

HASS RDC Technical Advisory Group Meeting

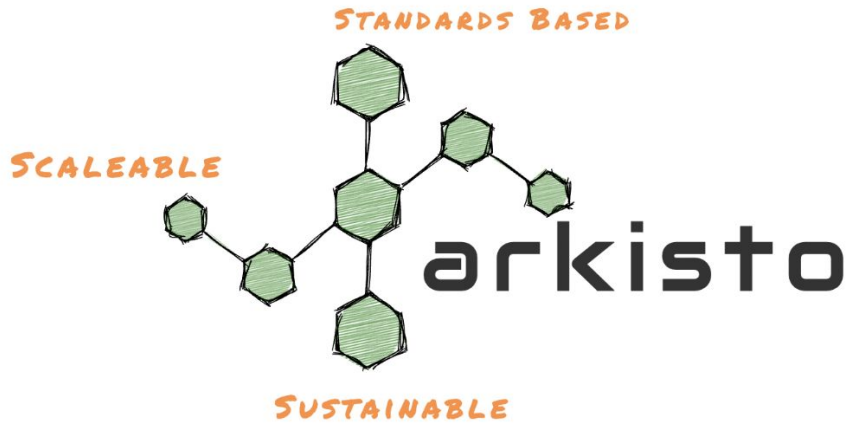
ATAP Architecture Discussion

Intro

Moises Sacal Bonequi - m.sacalbonequi@uq.edu.au

Ben Foley - b.foley@uq.edu.au

Peter Sefton - p.sefton@uq.edu.au



A scaleable, standards based platform
for sustainable data.

The basis of Arkisto is that the long-term preservability of well-described data is *always* the first consideration.

Data on an Arkisto deployment is always available on disc (or object storage) with a complete description *independently* of any services such as websites or APIs. Once the data is safe and well described, Arkisto has a flexible model for how data can be accessed using a variety of services.

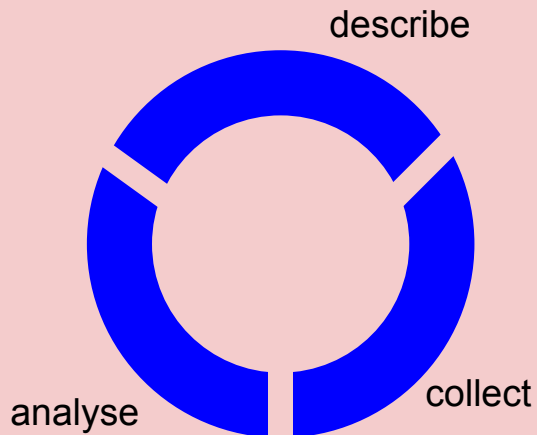
Arkisto is built on top of [Research Object Crate \(RO-Crate\)](#) and the [Oxford Common File System Layout \(OCFL\)](#).

With Arkisto there is no messy data migration.

Workspaces:

- working storage
- domain specific tools
- domain specific services

Research
Data
Management
Plan



Active cleanup processes
workspaces considered ephemeral

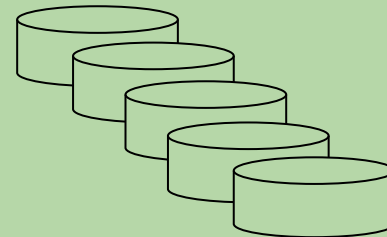
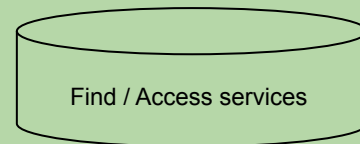
Reusable, Interoperable
data objects

- deposit early
- deposit often



reuse data objects

Repositories: institutional, domain or both



Findable, Accessible, Reusable
data objects



Policy based data management

Compute

HPC

Cloud

Desktop

Workspaces:

- working storage
- domain specific tools
- domain specific services



Active cleanup processes
workspaces considered ephemeral

Data Curation
& description

describe

analyse

collect

Data Cleaning

OCR / transcription
format migration

BYOData

STORAGE (including Cloudstor)

Identity Management

AAF / social media accounts

Archive & Preservation Repositories
institutional, domain or both

Harvested

external

PARADISEC

AU Nat. Corpus

AusLan (sign)

Sydney Speaks

ATAP Corpus
Reference, Training & BYO

ATAP Notebooks
Apps, Code, Workflows

... etc

Lang. portal(s)

