# code{4}lib
## JOURNAL

Mission     Editorial Committee     Process and Structure     Code4Lib

[                                                    ] [ Search ]

## The DSA Toolkit Shines Light Into Dark and Stormy Archives

*Themed web archive collections exist to make sense of archived web pages (mementos). Some collections contain hundreds of thousands of mementos. There are many collections about the same topic. Few collections on platforms like Archive-It include standardized metadata. Reviewing the documents in a single collection thus becomes an expensive proposition. Search engines help find individual documents but do not provide an overall understanding of each collection as a whole. Visitors need to be able to understand what individual collections contain so they can make decisions about individual collections and compare them to each other. The Dark and Stormy Archives (DSA) Project applies social media storytelling to a subset of a collection to facilitate collection understanding at a glance. As part of this work, we developed the DSA Toolkit, which helps archivists and visitors leverage this capability. As part of our recent International Internet Preservation Consortium (IIPC) grant, Los Alamos National Laboratory (LANL) and Old Dominion University (ODU) piloted the DSA toolkit with the National Library of Australia (NLA). Collectively we have made numerous improvements, from better handling of NLA mementos to native Linux installers to more approachable Web User Interfaces. Our goal is to make the DSA approachable for everyone so that end-users and archivists alike can apply social media storytelling to web archives.*

by Shawn M. Jones, Himarsha R. Jayanetti, Alex Osborne, Paul Koerbin, Martin Klein, Michele C. Weigle, Michael L. Nelson

*Editor's Note: This article makes use of Robust Links. Next to each hyperlink the reader will discover a menu that allows them to visit an archived version of the linked resource in case the current version has changed or is no longer available. Visit the Robust Links project for tools and more information on combating reference rot.*

## Web Archive Collections Are Too Large To Understand At A Glance

Web archives are invaluable for a variety of research studies. Historians ⤴ have analyzed how humans interacted on extinct websites, like Geocities. Social scientists ⤴ have used them to study the changes in social commerce over time. Journalists can use web archive evidence to bring attention to questionable medical practices ⤴ and document changes in government policy ⤴.

Some archivists create themed **web archive collections** by selecting web pages for preservation that support a topic. Each web page, or **original resource**, can change over time. Archivists capture these original resources at specific points in time, turning each observation into a **memento**. The date and time of capture is that memento's **memento-datetime**. A **TimeMap** contains the set of mementos for an original resource. Archive-It ⤴ is a popular platform for creating themed web archive collections. Themed collections also exist at the Library of Congress ⤴, Conifer ⤴, the Croatian Web Archive ⤴, the UK Web Archive ⤴, and the National Library of Australia's (NLA) PANDORA ⤴ and Trove ⤴ collections.
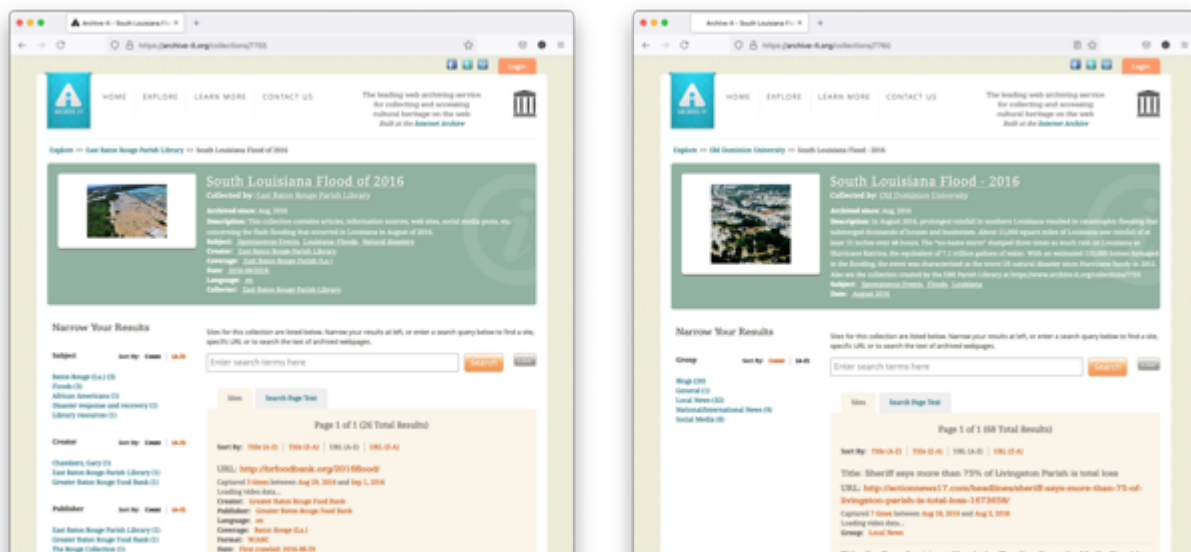
**Figure 1.** Archive-It has two collections about the South Louisiana Flood in 2016.
East Baton Rouge Parish Library created the one represented by the left screenshot 🔗
Old Dominion University created the one on the right 🔗. (Screenshots taken in June 2021.)

Such topical web archive collections are invaluable to those studying a topic, but it can be difficult for visitors to understand which collection they should begin exploring. Figure 1 shows screenshots of two Archive-It collections containing mementos about the South Louisiana Flood of 2016. Each collection contains different metadata fields. Which one should a researcher use in their project? Alternatively, if a visitor uses Archive-It's collection-level search engine to find collections about a topic like "human rights," they have 48 collections to review as of January 2022. Like *Government of Canada Publications* 🔗, some collections contain hundreds of thousands of documents but no metadata to help visitors understand the collection. In a 2019 study 🔗 (*preprint version* 🔗), we evaluated all collections at Archive-It and determined that collections with more original resources have less metadata to help visitors understand them.

Metadata is more consistent in the NLA's PANDORA collections than in Archive-It's collections, but the sheer size of some collections makes it difficult to understand them at a glance. Figure 2 shows a screenshot of the page for the PANDORA subject Politics 🔗. This subject is a collection in its own right and contains subcategories and collections so that visitors can view parts of the collection. It also includes a listing of 5,269 page titles. At a minimum, each PANDORA collection contains a collection title and page titles. Some are divided into sub-collections, but a human would still need to review many documents (or at least page titles) to understand the collection.
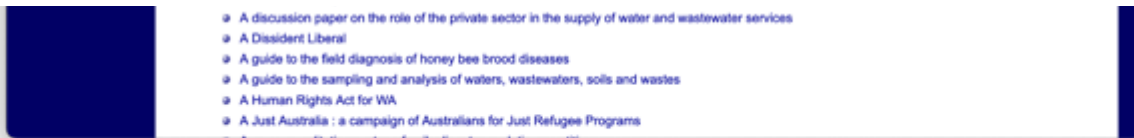
**Figure 2.** The PANDORA Subject *Politics* contains subcategories and collections, as well as a list of page titles.

## The DSA Toolkit For Summarizing Corpora

The Dark and Stormy Archives (DSA) Project ⚡ helps visitors achieve collection understanding by finding solutions that summarize web archive collections through visualizations that provide understanding at a glance. We apply social media storytelling as that visualization because most visitors are already familiar with the paradigm and thus require no additional training to understand how to consume these visualizations.

Social media storytelling consists of **surrogates** that summarize individual pages. Figure 3 shows the same page rendered as a surrogate in different web platforms. Most web users understand how to interpret the cards from search results. Many readers understand how to interpret social cards used by Facebook, Twitter, Tumblr, and other platforms. These surrogates take different shapes, but each attempts to summarize an individual document.
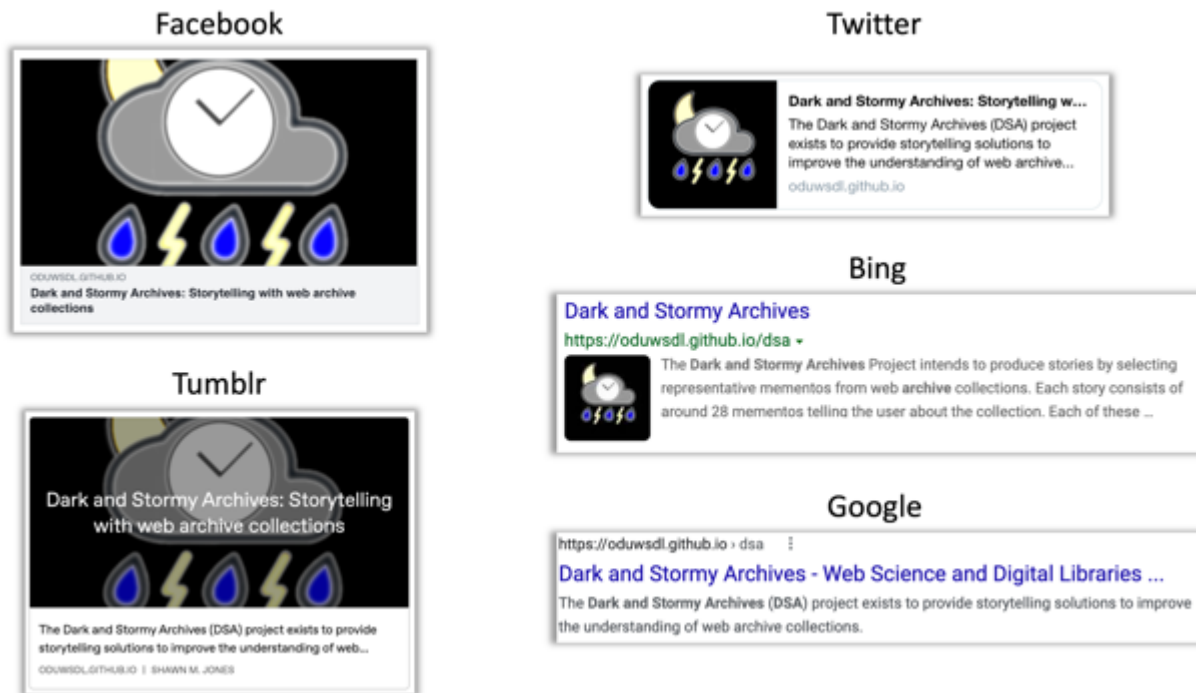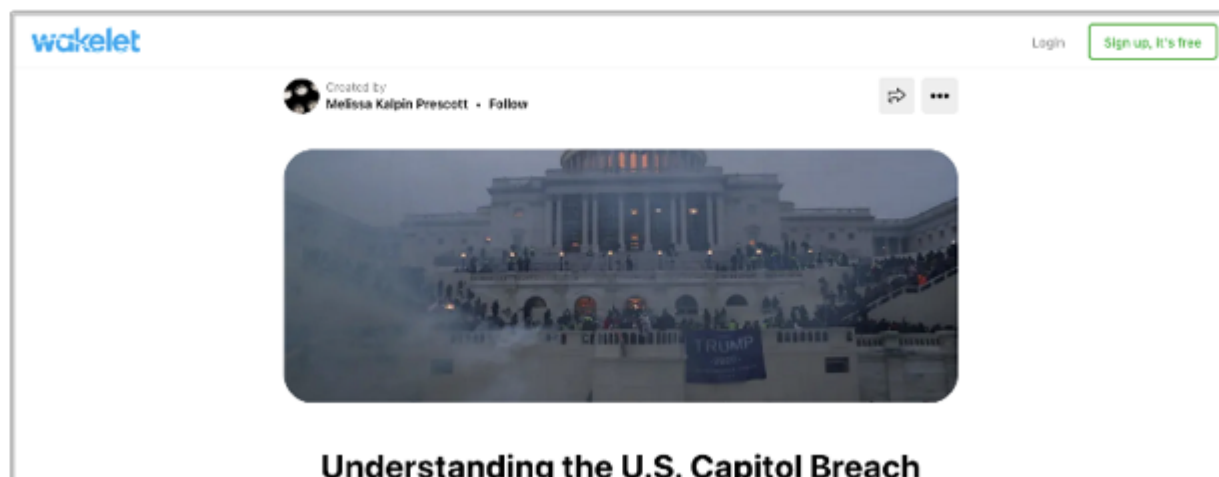
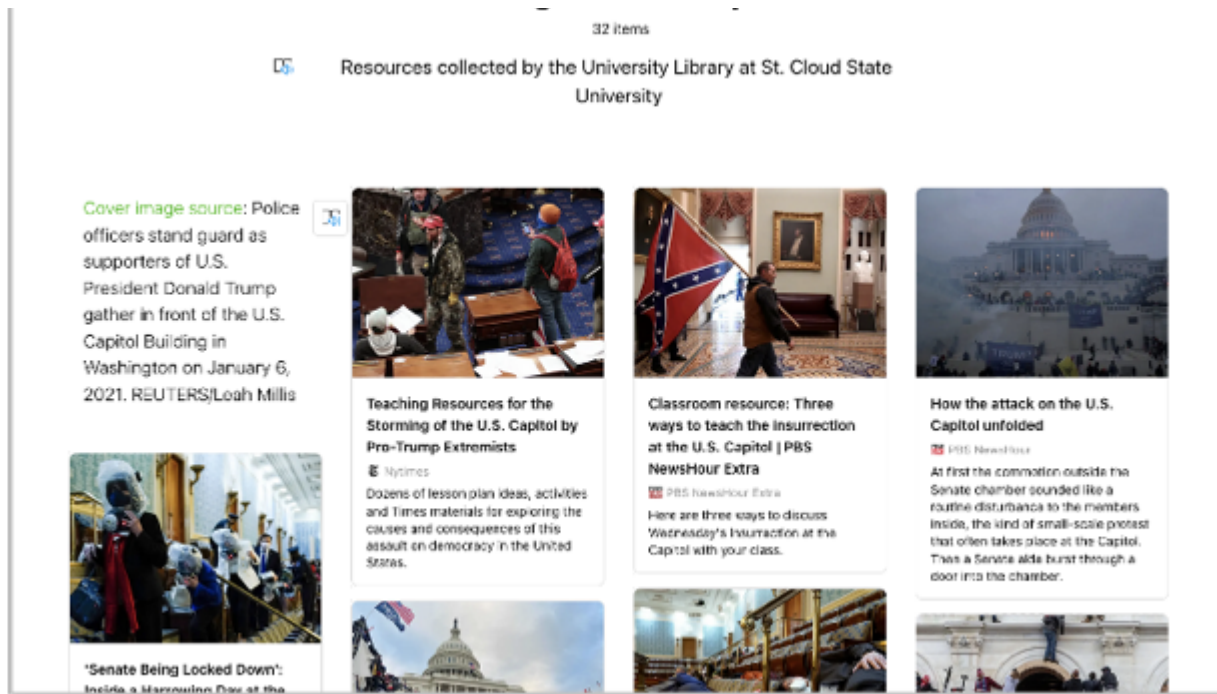**Figure 3.** Different surrogates for the same web page rendered by different platforms.

