# Construction of Interactive Websites for Remote Sensing Datasets
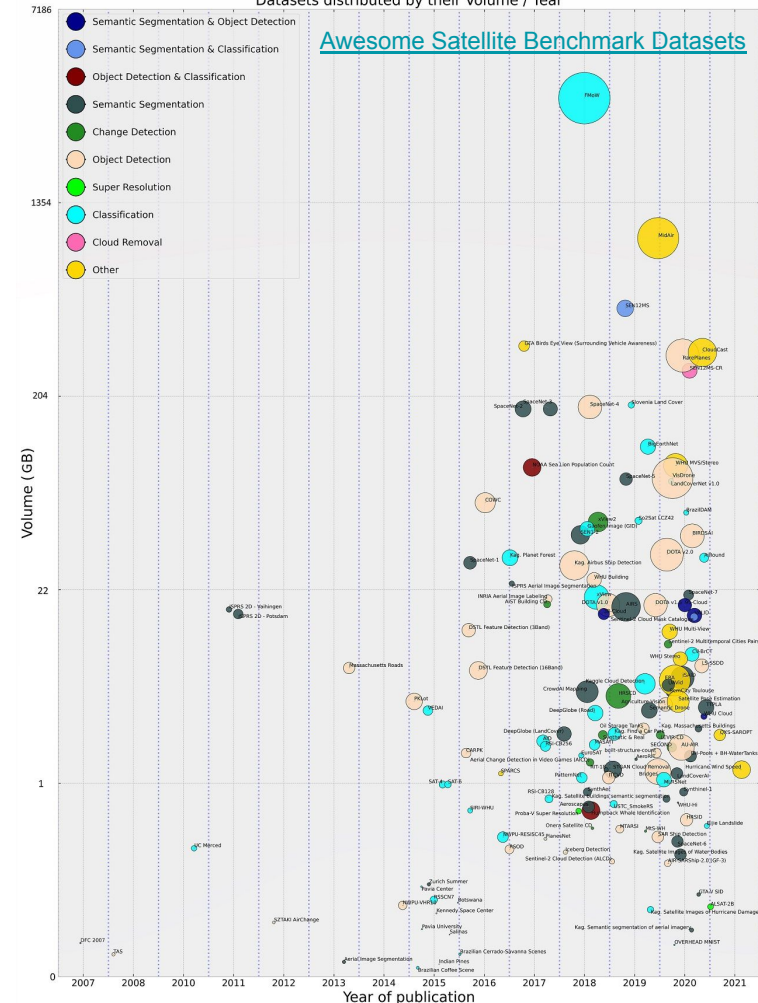
Kai Norman Clasen and Begüm Demir

E-mail: k.clasen@tu-berlin.de

# Remote Sensing Benchmark Datasets

- To train machine learning models large-scale training datasets are required.

- Thus, several well-designed and ready-to-use benchmark datasets have been recently introduced in remote sensing.

- The descriptions of the existing datasets are often published in scientific papers as PDF files.



Datasets distributed by their Volume / Year

Awesome Satellite Benchmark Datasets

# Limitations of Benchmark Dataset PDFs

- **Static Format**: PDF files have no interactive visualization capabilities.

- **Page Limit**: May limit the description of the dataset due to submission guidelines.

- **Hard to update**: Once published, it is difficult to update the dataset paper/description.

- **Limited communication**: Hard to interact with the creators and users of the dataset.

# Solution: Interactive Dataset Websites

- **Aim**: Increase accessibility through:
    - Interactive engagement of the community via commenting systems.

# Solution: Interactive Dataset Websites

- **Aim**: Increase accessibility through:
  - Allow community to
    contribute to the project.

# Solution: Interactive Dataset Websites

- **Aim**: Increase accessibility through:
  - Illustrative and interactive visualizations.



