

# NPLinker Community Meeting

2024-03-05

netherlands  
eScience center

# Agenda

- 16:00 Short Introduction
- 16:05 Announcements and upcoming events
- 16:10 NPLinker development
- 16:35 Monthly plan
- 16:40 Q&A
- 17:00 .



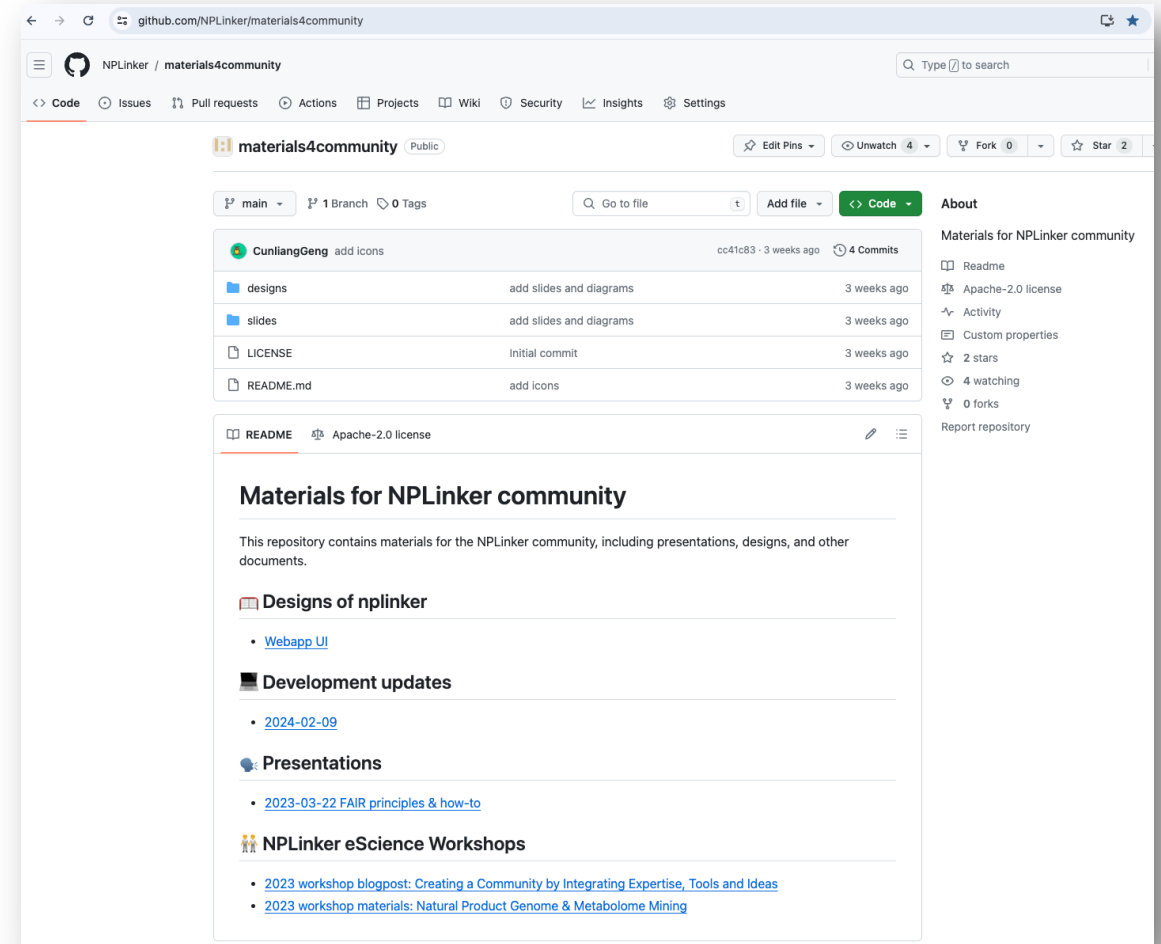
## Announcements and upcoming events

Arjan Draisma has joined the force of nplinker development and will first integrate BigScape v2 into NPLinker.

New repo for sharing the materials for NPLinker community

<https://github.com/NPLinker/materials4community>

The slides of community meetings are shared there.



