

Sailing the Marine Knowledge Landscape

ARDC RLA PROJECT FINAL REPORT

Meg Forestier-Wardley, Nicole Stubing, Kyaw Kyaw Soe Hlaing

02/12/2024

CONTENTS

1. PROJECT INFORMATION	2
1.1 Background	2
1.2 Achievement of project aims	3
1.3 Summary of value	5
2. DESCRIPTION OF PROJECT OUTPUTS	7
2.1. Achievements against project work packages/deliverables:	7
2.2 Project outputs and outcomes	11
2.2.1 Outputs	11
2.2.2 Outcomes	12
2.3 Outreach and training activities	13
3. PROJECT IMPACT	14
3.1. Impact stories	14
3.2. Impact metrics	16
3.2.1. Use, Outreach and Communications Metrics	16
3.2.2. Links to Publications and Other Communications	17
Blog posts/news articles that reference the platform/infrastructure	17
3.2.3. Details on users of the platform	17
4. LESSONS LEARNED	19
4.1. What Went Well?	19
4.2. What Could be Improved?	19

1. PROJECT INFORMATION

PROJECT TITLE AND CODE	Sailing the Marine Knowledge Landscape (RLA10)
PROJECT START AND END DATES	Start Date: 03/06/2024 End Date: 02/12/2024
LEAD CONTRACTING ORGANISATION	Fisheries Research and Development Corporation
SUBCONTRACTOR REPRESENTATIVE	Andrew Fletcher
PROJECT LEAD CONTACT PERSON	Meg Forestier-Wardley
PROJECT MANAGER	Nicole Stubing, Meg Forestier-Wardley
PARTNER ORGANISATIONS	Volve, National Marine Science Committee
FOCUS AREA and ACTIVITY	Research Link Australia

1.1 Background

Australia has 15 Research and Development Corporations (RDCs) across agriculture, fisheries, and forestry industries. RDCs are co-funded partnerships between the Australian government and industry that invest in research and development (R&D) initiatives pertaining to their respective sectors, collectively investing approximately \$800 M annually. The Fisheries Research and Development Corporation (FRDC) invests in R&D on behalf of the fishing and aquaculture sectors (including recreational and Indigenous) and, over its lifetime, has funded over 5000 research projects. This means that the FRDC has a vast collection of R&D knowledge to disseminate.

The FRDC is currently the only RDC operating primarily in the marine science space and, thus, has a unique research landscape to understand compared to other RDCs that have a primarily land-based

agricultural focus. To date, there have only been a handful of initiatives undertaken to synthesise cross-RDC research. While there are coordinating structures across Australian marine science initiatives (i.e., the National Marine Science Committee (NMSC), Research Provider Network), there have been limited attempts made to map Australian Government funded marine research across the diversity of research funders and institutions that operate in this space. Thus, there is a benefit for the FRDC to explore research specifically from the broader corpus of marine knowledge and to make its own research more discoverable in the process.

The current mechanisms through which the FRDC and other RDCs make their research discoverable include websites, cross-RDC platforms such as GrowAg, or initiatives such as the RDC Knowledge Transfer. The RLA platform provides an opportunity for FRDC-funded research to be made more discoverable and accessible, thereby complementing these existing mechanisms. There is also an opportunity to investigate ways to improve the searchability of research, which enhances discoverability further. Moreover, the RLA initiative provides an opportunity for the FRDC to explore and better understand the marine science landscape, explore the agricultural innovation landscape, and identify potential research partners (both traditional and non-traditional) or priority areas for research. Finally, RLA presents an opportunity to facilitate more cross-sectoral research opportunities through greater exposure on the RLA platform.

The project has three aims:

- 1) To share FRDC research project data to be made available on the RLA platform, thereby making it more discoverable, accessible, and usable.
- 2) To explore different knowledge tools (i.e., RLA-generated dashboards or applications of Large Language Models (LLMs)) to better understand the marine science space and explore the agricultural innovation landscape.
- 3) To make the RLA platform known within the wider marine science community and across different RDCs and encourage the adoption of the tool.

1.2 Achievement of project aims

The first project aim was to share FRDC research project data to be made available on the RLA platform, thereby making it more discoverable, accessible, and usable. To achieve this, four APIs containing FRDC research project data were developed and shared with the ARDC. Due to minor delays, the data are currently in the pipeline to be ingested into the RLA platform and made publicly available. This activity involved the following sub-activities:

- Initial research project data collection.
- Internal process/ procedural changes to enable near real-time data sharing with third parties.
- Data enrichment to maximise research project data shared to the RLA platform.
- Development of an internal database schema to align FRDC data to RLA data schema.
- Development of four API URLs in JSON format – data formatted in accordance with RLA data schema.
- Ensure data provision aligns with internal data-sharing processes/ procedures to enable near real-time data sharing with third parties.
- Provision of research project data via API URLs to the ARDC to be ingested into the RLA platform.
- Workshops with the Australian Research Data Commons (ARDC) to refine input to the RLA platform.

