
New Zealand Open Environmental Science Data sets

Albert Bifet
University of Waikato

Jacob Montiel
University of Waikato

Nick Lim
University of Waikato

Gregory Pearson
Metservice

Sebastian Delaux
Metocean

James McLeod
Waikato Regional Council

Phil Mourot
Waikato Regional Council

Abstract

1 Data Science on environmental spatio-temporal data is becoming a critical and
2 challenging research topic due to the changing nature and rapidly increasing volume
3 of available data. To this end, we would like to introduce TAIAO data repository,
4 comprising of over 30 datasets of various types including images, videos, textual
5 and tabular data.

6 1 Introduction

7 Environmental data science is strategically essential to New Zealand because it supports, leverages
8 and boosts research on climate change impacts, adaptation and conservation [6]. Effective data
9 science can take an essential role in the Government’s goals of improving the quality of freshwater
10 and to reach zero carbon by 2050. Environmental time series or data streams are found in many
11 practical applications in New Zealand. They can consist of monitoring observations or modelling
12 output of flow (e.g. wind, current, water level, ice flow, ice height), concentration (e.g. suspended
13 sediment, nutrients, contaminants), physical properties (e.g. temperature, density) and external
14 forcing (e.g. gravity, solar radiation).

15 Data Science on environmental spatio-temporal data is becoming a critical and challenging research
16 topic due to the changing nature and rapidly increasing volume of available data. New methods
17 are also required for automatic routine monitoring of biological variables (e.g. birdsong listening
18 stations, automated tree-ring measurements). Environmental time series data need specific processing
19 techniques because:

- 20 • **Decisions** (for example, with respect to management regimes or policies) are made over
21 time on the basis of partial information, and we do not have the time to collect perfect and
22 complete datasets;
- 23 • the properties of the information are likely to evolve over time (concept drift), violating the
24 assumptions of some standard statistical approaches;
- 25 • the information has a history that is difficult to delimit, yet incorporating history can
26 substantially improve predictive power; and
- 27 • the information can be multi-scale, ranging from broadscale satellite-derived data to
28 irregularly-spaced point measurements (e.g of temperature, wind velocity, water flow);

29 Remotely-sensed environmental data are often taken as archetypes of big data because they exhibit
30 three key properties:

- 31 • **Volume:** with multiple constellations of satellites offering daily global coverage at sub-metre
32 resolutions (e.g. over 25 PB of data per day from the ESA Sentinel satellites alone [5]), or
33 geostationary satellites such as Himawari 8, and airborne and terrestrial datasets such as
34 LiDAR and sonar datasets that capture centimetre resolution topographic information over
35 wide areas, the volume of remotely-sensed data is enormous;
- 36 • **Variety:** earth observation data are multimodal, comprising observations by both active and
37 passive sensors and across the electro-magnetic spectrum, and packaged in multiple formats
38 (raster/vector, structured/unstructured), but are typically sampled unselectively presenting
39 major challenges for pattern recognition and interpretation; and
- 40 • **Velocity:** image cadence has increased dramatically (e.g. the Himawari 8 satellite over NZ
41 produces images every 10 minutes), and growing archives of archival image stacks that are
42 ripe for temporal analysis [4].

43 **2 Time-Evolving Data Science / Artificial Intelligence for Advanced Open** 44 **Environmental Science (TAIAO)**

45 TAIAO[2] is a New Zealand government supported, multi-year, multi-million dollars, programme
46 aimed at improving the data capabilities of researchers in New Zealand. TAIAO was launched in
47 2020 and is currently in the early stages of the development of the platform. The motivation of
48 the TAIAO programme includes advancing the state-of-the-art in environmental data science by
49 developing new machine learning methods for time series and data streams with the capacity to deal
50 with large quantities of big data in real-time, with special emphasis on processing the data collected
51 on the New Zealand environment. TAIAO also aim to build an open-source framework to implement
52 machine learning on time series data, as well as provide an open available repository with datasets to
53 improve reproducibility in environmental data science. Ultimately, TAIAO aim to democratise and
54 build capability in fundamental and applied data science.

