

Guide for trident v0.28.0

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	7
3.3.2	Other options	8
3.4	Genoconvert command	9
3.5	Update command	10
4	Inspection commands	12
4.1	List command	12
4.2	Summarise command	13
4.3	Survey command	14
4.4	Validate command	14
5	Analysis commands	15

1 Poseidon package repositories

Trident generally requires Poseidon datasets to work with. Most trident subcommands therefore have a central parameter, called `--baseDir` or simply `-d` to specify one or more base directories to look for Poseidon 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 (-r|--inFormat ARG) (-g|--genoFile ARG) (-s|--snpFile ARG)  
87 (-i|--indFile ARG) [--snpSet ARG] (-o|--outPackagePath ARG)  
88 [-n|--outPackageName ARG] [--minimal]
```

89 Create a new Poseidon package from genotype data

90 Available options:

92	<code>-h,--help</code>	Show this help text
93	<code>-r,--inFormat ARG</code>	the format of the input genotype data: EIGENSTRAT or
94		PLINK
95	<code>-g,--genoFile ARG</code>	the input geno file path
96	<code>-s,--snpFile ARG</code>	the input snp file path
97	<code>-i,--indFile ARG</code>	the input ind file path
98	<code>--snpSet ARG</code>	the snpSet of the new package: 1240K, HumanOrigins or
99		Other. Default: Other
100	<code>-o,--outPackagePath ARG</code>	the output package directory path
101	<code>-n,--outPackageName ARG</code>	the output package name - this is optional: If no
102		name is provided, then the package name defaults to
103		the basename of the (mandatory) <code>--outPackagePath</code>
104		argument
105	<code>--minimal</code>	should only a minimal output package be created?

106 The command

```
107 trident init \  
108 -r EIGENSTRAT/PLINK \  
109 -g path/to/geno_file \  
110 -s path/to/snp_file \  
111 -i path/to/ind_file \  
112 --snpSet 1240K|HumanOrigins|Other \  
113
```

113 `-o path/to/new_package_name`

114 requires the format `-r` (`--inFormat`) of your input data (either `EIGENSTRAT` or `PLINK`), the paths to the
115 respective files in `-g` (`--genoFile`), `-s` (`--snpFile`), and `-i` (`--indFile`), and optionally the “shape”
116 of these files (`--snpSet`), so if they cover the `1240K`, the `HumanOrigins` or an `Other` SNP set.

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

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

121 3.2 Fetch command

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

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

124 Usage: trident fetch (-d|--baseDir DIR) [-f|--fetchString ARG] [--fetchFile ARG]
125 [--remoteURL ARG] [-u|--upgrade] [--downloadAll]

126 Download data from a remote Poseidon repository

127
128 Available options:

129 <code>-h,--help</code>	Show this help text
130 <code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (could be a Poseidon repository)
131	
132 <code>-f,--fetchString ARG</code>	List of packages to be downloaded from the remote 133 server. Package names should be wrapped in asterisks: 134 <code>*package_title*</code> . You can combine multiple values with 135 comma, so for example: <code>"*package_1*, *package_2*,</code> 136 <code>*package_3*"</code> . <code>fetchString</code> uses the same parser as 137 <code>forgeString</code> , but does not allow excludes. If groups 138 or individuals are specified, then packages which 139 include these groups or individuals are included in 140 the download.
141 <code>--fetchFile ARG</code>	A file with a list of packages. Works just as <code>-f</code> , but 142 multiple values can also be separated by newline, not 143 just by comma. <code>-f</code> and <code>--fetchFile</code> can be combined.
144 <code>--remoteURL ARG</code>	URL of the remote Poseidon 145 server (default: <code>"https://c107-224.cloud.gwdg.de"</code>)
146 <code>-u,--upgrade</code>	overwrite outdated local package versions
147 <code>--downloadAll</code>	download all packages the server is offering

148 It works with

```

149 trident fetch -d ... -d ... \
150   -f "*package_title_1*,*package_title_2*,*package_title_3*,group_name,<Individual1>" \
151   --fetchFile path/to/forgeFile

```

and the entities you want to download must be listed either in a simple string with comma-separated values (`-f / --fetchString`) or in a text file (`--fetchFile`). 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` causes fetch to ignore `-f` and download all packages from the server. The downloaded packages are added in the first (!) `-d` directory, 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.

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.

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

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.

3.3 Forge command

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

Click here for command line details