The second aim was to explore different knowledge tools (RLA-generated dashboards and application of LLMs) to better understand the marine science space and explore the agricultural innovation landscape. To achieve this, an RLA-generated dashboard was developed by the ARDC, which showcased Australian Research Council (ARC)-funded research activities within the fisheries research domain, was explored. To explore the application of LLMs, the APIs containing research project data were indexed by a third-party LLM and a chatbot proof of concept was built against that LLM. This activity involved the following sub-activities:

- Development and implementation of a policy on the safe and responsible use of generative Artificial Intelligence (AI) processes/tools that met the guidelines set out by the Australian Government/ Australian Signals Directorate (ASD).
- Indexation of the data from the four APIs developed by a third-party LLM.
- Development of a chatbot proof of concept against that LLM, refining the APIs and development process as challenges arose.
- Integration of learnings into FRDC business processes/organisational intelligence.

The third aim was to showcase the RLA platform to key members within the marine science community and members of other RDCs and gather user stories and input/feedback on the platform. This aim was partially achieved. The RLA platform was introduced to members across these spaces. This activity involved the following sub-activities:

- Showcase the RLA platform to RDC members.
- Showcase the RLA platform to National Marine Science Committee (NMSC) members.

Due to delays in the ingestion of FRDC research project data to the RLA platform, the following sub-activities could not be achieved within the timeframe of the project:

- Consult members of RDCs for user testing of the RLA platform and input/ feedback.
- Consult members of the NMSC for user testing of the RLA platform and input/ feedback.
- Provide user stories and input/ feedback on the RLA platform to the ARDC.

These activities will be achieved after the completion of the project once the data is made available on the RLA platform.

1.3 Summary of value

Improved research discoverability will yield value for current and future research funders (i.e., knowledge producers), those who use FRDC research, including researchers, policymakers, and stakeholders (i.e., knowledge users), the broader fishing and aquaculture community, and the RDC community. The project will also inform FRDCs organisational intelligence. The benefits yielded from the project are expected to be long-term and impact actors in and outside the marine science space.

Research Funders – Knowledge Producers

Sharing FRDC research projects on the RLA platform enables enhanced discoverability, facilitating knowledge transfer between different actors both in and outside the marine science space, including other research funders, industry, business, and government. Research becomes more accessible and usable, ensuring FRDC investment is maximised. Making FRDC research accessible through the RLA platform also enables it to be seamlessly integrated with research produced by other research funders in the marine science space, enabling them to:

- Identify where research has already been conducted, reducing the risk of duplicated investment,
- Identify knowledge gaps within the larger body of marine science knowledge, and
- Plan future research activities accordingly.

Moreover, sharing FRDC research to the RLA platform will link together different actors and foster collaboration, which is critical for innovation in the marine science space.

Knowledge Users

Exploring the application of LLMs can significantly enhance the searchability and discoverability of FRDC research, providing long-term value to diverse knowledge users. LLMs are advanced computational models designed to perform natural language processing (NLP) tasks, which enable computers to

understand, interpret, and generate human language based on vast amounts of data. One practical application of LLMs is in creating chatbots that take user inputs and generate contextually relevant responses. A user-friendly, conversational chatbot could improve upon existing search functionalities that have significant limitations, such as keyword dependencies and an overload of search results that require the user to sift through potentially a lot of deadwood to find relevant results, improving the information retrieval process overall. Knowledge users would be able to access research in an easier, more simplified manner. The lessons learned through the development of a chatbot based on FRDCs research project data could be extended and applied to other organisations, including RLA or GrowAg. The value of this would be in improving the searchability and discoverability of research from across the broader Australian research ecosystem.

Organisational Intelligence and Fishing and Aquaculture Community

By exploring RLA-generated knowledge tools, such as dashboards, the FRDC would be able to gain a better understanding of the broader marine science research landscape, informing future FRDC investments and R&D application assessments, and minimising the risk of duplicating investment that exists elsewhere. The FRDC would realise the value of this activity in the first instance as it would be expanding its current information base on the marine science research landscape. However, this exercise can have flow-on effects that are long lasting to benefit the broader fishing and aquaculture community. Furthermore, the learnings derived from this exploration could also inform national structures, such as the NMSC, which is an advisory body that sets the strategic direction for and assists in the coordination of marine science R&D activities in Australia, or the Department of Climate Change, Energy, the Environment, and Water's (DCCEEW) "Sustainable Oceans Plan", which aims to develop a strategic framework to guide the responsible management of Australia's marine areas.

