# Institutional partnerships and open standards: Unlocking your archive for digital scholarship

**Kevin Cunningham**
Developer, Cogapp, Brighton, UK.
E-mail address:  kevinc@cogapp.com


**Andy Cummins**
Director of Product Development, Cogapp, Brighton, UK.
 E-mail address:  andyc@cogapp.com

### Abstract:

*This paper will look at the tools available for digital scholars on several large digitised document archives. Each of these archives demonstrate how raw images and cataloguing metadata can be leveraged to create user-facing features based on the open standards provided by the International Image Interoperability Framework (IIIF).*

*We will explore how rich, responsive interfaces deepen engagement with users on mobile and desktop devices. Specifically, we will be look at the Endangered Archives Programme (https://eap.bl.uk/) and the Qatar Digital Library (http://www.qdl.qa) and the Arabian Gulf Digital Archive (https://www.agda.ae/) which between them comprise over eight million digitised pages.*

*These projects all demonstrate international inter-library cooperation and the use of common archival standards such as OAI-PMH, EAD and METS. We will demonstrate how these cataloguing standards allow us to present the metadata online in formats that are helpful to researchers and scholars. Not only this, but also how to use it to provide powerful, faceted search, and onward journeys for users to find relevant and related documents. The use of the IIIF Image and Presentation APIs provides advanced features to digital scholars including side-by side document comparison and image manipulation, all while relying on reusable, open-source software to keep custom coding to a minimum and allow experience in one online archive to be easily applied to another.*

*Finally, we will look at how OCR transcripts combined with positional information allow implementation of the IIIF Content Search API, providing users with fast and accurate highlighting of search terms within document fragments.*

**Keywords:** Digital archives, IIIF, OCR

**Introduction**

Archives have been a part of human culture for 5,000 years, preserving primary source records selected for their cultural, historical or evidentiary value. While we no longer keep records on stone tablets or papyrus, digital approaches to preserving these unique documents presents a number of opportunities and challenges in which archives can serve the scholars and academics who use them.

Certainly, these challenges are not entirely new - though the ever-increasing volume of documents stored by digital archives risks obscuring the signal with an abundance of noise. With more and more archives seeking to open up their archives, exposing documents that have been stored in basements and previously accessible only in reading rooms, we'd like to focus on some of these challenges and how we can lean into them for the benefit of digital scholars.

We are part of the team at Cogapp, a digital development agency based in Brighton, UK, that has worked alongside cultural institutions since the early 1980s. Since it's inception, Cogapp has worked alongside our partners to provide expertise in designing and deploying digital strategy, developing bespoke software and infrastructure solutions, and integrating disparate systems to provide unified, rich and engaging user experiences.

This paper will use examples from our own work to document and explore how we can to unlock archives in various ways. The archives we will draw examples from are the Qatar Digital Library (QDL - https://www.qdl.qa), the Endangered Archives Programme (EAP - https://eap.bl.uk) and the Arabian Gulf Digital Archive (AGDA - https://www.agda.ae/).

**Unlocking Mobile Devices - Responsive Design**

Archives are presenting large volumes of textual and visual information which can be hard to convey on mobile devices.

The percentage of mobile phone users who access the internet from their phone is increasing year on year. Globally, an estimated 63.4% of users will use their phone for internet access at least once every month in 2019[1]. Users expect that their browsing experience will not be meaningfully degraded due to them using a mobile device.

When considering archives, this could be a digital scholar quickly checking references or confirming sources while away from a laptop or desktop. It could be a researcher or student whose primary access is through a mobile device.

Each of these users need to have fully functional interfaces to access archival content from whichever device they are using.

The 'Explore the Archive' page on the QDL[2] exposes a lot of tools to our digital scholars. As well as the free text search box, we have detailed facets on the left, extracted from the cataloguing data, and a timeline date filter, which is initially hidden when a user arrives on the page.



The mobile user should, as much as possible, not be disadvantaged by their choice of device. However, exposing all of these options would be overwhelming on the screen. Instead, they are initially concealed behind the 'Refine Your Search' label while the free text search box and result summary are front and center.