# NPLinker Development

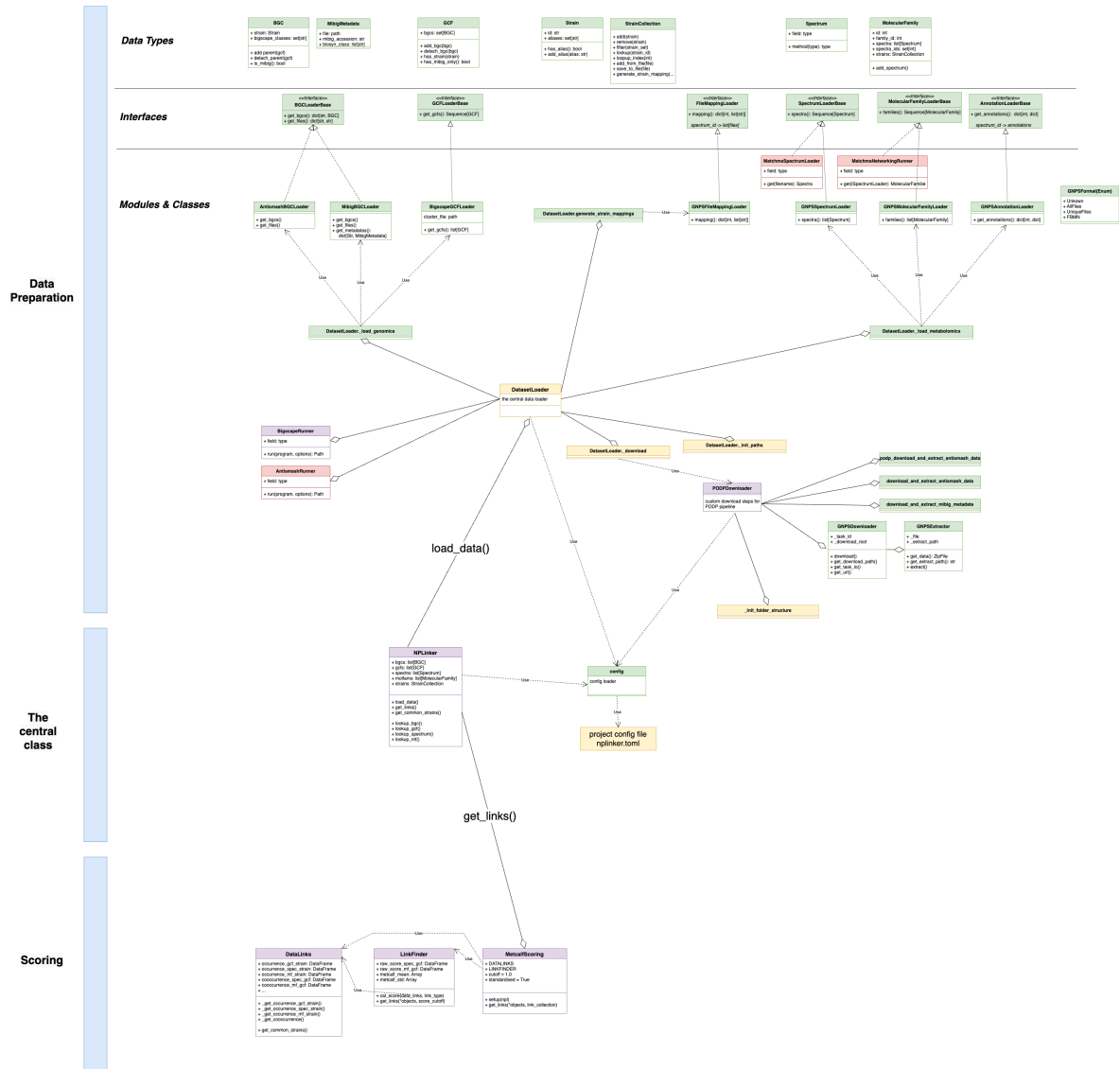
New features & changes

Contributors: Cunliang Geng, Giulia Crocioni

# NPLinker architecture

Updated: 2024-02-05

Ideas    Todo    In progress    Done



## Data Preparation

- **Arranging data:**
  - Initialize project folder
  - Provide or download/extract data
  - Validate data
- **Loading data:** load data from files to python objects
  - ☒ in February's update

## New project directory structure

The project directory structure was redesigned to make it more structured, well-defined and flexible.

It's also important for enabling the local mode.

- **local mode**: users provide input data directly to nplinker
- **podp mode**: users provide a PODP id and let nplinker get input data from PODP platform

With this new project dir structure, users are required to

- Create a `root\_dir` manually
- Create a config file `nplinker.toml`
- [Local mode] Provide
  - `strain\_mappings.json`
  - input data (gnps, antismash, bigscape...)
- Run nplinker in the `root\_dir`

```

root_dir
├── nplinker.toml
├── strain_mappings.json
├── strains_selected.json
├── downloads
│   ├── paired_datarecord_4b29ddc3-26d0-40d7-80c5-44fb6631dbf9.4.json
│   ├── GCF_000016425.1.zip
│   ├── GCF_0000514975.1.zip
│   ├── c22f44b14a3d450eb836d607cb9521bb.zip
│   ├── genome_status.json
│   └── mibig_json_3.1.tar.gz
├── gnps
│   ├── spectra.mgf
│   ├── molecular_families.tsv
│   ├── annotations.tsv
│   └── file_mappings.tsv
├── antismash
│   ├── GCF_000016425.1
│   │   ├── xxx.region001.gbk
│   │   └── ...
│   ├── GCF_000016425.1
│   │   ├── xxxx.region001.gbk
│   │   └── ...
│   └── ...
├── bigscape
│   ├── mix_clustering_c0.30.tsv
│   ├── bigscape_running_output
│   └── ...
├── mibig
│   ├── BGC0000001.json
│   ├── BGC0000002.json
│   └── ...
└── ...