RDC Community and Research Discoverability

The FRDC intends to provide a rich data source to the RLA platform – details of its funded research projects – thus enriching national information. Given the FRDC is the only RDC participating in RLA and, more broadly, dedicated funding source for fishing and aquaculture-related research, the input of such information to the RLA platform is inherently unique and enriches the existing body of marine knowledge. FRDCs involvement in the RLA initiative can be used to demonstrate best practice in research discoverability to other RDCs who collectively represent over \$800 M in research investment each year.

2. DESCRIPTION OF PROJECT OUTPUTS

2.1. Achievements against project work packages/deliverables:

Use the following criteria and colours to assign a status to each work package:

Status Key

Completed	Completed
Non Completion may impact on project outcomes	Not yet completed and non-completion may have some impact on the project outcomes. Delay has been accepted by the Steering Committee and responsibility for completion has been assigned.
Will not be completed and will impact on project outcomes	Will not be completed and has an impact on the project outcomes. The Steering committee has accepted the non-completion and understands the impact to the project outcomes.

DELIVERABLE / WORK PACKAGE	DETAILS INCLUDING EXPLANATION OF ANY VARIATION (INCLUDE LINKS TO RELEVANT DOCUMENTATION)	AGREED DUE DATE	ACTUAL or EXPECTED COMPLETION DATE	STATUS
WP 1 – Provision of FRDC research project data via API URL to the RLA platform.	<ul style="list-style-type: none"> - 1.1 – Initial research project data collection. - 1.2 – Internal consultation and process change to enable near real-time data sharing with third parties. - 1.3 – Data enrichment 1 - determine if additional fields/translations need to be included to maximise research project data shared to the RLA platform. - 1.4 – Data enrichment 2 - consult members of NMSC for feedback on research project data to be shared to RLA platform. - 1.5 – Development of an API URL in JSON format. - 1.6 – Ensure provision of API URL aligns with internal FRDC data sharing processes/ procedures to enable the near real-time sharing of research project data. - 1.7 – Provision of research project data via API URL to the RLA platform. - 1.8 – Workshop(s) with the ARDC RLA project team on refining input to the RLA platform. 	06/09/2024	06/09/2024	Completed
WP 2 – Explore knowledge	<ul style="list-style-type: none"> - 2.1 – Explore RLA-generated knowledge tools. 	02/12/2024	02/12/2024	Completed

tools (i.e., RLA-generated dashboards or application of LLMs) to explore and better understand the marine science landscape and the agricultural innovation space.	<ul style="list-style-type: none"> - 2.2 – Develop and implement the safe and responsible use of generative AI processes that meets the requirements set out by the ASD. - 2.3 – Index the data from the developed API by third-party LLM. - 2.4 – Build a chatbot proof of concept (POC) against that LLM. - 2.5 – Investigate integration of learnings into FRDC business processes/ organisational intelligence. 			
WP 3 – Showcase the RLA platform to key members within the marine science community and members of other RDCs and gather user stories and input/feedback on the RLA platform.	<ul style="list-style-type: none"> - 3.1 – Showcase the RLA platform to RDC members. - 3.2 – Showcase the RLA platform to NMSC members. 	02/12/2024	02/12/2024	Completed
	<p>The following activities were not completed due to delays in the ingestion of FRDC research project data into the RLA platform:</p> <ul style="list-style-type: none"> - 3.3 – Showcase the RLA platform at an open ARDC online forum meeting (upon request by the ARDC). - 3.4 – Consult members of the NMSC for user testing of the RLA platform and input/ feedback. - 3.5 – Consult members of RDCs for user testing of the RLA platform and input/ feedback. 	02/12/2024	TBD	Activities intended to occur once FRDC research project data are made available on the RLA platform.

