

Ripple – Using open-source tools to turn complex data into applied intelligence

Dr Petra Muellner Nick Snellgrove The challenge

Data

Analysis

Surveillance

Research

Bridging the GAP

Intelligence

The ability to read and respond effectively to a situation'. It's all about how you can gather together data in order to make faster, clearer decisions.

Decision advantage

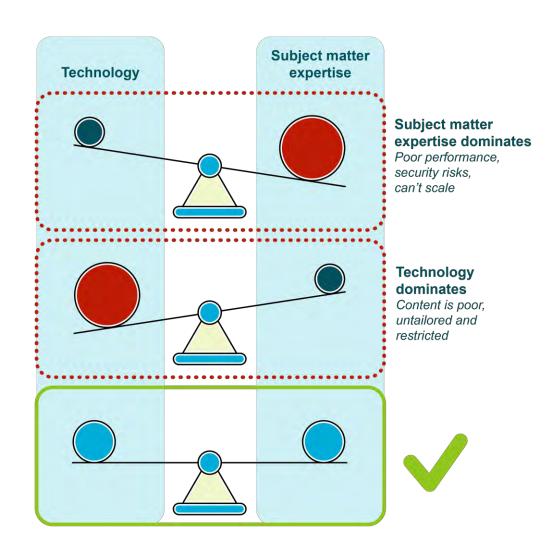
Decision advantage when intelligence enables a decision-maker to better understand and address an issue.



A fine balance

 Operational intelligence to improve human, animal and environmental health

 Highly bespoke solutions – driven by the joint power of subject matter expertise & technology





Environmental DNA (eDNA)

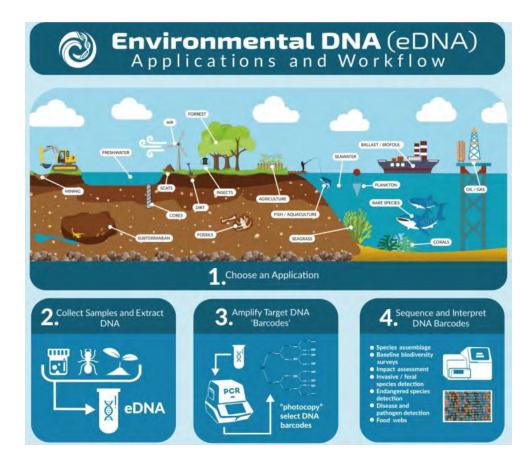
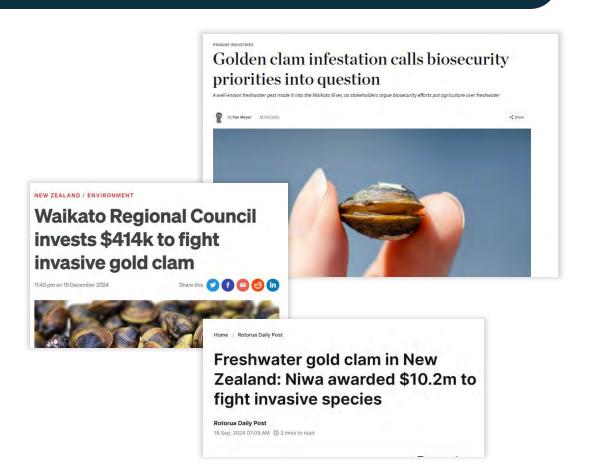


Image credit: B3 Better Border Biosecurity



Acknowledgements

We would like to thank the following organisations for their collaboration, data, insights and support







eDNA data – this is what we have heard

- Frustrations about difficulty to work with the data
- High cost of doing nothing e.g. only finding out of an incursion when it is too late; ecosystems under pressure, water quality
- Not able to work together with others
- Fragmentation by location, by teams within teams



BUT so much knowledge and passion.

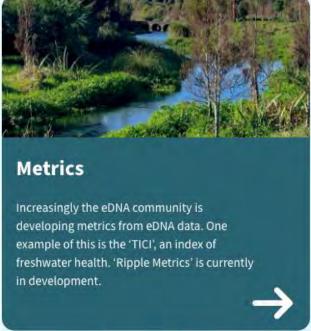




Connecting eDNA with action

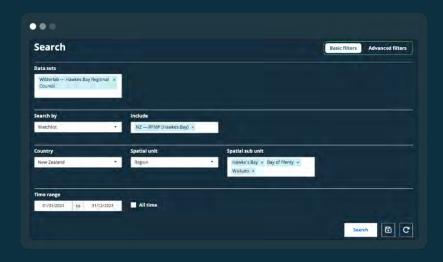
Ripple includes two dedicated modules. The first module, "Biosecurity and biodiversity" supports biodiversity surveys and biosecurity activities. For example, the identification and management of biosecurity incursions, and the discovery of new populations of nationally critical endemic species. Ripple's second module, "Metrics" is currently in development and focusses on water quality, including metrics for monitoring ecosystem health.