Working this way allows a user to choose the most convenient and appropriate method to search for the information they need. It may be a search begun on a mobile device could be continued on a desktop or vice versa. For our user, the move between devices would make sense and be perfectly familiar, and we ensure that all of our search results addresses are persistent and bookmarkable to allow exactly this sort of transition between device types.
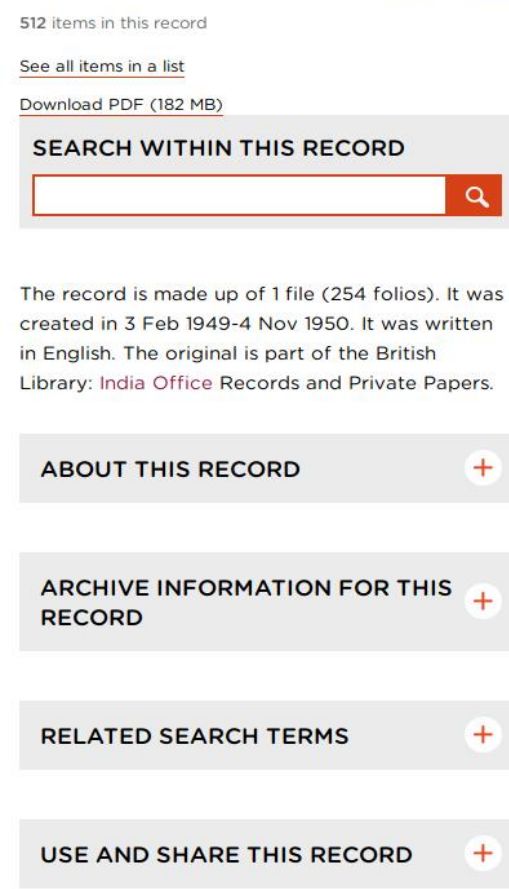
Pictured here is an archive item from the QDL site, pictured on both a desktop device and a mobile device[3]. Both versions of this page provide the same information that a scholar may need. The viewer is at the top to scroll through the archive item pages and the summary of the record is below.

In each view, the detailed archival information is initially concealed to stop the page being overloaded with text. Each contains the same information and could be equally useful to researchers accessing this item.

It might be thought that mobile use on archives would be not track with general trends. This, however, has not been our experience. With the archives we have built and continue to manage, the use of mobile and desktop are roughly equivalent. Users are also spending roughly the same amount of time on the archives whether they are using mobile or desktop devices.

These users are viewing a number of pages on each archive before they leave and they have all the tools they need to allow them to search and examine documents that may be of interest.

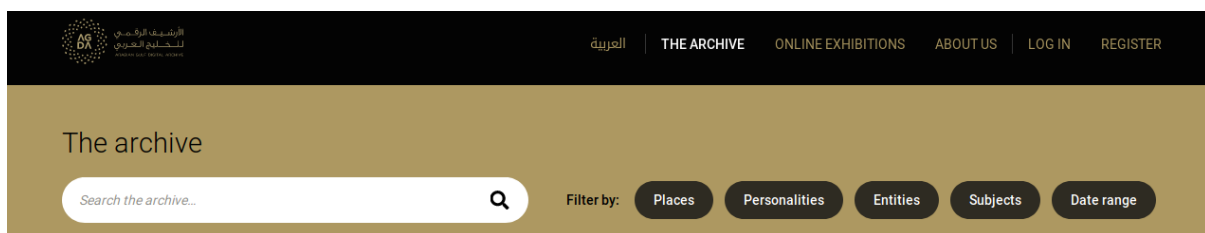**Unlocking the Metadata – harvesting and presenting**

Great archives depend on great cataloguing. A digital archive needs to be able to consume this cataloguing data, potentially from several sources, and allow it to be exposed to users in ways that make the discovery of relevant documents intuitive, straightforward and helpful.

For a digital archive to have the best chance of being successful, cataloguing partners need to use consistent standards. We have experience working with a number of standards as we work with a variety of different libraries and archive. In projects in which we have several digitising and cataloguing partners, we even support a variety of different formats simultaneously. Our observation is that the choice of cataloguing standard is less important than the consistency with which it is applied.

Having gathered this data, archives can make the choice to enrich this data with other sources. We also have experience using machine learning to aid in this effort[4]. Equally, using Wikidata or other authoritative, open-source, third-party sources to add more detail to an archive makes the job of the researcher easier.