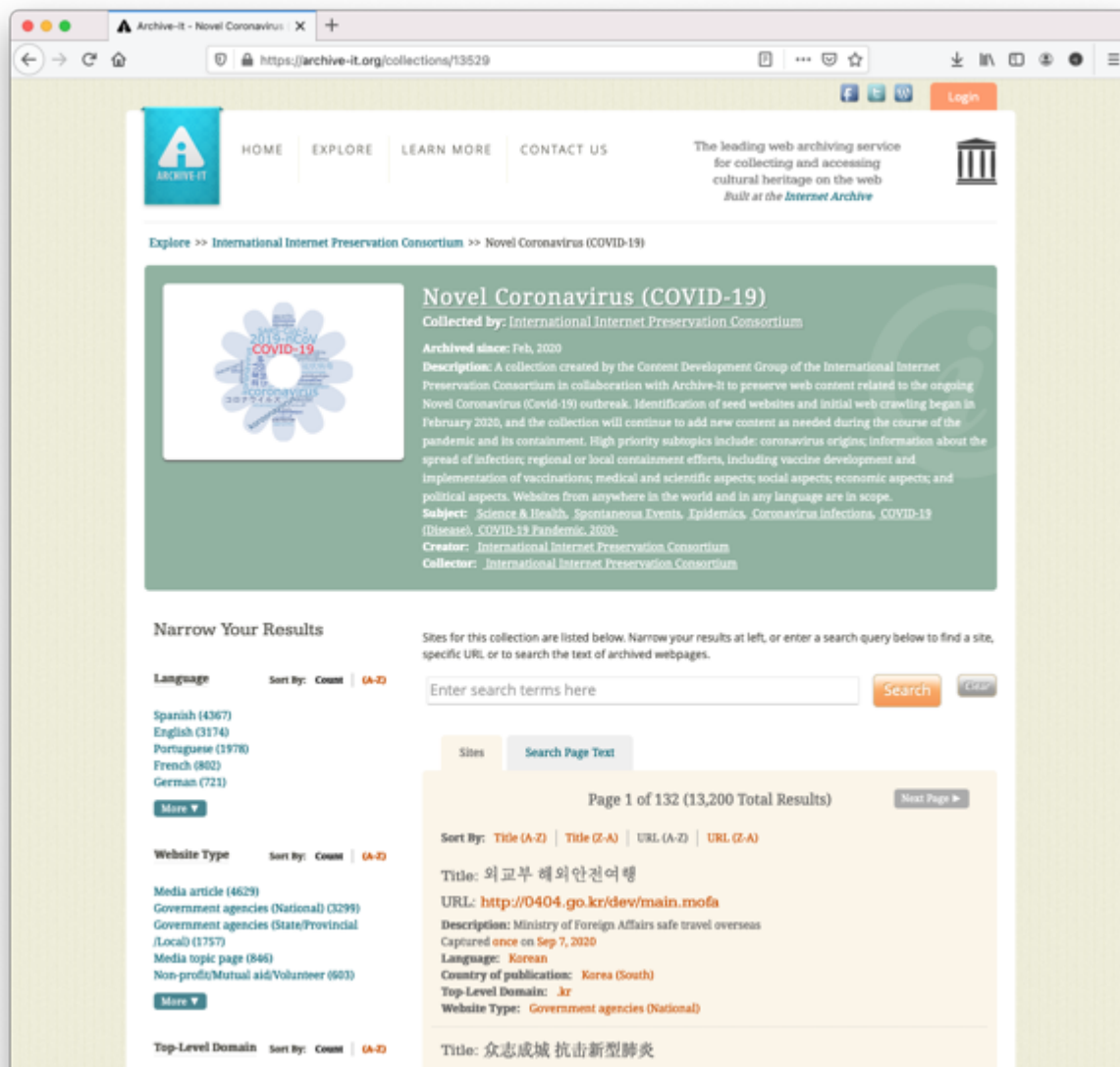**Figure 4.** A social media story on Wakelet about the 2021 US Capitol Breach.

**Figure 5.** A screenshot of the IIPC's Archive-It Collection *Novel Coronavirus (COVID-19)*.

When we combine a group of surrogates together as a unit, we create a **story** that summarizes a **topic**. Figure 4 shows a story that summarizes the 2021 US Capitol Breach. This story contains surrogates, in the form of social cards, each describing a different web resource discussing the attack. Archive-It and PANDORA have their own surrogates. In Archive-It's case, the surrogate consists of the URL captured and any metadata added by the archivist. In PANDORA's case, each surrogate is a page title. We want to extend this idea of the social media story further than is already accomplished by these web archiving platforms.

Figure 5 shows a 2019 screenshot of the International Internet Preservation Consortium's (IIPC) Archive-It Collection *Novel Coronavirus (COVID-19)* . This collection contains more than 23,000 mementos, far more than a single human can review, let alone understand at a glance. We applied the DSA toolkit  to produce a story summarizing this collection (Figure 6). The Archive-it collection contains metadata painstakingly provided by the archivist. We see page titles, countries of publication, languages, and, in some cases, descriptions of individual resources. There are even facets to help users explore different aspects of the collection. The story generated by the DSA Toolkit has people in masks, pictures of the virus, headlines, dates, sources, names, places, page summaries, maps showing the virus spreading across the world, and links back to the collection so visitors can explore that collection further.

**Figure 6.** A story by the DSA Toolkit summarizing IIPC's Archive-It Collection *Novel Coronavirus (COVID-19)* using sampling algorithms and social media storytelling.

The DSA Toolkit is built on our model of five storytelling processes 🔗, as shown in Figure 7. A user follows these processes to tell a story that summarizes a corpus.



**Figure 7.** The five processes for telling a story that summarizes a corpus.

As a summarization, our story provides a subset of the collection to a user. This subset consists of **exemplars** – documents that represent the collection well. Thus, we need to **select exemplars** from the collection to tell our story. Selecting exemplars can be done by humans or by various sampling algorithms. We can choose exemplars that attempt to summarize the collection as a whole or select those that feature a particular aspect of the collection, such as a specific time period, person, or source.

We **generate story metadata** to enrich our story for the visitor. Story metadata can consist of the collection name, who created the collection, and other collection metadata. The story from Figure 6 shows that story metadata can also contain entities, terms, images, and other content extracted from the collection.

Once our exemplars are selected, we **generate document metadata** to summarize each exemplar. In Figure 6, each of the social cards shown in the story visualizes this document metadata. Document metadata can take many forms, as needed by the storyteller.

We can then **visualize the story** in the desired medium by applying our story and document metadata. Most users will want a story that visitors can view on the web, and hence we have focused on providing visualizations that use HTML. We can also visualize stories in other media, such as video.

Finally, we can **distribute the story** for others to consume our visualization. Distribution can be from an author's website. It could also be via social media. Alternatively, a user could potentially print the story and manually hand it to someone. The point is that visitors cannot consume a story if they cannot access it.

As shown in Figure 8, the DSA Toolkit provides tools to help meet the goals of this model. **Hypercane** 🔗 helps users select exemplars and generate story metadata. **MementoEmbed** 🔗 generates document metadata. **Raintale** 🔗 helps users visualize the story by accepting input from Hypercane and MementoEmbed. Not shown in Figure 8 are two additional supporting packages. The library **AIU** 🔗 helps programmers identify the mementos and collection metadata from a web archive collection. The **OTMT** 🔗 (short for Off-Topic Memento Toolkit) identifies off-topic mementos in a collection. This paper highlights these tools and the results of their recent pilot with the NLA, thanks to a grant 🔗 from the IIPC.
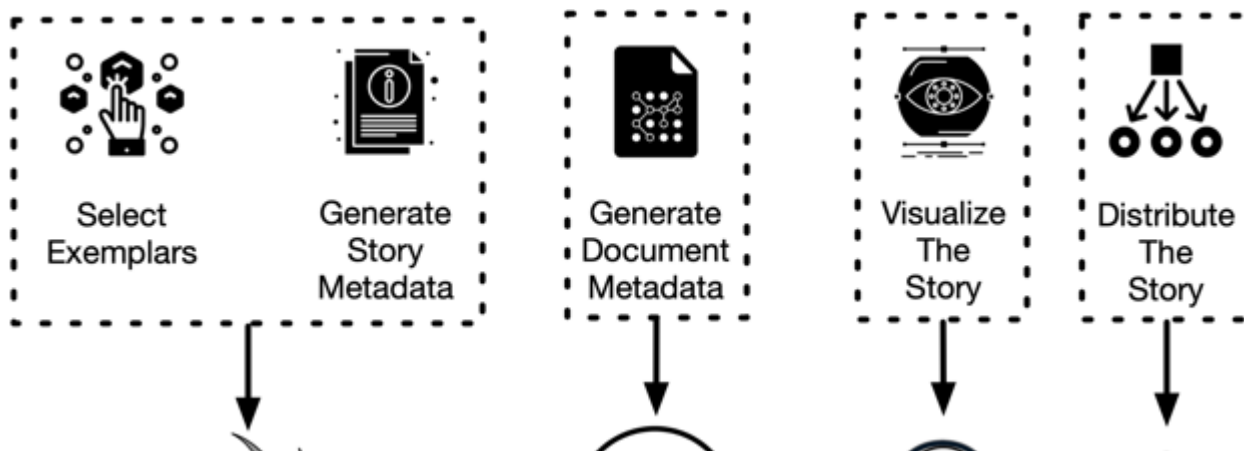
**Figure 8.** How the DSA Toolkit fits with the storytelling model shown in Figure 7.

