
Tentaclio

octoenergy

Apr 20, 2023

CONTENTS:

1	Tentacio	1
2	Quick Examples.	3
2.1	Read and write streams.	3
2.2	Copy streams	3
2.3	Delete resources	3
2.4	List resources	4
2.5	Authenticated resources.	4
2.6	Database connections.	4
2.7	Pandas interaction.	4
3	Installation	7
3.1	Developing.	7
3.2	How to use	7
3.3	Quick note on protocols structural subtyping.	11
4	Documentation	13
4.1	Building docs	13
5	Indices and tables	15

**CHAPTER
ONE**

TENTACLIO

Python library that simplifies:

- Handling streams from different protocols such as `file:`, `ftp:`, `sftp:`, `s3:`, ...
- Opening database connections.
- Managing the credentials in distributed systems.

Main considerations in the design:

- Easy to use: all streams are open via `tentaclio.open`, all database connections through `tentaclio.db`.
- URLs are the basic resource locator and db connection string.
- Automagic authentication for protected resources.
- Extensible: you can add your own handlers for other schemes.
- Pandas interaction.

QUICK EXAMPLES.

2.1 Read and write streams.

```
import tentaclio
contents = " "

with tentaclio.open("ftp://localhost:2021/upload/file.txt", mode="w") as writer:
    writer.write(contents)

# Using boto3 authentication under the hood.
bucket = "s3://my-bucket/octopus/hello.txt"
with tentaclio.open(bucket) as reader:
    print(reader.read())
```

2.2 Copy streams

```
import tentaclio

tentaclio.copy("/home/constantine/data.csv", "sftp://constantine:tentacle@ftp.
↪octoenergy.com/uploads/data.csv")
```

2.3 Delete resources

```
import tentaclio

tentaclio.remove("s3://my-bucket/octopus/the-9th-tentacle.txt")
```

2.4 List resources

```
import tentacio

for entry in tentacio.listdir("s3://mybucket/path/to/dir"):
    print("Entry", entry)
```

2.5 Authenticated resources.

```
import os

import tentacio

print("env ftp credentials", os.getenv("OCTOIO__CONN__OCTOENERGY_FTP"))
# This prints `sftp://constantine:tentacle@sftp.octoenergy.com/`

# Credentials get automatically injected.

with tentacio.open("sftp://sftp.octoenergy.com/uploads/data.csv") as reader:
    print(reader.read())
```

2.6 Database connections.

```
import os

import tentacio

print("env TENTACLIO__CONN__DB", os.getenv("TENTACLIO__CONN__DB"))

# This prints `postgresql://octopus:tentacle@localhost:5444/example`

# hostname is a wildcard, the credentials get injected.
with tentacio.db("postgresql://hostname/example") as pg:
    results = pg.query("select * from my_table")
```

2.7 Pandas interaction.

```
import pandas as pd  #
import tentacio  #

df = pd.DataFrame([[1, 2, 3], [10, 20, 30]], columns=["col_1", "col_2", "col_3"])

bucket = "s3://my-bucket/data/pandas.csv"

with tentacio.open(bucket, mode="w") as writer: # supports more pandas readers
```

(continues on next page)

(continued from previous page)

```
df.to_csv(writer, index=False)

with tentaclio.open(bucket) as reader:
    new_df = pd.read_csv(reader)

# another example: using pandas.DataFrame.to_sql() with tentaclio to upload
with tentaclio.db(
    connection_info,
    connect_args={'options': '-csearch_path=schema_name'}
) as client:
    df.to_sql(
        name='observations', # table name
        con=client.conn,
    )
```

CHAPTER THREE

INSTALLATION

You can get tentaclio using pip

```
pip install tentaclio
```

or pipenv

```
pipenv install tentaclio
```

3.1 Developing.

Clone this repo and install [pipenv](#):

In the `Makefile` you'll find some useful targets for linting, testing, etc. i.e.:

```
make test
```

3.2 How to use

This is how to use `tentaclio` for your daily data ingestion and storing needs.

3.2.1 Streams

In order to open streams to load or store data the universal function is:

```
import tentaclio

with tentaclio.open("/path/to/my/file") as reader:
    contents = reader.read()

with tentaclio.open("s3://bucket/file", mode='w') as writer:
    writer.write(contents)
```

Allowed modes are `r`, `w`, `rb`, and `wb`. You can use `t` instead of `b` to indicate text streams, but that's the default.

In order to keep tentaclio as light as possible, it only includes `file`, `ftp`, `sftp`, `http` and `https` schemes by default. However, many more are easily available by installing extra packages:

Default:

- /local/file
- file:///local/file
- ftp://path/to/file
- sftp://path/to/file
- http://host.com/path/to/resource
- https://host.com/path/to/resource

tentacio-s3

- s3://bucket/file

tentacio-gs

- gs://bucket/file
- gsc://bucket/file

tentacio-gdrive

- gdrive:/My Drive/file
- googledrive:/My Drive/file

tentacio-postgres

- postgresql://host/database::table will allow you to write from a csv format into a database with the same column names (note that the table goes after :: :warning:).