	<ul style="list-style-type: none">- 3.6 – Provide user stories and input/feedback on the RLA platform to the ARDC.			
--	--	--	--	--

2.2 Project outputs and outcomes

2.2.1 Outputs

OUTPUTS	RESULTS SUMMARY	EVIDENCE
Output 1 – FRDC research project data provisioned to the RLA platform.	<p>Four API URLs were developed in JSON format containing FRDC research project data and shared with the ARDC.</p> <p>The ARDC fetched the data and mapped it to RLAs data schema. The FRDC and ARDC worked to refine input to the platform.</p> <p>Due to delays, the data currently remains in the pipeline to be ingested into the RLA platform.</p>	<p>FRDC research project data provisioned to the ARDC.</p> <p>The data is currently in the pipeline to be ingested into the RLA platform and made accessible here: https://devl.researchlink.ardc.edu.au/.</p>
Output 2 – Exploration of RLA-generated knowledge tools.	<p>The ARDC developed a demonstration dashboard displaying ARC-funded activities in the fisheries research domain.</p> <p>The dashboard provided a broader awareness of research activities being conducted in the marine science research space and highlights the sorts of valuable insights that can be easily derived from data displayed in this format (collaborations, project/ grant details).</p>	<p>RLA-generated dashboard, “Collaborations in Fisheries Research”, accessible here: https://demo.researchlink.ardc.edu.au/dashboardList/published.</p>
Output 3 – Exploration of the application of LLMs, and development of an LLM-powered chatbot proof of concept.	<p>The data contained within the APIs were further refined and indexed by a third-party LLM and a chatbot proof of concept was built against the LLM.</p>	<p>The chatbot is in the development phase. Demonstration is available upon request to the FRDC. See Section 4.2 for further details.</p>

OUTPUTS	RESULTS SUMMARY	EVIDENCE
	The challenges, learnings, and broader implications for RLA of this activity are discussed further in Section 4.	
Output 4 – RLA platform showcased to RDC and NMSC members.	<p>The RLA initiative and RLA platform were introduced to RDC and NMSC members.</p> <p>The current project, “Sailing the Marine Knowledge Landscape”, was also introduced as a case study demonstrating the capacity for RLA to integrate knowledge across the marine science research landscape and situate it within the broader Australian research ecosystem.</p>	Increased awareness of the RLA initiative, the RLA platform, and the “Sailing the Marine Knowledge Landscape” project within RDC and NMSC members.

2.2.2 Outcomes

OUTCOME	RESULTS SUMMARY	EVIDENCE
Outcome 1 – Greater transfer of knowledge between research funders, industry, business, and government (lag outcome).	<p>Increase in awareness of FRDC research projects as measured through:</p> <ul style="list-style-type: none"> - Number of FRDC research projects shared to the RLA platform. - Number of page visits to FRDC research project pages on RLA platform. 	<p>Outcome 1 to be achieved once FRDC research project data is made available on the RLA platform.</p> <p>Currently, 4500 research projects and their associated researchers (10239), organisations (4371), and final reports (3652) have been made available to ARDC/RLA. These data are currently in the RLA ingestion backlog.</p> <p>Data for measure to be obtained from RLA platform user usage data.</p>

OUTCOME	RESULTS SUMMARY	EVIDENCE
Outcome 2 – Greater collaboration between actors in the marine science space (lag outcome).	<p>Increased number of interactions between members within the marine science research space as measured through:</p> <ul style="list-style-type: none"> - Number and frequency of inquiries to RLA about FRDC and/ or FRDC research by marine science researchers. 	<p>Outcome 2 to be achieved once FRDC research project data is made available on the RLA platform.</p> <p>Data for measure to be obtained from RLA platform user usage data.</p>

2.3 Outreach and training activities

ACTIVITY	ACTIVITY DESCRIPTION	NO. OF PARTICIPANTS	DATE OF ACTIVITY
RLA initiative, RLA platform, and current project introduced to members from other RDCs and the NMSC.	<p>The RLA initiative and RLA platform were introduced to RDC and NMSC members.</p> <p>The current project, “Sailing the Marine Knowledge Landscape”, was also introduced as a case study demonstrating the capacity for RLA to integrate knowledge across the marine science research landscape and situate it within the broader Australian research ecosystem.</p> <p>The purpose of this activity was to increase awareness of the platform and encourage adoption of the RLA infrastructure.</p>	55	02/12/2024

3. PROJECT IMPACT

The ARDC and the Government wish to demonstrate the impact on researchers, industry and the general public of the NCRIS investment. Please complete the following sections as relevant to your program.

3.1. Impact stories

The FRDC has invested in and produced an extensive amount of research across its lifetime (over 5000 projects). However, the life of this research does not end after it has been produced. Research can be lost if it is not developed, used, and effectively translated into practical, real-world applications. As shown in Figure 1, there is a critical juncture in the research translation journey after research has been produced where it either remains underutilised due to inadequate translation or becomes impactful through adoption and application. The risks of non-adoption, often stemming from inadequate translation, are significant and can result in lost societal benefits and wasted resources. Mechanisms that facilitate knowledge transfer and enable research to be made more visible, like the RLA platform, become vital to increase awareness of research – not just FRDC-produced research, but research generated by any institution, within any sector – thereby mitigating the risks associated with non-adoption.