## Other Efforts For Understanding Web Archive Collections

The DSA Toolkit is one of many tools for helping users understand web archive collections. For users with access to the Web ARChive (WARC) files that make up a web archive collection, ArchiveSpark 🔗 (by Holzmann et al. 🔗) and the Archives Unleashed Toolkit 🔗 (by Ruest et al. 🔗) offer an environment for exploring collections.

A user can load WARCs into ArchiveSpark and then perform a set of operations for generating data, such as extracting titles, calculating term distributions, discovering named entities, building a graph of outgoing hyperlinks, and extracting images. ArchiveSpark typically outputs a JSON file consisting of the desired data. The output of one operation in ArchiveSpark can feed into another, allowing users to chain them together and produce the desired result. It is up to the end-user to further process this data with a third-party tool.

Archives Unleashed Toolkit (AUT) is built on Warcbase 🔗 (Lin et al. 🔗). Based on work with web archive researchers, Lin et al. developed four steps for those trying to work with web archives: filter, analyze, aggregate, and visualize. AUT implements these four steps. With AUT, a user can filter based on various features, such as original resource URL, language, domain name, memento-datetime, or URL pattern. An AUT user can also generate reports on top-level domains, named entities, or link structure. Archives Unleashed Cloud (AUK) integrates a subset of this functionality with Archive-It, providing different reports for archivists to apply to their own collections. If an AUT user wishes to visualize the output, they must use a third-party tool.

These tools inspired the development of the DSA Toolkit. Hypercane allows users to chain together several web archive collection operations, just like ArchiveSpark and AUT. Raintale accepts the output of Hypercane's operations and visualizes it.

AUT and ArchiveSpark can potentially help archivists select exemplars and generate story metadata. They can also extract or generate some document metadata. Their focus, however, is not storytelling, so they rely on third-party tools for visualization and distribution.

The input type is another difference between the DSA Toolkit and these tools. Where ArchiveSpark and AUT accept WARCs as input, the DSA Toolkit does not. Because its goal is storytelling through summarization and visualization, the DSA Toolkit expects that all resources to be shared are accessible to the visitor of the story. Also, we wanted the DSA Toolkit to be usable by those who did not have access to a collection's WARCs. Thus, the DSA Toolkit uses memento URLs instead of WARC files as its input.

Jatowt et al. 🔗 may have been the first to visualize web pages over time. Their Page History Explorer downloaded mementos from web archives and generated screenshots ordered by memento-datetime as part of a timeline. Similarly, the TMVis project 🔗 (by Mabe et al. 🔗) accepts a single original resource URL, applies AlSum's Algorithm 🔗 to find the most novel mementos of that page, allows the user to reduce that set further manually, and then renders those mementos as a set of screenshots. These efforts work with a single original resource URL, not a collection of many different ones. However, they satisfy parts of our storytelling model by selecting exemplars, generating story metadata, generating document metadata (in the form of a screenshot), and visualizing the story – with Jatowt using a timeline and TMVis providing several different visual arrangements. With TMVis, users can distribute parts of their stories in various forms such as embeds, animated GIFs, and image sliders.

Padia et al. 🔗 generated a set of visualizations for web archive collections. Padia's visualizations did not focus on a single original resource URL but instead tried to consider different features of the mementos in the collection. They produced visualizations consisting of heat maps, timelines, word clouds, bubble charts, and treemaps. When possible, they leveraged the metadata provided by the collection, allowing them to cluster individual mementos as part of these visualizations. While they successfully created a set of visualizations to convey meaning about aspects of the collection, these visualizations required training for the viewer to understand them. Thus, they did not satisfy our goal of conveying meaning at a glance.
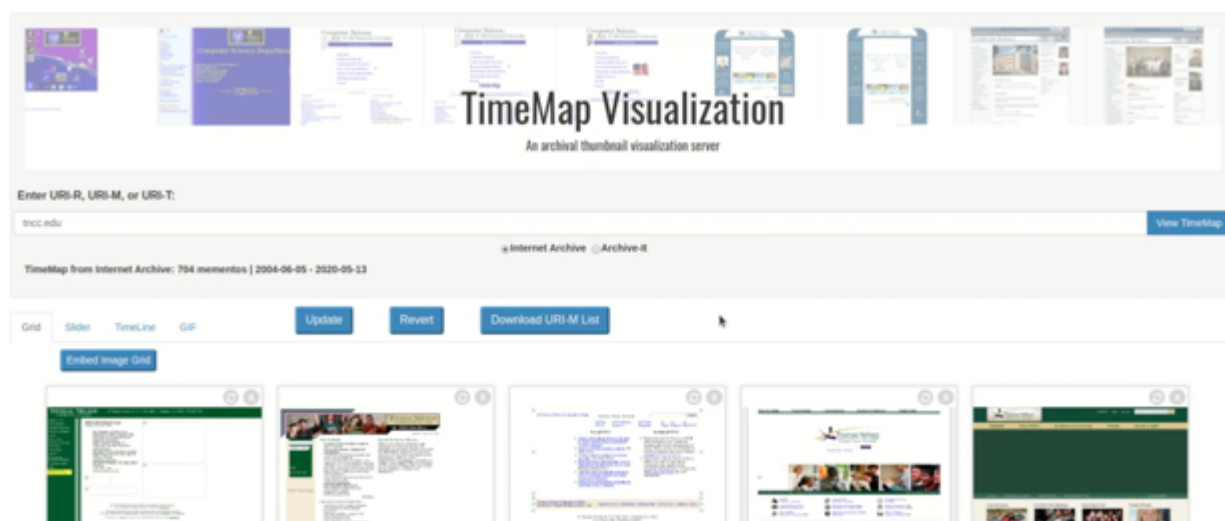
**Figure 9.** With TMVis, users can explore a single original resource over time and visualize thumbnails for that resource.

AlNoamany et al. pioneered combining storytelling with web archive collections. At the time of her work, the most popular social media storytelling service was the now-defunct Storify. She analyzed Storify stories and compared them with Archive-It collections to understand their similarities and differences. She identified that popular stories contain a median of 28 links, giving automated storytelling platforms a target number of items to select from a collection.

She identified some of the operations necessary to select exemplars automatically. Her algorithm filtered off-topic mementos, then duplicate mementos, and then non-English content from the collection. She also applied clustering, both by memento-datetime and by content. Finally, she scored the resulting mementos by Padia's web page category, McCown's URL path depth, and Brunelle's memento damage. Her algorithm, which the DSA Toolkit implements as *DSA1*, is the first of many possible methods for selecting exemplars. She demonstrated (*preprint version*) that participants *could* determine the difference between randomly generated stories and those created by human archivists through a user study. Equally important, she confirmed that participants *could not* distinguish between stories produced by her visualization algorithm and those created by human archivists.

Figure 10 shows a story generated by AlNoamany's proof-of-concept with Storify. Her story metadata was the collection name. Her proof-of-concept generated document metadata to fit into this platform. Storify handled distributing the story to visitors. Her research project was named *Dark and Stormy Archives*, and the current DSA project is a continuation of her work.

**Figure 10.** A story produced by AlNoamany's Dark and Stormy Archives proof-of-concept from the Archive-It collection *Russia Plane Crash* ⤸.

**Table 1.** Different web archive collection analysis efforts and suitability for summarizing collections with storytelling.

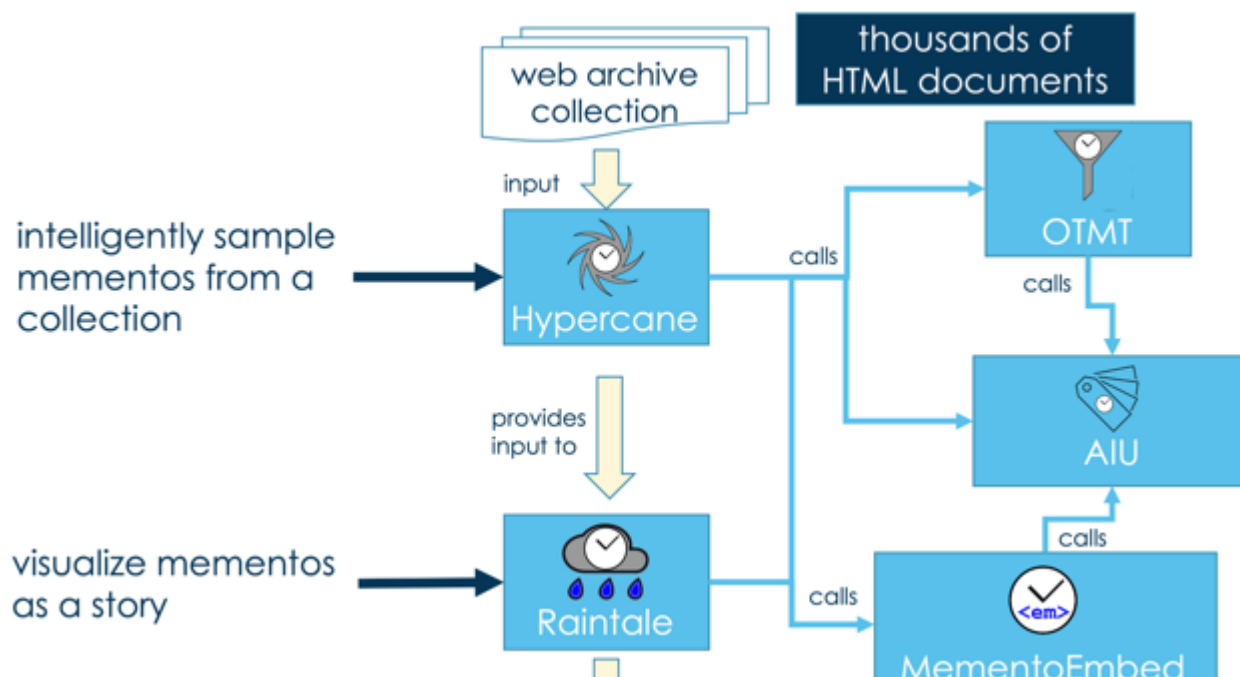| | | ArchiveSpark | Archives Unleashed Toolkit | Page History Explorer | TMVis | Padia's Visualizations | AlNoamany's Dark and Stormy Archives | DSA Toolkit (Current) |
|---|---|---|---|---|---|---|---|---|
| **Storytelling Model Support** | **Select Exemplars** | yes | yes | yes | yes | yes | yes | yes |
| | **Generate Story Metadata** | yes | yes | yes | yes | yes | yes | yes |
| | **Generate Document Metadata** | yes | yes | yes | yes | yes | yes | yes |
| | **Visualize the Story** | no | no | yes | yes | yes | yes | yes |
| | **Distribute the Story** | no | no | no | yes | not currently public, but could be | yes | yes |
| **Can User Customize Part of Storytelling Model?** | **Select Exemplars** | yes | yes | no | yes | no | yes | yes |
| | **Generate Story Metadata** | yes | yes | no | no | no | no | yes |
| | **Generate Document Metadata** | yes | yes | no | no | no | no | yes |
| | **Visualize the Story** | N/A | N/A | no | yes | no | no | yes |
| | **Distribute the Story** | N/A | N/A | no | no | no | no | yes |
| **Supports More Than One Original Resource** | | yes | yes | no | no | yes | yes | yes |
| **Allows Analysis Without WARCs** | | no | no | yes | yes | yes | yes | yes |
| **Web Archive Platforms Supported** | | any WARC files | any WARC files | 1 | 3 | 1 | 1 | any Memento-compliant archive |

**The DSA Toolkit**

**Figure 11.** The DSA Toolkit workflow for producing a story relies heavily on Hypercane and Raintale.

The DSA Toolkit consists of five different software projects that satisfy our storytelling processes. Figure 11 provides an overview of the workflow for these tools. A user executes Hypercane to select exemplars and generate story metadata. Then they feed that output into Raintale to visualize their story. Both Hypercane and Raintale consult the other components as needed. Here we introduce each tool and highlight our recent improvements thanks to the pilot with NLA.

