

---

**datazimmer**

**Social Science Computing Unit Budapest**

**Feb 04, 2024**



## CONTENTS:

<b>1</b>	<b>To create a new project</b>	<b>3</b>
1.1	Scheduling . . . . .	3
1.2	Test projects . . . . .	3
<b>2</b>	<b>Lookahead</b>	<b>5</b>
2.1	W3C compliancy plan . . . . .	5
<b>3</b>	<b>Installation:</b>	<b>7</b>
3.1	using pip . . . . .	7
<b>4</b>	<b>Glossary</b>	<b>9</b>
4.1	Namespace . . . . .	9
4.2	Data Project . . . . .	9
4.3	Registry . . . . .	10
4.4	Metadata . . . . .	10
4.5	Config . . . . .	11
4.6	Environment . . . . .	11
4.7	Tabular data related phrases . . . . .	12
<b>5</b>	<b>Mock Projects</b>	<b>13</b>
<b>6</b>	<b>Naming Conventions and Restrictions</b>	<b>15</b>
6.1	Metadata . . . . .	15
<b>7</b>	<b>Rules</b>	<b>17</b>
<b>8</b>	<b>CLI (WIP)</b>	<b>19</b>
<b>9</b>	<b>API</b>	<b>21</b>
9.1	datazimmer Package . . . . .	21
<b>10</b>	<b>Release Notes</b>	<b>29</b>
10.1	v0.1.0 . . . . .	29
10.2	v0.1.1 . . . . .	29
10.3	v0.1.2 . . . . .	29
10.4	v0.1.3 . . . . .	29
10.5	v0.1.4 . . . . .	29
10.6	v0.1.5 . . . . .	29
10.7	v0.2.1 . . . . .	29
10.8	v0.2.2 . . . . .	29
10.9	v0.2.3 . . . . .	30

10.10 v0.2.4 . . . . .	30
10.11 v0.2.5 . . . . .	30
10.12 v0.2.6 . . . . .	30
10.13 v0.2.7 . . . . .	30
10.14 v0.3.0 . . . . .	30
10.15 v0.3.1 . . . . .	30
10.16 v0.3.10 . . . . .	30
10.17 v0.3.2 . . . . .	31
10.18 v0.3.3 . . . . .	31
10.19 v0.3.4 . . . . .	31
10.20 v0.3.5 . . . . .	31
10.21 v0.3.6 . . . . .	31
10.22 v0.3.7 . . . . .	31
10.23 v0.3.8 . . . . .	31
10.24 v0.3.9 . . . . .	31
10.25 v0.4.0 . . . . .	32
10.26 v0.4.1 . . . . .	32
10.27 v0.4.10 . . . . .	32
10.28 v0.4.11 . . . . .	32
10.29 v0.4.12 . . . . .	32
10.30 v0.4.13 . . . . .	32
10.31 v0.4.14 . . . . .	32
10.32 v0.4.15 . . . . .	32
10.33 v0.4.2 . . . . .	33
10.34 v0.4.3 . . . . .	33
10.35 v0.4.4 . . . . .	33
10.36 v0.4.5 . . . . .	33
10.37 v0.4.6 . . . . .	33
10.38 v0.4.7 . . . . .	33
10.39 v0.4.8 . . . . .	33
10.40 v0.4.9 . . . . .	33
10.41 v0.5.0 . . . . .	34
10.42 v0.5.1 . . . . .	34
10.43 v0.5.2 . . . . .	34
10.44 v0.5.3 . . . . .	34
10.45 v0.5.4 . . . . .	34
<b>11 Indices and tables</b>	<b>35</b>
<b>Python Module Index</b>	<b>37</b>
<b>Index</b>	<b>39</b>





## TO CREATE A NEW PROJECT

- make sure that `python` points to `python>=3.8` and you have `pip` and `git` then `pip install datazimmer`
- run `dz init project-name`
  - pulls `project-template`
- add a remote
  - both to `git` and `dvc` (can run `dz build-meta` to see available `dvc` remotes)
  - `git` remote can be given with `dz init`
- create, register and document steps in a pipeline you will run in different *environments*
- build metadata to exportable and serialized format with `dz build-meta`
  - if you defined importable data from other artifacts in the config, you can import them with `load-external-data`
  - ensure that you import envs that are served from sources you have access to
- build and run pipeline steps by running `dz run`
- validate that the data matches the *datascript* description with `dz validate`

### 1.1 Scheduling

- a project as a whole has a cron expression in `zimmer.yaml` to determine the schedule of reruns
- additionally, aswan projects within the `dz` project can have different cron expressions for scheduling new runs of the aswan projects