55 This programme is a multi-institute, multi-domain and includes data scientists, data engineers,
56 environmental scientists, and machine learning researchers from undergraduate to post-graduate level.
57 Moreover, collaboration is expected to extend beyond technical aspects to include regional councils,
58 iwi¹ and co-governance entities to implement the methods we develop to support governance and
59 management decisions with analyses based on large volumes of data that they cannot currently
60 process.

61 **2.1 Reproducible Notebooks**

62 TAIAO uses Jupyter Notebooks to document and visualize the codes for better reproducibility and
63 documentation. Each notebook is associated to a task and describe how the data can be accessed. The
64 notebooks also documents the application of the dataset as well as the questions the data provider and
65 researcher seek to answer with the data. As part of TAIAO's commitment to open-source platforms,
66 we use Jupyter notebooks as it is open-sourced, light-weight while being capable. We envision
67 that as more notebooks in the platform are developed, we can improve the transparency and the
68 reproducibility of the findings, as well as improve the accessibility of data science research. Figure 1
69 describes a conceptual view of the platform and how the corresponding components of the platform
70 interact.

71 **2.2 Indigenous Data Science**

72 TAIAO and the New Zealand government is committed to “Vision Mātuaranga” [7] which aims to
73 unlock the potential of traditional indigenous knowledge and recognize the value of the generations

¹largest social units in Aotearoa Māori society.

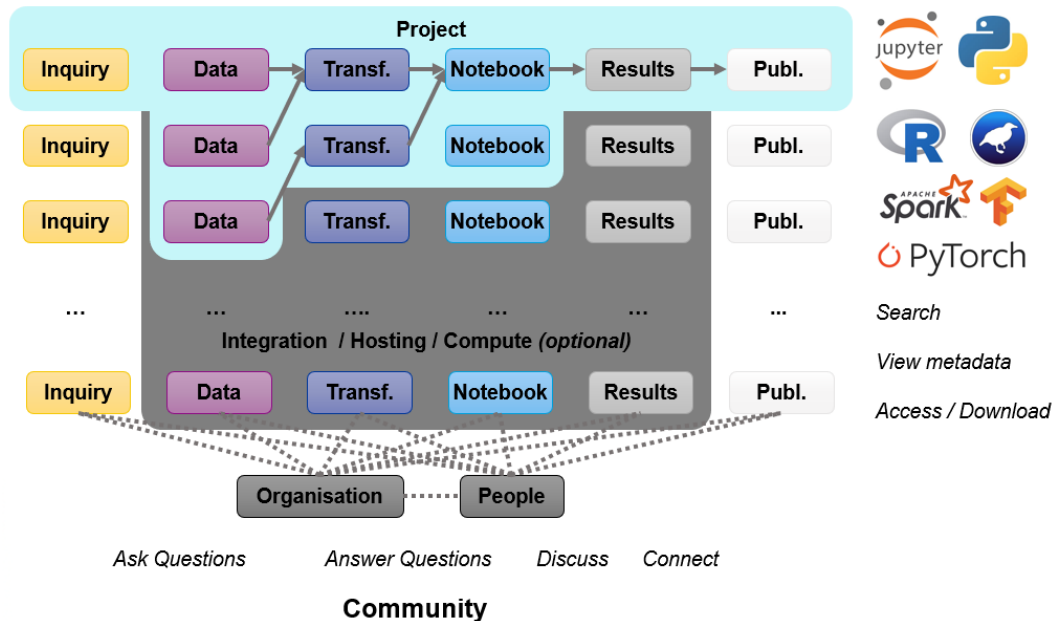


Figure 1: Conceptual view of the platform and the interaction of the corresponding components the TAI AO platform

74 of tradition and knowledge passed down through the people living in the land. As part of this
 75 commitment, TAI AO recognizes the community as partners in science and innovation and guardians
 76 of the natural resources and indigenous knowledge.