Research Translation Journey

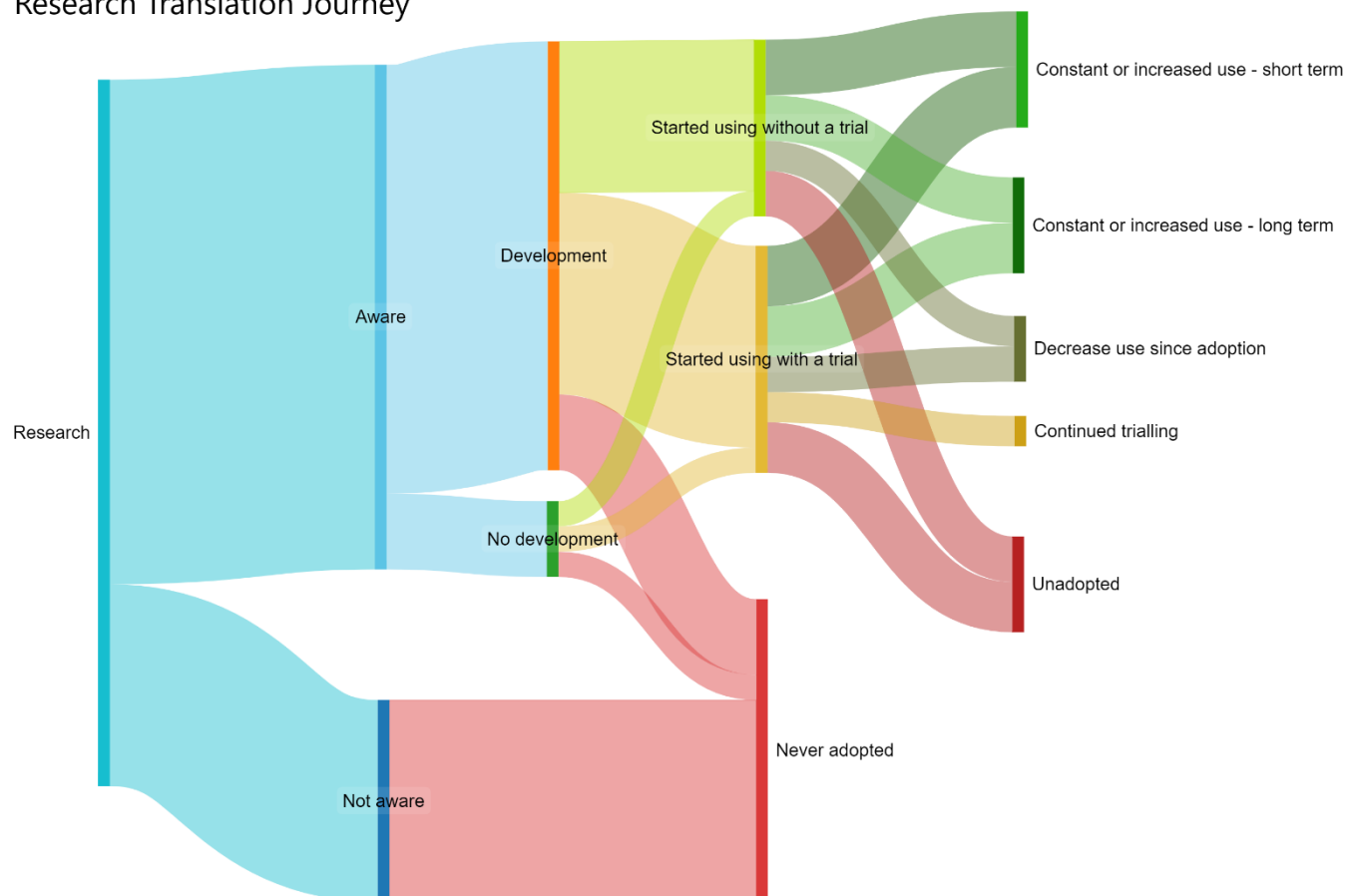


Figure 1: Stages of the research translation journey.

Particularly in the marine science space, RLA is only one among very few mechanisms that currently exist that aims to improve research awareness and translation. The impact of the RLA initiative will enhance the discoverability of FRDC research, which, in turn, will increase awareness and likelihood of adoption, which is fundamental to ensure research is maximally utilised.

