# FINAL REPORT
# FOUNDATION PROJECTS

## HERITAGE CONNECTOR
### Transforming text into data to extract meaning and make connections

**PI: John Stack, Science Museum Group**

Science Museum Group | Victoria and Albert Museum |
School of Advanced Study, University of London

**FEBRUARY 2022**

# TABLE OF CONTENTS

## AUTHORS

*Jane Winters*, *John Stack*, *Kalyan Dutia*, *Jamie Unwin*, *Rhiannon Lewis*, *Richard Palmer*, & *Angela Wolff*

# Executive Summary

The aims of the Heritage Connector (HC) project were to make a substantial contribution to enable realisation of the ambitions within the AHRC's Towards a National Collection (TaNC) programme to make collections accessible for research and public engagement purposes.

Bringing multiple cultural heritage collections together is fundamentally about building links. Online, these can be manifested as hypertext links which create a rich web of deep and broad user journeys between related content and information. These links also have the potential for computational analysis and visualisation enabling new forms of digital humanities research into collections.

The project explored three technologies that together have the potential to provide a step-change in access and discoverability, research and public engagement by augmenting traditional catalogue data and associated keyword search through generation of a vast number of interlinked resources and content. The three technologies Heritage Connector explored were:

- artificial intelligence (AI) – specifically, natural language processing (NLP), named entity recognition (NER) and entity linking (EL) – to build links at scale from thin collection records;
- linked open data (LOD) as a scalable and flexible structuring methodology;
- knowledge graphs to store links and make them accessible.

The project sought to demonstrate that generation of a rich web of links could be built and made available using these technologies on the following source datasets:

- Science Museum Group (SMG) Collection catalogue,
- Victoria and Albert Museum (V&A) Collection catalogue,
- Wikidata,
- Science Museum Group Journal,
- Science Museum blog.

The final web of links (structured in the knowledge graph) has 1,208,256 entities and 53 relations. The techniques used to generate the links was tuned in ways which were able to provide high quality links and even though the accuracy of these links in some cases falls short of those generated manually, a greater wealth of associated material is surfaced which has practical benefits.

The project findings are:

- Overall, the Heritage Connector project demonstrated that the methods used can be used to build links at scale between and within collections.
- The approach taken by the project provides avenues to better solve existing challenges of discovery and exploration for large collection datasets and enable new forms of data-based analysis.
- In common with other linked data projects in the GLAM sector, the project found that the use of linked content in the cultural heritage sector has the potential to make collections more visible, expose hidden aspects, enrich existing catalogues, allow data reuse in new contexts, and enable improved user experiences.
- The availability of museum collection catalogues and other data sources as well-documented APIs (application programming interfaces) significantly speeded up the project's technical work.

However, none of these APIs are designed for large-scale 'bulk' use and can be slow and/or unstable for this application. So, for this kind of project, data extracts are ideally required.

- As link building is undertaken on the collection catalogue datasets, the resulting dataset rapidly becomes extremely large. Given the size of the source datasets, there is significant benefit in being selective about the data and content processed and making sure it is focussed on specific outcomes.

- Aligning controlled and free-text dataset fields (such as those found in collection catalogues) to entities can take a significant time using existing tools, and using the NLP, NER and EL methods trailed by the project generated vast numbers of links: within collections, between collections, and to and from other content sources (Wikidata, journal articles and other texts such as blog posts).

- Working with the collections and NLP, NER and EL, a pipeline approach with various stages was demonstrated to work well.

- If the source collection data includes persistent identifiers (PIDs) and links to sources (e.g. Wikipedia or biographical sources), these can be used to build links with a high degree of confidence.

- False positives ('mistakes' by the machine learning) were a tiny minority of the links and are usually readily apparent, even to non-specialist audiences.

- For the best results, it is not a question of if human intervention and curation is needed, but when it should be used and how it may be most usefully focused. Providing descriptive texts for collection objects was proven to be valuable not only to human users but also for entity extraction by machine. Subject matter expertise is required to select the appropriate source datasets to ensure that they were manageable and to review machine learning outputs to improve the outputs in subsequent stages.

- Although it would be desirable to display a 'degrees of confidence' rating for each link in user interfaces, these were impossible to calculate in this project, as links were created by multiple machine learning models used in a pipeline approach.

- Using the project's methodologies generates an extremely large dataset of linked open data (LOD). Handling these output linked data sets of links in a knowledge graph was demonstrated to work well. Knowledge graphs make it relatively easy to visually map disparate data and apply mathematical functions across large volumes of data.

- Barriers to LOD in the cultural heritage sector fall under four broad headings: technical, conceptual, legal and financial. Working with LOD at any kind of scale is both time consuming and resource intensive. A great deal of LOD work to date has focused on records for people rather than objects. Many LOD projects involve only one or at most two institutions, and international collaboration is relatively rare.

- Wikidata is a rich and diverse data source that was shown to be valuable to enrich museum collection entities with new data and act as a bridging point between collections.

- Easy to use interfaces are important for non-technical users to visualise, explore and navigate output data (which is potentially colossal in scale). Rather than creating single monolithic user interfaces, there is significant potential to create multiple light-weight interfaces and tools of different kinds onto the same output dataset.

- When considering further projects of this type, it is important that consideration is given to framing and contextualisation for machine learning generated outputs. The approach taken challenges traditional cultural heritage notions of the 'canonical' collection catalogue.

# Abstract

As with almost all data, museum collection catalogues are largely unstructured, variable in consistency and overwhelmingly composed of thin records. This is largely a legacy of the development of these catalogues from handwritten paper records used primarily for managing collections rather than public access. The form of the catalogues means that the potential for new kinds of digital research, access and scholarly enquiry remain dormant. Searching across collections is currently possible only through aggregation, which is labour-intensive to implement, or by third-party search engines where results are variable and unreliable. In this project, we applied artificial intelligence techniques to connect similar, identical and related items within and across collections. Our primary research question was "How can existing digital tools and methods be used to build relationships at scale between poorly and inconsistently catalogued digitised collection objects and other content sources?"

The Heritage Connector created a linked data knowledge graph that enables new forms of research and exploration. Furthermore, it explored the opportunity for computer generated links with Wikidata to provide new levels of structure and machine-readable data that can form the foundation of new types of discovery and access. The Heritage Connector uses a range of technologies including named entity recognition; entity linking; open data; and knowledge graphs. These methods created a large-scale data source of links. Computational inquiry to generate links via an application programming interface (API) enabled the creation of a range of proof-of-concept demonstrator research and discovery tools.

# Aims and Objectives

The project aims were to:

- make a substantial contribution to enable realisation of the ambitions within the UKRI Research Infrastructure Roadmap to make collections accessible for research purposes;
- provide proof of concept for the application of entity-linking approaches to making connections between online representations of heritage objects, including catalogue records, images and other assets;
- share this new and developing understanding with stakeholders and all interested parties.

 The project's objectives were:

- to conduct a review of relevant literature and digital tools;
- to construct a dense web of links between object records and Wikidata for evaluation and experimentation;
- to apply a series of digital tools / computational methods to create speculative identifications between different records within the test dataset;
- to work on successively larger and varied datasets as the project proceeds in three stages, starting with SMG records, then adding those of the V&A and then adding textual sources including the Science Museum's blogs and the Science Museum Journal articles;
- at every stage to hold focussed workshops for groups relevant to each of the project's phases;
- to publish the resulting code as open source software and data as open data;
- to write a report to Arts and Humanities Research Council (AHRC) on our findings congruent with the demands of the TaNC Programme Director;
- to write articles for peer reviewed journals.

# Partnership structure

The project was a collaboration between the Science Museum Group (SMG); Victoria and Albert Museum (V&A); and School of Advanced Study, University of London (SAS). Wikimedia UK also joined the project's working group.

The SMG is made up of the Science Museum, London; National Science and Media Museum, Bradford; Science and Industry Museum, Manchester; National Railway Museum, York; and Locomotion, Shildon. SMG holds a collection of approximately 425,000 objects and seven million archival items. SMG was responsible for undertaking project leadership; project coordination, administration and project meetings; event hosting and organisation, project website, blog and YouTube channel hosting; and GitHub software code repository management. SMG led on designing the technical architecture, development of the project software, and documentation of these components. In parallel with the Heritage Connector project, SMG undertook a partnership with Wikimedia UK (Wikipedia's UK body) and hosted a Wikimedian in Residence and that partnership's activities were being aligned with Heritage Connector.