Once the data has been ingested by a digital archive, it needs to be presented and discovered by users. Some of the ways an archive can meet the needs of digital scholars are:

- Making the search fast and intuitive. We tend to prefer open-source solutions and have experience using both Apache Solr and Elastisearch to serve our clients. Each of these software solutions allow for complex searches against a large selection of documents in a fast and reliable way.

- Exposing the cataloguing information the form of detailed search facets allows for the user to get to the results they need more easily.



The search filters on the AGDA are categorised and displayed discretely beside the free text search box[5]. A scholar can click on any of these buttons and receive a comprehensive set of filters based on the results returned by the archive.

**Archive type**

Project (369)

Collection (4,805)

File (155,182)

Item (207,037)

Web page (20)

**Content type**

Photograph (175,112)

Newspaper (57,552)

Periodical (30,672)

Manuscript (24,895)

Document (11,966)

Book (11,905)

Correspondence (6,096)

Volume (3,886)

Report (2,904)

Register (2,228)

Show more

**Related places**

India, Asia (58,200)

FILTER BY DOCUMENT SOURCE:

(-) Items in the archive

FILTER BY SEARCH TERM LOCATION:

Archival description (16,867)
Transcription (1,775,457)

FILTER BY SUBJECT:

Second World War (1939-1945) (236)
Diplomacy (216)
Petroleum industry (189)
Reports (179)
Treaties (167)

Show more

FILTER BY PLACE:

Persia (540)
Bahrain (357)
Persian Gulf (348)
Muscat (263)
Kuwait (246)

Show more

FILTER BY TYPE:

The facets for both the QDL and the EAP are presented to the left of the search results page. A small selection of the top results are shown with the option for the user to expand to see more.

- When a scholar has discovered a relevant document, they need to be able to cite that record and trust that it will be there for future reference. For any item on the QDL, under the 'Use and Share' section, there is information on how to cite the record and how to link to it. We work hard to ensure that the link never changes for  an item so that scholars can return to a given item when they need to in the future.

- The presentation of the archival information should be clear and consistent, following the standards of the institution and be consistent across the archive.

**ABOUT THIS RECORD**                                                    ✕

**Content**

This file contains correspondence between British officials regarding Britain supplying Qatar with arms and ammunition.

The correspondence is primarily between officials at the Political Agency in Bahrain, the Political Residency in Bahrain and the British Agency in Doha. The file also contains letters sent to the Political Agency in Bahrain by the ruler of Qatar, Shaikh ʿAbdullāh bin Jāsim Āl Thānī as well as his son (and successor) ʿAli bin ʿAbdullāh Āl Thānī. These letters are in Arabic and accompanied by English translations.

The file contains correspondence that discusses Shaikh ʿAbdullāh's desire for arms, exactly what type of arms and ammunition should be provided and the importation of arms into the country, as well as the broader political context of these events.

**Extent and format**

1 file (160 folios)

**Arrangement**

The papers are arranged in approximate chronological order from the front to the rear of the file.

**Physical characteristics**

Foliation: the main foliation sequence (used for referencing) commences at the front cover with 1 and terminates at the inside back cover with 160; these numbers are written in pencil, are circled, and are located in the top centre of the recto side of each folio. An additional foliation sequence is also present in parallel between ff 2-142; these numbers are written in a combination of blue ink and pencil, but are not circled.

**Written in**

English and Arabic in Latin and Arabic script

**Type**

Archival file

- It is not a great feeling to reach a dead-end in a search. The metadata can be used to provide suggestions and onward journeys for the researcher, surfacing connections that may have otherwise remained hidden. This could be leveraging the cataloguing, recent searches or even other media beyond the archived items.

**RELATED SEARCH TERMS**                                                 ✕

**Subjects**

Weapons  Arms sales

**Places**

Qatar

## Unlocking the images – IIIF

Once the metadata has been unlocked, it can be used to surface archive items. The richest resource is the digital images of the archive items themselves. A lot of image-based archive data on the Internet, be that photographs, books, newspapers or manuscripts, are stored behind paywalls with access restricted to bespoke, locally built applications. The International Image Interoperability Framework (IIIF)[6] is a set a shared application programming interface (API) specifications that aims to provide ways to consolidate this data and make research, scholarship and the transmission of cultural knowledge more open.