```

170 Usage: trident forge [-d|--baseDir DIR]
171                   [(-r|--inFormat ARG) (-g|--genoFile ARG) (-s|--snpFile ARG)
172                   (-i|--indFile ARG) [--snpSet ARG]]
173                   [--forgeFile ARG | (-f|--forgeString ARG)]
174                   [--selectSnps ARG] [--intersect] [--outFormat ARG]
175                   [--minimal] [--onlyGeno] (-o|--outPackagePath ARG)
176                   [-n|--outPackageName ARG] [-w|--warnings] [--no-extract]
177   Select packages, groups or individuals and create a new Poseidon package from
178   them

```

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>-r,--inFormat ARG</code>	the format of the input genotype data: EIGENSTRAT or PLINK
<code>-g,--genoFile ARG</code>	the input geno file path
<code>-s,--snpFile ARG</code>	the input snp file path
<code>-i,--indFile ARG</code>	the input ind file path
<code>--snpSet ARG</code>	the snpSet of the new package: 1240K, HumanOrigins or

190		Other. Default: Other
191	--forgeFile ARG	A file with a list of packages, groups or individual
192		samples. Works just as -f, but multiple values can
193		also be separated by newline, not just by comma.
194		Empty lines are ignored and comments start with "#",
195		so everything after "#" is ignored in one line.
196	-f,--forgeString ARG	List of packages, groups or individual samples to be
197		combined in the output package. Packages follow the
198		syntax *package_title*, populations/groups are simply
199		group_id and individuals <individual_id>. You can
200		combine multiple values with comma, so for example:
201		"*package_1*, <individual_1>, <individual_2>,"
202		group_1". Duplicates are treated as one entry.
203		Negative selection is possible by prepending "-" to
204		the entity you want to exclude (e.g. "*package_1*,
205		-<individual_1>, -group_1"). forge will apply
206		excludes and includes in order. If the first entity
207		is negative, then forge will assume you want to merge
208		all individuals in the packages found in the baseDirs
209		(except the ones explicitly excluded) before the
210		exclude entities are applied. An empty forgeString
211		will therefore merge all available individuals.
212	--selectSnps ARG	To extract specific SNPs during this forge operation,
213		provide a Snp file. Can be either Eigenstrat (file
214		ending must be '.snp') or Plink (file ending must be
215		'.bim'). When this option is set, the output package
216		will have exactly the SNPs listed in this file. Any
217		SNP not listed in the file will be excluded. If
218		option '--intersect' is also set, only the SNPs
219		overlapping between the SNP file and the forged
220		packages are output.
221	--intersect	Whether to output the intersection of the genotype
222		files to be forged. The default (if this option is
223		not set) is to output the union of all SNPs, with
224		genotypes defined as missing in those packages which
225		do not have a SNP that is present in another package.
226		With this option set, the forged dataset will
227		typically have fewer SNPs, but less missingness.
228	--outFormat ARG	the format of the output genotype data: EIGENSTRAT or
229		PLINK. Default: PLINK
230	--minimal	should only a minimal output package be created?
231	--onlyGeno	should only the resulting genotype data be returned?
232		This means the output will not be a Poseidon package
233	-o,--outPackagePath ARG	the output package directory path
234	-n,--outPackageName ARG	the output package name - this is optional: If no

235 name is provided, then the package name defaults to
 236 the basename of the (mandatory) `--outPackagePath`
 237 argument
 238 `-w,--warnings` Show all warnings for merging genotype data
 239 `--no-extract` Skip the selection step in forge. This will result in
 240 outputting all individuals in the relevant packages,
 241 and hence a superset of the requested
 242 individuals/groups. It may result in better
 243 performance in cases where one wants to forge entire
 244 packages or almost entire packages. Note that this
 245 will also ignore any ordering in the output
 246 groups/individuals. With this option active,
 247 individuals from the relevant packages will just be
 248 written in the order that they appear in the original
 249 packages.

250 `forge` can be used with

```

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

255 where the entities (packages, groups/populations, individuals/samples) you want in the output package can
 256 be denoted either as as simple string with comma-separated values (`-f / --forgeString`) or in a text file
 257 (`--forgeFile`).

