

# CrocoLakeTools: A Python package to convert ocean observations to the parquet format

Enrico Milanese<sup>1</sup>, David Nicholson<sup>1</sup>, Gaël Forget<sup>2</sup>, and Susan Wijffels<sup>1</sup>

<sup>1</sup> Woods Hole Oceanographic Institution, United States of America <sup>2</sup> Massachusetts Institute of Technology, United States of America

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Open Journals](#) ↗

## Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

Investigations of the ocean state are possible thanks to the ever growing number of measurements performed with multiple instruments by different research missions. The vast and variegated efforts have brought the community to define data storage conventions (e.g. CF-netCDF) and to assemble collections of datasets (e.g. the World Ocean Database). Yet, accessing these datasets often requires the usage of multiple tools, is inefficient over the cloud, and presents an overall high entrance barrier in terms of knowledge and time required to effectively access these resources. CrocoLakeTools is a Python package that address those shortcomings by providing workflows to convert several datasets from their original format to a uniform parquet dataset with a shared schema.

## Statement of need

CrocoLakeTools is a Python package to build workflows that convert ocean observations from different formats (e.g. netCDF, CSV) to parquet. CrocoLakeTools take advantage of Python's well-established and growing ecosystem of open tools: it uses dask's parallel computing capabilities to convert multiple files at once and to handle larger-than-memory data. dask is already well-integrated with xarray and pandas, two widely used Python libraries for the treatment of array and tabular data, respectively, and with pyarrow, the API to the Apache Arrow library which is used to generate the parquet dataset.

Parquet is a data storage format for big tidy data which presents several advantages: it is language agnostic (it can be accessed with Python, Matlab, Julia and web development technologies); it offers faster reading performances than other tabular formats such as CSV; it is optimized for cloud systems storage and operations; it is widespread in the data science community, leading to a multitude of freely accessible tools and educational material.

CrocoLakeTools was developed with the goal of building and serving CrocoLake, a regularly refreshed database of oceanographic observations that are pre-filtered to contain only quality-controlled measurements. CrocoLakeTools was designed to be used by researchers and data scientists and engineers in oceanography. CrocoLake was designed to be accessed by the wider oceanographic community.

## Code architecture

### Converters

The core task of CrocoLakeTools is to take one or more files from a dataset and convert them to parquet, ensuring that CrocoLake's schema is followed. This is achieved through the methods contained in the Converter class and its subclasses. While the conversion of all

39 datasets requires some general functionality (e.g. renaming the original variables to the final  
40 schema), each conversion requires specific tools for the specific dataset (e.g. the map used to  
41 rename the variable). CrocoLakeTools then hosts a converter for each dataset that implement  
42 the specific needs of that datasets and inherits from Converter, which contains the shared  
43 methods.

## 44 Workflow

45 The first step in the workflow is to retrieve the original files. This is generally left to the user,  
46 although we provide methods to download Argo data and we wish to be able in the future to  
47 provide this type of tools for other datasets as well (Figure 1). The second step is to convert  
48 the data to parquet, and finally merge the datasets into CrocoLake (Figure 2).

## 49 Original sources

50 These are the original data provided by the project, mission, scientist, etc. The format, schema,  
51 nomenclature and conventions are those defined by the individual project and are unaware of  
52 CrocoLake's workflow.

## 53 Local mirrors

54 Modules to download the original data are optional. They should inherit from the Downloader  
55 class and be called downloader<ProjectName>, e.g. downloaderArgoGDAC. Whether a down-  
56 loader module exists or the user downloads the data themselves, the original data is stored on  
57 disk and this is the starting point for the converter.