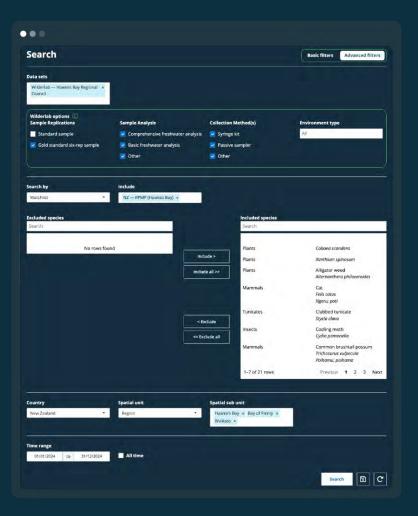




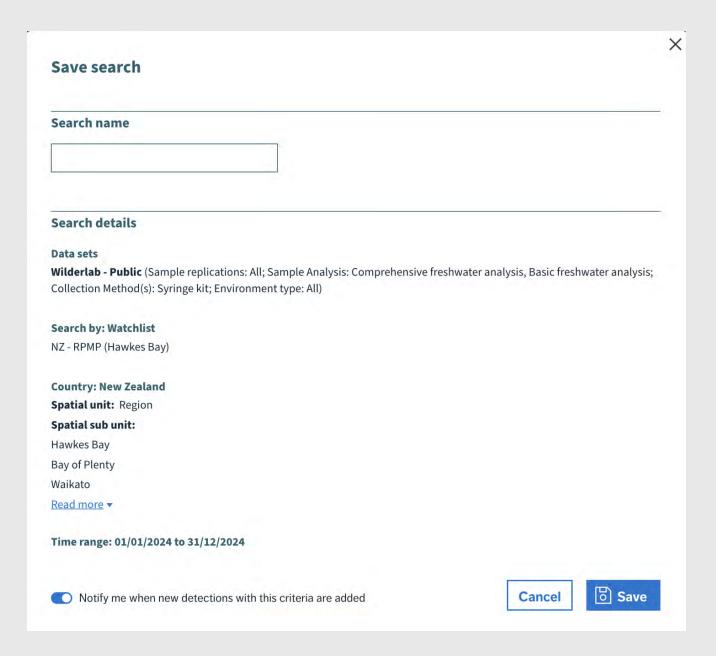


Highly customised searches





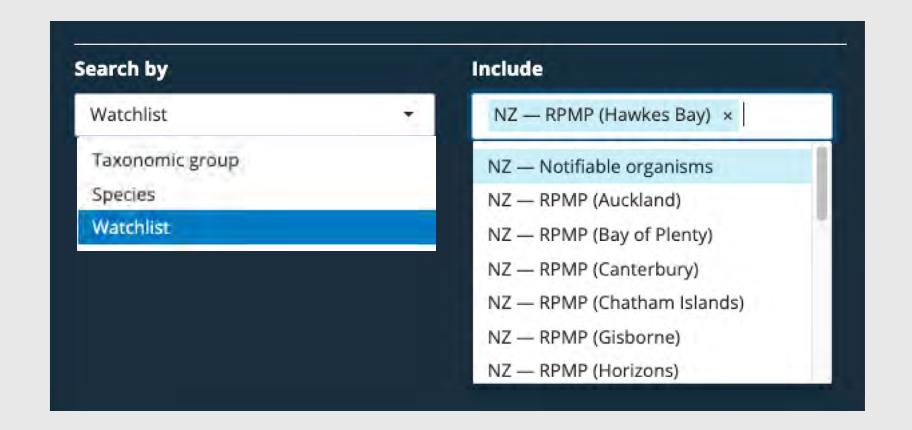
Saved searches



Lab specific search options



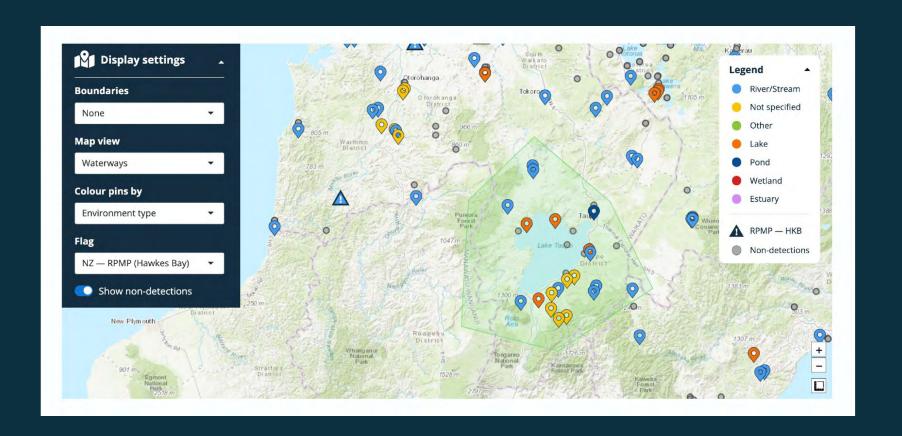
Watchlists



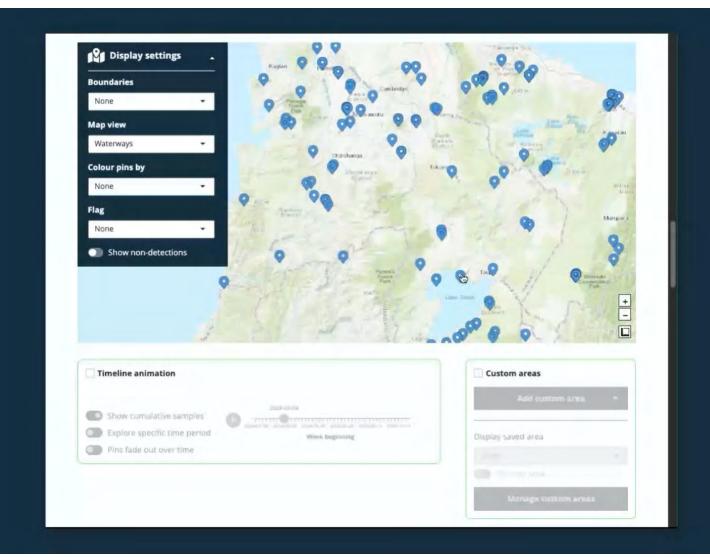
Spatial units



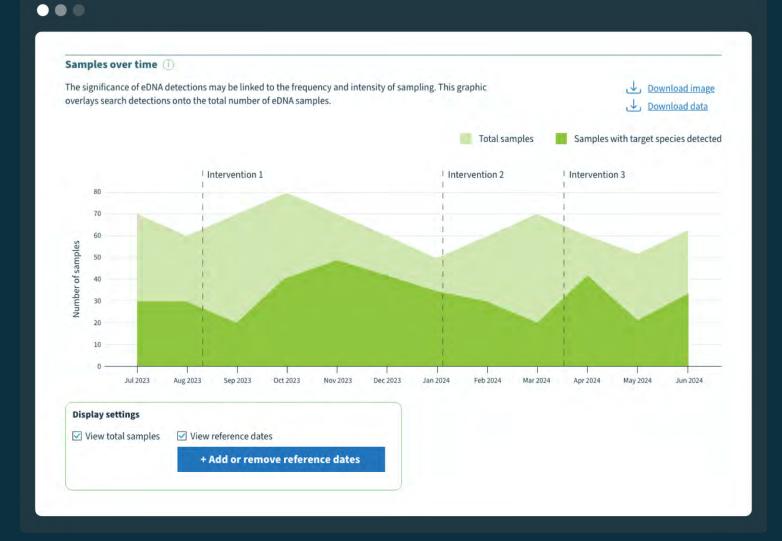
Spatial data



Animation & custom areas



Detection vs sampling effort



Downloads



How it works

- **Subscription model** to make it affordable for everyone.
 - Private data
 - Share your data with everyone or some
 - Individuals or organisations can join
- Bespoke solutions, where needed –
 e.g. Minderoo / Parks Australia
 public dashboards; research or
 project dashboards





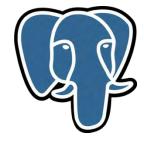


Under the hood

Under the hood

- SaaS platform built from the ground up (primarily) with R!
- Lots of moving parts under the surface of Ripple
- Interface built with R Shiny, JavaScript and CSS
- PostgreSQL relational database to store application metadata
- External relational databases for data providers
- AWS cloud hosting infrastructure









How it all fits together

• Some challenging tech problems, to name a few:

- Secure user authentication & dataset permission system
- Live data access to samples from multiple labs, metadata types
- Performant queries with large / complex searches

User authentication and permissions

- We needed a system for user management that is:
 - Available within cloud environment
 - Trusted identity provider
 - Secure features like 2FA and session management
- AWS Cognito ticks all of our boxes here!
- Users have a variety of user specific data, including:
 - Saved searches
 - Reference dates, custom areas
 - Dataset permissions
- User metadata managed within our relational database



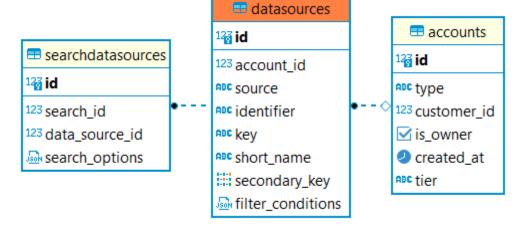
Cognito authentication in R

- Existing packages (e.g. CognitoR) opinionated and prescriptive
- PAWS packages require a Lot of manual configuration
- Secure authentication is becoming a more common request
- Created gRdian (coming soon) to standardise authentication in R
- Using R6 classes in gRdian to allow extendable / modifiable functionality
- In Ripple, we extend gRdian user class to contain extra user data
- gRdian user object contains all user app data, including dataset permissions



The Ripple dataset model

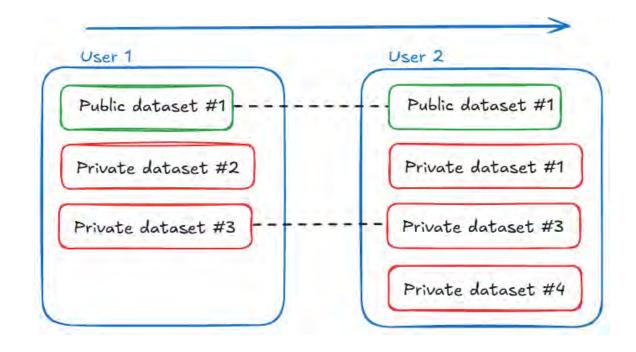
- Stored in the Ripple database with 'principle of least privilege' model
- Curated set of 'public' datasets available to all Ripple users
- Individuals can request access to a private dataset (with proven rights to that data) from our data providers
- JSON filter conditions to control data retrieval, e.g.
 - Job ID
 - Region
 - Owner





The Ripple dataset model

- Secure access to public / private eDNA data is critical to Ripple's operating model
- Each user has a set of dataset permissions explicitly assigned, separate DB entries
- Public datasets open access to any Ripple user



Unifying samples from multiple labs

- Multiple labs contributing data
 - Wilderlab over 20,000 primarily freshwater samples
 - Minderoo over 1200 marine samples
 - GBIF coming soon
- Some common attributes, some distinct to the lab
 - Wilderlab: six rep gold standard, environment type ...
 - Minderoo: site / sample depth
- Different data providers may use different taxonomy databases! (E.g. NCBI, APHIA, Fishbase)
- How do we unify all these together?









Unifying samples from multiple labs

- Create a DataProvider R6 object for generic 'data provider'
- R6 class contains default definitions for sample retrieval, e.g.:
 - Database access
 - SQL query construction
 - Data cleaning / processing
- Create lab-specific DataProvider object with customised functionality for each lab

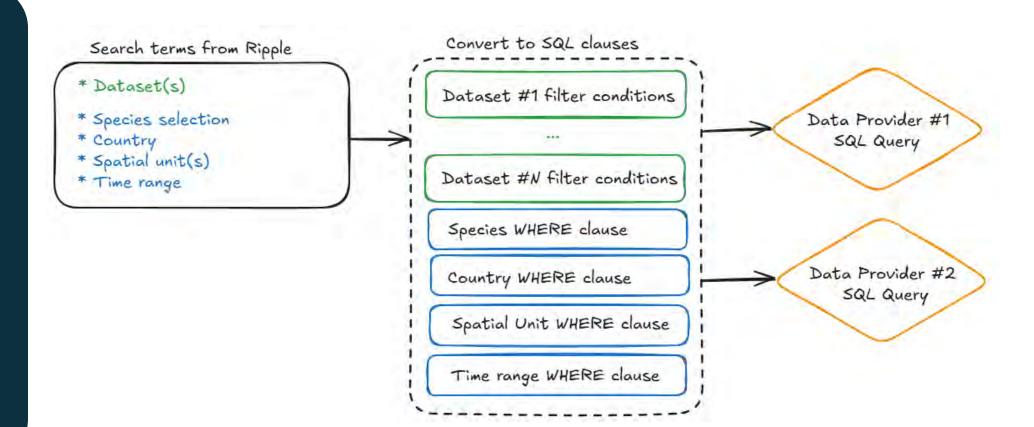


Unifying samples from multiple labs

- For each data provider, we:
 - Construct a SQL query that includes the filter conditions for each dataset
 - Retrieve the data one query per data provider
 - Clean the data into a consistent format
- Join results from each data provider into a single data frame in R
- Visualise unified data in Ripple outputs



