

Spatio-Temporal Time Series Modelling *20/12/2024*

With a Focus on Coral Reef Benthic Group Shifting in the Great Barrier Reef

By: Daud Waqas

I started off this project to tackle a fairly obscure and unexplored topic: *understanding and predicting shifts in coral reef benthic cover using data-driven approaches*. This was also my attempt to try and “handle” a much larger and unoptimised dataset. The data, initially provided at yearly intervals, represented the distribution of coral reef covers (*specifically hard coral cover, soft coral cover and algae cover*) across various reef systems throughout the Great Barrier Reef; due to the massive ecological size of the reef, I made an assumption that there would be distinct temporal patterns based on the spatial location that the data was sampled from. This was the basis of a study that was not just focused on fitting **time series**, but also on building a **spatio-temporal model** that would take those location-based dependencies into account.

The data was compiled from the **AIMS** and **eReefs** platform, both of which are managed by the Australian government; **AIMS** was used to source the primary data, and **eReefs** was used to source the auxiliary data. However, time constraints only allowed me to use the data from the **AIMS** platform (*ie. the primary data*). The data was made up of spatio-temporal indexes (*latitude, longitude, date, etc*) alongside 4 cover values:

- **SOFT CORAL_COVER**
- **HARD CORAL_COVER**
- **ALGAE_COVER**
- **OTHER_COVER**

Most of the data came in an unusable format, and had to be *heavily augmented* to be used properly by the subsequent models. You should be able to find most of the context regarding these augmentations in the ``data_collection.ipynb`` notebook, which I wasn't able to knit due to the fact that it includes a HTML map (labelled as ``folium_click_map.html``). This click map has also been included alongside this report, and you can observe how the details for each selected spatial location was given through the console (*within browser DevTools*) of that HTML file.

All of the modelling took place with the **primary data**, which *originally came sampled in the yearly format* from the **1990s** to **2004**. Since yearly data lacked the granularity needed for temporal analysis, I resampled it into weekly intervals using spline-based interpolation. This was a fairly safe augmentation, since biological metrics over square kilometres of area (like benthic group coverage) don't tend to have massive fluctuations or show volatility in that regard. To mimic **natural variability** and prevent models from fitting from the interpolated data, I added **Gaussian noise** and **scaled the data randomly** within a controlled range. This made sure that any models (especially neural networks) won't derive patterns from the interpolation.

To test temporal dependencies, I started with **univariate models** like **Auto ARIMA** and **TBATS**. These methods failed in capturing any trends at all; the average **nRMSE** result, which was close to **1** for both models (*exactly 0.98 and 0.94 respectively*) made it so that the model was completely off (*note that an acceptable nRMSE starts at around 0.2 or lower*). The heat map showed how there were massive variations, with some values **above 1** and some closer to **0.05** (*which would have been considered a good fit if those results across spatial locations were consistent*). Overall, the models were simply lacking and overly inconsistent.

My choice for **multivariate time series model** was a **Temporal Convolutional Network (TCN)**. This neural network model leveraged the values of all 4 cover values to predict for the next time iteration. The **TCN** outperformed the **univariate models**, and quite aggressively so achieving an