### 1.2 Test projects

TODO: document dogshow and everything else much better here





## LOOKAHEAD

- overlapping names convention
- resolve naming confusion with colassigner, colaccessor and table feature / composite type / index base classes
- abstract composite type + subclass of entity class
  - import ACT, inherit from it and specify
  - importing composite type is impossible now if it contains foreign key :(
- add option to infer data type of assigned feature
  - can be problematic b/c pandas int/float/nan issue
- create similar sets of features in a dry way
- overlapping in entities
  - detect / signal the same type of entity
- exports: postgres, postgis , superset

### 2.1 W3C compliancy plan

- test suite for compliance: <https://w3c.github.io/csvw/publishing-snapshots/PR-earl/earl.html>
- <https://github.com/w3c/csvw>
  - <https://www.w3.org/TR/2015/REC-tabular-data-model-20151217/>
  - <https://www.w3.org/TR/tabular-metadata/>

```
@article{tennison2015model,  
  title={Model for tabular data and metadata on the web},  
  author={Tennison, Jeni and Kellogg, Gregg and Herman, Ivan},  
  year={2015}  
}
```

```
@article{pollock2015metadata,  
  title={Metadata vocabulary for tabular data},  
  author={Pollock, Rufus and Tennison, Jeni and Kellogg, Gregg and Herman, Ivan},  
  journal={W3C Recommendation},  
  volume={17},  
  year={2015}  
}
```



## INSTALLATION:

### 3.1 using pip

```
pip install datazimmer
```



## GLOSSARY

### 4.1 Namespace

The atomic unit of the knowledge system containing data and metadata

- defines
  - tables
  - composite types
  - entity classes
  - code to build data based on these
- represented by
  - a module in a data project (as datascrip) - nested right below the main (src) module
  - a set of YAML files in {namespace name}/\*\*/\*.yaml as serialized metadata in the released sdist
    - \* automatically generated from the code
  - an exported .py file with basic datascrip in {namespace name}/\_\_init\_\_.py in the released sdist
- can import other namespaces, either to use
  - data (even for foreign keys in tables)
  - defined composite types / entity classes

### 4.2 Data Project

A versioned set of interconnected namespaces with metadata and different environments

- defines
  - namespaces
  - different environments where (usually) the same code runs for different data
- represented by
  - a git repository
    - \* is a DVC repository
    - \* based on a [template](#)
    - \* has fixed form tags representing the releases and data versions

## 4.3 Registry

A repository containing data about the releases and dependencies of projects to make importing namespaces straightforward

- represented by
  - a git repository (either local or remote)
  - write access needed to the repo to release to it
- contains data about
  - (named) projects
    - \* URI
    - \* versions
    - \* environment->dvc remote mapping
- contains sdists forms of metadata of projects release there
  - to set up a special PyPI index so that installation and dependency resolution is outsourced

## 4.4 Metadata

Information about the data contained in projects

- defines
  - for each namespace
    - \* tables
    - \* composite types
    - \* entity classes
- represented
  - in a project repository
    - \* defined in code (datascript object)
      - scrutable
      - entitybase
      - compositetypebase
    - \* serialized (generated from code)
      - YAML files
  - in runtime
    - \* converted as soon as possible to dataclasses in bedrock module
  - some even in data output in parquet

## 4.5 Config

- defines
  - name
  - version (this is the metadata version, the data version is determined at release)
  - default-environment name (the first environment in envs config by default)
  - validation-environments (the default-environment by default)
  - registry address (the SSCUB registry by default) TODO: link
  - imported\_projects
    - \* either a list of project names to be imported, where other than name, all default values are used
    - \* or a dictionary, where the key is the project name (in the registry), and values are:
      - version (metadata version)
      - data\_namespaces - the namespaces where loading the data is required
  - in envs for each environment (one empty env named complete by default)
    - \* params for all local namespaces and global params (namespace params default to these if not defined)
      - logged to DVC from here
    - \* environments of imported projects (where data is needed)
    - \* specific DVC remote
      - where to push data generated as outputs of running the code from namespaces (TODO - find a proper name, e.g. namespace processor) - identified by the name of a remote defined in DVC config
    - \* parent env (default-environment by default)
      - all missing keys of parameters or imported ns
- represented as `zimmer.yaml` in project root

## 4.6 Environment

A complete run of the code in an project with its values for parameters and environments for imported data

- defined by config
- ...

## 4.7 Tabular data related phrases

- feature: *A named set of columns in a table*
  - can be a primitive feature, foreign key or composite feature
- the subject of records: entity class that is represented in a table



## MOCK PROJECTS

The *dogshow standard* projects and explorer are intended to showcase all features of the package.

