# GLAM Collections as Data

Tim Sherratt · @wragge · timsherratt.org

These slides: https://slides.com/wragge/hussdata/

Watch the video: https://vimeo.com/350981680

# Trove

| All | Books | Pictures, photos, objects | Journals, articles and data sets | **Digitised newspapers and more** | Government Gazettes | Music, sound and video | Maps | Diaries, letters, archives | Archived websites (1996 – now) | People and organisations | Lists |

radio

**Search**

☐ Available online   ☐ Australian content   ☐ In my libraries   *Advanced Search*

## Refine your results:

▼ **Place**

New South Wales (1,161,749)

Queensland (689,006)

Victoria (382,009)

Western Australia (345,487)

South Australia (310,867)

ACT (180,429)

Tasmania (176,568)

International (42,504)

National (10,934)

Northern Territory (10,537)

▼ **Title**

The Canberra Times (AC...
(179,315)

The Sydney Morning Her...
(100,430)

The West Australian (P...
(97,400)

The Age (Melbourne, Vi...
(93,010)

The Argus (Melbourne, ...
(87,878)

## Digitised newspapers and more

Showing: 1 - 20 of at least **3,310,090** Refine search

Sort by: Relevance   Sort

**RADIO.**

*The West Australian (Perth, WA : 1879 - 1954)* **Tuesday 15 September 1925** p 6 Article

... *RADIO*. ?Mr. H. Brooghton Jensen, who was recently engaged to make an examination of the *Radio* mime ... 68 words

**RADIO**

*The Mail (Adelaide, SA : 1912 - 1954)* **Saturday 2 August 1930** p 17 Article Illustrated

... *RADIO* Variety will be the keynote of programmes by the Australian Broadcasting Company this week ... 128 words

**RADIO.**

*The Longreach Leader (Qld. : 1923 - 1954)* **Friday 26 September 1924** p 14 Article

... licensed: dealers'in all *radio* goods, and can quote Sydney prices.,If contemplating installing a set, why ... 111 words

**Radio**

*Queensland Figaro (Brisbane, Qld. : 1901 - 1936)* **Saturday 29 October 1927** p 9 Article

... the retail *Radio* shops in Brisbane. o Ireland has experienced three eras— The Pagan Era. The Christian ... 458 words

# Trove

| All | Books | Pictures, photos, objects | Journals, articles and data sets | Digitised newspapers and more | Government Gazettes | Music, sound and video | Maps | Diaries, letters, archives | Archived websites (1996 – now) | People and organisations | Lists |

radio

**Search**

☐ Available online    ☐ Australian content    ☐ In my libraries    Advanced Search

## Refine your results:

▼ **Place**

New South Wales (1,161,749)
Queensland (689,006)
Victoria (382,009)
Western Australia (345,487)
South Australia (310,867)
ACT (180,429)
Tasmania (176,568)
International (42,504)
National (10,934)
Northern Territory (10,537)

▼ **Title**

The Canberra Times (AC...
(179,315)
The Sydney Morning Her...
(100,430)
The West Australian (P...
(97,400)
The Age (Melbourne, Vi...
(93,010)
The Argus (Melbourne, ...
(87,878)

## Digitised newspapers and more

Showing: 1 - 20 of at least **3,310,090** Refine search

Sort by: Relevance    Sort

**RADIO.**
*The West Australian (Perth, WA : 1879 - 1954)* **Tuesday 15 September 1925** p 6 Article
... *RADIO*. ?Mr. H. Brooghton Jensen, who was recently engaged to make an examination of the *Radio* mime ... 68 words

**RADIO**
*The Mail (Adelaide, SA : 1912 - 1954)* **Saturday 2 August 1930** p 17 Article Illustrated
... *RADIO* Variety will be the keynote of programmes by the Australian Broadcasting Company this week ... 128 words