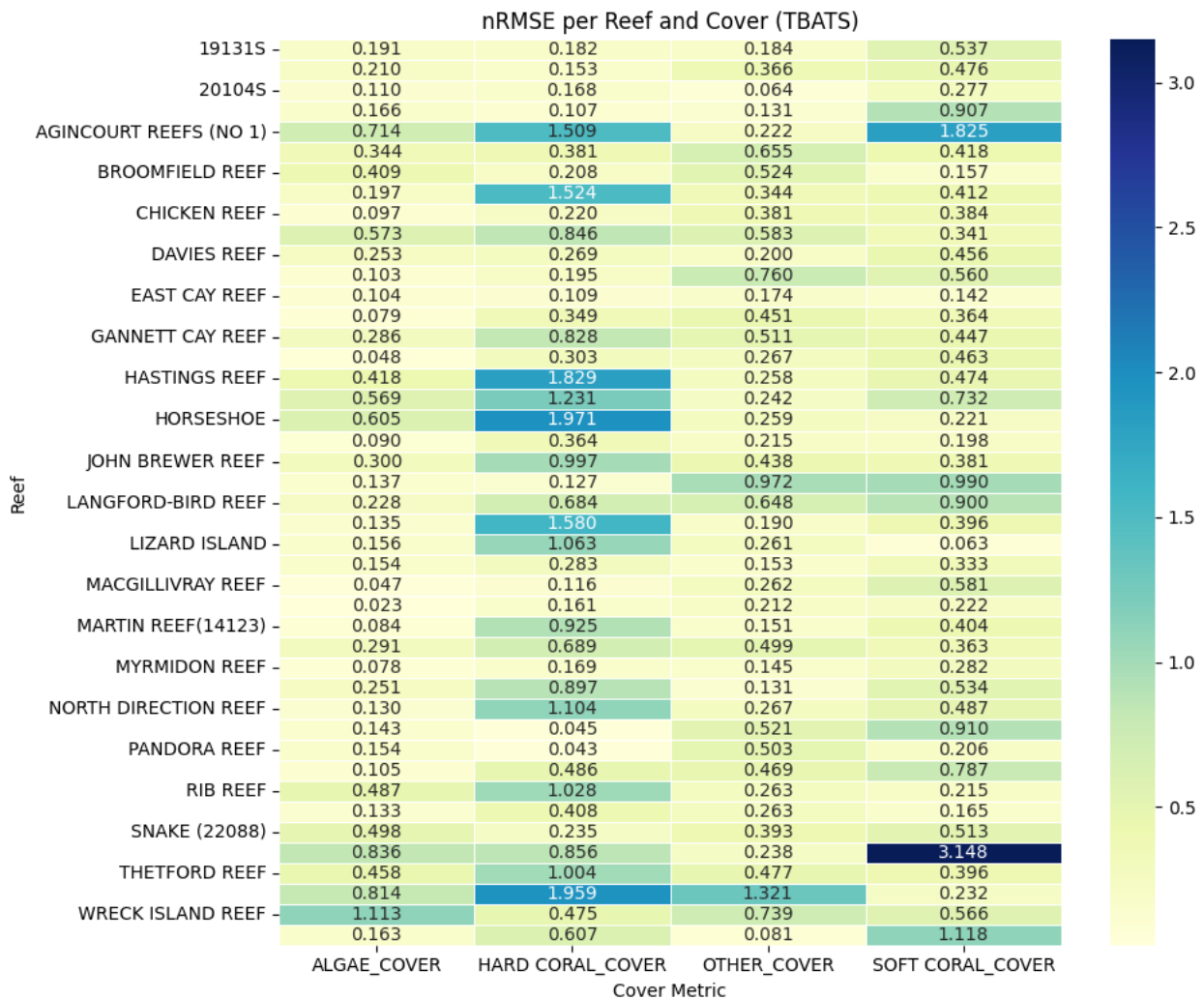
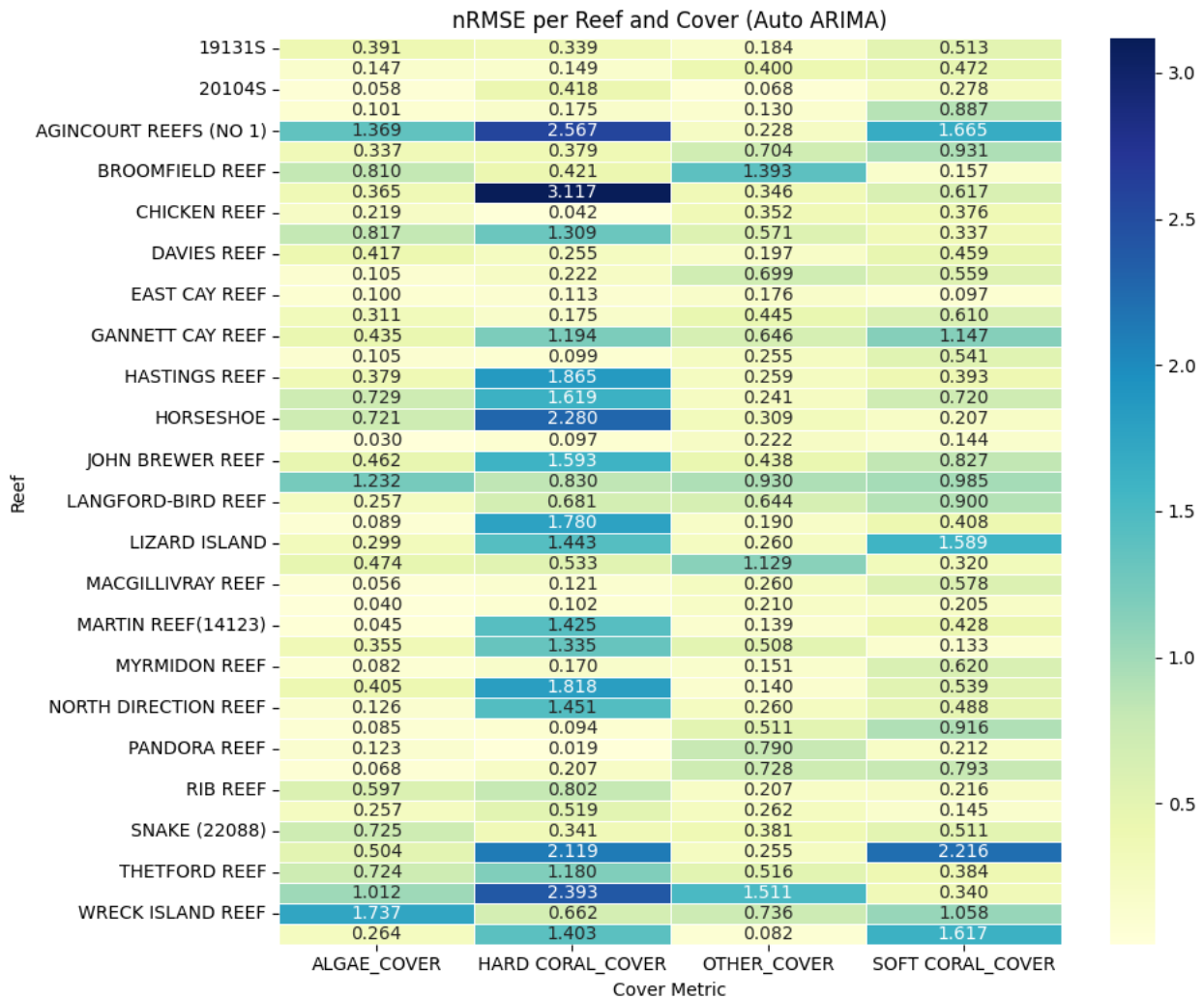
average **nRMSE** of **0.116**. This is a score that is considered a “**moderate fit**”, and goes to show that the granularity of the **neural network** within **TCN** did make up for some of the weaknesses of the data, such as its **multicollinearity**. Please note that the score was derived on an **80/20 split** on the data, where the **latter 20% of the temporal range is the test range**.