They have their test registry (TODO: should have sscu-budapest GitHub repo address)

- dog-show
- dograce
- dogsuccess
- dogcombine

In the code, changes from versions v0.0, v0.1 and v1.0 are noted with comments - `# add in v0.1:` at the beginning of a line or `# remove in v1.0` at the end

- to test for
- rebuilding the code with the same version
- trying (and failing) to publish a different codebase for a previously published version
- test for changed dependencies



## NAMING CONVENTIONS AND RESTRICTIONS

- project names:
  - lower case letters and non-duplicated dashes (-) not at either end
- environment and namespace names:
  - lower case letters and non-duplicated underscores (\_) not at either end

### 6.1 Metadata

- [a-b]\*\_table table names based on singular form of entity e.g. dog\_table
- see all in dogshow standard



## **RULES**

Notes of rules that are necessary for operation. Too strict or stupid rules need to change!

- Only one registered function per namespace
  - unless for separate environments, like with the helper function of data loaders and environment creation
- usually, an env run of a namespace processor (is it called this? TODO) reads and writes to one env, the one it corresponds to, but:
  - can only write to its env
  - can read from a different one, but then can only read from that one
  - this allows for registering a function that creates its environment from a base/complete set
- TableNameFeatures should be the name of table feature classes if the table name is to be inferred
- No composite features with the same prefix in the same table
- Feature name can't contain \_\_ (dunder)



## CLI (WIP)

- csv format default
- core namespace default
- datazimmer registry repository default registry

```
datazimmer pull dog-show --env top_comps datazimmer pull dogssuccess/sex_matches --parquet  
datazimmer pull something --registry uri://sg-registry/...
```





## 9.1 datazimmer Package

sscu-budapest utilities for scientific data engineering

### 9.1.1 Functions

<code>dump_dfs_to_tables(df_structable_pairs[, ...])</code>	helper function to fill the detected env of a dataset
<code>get_raw_data_path(leaf_name[, project])</code>	if project is None, raw data output path is given, otherwise imported
<code>parse_df(df, entity[, verbose])</code>	
<code>register([procfun, dependencies, outputs, ...])</code>	registers a function to the pipeline the names of parameters will matter and will be looked up in conf/envs.yaml params
<code>register_data_loader([fun, extra_deps])</code>	Convenience functions to use register with typical parameters
<code>register_env_creator([fun, extra_deps])</code>	Convenience functions to use register with typical parameters

#### **dump\_dfs\_to\_tables**

`datazimmer.dump_dfs_to_tables(df_structable_pairs: list[tuple[pd.DataFrame, ScruTable]], parse=True, skip_empty=False, **kwargs)`

helper function to fill the detected env of a dataset

#### **get\_raw\_data\_path**

`datazimmer.get_raw_data_path(leaf_name: str, project: str | None = None) → Path`

if project is None, raw data output path is given, otherwise imported

## parse\_df

`datazimmer.parse_df(df: DataFrame, entity: AbstractEntity, verbose=False)`

## register

`datazimmer.register(procfun=None, *, dependencies: list | None = None, outputs: list | None = None, outputs_nocache: list | None = None, outputs_persist: list | None = None)`

registers a function to the pipeline the names of parameters will matter and will be looked up in `conf/envs.yaml` params

## register\_data\_loader

`datazimmer.register_data_loader(fun=None, *, extra_deps=None)`

Convenience functions to use register with typical parameters

## register\_env\_creator

`datazimmer.register_env_creator(fun=None, *, extra_deps=None)`

Convenience functions to use register with typical parameters

## 9.1.2 Classes

---

*AbstractEntity()*

*CompositeTypeBase()*

*DzAswan*([global\_run])

*EntityClass*(name, identifiers, ...)

*Index*()

*Nullable*(dtype)

*PersistentState*()

*ReportFile*(filename)

*ScruTable*(entity[, entity\_key\_table\_map, ...])