## AIU

AIU (formerly Archive-It Utilities) is a Python library containing several classes that provide information about web archive collections. A programmer can invoke one of these classes and supply a collection identifier. The class includes several methods providing information like collection title, the URLs of the original resources in the collection, and other metadata, if present. AIU applies APIs and scraping HTML to gather its information.

AIU initially provided this information only for Archive-It collections. Archive-It collections contain metadata and a list of preserved original resources. From these original resources, we find TimeMaps. From these TimeMaps, we discover mementos. AIU follows this path to find the mementos in an Archive-It collection.

Below is an example iPython session where a user can extract information about an Archive-It collection (we measure the length of the list returned by the `list_seed_uris()` method for brevity, but it does produce a list of original resource URLs for the collection):

```
 1   In [1]: from aiu import ArchiveItCollection
 2
 3   In [2]: aic = ArchiveItCollection(5728)
 4
 5   In [3]: aic.get_collection_name()
 6   Out[3]: 'Social Media'
 7
 8   In [4]: aic.get_collectedby()
 9   Out[4]: 'Willamette University'
10
11   In [5]: aic.get_description()
12   Out[5]: 'Social media content created by Willamette University.'
13
14   In [6]: aic.get_collection_uri()
15   Out[6]: 'https://archive-it.org/collections/5728'
16
17   In [7]: len(aic.list_seed_uris())
18   Out[7]: 113
```
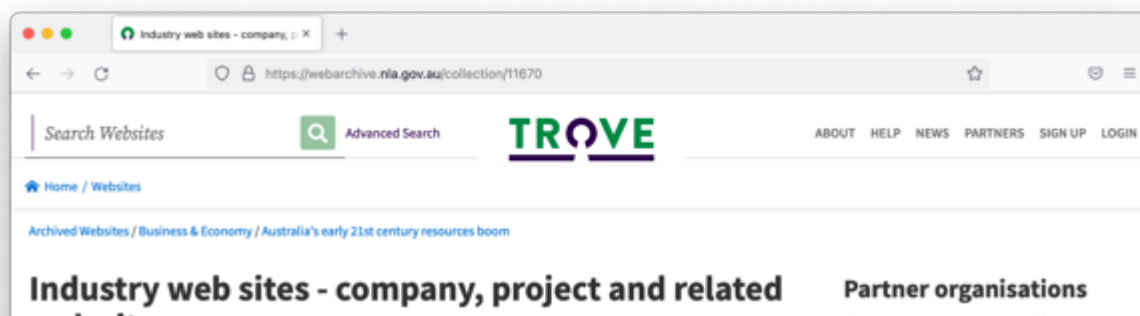
As part of our IIPC grant, we extended AIU to service the three types of collections at the NLA. as presented through their PANDORA and Trove websites. The PANDORA Archive is the curated selective web archive component of the larger Australian Web Archive (AWA). The NLA's entire web archive corpus (including PANDORA) is accessible through the Trove discovery service. While the AWA includes PANDORA's thematic sub-collections, the more extensive PANDORA subject listings are only viewable through the stand-alone PANDORA Archive website.

PANDORA subjects contain other PANDORA subjects listed as "subcategories." Each PANDORA subject includes a list of page titles, as seen in Figure 2. Each page title represents an original resource and links to a Title Entry Page (TEP). The TEP serves as a TimeMap that links to all mementos for that original resource. AIU will follow this chain to find similar metadata to Archive-It, including original resource and memento URLs.

PANDORA subjects also contain PANDORA collections. Like subjects, each PANDORA collection can include collections and page titles. Each page title links to a TEP which functions the same as with PANDORA subjects. AIU will follow these links to find metadata and URLs for PANDORA collections.

Trove collections are slightly different. Each Trove collection fans out into subcollections. Subcollections contain direct links to mementos, as shown in Figure 12. Now AIU includes a class that will follow these links, gather metadata, and capture URLs as needed.

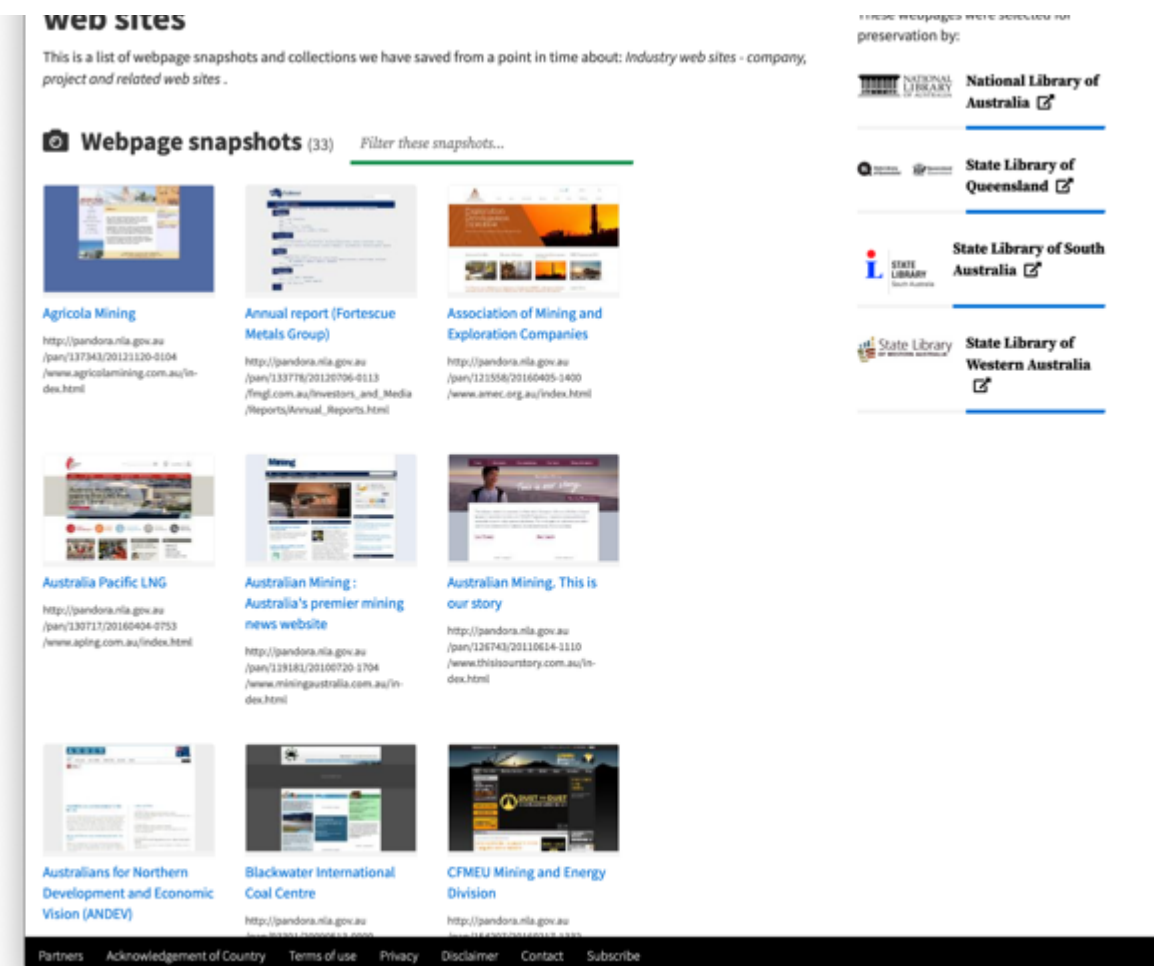Please view AIU's documentation ⤢ for more information on applying it to gather collection information.

**Figure 12.** A screenshot of a Trove collection.

### OTMT

The Off-Topic Memento Toolkit ⤢ (OTMT) is an application and library that supports different text similarity measures for identifying which mementos in a collection are on-topic or off-topic. It is used as a library by Hypercane to perform this same task. The OTMT received bug fixes as part of our IIPC grant work.

### Hypercane

Hypercane ⤢ is a complex application that provides the user with many ways to select exemplars and generate story metadata automatically. Hypercane exists in the command line application `hc`. A user selects exemplars with the `sample` action of that application. For example, to select exemplars from Archive-It collection 694 using the DSA1 algorithm, a user types:

```
1  # hc sample dsa1 -i archiveit -a 694 -o story-mementos.tsv
```

Generating story metadata is handled by Hypercane's `report` action. For example, to generate a list of named entities from a file containing a list of mementos, a user types:

```
1  # hc report entities -i memento -a story-mementos.tsv \
2      -o entity-report.tsv
```

Finally, to create a rich story, a user can `synthesize` the output of other Hypercane commands into other formats. To create a rich JSON story file for Raintale, a user types:

```
1  # hc synthesize raintale-story -i mementos -a story-mementos.tsv \
2      --imagedata ${image_report} \
3      --title "My Story for Collection X" \
4      --termdata terms-report.tsv \
5      --entitydata entity-report.tsv \
6      -o mystory.json
```

The `sample`, `report`, and `synthesize` actions provide Hypercane's top-level functionality. Their capabilities are not limited to what we show in the examples

above. The `sample` action currently supports fourteen different sampling algorithms, denoted by their name – e.g., a user can execute stratified random sampling with `hc sample stratified-random`, or they can run the DSA3 algorithm from our previous work with `hc sample dsa3`. Likewise, the `report` action supports ten different types of reports that can produce helpful metadata for storytelling, like extracted collection metadata, image analysis, or the list of top phrases for the input. The `synthesize` action gives the user different output format options, like a Raintale story file or WARCs.

As seen above, each Hypercane command supports the arguments `-i` (for input type) and `-a` (for input argument). The input type instructs Hypercane on the nature of the input argument. As seen above, if the input type is `archiveit`, then the input argument must be an Archive-It collection identifier. If the input type is `mementos`, then the input argument must be a file containing a list of memento URLs.