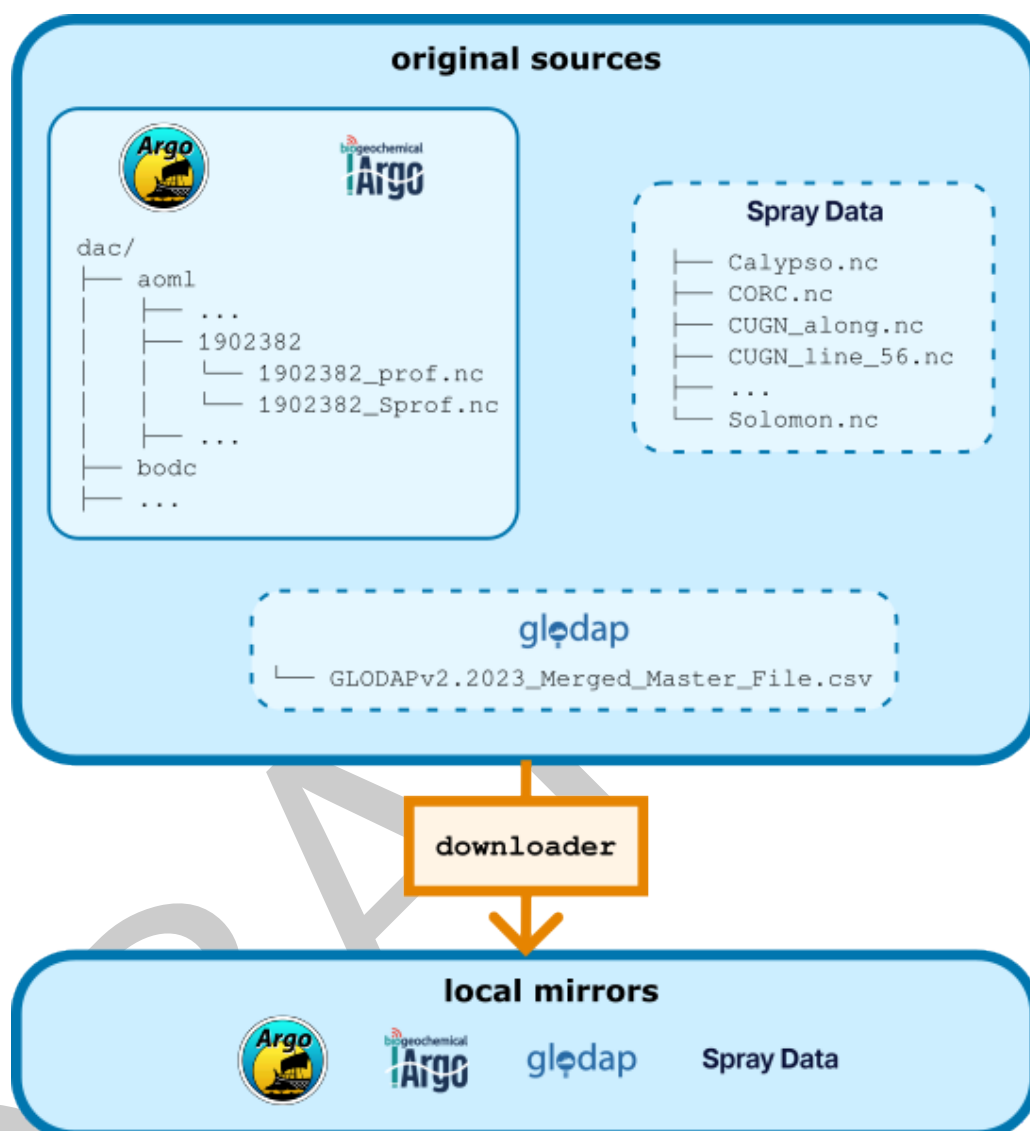
## 58 Parquet datasets

59 The core of CrocoLakeTools are the modules in the Converter class and its subclasses. Each  
60 project has its own subclass called converter<ProjectName>, e.g. converterGLODAP; further  
61 specifiers can be added as necessary (e.g. at this time there are a few different converters for Argo  
62 data to prepare different datasets). The need for a dedicated converter for each project despite  
63 the usage of common data formats (e.g. netCDF, CSV) is due to differences in the schema,  
64 e.g. variable names, units, etc., while the steps that can be generalized for each format are  
65 usually already included in other libraries that CrocoLakeTools rely on (e.g. pandas, dask,  
66 pyarrow). Depending on the dataset, multiple converters can be applied. For example, to create  
67 CrocoLake, Argo data goes through two converters: 1. converterArgoGDAC, which converts  
68 the original Argo GDAC preserving most of its original conventions; 2. converterArgoQC,  
69 which takes the output of the previous step and applies some filtering based on Argo's QC  
70 flags and makes the data conforming to CrocoLake's schema.

