

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

Enrico Milanese¹, David Nicholson¹, Gaël Forget², and Susan Wijffels¹

¹ Woods Hole Oceanographic Institution, United States of America ² 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Published: unpublished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 addresses 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 takes 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, engineers and data scientists in oceanography and 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 datasets requires some general functionality (e.g. renaming the original variables to the final

39 schema), each conversion requires specific tools for the specific dataset (e.g. the map used to
40 rename the variables). CrocoLakeTools then contains the Converter class, which contains the
41 methods shared across datasets, and from which converter subclasses inherit and implement
42 the specific needs of each dataset.

43 Workflow

44 Local mirrors

45 The first step in the workflow is to retrieve the original files (Figure 1). The original sources
46 follow the format, schema, nomenclature and conventions defined by the individual project
47 (mission, scientist, etc.) the generated them and are unaware of CrocoLake's workflow. Modules
48 to download the original data are optional. They should inherit from the Downloader class
49 and be called downloader<DatasetName> (e.g. downloaderArgoGDAC). At the time of writing
50 CrocoLakeTools is released with a downloader to build a local mirror of the Argo GDAC, and
51 we hope to support more in the future. Whether a downloader module exists or the user
52 downloads the data themselves, the original data is stored on disk and this is the starting point
53 for the converter.

54 Parquet datasets

55 The second step is to convert the data to parquet, and finally merge the datasets into
56 CrocoLake (Figure 2). The core of CrocoLakeTools are the modules in the Converter class
57 and its subclasses. Each original dataset has its own subclass called converter<DatasetName>,
58 e.g. converterGLODAP; further specifiers can be added as necessary (e.g. at this time there a
59 few different converters for Argo data to prepare different datasets). The need for a dedicated
60 converter for each project despite the usage of common data formats (e.g. netCDF, CSV) is
61 due to differences in the schema, e.g. variable names, units, etc. Depending on the dataset,
62 multiple converters can be applied. For example, to create CrocoLake, Argo data goes through
63 two converters: 1. converterArgoGDAC, which converts the original Argo GDAC preserving
64 most of its original conventions; 2. converterArgoQC, which takes the output of the previous
65 step and applies some filtering based on Argo's QC flags and makes the data conforming to
66 CrocoLake's schema.

67 CrocoLake

68 CrocoLake is one parquet dataset that contains each converted dataset merged together. This
69 can be achieved with the script merge_crocolake.py. The script first creates a directory con-
70 taining symbolic links to each converted dataset. It then uses the submodule CrocoLakeLoader
71 to seamlessly load all the converted datasets into memory as one dask dataframe with a
72 uniform schema and merge them into CrocoLake and stores it back to disk.