Hypercane currently supports the following input types:

- `mementos` – for a file containing a list of memento URLs

- `timemaps` – the input argument is a file containing a list of TimeMap URLs

- `original-resources` – the input argument is a file containing a list of original resource URLs

- `archiveit` – the input argument is an Archive-It collection identifier

- `pandora-subject` – the input argument is a PANDORA subject identifier, added as part of the IIPC grant work

- `pandora-collection` – the input argument is PANDORA collection identifier, added as part of the IIPC grant work

- `trove` – the input argument is a Trove collection identifier, added as part of the IIPC grant work

Hypercane supports all Memento-compliant archives. Thus, a user need not rely upon mementos from Archive-It or NLA. The `mementos` and `timemaps` input types allow users to submit URLs from other archives, such as the Internet Archive's Wayback Machine, Archive.Today, or Arquivo.pt.
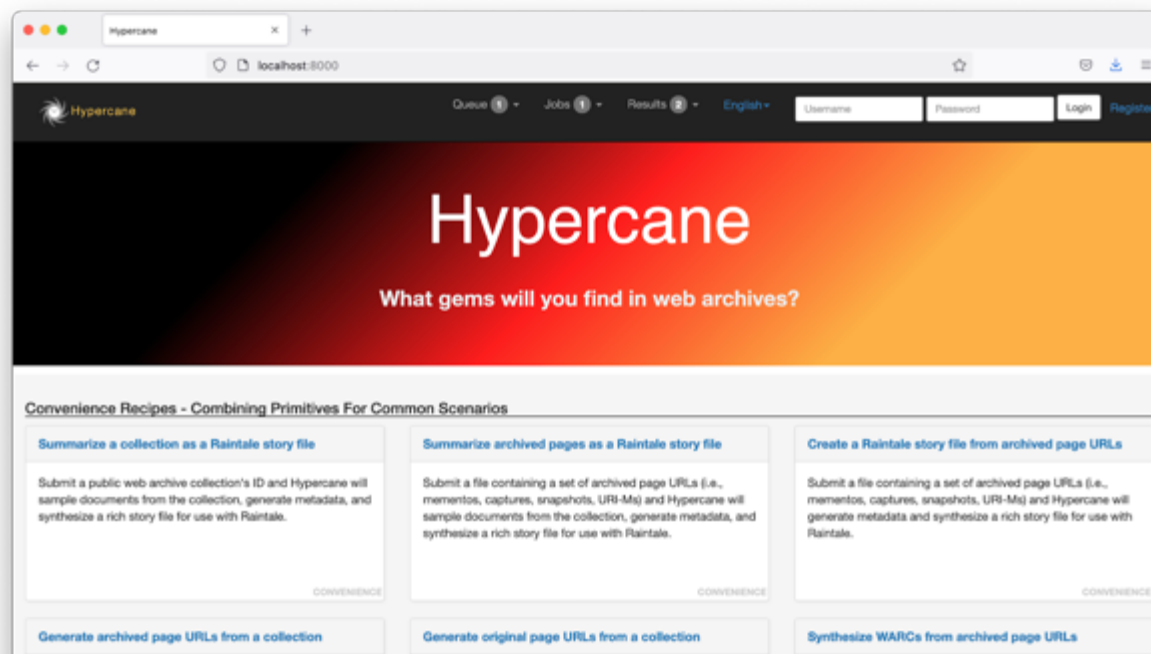
Every Hypercane command also supports the `-o` argument to specify an output file. In most cases, the output file is a list of memento URLs. This feature allows the user to chain Hypercane commands together so that the output of one command can feed into another.

Hypercane's sample action helps users execute existing sampling algorithms, but users also have access to the following algorithmic primitives that allow them to build their samples:

- `identify` – identifies the memento, TimeMap, or original-resource URLs for a given input type

- `filter` – produces a list of mementos that match the provided criteria (e.g., are on-topic, have unique content, are written in a specific language)

- `cluster` – clusters the collection via a given algorithm, such as through LDA topic modeling or K-means clustering by memento-datetime

- `score` – scores each memento in the input by a given function

- `order` – sorts the input by a given feature, like memento-datetime

By combining these primitives, users can create powerful scripts for sampling from collections, allowing them to tell many types of stories. We will not describe all of these primitives in detail for brevity but instead, direct interested readers to Hypercane's official documentation and our recent publications.

As part of the IIPC grant, we also sought to make Hypercane more approachable to administrators and end-users. Hypercane is a Python application and also depends on MongoDB for caching. We are evaluating native Linux installers for CentOS 8 and Ubuntu 21.10 that handle installing dependencies and provide convenience scripts for administering the application.
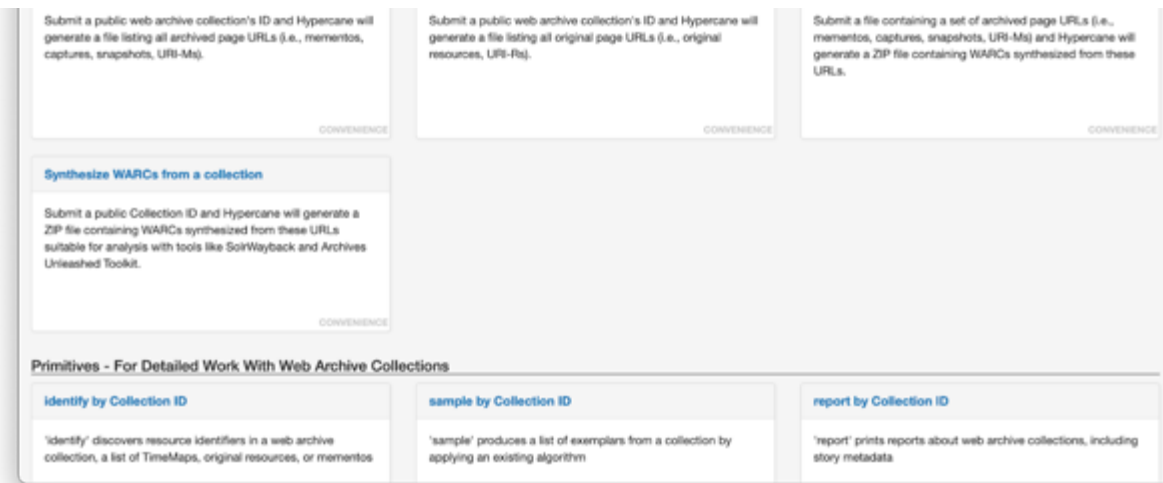
**Figure 13.** A screenshot of Hypercane's Web User Interface landing page.

We developed a new Web User Interface (WUI) so that Hypercane was approachable for users not familiar with scripting or working with command-line interfaces. Figure 13 shows a screenshot of this WUI. We wanted to create a web user interface while still preserving the existence and capabilities of the command-line application for scripting. We built Hypercane's WUI with the Wooey 🔗 library, which creates web user interfaces from command-line applications, satisfying our needs. With Wooey, a user can fill out a web form to submit a Hypercane command as a job and come back later to retrieve the results. Wooey also provides authentication and separation of jobs to protect users' privacy as a Django application.

Hypercane's `sample` action allows users to run pre-existing sampling algorithms to get a list of mementos for their story. As noted above, it supports ten sampling algorithms, including AlNoamany's. But what if a user wants to run `sample`, execute all `report` actions necessary to generate different types of story metadata, and finally `synthesize` a JSON file for use in a rich Raintale story? Before the WUI, we suggested that they write a script that executes these commands in order. The WUI presents **convenience recipes** that help users run complex tasks with minimal input and no scripting. Hypercane's WUI currently supports seven convenience recipes as web forms.
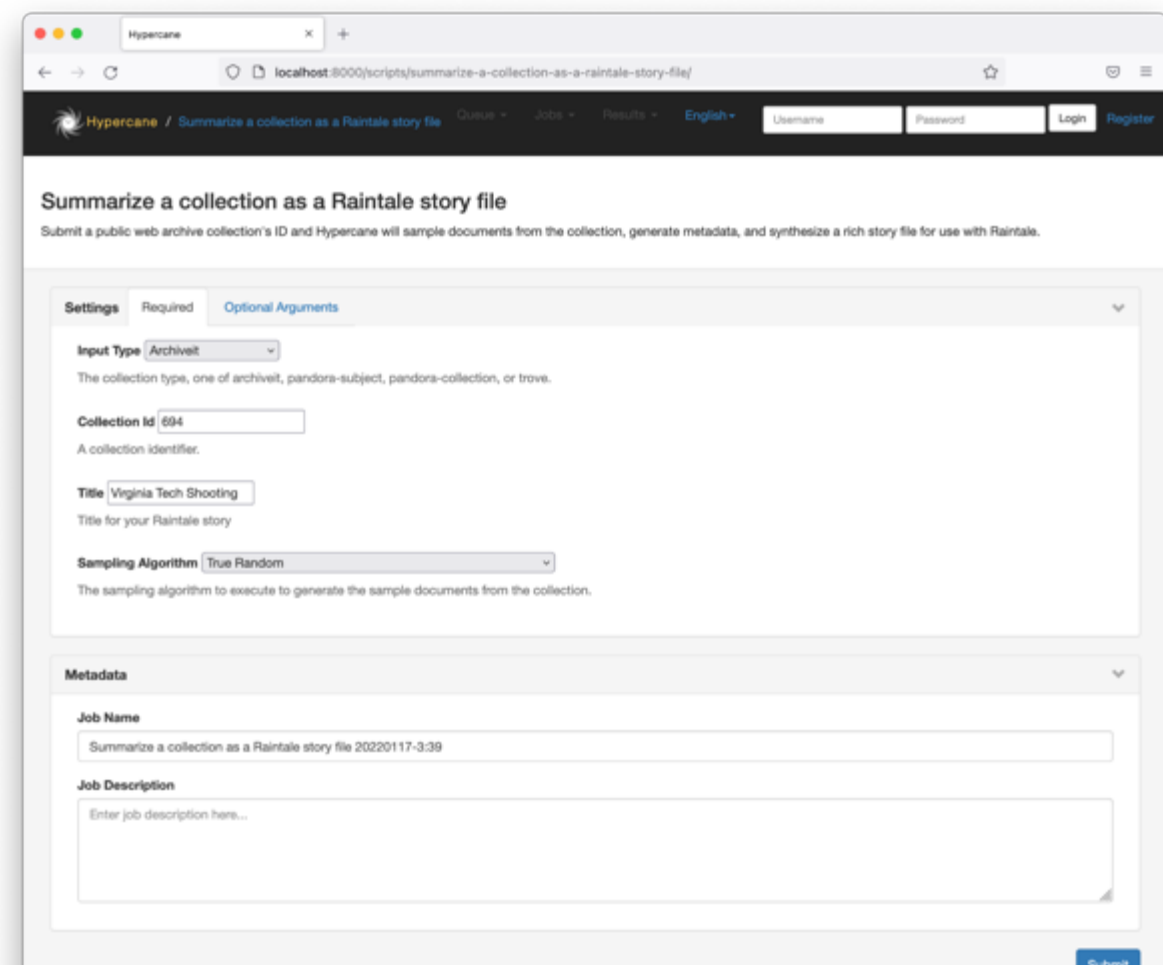
**Figure 14.** A screenshot of the Hypercane web form for the recipe "Summarize a collection as a Raintale story file."
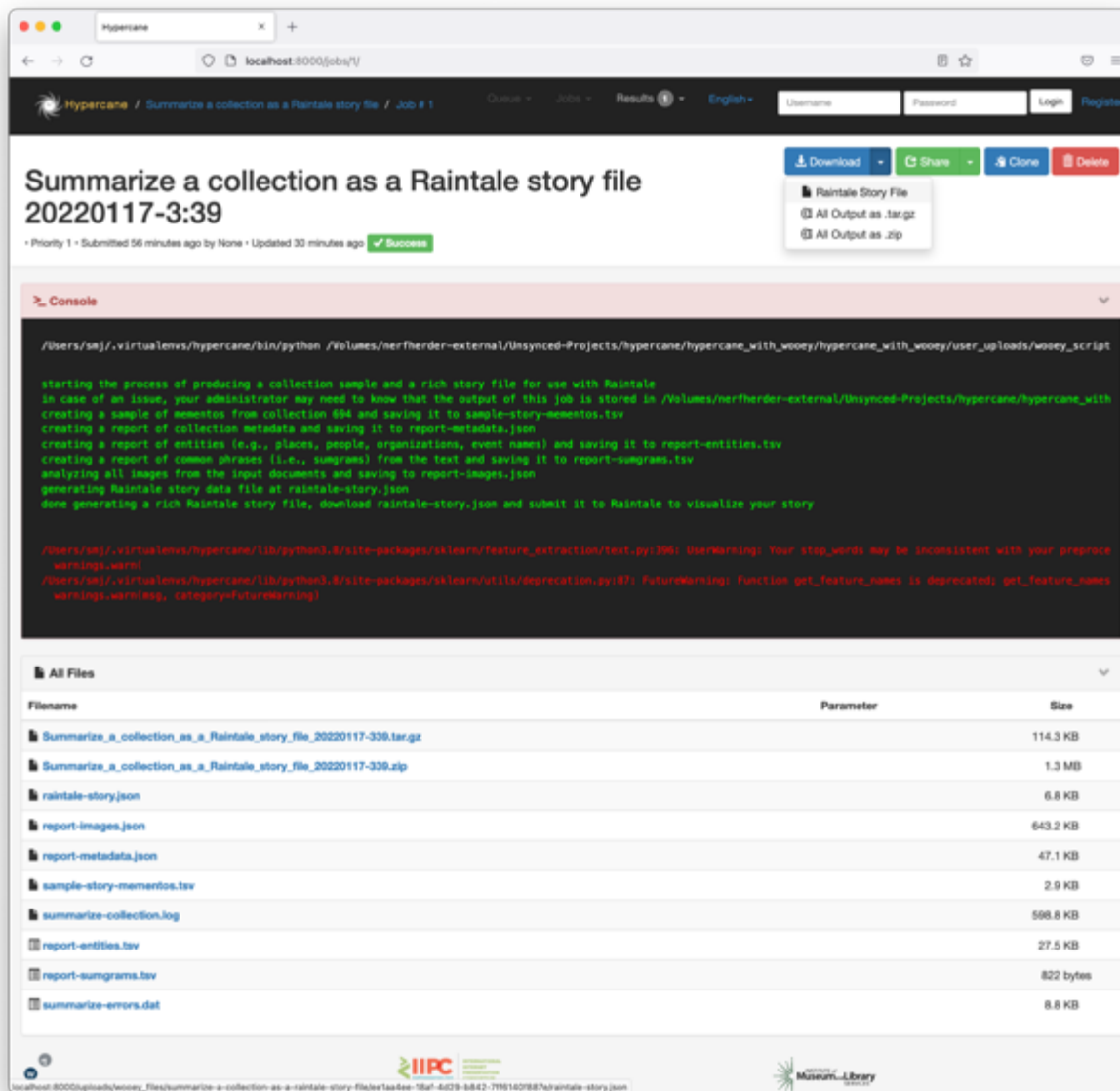