**BigEarthNet-S2**

The general contents of the BigEarthNet-S2 archive looks as follows:

```
BigEarthNet-S2-Example
  S2A_MSIL2A_20170613T101031_87_48
    S2A_MSIL2A_20170613T101031_87_48_B01.tif
    S2A_MSIL2A_20170613T101031_87_48_B02.tif
    S2A_MSIL2A_20170613T101031_87_48_B03.tif
    S2A_MSIL2A_20170613T101031_87_48_B04.tif
    S2A_MSIL2A_20170613T101031_87_48_B05.tif
    S2A_MSIL2A_20170613T101031_87_48_B06.tif
    S2A_MSIL2A_20170613T101031_87_48_B07.tif
    S2A_MSIL2A_20170613T101031_87_48_B08.tif
    S2A_MSIL2A_20170613T101031_87_48_B8A.tif
    S2A_MSIL2A_20170613T101031_87_48_B09.tif
    S2A_MSIL2A_20170613T101031_87_48_B11.tif
    S2A_MSIL2A_20170613T101031_87_48_B12.tif
    S2A_MSIL2A_20170613T101031_87_48_labels_metadata.json
  S2A_MSIL2A_20170617T113321_4_55
    S2A_MSIL2A_20170617T113321_4_55_B01.tif
    S2A_MSIL2A_20170617T113321_4_55_B02.tif
```



RSiM

6

# Solution: Interactive Dataset Websites

- **Aim**: Increase accessibility through:
  - Example code for using, loading and visualizing the data;
  - Providing links to useful tools/libraries to work with the datasets.

## Helpful Libraries

The following is a short list of *unofficial* BigEarthNet-related libraries:

### BigEarthNet Common

The BigEarthNet Common library provides a collection of high-level tools to better work with the BigEarthNet dataset. Use this library to:

- Use any BigEarthNet related constants
  - Quickly print constants by using a CLI tool
- Safely read JSON files
- Deterministically multi-hot encode/decode 19/43-class labels
- Quickly accessing metadata from a patch for filtering
  - Country

```python
import lmdb

import numpy as np

# readahead should be True if dataset fits in RAM
# otherwise it may be faster to set readahead = False
# as readonly=True no need for `locking` which _should_ take longer if lock=True
env = lmdb.open(str(p), readonly=True, readahead=True, lock=False)
# possible optimization use single call to
# getmulti(keys) instead of a new thread with a single element as transaction?

with env.begin() as txn:
    byteflow = txn.get(example_patch.encode("utf-8"))
    s2_patch = BigEarthNet_S2_Patch.loads(byteflow)

bands_10m = s2_patch.get_stacked_10m_bands()
bands_20m = s2_patch.get_stacked_20m_bands()

# interpolate to 10m dimension
```

```python
import matplotlib.pyplot as plt

bands_10m_torch = Tensor(np.float32(bands_10m)).unsqueeze(dim=0)
bands_20m_torch = Tensor(np.float32(bands_20m)).unsqueeze(dim=0)

bands_20m_interp = interpolate(bands_20m_torch, bands_10m.shape[-2:], mode="bicubic")
plt.imshow(bands_20m_interp[0][0], cmap="gray")
plt.title("Torch interpolate (bicubic)")
plt.axis("off");
```

Torch interpolate (bicubic)

# Suggested Workflow with Open-Source Tools

Source

Transform

Publish

**Jupyter Notebook**

- Code cells for examples & to produce images
- Markdown cells for prose

**Sphinx – Python Documentation Generator**

- MyST extension for Markdown support
- MyST-NB for Notebook support

**GitHub**

- GitHub Pages for free website hosting
- GitHub Issues for community engagement & Feedback

Interactive Dataset Websites like docs.kai-tub.tech/ben-docs

=

More accessible and engaging science

SCAN ME

# Acknowledgement

# Additional Material

# Jupyter Project

- A Jupyter Notebook is a web-based interactive computing platform that allows users to create and share documents that combine live code, equations, visualizations, and narrative text.

- Notebooks provide an easy way to prototype, experiment, and iterate on data analysis and machine learning models.