**RADIO.**
*The Longreach Leader (Qld. : 1923 - 1954)* **Friday 26 September 1924** p 14 Article
... licensed: dealers'in all *radio* goods, and can quote Sydney prices.,If contemplating installing a set, why ... 111 words
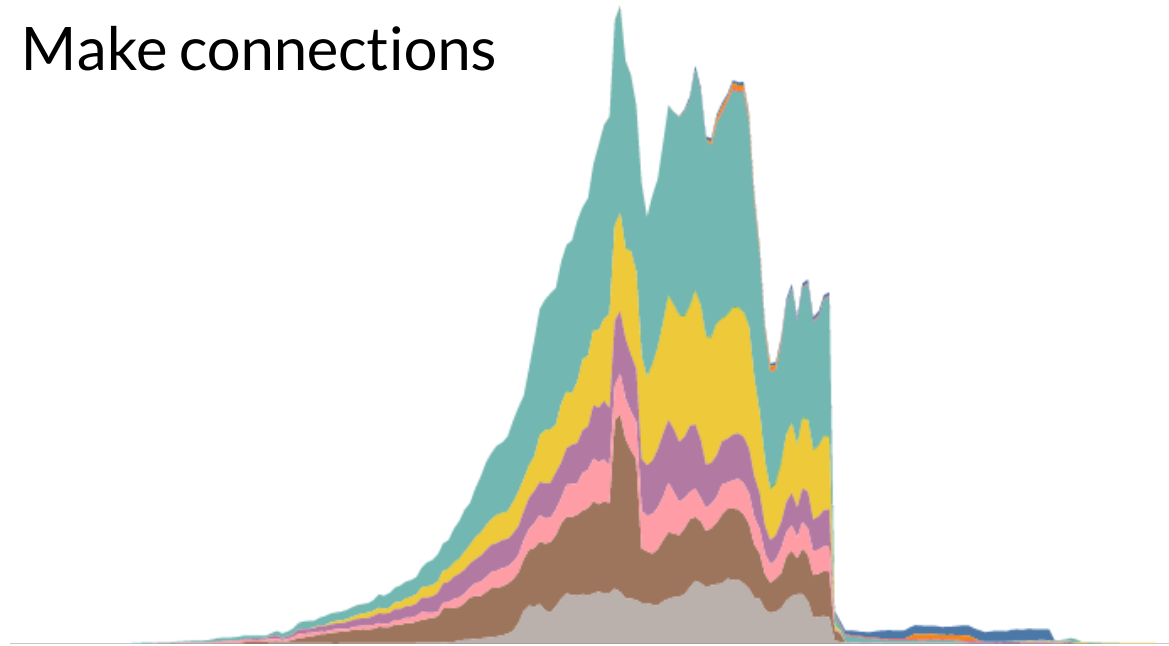
**Radio**
*Queensland Figaro (Brisbane, Qld. : 1901 - 1936)* **Saturday 29 October 1927** p 9 Article
... the retail *Radio* shops in Brisbane. o Ireland has experienced three eras— The Pagan Era. The Christian ... 458 words

# Collections as data

- Shift scales
- Find patterns
- Extract features
- Make connections

# Visualise Trove newspaper searches over time

You know the feeling. You enter a query into Trove's digitised newspapers search box and...

