

Guide for trident v1.1.6.0

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
3.3.3	Treatment of the .janno file while merging	9
3.4	Genoconvert command	10
3.5	Update command	11
4	Inspection commands	13
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 “unpackaged” 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 by
55 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` and
 74 `my_project.janno`. Of course you can add them manually - `init` automatically creates empty dummy versions.
 75 Once you have set up your own “Poseidon” package (which is really only a skeleton so far), you can add it to
 76 your `trident` analysis, by simply adding your project directory to the command using `-d`:

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

80 3 Package creation and manipulation commands

81 3.1 Init command

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

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

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

89 Create a new Poseidon package from genotype data

91 Available options:

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

111 The command

```
112 trident init \
```

```

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

```

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

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

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

126 3.2 Fetch command

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

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

```

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

```

133 Download data from a remote Poseidon repository

134 Available options:

136 <code>-h,--help</code>	Show this help text
137 <code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (could be a Poseidon repository)
138 <code>--downloadAll</code>	download all packages the server is offering
139 <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.
140 <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*, *package_3*"</code> . <code>fetchString</code> uses the same parser as <code>forgeString</code> , but does not allow excludes. If groups

```

149         or individuals are specified, then packages which
150         include these groups or individuals are included in
151         the download.
152     --remoteURL ARG        URL of the remote Poseidon server
153                           (default: "https://c107-224.cloud.gwdg.de")
154     -u,--upgrade           overwrite outdated local package versions

```

155 It works with

```

156 trident fetch -d ... -d ... \
157     -f "*package_title_1*,*package_title_2*,*package_title_3*,group_name,<Individual1>" \
158     --fetchFile path/to/forgeFile

```

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

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

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

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

174 3.3 Forge command

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

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

```

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

```

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

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

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

199 --genoFile ARG the input geno file path

200 --snpFile ARG the input snp file path

201 --indFile ARG the input ind file path

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

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

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

229 --selectSnps ARG To extract specific SNPs during this forge operation,
230 provide a Snp file. Can be either Eigenstrat (file
231 ending must be '.snp') or Plink (file ending must be
232 '.bim'). When this option is set, the output package
233 will have exactly the SNPs listed in this file. Any
234 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

279 merge data from one Poseidon package and two genotype datasets to get a new EIGENSTRAT dataset.

```
280 trident forge \  
281   -d 2017_GonzalesFortesCurrentBiology \  
282   -p 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \  
283   --inFormat PLINK \  
284   --genoFile 2017_HaberAJHG/2017_HaberAJHG.bed \  
285   --snpFile 2017_HaberAJHG/2017_HaberAJHG.bim \  
286   --indFile 2017_HaberAJHG/2017_HaberAJHG.fam \  
287   -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \  
288   -o testpackage \  
289   --outFormat EIGENSTRAT \  
290   --onlyGeno
```

291 3.3.1 The forge selection language

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

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

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

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

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

```
303 # Packages  
304 *package1*, *package2*  
305  
306 # Groups and individuals from other packages beyond package1 and package2  
307 group1, <individual1>, group2, <individual2>, <individual3>  
308  
309 # group2 has two outlier individuals that should be ignored  
310 -<bad_individual1> # This one has very low coverage  
311 -<bad_individual2> # This one is from a different time period
```

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

3.3.2 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 might be especially useful for data analysis pipelines, where only the genotype data is required. Even more basic output comes with `--onlyGeno`, which means that only the genotype data is returned without any Poseidon package.

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

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

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

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

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

3.3.3 Treatment of the `.janno` file while merging

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

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

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

The following example illustrates the described behaviour:

A.janno

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

B.janno

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

A.janno + B.janno

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

3.4 Genoconvert command

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

[Click here for command line details](#)

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

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

Available options:

-h,--help	Show this help text
-d,--baseDir DIR	a base directory to search for Poseidon Packages (could be a Poseidon repository)
-p,--genoOne ARG	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

```

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

```

391 With the default setting

```

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

```

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

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

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

```

401 trident genoconvert \
402   -p 2018_Mittnik_Baltic/Mittnik_Baltic.bed \
403   --outFormat EIGENSTRAT
404   -o my_directory

```

405 3.5 Update command

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

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

```

409 Usage: trident update (-d|--baseDir DIR) [--poseidonVersion ARG]
410                          [--ignorePoseidonVersion] [--versionComponent ARG]
411                          [--noChecksumUpdate] [--newContributors ARG]
412                          [--logText ARG] [--force]
413 Update POSEIDON.yml files automatically