- Additional links:
  - Official Jupyter Project Documentation
  - Try Jupyter on the web without installing anything
  - Example Jupyter Lab Notebook

# Jupyter NB source to HTML page I

## Source File

With the following conventions:

- Each folder corresponds to a single patch
- The `patch_name` is encoded as the name of the folder
- Each patch folder contains a GeoTIFF file for each of the 12 bands.
  - The name of the GeoTIFF file is encoded as `<patch_name>_<band>.tif`.
- The JSON file, named `<patch_name>_labels_metadata.json`, contains the metadata

The prettified contents of a metadata file is:

```python
# remove-input

from rich import print_json
from copy import copy
import json

ben_s2_json_file_paths = list(Path(ben_s2_path).rglob("*.json"))
ben_s2_json_fp = ben_s2_json_file_paths[0]
text = ben_s2_json_fp.read_text()
j = json.loads(text)
simple_j = copy(j)
simple_j["projection"] = simple_j["projection"][:75] + "..."

print_json(data=simple_j)
```

## HTML Page

With the following conventions:

- Each folder corresponds to a single patch
- The `patch_name` is encoded as the name of the folder
- Each patch folder contains a GeoTIFF file for each of the 12 bands.
  - The name of the GeoTIFF file is encoded as `<patch_name>_<band>.tif`.
- The JSON file, named `<patch_name>_labels_metadata.json`, contains the metadata

The prettified contents of a metadata file is:

```
{
  "labels": [
    "Pastures"
  ],
  "coordinates": {
    "ulx": 604800,
    "uly": 5834040,
    "lrx": 606000,
    "lry": 5832840
  },
  "projection": "PROJCS[\"WGS 84 / UTM zone 29N\",GEOGCS[\"WGS 84\",DATUM[\"W
  "tile_source": "S2A_MSIL1C_20170617T113321_N0205_R080_T29UPU_20170617T1133
  "acquisition_date": "2017-06-17 11:33:21"
}
```

**RSiM**

# Jupyter NB source to HTML page II

## Source File

Example output

```
In [ ]:   # scroll-output
          from bigearthnet_gdf_builder.builder import get_gdf_from_s2_patch_dir

          # gdf_builder also has a CLI tool to convert the entire archive into a single
          # parquet file!
          # Example "raw" subset
          gdf = get_gdf_from_s2_patch_dir(ben_s2_path)
          # showing first row as tables have display issues
          gdf
```

Parquet files allow for easy data-processing and visualization. These files work particularly well with geopandas:

## HTML Page

**Example output**

```
from bigearthnet_gdf_builder.builder import get_gdf_from_s2_patch_dir

# gdf_builder also has a CLI tool to convert the entire archive into a single
# parquet file!
# Example "raw" subset
gdf = get_gdf_from_s2_patch_dir(ben_s2_path)
# showing first row as tables have display issues
gdf
```

|   | labels | tile_source | acquisition_date |
|---|---|---|---|
| 0 | [Pastures] | S2A_MSIL1C_20170617T113321_N0205_R080_T29UPU_2... | 2017-06-17 11:33:21 |
| 1 | [Coniferous forest, Mixed forest, Transitional... | S2B_MSIL1C_20170924T93020_N0205_R136_T35VPK_20... | 2017-09-24 09:30:20 |
| 2 | [Non-irrigated arable land, Land principally o... | S2A_MSIL1C_20170613T101031_N0205_R022_T33UUP_2... | 2017-06-13 10:10:31 |
| 3 | [Non-irrigated arable land, Coniferous forest,... | S2B_MSIL1C_20180204T94161_N0206_R036_T35VPK_20... | 2018-02-04 09:41:56 |

# Jupyter NB source to HTML page III

## Source File

:::{note}

The polygons have been merged together to minimize page-load time and storage requirements