```

## New template for config file

The settings for NPLinker were redesigned to make the config file as simple as possible.

Now users usually only need to set values for `root\_dir`,  
`mode` and/or `podp\_id`.

```
1 #####
2 # NPLinker configuration file
3 #####
4
5 root_dir = "<NPLinker root directory>"
6 mode = "podp"
7 podp_id = ""
8
9 [log]
10 level = "INFO"
11 file = "path/to/logfile"
12 to_stdout = true
13
14 [mibig]
15 to_use = true
16 version = "3.1"
17
18 [bigscape]
19 parameters = "--mibig --clans-off --mix --include_singletons --cutoffs 0.30"
20 cutoff = "0.30"
21
22 [scoring]
23 methods = ["metcalf"]
```



## New config loader and validators

We added the config loader and validators by taking advantage of the library Dynaconf.

So users are able to

- specify the path to config file with the env variable  
`NPLINKER\_CONFIG\_FILE`
- Or directly put the config file `nplinker.toml` in the root dir

Validation of the config will be triggered by python's import, e.g. `import nplinker`, so validation is done even before running any code of nplinker.

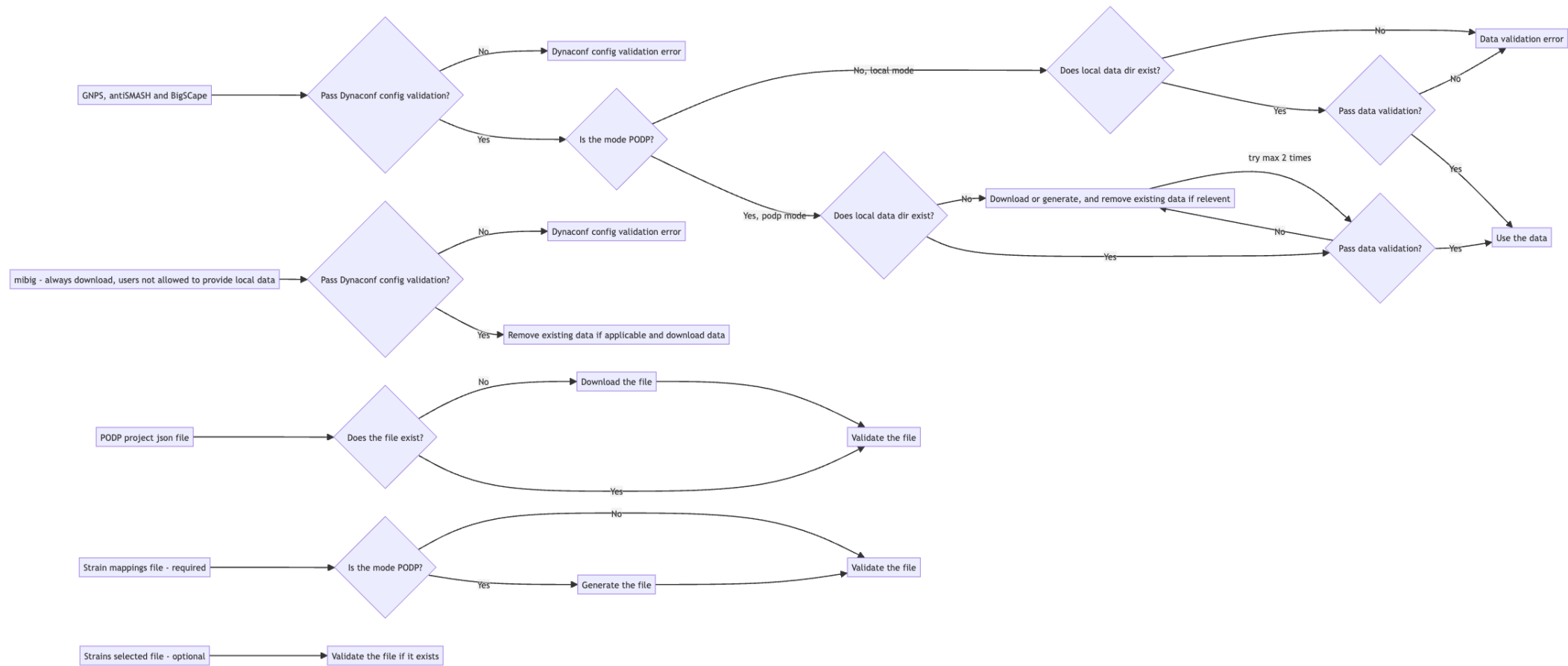
```
16 # Locate the user's config file
17 user_config_file = os.environ.get("NPLINKER_CONFIG_FILE", "nplinker.toml")
18 if not os.path.exists(user_config_file):
19     raise FileNotFoundError(f"Config file '{user_config_file}' not found")
20
21 # Locate the default config file
22 default_config_file = Path(__file__).resolve().parent / "nplinker_default.toml"
23
24 # Load config files
25 config = Dynaconf(settings_files=[user_config_file], preload=[default_config_file])
26
27
28
29
30
31 validators = [
32     # General settings
33     # root_dir value is transformed to a pathlib.Path object and must be a directory.
34     Validator(
35         "root_dir", required=True, cast=transform_to_full_path, condition=lambda v: v.is_dir()
36     ),
37     Validator("mode", required=True, cast=lambda v: v.lower(), is_in=["local", "podp"]),
38     # root_dir must be set if mode is "podp"; must not be set if mode is "local".
39     Validator("podp_id", required=True, when=Validator("mode", eq="podp")),
40     Validator("podp_id", required=False, when=Validator("mode", eq="local")),
41     # Log
42     # loglevel must be a string and must be one of the supported levels. It is transformed to
43     # uppercase to avoid case sensitivity.
44     Validator(
45         "log.level",
46         is_type_of=str,
47         cast=lambda v: v.upper(),
48         is_in=["NOTSET", "DEBUG", "INFO", "WARNING", "ERROR", "CRITICAL"],
49     ),
50     Validator("log.file", is_type_of=str, cast=Path),
51     Validator("log.stdout", is_type_of=bool),
52     # Mibig
53     Validator("mibig.to_use", required=True, is_type_of=bool),
54     Validator(
55         "mibig.version",
56         required=True,
57         is_type_of=str,
58         when=Validator("mibig.to_use", eq=True),
59     ),
60     # BigScape
61     Validator("bigscape.parameters", required=True, is_type_of=str),
62     Validator("bigscape.cutoff", required=True, is_type_of=str),
63     # Scoring
64     # scoring.methods must be a list of strings and must contain at least one of the
65     # supported scoring methods.
66     Validator(
67         "scoring.methods",
68         required=True,
69         cast=lambda v: [i.lower() for i in v],
70         is_type_of=list,
71         len_min=1,
72         condition=lambda v: set(v).issubset({"metcalf", "rosetta"}),
73     ),
74 ]
```