```

```

414
415 Available options:
416 -h,--help          Show this help text
417 -d,--baseDir DIR    a base directory to search for Poseidon Packages
418                    (could be a Poseidon repository)
419 --poseidonVersion ARG Poseidon version the packages should be updated to:
420                    e.g. "2.5.3" (default: Nothing)
421 --ignorePoseidonVersion Read packages even if their poseidonVersion is not
422                    compatible with the trident version. The assumption
423                    is, that the package is already structurally adjusted
424                    to the trident version and only the version number is
425                    lagging behind.
426 --versionComponent ARG Part of the package version number in the
427                    POSEIDON.yml file that should be updated: Major,
428                    Minor or Patch (see https://semver.org)
429                    (default: Patch)
430 --noChecksumUpdate    Should update of checksums in the POSEIDON.yml file
431                    be skipped
432 --ignoreGeno          ignore SNP and GenoFile
433 --newContributors ARG Contributors to add to the POSEIDON.yml file in the
434                    form "[Firstname Lastname](Email address);..."
435 --logText ARG         Log text for this version jump in the CHANGELOG file
436                    (default: "not specified")
437 --force              Normally the POSEIDON.yml files are only changed if
438                    the poseidonVersion is adjusted or any of the
439                    checksums change. With --force a package version
440                    update can be triggered even if this is not the case.

```

441 It can be called with a lot of optional arguments

```

442 trident update -d ... -d ... \
443 --poseidonVersion "X.X.X" \
444 --versionComponent Major/Minor/Patch \
445 --noChecksumUpdate
446 --ignoreGeno
447 --newContributors "[Firstname Lastname](Email address);..."
448 --logText "short description of the update"
449 --force

```

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

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

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

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

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

476 4 Inspection commands

477 4.1 List command

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

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

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

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

485 Available options:

487 -h,--help	Show this help text
488 -d,--baseDir DIR	a base directory to search for Poseidon Packages 489 (could be a Poseidon repository)
490 --remote	list packages from a remote server instead the local 491 file system
492 --remoteURL ARG	URL of the remote Poseidon server 493 (default: "https://c107-224.cloud.gwdg.de")
494 --packages	list all packages
495 --groups	list all groups, ignoring any group names after the 496 first as specified in the Janno-file
497 --individuals	list individuals
498 -j,--jannoColumn JANNO_HEADER	

```

499         list additional fields from the janno files, using
500         the Janno column heading name, such as Country, Site,
501         Date_C14_Uncal_BP, Endogenous, ...
502     --raw                output table as tsv without header. Useful for piping
503                          into grep or awk
504     --ignoreGeno         ignore SNP and GenoFile

```

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

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

This will yield a table like this

```

508 .------.------.------.
509 |                Title                |    Date    | Nr Individuals |
510 :=====:=====:=====:
511 | 2015_1000Genomes_1240K_haploid_pull | 2020-08-10 | 2535           |
512 | 2016_Mallick_SGDP1240K_diploid_pull | 2020-08-10 | 280             |
513 | 2018_BostonDatashare_modern_published | 2020-08-10 | 2772           |
514 | ...                                | ...        |                 |
515 '-----'-----'-----'

```

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

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

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

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

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

```

523 trident list -d ... -d ... --groups
524 trident list -d ... -d ... --individuals

```

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

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

4.2 Summarise command

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

[Click here for command line details](#)

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

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

537

538 Available options:

539 -h,--help Show this help text

540 -d,--baseDir DIR a base directory to search for Poseidon Packages

541 (could be a Poseidon repository)

542 --raw output table as tsv without header. Useful for piping

543 into grep or awk

544 You can run it with

545 `trident summarise -d ... -d ...`

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

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

551 4.3 Survey command

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

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

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

555 Survey the degree of context information completeness for Poseidon packages

556

557 Available options:

558 -h,--help Show this help text

559 -d,--baseDir DIR a base directory to search for Poseidon Packages

560 (could be a Poseidon repository)

561 --raw output table as tsv without header. Useful for piping

562 into grep or awk

563 Running

564 `trident survey -d ... -d ...`

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

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

568 4.4 Validate command

569 `validate` checks poseidon datasets for structural correctness.

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

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

572 Check one or multiple Poseidon packages for structural correctness

573

574 Available options:

575	<code>-h,--help</code>	Show this help text
576	<code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages
577		(could be a Poseidon repository)
578	<code>--ignoreGeno</code>	ignore SNP and GenoFile
579	<code>--noExitCode</code>	do not produce an explicit exit code

580 You can run it with

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

582 and it will either report a success (**Validation passed**) or failure with specific error messages to simplify fixing

583 the issues.

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

585 what is checked:

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

591 In fact much of this validation already runs as part of the general package reading pipeline invoked for many

592 trident subcommands (e.g. **forge**). **validate** is meant to be more thorough, though, and will explicitly fail if

593 even a single package is broken.

594 Remember to run it with `--logMode VerboseLog` to get more information if the output is not sufficient to debug

595 an issue.