## Digitised newspapers and more

Showing: 1 - 20 of at least **3,308,203** Refine search

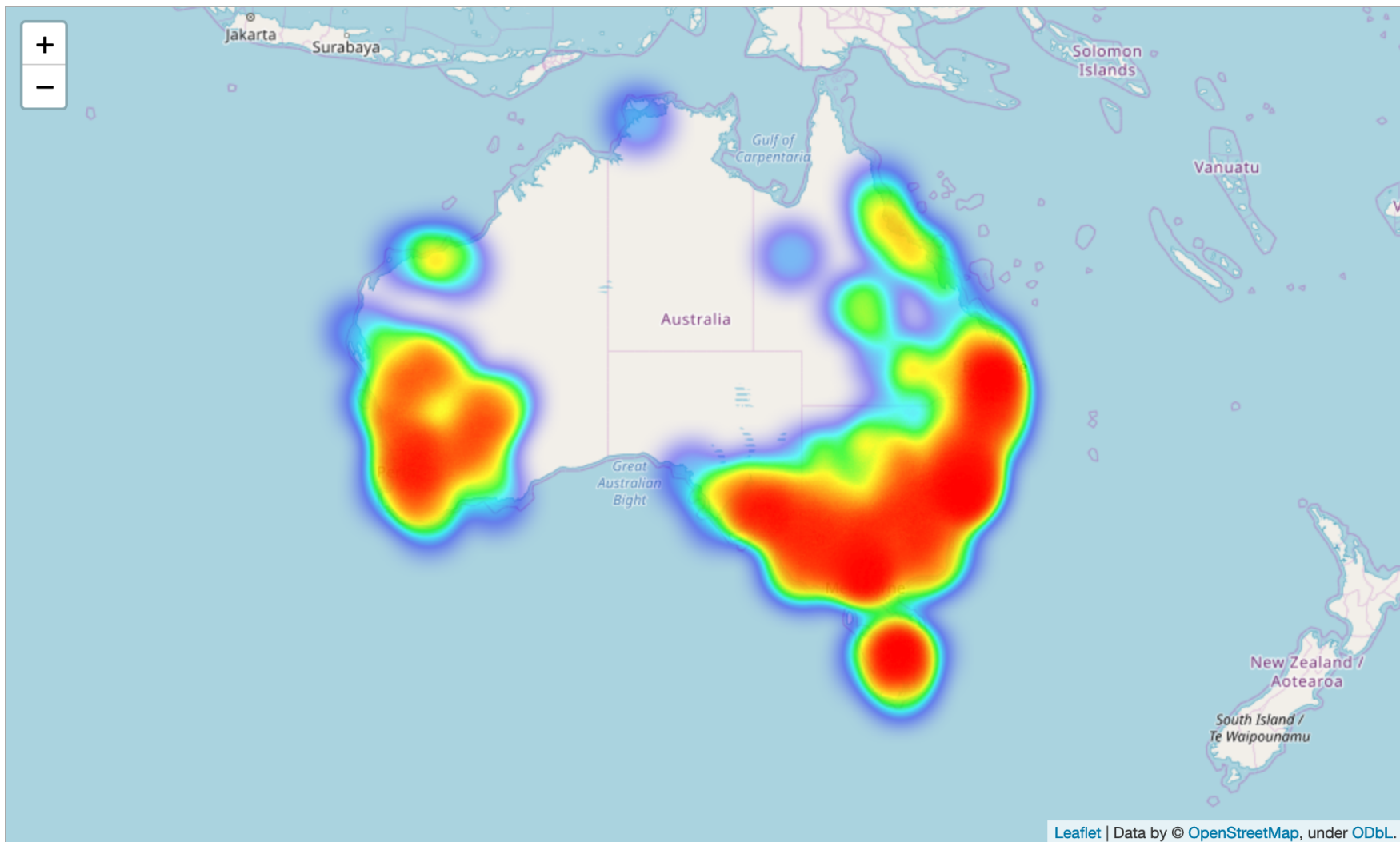Hmmm, **3 million results**, how do you make sense of that..?

Trove tries to be as helpful as possible by ordering your results by relevance. This is great if you aim is to find a few interesting articles. But how can you get a sense of the complete results set? How can you *see* everything? Trove's web interface only shows you the first 2,000 articles matching your search. But by getting data directly from the Trove API we can go bigger.

This notebook helps you zoom out and explore how the number of newspaper articles in your results varies over time by using the `decade` and `year` facets. We'll then combine this approach with other search facets to see how we can slice a set of results up in different ways to investigate historical changes.

1. Setting things up
2. Find the number of articles per year using facets
3. How many articles in total were published each year?
4. Charting our search results as a proportion of total articles
5. Comparing multiple search terms over time
6. Comparing a search term across different states
7. Comparing a search term across different newspapers
8. Chart changes in illustration types over time
9. But what are we searching?
10. Next steps
11. Related resources
12. Further reading

https://glam-workbench.github.io/trove-newspapers/

# Trove newspaper search by place of publication



https://glam-workbench.github.io/trove-newspapers/

# Newspaper articles with 'White Australia Policy' in their title

https://easyzoom.com/embed/139535

https://glam-workbench.github.io/trove-newspapers/

# Harvest newspaper articles in bulk

## Enter your search query

Use the Trove web interface to construct your search. Remember that the harvester will get **all** of the matched results, not just the first 2,000 you see in the web interface. Once you're happy with your search, just copy the url and paste it below.

