

Guide for trident v1.1.0.0 to v1.1.4.2

Contents

1	Poseidon package repositories	1
2	Analysing your own dataset outside of the main repository	2
3	Package creation and manipulation commands	3
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
3.4	Genoconvert command	9
3.5	Update command	11
4	Inspection commands	12
4.1	List command	12
4.2	Summarise command	14
4.3	Survey command	14
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
  /ancient
```

```

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

```

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

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

```

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

```

46 2 Analysing your own dataset outside of the main repository

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

```

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

```

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

```

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58 description: Unpublished genetic data from my awesome project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

```

71 Two remarks: 1) all file paths are considered *relative* to the directory in which `POSEIDON.yml` resides. Here I
 72 assume that you put this file into the same directory as the three genotype files. 2) Besides the genotype data

73 files there are two (technically optional) files referenced by this example `POSEIDON.yml` file: `sources.bib`
74 and `my_project.janno`. Of course you can add them manually - `init` automatically creates empty dummy
75 versions.

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

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

81 3 Package creation and manipulation commands

82 3.1 Init command

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

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

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

90 Create a new Poseidon package from genotype data

92 Available options:

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

112 The command

```

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

```

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

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

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

128 3.2 Fetch command

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

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

```

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

```

135 Download data from a remote Poseidon repository

136 Available options:

138 <code>-h,--help</code>	Show this help text
139 <code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (could be a Poseidon repository)
140 <code>--downloadAll</code>	download all packages the server is offering
141 <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.
142 <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> ,

```

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

```

157 It works with

```

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

```

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

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

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

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

176 3.3 Forge command

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

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

```

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

```

190 Available options:

191 -h,--help Show this help text

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

194 -p,--genoOne ARG one of the input genotype data files. Expects .bed or
195 .bim or .fam for PLINK and .geno or .snp or .ind for
196 EIGENSTRAT. The other files must be in the same
197 directory and must have the same base name

198 --inFormat ARG the format of the input genotype data: EIGENSTRAT or
199 PLINK (only necessary for data input with --genoFile
200 + --snpFile + --indFile)

201 --genoFile ARG the input geno file path

202 --snpFile ARG the input snp file path

203 --indFile ARG the input ind file path

204 --snpSet ARG the snpSet of the new package: 1240K, HumanOrigins or
205 Other. Default: Other

206 --forgeFile ARG A file with a list of packages, groups or individual
207 samples. Works just as -f, but multiple values can
208 also be separated by newline, not just by comma.
209 Empty lines are ignored and comments start with "#",
210 so everything after "#" is ignored in one line.
211 Multiple instances of -f and --forgeFile can be
212 given. They will be evaluated according to their
213 input order on the command line.

214 -f,--forgeString ARG List of packages, groups or individual samples to be
215 combined in the output package. Packages follow the
216 syntax *package_title*, populations/groups are simply
217 group_id and individuals <individual_id>. You can
218 combine multiple values with comma, so for example:
219 "*package_1*, <individual_1>, <individual_2>,"
220 group_1". Duplicates are treated as one entry.
221 Negative selection is possible by prepending "-" to
222 the entity you want to exclude (e.g. "*package_1*,"
223 -<individual_1>, -group_1"). forge will apply
224 excludes and includes in order. If the first entity
225 is negative, then forge will assume you want to merge
226 all individuals in the packages found in the baseDirs
227 (except the ones explicitly excluded) before the
228 exclude entities are applied. An empty forgeString
229 (and no --forgeFile) will therefore merge all
230 available individuals.

231 --selectSnps ARG To extract specific SNPs during this forge operation,
232 provide a Snp file. Can be either Eigenstrat (file
233 ending must be '.snp') or Plink (file ending must be
234 '.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

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

322 3.3.2 Other options

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

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

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

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

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

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

342 3.4 Genoconvert command

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

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

