

Unsupervised learning in in shallow marine environments using satellite imagery

Zayad AlZayer^{1,†,‡} , Cedric John^{2,‡} and Philippa Mason^{2,*}

¹ Imperial College 1; zba21@ic.ac.uk

² Affiliation 2; e-mail@e-mail.com

* Correspondence: e-mail@e-mail.com; Tel.: (optional; include country code; if there are multiple corresponding authors, add author initials) +xx-xxxx-xxx-xxxx (F.L.)

† Current address: Affiliation 3.

‡ These authors contributed equally to this work.

Abstract: A single paragraph of about 200 words maximum. For research articles, abstracts should give a pertinent overview of the work. We strongly encourage authors to use the following style of structured abstracts, but without headings: (1) Background: place the question addressed in a broad context and highlight the purpose of the study; (2) Methods: describe briefly the main methods or treatments applied; (3) Results: summarize the article's main findings; (4) Conclusions: indicate the main conclusions or interpretations. The abstract should be an objective representation of the article, it must not contain results which are not presented and substantiated in the main text and should not exaggerate the main conclusions.

Keywords: keyword 1; keyword 2; keyword 3 (List three to ten pertinent keywords specific to the article; yet reasonably common within the subject discipline.)

1. Introduction

Over the past 40 years, the crucial yet challenging task of monitoring coral reefs has been undertaken, with data gathering initiatives tracing back to as early as the 1960s [1], and more comprehensive databases from the 1980s to 2022 [2], as well as citizen science datasets [3]. With these data sets encompassing sub-mapping scale information, remote sensing studies encompass a broad range of objectives, from local ecological surveillance to tracking carbon budgets [4]. In light of the threats imposed by climate change and anthropogenic activities [5], and the rapid temperature rise that has led to a reduction in both coral cover and diversity [6], there exists an immediate and pressing need for accurate and swift global coral reef monitoring and data fusion techniques.

[?] provide a comprehensive overview of sensor limitations and uses

Much of the recent research has centered around supervised learning algorithms, a form of machine learning that utilizes labeled data to train a model, thereby enabling it to predict labels for new data. However, this approach often entails certain assumptions about the labeled data, including a uniform quality of labels among all labelers [7], thereby necessitating expert verification. This methodology has been applied to categorize images of coral reefs into various classes such as coral, sand, algae, and rubble [8]. Nevertheless, such a process requires labeled data, which can be challenging to procure and process. Moreover, it can be outright impossible in cases dealing with historical satellite imagery, where the ground truth may not always be accessible in an environment that is living and adapting.

Unsupervised learning is a type of machine learning that uses unlabelled data to train a model to find patterns in the data itself, helping unlock bottlenecks that exist within labelled data [9].

In this study we aim to use a combination of unsupervised and supervised learning to classify coral reefs into different classes. Using a combination of more traditional clustering

Citation: AlZayer, Z.; Mason, P.; John, C. Title. *Journal Not Specified* **2023**, *1*, 0. <https://doi.org/>

Received:

Revised:

Accepted:

Published:

Copyright: © 2023 by the authors. Submitted to *Journal Not Specified* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

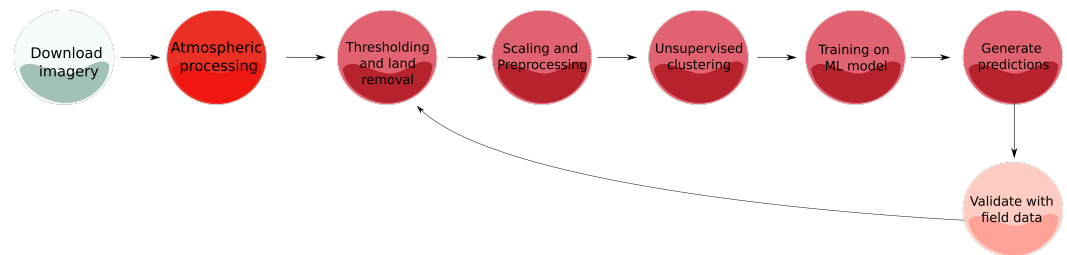


Figure 1. Basic overall workflow in the study of coral reefs using Sentinel-2 imagery

methods and various color spaces. We then use a supervised learning algorithms to provide additional insight into the clustering and retrieve understandable results from the data and its clusters, including using simple logistic regression to gain insight into the data itself.