Making the SQL query





Performance with complex searches

Ripple search is combining and processing a lot of complex data!

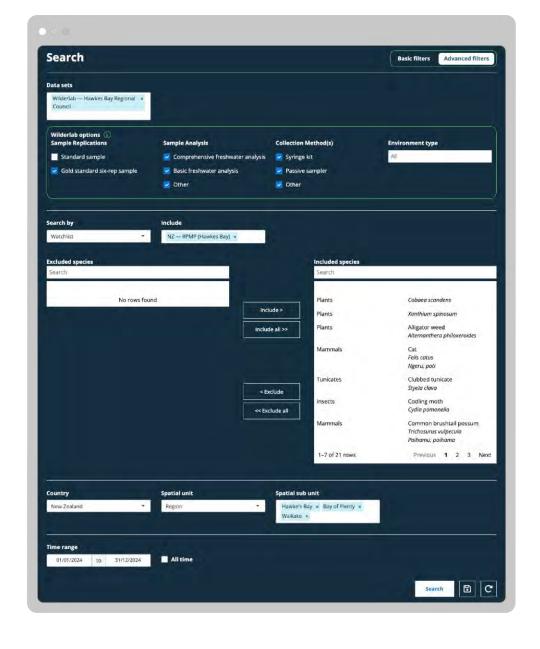
- eDNA Samples (and sample ownership)
- Species detection records (potentially thousands per sample)
- Taxonomy information (NCBI / APHIA IDs, naming conventions)
- Multiple data providers
- Spatial filtering (large river catchments, custom areas)



Performance with complex searches

- Lots of possible search parameters
- UX requirements (need to query all samples, then by species)
- Needs to query fast for most common searches!

Worst case scenario search:
 20,000 samples, 400,000 records





Performance with complex searches

- Indexing most common search terms:
 - Region
 - Owner
 - Public
- Samples 1:N Records, create a temp table of sample UIDs to speed up record query
- Interface design: encourage more specific search terms, enable saved searches to make it easy to re-run a search.



Maintaining Ripple

- Early Access app is still being actively developed and improved
- DevOps pipeline and workflow simplifies version management across multiple branches









New packages incoming!

- Rchv
 Managing database migrations within a complex codebase (coming soon)
- gRdian
 User authentication framework to connect identity providers with R and Shiny applications (coming soon)







Summary

- Excellent use case for power and flexibility of open source technology
- Very close collaboration with eDNA subject matter experts
- Ripple provides solutions for interpretation of high impact eDNA data
- Future plans: Bespoke instances (Parks Australia, Al integration, global data providers (GBIF) and more!





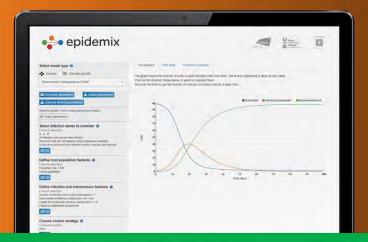
Democracy Tools

Democratic Performance

Index

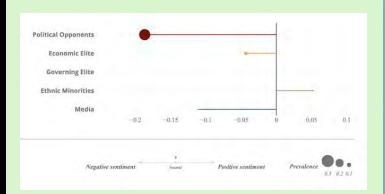
Tracking global democratic performance
 Analyzing the rhetoric of politicians
 Cataloging national-level referendums in democracies

Illuminating key areas of democracy









Political Rhetoric

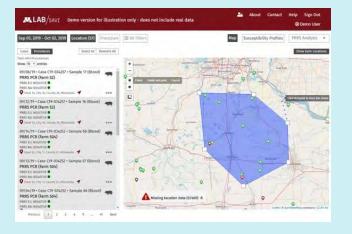
Explore ->

Analysis

 $Explore \longrightarrow$









Thank you!

Any questions, please get in touch! petra@epi-interactive.com nick@epi-interactive.com

Check out Ripple: https://rippledna.com/