77 TAI AO works closely with the local co-governance entities and the iwi to understand the interests and
 78 issues of the community. In addition, TAI AO also hold regular dialogues and meetings to review the
 79 environmental findings and review the capability developed as well as to democratize the platform
 80 and to build the data science capabilities of the community.

81 In addition to that, TAI AO is also committed to principals of “Te Mana Raruanga” [1] which
 82 recognizes Indigenous data rights and sovereignty. Here, TAI AO recognizes the rights and ownership
 83 of the data and are have regular dialogue with the owners of the data regarding how the data is being
 84 used and the rights and accessibility to the data.

85 3 Available datasets

86 Over the last year (2020), we have been building a repository of environmental data and hosting the
 87 data within the TAI AO platform. Currently, the TAI AO environmental dataset repository currently
 88 contains a collection of over 30 datasets of various type including images, videos, textual data, and
 89 multi-variate tabular data

90 Table 1 is a summarises of some of the datasets in the platform. Note that additional datasets are still
 91 being added and the list in the table may not include some of the more recent datasets.

Table 1: Summary of the datasets in the TAI AO data repository

Name	Description	Category/Tag	Type
Aggregated wave and atmospheric forecast derived from GFS guidance	Hindcast and forecast wave and atmosphere model	Model, Forecast, Sea, Wind, Wave	Tabular

Continued on next page

Table 1 – continued from previous page

Name	Description	Category/Tag	Type
Ardmore airport automatic weather station	Nowcast percipitation, temperature and barometric, updated every 7.5 minutes	Observation, Land, Wind, Temperature	Textual
Ashburton Aerodrome automatic weather station	Nowcast percipitation, temperature and barometric, updated every 7.5 minutes	Observation, Land, Wind, Precipitation, Temperature	Textual
Auckland airport automatic weather station	Nowcast percipitation, temperature and barometric, updated every 7.5 minutes	Observation, Land, Wind, Precipitation, Temperature	Textual
Auckland-Hamilton Corridor Aerial Photography	Aerial photograph on the Auckland-Hamilton Corridor region	Observation Image	Image
Birchwood automatic weather station	Nowcast percipitation, temperature and barometric, updated every 7.5 minutes	Observation, Land, Wind, Precipitation, Temperature	Textual
Blitzen lightning feed	Archive of lightning records and live feed	Observation, lightning, location, time, intensity, type	Tabular
Campbell Island Wave buoy	Nowcast Wave heighth, Wave direction, Wave period	Observation, Sea, Wave, Buoy, Wave heighth, Wave direction, Wave period	Tabular
Coromandel river + rain gauge time series	Time Series of rivers and ran gauages in Coromandel Peninsular dated from 2010-2020 at 5 minutes resolution	Observation, time, intensity	Tabular
Flat Hills Automatic weather station	Nowcast percipitation, temperature and barometric, updated every 7.5 minutes	Observation, Land, Wind, Precipitation, Temperature	Textual
Global topography (elevation, slope, and aspect)	Topography metrics; ARD tile format. Extracted from the NASA Making Earth System Data Records for Use in Research Environments (MEaSUREs) version of the Shuttle Radar Topography Mission (SRTM) global elevation dataset	Land, topography	Image
Google Earth Engine global geospatial data sets	Planetary scale satellite imagery and derived data products	Cimate, weather, geophysical data products, satellite imagery (remote sensing)	Image
GPATS Oceania lightning feed	Archive of lightning records and live feed for Oceania area	Observation, lightning, location, time, intensity, type	Tabular