The V&A took part in project meetings (weekly and monthly) to input into project development and from early 2021 provided the second collection dataset (after SMG's) to be analysed, processed and ingested into the Heritage Connector knowledge graph.

SAS led on the literature review, aligning work with the wider digital humanities field and research interests, and management of the project's Zotero literature library of publications, presentations and case studies. SAS attended the monthly project meetings.

Although not funded through the project, Wikimedia UK joined the monthly project team meetings to provide guidance on their Wikidata project and identify opportunities to join up with other Wikidata projects in the GLAM (galleries, libraries, archives and museums) sector.

# Staffing structure

**John Stack** (Digital Director, Science Museum Group) was Principal Investigator and responsible for overall delivery of the project and management of the Project Coordinator.

**Jamie Unwin** (Technical Architect: Collections Online, Science Museum Group) was Co-Investigator and oversaw software development and approach and was manager of the Research Developer.

**Professor Jane Winters** (Professor of Digital Humanities and Pro-Dean for Libraries, School of Advanced Study, University of London) was Co-Investigator and led on literature review and will author one of the project papers.

**Kalyan Dutia** (Research Developer) led on software development and technical implementation.

**Rhiannon Lewis** (Project Coordinator and Doctoral research student, School of Advanced Study, University of London) undertook event logistics, blog authoring, project enquiries, meeting organisation and representation of the interests of digital humanities researchers.

**Angela Wolff** (Full Stack Developer, Digital Media, V&A) led on the addition of the V&A's collection catalogue data and inputting into the technical approach to enable this.

**Richard Palmer** (Senior Web Developer, Digital Media, V&A) provided technical input from the V&A and oversaw V&A collection data ingestion into Heritage Connector.

**Stuart Prior** (Project Coordinator, Wikimedia UK) provided guidance on the use of Wikidata and introductions to Wikidata specialists and related projects.

**Hope Miyoba** (Wikimedian in Residence, Science Museum Group) led a parallel project at SMG working with Wikipedia and Wikimedia Commons and provided the link between that project and Heritage Connector.

**Tim Boon** (Head of Research, Science Museum Group) advised on links to curatorial culture and practice.

# Revised overall programme

| Summary | 2020 | | | | | | | | | | | 2021 | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Feb | Mar | Apr | May | Jul | Jun | Aug | Sept | Oct | Nov | Dec | Jan | Feb | Mar | Apr | May | Jul | Jun | Aug | Sept | Oct | Nov | Dec |
| **Phase 1** | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | |
| Start up | ░ | ░ | | | | | | | | | | | | | | | | | | | | | |
| Review of literature | ░ | ░ | ░ | ░ | ░ | | | | | | | | | | | | | | | | | | |
| Review of technologies | ░ | ░ | ░ | ░ | ░ | | | | | | | | | | | | | | | | | | |
| Convening 1: Webinar | | | | | X | | | | | | | | | | | | | | | | | | |
| **Phase 2** | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| Software development | | | | | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ |
| Workshop at Linked Pasts 6 | | | | | | | | | | | X | | | | | | | | | | | | |
| Convening 2: Workshop | | | | | | | | | | | | | X | | | | | | | | | | |
| **Phase 3** | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| Ingestion of V&A data | | | | | | | | | | | | ░ | ░ | | | | | | | | | | |
| Stage 2: Ingestion of other data | | | | | | | | | | | | | | | ░ | | | | | | | | |
| Hackathon event | | | | | | | | | | | | | | | | | | | | | | X | |
| Convening 3: Webinar | | | | | | | | | | | | | | | | | | | | | | | X |
| **Phase 4** | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| User evaluation | | | | | | | | | | | | | | | | ░ | ░ | ░ | ░ | ░ | ░ | ░ | ░ |
| Write-up: AAIL journal | | | | | | | | | | | | | X | | | | | | | | | | |
| Milestone: Publication of software | | | | | | | | | | | | | | | | | | | | | X | | |
| Write-up: Report to AHRC | | | | | | | | | | | | | | | | | | | | | | | X |

# Events and consultations

| Event | Date | Attendees | Outputs |
|---|---|---|---|
| **Convening 1: Webinar: Wikidata and Cultural Heritage Collections** | 19 June 2020 | 296 webinar attendees<br>231 responses to event survey | Blog post<br>Event recordings<br>Survey of attendees |
| **Extra event added to project: workshop at *Linked Pasts 6* conference, British Library** | 7–16 December 2020 | 39 participants | Demo |
| **Convening 2: Project Workshop** | 22 March 2021 | 50 participants[1] | Blog post |
| **Hackathon in partnership with Cogapp** | 19 November 2021 | 20 participants[2] | Blog post |
| **Convening 3: Webinar: Heritage Connector: Findings, Demonstrators and Potential** | 3 December 2021 | 114 webinar attendees | Blog post<br>Event recordings |

# Talks and presentations

- RLUK Digital Shift Forum, 20 January 2021

- ai4lam Community Call, 16 February 2021

- Towards a National Collection Webinar: Heritage Connector (John Stack) & Deep Discoveries (Lora Angelova), 22 February 2021

- AEOLIAN Network: Employing Machine Learning and Artificial Intelligence in Cultural Institutions, 7 July 2021

- Towards a National Collection Webinar Opening UK Heritage to the World - Introduction Webinar: Connecting the UK's Cultural Heritage - Prof. Jane Winters, SAS, University of London, 18 August 2021

- Dealing with complexity: Collections Trust conference 2021, 14 October 2021

---

[1] Curators, collections management professionals, Wikipedia professionals; academics from digital humanities, history and other disciplines; community-based historians and practitioners
[2] Software developers, designers, Wikipedians, digital humanities professionals, museum professionals

# Research approach

Heritage Connector used artificial intelligence techniques on existing collection catalogues (SMG's and V&A's) and other datasets (Wikidata and SMG's text-based content) to build a knowledge graph holding a large volume of linked open data.

Knowledge graphs were popularised by Google in their 2012 blog post *Introducing the Knowledge Graph: Things, not strings*, in which they described how data structured in a graph (rather than a table which is how museum collection catalogues are generally stored) can help users receive better responses to their search queries, retrieve context around a specific object, and make new discoveries via serendipitous connections. These outputs can be achieved using knowledge graphs as they can create new links (relations) between items (entities) without worrying about creating pre-defined and specific database tables to hold such information.

The structure of data in a knowledge graph takes the form: a subject (an entity) has a defined relationship (a relation) to an object (another entity). Entities and relations can take the form of URLs (Universal Resource Locators), enabling further computational analysis such as examining all links to and from an entity, all relations of a certain type, or a series of such links both within a knowledge graph but also beyond (where the entities or relations are expressed as URLs. This network of relationships ('Triples' in LOD parlance) connects together to form a 'network graph' that can not only be traversed, but to which mathematical functions can be applied to find clusters (clustering) and nearby nodes (graph embeddings).

For example, Robert Stephenson and Company (entity) made (relation) The Rocket (entity) can be expressed in RDF and stored in knowledge graph as:

- https://collection.sciencemuseumgroup.org.uk/objects/co8084947
- http://xmlns.com/foaf/0.1/maker
- https://collection.sciencemuseumgroup.org.uk/people/cp2736

The value in this approach is that – unlike current approaches to collection catalogue data – a knowledge graph can:

- hold more diverse kinds of relationships;

- string relationships together;

- infer relationships;

- change and grow as required;

- easily cluster related content in different ways.

Heritage Connector also sought to demonstrate how such approaches might sit alongside existing methods for collection digital access. Among the areas explored by the Heritage Connector were ways in which this approach can:

- improve browsing interfaces where those are currently limited by the available collection catalogue data;

- enhance search by keyword where that is currently limited and where presentation of results in an ordered list is problematic;

- extract information (datapoints) from the free-text description and interpretation/description fields;

- create thematic and topic entry points for users of the collection(s);

- enable cross-collection linking and discovery so users can rapidly explore larger and more diverse volumes of material;

- offer onward links to related material in other collections and related sources;

- increase usage of collections through surfacing deep content more effectively;

- facilitate serendipitous discovery of material by providing a wealth of surrounding material and contextual content;

- generate links into knowledge graphs and third-party datasets that surface new data and facilitate new forms of cross-disciplinary research;

- explore playful and experimental approaches to collection access which will broaden audiences.