**Figure 15.** A screenshot of a completed Hypercane job. The user can click the Download button (top right) and save the Raintale story file. (The red text only indicates a warning.)

Most users just wish to leverage Hypercane's storytelling capabilities. The recipe "Summarize a collection as a Raintale story file" accepts a collection ID for Archive-It, PANDORA Subject, PANDORA Collection, or Trove and executes the necessary `sample`, `report`, and `synthesize` commands to produce a Raintale story file. A user will need to supply the collection type, collection ID, sampling algorithm, and story title, as shown in Figure 14. Figure 15 shows the interface once the job is complete, allowing the user to download their Raintale story file.

Some users already have a list of memento URLs or wish to use mementos from a Memento-compliant archive not listed above. The recipe "Summarize

archived pages as a Raintale story file" accepts an input file containing a newline-separated list of memento URLs and does the same as the recipe above.

Other users have already selected exemplars from web archives and do not need to execute the `sample` step. The recipe "Create a Raintale story file from archived page URLs" accepts an input file containing a newline-separated list of memento URLs and executes the necessary `report` and `synthesize` commands to produce a Raintale story file.

Sometimes archivists or researchers just need a list of the memento page URLs in a collection. The "Generate archived page URLs from a collection" recipe accepts a collection ID for Archive-It, PANDORA Subject, PANDORA Collection, or Trove and executes the necessary `identify` command to find all memento page URLs in the collection. Similarly, a user may want the list of original resource URLs in a collection. The "Generate original page URLs from a collection" recipe performs the same steps but outputs a list of original resource URLs. The user could also run the "identify by Collection ID" web form, but these recipes present simplified language and options.

Additionally, there are cases where a user will desire WARCs for a third-party application. The recipe "Synthesize WARCs from a collection" accepts a collection ID and will execute the necessary `identify` and `synthesize` steps to convert the mementos in the collection into WARCs that closely resemble what the archive initially crawled. Hypercane leverages the Memento protocol and its knowledge of accessing raw memento content from different web archives to make this happen. These WARCs are an approximation built from what Hypercane can access via the web and are not identical to the web archive's holdings. Still, they are sufficient for experimentation and analysis with tools like AUT. Likewise, a user can produce WARCs from a list of memento URLs with the recipe "Synthesize WARCs from archive page URLs."

These recipes continue to make Hypercane more approachable for end-users. Though Hypercane primarily exists to feed content to Raintale, it has evolved into a much more promising toolkit of its own.

### MementoEmbed

MementoEmbed is a surrogate generation tool that can create four different types of surrogates. It produces social cards like those seen on Facebook and Twitter, page screenshots (also called browser thumbnails), word clouds, and **imagereels**. Imagereels are animated GIFs (GIF version 89a) of the top images found in a memento. Figure 16 displays one of MementoEmbed's social cards.

MementoEmbed also supports a web API for generating document metadata. Through this web API, a client can create surrogates. Clients can also request specific document metadata, such as page titles, memento-datetimes, image ranking, and automatic text summaries.

We developed MementoEmbed after analyzing more than 50 platforms and finding that none support mementos correctly. Some platforms do not understand how to extract the original resource URL from the memento. Others could not differentiate between a web archive's content and the memento content. Still, others refused to generate surrogates for mementos.
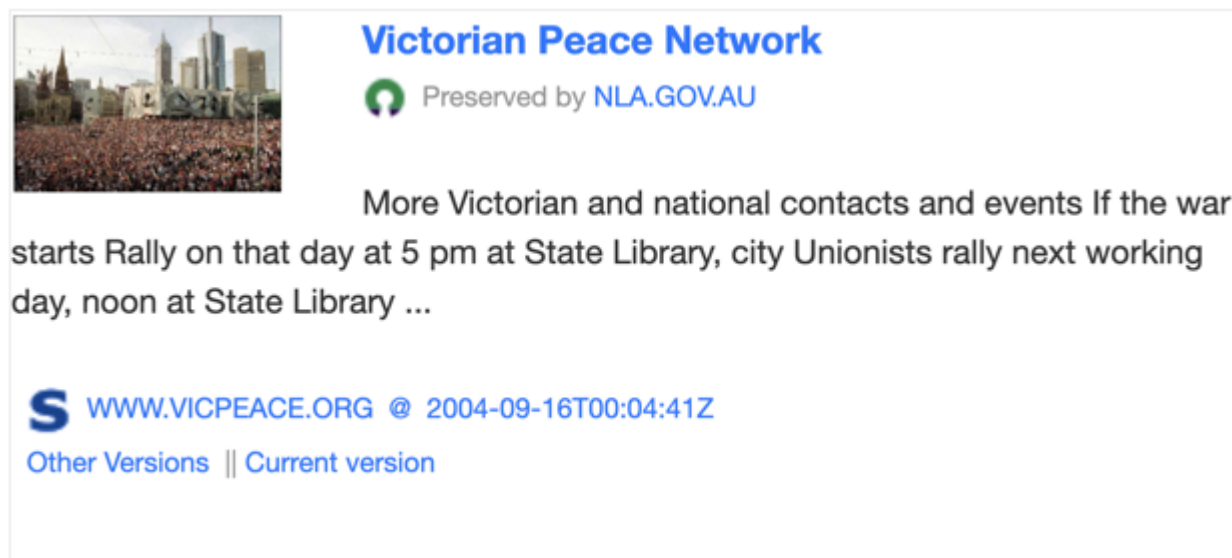
MementoEmbed is critical to the DSA Toolkit. Hypercane applies MementoEmbed's functions that discover images (preprint version) and raw memento content. Raintale is a client of MementoEmbed's web API.

Most web archives augment their mementos with navigation elements and branding to aid visitors. Unfortunately, these augmentations confuse natural language processing technology, like that implemented by Hypercane and MementoEmbed. For example, when comparing frequent terms among a set of mementos, the word "Trove" appears to be the most frequent word, even though it does not represent the actual memento content. For Wayback and OpenWayback web archives, like Archive-It, MementoEmbed simply inserts an `id_` flag into their memento URL to visit a page lacking these augmentations. Trove mementos required different handling.

Figure 17a displays a screenshot of a Trove memento. Figure 17b blurs its augmentations to emphasize the actual memento's content. As part of our pilot with the NLA, we had to update MementoEmbed to discard these augmentations when processing Trove's mementos. Finding the actual memento content required the following steps:

1. change the domain name in the memento URL from webarchive.nla.gov.au to web.archive.org.au

2. insert the `id_` flag into the memento URL

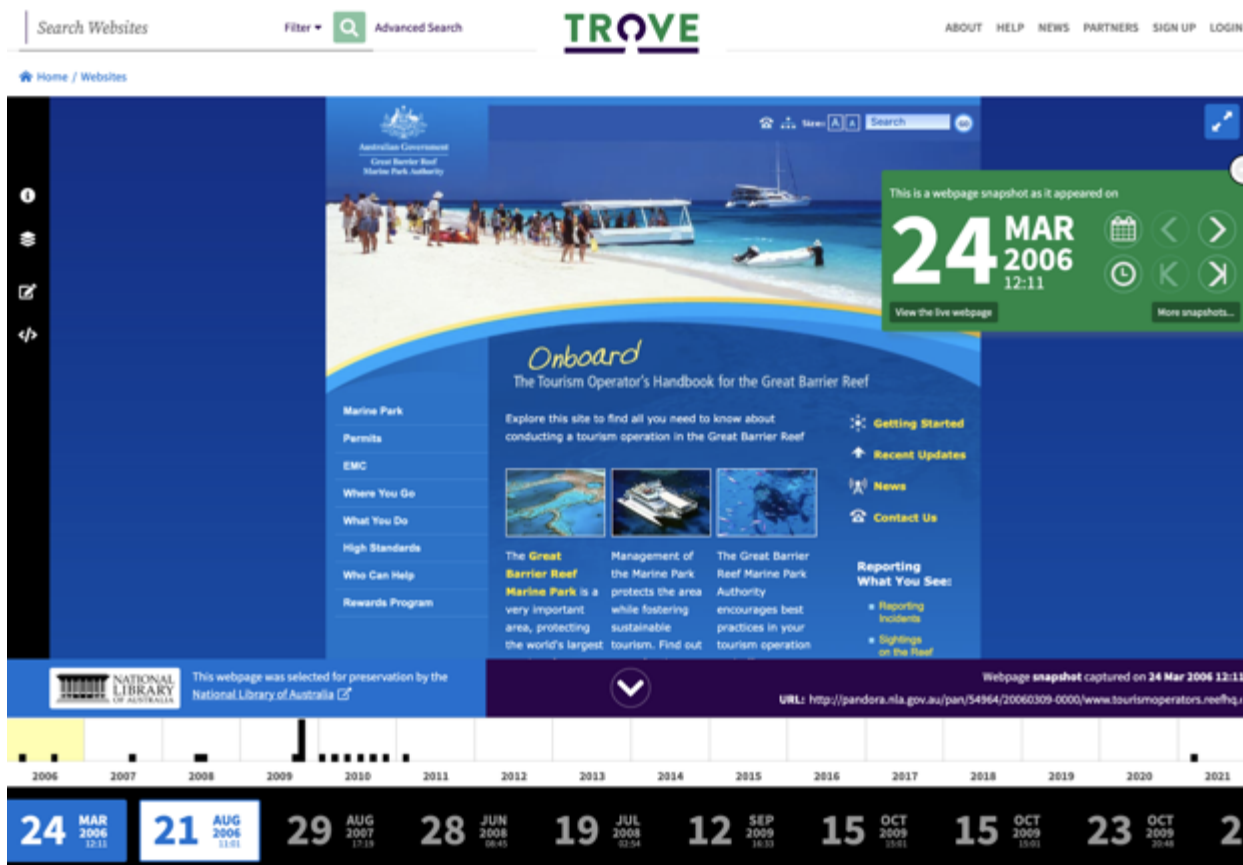From here, MementoEmbed can now process the raw memento content.