*SourceUrl*(\_)

---

## AbstractEntity

```
class datazimmer.AbstractEntity
    Bases: ColAssigner
```

## CompositeTypeBase

```
class datazimmer.CompositeTypeBase
    Bases: ColAssigner
```

## DzAswan

```
class datazimmer.DzAswan(global_run=False)
    Bases: object
```

### Attributes Summary

<i>cron</i>
<i>name</i>
<i>starters</i>

### Methods Summary

<i>get_all_events</i> (handler[, only_latest])	
<i>get_aswan_status</i> ()	
<i>get_unprocessed_events</i> (handler[, only_latest])	
<i>prepare_run</i> ()	this runs prior to running the project
<i>run</i> ()	

### Attributes Documentation

```
cron:  str = ''
name:  str = None
starters: dict[type[ANY_HANDLER_T], list[str]] = {}
```

## Methods Documentation

**get\_all\_events**(*handler: ANY\_HANDLER\_T, only\_latest=True*)

**get\_aswan\_status**()

**get\_unprocessed\_events**(*handler: ANY\_HANDLER\_T, only\_latest=False*)

**prepare\_run**()

this runs prior to running the project

**run**()

## EntityClass

```
class datazimmer.EntityClass(name: str, identifiers:
    List[Union[datazimmer.metadata.atoms.PrimitiveFeature,
datazimmer.metadata.atoms.CompositeFeature,
datazimmer.metadata.atoms.ObjectProperty]] = <factory>, properties:
    List[Union[datazimmer.metadata.atoms.PrimitiveFeature,
datazimmer.metadata.atoms.CompositeFeature,
datazimmer.metadata.atoms.ObjectProperty]] = <factory>, parents:
    List[ForwardRef('EntityClass')] = <factory>, description: Optional[str] =
    None)
```

Bases: `_AtomBase`

## Attributes Summary

<i>description</i>
<i>table_all_columns</i>
<i>table_feature_cols</i>
<i>table_feature_dt_map</i>
<i>table_full_dt_map</i>
<i>table_index_cols</i>
<i>table_index_dt_map</i>

### Attributes Documentation

`description: str | None = None`

`table_all_columns`

`table_feature_cols`

`table_feature_dt_map`

`table_full_dt_map`

`table_index_cols`

`table_index_dt_map`

### Index

`class datazimmer.Index`

Bases: object

### Nullable

`class datazimmer.Nullable(dtype)`

Bases: type

### PersistentState

`class datazimmer.PersistentState`

Bases: object

### Methods Summary

---

*`get_conf()`*

*`get_full_name()`*

*`load()`*

*`save()`*

---

### Methods Documentation

`classmethod get_conf()`

`classmethod get_full_name()`

`classmethod load()`

`save()`

### ReportFile

`class datazimmer.ReportFile(filename: str)`

Bases: `object`

### Attributes Summary

<code>current_path</code>
---------------------------

### Methods Summary

<code>env_path(env)</code>
----------------------------

<code>env_posix(env)</code>
-----------------------------

<code>write_bytes(blob)</code>
--------------------------------

<code>write_text(text)</code>
-------------------------------

### Attributes Documentation

`current_path`

### Methods Documentation

`env_path(env)`

`env_posix(env)`

`write_bytes(blob)`

`write_text(text)`

## ScruTable

```
class datazimmer.ScruTable(entity: type[AbstractEntity], entity_key_table_map: dict[str, ScruTable] | None =
                             None, partitioning_cols: list[str] | None = None, max_partition_size: int | None
                             = None)
```

Bases: object

### Attributes Summary

<i>dfs</i>
<i>paths</i>

### Methods Summary

<i>env_ctx</i> (env)
<i>get_partition_paths</i> (partition_col[, env])
<i>purge</i> ()