3.2. Impact metrics

3.2.1. Use, Outreach and Communications Metrics

ARDC realises that a number of possible impact metrics are lag indicators.

IMPACT METRIC	PROJECT TOTAL TO DATE	DESCRIPTION, LINKS and DOIs
FRDC project site created for “Sailing the Marine Knowledge Landscape”, enhancing awareness of the current project and FRDCs contribution to the RLA platform.	26 page views since page was published.	https://www.frdc.com.au/project/2023-204
ARDC project site created for “Sailing the Marine Knowledge Landscape”, enhancing awareness of the current project and FRDCs contribution to the RLA platform.	43 page views since page was published.	https://ardc.edu.au/project/sailing-the-marine-knowledge-landscape-with-rla/
Project mentioned at launch of Research Link Australia held on 19 September 2024, enhancing awareness of the current project and FRDCs contribution to the RLA platform.	N/A	https://ardc.edu.au/article/driving-innovation-with-research-link-australia/#:~:text=The%20platform%20supports%20research%20discoveries,Canberra%20on%2019%20September%202024.
Project mentioned at “ARDC Leadership Series: Enabling Research Translation” panel held on 25 September 2024, enhancing awareness of the current project and FRDCs contribution to the RLA platform.	N/A	https://ardc.edu.au/event/ardc-leadership-series-enabling-research-translation/
Future impact metric: A forthcoming news article developed by the FRDC discussing the “Sailing the Marine Knowledge Landscape” project.	Number of page visits.	Article in development. To be published early 2025. The article will be published to the FRDC website and disseminated via FRDCs

		monthly newsletter, which is received by a variety of stakeholders across fishing and aquaculture – sent to 8,800+ recipients.
Future impact metric: FRDC research projects accessed on the RLA platform.	Number of page visits.	Once FRDC research project data has been made available on the RLA platform, an impact metric will be the number of page visits for FRDC research projects pages.

3.2.2. Links to Publications and Other Communications

Blog posts/news articles that reference the platform/infrastructure

DATE	LINK
FRDC news article in development. To be published early 2025.	https://www.frdc.com.au/frdc-news

3.2.3. Details on users of the platform

Not applicable.

ORGANISATION TYPE	NUMBER OF USERS
.edu.au	
.gov.au	

.org.au	
.edu	
.com	
other	

4. LESSONS LEARNED

4.1. What Went Well?

Greater Research Discoverability

There are very few mechanisms in place that attempt to synthesise and centralise the vast amount of knowledge produced by different RDCs. GrowAg is the primary platform to connect stakeholders across the Australian agrifood sectors. However, the contributors to this platform are predominantly from land-based agricultural backgrounds and exclusively RDCs. The FRDC stands out as the only RDC focused primarily on marine-related R&D. Despite its unique position, there have been very few attempts to synthesise RDC-produced marine research with research from other institutions. This lack of integration has hindered the ability to create a cohesive and comprehensive understanding of work conducted/underway in the space, limiting the potential for collaboration and innovation.

The RLA platform addresses this issue by centralising research data and linking actors across different sectors. RLA facilitates the synthesis of diverse research outputs, making it easier for knowledge users to leverage this information. This integration not only enhances the discoverability of FRDCs research but also situates it within the broader ecosystem of marine research. Consequently, the FRDC can more effectively contribute to and benefit from the collective knowledge, thereby fostering greater collaboration and innovation in the marine science space.

4.2. What Could be Improved?

Harmonising Datasets and Project Tagging/ Categorisation

To share FRDC research projects to the RLA platform, the data needed to align with the RLA data schema. Largely, reconciling the data with the schema was a straightforward task. However, there were instances where the alignment of data fields was not easy. For example, project keywords. The FRDC tags its projects with keywords from a controlled vocabulary list highly specific to its domain (marine research). RLA uses Australian and New Zealand Standard Research Classification (ANZSRC) codes. ANZSRC codes capture broader research categories that extend beyond marine research but do not capture the details necessary for FRDC-specific usage.

A challenge within this existing approach is that harmonising distinct keyword lists requires manually mapping them together to make them ‘fit’. However, as shown through this project, this can be difficult. A further area for investigation beyond this project could be exploring the role of AI-driven keyword tagging and project categorisation. This approach would remove the need for mapping or aligning distinct keyword lists. Another benefit of such an approach is making the keyword tagging process more efficient overall, reducing human workload and human error. Furthermore, AI-driven keyword tagging could improve the discoverability and searchability of research through more accurate keywords that better reflect the research, which relates to one of the core objectives of this project.