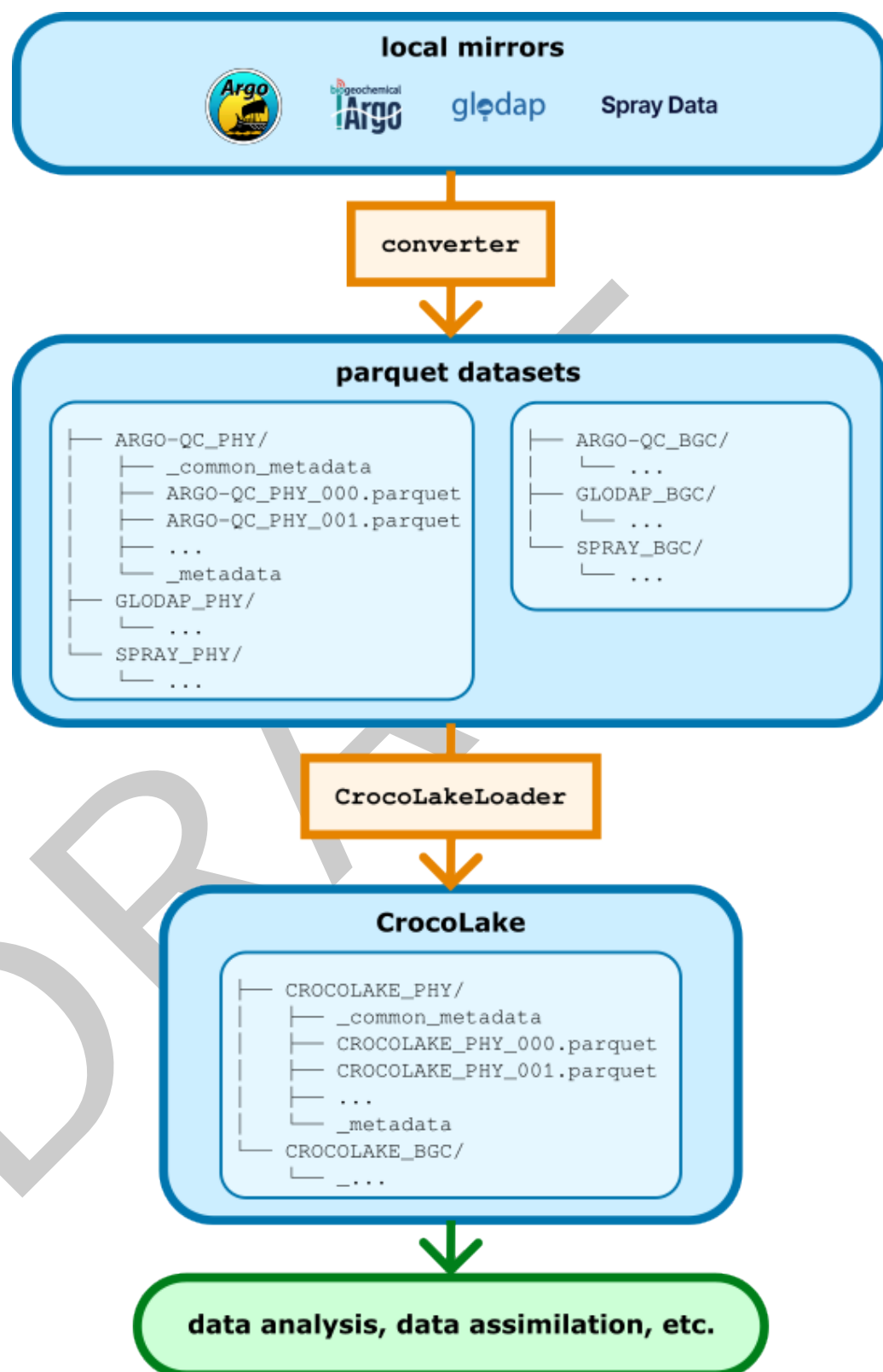
## 71 CrocoLake

72 CrocoLake contains each converted dataset. The first step to build it is to create a direc-  
73 tory containing symbolic links to the converted datasets (an example script is provided).  
74 The submodule CrocoLakeLoader then allows to seamlessly load all the converted datasets  
75 into memory as one dask dataframe with a uniform schema, using just a few lines. The  
76 script merge\_crocolake.py exploits CrocoLakeLoader's capabilities to generate one merged  
77 CrocoLake dataset that contains all the converted datasets and is stored back to disk (in  
78 parquet).

79 CrocoLake can be accessed with several programming languages with just a few lines of codes:  
80 the submodules CrocoLake-Python, CrocoLake-Matlab, and CrocoLake-Julia contain tools and  
81 examples in some languages.



**Figure 1:** CrocoLake's workflow: 'downloader's. 'CrocoLakeTools' is set up to host modules that are dedicated to download the desired datasets from the web. It currently supports the download only of Argo data (solid line subset), and other datasets require the user to download them manually (dashed border subsets).



**Figure 2:** CrocoLake's workflow. 'converter's read the data in their original format, transform it following CrocoLake's conventions, converts it to parquet, and stores it back to disk. Each dataset is converted to its own parquet version. Thanks to the submodule 'CrocoLakeLoader', multiple parquet datasets are merged into a uniform dataframe which is save to disk as CrocoLake. For each dataset, a version containing only physical variables ('PHY') or also biogeochemical variables ('BGC') can be generated.

82 **Documentation**

83 Documentation is available at [TD add link]. It describes the CrocoLake dataset, the sub-  
84 datasets that it is made, and how they are obtained. It provides references to the original files  
85 origins for download before converting them.

86 **Citation**

87 If you use CrocoLakeTools and/or CrocoLake, do not limit yourself to citing this manuscript.  
88 Remember to cite the datasets that you have used as indicated in the documentation. For  
89 example, if your work relies on Argo measurements, acknowledge Argo (Wong et al., 2020).  
90 This is important for each product to track their impact.

91 **Acknowledgements**

92 We acknowledge funding from [NSF CSSI CROCODILE details]

93 **References**

94 Wong, A. P., Wijffels, S. E., Riser, S. C., Pouliquen, S., Hosoda, S., Roemmich, D., Gilson,  
95 J., Johnson, G. C., Martini, K., Murphy, D. J., & others. (2020). Argo data 1999–2019:  
96 Two million temperature-salinity profiles and subsurface velocity observations from a global  
97 array of profiling floats. *Frontiers in Marine Science*, 7, 700.

DRAFT