Using IIIF has a huge number of benefits for online archives and their users.

For the end user, it is much more straightforward to provide advanced, interactive functionality such as zooming, selecting and comparing documents.

While the institution can be confident that they are building a resilient and flexible archive that can grow and change in step with their needs. Since IIIF is open-source there is no proprietary lock-in and there is an ability to combine content from different sources for comparison and consolidation.

There is also a reduction of long term costs as licensing and operational costs are low and predictable over time.

Leveraging IIIF and cloud-based solutions, we were able to build and deliver the system for EAP, which initially consisted of over 200Tb of TIFF-format master images, within two months of starting work[7].

## Unlocking the viewers – Presentation and Image API

Now the images are unlocked and available for manipulation and sharing, other tools are now open to us. One of the benefits of using a well-supported open-source framework is that we can build on the work of others within the community.

## Viewers



The Universal Viewer[8], seen here on the EAP project[9], is a viewer which can be used to display a variety of file formats, one of which is IIIF image manifests. This rich viewer serves zoomable images that can be easily embedded elsewhere. It can be styled to fit with the brand of your archive, and it can be translated to any language including right-to-left scripts.



The viewer is used on a number of archives already, so the benefit of using it on a new archive would be familarity for the users. If you are a scholar who has gained proficiency with this viewer on one archive, you can apply that skill set to this new archive.

There are other viewers available which you could implement and test. The Mirador[10] viewer is specifically designed to be a IIIF client and has some powerful tools to allow scholars to be productive.



A user can compare different elements of different documents side-by-side and zoom each of these independently. There are also advanced capabilities to alter the brightness and contrast of the image which sometimes allows previously unobserved elements of documents to stand out.

Each of these viewers can be integrated into a new archive with a minimum of code and configuration, once the digital assets are being served using a IIIF manifest.

Now that we can access high-resolution, specific zooms of images, we can use these to tell richer stories about the content.

This is a cover page from a medical textbook that is archived in the QDL[11]. During a hack day in March 2019, a team explored ways to detail the inscriptions around the outside of this page.

The writing around the outside of this panel details the ownership of the book over the years. This provenance information had previously been captured and added to detail of the document.

**Content**

The Canon of Medicine (القانون في الطبّ) by Abū ʿAlī al-Ḥusayn ibn ʿAbd Allāh (أبو علي الحسين بن عبد الله), known as Ibn Sīnā (ابن سينا), latinised as Avicenna, 980-1037).

An unsigned note, without date, in the upper margin of f. 1r states that the physician Ibn al-Nafīs (ابن النفيس, d. 1288), here referred to by his full name Abū al-Hasan ʿAlāʾ al-Dīn ʿAlī ibn Abī al-Ḥazm al-Qurashī (أبو الحسن علاء الدين علي بن أبي الحزم الشهير بالقرشى), claimed that this copy was produced by Avicenna himself.

Begins (f. 1v, lines 1-3):

وبعد فقد التمس ...

منى بعض خلص إخوانى ومن يلزمنى إسعافه فيما يسمح به وسعى أن أضف فى الطب كتاتا مشتملا على قوانينه الكلية والجزؤية اشتمالا يجمع إلى الشرح الاختصار

وإلى ايفا الأكثر حقه من البيان

Ends (f. 274r, line 9):