Sea roughness due to wind is also a problem as very large changes in reflectance such as sun glint also affect the imagery negatively and should generally be discarded and or masked out [10]

2. Methods

In this study we setup a series of systematic experiments using the workflow described briefly below.

- Data collection and preprocessing: We gather Sentinel-2 L1C data using the API provided by the Copernicus Open Access Hub. We then preprocess the data to remove clouds and other noise. We then use the data to create a time series of images for each location. We then use the time series to create an image stack for each location.
- Stack processing: For each stack we remove the land using a combination of band 8 and 11 to create a mask. We then use the mask to remove the land from the stack. We also include additional features such as NDCI, BGR and and pseudo-bathymetry
- Unsupervised learning: We then unstack the array to create individual data points of each pixel. We then use a combination of clustering methods to cluster the data points into different classes.
- Supervised learning: We then use a combination of supervised learning methods to classify the data points into the data classes previously defined by the clustering algorithm.

Data collection, preprocessing and Stacking

Image correction for L1C data was done using the sen2cor processor for all the imagery to remove atmospheric effects. This was followed by cloud masking using the Fmask algorithm [11]. The Fmask algorithm uses the blue, red, near-infrared and shortwave infrared bands to create a cloud mask. The cloud mask is then used to remove the clouds from the image. The Fmask algorithm was chosen as it is a widely used and tested algorithm for cloud masking Sentinel-2 imagery and it also provides a mask for water, a ratio of water to clouds was used to filter out the imagery which resulted in a relatively cloud free dataset for the study area (Lizard Island Australia) with a total of 56 images that were cloud free, one cloudy image was also included in the dataset in order to cover a wider variety of data in the training set.

Color correction and color spaces

As bands 4,3,2 (centered at 664, 559 and 492 nm) roughly represent the RGB color space, we use this as a starting point for our color correction. We then use the following color spaces to create additional features for the data, from the color spaces examined tests were run on the LAB [12], HSV and HSI [13] color spaces respectively. This was done to preserve the overall color scheme and ensure the images are stretched correctly.

The images were then stacked into time series for the lizard island location, and feature generation for each individual time slice was done using the following indices:

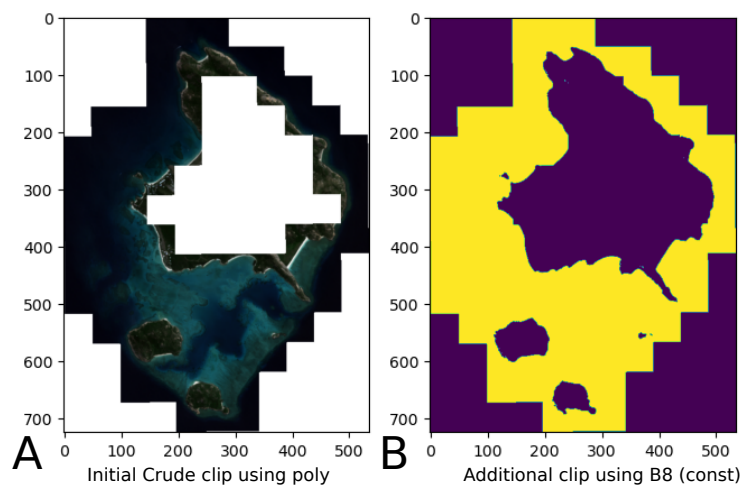


Figure 2. Clipping Image data with NIR mask, (A) showing original clipped Image in dataset (B) showing clip mask with specified threshold

- Chlorophyll Index (CI): Used to estimate chlorophyll content in vegetation. This information can give insights into the health and vigor of plants. 78
- Ocean Color Index (OCI): Used to assess ocean color properties, particularly the presence of chlorophyll. This index can help in studying phytoplankton abundance and water quality in marine environments. 79
- Suspended Sediment Index (SSI): Used to estimate the concentration of suspended sediments in water bodies. This index is helpful in monitoring water quality, sediment transport, and erosion processes. 80
- Turbidity Index (TI): Used to estimate the turbidity in water bodies. Like the SSI, this index is also useful in monitoring water quality, sediment transport, and erosion processes. 81
- Water Quality Index (WQI): Used to assess water quality based on multiple parameters. It provides a comprehensive measure of water health, considering the contributions of various spectral bands to the index computation. 82
- Normalized Difference Chlorophyll Index (NDCI): Used to estimate chlorophyll content in vegetation. The NDCI provides a normalized measure of the difference between green reflectance and red-edge reflectance, indicating vegetation health. 83
- Blue to Green Ratio (BGR): Used to assess water quality by comparing the blue and green reflectance values. This index provides information about the concentration of chlorophyll and suspended sediments in water bodies. 84
- In addition to these indices, the code contains a function for masking out land areas in an image (`mask_land`) using the NIR band and threshold, generally named the black pixel approximation [14]. 85