In addition to content from the V&A and SMG, the Heritage Connector knowledge graph was linked to Wikidata. Wikidata is the free, open, linked, multilingual and structured database which underpins Wikipedia but which is also a project in its own right. Today, Wikidata contains over 96.4 million items (December 2021) structured as linked data and includes references to numerous external data points in cultural heritage collections and to other data sources. Because of its size and origins, Wikidata covers a vast range of subject domains and extends far beyond the areas traditionally covered by museum collection catalogues. Some potential opportunities of linking collection catalogues to Wikidata explored were:

- extending the Heritage Connector knowledge graph to include 'facts' not in the collection catalogue and presenting them to users;

- using Wikidata as a 'bridging service' to provide onward links to other museum collections and data sources;

- ingesting data from Wikidata – or from other sources via Wikidata – into an index to improve discovery of collection objects, people, companies, organisations, etc.;

- using Wikidata points in the knowledge graph to infer and present new entry points into the collections such as themes, events and topics;

- using Wikidata as a route to Wikipedia entries for articles associated with collections.

Alongside building this knowledge graph and extending it through the addition of different datasets, we also built a number of prototypes to evaluate the quality of the knowledge graph connections and its potential affordances.

# Research results

The project had four phases:

1. project set up and literature review;

2. initial software development with Science Museum Group (SMG) collection and Wikidata;

3. extended software development with Victoria and Albert Museum (V&A) collection and other data sources, and development of a suite of demonstrator interfaces;

4. evaluation, write-up and dissemination.

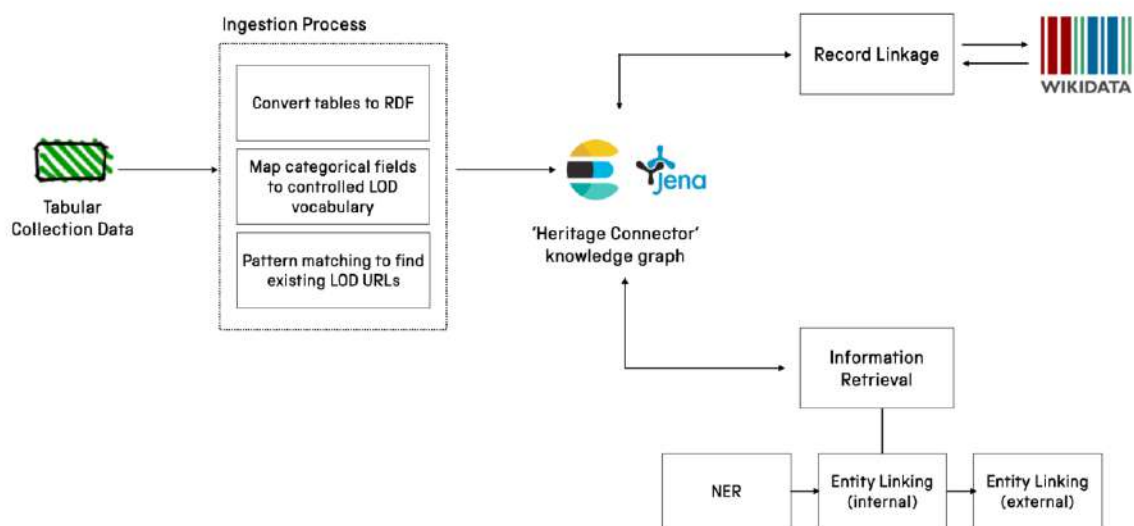## Phase 1

### Literature review

Findings from literature review:

- Motivations for GLAM (galleries, libraries archives and museums) institutions working with Linked Open Data (LOD) include: a concern to make cultural heritage more visible; an interest in exposing 'hidden' collections, or 'hidden' aspects of relatively well-known collections; the enrichment of catalogues and metadata; the encouragement of data reuse in new contexts; the desire to create a better user experience; the challenges of dealing with large volumes of data when resources are scarce.

- Many projects involve only one or at most two institutions, and international collaboration is relatively rare.

- Cultural heritage databases are rich, large and complex, there is limited standardisation, and institutional histories and cultures can make standardisation challenging.

- Barriers to LOD in the cultural heritage sector fall under four broad headings: technical, conceptual, legal and financial.

- Working with LOD at any kind of scale is both time consuming and resource intensive.

- A great deal of LOD work to date has focused on records for people rather than objects.

- It is not a question of *if* human intervention and curation is needed, but at what point in the pipeline it should be introduced and how it may be most usefully focused.

- Many LOD projects envisage personalisation as an important outcome, but this remains a mid- to long-term goal.

- Quality, authority and trust are crucial for cultural heritage organisations, but these can hold back experimentation and present a challenge for scalability.

- It is rare for promising experimental projects to move beyond the prototype stage.

# Phases 2 and 3

These phases of Heritage Connector involved building a system to create new links within and between the SMG and V&A collections and Wikidata, using the existing collections metadata and catalogue text, as well as the Science Museum Blog and Science Museum Group Journal (as exemplars of two forms of writing about collections: informal and scholarly). Existing metadata and these new links were stored in a Resource Description Framework (RDF) triplestore for the duration of the project, with a publicly accessible SPARQL (a query language for RDF) endpoint and user interface. Alongside developing the knowledge graph, various experimental interfaces were built to enable non-developers to view and interrogate the knowledge graph.

## Knowledge graph development

To develop our knowledge graph a range of machine learning and natural language processing (NLP) techniques were used, as shown in the schematic below.



First, an ingestion process was developed, which involved converting relational tables from SMG's Collections Online to RDF, mapping categorical fields such as places and object types (which were held in the SMG collections data as text) to Wikidata entities, and finding existing URLs (universal resource locators) that had previously been manually added to the dataset and which we could resolve to Wikidata entities through the Linked Open Data (LOD) cloud. The primary database and format for our data was Elasticsearch and JSON-LD, where each data source was held in a separate index. Later on in the project script was written to convert this JSON-LD to RDF, which was hosted in an Apache Jena Fuseki triplestore hosted on Amazon Web Services (AWS).

Noticing that there were no flexible open-source approaches for record linkage to Wikidata, a bespoke approach for this task for low training data and low metadata uses was developed. In summary, this approach searched for a record on Wikidata using its title, compared fields in the collection record to each candidate Wikidata record, and used a binary classifier to predict whether the collection record should link to each of the candidate Wikidata records (detailed information is available in our technical paper). To enable suitably fast text search on Wikidata, a tool was developed to create an Elasticsearch index from a Wikidata data dump, which we open-sourced as elastic-wikidata. Using this classifier we made 551

connections from SMG objects, 5,343 connections from SMG people, and 1,692 connections from SMG objects to Wikidata.

The final component of the system to build knowledge graphs used information retrieval methods to create links from text (e.g. object descriptions within a collection catalogue entry to Wikidata). We split this into a pipeline of three sequential methods:

- named entity recognition (NER) to identify mentions of people, places, events etc in text;

- entity linking (EL) to predict connections between any of these entities and SMG collection records;

- EL to predict connections from any remaining entities to Wikidata records.

For NER some extensions to the popular spaCy natural language processing (NLP) library were developed, for EL to the collection custom code based on existing research was written, and for EL to Wikidata a large pretrained model from Facebook Research was used. More information on the methods used can be found in the relevant technical paper and blog post.

After the work to create a knowledge graph from the SMG collection was complete, the same pipeline was used to process some other data sources, ultimately creating one knowledge graph containing all these data sources. Using the information retrieval components, Science Museum blog and Science Museum Group Journal were then processed into the knowledge graph. The stories and content about the collection in the sources made connections which would have not existed if we had processed only the collection catalogue. A subset of the V&A's digital collections was then identified which it was expected would have an overlap in themes to SMG's collection, and this was processed using the ingestion process and information retrieval components.