الثلث أوثولات سبع قراريط المواتوس أوقية ونصف المواتوسين ثلث أواق تمنه أوثولو أربعة وغرمى قيراطان
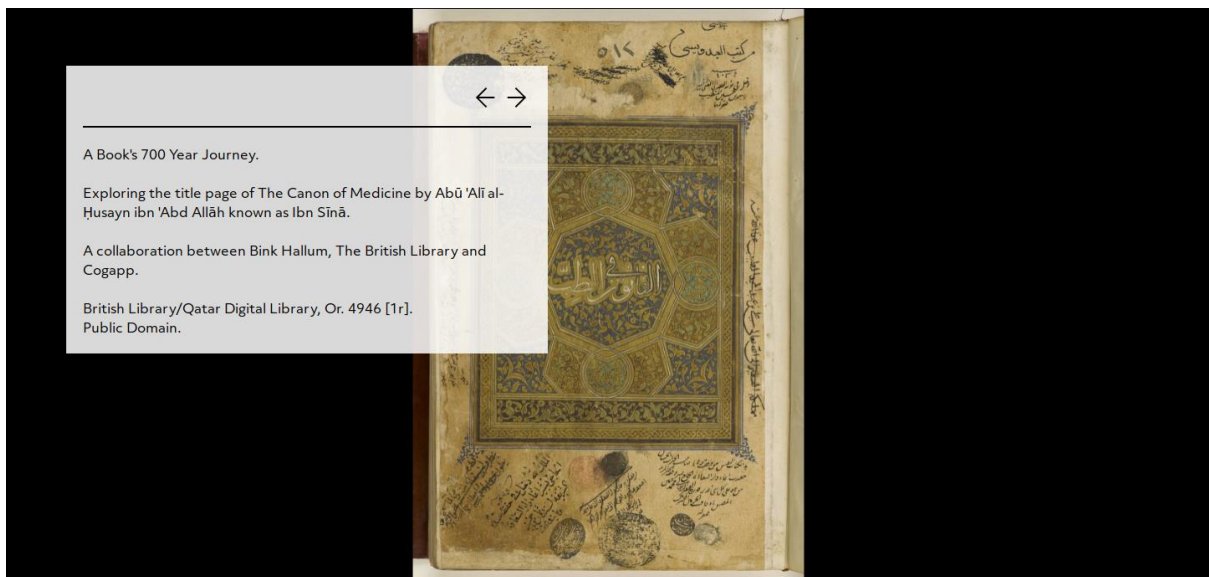
Colophon (f. 274r, line 10):

تم كتاب القانون فى علم الطب

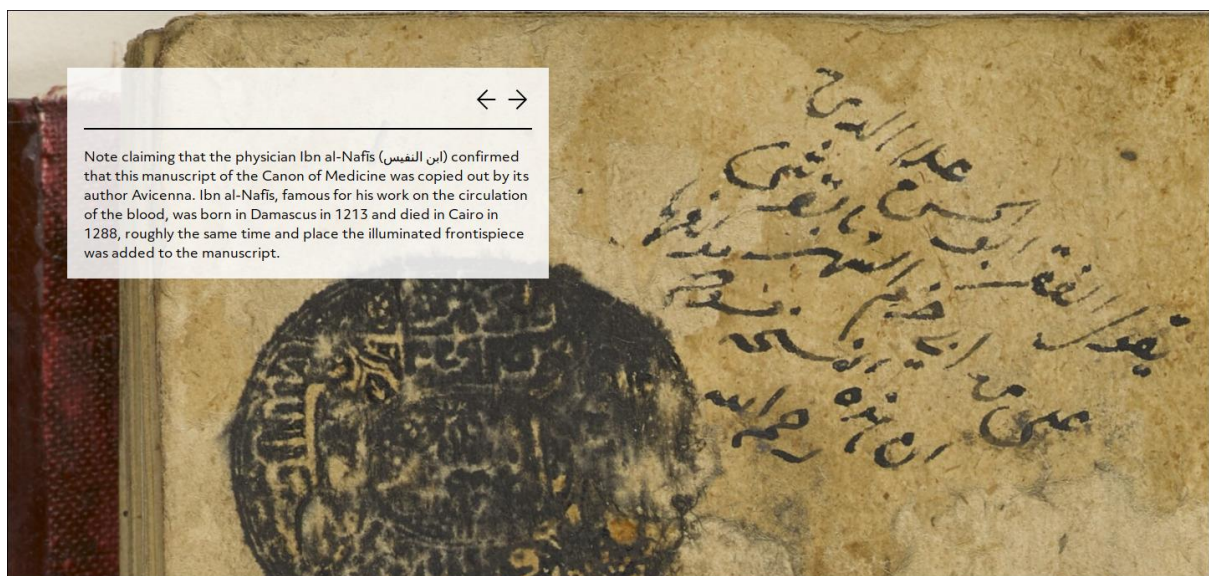after these words another hand writing with another ink writes

على يد مؤلفه

The details as they are written here can be difficult to grasp and visualise.

It is not immediately clear which piece of writing is being referred to or how interesting and engaging the details actually are.

A Book's 700 Year Journey.

Exploring the title page of The Canon of Medicine by Abū 'Alī al-Ḥusayn ibn 'Abd Allāh known as Ibn Sīnā.