Resulting in a total of approximately 1 million unique data points covering the range of the time series containing the original 13 bands and 7 additional features. 86

Unsupervised learning 87

These are then entered into 3 dimensionality reduction algorithms, PCA [15], t-SNE [16] and UMAP [17]. These dimensionality reduction algorithms are then used to reduce the dimensionality of the data to 2 dimensions. These are then used to cluster the data using a combination of K-means, DBSCAN and HDBSCAN. The clusters retrieved from these algorithms are then visualised and analysed to create psuedo-labels using k-means [18] and gaussian mixture models were also tested [19]. After hyper parameter tuning and optimisation, we then use these as labels for the supervised learning algorithm classifier which provides additional scope for creating probability maps and testing the accuracy of the clusters themselves. 88

3. Results

This section may be divided by subheadings. It should provide a concise and precise description of the experimental results, their interpretation as well as the experimental conclusions that can be drawn.

4. Discussion

In this study we show that it is possible to achieve repeatable results using a combination of unsupervised and supervised learning methods to classify shallow water imagery in Sentinel-2 imagery. We also show that it is possible to use these methods to create a probability map of a variety of shallow water classes (will expand this) with minimal preprocessing to cluster various times in different scenes and that this methodology can be applied to different areas whilst using simple explainable algorithms.

4.1. Conclusions

4.1.1. Subsubsection

Bulleted lists look like this:

- First bullet;
- Second bullet;
- Third bullet.

Numbered lists can be added as follows:

1. First item;
2. Second item;
3. Third item.

The text continues here.

4.2. Figures, Tables and Schemes

All figures and tables should be cited in the main text as Figure 3, Table 1, etc.



Figure 3. This is a figure. Schemes follow the same formatting. If there are multiple panels, they should be listed as: (a) Description of what is contained in the first panel. (b) Description of what is contained in the second panel. Figures should be placed in the main text near to the first time they are cited. A caption on a single line should be centered.

Table 1. This is a table caption. Tables should be placed in the main text near to the first time they are cited.

Title 1	Title 2	Title 3
Entry 1	Data	Data
Entry 2	Data	Data ¹

¹ Tables may have a footer.

The text continues here (Figure 4 and Table 2).

137



Figure 4. This is a wide figure.

Table 2. This is a wide table.

Title 1	Title 2	Title 3	Title 4
Entry 1 *	Data	Data	Data
	Data	Data	Data
	Data	Data	Data
Entry 2	Data	Data	Data
	Data	Data	Data
	Data	Data	Data
Entry 3	Data	Data	Data
	Data	Data	Data
	Data	Data	Data
Entry 4	Data	Data	Data
	Data	Data	Data
	Data	Data	Data

* Tables may have a footer.

Text.
Text.

138

139

4.3. Formatting of Mathematical Components140

This is the example 1 of equation:141

$$a = 1,$$
(1)

the text following an equation need not be a new paragraph. Please punctuate equations as142
regular text.143

This is the example 2 of equation:144

$$a = b + c + d + e + f + g + h + i + j + k + l + m + n + o + p + q + r + s + t + u + v + w + x + y + z$$
(2)

Please punctuate equations as regular text. Theorem-type environments (including145
propositions, lemmas, corollaries etc.) can be formatted as follows:146

Theorem 1. *Example text of a theorem.*147

The text continues here. Proofs must be formatted as follows:148

Proof of Theorem 1. Text of the proof. Note that the phrase “of Theorem 1” is optional if it149
is clear which theorem is being referred to. □150

The text continues here.151

5. Discussion152

Authors should discuss the results and how they can be interpreted from the perspec-153
tive of previous studies and of the working hypotheses. The findings and their implications154
should be discussed in the broadest context possible. Future research directions may also155
be highlighted.156