The final knowledge graph from the project contains four data sources (V&A collection, SMG collection, blogs and articles, and Wikidata) and has 1,208,256 entities (subjects and objects) and 53 relations (types of relationship between subjects and objects).

### Knowledge graph embeddings

To take advantage of knowledge graphs' capability for mathematical analysis across huge volumes of entities and relations, knowledge graph embedding models were explored to attempt to cluster and/or classify related content, visualise the entire contents of the knowledge graph, and use the knowledge graph to develop 'related items' recommendations. These models produced a vector for each entity and relation in the knowledge graph, while seeking to preserve their semantic meaning (i.e. their neighbours and connections) in the space in which these vectors lie.

After some options analysis, the RotatE model using the DGL-KE library was used. This was then implemented once for a knowledge graph formed from only SMG data and once for a knowledge graph formed from SMG and V&A data. Our demonstrators (next section) describe some successful applications of these embeddings.

### Demonstrators

Alongside the building of the knowledge graph several demonstrators were built. These served, and still serve, two purposes: to enable the Heritage Connector team and wider stakeholders to interrogate the knowledge graph and methods used to build it; and to investigate the affordances created by combining this dataset with new forms of interface. A webpage presents all these demonstrators.

Whilst the models were being developed, Streamlit was used to build a simple interface for audiences to inspect the NER and EL models. This was particularly valuable from a development perspective to allow team members with curatorial, digital humanities and museum collection data expertise to feed into the software development process.

Since an RDF store was being built as an output of this project, it was felt essential to implement a SPARQL query user interface (UI) onto this. An existing open source UI was forked, and a force-directed graph visualisation mode and our own sample queries were added to it.

Three different interfaces were then built, iterating and improving on them as a team:

- A 'connections' browser bookmarklet, which opens a sidebar from any SMG or V&A collections page, SMG Journal article and Science Museum blog post and shows the connections and similar items to the item on that page produced by Heritage Connector. This was based on an idea and some code from the Towards a National Collection's Heritage PIDs project.

- A visualisation, showing a map of the knowledge graph, in which each record is presented as a dot in 2D space. Two versions of this were built: one with the SMG collection, Science Museum blog, and Science Museum Group Journal, and one which also included the parts of the V&A's digital collection that were processed.

- A 'metadata explorer', which enabled guided exploration of the connections and similar items in the knowledge graph by clicking through links between webpages representing the knowledge graph.

## Hackathon

To further explore the resulting knowledge graph and the potential for new forms of interface for public engagement, exploration and visualisation, an end-of-project hackathon was held in collaboration with the digital agency Cogapp, which developed an additional suite of proof-of-concept interfaces:

- Augmenting From The Outside

- 3D Space Curator

- Heritage Connector Link Race

- Good Neighbours

- RHiZOME

- Timeline Interface

- Map Interface

# Projects outputs

The software developed by the project can be found in the list of [GitHub code repositories](#) in the Annexes to this report. The project's output datasets, including the knowledge graph and embeddings, are available at [https://doi.org/10.5281/zenodo.5752010](https://doi.org/10.5281/zenodo.5752010). Demonstrator interfaces using the project's knowledge graph are listed in the Annexes. These include those developed by the project team as well as those developed through the project hackathon event. Recordings of the project's events are posted on the [project YouTube channel](#) and on the [Towards a National Collection YouTube channel](#). The project's [blog](#) is available online. The project's [Zotero library](#) is publicly available. The project's [interim report](#) is available on the Towards a National Collection website.

The published paper Dutia, K., Stack, J. Heritage connector: A machine learning framework for building linked open data from museum collections. *Applied AI Letters*. 2021; 2( 2):e23 is available at [https://doi.org/10.1002/ail2.23](https://doi.org/10.1002/ail2.23) with the abstract:

> As with almost all data, museum collection catalogues are largely unstructured, variable in consistency and overwhelmingly composed of thin records. The form of these catalogues means that the potential for new forms of research, access and scholarly enquiry that range across multiple collections and related datasets remains dormant. In the project *Heritage Connector: Transforming text into data to extract meaning and make connections*, we are applying a battery of digital techniques to connect similar, identical and related objects within and across collections and other publications. In this article, we describe a framework to create a Linked Open Data knowledge graph from digital museum catalogues, perform record linkage to Wikidata, and add new entities to this graph from textual catalogue record descriptions (information retrieval). We focus on the use of machine learning to create these links at scale with a small amount of labelled data, and models which are small enough to run inference on datasets the size of museum collections on a mid-range laptop or a small cloud virtual machine. Our method for record linkage against Wikidata achieves 85%+ precision with the Science Museum Group (SMG) collection, and our method for information retrieval is shown to improve NER performance compared with pretrained models on the SMG collection with no labelled training data. We publish open-source software providing tools to perform these tasks.

# Recommendations for the programme

Working with two cultural heritage collection catalogues (SMG and V&A), the Heritage Connector project established that – as expected – these were rich in content; substantial in size; complex in structure; variable in content depth; are limited in standardisation; have few internal links and structures; and exist mostly in silos with few external references, persistent identifiers (PIDs) and links. However, the project demonstrated that as source datasets for machine-learning-based link-building, the catalogues are rich sources and could potentially be even richer, as supplementing object records with even short descriptive texts will often add keywords that can be transformed into entities with links.

The availability of museum collection catalogues and other data sources as well-documented application programming interfaces (APIs) significantly speeded up the project's technical work. Indeed, the project's source datasets (V&A collection, SMG collection, Wikidata, etc.) were in part selected because access to data was not an issue with them, though it would likely have been in most other instances. However, none of these APIs are designed for large-scale 'bulk' use and can be slow and/or unstable for this application. So, for this kind of project, data extracts are ideally required. Once these 'bulk' datasets are available, work is required to make them work as desired: to pre-process them, to load them into software in a sensible way,

to mould them to answer different questions, and to display answers to those questions in a variety of meaningful ways.

As link-building is undertaken on the collection catalogue datasets, the resulting dataset rapidly becomes extremely large. Given the size of the source datasets, there is significant benefit in being selective about the data and content processed and making sure it is focussed on specific outcomes. Since the outcomes desired are likely to be diverse, there is mileage in providing access to cultural heritage datasets not only for bulk download but also enabling this to be segmented in comprehensive and coherent ways to allow subsequent processing for given research and engagement needs. Since such datasets are not static – arguably, they are forever a work in progress – consideration needs to be given not only to 'snapshot' bulk and segmented access, but also enabling ongoing access to the updated datasets over time in a logical way.

The objective of the project – to generate a wealth of new links through a range of machine learning methods – was shown to be viable and that a huge number of high-quality links could be produced with the methods used.

Aligning controlled and free-text dataset fields (such as those found in collection catalogues) to entities can take a significant time using existing tools, and using the natural language processing, named entity recognition and entity-linking (NLP, NER and EL) methods trialled by the project generated vast numbers of links: within collections, between collections, and to and from other content sources (Wikidata, journal articles and other texts such as blog posts). Current NLP, NER and EL methods are strongest at identifying people, organisations, companies, events, dates, nationalities, facilities (e.g. buildings, bridges, etc.), laws, and some object types (e.g. ships, mass-produced products, etc.). Working with the collections and NLP, NER and EL, a pipeline approach with various stages was demonstrated to work well. If the source collection data includes persistent identifiers (PIDs) and links to sources (e.g. Wikipedia or biographical sources), these can be used to build links with a high degree of confidence. Free text fields were shown to be valuable sites for NER-based link building, even thinly populated fields if they contain identifiable entities.

False positives ('mistakes' by the machine learning) were a tiny minority of the links and are usually readily apparent, even to non-specialist audiences. Although it would be desirable to display a 'degrees of confidence' rating for each link in user interfaces, these were impossible to calculate in this project, as links were created by multiple machine learning models used in a pipeline approach, some of which don't output a confidence score.

For the best results, it is not a question of *if* human intervention and curation is needed, but *when* it should be used and how it may be most usefully focused. Providing descriptive texts for collection objects was proved to be valuable not only to human users but also for entity extraction by machine. In projects of this type, subject matter expertise is required to select the appropriate source datasets to ensure that they are manageable, to provide technical teams with insights into the form and structure of the source datasets, and to review machine learning outputs to improve the outputs in subsequent stages.