258 Including one or multiple Poseidon packages with `-d` is not the only way to include data for a forge operation.
 259 It is also possible to include unpackaged genotype data directly with `-r + -g + -s + -i (+ --snpSet)` . This
 260 makes the following example possible, where we merge data from one Poseidon package and two genotype
 261 datasets.

```

262 trident forge \
263   -d 2017_GonzalesFortesCurrentBiology \
264   -r PLINK \
265   -g 2017_HaberAJHG/2017_HaberAJHG.bed \
266   -s 2017_HaberAJHG/2017_HaberAJHG.bim \
267   -i 2017_HaberAJHG/2017_HaberAJHG.fam \
268   -r PLINK \
269   -g 2018_VeeramahPNAS/2018_VeeramahPNAS.bed \
270   -i 2018_VeeramahPNAS/2018_VeeramahPNAS.fam \
271   -s 2018_VeeramahPNAS/2018_VeeramahPNAS.bim \
272   -f "<STR241.SG>,<ERS1790729.SG>,Iberia_HG.SG" \
273   -o testpackage \
274   --onlyGeno
  
```

275 3.3.1 The forge selection language

276 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 either use `-f / --forgeString` or `--forgeFile`. 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 packages found in the `baseDirs` (except the ones explicitly excluded, of course). An empty `forgeString` will 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 an optional flag `--intersect`, that defines, if the genotype data from different packages should be merged with an **union** or an **intersect** operation. The default (if this option is not set) is to output the union of all SNPs, with genotypes defined as missing in samples from packages which do not have a SNP that is present in another package. With this option set, on the other hand, the forged dataset will typically have fewer SNPs, but less missingness.

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

Input snpSet A	Input snpSet B	<code>--intersect</code>	Output snpSet
Other	*	*	Other

Input snpSet A	Input snpSet B	<code>--intersect</code>	Ouput snpSet
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 `-w / --warnings` flag is set.

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]
                        [(-r|--inFormat ARG) (-g|--genoFile ARG)
                        (-s|--snpFile ARG) (-i|--indFile ARG)
                        [--snpSet ARG]] --outFormat ARG [--onlyGeno]
                        [--removeOld]
  
```

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

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>-r,--inFormat ARG</code>	the format of the input genotype data: EIGENSTRAT or PLINK
<code>-g,--genoFile ARG</code>	the input geno file path
<code>-s,--snpFile ARG</code>	the input snp file path
<code>-i,--indFile ARG</code>	the input ind file path
<code>--snpSet ARG</code>	the snpSet of the new package: 1240K, HumanOrigins or Other. Default: Other
<code>--outFormat ARG</code>	the format of the output genotype data: EIGENSTRAT or PLINK.
<code>--onlyGeno</code>	should only the resulting genotype data be returned? This means the output will not be a Poseidon package
<code>--removeOld</code>	Remove the old genotype files when creating the new ones

With the default setting

```

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

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

Instead of `-d` to change Poseidon packages, the combination `-r + -g + -s + -i (+ --snpSet)` allows to directly convert genotype data that is not wrapped in a Poseidon package. See this example:

```
trident genoconvert \  
-r PLINK \  
-g 2018_Mittnik_Baltic/Mittnik_Baltic.bed \  
-s 2018_Mittnik_Baltic/Mittnik_Baltic.bim \  
-i 2018_Mittnik_Baltic/Mittnik_Baltic.fam \  
--outFormat EIGENSTRAT
```

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

Remember that the `POSEIDON.yml` file can also be edited by hand if you want to replace the genotype data in a package.

3.5 Update command

`update` automatically updates `POSEIDON.yml` files of one or multiple packages if the packages were changed.

