion ARG]
414                 [--ignorePoseidonVersion] [--versionComponent ARG]

```

```

415             [--noChecksumUpdate] [--newContributors ARG]
416             [--logText ARG] [--force]
417 Update POSEIDON.yml files automatically
418
419 Available options:
420 -h,--help          Show this help text
421 -d,--baseDir DIR    a base directory to search for Poseidon Packages
422                     (could be a Poseidon repository)
423 --poseidonVersion ARG Poseidon version the packages should be updated to:
424                     e.g. "2.5.3" (default: Nothing)
425 --ignorePoseidonVersion Read packages even if their poseidonVersion is not
426                     compatible with the trident version. The assumption
427                     is, that the package is already structurally adjusted
428                     to the trident version and only the version number is
429                     lagging behind.
430 --versionComponent ARG Part of the package version number in the
431                     POSEIDON.yml file that should be updated: Major,
432                     Minor or Patch (see https://semver.org)
433                     (default: Patch)
434 --noChecksumUpdate    Should update of checksums in the POSEIDON.yml file
435                     be skipped
436 --ignoreGeno          ignore SNP and GenoFile
437 --newContributors ARG Contributors to add to the POSEIDON.yml file in the
438                     form "[Firstname Lastname](Email address);..."
439 --logText ARG         Log text for this version jump in the CHANGELOG file
440                     (default: "not specified")
441 --force              Normally the POSEIDON.yml files are only changed if
442                     the poseidonVersion is adjusted or any of the
443                     checksums change. With --force a package version
444                     update can be triggered even if this is not the case.
445
446 It can be called with a lot of optional arguments
447
448 trident update -d ... -d ... \
449   --poseidonVersion "X.X.X" \
450   --versionComponent Major/Minor/Patch \
451   --noChecksumUpdate
452   --ignoreGeno
453   --newContributors "[Firstname Lastname](Email address);..."
454   --logText "short description of the update"
455   --force
456
457 By default update will not edit a package's POSEIDON.yml file, even when arguments like --versionComponent ,
458 --newContributors or --logText are explicitly set. This default exists to run the function on a large set of
459 packages where only few of them were edited and need an active update. A package will only be modified by
460 update if either
461
462 • any of the files with checksums (e.g. the genotype data) in it were modified,

```

- the `--poseidonVersion` argument differs from the `poseidonVersion` in the package's POSEIDON.yml file
- or the `--force` flag was set in `update`.

If any of these applies to a package in the search directory (`--baseDir / -d`), it will be updated. This includes the following steps:

- If `--poseidonVersion` is different from the `poseidonVersion` field in the package, then that will be updated.
- The `packageVersion` will be incremented. If `--versionComponent` is not set, then it falls back to `Patch`, so a change in the last position of the three digit version number. `Minor` increments the middle, and `Major` the first position (see [semantic versioning](#)).
- The `lastModified` field will be updated to the current day (based on your computer's system time).
- The contributors in `--newContributors` will be added to the `contributor` field if they're not there already.
- If any checksums changed, then they will be updated. If certain checksums are not set yet, then they will be added. The checksum update can be skipped with `--noChecksumUpdate` or partially skipped for the genotype data with `--ignoreGeno`.
- The CHANGELOG.md file will be updated with a new row for the new version and the text in `--logText` (default: "not specified"), which will be appended as the first line of the file. If no CHANGELOG.md file exists, then it will be created and referenced in the POSEIDON.yml file.

:heavy\_exclamation\_mark: As `update` reads and rewrites POSEIDON.yml files, it may change their inner order, layout or even content (e.g. if they have fields which are not in the [Poseidon package definition](#)). Create a backup of the POSEIDON.yml file before running `update` if you are uncertain.

## 4 Inspection commands

### 4.1 List command

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

Click here for command line details

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

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

Available options:

<code>-h,--help</code>	Show this help text
<code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (could be a Poseidon repository)
<code>--remote</code>	list packages from a remote server instead the local file system
<code>--remoteURL ARG</code>	URL of the remote Poseidon server (default: "https://c107-224.cloud.gwdg.de")
<code>--packages</code>	list all packages

```

500  --groups                list all groups, ignoring any group names after the
501                          first as specified in the Janno-file
502  --individuals           list individuals
503  -j,--jannoColumn JANNO_HEADER
504                          list additional fields from the janno files, using
505                          the Janno column heading name, such as Country, Site,
506                          Date_C14_Uncal_BP, Endogenous, ...
507  --raw                   output table as tsv without header. Useful for piping
508                          into grep or awk
509  --ignoreGeno            ignore SNP and GenoFile

```

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

```
511 trident list -d ... -d ... --packages
```

512 This will yield a table like this

```

513 .------.------.------.
514 |                Title                |    Date    | Nr Individuals |
515 :=====:=====:=====:
516 | 2015_1000Genomes_1240K_haploid_pull | 2020-08-10 | 2535           |
517 | 2016_Mallick_SGDP1240K_diploid_pull | 2020-08-10 | 280            |
518 | 2018_BostonDatashare_modern_published | 2020-08-10 | 2772           |
519 | ...                                | ...        |                |
520 '-----'-----'-----'

```

521 so a nicely formatted table of all packages, their last update and the number of individuals in it.

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

```
524 trident list --packages --remote
```

525 will result in a view of all published packages in our public online repository.

526 You can also list groups, as defined in the third column of EIGENSTRAT `.ind` files (or the first column of a  
527 PLINK `.fam` file), and individuals:

```

528 trident list -d ... -d ... --groups
529 trident list -d ... -d ... --individuals

```

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

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

## 537 4.2 Summarise command

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

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

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

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

542

543 Available options:

544 -h,--help Show this help text

545 -d,--baseDir DIR a base directory to search for Poseidon Packages  
546 (could be a Poseidon repository)

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

549 You can run it with

550 trident summarise -d ... -d ...

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

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

## 556 4.3 Survey command

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

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

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

560 Survey the degree of context information completeness for Poseidon packages

561

562 Available options:

563 -h,--help Show this help text

564 -d,--baseDir DIR a base directory to search for Poseidon Packages  
565 (could be a Poseidon repository)

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

568 Running

569 trident survey -d ... -d ...

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

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

## 573 4.4 Validate command

574 `validate` checks poseidon datasets for structural correctness.

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

```

576 Usage: trident validate (-d|--baseDir DIR) [--verbose]
577     Check one or multiple Poseidon packages for structural correctness
578
579 Available options:
580     -h,--help                Show this help text
581     -d,--baseDir DIR         a base directory to search for Poseidon Packages
582                             (could be a Poseidon repository)
583     --ignoreGeno             ignore SNP and GenoFile
584     --noExitCode             do not produce an explicit exit code
585
586 You can run it with
587
588 trident validate -d ... -d ...
589
590 and it will either report a success ( Validation passed ) or failure with specific error messages to simplify
591 fixing the issues.
592
593 validate tries to ensure that each package in the dataset adheres to the schema definition. Here is a list of
594 what is checked:
595
596 • Presence of the necessary files
597 • Full structural correctness of .bib and .janno file
598 • Superficial correctness of genotype data files. A full check would be too computationally expensive
599 • Correspondence of BibTeX keys in .bib and .janno
600 • Correspondence of individual and group IDs in .janno and genotype data files
601
602 In fact much of this validation already runs as part of the general package reading pipeline invoked for many
603 trident subcommands (e.g. forge). validate is meant to be more thorough, though, and will explicitly fail if
604 even a single package is broken.
605
606 Remember to run it with --logMode VerboseLog to get more information if the output is not sufficient to
607 debug an issue.

```