Using the project's methodologies generates an extremely large dataset of links. Handling these output linked datasets of links in a knowledge graph (KG) was demonstrated to work well. KGs make it relatively easy to visually map disparate data and apply mathematical functions across large volumes of data. Knowledge graph embeddings models were shown to be useful when trying to run algorithms across large knowledge graphs (e.g. nearest neighbours/related item engine/visualisation). These embeddings can then be used in a variety of user interfaces including those that visualise clustering of related content.

In common with other linked data projects in the GLAM sector, the project found that the use of linked content in the cultural heritage sector has the potential to make collections more visible, expose hidden aspects, enrich existing catalogues, allow data reuse in new contexts, and enable improved user experiences.

Wikidata is a rich and diverse data source that was shown to be valuable to enrich museum collection entities with new data and act as a bridging point between collections. However, because of its origins, Wikidata is occasionally patchy in implementation. So, comprehensiveness of Wikidata should not be relied upon at this point in time. Depending on their type, varying degrees of success can be expected when disambiguating (attempting to link) collection records and identified entities with Wikidata. Wikidata is, for example, strong on people, organisations and events. Objects are naturally weaker, as the scope of what an 'object' is is somewhat loosely defined in Wikidata (compared to museum collections), and specific instances of objects are not always in Wikidata (nor should they be: e.g. 'Olympic torch carried by person X'). As with its sister project Wikipedia, Wikidata contains many of the expected biases: gender, race, under-presented minorities, 'Western' skewed, etc.

Easy to use interfaces are important for non-technical users to visualise, explore and navigate output data (which is potentially colossal in scale). SPARQL in particular, while powerful, is challenging for such users. The needs of users are likely to be disparate and need to be defined carefully as it can be costly and time consuming to build, host and maintain such interfaces. Rather than creating single monolithic user interfaces, there is significant potential to create multiple lightweight interfaces and tools of different kinds onto the same output dataset using a range of APIs which allow software to access the data. In addition, allowing large-scale data dumps which can then be processed by users with the skills to apply digital humanities and visualisation tools is valuable.

The project's demonstrator interfaces explored a variety of approaches which enabled exploration via an overlay on existing collection interfaces (via a browser bookmarklet), using a dedicated metadata explorer interface, a macro visualisation and a number of more playful interfaces (games, maps, Twitter integrations, etc.). These interfaces showed that it was possible to build interfaces on top of the knowledge graph, that the knowledge graph dataset was valuable to a range of user types, and that the quality of links generated was high.

Integration of content beyond collections – blog posts and journal articles – showed that linking to and from these was viable and that integrating these into the knowledge graph was valuable as they became part of the wider ecosystem of content links enabling user journeys to and from related content.

When considering further projects of this type, it is important that consideration is given to framing and contextualisation for machine learning generated outputs. The approach taken challenges traditional cultural heritage notions of the 'canonical' collection catalogue. Rich user interfaces of the kind that the Heritage Connector approach enables, risk giving the impression of comprehensiveness and authoritativeness, whereas the outputs are variable in breadth, depth and quality depending on the content sources and the methods used. The project established that it is therefore desirable to approach the use of machine learning critically as an adjunct to manual cataloguing and link building, so further exploration of how best to frame machine-generated content alongside human-generated content would be valuable. Furthermore, consideration should be given to using crowdsourcing to flag 'false positive' links within user interfaces which could then be suppressed for subsequent users.

The software developed through the project is not a standalone tool that can be easily transformed into an ongoing product for processing cultural heritage content. Rather, Heritage Connector demonstrates significant potential for NER, NLP, EL and KGs approaches, and is a set of tools, techniques, approaches and software libraries that can be used to achieve link-building at scale. Considering future technical developments using these methods with such content, the project identifies that there is potential to move towards approaches for specific subject domains, aspects of the collections or user needs that could be valuable to those exploring the content. For example, a tool that allowed users to visually explore a collection (or set of collections) through a knowledge graph; development of more specific language models that were trained for a more specific subject domain; or a standalone 'disambiguation engine' which could take a webpage record for a person, place or object and disambiguate that entity to a known source with an accuracy score using all the other information on that page.

The approach taken by the project provides avenues to better solve existing challenges of discovery and exploration for large collection datasets and enable new forms of data-based analysis. Through the generation of a vast number of structured links within, between and beyond collections, the project's approach is able to highlight connections and suggest related content to users which would be impractical to generate manually; it is able to visualise and highlight the range and diversity of collections for users; and it is able to provide links to and from external content (e.g. Wikidata, journal articles and blog posts). While a limited number of datasets were used in the project, the methods used could achieve similar results with a larger number of datasets, and the project therefore identified the potential alternatives to or improvements to large-scale aggregator approaches currently constructed by manual linking and/or collection data alignment through transformation.

# Contacts

**John Stack (PI)**
Digital Director
Science Museum Group
John.Stack@ScienceMuseum.ac.uk


**Jamie Unwin (Co-I)**
Technical Architect: Collections Online
Science Museum Group
Jamie.Unwin@sciencemuseum.ac.uk


**Professor Jane Winters (Co-I)**
Professor of Digital Humanities
School of Advanced Study, University of London
jane.winters@sas.ac.uk
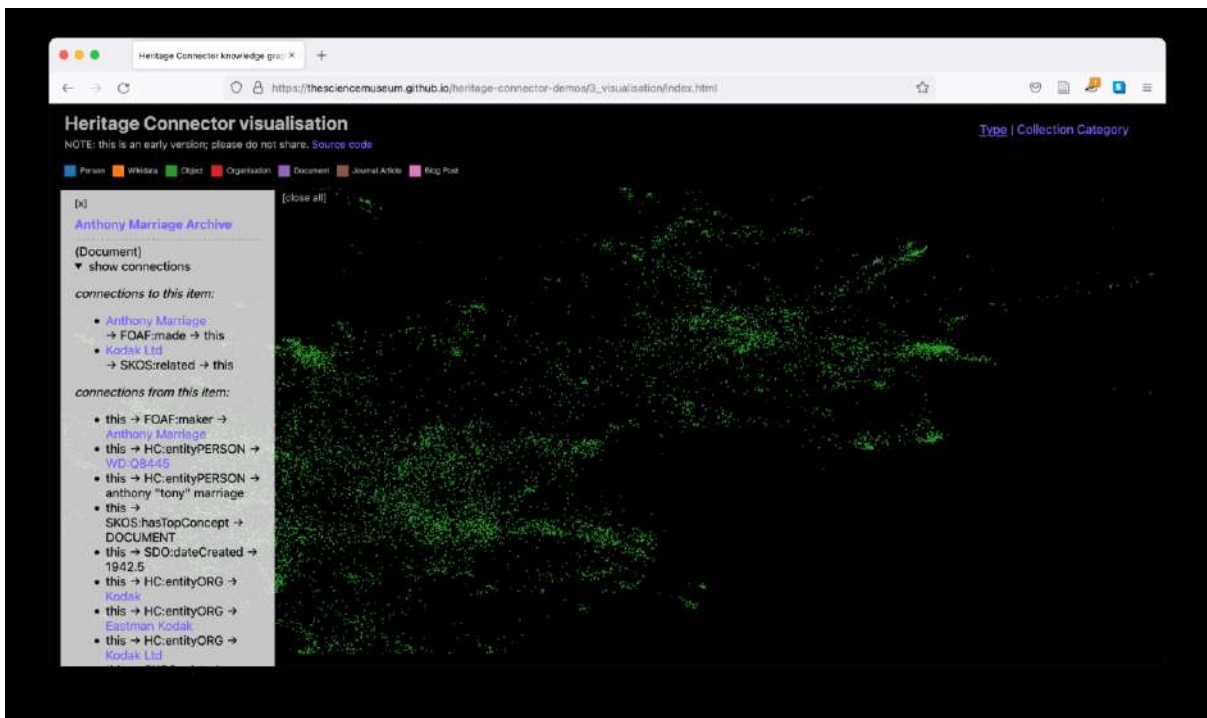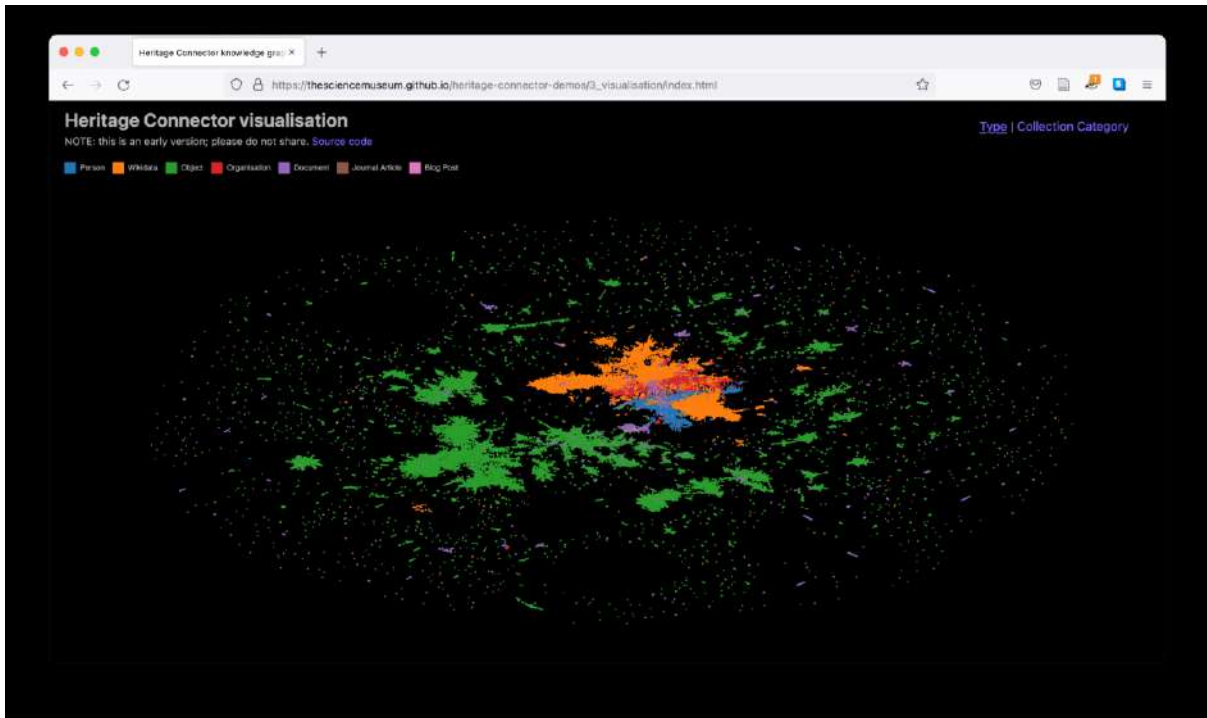
# Annexes and links