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:

<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>--poseidonVersion ARG</code>	Poseidon version the packages should be updated to: e.g. "2.5.3" (default: Nothing)
<code>--ignorePoseidonVersion</code>	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.
<code>--versionComponent ARG</code>	Part of the package version number in the <code>POSEIDON.yml</code> file that should be updated: Major, Minor or Patch (see https://semver.org) (default: Patch)
<code>--noChecksumUpdate</code>	Should update of checksums in the <code>POSEIDON.yml</code> file be skipped

```

392 --ignoreGeno          ignore SNP and GenoFile
393 --newContributors ARG Contributors to add to the POSEIDON.yml file in the
394                        form "[Firstname Lastname](Email address);..."
395 --logText ARG         Log text for this version jump in the CHANGELOG
396                        file (default: "not specified")
397 --force               Normally the POSEIDON.yml files are only changed if
398                        the poseidonVersion is adjusted or any of the
399                        checksums change. With --force a package version
400                        update can be triggered even if this is not the case.

```

401 It can be called with a lot of optional arguments

```

402 trident update -d ... -d ... \
403     --poseidonVersion "X.X.X" \
404     --versionComponent Major/Minor/Patch \
405     --noChecksumUpdate
406 --ignoreGeno
407 --newContributors "[Firstname Lastname](Email address);..."
408 --logText "short description of the update"
409 --force

```

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

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

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

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

434 :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:

```
-h,--help          Show this help text
-d,--baseDir DIR   a base directory to search for Poseidon Packages
                   (could be a Poseidon repository)
--remote           list packages from a remote server instead the local
                   file system
--remoteURL ARG    URL of the remote Poseidon
                   server (default: "https://c107-224.cloud.gwdg.de")
--packages         list all packages
--groups           list all groups, ignoring any group names after the
                   first as specified in the Janno-file
--individuals      list individuals
-j,--jannoColumn JANNO_HEADER
                   list additional fields from the janno files, using
                   the Janno column heading name, such as Country, Site,
                   Date_C14_Uncal_BP, Endogenous, ...
--raw             output table as tsv without header. Useful for piping
                   into grep or awk
--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

```
.------.------.------.
|          Title          |    Date    | Nr Individuals |
:=====:=====:=====:
| 2015_1000Genomes_1240K_haploid_pulldown | 2020-08-10 | 2535          |
| 2016_Mallick_SGDP1240K_diploid_pulldown | 2020-08-10 | 280           |
| 2018_BostonDatashare_modern_published   | 2020-08-10 | 2772          |
```

```

475 | ... | ... |
476 '-----'-----'-----'

```

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

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

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

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

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

```
484 trident list -d ... -d ... --groups
```

```
485 trident list -d ... -d ... --individuals
```

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

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

493 4.2 Summarise command

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

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

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

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

```

498
499 Available options:
500   -h,--help           Show this help text
501   -d,--baseDir DIR    a base directory to search for Poseidon Packages
502                       (could be a Poseidon repository)
503   --raw               output table as tsv without header. Useful for piping
504                       into grep or awk

```

505 You can run it with

```
506 trident summarise -d ... -d ...
```

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

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

512 4.3 Survey command

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

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

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

516 Survey the degree of context information completeness for Poseidon packages

517 Available options:

519 <code>-h,--help</code>	Show this help text
520 <code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (could be a Poseidon repository)
522 <code>--raw</code>	output table as tsv without header. Useful for piping 523 into <code>grep</code> or <code>awk</code>

524 Running

525 `trident survey -d ... -d ...`

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

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

529 4.4 Validate command

530 `validate` checks poseidon datasets for structural correctness.

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

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

533 Check one or multiple Poseidon packages for structural correctness

534 Available options:

536 <code>-h,--help</code>	Show this help text
537 <code>-d,--baseDir DIR</code>	a base directory to search for Poseidon Packages (could be a Poseidon repository)
539 <code>--verbose</code>	print more output to the command line
540 <code>--ignoreGeno</code>	ignore SNP and GenoFile
541 <code>--noExitCode</code>	do not produce an explicit exit code

542 You can run it with

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

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

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

- 548 • Presence of the necessary files
- 549 • Full structural correctness of `.bib` and `.janno` file

- 550 • Superficial correctness of genotype data files. A full check would be too computationally expensive
- 551 • Correspondence of BibTeX keys in .bib and .janno
- 552 • Correspondence of individual and group IDs in .janno and genotype data files

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

556 **5 Analysis commands**

557 All analysis commands (e.g. `trident fstats`) have been moved from trident to the analysis tool `xerxes`.