A collaboration between Bink Hallum, The British Library and Cogapp.

British Library/Qatar Digital Library, Or. 4946 [1r]. Public Domain.

By leveraging IIIF and some work we have been doing in other projects, it is possible to expose this information, zooming to and explaining parts of the text in detail[12]. In zooming to the relevant part of the image, the commentary and description makes more sense and is more engaging.



Note claiming that the physician Ibn al-Nafīs (ابن النفيس) confirmed that this manuscript of the Canon of Medicine was copied out by its author Avicenna. Ibn al-Nafīs, famous for his work on the circulation of the blood, was born in Damascus in 1213 and died in Cairo in 1288, roughly the same time and place the illuminated frontispiece was added to the manuscript.

## Unlocking the OCR transcripts - Content Search API

For typewritten script, OCR has a fairly high degree of accuracy. When combined with positional data, often through an ALTO[13] file, and IIIF images, more tools can be unlocked for your scholars.



Here, in the Qatar Digital Library, we see a search for the word 'elephant' within a specific archive item[14]. As well as displaying the matching text fragment, we have a 'Quick look' selector. When a user hovers over this, the relevant snippet appears with the word highlighted.



Here, in the Arabian Gulf Digital Archive, we used Content Search API and the positional data to highlight both the OCR text and the digitised document at the same time.
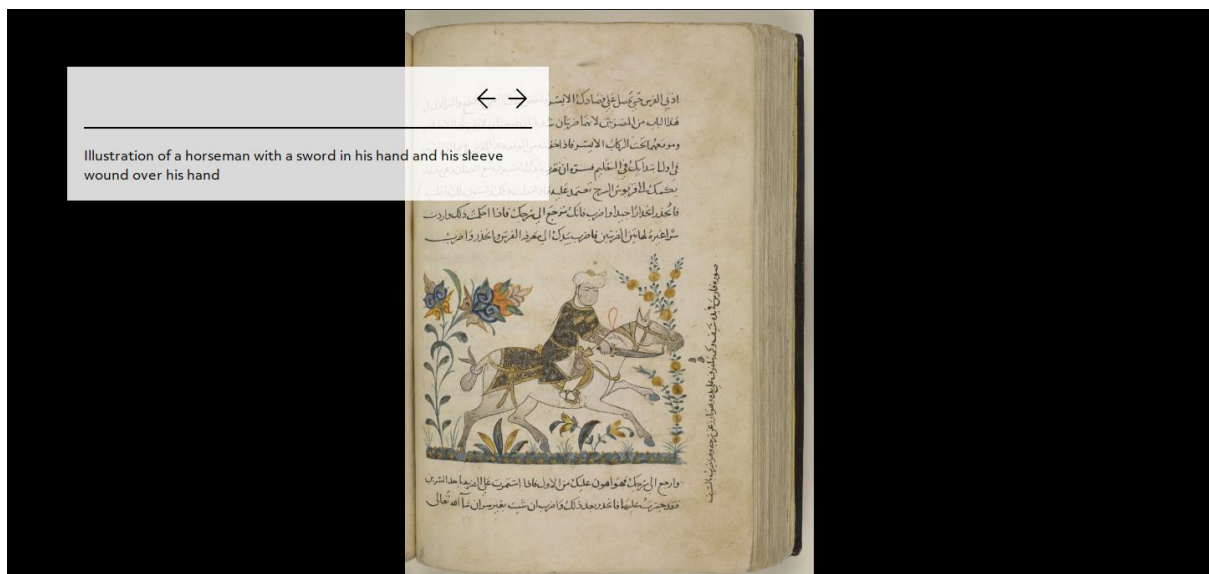
**Unlocking the Fun with IIIF**