# Pipeline of arranging data

The pipeline was designed to keep arranging data independently for each type of data while keeping them sharing common logics.



## DatasetArranger

The data arranging was coupled with data loading before.

By following the new pipeline of arranging data, a new class

`DatasetArranger` and some validation functions were developed to be independently responsible for the steps of arranging data.

It supports both the local mode and podp mode.

**It's ready for you to try and test the refactored nplinker using your own data when this PR is reviewed and merged to dev branch.**

```
140     ... def load_data(self):  
141     ...     """Loads the basic components of a dataset."""  
142     ...     arranger = DatasetArranger()  
143     ...     arranger.arrange()  
144     ...     self._loader.load()
```



# Monthly plan

What will happen before next meeting

New Kanban board enhanced with planning feature

<https://github.com/orgs/NPLinker/projects/9>

github.com/orgs/NPLinker/projects/9/views/6

NPLinker / Projects / dev

Q Type to search

dev

OverviewCurrent monthNext MonthRoadmapKanban AllNew view

Filter by keyword or by field

Title	Month	Status	Size	Labels	Linked pull requests
> 2024 Feb 12 Feb 16 - Mar 03					
v 2024 March 7 Mar 04 - Mar 31 Current					
13 writing tutorials and docs for preparing data #204	2024 March	In progress	L	docs	
14 Inform users to start trying refactored nplinker	2024 March	Backlog			
15 Design the architecture of Webapp	2024 March	Backlog			
16 Abstraction of run_bigscape function #115	2024 March	Ready		GEN	
17 support bigscape v2 #216	2024 March	Backlog		GEN	
18 Further refactor downloaders including PODODownloader	2024 March	Backlog			
19 Remove unused GNPSLoader	2024 March	Backlog			
+ Add item					



# Q&A

Questions? Ideas? Feedback?