Challenges in the Application of Large Language Models

One of the aims of the project was to explore the application of LLMs and develop a chatbot proof of concept. The main purpose of this investigation was to see if and how an LLM-powered chatbot could improve the searchability and overall information retrieval process of FRDC research given the limitations of existing search functionalities. Figure 2 shows the chatbot user-interface that was built, which is currently embedded within FRDCs User Acceptance Testing (UAT) environment. The chatbot is not yet publicly available as it is being further developed/ refined.

Tell me about bycatch reduction in fishing

Response

Bycatch Reduction in Fishing: A Comprehensive Overview

Bycatch, the unintentional capture of non-target species during fishing operations, represents a significant challenge in fisheries management. Reducing bycatch is critical for maintaining marine biodiversity, adhering to conservation goals, and improving the sustainability of fishing practices. However, based on the provided context, there is limited direct information regarding specific bycatch mitigation measures in Australian fisheries research projects. The available data highlights related themes, such as ecological sustainability, the utilisation of bycatch species, and the role of advanced technologies, but lacks explicit focus on measures like gear innovations or regulatory frameworks aimed purely at reducing bycatch.

Key Insights Relevant to Bycatch Reduction

1. **Sustainable Use of Bycatch Species: Shark Bycatch in the Northern Territory**
 - **Project 2020-036** investigates the biology and population connectivity of shark bycatch species, specifically the Whitecheek and Milk sharks, in the Offshore Snapper Fishery (OSF). These species account for approximately 10% of the total bycatch in the OSF, which primarily targets Saddletail and

What do you want to know?



Figure 2: Chatbot user-interface embedded in FRDCs UAT.

One positive outcome of the investigation was the implementation of semantic search capabilities, which improved the relevance of search results and made the information retrieval process of FRDC research easier. Another outcome was the learnings gained from the challenges encountered during the development of the LLM-powered chatbot (see Figure 3 for development workflow).