Corpus discovery
Item discovery
Authenticated API
Create virtual corpora

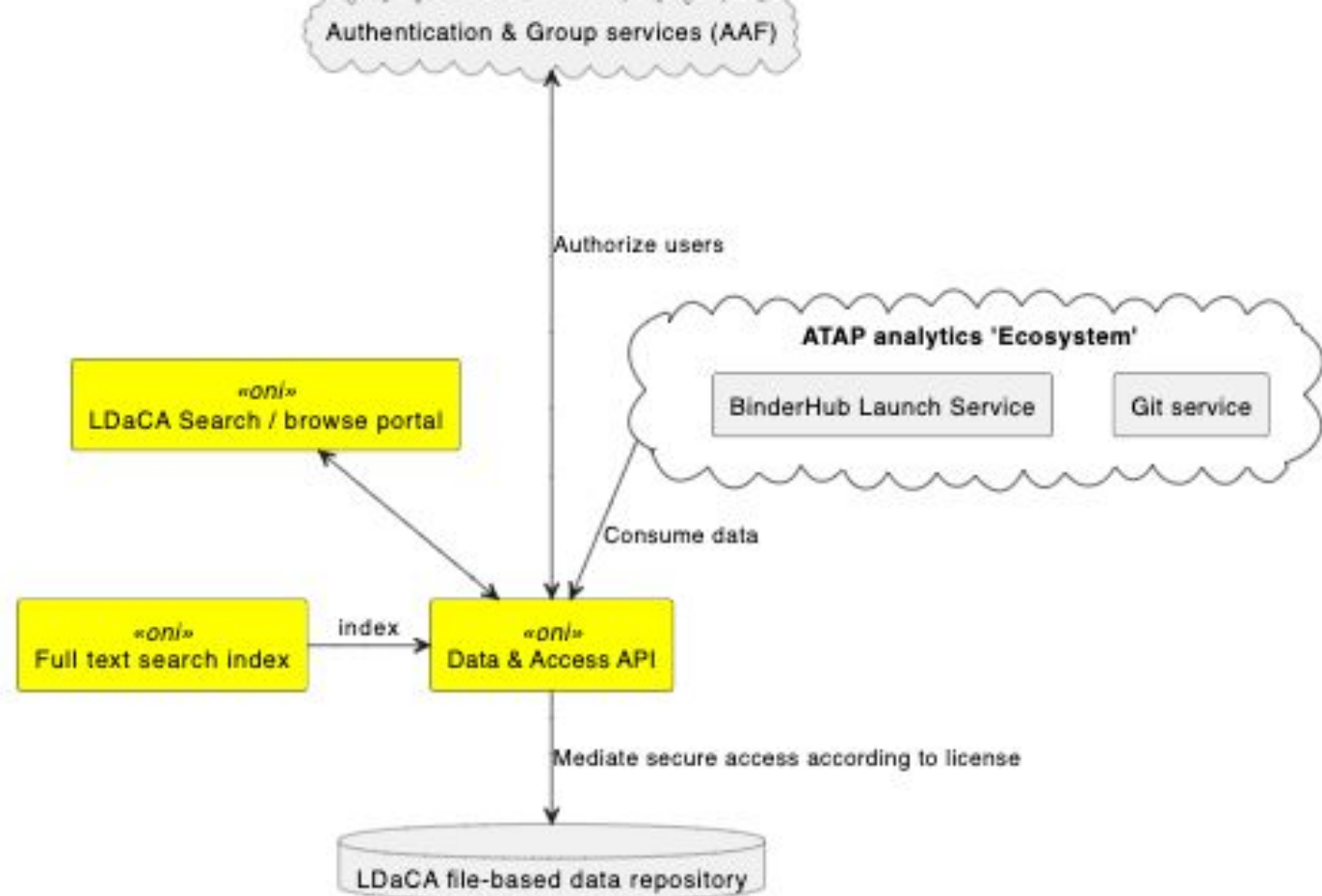
Deposit / Publish

Licence
Server

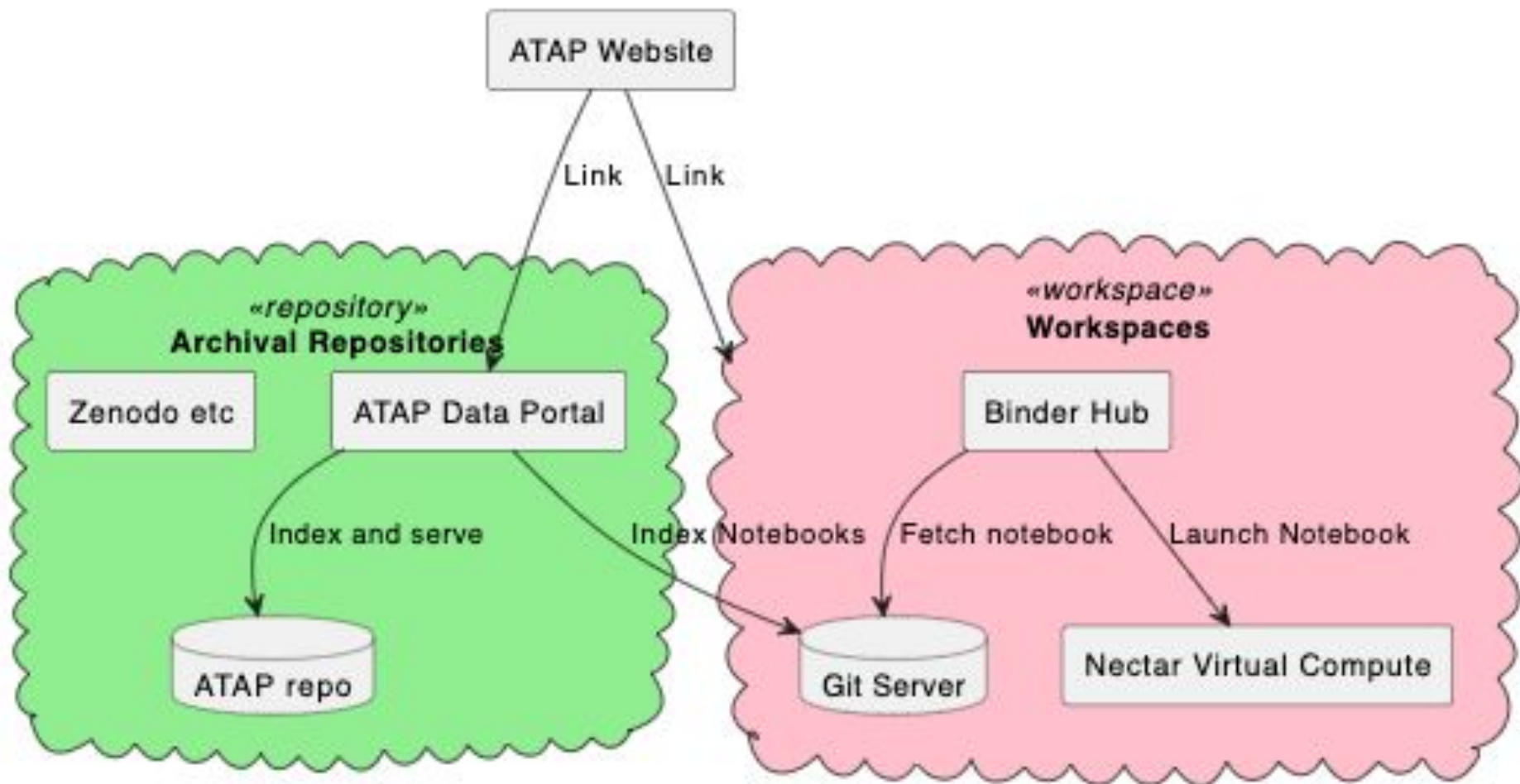
Reuse

Analytics
Portal

Code discovery
Launch / Rerun
Data Discovery
Authenticated API



ATAP Ecosystem - components



Demo

The screenshot shows a web browser displaying the Australian TextAnalytics Platform. The page features a search bar and navigation tabs for 'Collections' and 'Resources'. Two dataset cards are visible:

- A Corpus of Oz Early English (COOEE)**
 - Description: Material to be included had to meet with a regional and a temporal criterion. The latter required texts to have been produced between 1788 and 1900 in order to become eligible for COOEE. It was mandatory for a text to have been written in Australia, New...
 - Language: English: 4671
 - Linguistic Genre: Private Written: 810, Public Written: 465, Government English: 186, Speech Based: 147
 - Modality: Diagraphy: 4071
 - File Formats: text/plain: 2174
 - Access: [ES39u001.0239M36989JCCJBY](#), DOI: [10.26434/chemrxiv-2024-00000](#), Public responses: indexed
 - More
- Farms to Freeways Example Dataset**
 - Description: This data set was exported from an OpenAI Repository as an example of a DataChat. It contains the Collections and Items from the repository but does NOT have the attributions. The DOI resolved to an archive of the data elsewhere.
 - Language: English: 138
 - Linguistic Genre: Interview: 34
 - Modality: Diagraphy: 88, Speech: 34
 - File Formats: audio/mpeg: 68, application/pdf: 34, text/plain: 34
 - Access: [ES39u001.0239M36989JCCJBY](#), DOI: [10.26434/chemrxiv-2024-00000](#), Public responses: indexed
 - More

© 2022 LCoCA Program ATAP