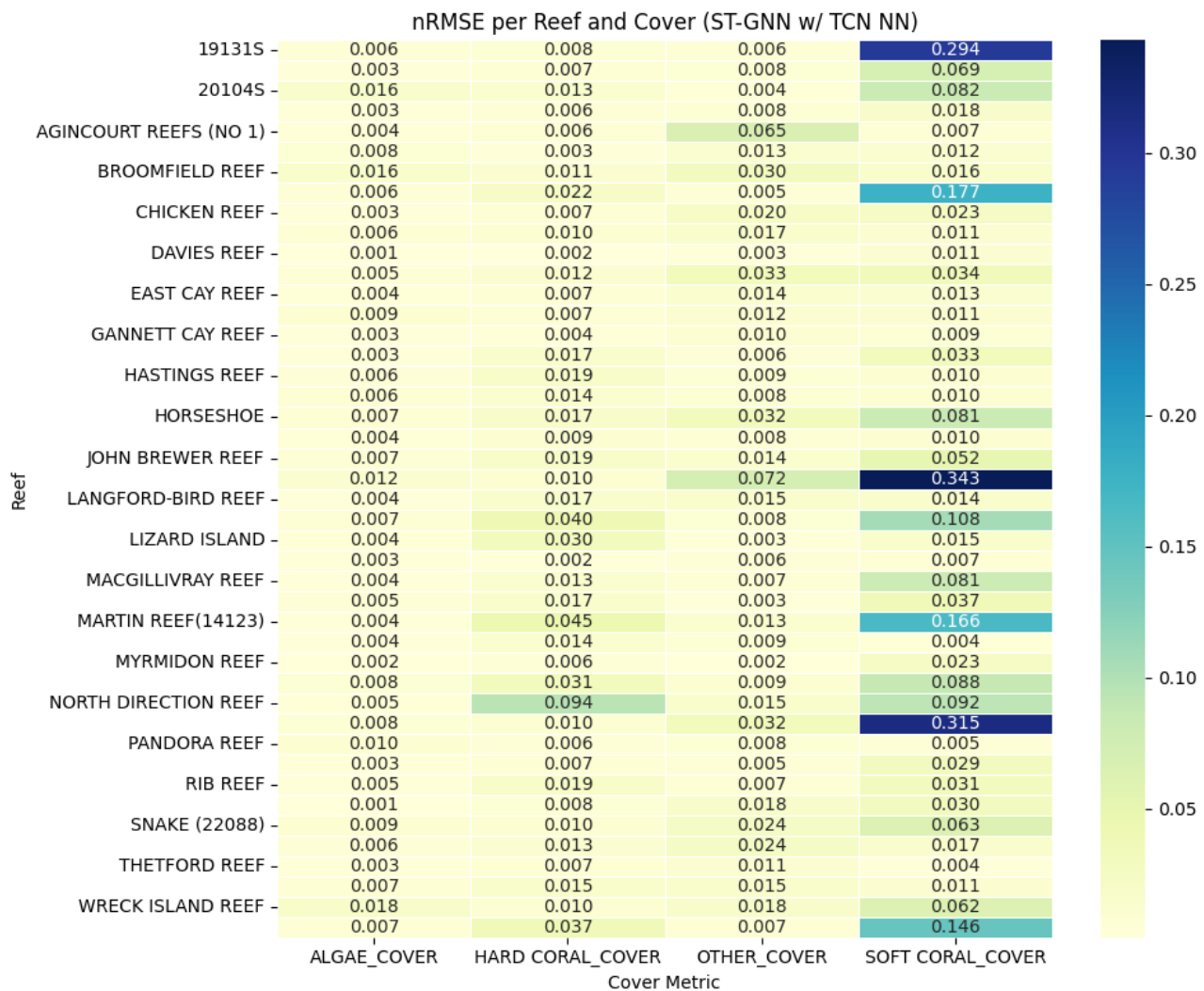
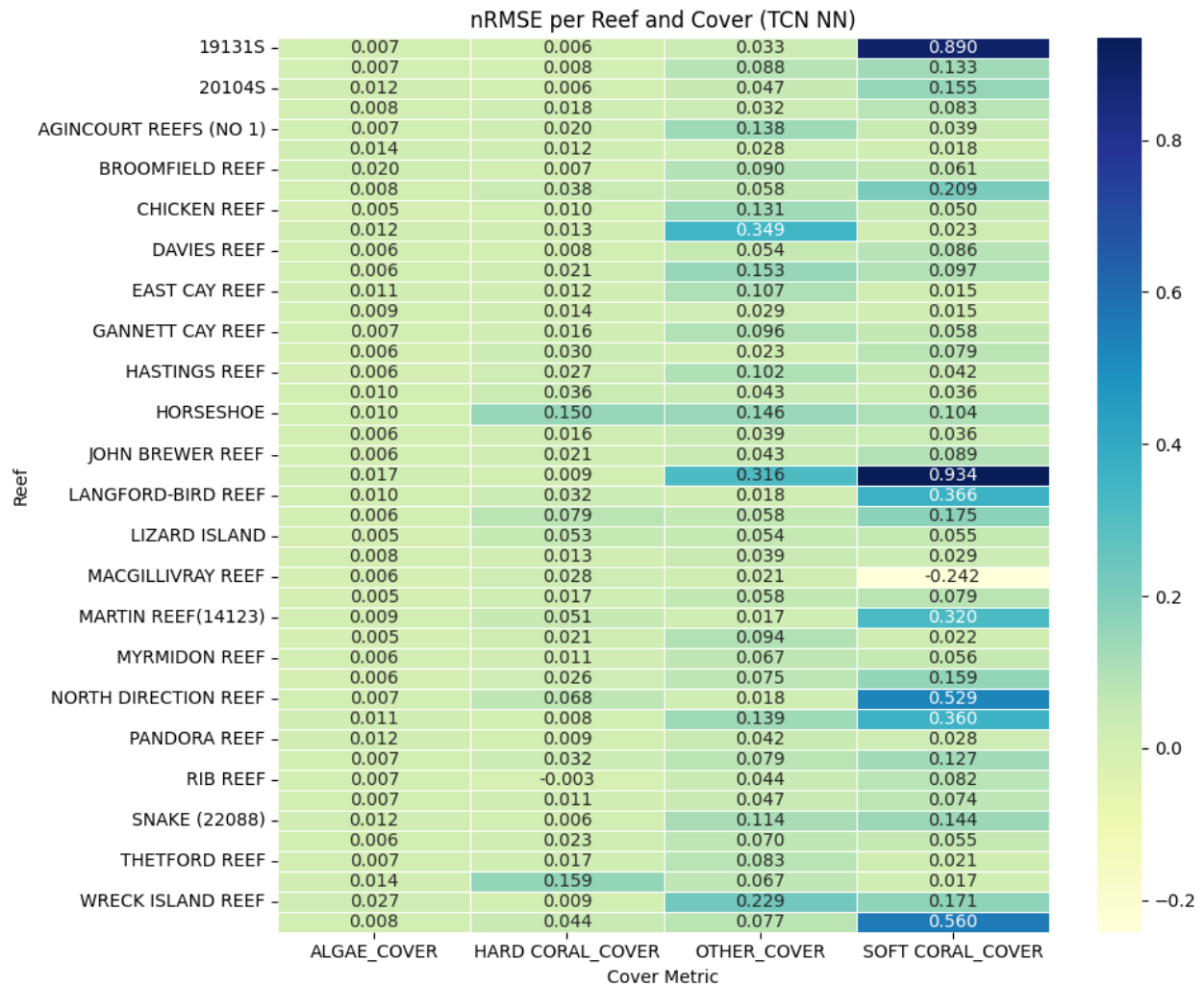
The main part of this project was the implementation of a **Spatio-Temporal Graph Neural Network (ST-GNN)**. By *building a spatial graph where reef locations were connected based on proximity*, I would have been (in theory) able to capture the temporal relationships between the cover values **relative to different reef locations** along the spatial map. This part of the project was quite a struggle, since I couldn't find any proper documentation on how to best fit ST-GNNs. There were many papers out there for the fitting of ST-GNNs for traffic monitoring, but not that many for sensitive environmental-based modelling (*admittedly, I spent too long trying to straighten out this part of the notebook, and had to redesign multiple times to get the results I wanted. I had fairly inadequate results, but this the afternoon before the submission of this assignment I managed to figure out how to implement this properly*).