Query url: `Enter the url of your search`

## Set harvest options

By default the harvester only saves the metadata (date, page, title, newspaper etc) from the search results. If you want to save the full text content of each article, just check the `Text` box. You can also save PDF copies of every article by checking the `PDF` option, but be warned that this will slow down your harvest and generate large download files. If you want to save PDFs from large harvests, you're probably better off installing and running the harvester on your own computer.

☐ Save full text

☐ Save PDFs (this can be slow)

**Start harvest**

Once your harvest is complete a link will appear to download the results as a single, zipped file. See this notebook for more information about the contents and format of the results folder.

You can also start to explore your results using this notebook.

Created by Tim Sherratt (@wragge) as part of the GLAM Workbench project.

If you think this project is worthwhile you can support it on Patreon.

https://glam-workbench.github.io/trove-harvester/

# Words before 'aliens' in Australian newspapers



http://timsherratt.org/blog/who-belongs/

# Jupyter notebooks

- computational narratives
- experimentation & learning
- tools & tutorials

# In the cloud or on your desktop

- Binder
- Tinker
- SWAN (part of Cloudstor)
- CoLab
- nteract

# GLAM Workbench



https://glam-workbench.github.io/

# Static or live



Visualise Trove newspaper searches over time

This notebook helps you zoom out and explore how the number of Trove newspaper articles in your search results varies over time by using the `decade` and year `facets`. We then combine this approach with other search facets to see how we can slice a set of results up in different ways to investigate historical changes.

view static

- Download from GitHub
- View using NBViewer
- Run live on Binder

run live

https://glam-workbench.github.io/trove-newspapers/

# What data?

# Sources of GLAM data

## GLAM datasets on data.gov.au

- Human readable list of GLAM datasets harvested from data.gov.au (July 2019)
- CSV formatted list of GLAM datasets harvested from data.gov.au (July 2019)
- CSV formatted list of GLAM datasets (CSVs only) harvested from data.gov.au (July 2019)

## Other downloadable datasets

### Full text

- Commonwealth Parliamentary Debates (Hansard), 1901-1980 (harvested from Parliamentary Library)
- Hansard interjections
- Australian Government Gazettes (1832-1968) (Trove)
- Federal Election speeches (Museum of Australian Democracy)
- Prime Ministers transcripts (harvested from DPMC)
- OCRd text from Trove digitised books (and ephemera) (harvested from Trove)
- OCRd text from the Internet Archive of 'Australian' books listed in Trove (harvested from Trove and Internet Archive)
- OCRd text of Trove digitised journals (harvested from Trove)
- Parliamentary press releases relating to immigrants and refugees (harvested from Trove and Parliamentary Library)
- Real Face of White Australia data (transcribed from National Archives of Australia: ST84/1)

### Images

- Editorial cartoons from The Bulletin, 1886 to 1952 (harvested from Trove)
- DIY #redactionart (harvested from NAA)
- Faces extracted from Trove newspaper photographs