You can add the credentials for any of the urls in order to access protected resources.

You can use these readers and writers with pandas functions like:

```
import pandas as pd
import tentacio

with tentacio.open("/path/to/my/file") as reader:
    df = pd.read_csv(reader)

[...]

with tentacio.open("s3::/path/to/my/file", mode='w') as writer:
    df.to_parquet(writer)
```

Readers, Writers and their closeable versions can be used anywhere expecting a file-like object; pandas or pickle are examples of such functions.

Notes on writing files for Spark, Presto, and similar downstream systems

The default behaviour for the `open` context manager in python is to create an empty file when opening it in writable mode. This can be annoying if the process that creates the data within the `with` clause yields empty dataframes and nothing gets written. This will make Spark and Presto panic.

To avoid this we can make the stream *empty safe* so the empty buffer won't be flushed if no writes have been performed so no empty file will be created.

```
with tio.make_empty_safe(tio.open("s3://bucket/file.parquet", mode="wb")) as writer:
    if not df.empty:
        df.to_parquet(writer)
```

3.2.2 File system like operations to resources

Listing resources

Some URL schemes allow listing resources in a pythonic way:

```
import tentacio

for entry in tentacio.listdir("s3://mybucket/path/to/dir"):
    print("Entry", entry)
```

Whereas `listdir` might be convinient we also offer `scandir`, which returns a list of `DirEntries`, and, `walk`. All functions follow as closely as possible their standard library definitions.

3.2.3 Database access

In order to open db connections you can use `tentacio.db` and have instant access to postgres, sqlite, athena and mssql.

```
import tentacio

[...]

query = "select 1";
with tentacio.db(POSTGRES_TEST_URL) as client:
    result = client.query(query)
[...]
```

The supported db schemes are:

Default:

- `sqlite://`
- `mssql://`
- – Any other scheme supported by sqlalchemy.

`tentacio-postgres`

- `postgresql://`

`tentacio-athena`

- awsathena+rest://
tentacio-databricks
 - databricks+thrift://
tentacio-snowflake
 - snowflake://

Extras for databases

For postgres you can set the variable TENTACLI0__PG_APPLICATION_NAME and the value will be injected when connecting to the database.

3.2.4 Automatic credentials injection

1. Configure credentials by using environmental variables prefixed with TENTACLI0__CONN__ (i.e. TENTACLI0__CONN__DATA_FTP=sftp://real_user:1321dsf@ftp.octoenergy.com).
2. Open a stream:

```
with tentacio.open("sftp://ftp.octoenergy.com/file.csv") as reader:  
    reader.read()
```

The credentials get injected into the url.

3. Open a db client:

```
import tentacio  
  
with tentacio.db("postgresql://hostname/my_data_base") as client:  
    client.query("select 1")
```

Note that hostname in the url to be authenticated is a wildcard that will match any hostname. So authenticate("http://hostname/file.txt") will be injected to http://user:pass@octo.co/file.txt if the credential for http://user:pass@octo.co/ exists.

Different components of the URL are set differently:

- Scheme and path will be set from the URL, and null if missing.
- Username, password and hostname will be set from the stored credentials.
- Port will be set from the stored credentials if it exists, otherwise from the URL.
- Query will be set from the URL if it exists, otherwise from the stored credentials (so it can be overridden)

Credentials file

You can also set a credentials file that looks like:

```
secrets:
  db_1: postgresql://user1:pass1@myhost.com/database_1
  db_2: mssql://user2:pass2@otherhost.com/database_2?
  ↵driver=ODBC+Driver+17+for+SQL+Server
  ftp_server: ftp://fuser:ffpass@ftp.myhost.com
```

And make it accessible to tentacio by setting the environmental variable TENTACLIO__SECRETS_FILE. The actual name of each url is for traceability and has no effect in the functionality.

(Note that you may need to add ?driver={driver from /usr/local/etc/odbcinst.ini} for mssql database connection strings; see above example)

Alternatively you can run curl https://raw.githubusercontent.com/octoenergy/tentacio/master/extras/init_tentacio.sh to create a secrets file in `~/.tentacio.yml` and automatically configure your environment.

3.3 Quick note on protocols structural subtyping.

In order to abstract concrete dependencies from the implementation of data related functions (or in any part of the system really) we use typed [protocols](#). This allows a more flexible dependency injection than using subclassing or [more complex approaches](#). This idea is heavily inspired by how this exact thing is done in [go](#). Learn more about this principle in our [tech blog](#).

**CHAPTER
FOUR**

DOCUMENTATION

Tentacio uses reStructuredText (Restructured Text), MyST (Markedly Structured Text) and the Sphinx documentation system. This allows it to be built into other forms for easier viewing and browsing.

4.1 Building docs

To create an HTML version of the docs, use:

```
$ cd docs  
$ make html
```

This will generate a static set of HTML files with root *docs/_build/html/index.html*. This can be viewed in a browser

**CHAPTER
FIVE**

INDICES AND TABLES

- genindex