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

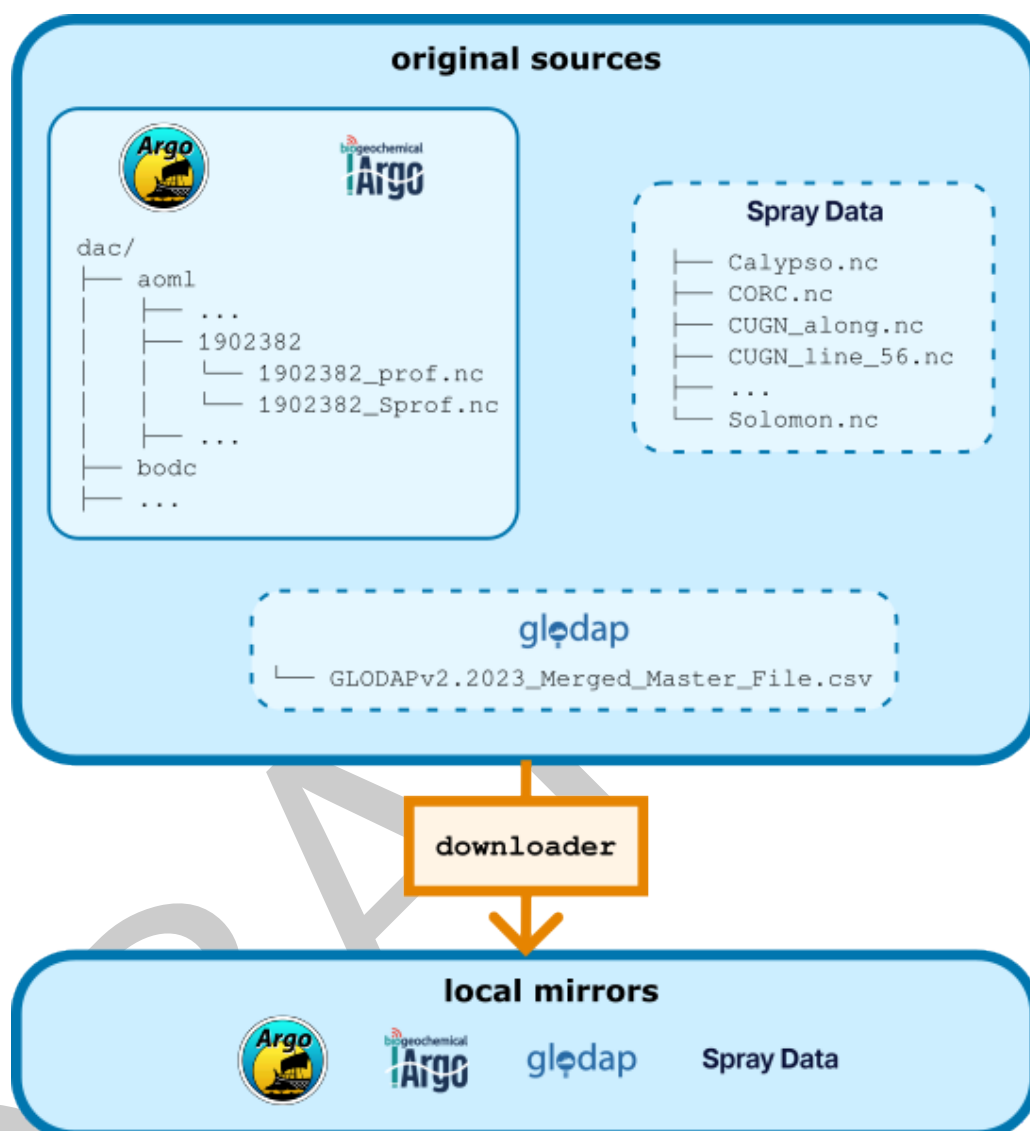


Figure 1: CrocoLake's workflow: 'downloader's. 'CrocoLakeTools' is set up to host modules that are dedicated to download the desired datasets from the web. It currently supports the download only of Argo data (solid line subset), and other datasets require the user to download them manually (dashed border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