https://github.com/GLAM-Workbench/glam-data-list

- Machine-readable but coding required for download

- Downloadable but not easily findable

- Structured but not downloadable

- Downloadable but in many separate pieces

# Machine-readable but coding required for download

https://glam-workbench.github.io/trove-lists/

# Machine-readable but coding required for download

https://glam-workbench.github.io/digitalnz/

# Downloadable but not easily findable

https://glam-workbench.github.io/glam-data-portals/

# Downloadable but not easily findable



https://glam-workbench.github.io/csv-explorer/

# Structured but not downloadable

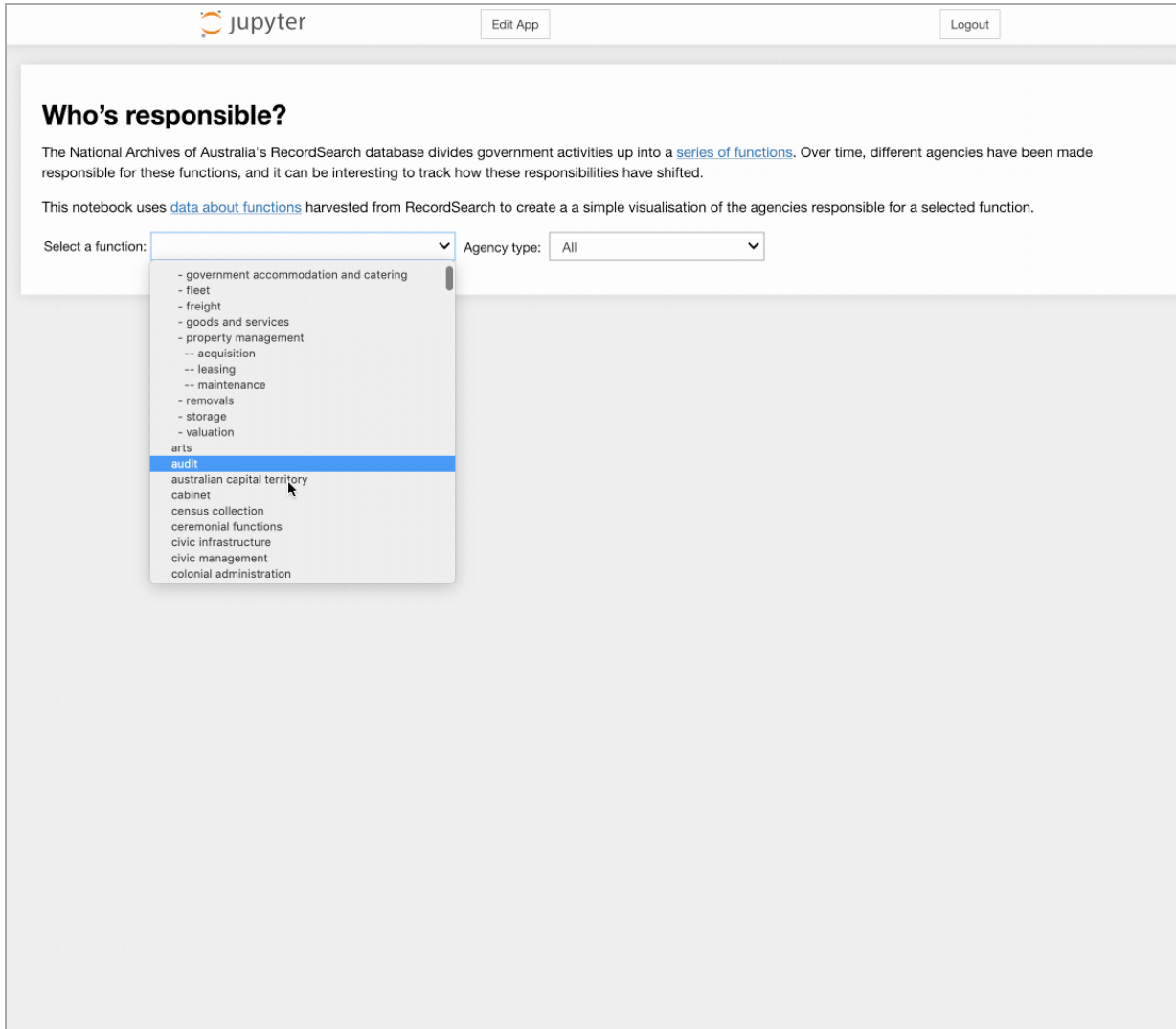https://glam-workbench.github.io/recordsearch/