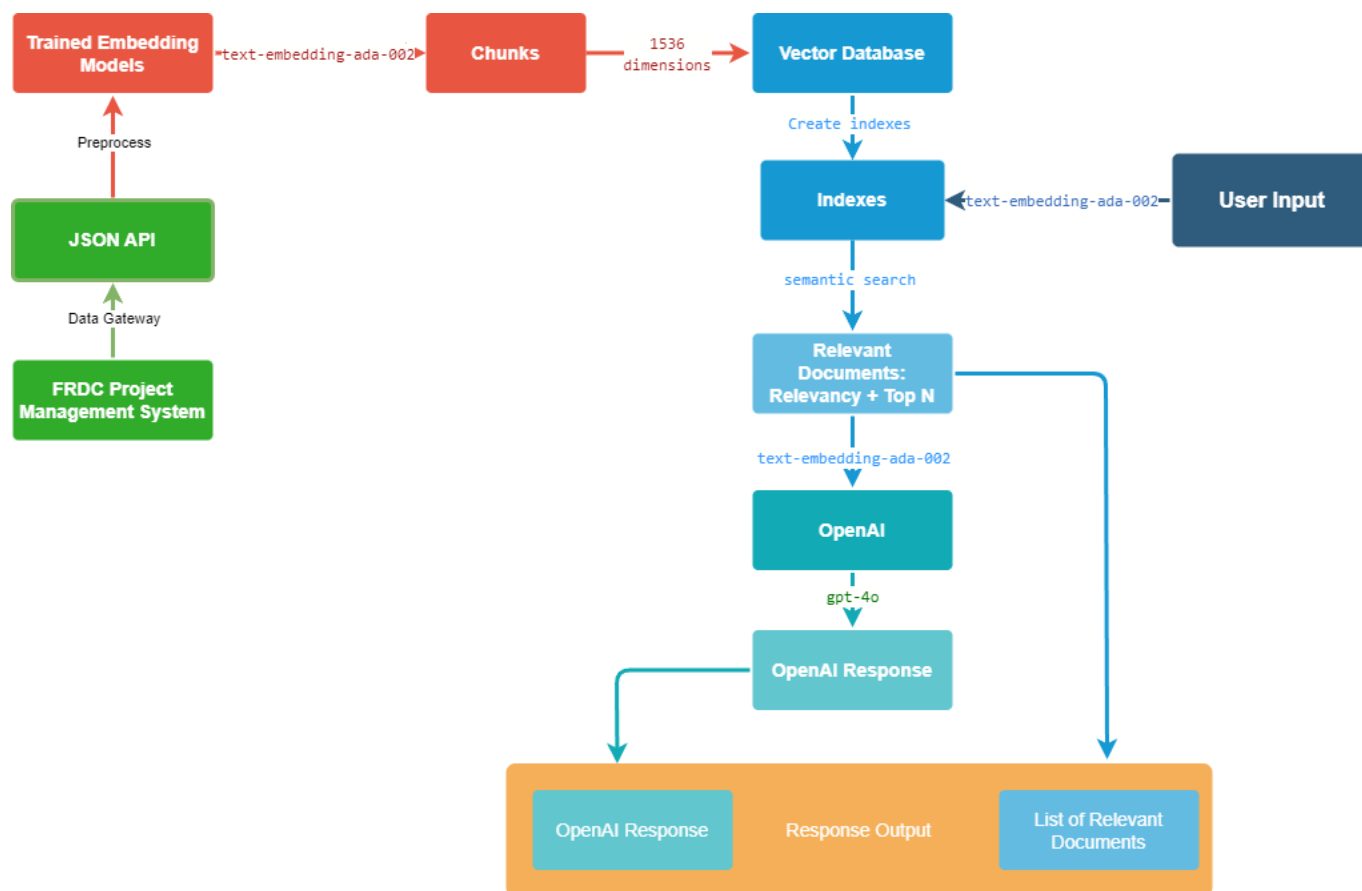


Figure 3: System workflow for the development of an LLM-powered chatbot.

Development of Safe and Responsible Use of AI Applications Policy

Before developing the LLM-powered chatbot, it was necessary to develop and establish a policy outlining the safe and responsible use of AI applications within the FRDC. Developing such a policy was challenging as the use of AI in organisations is still in its nascent stages, and guidelines outlining the safe and secure use and implementation of AI are still being established. Developing the policy required drawing on current Australian government guidelines, including those established by the ASD, to create a policy robust enough to address immediate concerns while remaining flexible to adapt to future advancements and regulatory changes.

Need for Large and Comprehensive Datasets

The first challenge faced in developing the LLM-powered chatbot was the limited amount of information contained within the APIs that were developed as part of the project, which were intended to train the LLM. These APIs didn't contain sufficient data, which constrained the breadth of knowledge the LLM

could process and understand, resulting in a narrow knowledge scope that restricted its ability to generate meaningful and contextually relevant responses to different queries. To address this, new APIs needed to be created, which pulled a broader range of project details from FRDCs project management system. By expanding the dataset, the LLM was exposed to a more comprehensive set of information, significantly improving its ability to produce relevant outputs.

Chunking Large Datasets

Handling such large datasets required splitting the data into small, manageable chunks before embedding them into vector representations and storing them in a vector database. However, the chunking process introduced complexities, and errors in chunking propagated through the system workflow, ultimately impacting the ability of the chatbot to return accurate information. Deciding how to split the data was non-trivial, as it was essential to retain semantic coherence to ensure chunks remained contextually relevant. Additionally, testing these chunks against models for accuracy could be computationally expensive and time-consuming, requiring efficient preprocessing pipelines and robust evaluation metrics.

Testing Different Embedding Models

A significant challenge was in finding the right embedding model that provided the necessary semantic depth, offered good domain specificity, and was high-performing yet not overly computationally demanding – balancing processing time speed with the quality of responses. This balance was crucial for ensuring the chatbot could generate highly relevant and contextually rich responses in a timely manner. Many models were tested to find the best fit, including “text-embedding-ada-002” developed by OpenAI. However, the model that was found to suit our specific needs was “nomic-embed-text-v1.5”, developed by Nomic. The model was chosen for its flexibility, performance, and ability to handle the domain-specific nuances of the data. Furthermore, it was able to process queries in a timely manner while still maintaining contextual relevance in responses.

Implications for Other Organisations

The LLM-powered chatbot developed through the project demonstrates how the searchability and discoverability of research, or knowledge more broadly, can be improved, going beyond the capabilities of current search functionalities. The result of this investigation highlights the potential to change how research data is accessed by knowledge users. Moreover, the learnings from this attempt can be utilised in other attempts to develop an LLM-powered chatbot. Specifically, the need for comprehensive datasets to train an LLM on, consideration for chunking large datasets and how to chunk the data to

retain semantic coherence, and the need to find the right embedding model for your specific needs, which balances quality with computational cost. Integrating this capability into RLA or other platforms such as GrowAg, could further improve the discoverability and accessibility of research from across the broader Australian research ecosystem.