Specifications are not always the most inspiring and interesting of documents. It can be hard to workout how these standards could have a positive impact on your archive and institutions.
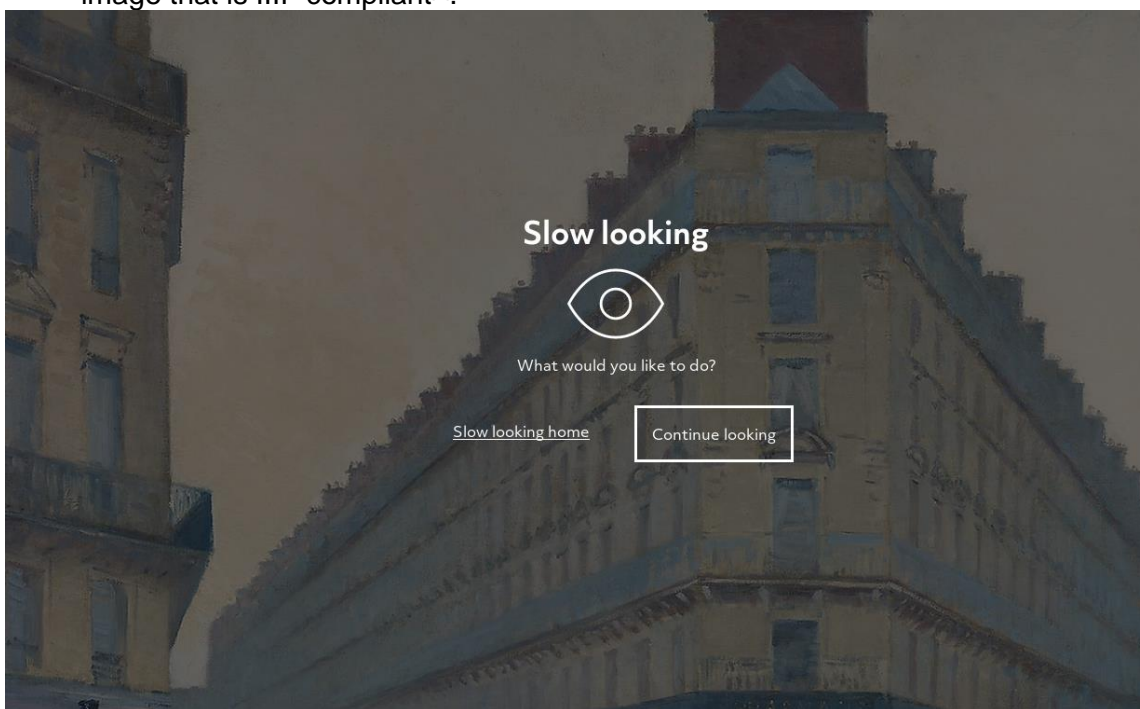
Approaching these in a playful way, however, can help to unlock some flashes of brilliance that can have unexpected and positive impacts on user experience.

We have regular hack days both internally and with partners such as the British Library and The National Archives. We often will explore projects based on ideas we have about IIIF. Some of those could be lifted and developed as they are into elements of an archive, others are more fun but have seeds of understanding and exploiting the technology more.
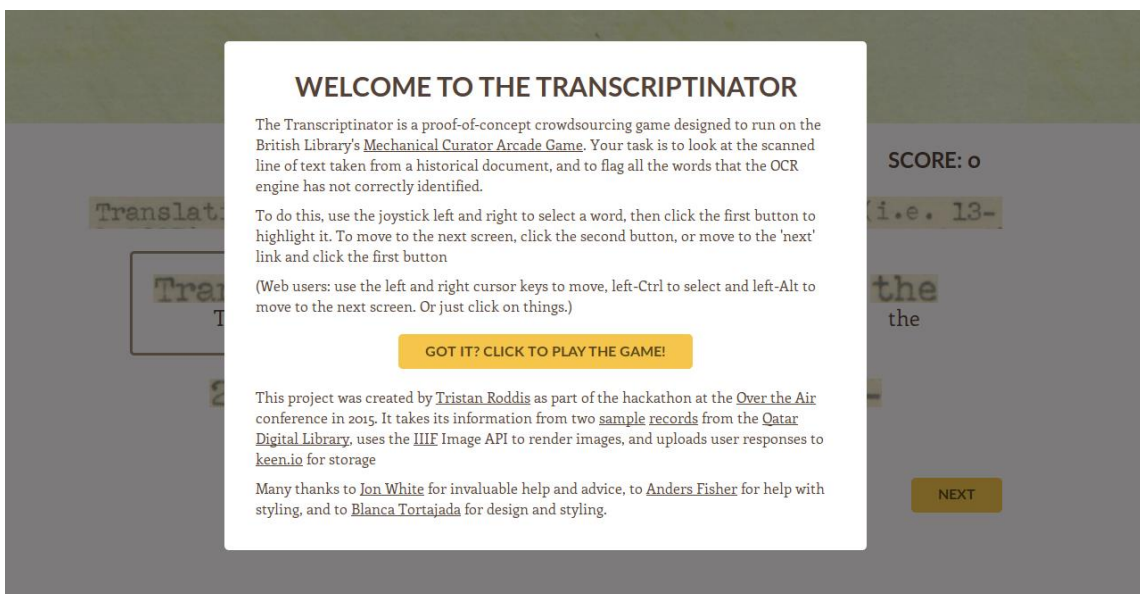
- Storiiies is a way to tell stories using images and words. Since it is using IIIF rather than a closed system these storiiies can be shared and developed by others[15].