**Figure 16.** A screenshot of a MementoEmbed social card containing a striking image, title, description, original resource domain, and favicons, as well as information about the archive, when the archive preserved the page, and links to other mementos.

a. memento ♻ from Trove



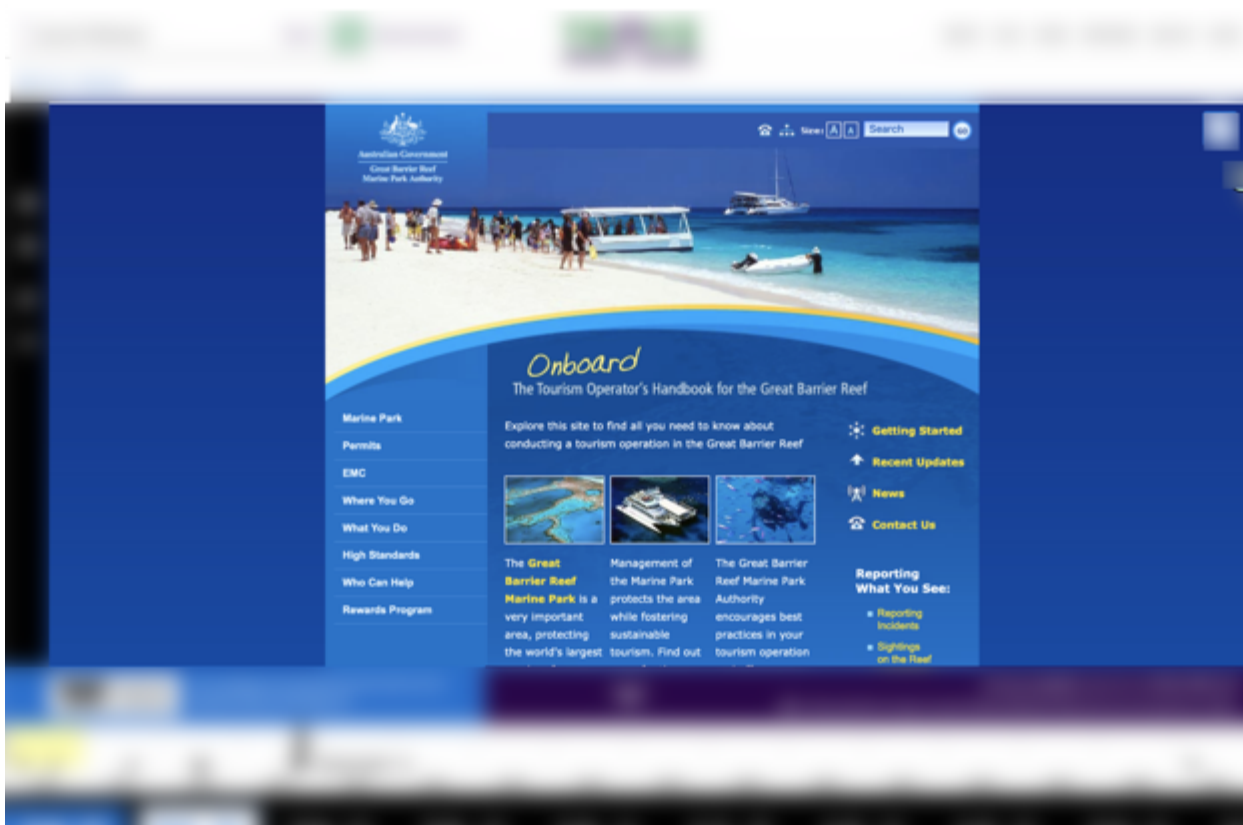b. The content of the memento exists in the unblurred area.

**Figure 17.** Trove mementos consist of augmentations that aid users in navigation and provide additional information that confuse natural language processing systems.

We also updated MementoEmbed to better handle creating screenshot surrogates of Trove's mementos. The additional augmentations required changing the timeouts on MementoEmbed's screenshot capabilities.

As with Hypercane, we created native Linux installers for MementoEmbed for CentOS 8 and Ubuntu 21.10 to help system administrators quickly stand up MementoEmbed. For more information on MementoEmbed, please see its documentation.

### Raintale

Hypercane focuses on selecting exemplars and generating story metadata. MementoEmbed focuses on generating document metadata for a single memento. Raintale takes in all of this information and produces a story in a format that the storyteller desires. Raintale is a command-line application executed by the `tellstory` command. To apply the rich story file generated in the Hypercane section above, a user types:

```
1   # tellstory -i mystory.json --storyteller template \
2       --story-template mytemplate.html -o mystory.html
```

The `-i` argument specifies the input file containing the information for the story. The `--storyteller` argument indicates the type of story being told. Because the above example suggests that we wish to use a template, the `--story-template` argument is needed to specify a file containing a template that Raintale will use to format the story. Finally, the `-o` argument indicates where to store the resulting story.

The user heavily controls the format of Raintale stories through templates. Figure 18 shows a story that applies a Bootstrap carousel template to a set of mementos to summarize the Trove collection *Tourism*. Figure 18 shows a story that applies an NLA-authored template to a group of hand-selected mementos about Australian Zoos, Wildlife Sanctuaries & Aquariums.
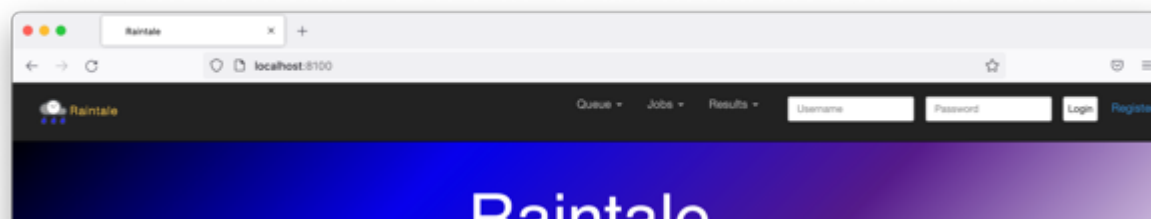
**Figure 18.** A Raintale story generated using a Bootstrap carousel.



**Figure 19.** A Raintale story generated using a template developed by the National Library of Australia.

Because Raintale supports templates, users can construct stories that take many forms, from the COVID-19 story in Figure 6 to the stories shown in Figures 17 and 18. Raintale applies a modified version of the Jinja2 template language 🔗 that supports additional options for extracting specific document metadata for the output. This template language does not limit Raintale to HTML. It permits many textual formats, like Markdown, XML, JSON, and Jekyll.

As with Hypercane, we recognized that a command-line application is not always approachable for all software users. As part of our pilot, we leveraged Wooey to create a WUI for Raintale as well. Figures 20–22 show screenshots of a user executing "Create Story From Template," which produces a templated story just like we demonstrated with the `tellstory` command above. The "Create Story From Preset" option lets users forgo a template in favor of built-in story templates. With "Tell Story With Twitter," Raintale will create a Twitter thread from the mementos in the input. Finally, "Create Video Story" generates an MP4 file consisting of the top-ranked images and text from each memento in the story.
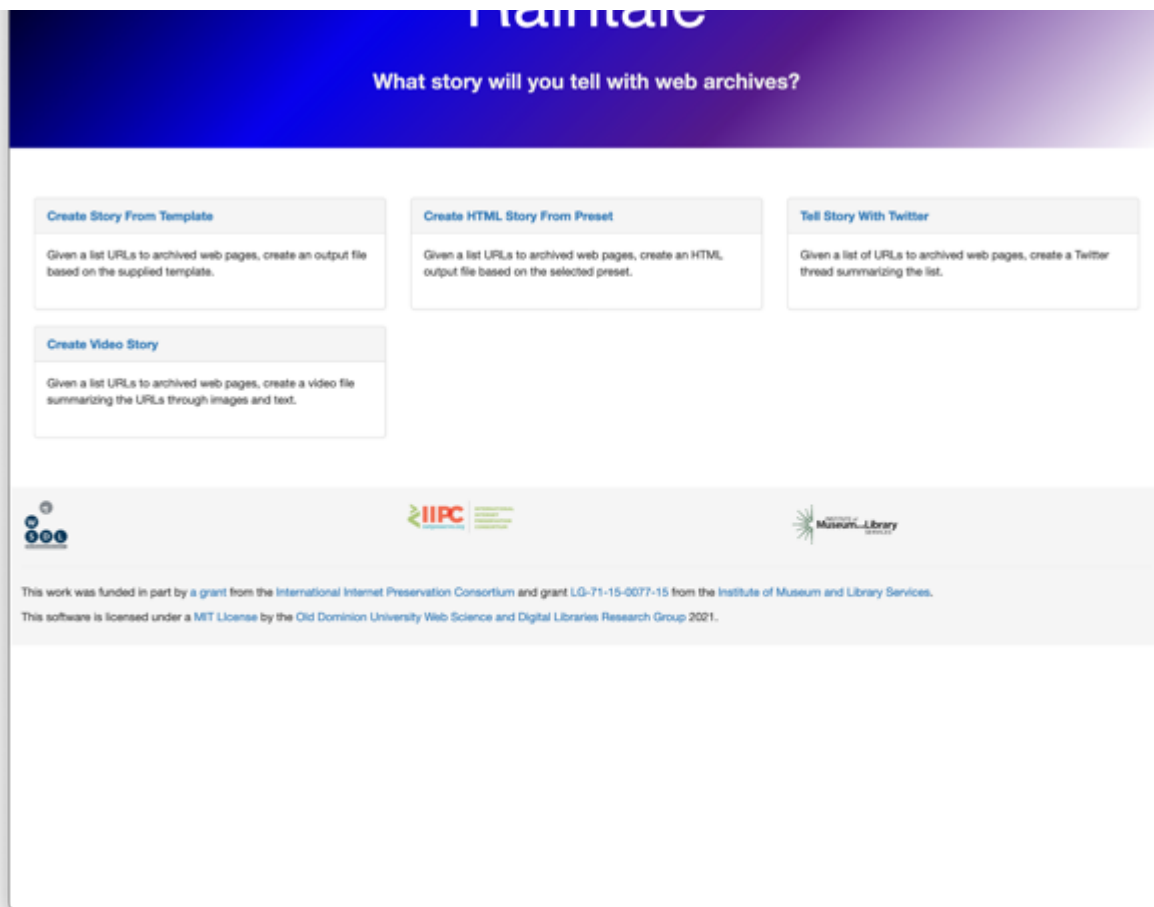
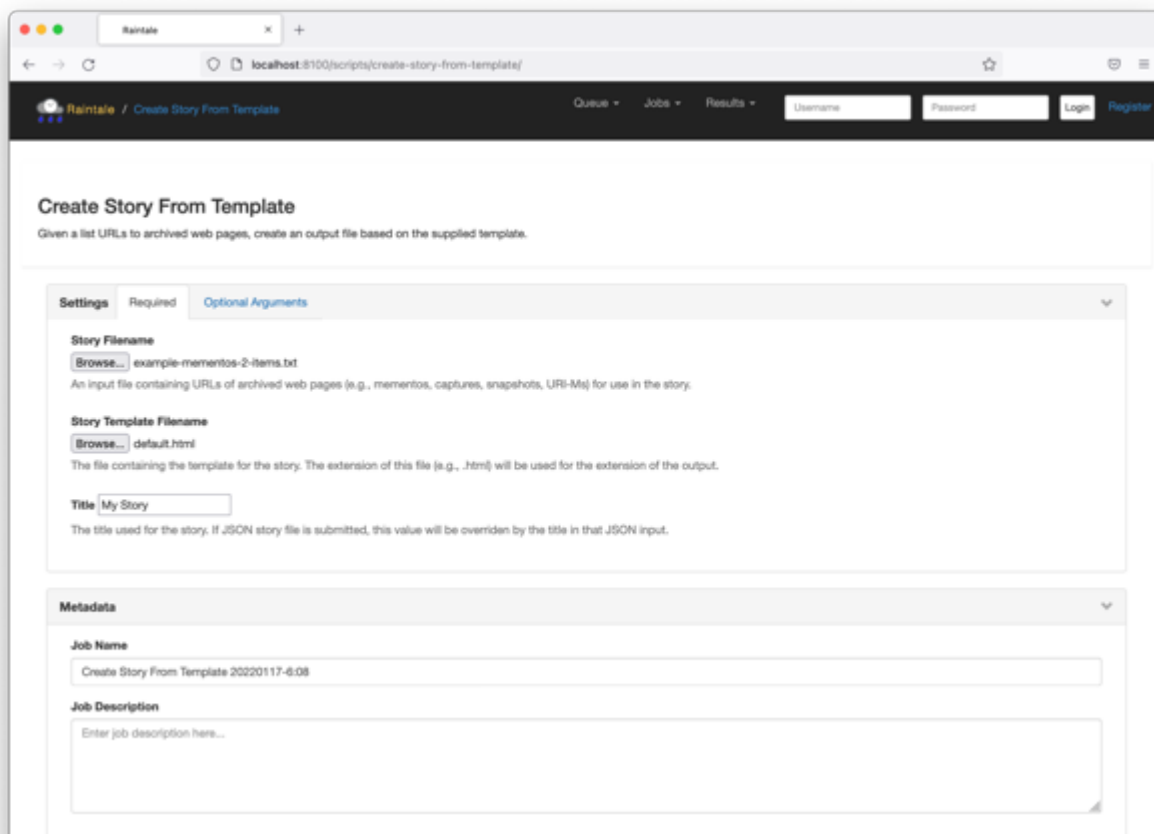**Figure 20.** A screenshot of the landing page of the Raintale WUI.