:::

```python
# remove-input
import warnings

warnings.filterwarnings("ignore")
import geopandas
import folium

# import folium.plugins


def draw_fast_marker_cluster(gdf):
    marker_gdf = gdf.copy()
    marker_gdf = marker_gdf.to_crs("EPSG:4326")
    m = folium.Map(tiles="Stamen Terrain")
    data = marker_gdf.representative_point().apply(lambda point: [point.y, point.x])
    map_data = folium.plugins.FastMarkerCluster(data)
    m.add_child(map_data)
    return m


def simplify_gdf(gdf, tolerance=100):
    geo_series = gdf.geometry.unary_union
    g_series_simp = geo_series.simplify(tolerance)
    return geopandas.GeoDataFrame(
        {"name": ["BigEarthNet-simplified"]}, geometry=[g_series_simp], crs=gdf.crs
    )


# gdf = geopandas.read_parquet("../gdf/raw_ben_s2_gdf_3035.parquet")
# gdf_simp = simplify_gdf(gdf)
# gdf_simp.to_parquet("_static/ben_simple_union.parquet")

gdf = geopandas.read_parquet("_static/ben_simple_union.parquet")
```
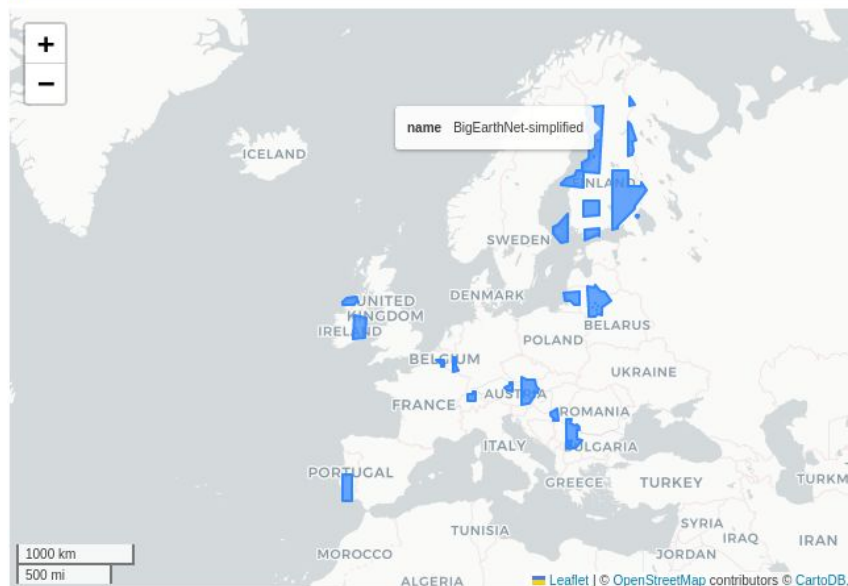
## HTML Page

# Executable Books Project

- The Executable Books Project is an open-source project that aims to improve the sustainability of scientific research by developing tools that facilitate publishing computational narratives using the Jupyter ecosystem, such as:
    - Jupyter Books
    - MyST-NB project

- Additional links:
    - Official Executable Books Documentation
    - Jupyter Book Gallery
    - MyST (Markedly Structured Text) – A superset of the CommonMark language

- Sphinx is an open-source documentation generator that:
  - is widely used in the Python community, but can also be used for other programming languages;
  - is used as the foundation for many Executable Book Projects/Tools.

- Additional Links:
  - [Official Sphinx Documentation](#)
  - [Furo – A popular Sphinx Theme](#) (also used for the BigEarthNet Guide)
  - [Using Markdown (MyST) in Sphinx](#) instead of *[reStructured Text (reST)](#)*

# GitHub Pages

- *GitHub Pages* is a free web hosting service provided by GitHub that:
    - allows users to create static websites and host them on GitHub's servers;
    - supports a variety of static site generators, such as Sphinx and Jekyll

- Additional links:
    - [Official GitHub Pages Documentation](#)
    - [Using a custom domain for your GitHub Pages site](#)
    - [GitHub Pages usage limits](#)