## Annexe 1: Glossary

- **AHRC** - Arts and Humanities Research Council

- **AI** - artificial intelligence

- **API** - application programming interface

- **AWS** - Amazon Web Services

- **EL** - entity linking

- **GLAM** - galleries, libraries, archives and museums

- **JSON-LD** - linked data format based on JSON (JavaScript Object Notation)

- **KG** - knowledge graph

- **LOD** - linked open data

- **NER** - named entity recognition

- **NLP** - natural language processing

- **PIDs** - persistent identifiers

- **RDF** - Resource Description Framework

- **SAS** - School of Advanced Study, University of London

- **SMG** - Science Museum Group

- **SPARQL** - RDF query language

- **TaNC** - Towards a National Collection programme

- **UI** - user interface

- **UKRI** - United Kingdom Research and Innovation

- **URL** - universal resource locator (web address)

- **V&A** - Victoria and Albert Museum

## Annexe 2: Links

- Project webpage: https://www.sciencemuseumgroup.org.uk/project/heritage-connector/

- Project blog: https://thesciencemuseum.github.io/heritageconnector/

- Demonstrators: https://thesciencemuseum.github.io/heritage-connector-demos/

- Project Zotero library: https://www.zotero.org/groups/2439363/heritage_connector/library

- Project YouTube channel: https://www.youtube.com/channel/UCzO6jroIvj-JbFuiQ9BpZdQ

- Project Github software repositories:

  - https://thesciencemuseum.github.io/heritageconnector/post/2022/01/24/Code-Repositories/

  - https://github.com/TheScienceMuseum/heritage-connector

  - https://github.com/TheScienceMuseum/heritage-connector-nlp

  - https://github.com/TheScienceMuseum/heritage-connector-apis

  - https://github.com/TheScienceMuseum/heritage-connector-deployment

  - https://github.com/TheScienceMuseum/heritage-connector-vectors

  - https://github.com/TheScienceMuseum/heritage-connector-data

  - https://github.com/TheScienceMuseum/heritage-connector-demos/

  - https://github.com/TheScienceMuseum/thor-cors-proxy

  - https://github.com/TheScienceMuseum/fuseki-docker

  - https://github.com/TheScienceMuseum/elastic-wikidata

  - https://github.com/LinkedPasts/LaNC-workshop/tree/main/heritageconnector

- Project data outputs https://doi.org/10.5281/zenodo.5752010

# Annexe 2: Illustrations





*Screengrabs of Visualisation of knowledge graph*

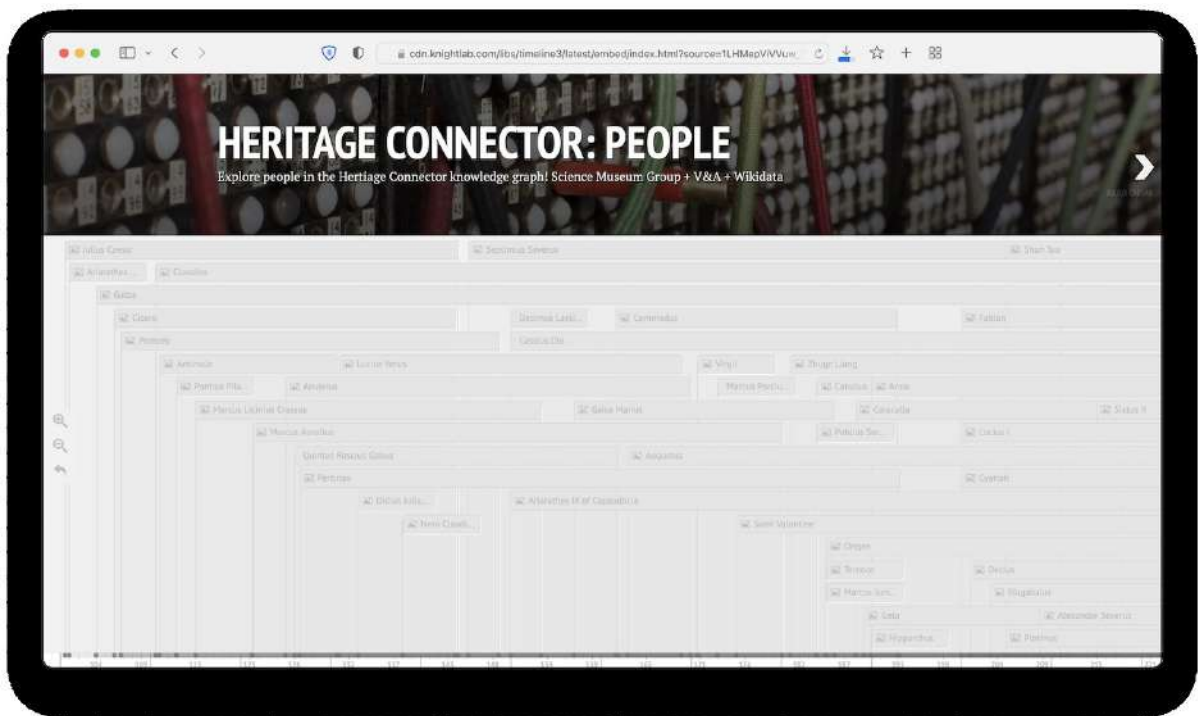*Screengrabs of Metadata Explorer interface demonstrator*