**Figure 21.** A screenshot of the web form allowing a user to "Create Story From Template."
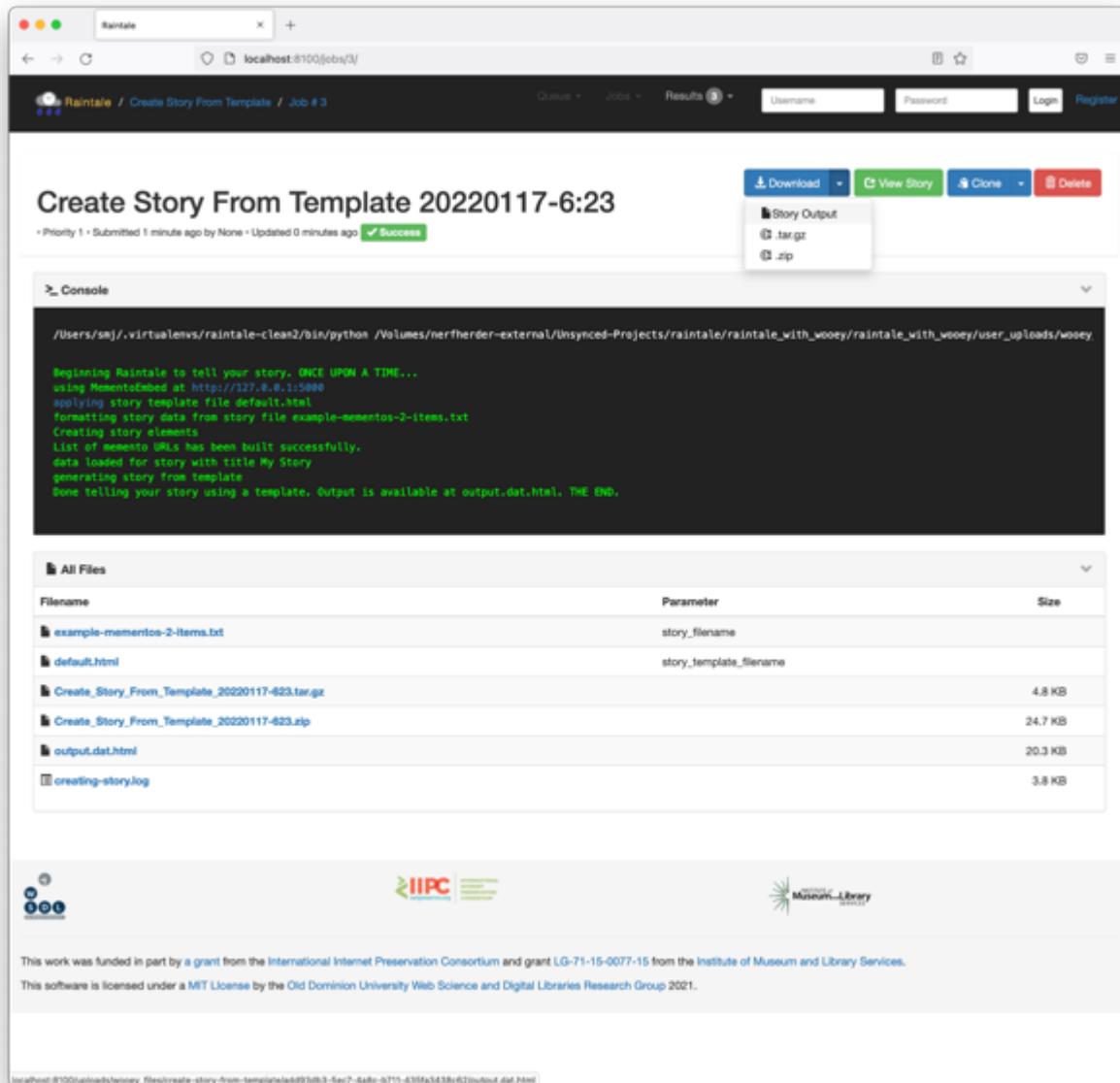


**Figure 22.** A screenshot showing a completed Raintale job where the user can download the story.

As with Hypercane and MementoEmbed, we developed native CentOS 8 and Ubuntu 21.10 installers for Raintale. We presented Raintale at the WADL 2020 workshop and IIPC WAC 2021. An early version of Raintale helped us conduct user testing that showed that social cards probably provide a better understanding of web archive collections when compared to alternatives like page screenshots. For more information on installing and leveraging Raintale, please consult its documentation.

**The Future of Storytelling With Web Archives**

We see a lot of potential for further development of these tools. This pilot project has helped us better understand the needs of our users. We continue to improve the error handling of these tools as we learn more from the community. Overall, we want to improve the approachability of these tools.

Our native installers and Web User Interfaces are a start but only apply to Linux. Thanks to the experience of creating an installer for Linux, we have some ideas for a macOS DMG file that supports the applications of Hypercane, MementoEmbed, and Raintale.

With Wooey, we could preserve our existing command-line interfaces for scripting while also providing new web user interfaces. The authors of Gooey promise this same functionality using the same methods but with a native non-browser graphic user interface. We are hopeful that our efforts with Wooey might make it easier to apply Gooey and provide a proper native GUI interface.

As we build upon these technologies, we are charting a path toward making the DSA Toolkit a set of Microsoft Windows desktop applications. As a set of Python applications, the DSA Toolkit will have to contend with library installation issues on Windows. Once we resolve those issues, we will have reached the maximum number of users with the current code.

We are also exploring additional functionality. AIU, and, by proxy, Hypercane now support Archive-It, Trove, and PANDORA. We can expand it to support web archive collections from the Croatian Web Archive, the UK Web Archive, the Library of Congress, and Conifer. This way, users can apply our tools to summarize more collection types. Of course, we would continue to support sets of memento URLs as input because that would extend the reach to web archives that do not currently support collections.

The Hypercane recipes inspired by our collaboration are just the beginning. We are exploring the idea of a "Recipe Builder" where users can create their own Hypercane scripts but in a much more user-friendly fashion than shell scripting. Once this is in place, we hope that users will share their recipes. We could even facilitate a recipe exchange for interested archives that want to share Hypercane solutions.

Raintale supports templates so that archivists can customize their stories. We see a future where archivists can share the look and feel of their stories with each other. In addition to templates, we are considering the idea of a "Story Builder" allowing users to create stories graphically in real-time, much like bloggers do with WordPress blocks.

Our collaboration between LANL, ODU, and the NLA has been fruitful but was only the beginning. Either by executing these tools on your own or visiting someone else's stories, what gems will you discover in a web archive collection?

## Acknowledgements

## About the Authors

Shawn M. Jones (0000-0002-4372-870X), Los Alamos National Laboratory, is an Information Science & Technology Institute Postdoctoral Fellow working for LANL's Information Sciences (CCS-3) division. Shawn recently received his Ph.D. from Old Dominion University and was advised by Dr. Michael L. Nelson. He is also an alumnus of ODU's Web Science and Digital Libraries research group. He has contributed tools, data, and analysis to projects and frameworks such as Memento, Signposting, and Robust Links. In addition, Shawn is the lead of the Dark and Stormy Archives project, an initiative to summarize web archive collections through visualization techniques common in social media. More information about Shawn is available at: https://www.shawnmjones.org/.

Himarsha R. Jayanetti (0000-0003-4748-9176) is a Ph.D. student at Old Dominion University working under the supervision of Dr. Michele C. Weigle. She is also a member of the Web Science and Digital Libraries research group as a graduate research assistant. Her research interests are in digital preservation, digital libraries, and social media. She graduated from Gujarat Technological University in India with a Bachelor's degree in Computer Engineering in 2017. More information about Himarsha is available at: https://himarshaj.github.io/.

Alex Osborne develops open source tools for web archiving and maintains the infrastructure of the Australian Web Archive at the National Library of Australia.

Paul Koerbin is Assistant-Director Web Archiving at the National Library of Australia. He has been involved with the development and operation of the NLA's web archiving program since its inception in the 1990s, including being part of the team that developed one of the first workflow systems for selective web archiving. He has published papers and given many presentations, in Australia and overseas, on web archiving operations and practice. He holds a graduate qualification in library and information studies from the University of Tasmania and a Ph.D. from the University of Western Sydney.

Martin Klein (0000-0003-0130-2097), Los Alamos National Laboratory, is a scientist and lead of the Prototyping Team in LANL's Research Library. In this role, he focuses on research and development efforts in the realm of web archiving, scholarly communication, digital system interoperability, and data management. He is involved in standards and frameworks such as Memento, ResourceSync, Signposting, and Robust Links. Martin holds a Diploma in Computer Science from the University of Applied Sciences Berlin, Germany, and a Ph.D. in Computer Science from Old Dominion University.

Michele C. Weigle (0000-0002-2787-7166) is a Professor of Computer Science at Old Dominion University. Her research interests include web science, social media, digital preservation, and information visualization. She has published over 115 articles in peer-reviewed conferences and journals and has served as PI or Co-PI on external research grants totaling almost $6M from a wide range of funders, including the National Science Foundation, the National Endowment for the Humanities, the Institute of Museum and Library Services, and the Andrew W. Mellon Foundation. She currently serves on the editorial boards of the *Journal of the Association for Information Science and Technology* (JASIST) and the *International Journal on Digital Libraries* (IJDL). Dr. Weigle received her Ph.D. in computer science from the University of North Carolina in 2003.

Michael L. Nelson (0000-0003-3749-8116) is a professor at Old Dominion University and the Virginia Modeling, Analysis, and Simulation Center (VMASC), and he co-leads the Web Science and Digital Libraries Research Group.  Prior to joining ODU, he worked at NASA Langley Research Center from 1991-2002, where he created the NASA Technical Report Server (NTRS).  More information about Dr. Nelson can be found at www.cs.odu.edu/~mln/ and twitter.com/phonedude_mln.

Subscribe to comments: For this article | For all articles