### Attributes Documentation

**dfs**

**paths**

### Methods Documentation

**env\_ctx**(env)

**get\_partition\_paths**(partition\_col, env=None)

**purge**()

## SourceUrl

```
class datazimmer.SourceUrl(_)
```

Bases: str

### **9.1.3 Class Inheritance Diagram**



## RELEASE NOTES

### 10.1 v0.1.0

start datazimmer from sscutils

### 10.2 v0.1.1

fix build - add some cron logic

### 10.3 v0.1.2

revert build

### 10.4 v0.1.3

### 10.5 v0.1.4

### 10.6 v0.1.5

### 10.7 v0.2.1

init explorer frame

### 10.8 v0.2.2

auth improvement

## **10.9 v0.2.3**

tighter explorer

## **10.10 v0.2.4**

tighter explorer

## **10.11 v0.2.5**

multitable dec

## **10.12 v0.2.6**

minimal to dataset level

## **10.13 v0.2.7**

scrutable extension

## **10.14 v0.3.0**

simplification and unification

## **10.15 v0.3.1**

incremental fixes

## **10.16 v0.3.10**

flat index and improved secrets

## 10.17 v0.3.2

requirement and explorer updates

## 10.18 v0.3.3

init explorer and os fix

## 10.19 v0.3.4

nb requirement fix

## 10.20 v0.3.5

title detection fix

## 10.21 v0.3.6

explorer upgrades

## 10.22 v0.3.7

freeze jb version

## 10.23 v0.3.8

flexible pulling

## 10.24 v0.3.9

sql loader extension

## **10.25 v0.4.0**

aswan integration init

## **10.26 v0.4.1**

here goes nothing

## **10.27 v0.4.10**

dvc realign, minor extensions and reorg

## **10.28 v0.4.11**

stage non cached outputs

## **10.29 v0.4.12**

private zenodo

## **10.30 v0.4.13**

aswan reset plus scrutable pickle fix

## **10.31 v0.4.14**

minimal alignment, pre-dependency fix

## **10.32 v0.4.15**

move import

## **10.33 v0.4.2**

collect dependency fix

## **10.34 v0.4.3**

minor iterations, redundancy removal

## **10.35 v0.4.4**

git fixes, unifications

## **10.36 v0.4.5**

aswan alignment

## **10.37 v0.4.6**

aswan align

## **10.38 v0.4.7**

simplifications and raw-data introduction

## **10.39 v0.4.8**

properly citeable

## **10.40 v0.4.9**

init zenodo integration

## **10.41 v0.5.0**

fix dependencies and drop explorer

## **10.42 v0.5.1**

many many bugs fixed

## **10.43 v0.5.2**

external dvc zimmauth

## **10.44 v0.5.3**

properly rm dvc as dependency

## **10.45 v0.5.4**

metaclass def bugfix

## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`





## PYTHON MODULE INDEX

### d

`datazimmer`, [21](#)



## A

`AbstractEntity` (class in `datazimmer`), 23

## C

`CompositeTypeBase` (class in `datazimmer`), 23

`cron` (`datazimmer.DzAswan` attribute), 23

`current_path` (`datazimmer.ReportFile` attribute), 26

## D

`datazimmer`  
module, 21

`description` (`datazimmer.EntityClass` attribute), 25

`dfs` (`datazimmer.Scrutable` attribute), 27

`dump_dfs_to_tables()` (in module `datazimmer`), 21

`DzAswan` (class in `datazimmer`), 23

## E

`EntityClass` (class in `datazimmer`), 24

`env_ctx()` (`datazimmer.Scrutable` method), 27

`env_path()` (`datazimmer.ReportFile` method), 26

`env_posix()` (`datazimmer.ReportFile` method), 26

## G

`get_all_events()` (`datazimmer.DzAswan` method), 24

`get_aswan_status()` (`datazimmer.DzAswan` method), 24

`get_conf()` (`datazimmer.PersistentState` class method), 26

`get_full_name()` (`datazimmer.PersistentState` class method), 26

`get_partition_paths()` (`datazimmer.Scrutable` method), 27

`get_raw_data_path()` (in module `datazimmer`), 21

`get_unprocessed_events()` (`datazimmer.DzAswan` method), 24

## I

`Index` (class in `datazimmer`), 25

## L

`load()` (`datazimmer.PersistentState` class method), 26

## M

module  
    `datazimmer`, 21

## N

`name` (`datazimmer.DzAswan` attribute), 23

`Nullable` (class in `datazimmer`), 25

## P

`parse_df()` (in module `datazimmer`), 22

`paths` (`datazimmer.Scrutable` attribute), 27

`PersistentState` (class in `datazimmer`), 25

`prepare_run()` (`datazimmer.DzAswan` method), 24

`purge()` (`datazimmer.Scrutable` method), 27

## R

`register()` (in module `datazimmer`), 22

`register_data_loader()` (in module `datazimmer`), 22

`register_env_creator()` (in module `datazimmer`), 22

`ReportFile` (class in `datazimmer`), 26

`run()` (`datazimmer.DzAswan` method), 24

## S

`save()` (`datazimmer.PersistentState` method), 26

`Scrubable` (class in `datazimmer`), 27

`SourceUrl` (class in `datazimmer`), 27

`starters` (`datazimmer.DzAswan` attribute), 23

## T

`table_all_columns` (`datazimmer.EntityClass` attribute), 25

`table_feature_cols` (`datazimmer.EntityClass` attribute), 25

`table_feature_dt_map` (`datazimmer.EntityClass` attribute), 25

`table_full_dt_map` (`datazimmer.EntityClass` attribute), 25

`table_index_cols` (`datazimmer.EntityClass` attribute), 25

`table_index_dt_map` (`datazimmer.EntityClass` attribute), 25

## W

`write_bytes()` (*datazimmer.ReportFile* method), [26](#)

`write_text()` (*datazimmer.ReportFile* method), [26](#)