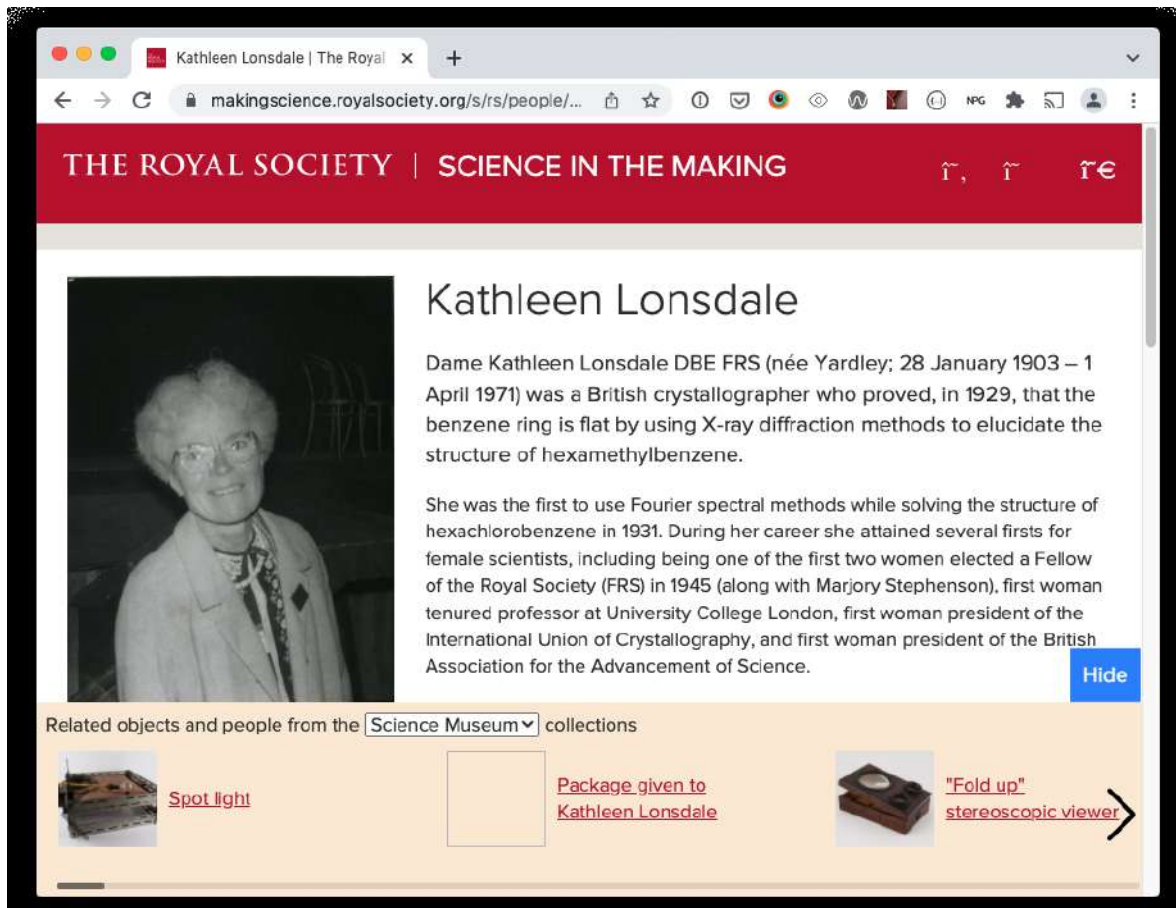*Screengrabs of browser bookmarklet on Science Museum Group Collection Online*

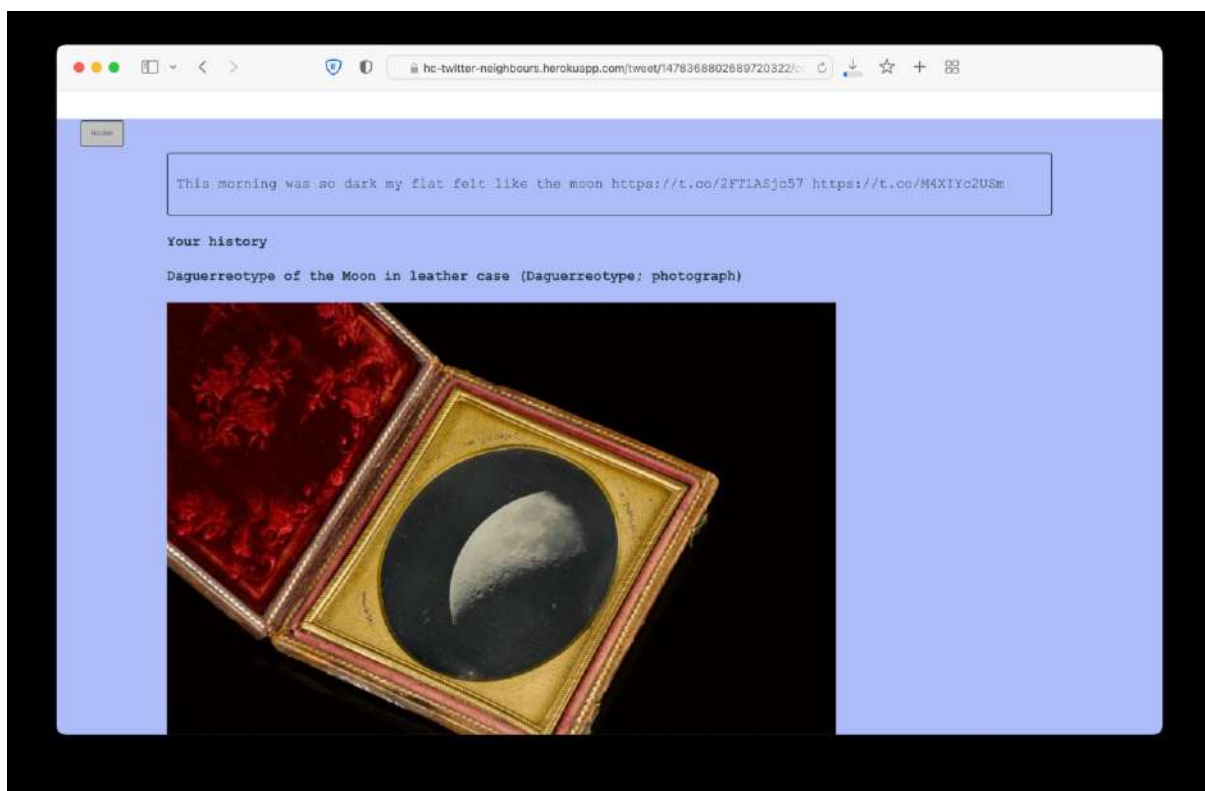*Screengrabs of browser bookmarklet on Science Museum Blog*