Continued on next page

Table 1 – continued from previous page

Name	Description	Category/Tag	Type
Haast automatic weather station	Nowcast percipitation, temperature and barometric, updated every 7.5 minutes	Observation, Land, Wind, Precipitation, Temperature	Textual
Himawari-8 2km half disc archive	Archive of himawari-8 satellite data	Observation, Remote-sensing, visible, in-fared	Image
Himawari-8 500m channel 3 half disc archive	Archive of full resolution channel 3 data from Himawari-8 satellite	Observation, Remote-sensing, visible	Image
Himawari-8 AWS NOAA archive	Full resolution archive of Himawari-8 satellite (data from end of 2019 only)	Observation, Remote-sensing, visible	Image
Landsat 8 remote sensing data	Global analysis ready multispectral satellite data from Landsat 8 OLI sensor	Land, remote sensing, multi-spectral	Image
LILA Wellington Camera Trap	270450 Images of wildlife from 187 camera traps locations[3]	Observation Image	Image
Moana New Zealand Hydrodynamics Re-analysis v1.9	Hydrodynamic reanalysis of new zealand waters.	Model, Hind-cast, Sea, Current, SST	Tabular
Mt Karioi predator camera video feed	2101 videos of wildlife from 20 camera trap locations. Note that this dataset is available only on request and with permission from the Hapu data owners	Observation, Video, Image	Video
Mt Karioi predator trap logs	Table of status of traps, including description of bait used, deployment date and date last checked	Observation, Manual logging	Tabular, Textual
NZ rain radar archive - RAW data	Archive of data from MetServices doppler radar network.	Observation, Atmosphere, Reflectivity	Image
Regional Council Water quality and Discharge data		Water, rainfall	Tabular
Sentinel 1/2 snapshot of waikato region	Hyperspectral satellite image from Sentinel 1/2	Observation Image	Image
Southern Ocean Waverider buoy	Nowcast Wave heigth, Wave direction, Wave period	Observation, Sea, Wave, Buoy, Wave heigth, Wave direction, Wave period	Tabular
TOA lightning feed	Archive of lightning records and live feed for Oceania	Observation, lightning, location, time, intensity, type	Tabular
Tropical cyclone archive	Archive of tropical cyclone trajectory	Observation, cyclone, trajectories	Tabular

Continued on next page

Table 1 – continued from previous page

Name	Description	Category/Tag	Type
Horizons Air Quality	PM10 and PM100 particulate data taken at 5 minute samples	Air Quality	Tabular
Hawke’s Bay Air Quality	PM10 and PM100 particulate data taken at 5 minute samples	Air Quality	Tabular
Hawke’s Bay Air Quality - Raw	PM10 and PM100 particulate data taken at 5 minute samples	Air Quality	Tabular
Waikato Region Aerial Photography	Orthorectilinearized projection of aerial photography of the Waikato region taken at 0.03m resolution	Observation Image	Image
GFS	Historical percipitation forecast of the Coromandel region from 2015-2018 [8]	Water, Rain-fall, Climate	Tabular

Table 1: Summary of the datasets in the TAI AO data repository

92 3.1 Strategic Use of Big Data Sets

93 Unlike other centres of data science expertise such as Europe, North America and Asia, New Zealand
 94 is an island with a low population density and a high level of urbanisation. Its weather changes quickly
 95 and in ways that are difficult to predict, and it has a low density of on-the-ground environmental data
 96 measurements (very low in surrounding oceans and on land away from population centres), so it
 97 relies heavily on satellite measurements and numerical modelling predictions, and has to combine
 98 broad-scale satellite data with sparse on-the-ground data. The TAI AO project particularly focuses on
 99 the challenges of making that combination effective.

100 Moreover, NZ climate is at the interface between tropical and polar air masses, and its coast is
 101 connected to the Southern and Pacific Oceans and the Tasman Sea, which is an ocean-warming
 102 hotspot. There is a need for fit-for purpose tools that are particularly tailored to this complex
 103 environment since existing methods developed overseas are often not suitable for nor transferable to
 104 New Zealand conditions.