6. Conclusions157

This section is not mandatory, but can be added to the manuscript if the discussion is158
unusually long or complex.159

7. Patents160

This section is not mandatory, but may be added if there are patents resulting from the161
work reported in this manuscript.162

Author Contributions: For research articles with several authors, a short paragraph specifying their163
individual contributions must be provided. The following statements should be used “Conceptualiza-164
tion, X.X. and Y.Y.; methodology, X.X.; software, X.X.; validation, X.X., Y.Y. and Z.Z.; formal analysis,165
X.X.; investigation, X.X.; resources, X.X.; data curation, X.X.; writing—original draft preparation,166
X.X.; writing—review and editing, X.X.; visualization, X.X.; supervision, X.X.; project administration,167
X.X.; funding acquisition, Y.Y. All authors have read and agreed to the published version of the168
manuscript.”, please turn to the [CRediT taxonomy](#) for the term explanation. Authorship must be169
limited to those who have contributed substantially to the work reported.170

Funding: Please add: “This research received no external funding” or “This research was funded171
by NAME OF FUNDER grant number XXX.” and and “The APC was funded by XXX”. Check172
carefully that the details given are accurate and use the standard spelling of funding agency names at173
<https://search.crossref.org/funding>, any errors may affect your future funding.174

Institutional Review Board Statement: In this section, you should add the Institutional Review175
Board Statement and approval number, if relevant to your study. You might choose to exclude this176
statement if the study did not require ethical approval. Please note that the Editorial Office might ask177
you for further information. Please add “The study was conducted in accordance with the Declaration178
of Helsinki, and approved by the Institutional Review Board (or Ethics Committee) of NAME OF179
INSTITUTE (protocol code XXX and date of approval).” for studies involving humans. OR “The180
animal study protocol was approved by the Institutional Review Board (or Ethics Committee) of181

NAME OF INSTITUTE (protocol code XXX and date of approval).” for studies involving animals. OR
“Ethical review and approval were waived for this study due to REASON (please provide a detailed
justification).” OR “Not applicable” for studies not involving humans or animals.

Informed Consent Statement: Any research article describing a study involving humans should
contain this statement. Please add “Informed consent was obtained from all subjects involved in the
study.” OR “Patient consent was waived due to REASON (please provide a detailed justification).”
OR “Not applicable” for studies not involving humans. You might also choose to exclude this
statement if the study did not involve humans.

Written informed consent for publication must be obtained from participating patients who can
be identified (including by the patients themselves). Please state “Written informed consent has been
obtained from the patient(s) to publish this paper” if applicable.

Data Availability Statement: We encourage all authors of articles published in MDPI journals to
share their research data. In this section, please provide details regarding where data supporting
reported results can be found, including links to publicly archived datasets analyzed or generated
during the study. Where no new data were created, or where data is unavailable due to privacy or
ethical restrictions, a statement is still required. Suggested Data Availability Statements are available
in section “MDPI Research Data Policies” at <https://www.mdpi.com/ethics>.

Acknowledgments: In this section you can acknowledge any support given which is not covered by
the author contribution or funding sections. This may include administrative and technical support,
or donations in kind (e.g., materials used for experiments).

Conflicts of Interest: Declare conflicts of interest or state “The authors declare no conflict of interest.”
Authors must identify and declare any personal circumstances or interest that may be perceived as
inappropriately influencing the representation or interpretation of reported research results. Any role
of the funders in the design of the study; in the collection, analyses or interpretation of data; in the
writing of the manuscript; or in the decision to publish the results must be declared in this section. If
there is no role, please state “The funders had no role in the design of the study; in the collection,
analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the
results”.

Sample Availability: Samples of the compounds ... are available from the authors.

Abbreviations

The following abbreviations are used in this manuscript:

MDPI	Multidisciplinary Digital Publishing Institute	
DOAJ	Directory of open access journals	
TLA	Three letter acronym	
LD	Linear dichroism	

Appendix A

Appendix A.1

The appendix is an optional section that can contain details and data supplemental to
the main text—for example, explanations of experimental details that would disrupt the
flow of the main text but nonetheless remain crucial to understanding and reproducing
the research shown; figures of replicates for experiments of which representative data are
shown in the main text can be added here if brief, or as Supplementary Data. Mathematical
proofs of results not central to the paper can be added as an appendix.

Table A1. This is a table caption.