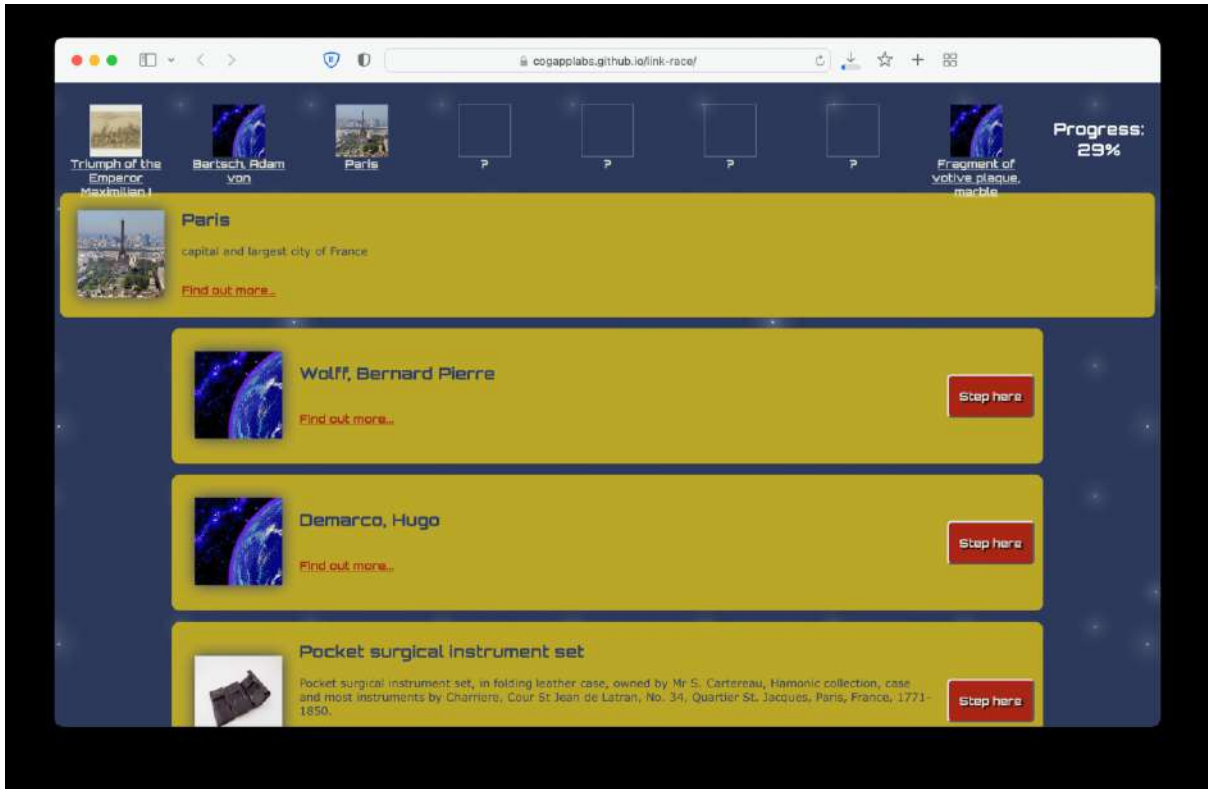*Screengrab of map interface demonstrator*
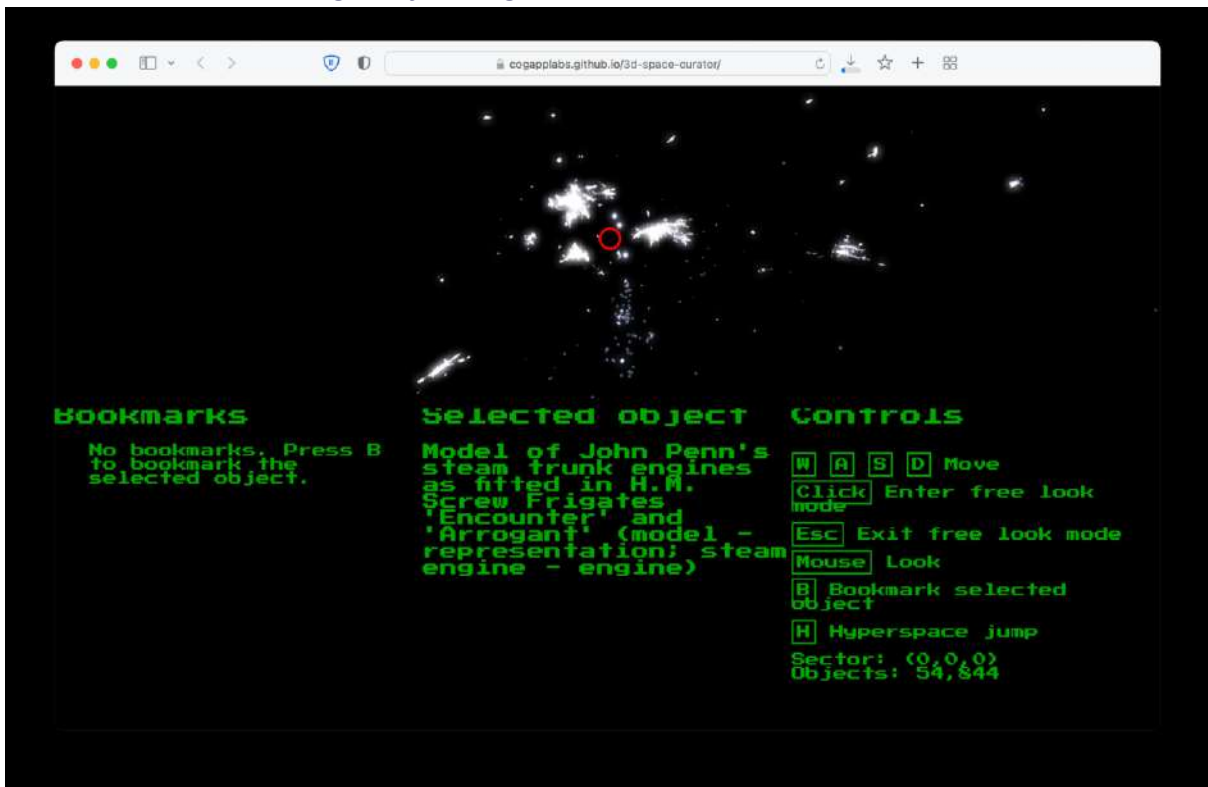


*Screengrab of timeline interface demonstrator*

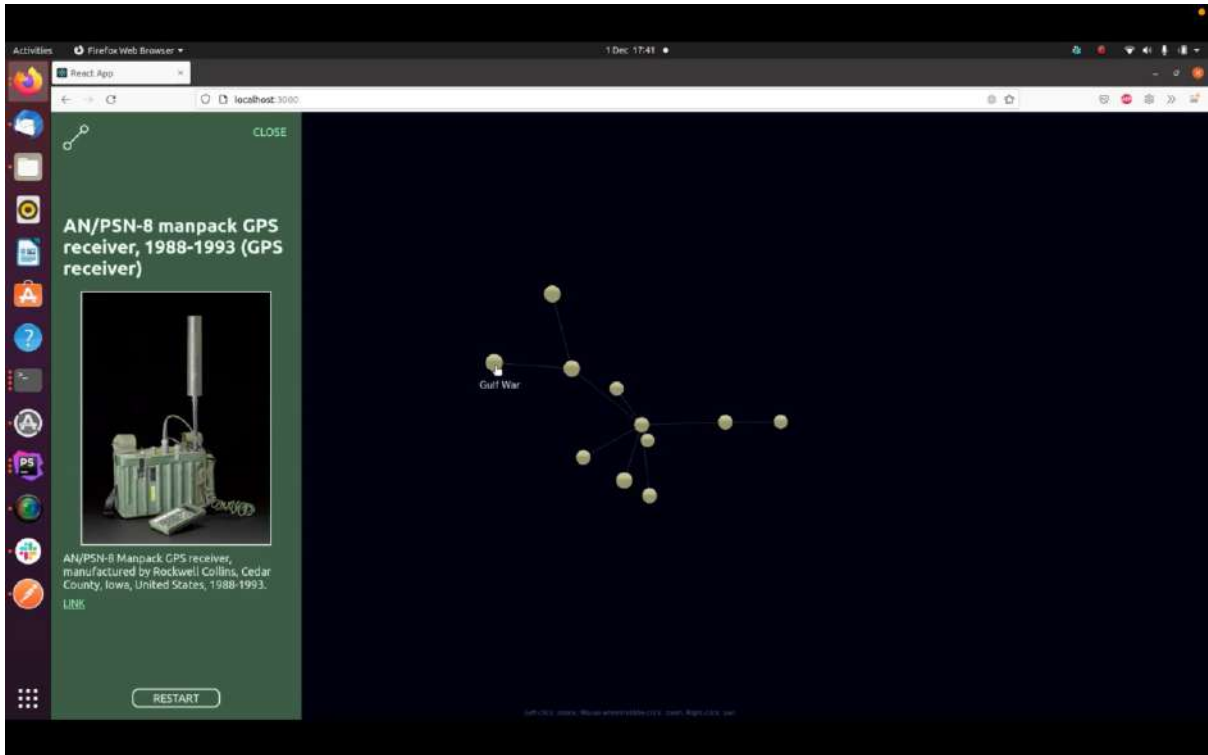*Screengrab of Augmentanator interface demonstrator*



*Screengrab of Good Neighbours demonstrator*

*Screengrab of Heritage Connector Link Race demonstrator*



*Screengrab of 3D Space Curator demonstrator*

*Screengrab of RHiZOME demonstrator*