# Structured but not downloadable

| series | total_items | date_from | date_to | Open | OWE | NYE | Closed | digitised_files | digitised_pages | % open | % digitised |
|--------|------------|-----------|---------|------|-----|-----|--------|-----------------|-----------------|--------|-------------|
| B13 | 20,194 | 1800 | 2005 | 19,786 | 8 | 400 | 0 | 354 | 5,043 | 97.98% | 1.75% |
| B6003 | 3 | 1904 | 1959 | 3 | 0 | 0 | 0 | 0 | 0 | 100.00% | 0.00% |
| BP343/15 | 2,571 | 1916 | 1955 | 2,566 | 0 | 5 | 0 | 85 | 176 | 99.81% | 3.31% |
| D2860 | 1 | 1902 | 1957 | 0 | 1 | 0 | 0 | 0 | 0 | 0.00% | 0.00% |
| D5036 | 1 | 1906 | 1935 | 1 | 0 | 0 | 0 | 0 | 0 | 100.00% | 0.00% |
| D596 | 11,395 | 1871 | 1971 | 2,983 | 31 | 8,381 | 0 | 185 | 3,031 | 26.18% | 1.62% |
| E752 | 722 | 1905 | 1941 | 719 | 0 | 3 | 0 | 717 | 9,310 | 99.58% | 99.31% |
| J2481 | 858 | 1897 | 1903 | 858 | 0 | 0 | 0 | 858 | 2,031 | 100.00% | 100.00% |
| J2482 | 799 | 1902 | 1912 | 799 | 0 | 0 | 0 | 798 | 3,153 | 100.00% | 99.87% |
| J2483 | 14,438 | 1903 | 1956 | 14,436 | 0 | 2 | 0 | 14,436 | 79,210 | 99.99% | 99.99% |
| J3115 | 161 | 1899 | 1928 | 161 | 0 | 0 | 0 | 161 | 1,344 | 100.00% | 100.00% |
| K1145 | 4,816 | 1900 | 1955 | 4,791 | 0 | 25 | 0 | 175 | 874 | 99.48% | 3.63% |
| P437 | 4,958 | 1901 | 1940 | 4,945 | 10 | 2 | 1 | 18 | 442 | 99.74% | 0.36% |
| P526 | 2 | 1909 | 1918 | 1 | 0 | 1 | 0 | 0 | 0 | 50.00% | 0.00% |
| PP4/2 | 613 | 1903 | 1947 | 610 | 0 | 3 | 0 | 28 | 1,512 | 99.51% | 4.57% |
| PP6/1 | 6,010 | 1906 | 1978 | 1,863 | 33 | 4,109 | 5 | 245 | 6,461 | 31.00% | 4.08% |
| SP11/26 | 27 | 1902 | 1902 | 27 | 0 | 0 | 0 | 5 | 84 | 100.00% | 18.52% |
| SP11/6 | 191 | 1902 | 1947 | 101 | 0 | 90 | 0 | 1 | 323 | 52.88% | 0.52% |
| SP115/1 | 1,787 | 1884 | 1943 | 1,787 | 0 | 0 | 0 | 9 | 285 | 100.00% | 0.50% |
| SP115/10 | 6 | 1884 | 1888 | 6 | 0 | 0 | 0 | 0 | 0 | 100.00% | 0.00% |
| SP42/1 | 16,256 | 1881 | 1960 | 15,525 | 0 | 731 | 0 | 3,253 | 45,862 | 95.50% | 20.01% |
| SP726/1 | 6 | 1902 | 1959 | 6 | 0 | 0 | 0 | 0 | 0 | 100.00% | 0.00% |
| ST84/1 | 2,765 | 1855 | 1975 | 2,758 | 0 | 7 | 0 | 434 | 13,979 | 99.75% | 15.70% |

https://glam-workbench.github.io/naa-wap/

# Structured but not downloadable



https://glam-workbench.github.io/recordsearch/

# Structured but not downloadable

https://glam-workbench.github.io/nsw-state-archives/

# Downloadable but in many separate pieces

https://player.vimeo.com/video/321657685?api=1

https://glam-workbench.github.io/hansard/

# Downloadable but in many separate pieces

https://glam-workbench.github.io/trove-journals/

# Downloadable but in many separate pieces