The **temporal component** of the ST-GNN was handled by a **TCN**, keeping consistent to the multivariate model done beforehand. This **ST-GNN** combined spatial correlations between reefs with temporal patterns provided by the **TCN** to provide a holistic view of the system. Once I did manage to figure out how to implement all this properly, the **nRMSE** values **improved significantly** for several cover types, and managed to not just beat **TCN** but also reduce the sensitivity of the more hard-to-predict cover types in particular reefs. The average **nRMSE** ended up being **0.033**, which is considered a good fit according to nRMSE standards (*and far better than the 0.116 of the TCN*).

While the deep learning models showed promising results, there are still some areas of interest left around this particular topic. A main issue was the implementation of the **auxiliary data**; it ends up that implementing varying feature sets into **neural networks** is a large ongoing subject with not that much knowledge surrounding it, *especially for ST-GNNs where documentation was already sparse*. The paper which did initially fuel my interest, **Spatial-temporal graph neural networks for groundwater data**, uses a quite a complicated setup for their auxiliary data; and for the most part, it involves the combination of **multiple different ST-GNNs knitted together for a custom-purposed solution**. After learning how to interpret and build models at a more mathematical and case-specific level, I may come back to this project in the future.

All the **heat maps** of the **nRMSE results** sorted by **reef** and **cover type** may be found in the below section, including the **graph network** of the **ST-GNN**.





Reef Graph Visualization