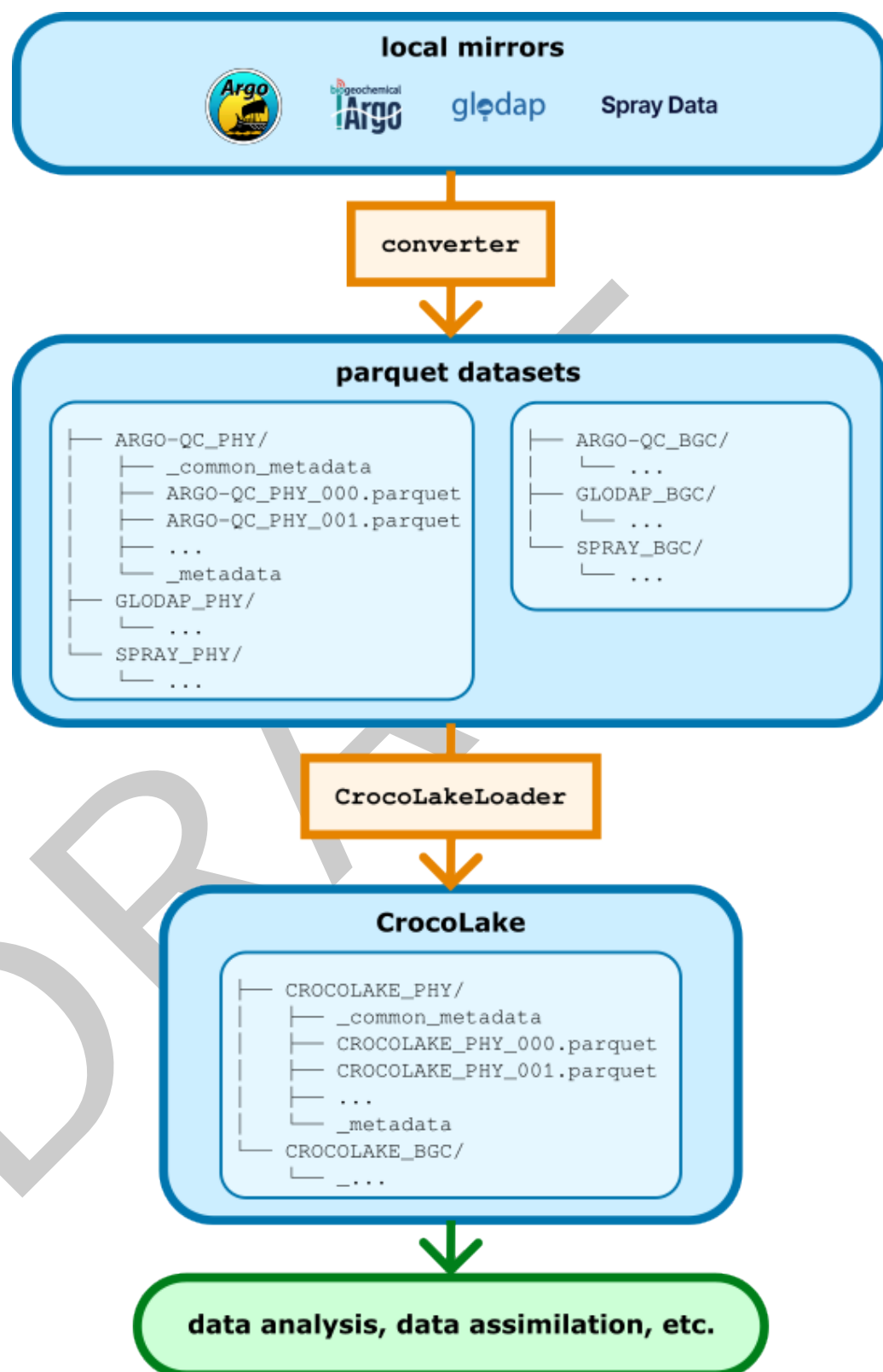


Figure 2: CrocoLake's workflow. 'converter's read the data in their original format, transform it following CrocoLake's conventions, converts it to parquet, and stores it back to disk. Each dataset is converted to its own parquet version. Thanks to the submodule 'CrocoLakeLoader', multiple parquet datasets 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.

76 Schema

77 The nomenclature, units and data types are generally based on Argo's (<https://vocab.nerc.ac.uk/collection/R03/current/>). For CrocoLake's variables that are not present in the
78 Argo program, we provide new names maintaining consistency with Argo's style.
79

80 Profile numbering

81 Ocean data is often accessed by profiles and we provide this functionality for CrocoLake too:
82 the user can retrieve the profiles through the CYCLE_NUMBER variable, which is unique and
83 progressive for each PLATFORM_NUMBER of each subdataset (DB_NAME). As each original product
84 uses its own conventions, the CYCLE_NUMBER of some datasets is generated ad hoc during their
85 conversion if no obvious match with CYCLE_NUMBER exists. The procedure for each dataset is
86 explained in the online documentation.

87 Quality control

88 CrocoLake contains only quality-controlled measurements. We rely exclusively on the quality-
89 controls performed by the data providers and at the time we do not perform any control ourselves
90 (although this might change in the future). Each parameter <PARAM> has a corresponding
91 <PARAM>_QC flag that is generally set to 1 to indicate that the data is reliable. For Argo
92 measurements, the original QC value is preserved, and only measurements with QC values of
93 1, 2, 5, and 8 considered.

94 Measurement errors

95 Each parameter <PARAM> has a corresponding <PARAM>_ERROR that indicates a measurement's
96 error as provided in the original dataset. When no error is provided, <PARAM>_ERROR is set to
97 null.

98 Documentation and updates

99 Documentation is available at (<https://crocolakedocs.readthedocs.io/en/latest/index.html>).
100 It describes the specifics of each dataset (e.g. what quality-control filter we apply to each
101 dataset, the procedure to generate the profile numbers, etc.), and it is updated every time a
102 new feature is available.

103 Citation

104 If you use CrocoLakeTools and/or CrocoLake, do not limit yourself to citing this manuscript
105 and remember to cite the datasets that you have used as indicated in the documentation.
106 For example, if your work relies on Argo measurements, acknowledge Argo (Wong et al.,
107 2020). This is important both for the maintainers of each product to track their impact and
108 to acknowledge their efforts that made your work possible.

109 Acknowledgements

110 We acknowledge funding from [NSF CSSI CROCODILE details]

111 References

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