Title 1	Title 2	Title 3
Entry 1	Data	Data
Entry 2	Data	Data

Appendix B

All appendix sections must be cited in the main text. In the appendices, Figures, Tables, etc. should be labeled, starting with “A”—e.g., Figure A1, Figure A2, etc.

References

1. Goreau, T.F. Mass expulsion of zooxanthellae from Jamaican reef communities after Hurricane Flora. *Science* **1964**, *145*, 383–386.
2. van Woesik, R.; Kratochwill, C. A global coral-bleaching database, 1980–2020. *Scientific Data* **2022**, *9*. <https://doi.org/10.1038/s41597-022-01121-y>.
3. Belbin, L.; Wallis, E.; Hobern, D.; Zerger, A. The Atlas of Living Australia: History, current state and future directions. *Biodiversity Data Journal* **2021**, *9*, e65023, [<https://doi.org/10.3897/BDJ.9.e65023>]. <https://doi.org/10.3897/BDJ.9.e65023>.
4. Duarte, C.M. Reviews and syntheses: Hidden forests, the role of vegetated coastal habitats in the ocean carbon budget. *Biogeosciences* **2017**, *14*, 301–310.
5. Hughes, T.P.; Graham, N.A.; Jackson, J.B.; Mumby, P.J.; Steneck, R.S. Rising to the challenge of sustaining coral reef resilience. *Trends in ecology & evolution* **2010**, *25*, 633–642.
6. Bruno, J.F.; Selig, E.R.; Casey, K.S.; Page, C.A.; Willis, B.L.; Harvell, C.D.; Sweatman, H.; Melendy, A.M. Thermal stress and coral cover as drivers of coral disease outbreaks. *PLoS biology* **2007**, *5*, e124.
7. Sheng, V.S.; Provost, F.; Ipeirotis, P.G. Get another label? improving data quality and data mining using multiple, noisy labelers. In Proceedings of the Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining, 2008, pp. 614–622.
8. Li, J.; Knapp, D.E.; Fabina, N.S.; Kennedy, E.V.; Larsen, K.; Lyons, M.B.; Murray, N.J.; Phinn, S.R.; Roelfsema, C.M.; Asner, G.P. A global coral reef probability map generated using convolutional neural networks. *Coral Reefs* **2020**, *39*, 1805–1815. <https://doi.org/10.1007/s00338-020-02005-6>.
9. Usama, M.; Qadir, J.; Raza, A.; Arif, H.; Yau, K.L.A.; Elkhatib, Y.; Hussain, A.; Al-Fuqaha, A. Unsupervised machine learning for networking: Techniques, applications and research challenges. *IEEE access* **2019**, *7*, 65579–65615.
10. Gordon, H.R. Atmospheric correction of ocean color imagery in the Earth Observing System era. *Journal of Geophysical Research: Atmospheres* **1997**, *102*, 17081–17106.
11. Zhu, Z.; Woodcock, C. Object-based cloud and cloud shadow detection in Landsat imagery. *Remote Sensing of Environment* **2012**, *118*, 83–94. <https://doi.org/10.1016/j.rse.2011.10.028>.
12. Wyszecki, G.; Stiles, W.S. *Color science: concepts and methods, quantitative data and formulae*; Vol. 40, John Wiley & sons, 2000.
13. Gonzalez R, R.; Woods, E. Digital Image Processing, 3rd ed Prentice-Hall Inc. Upper Saddle River, New Jersey **2006**.
14. Siegel, D.A.; Wang, M.; Maritorena, S.; Robinson, W. Atmospheric correction of satellite ocean color imagery: the black pixel assumption. *Applied optics* **2000**, *39*, 3582–3591.
15. Pearson, K. LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin philosophical magazine and journal of science* **1901**, *2*, 559–572.
16. Van der Maaten, L.; Hinton, G. Visualizing data using t-SNE. *Journal of machine learning research* **2008**, *9*.
17. McInnes, L.; Healy, J.; Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426* **2018**.
18. MacQueen, J.; et al. Some methods for classification and analysis of multivariate observations. In Proceedings of the Proceedings of the fifth Berkeley symposium on mathematical statistics and probability. Oakland, CA, USA, 1967, Vol. 1, pp. 281–297.
19. Rasmussen, C. The infinite Gaussian mixture model. *Advances in neural information processing systems* **1999**, *12*.

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.