```
346 Usage: trident genoconvert ((-d|--baseDir DIR) |
347                             ((-p|--genoOne ARG) | --inFormat ARG --genoFile ARG
348                             --snpFile ARG --indFile ARG) [--snpSet ARG])
349                             --outFormat ARG [--onlyGeno]
350                             [-o|--outPackagePath ARG] [--removeOld]
```

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

352

353 Available options:

354	<code>-h,--help</code>	Show this help text
355	<code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (could be a Poseidon repository)
356		
357	<code>-p,--genoOne ARG</code>	one of the input genotype data files. Expects .bed or .bim or .fam for PLINK and .geno or .snp or .ind for EIGENSTRAT. The other files must be in the same directory and must have the same base name
358		
359		
360		
361	<code>--inFormat ARG</code>	the format of the input genotype data: EIGENSTRAT or PLINK (only necessary for data input with --genoFile + --snpFile + --indFile)
362		
363		
364	<code>--genoFile ARG</code>	the input geno file path
365	<code>--snpFile ARG</code>	the input snp file path
366	<code>--indFile ARG</code>	the input ind file path
367	<code>--snpSet ARG</code>	the snpSet of the new package: 1240K, HumanOrigins or Other. Default: Other
368		
369	<code>--outFormat ARG</code>	the format of the output genotype data: EIGENSTRAT or PLINK.
370		
371	<code>--onlyGeno</code>	should only the resulting genotype data be returned? This means the output will not be a Poseidon package
372		
373	<code>-o,--outPackagePath ARG</code>	the output package directory path - this is optional: If no path is provided, then the output is written to the directories where the input genotype data file (.bed/.geno) is stored
374		
375		
376		
377	<code>--removeOld</code>	Remove the old genotype files when creating the new ones
378		

379 With the default setting

380 `trident genoconvert -d ... -d ... --outFormat EIGENSTRAT|PLINK`

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

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

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

```
389 trident genoconvert \  
390   -p 2018_Mittnik_Baltic/Mittnik_Baltic.bed \  
391   --outFormat EIGENSTRAT  
392   -o my_directory
```

3.5 Update command

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

Click here for command line details

```
Usage: trident update (-d|--baseDir DIR) [--poseidonVersion ARG]
        [--ignorePoseidonVersion] [--versionComponent ARG]
        [--noChecksumUpdate] [--newContributors ARG]
        [--logText ARG] [--force]

Update POSEIDON.yml files automatically

Available options:
  -h,--help                Show this help text
  -d,--baseDir DIR          a base directory to search for Poseidon Packages
                           (could be a Poseidon repository)
  --poseidonVersion ARG     Poseidon version the packages should be updated to:
                           e.g. "2.5.3" (default: Nothing)
  --ignorePoseidonVersion   Read packages even if their poseidonVersion is not
                           compatible with the trident version. The assumption
                           is, that the package is already structurally adjusted
                           to the trident version and only the version number is
                           lagging behind.
  --versionComponent ARG    Part of the package version number in the
                           POSEIDON.yml file that should be updated: Major,
                           Minor or Patch (see https://semver.org)
                           (default: Patch)
  --noChecksumUpdate        Should update of checksums in the POSEIDON.yml file
                           be skipped
  --ignoreGeno              ignore SNP and GenoFile
  --newContributors ARG     Contributors to add to the POSEIDON.yml file in the
                           form "[Firstname Lastname](Email address);..."
  --logText ARG             Log text for this version jump in the CHANGELOG file
                           (default: "not specified")
  --force                   Normally the POSEIDON.yml files are only changed if
                           the poseidonVersion is adjusted or any of the
                           checksums change. With --force a package version
                           update can be triggered even if this is not the case.
```

It can be called with a lot of optional arguments

```
trident update -d ... -d ... \
  --poseidonVersion "X.X.X" \
  --versionComponent Major/Minor/Patch \
  --noChecksumUpdate
  --ignoreGeno
  --newContributors "[Firstname Lastname](Email address);..."
```

```
436 --logText "short description of the update"
437 --force
```

438 By default `update` will not edit a package's POSEIDON.yml file, even when arguments like `--versionComponent`,
439 `--newContributors` or `--logText` are explicitly set. This default exists to run the function on a large set of
440 packages where only few of them were edited and need an active update. A package will only be modified by
441 `update` if either

- 442 • any of the files with checksums (e.g. the genotype data) in it were modified,
- 443 • the `--poseidonVersion` argument differs from the `poseidonVersion` in the package's POSEIDON.yml
444 file
- 445 • or the `--force` flag was set in `update`.

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

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

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

465 4 Inspection commands

466 4.1 List command

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

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

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

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

474
475 Available options:

```

476 -h,--help          Show this help text
477 -d,--baseDir DIR   a base directory to search for Poseidon Packages
478                   (could be a Poseidon repository)
479 --remote           list packages from a remote server instead the local
480                   file system
481 --remoteURL ARG     URL of the remote Poseidon server
482                   (default: "https://c107-224.cloud.gwdg.de")
483 --packages         list all packages
484 --groups           list all groups, ignoring any group names after the
485                   first as specified in the Janno-file
486 --individuals      list individuals
487 -j,--jannoColumn JANNO_HEADER
488                   list additional fields from the janno files, using
489                   the Janno column heading name, such as Country, Site,
490                   Date_C14_Uncal_BP, Endogenous, ...
491 --raw             output table as tsv without header. Useful for piping
492                   into grep or awk
493 --ignoreGeno       ignore SNP and GenoFile

```

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

```
495 trident list -d ... -d ... --packages
```

496 This will yield a table like this

```

497 .------.------.------.
498 |           Title           |    Date    | Nr Individuals |
499 :=====:=====:=====:
500 | 2015_1000Genomes_1240K_haploid_pulldown | 2020-08-10 | 2535          |
501 | 2016_Mallick_SGDP1240K_diploid_pulldown | 2020-08-10 | 280           |
502 | 2018_BostonDatashare_modern_published   | 2020-08-10 | 2772          |
503 | ...                                     | ...         |                |
504 '-----'-----'-----'

```

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

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

```
508 trident list --packages --remote
```

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

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

```
512 trident list -d ... -d ... --groups
```

```
513 trident list -d ... -d ... --individuals
```

514 The `--individuals` flag also provides a way to immediately access information from the `.janno`
515 files on the command line. This works with the `-j / --jannoColumn` option. For example adding

516 `--jannoColumn Country --jannoColumn Date_C14_Uncal_BP` to the commands above will add the `Country`
517 and the `Date_C14_Uncal_BP` columns to the respective output tables.

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

521 4.2 Summarise command

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

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

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

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

526
527 Available options:
528 `-h,--help` Show this help text
529 `-d,--baseDir DIR` a base directory to search for Poseidon Packages
530 (could be a Poseidon repository)
531 `--raw` output table as tsv without header. Useful for piping
532 into `grep` or `awk`

533 You can run it with

534 `trident summarise -d ... -d ...`

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

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

540 4.3 Survey command

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

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

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

544 Survey the degree of context information completeness for Poseidon packages

545
546 Available options:
547 `-h,--help` Show this help text
548 `-d,--baseDir DIR` a base directory to search for Poseidon Packages
549 (could be a Poseidon repository)
550 `--raw` output table as tsv without header. Useful for piping
551 into `grep` or `awk`

552 Running

553 `trident survey -d ... -d ...`

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

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

4.4 Validate command

`validate` checks poseidon datasets for structural correctness.

[Click here for command line details](#)

Usage: `trident validate (-d|--baseDir DIR) [--verbose]`

Check one or multiple Poseidon packages for structural correctness

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>--verbose</code>	print more output to the command line
<code>--ignoreGeno</code>	ignore SNP and GenoFile
<code>--noExitCode</code>	do not produce an explicit exit code

You can run it with

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

and it will either report a success (`Validation passed`) or failure with specific error messages to simplify fixing the issues.

`validate` tries to ensure that each package in the dataset adheres to the [schema definition](#). Here is a list of what is checked:

- Presence of the necessary files
- Full structural correctness of `.bib` and `.janno` file
- Superficial correctness of genotype data files. A full check would be too computationally expensive
- Correspondence of BibTeX keys in `.bib` and `.janno`
- Correspondence of individual and group IDs in `.janno` and genotype data files

In fact much of this validation already runs as part of the general package reading pipeline invoked for many trident subcommands (e.g. `forge`). `validate` is meant to be more thorough, though, and will explicitly fail if even a single package is broken.