105 Topographic datasets are also relatively sparse compared to more populated land-masses; this gap
 106 is being corrected by (for example) the national LiDAR survey that is underway, funded by the
 107 Provincial Growth Fund and regional councils², but integrating the LiDAR data with other spatial
 108 data present challenges that the TAI AO project is well-placed to address. To meet these challenges,
 109 the TAI AO project will use large-scale datasets facilitated by research partner MetService and build on
 110 ongoing work of the environmental scientists within its team. The goal is to build on state-of-the-art
 111 modelling datasets of coastal ocean circulation, connectivity and marine temperature being developed
 112 in MetService’s Moana project³ and weather radar images, weather station data and high-resolution
 113 satellite imagery archived by MetService. It is expected to use existing and newly-acquired LiDAR
 114 datasets.

115 Going beyond physical data to biological data, a joint project between the University of Waikato and
 116 Xerra uses estuarine colour indices from Sentinel II satellites to detect estuarine ecosystem tipping
 117 points⁴.

118 4 Future plans and capabilities

119 TAI AO is currently in the early stages of development, and over the next few years, we plan to
 120 improve the platform to better index the datasets. We also plan to increase the number of notebooks
 121 as well as to include more complex, cross-domain, multi-dataset examples showing the application of

²<https://www.linz.govt.nz/data/linz-data/elevation-data>

³<https://www.moanaproject.org/>

⁴<https://www.xerra.nz/2019/06/11/calibrating-satellite-imagery-using-ground-based-data-collection/>

122 cross-domain datasets. Additionally, we plan to improve the process of adding additional datasets
123 and notebooks to improve the accessibility and contribution from the TAI AO community

124 **5 Conclusion**

125 While the TAI AO project is relatively young, we have compiled a repository of varied and unique
126 datasets that are pertinent to environmental research. We are confident that the TAI AO project can
127 improve the accessibility of data science research especially in the field of environmental science

128 **References**

- 129 [1] Principles of māori data sovereignty. Retrieved from [https://www.temanararaunga.maori.](https://www.temanararaunga.maori.nz/nga-rauemi)
130 [nz/nga-rauemi](https://www.temanararaunga.maori.nz/nga-rauemi), 2018.
- 131 [2] Taiao, time-evolving data science / artificial intelligence for advanced open environmental science.
132 Retrieved from <https://taiao.ai>, 2020.
- 133 [3] V. Anton, S. Hartley, A. Geldenhuis, and H. U. Wittmer. Monitoring the mammalian fauna of
134 urban areas using remote cameras and citizen science. *Journal of Urban Ecology*, 4(1):juy002,
135 2018.
- 136 [4] A. Bifet and J. Read. Ubiquitous artificial intelligence and dynamic data streams. In *Proceedings*
137 *of the 12th ACM International Conference on Distributed and Event-based Systems*, pages 1–6,
138 2018.
- 139 [5] W. Fan and A. Bifet. Mining big data: Current status, and forecast to the future. *SIGKDD Explor.*
140 *NewsL.*, 14(2):1–5, Apr. 2013.
- 141 [6] K. Gibert, J. S. Horsburgh, I. N. Athanasiadis, and G. Holmes. Environmental data science.
142 *Environmental Modelling & Software*, 106:4–12, 2018. Special Issue on Environmental Data
143 Science. Applications to Air quality and Water cycle.
- 144 [7] L. H. Kaiser and W. S. A. Saunders. Vision mātauranga research directions: opportunities for
145 iwi and hapū management plans. *Kōtuitui: New Zealand Journal of Social Sciences Online*,
146 0(0):1–13, 2021.
- 147 [8] National Centers for Environmental Prediction, National Weather Service, NOAA, U.S. Depart-
148 ment of Commerce. Ncep gfs 0.25 degree global forecast grids historical archive, 2015.