Illustration of a horseman with a sword in his hand and his sleeve wound over his hand
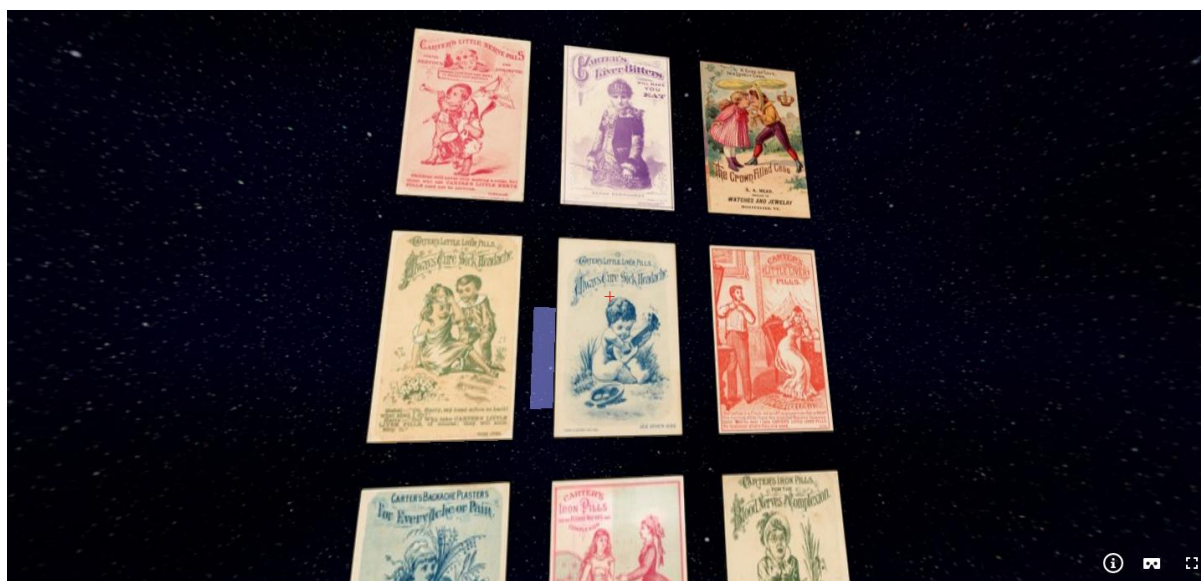
- Slow looking is a way to experience an image in a more contemplative and measured way. It is a relaxing and full-screen experience and can be used with any image that is IIIF compliant[16].



- Transcriptinator is a proof-of-concept crowdsourcing game. It asks the user to identify if the OCR has correctly transcribed a word through a game-like interface[17].

- Trade Cards Explorer places digitised Victorian trade cards into a 3D environment, viewable in-browser or by using VR goggles[18].



These are just some of the outcomes from our hackdays. The output feeds into our work and helps keep our creativity unlocked.

**Conclusion**

We have looked at some of the ways to unlock archives in this paper. Archives provide an essential role in helping us understand our historical and sociological contexts. When archives are well catalogued and helpfully unlocked, with clutter-free and intuitive designs, scholars of all backgrounds can dive deeply, make connections and move their respective fields forward.

Using open standards such as IIIF allows for co-operation between institutions in an unparallelled way.

Furthermore, using playful approaches in development and design can surface richer and more engaging solutions for end-users.

## Acknowledgments

Special thanks to our partners and colleagues whose work we have built on in this paper.

## References

[1] <https://www.statista.com/statistics/284202/mobile-phone-internet-user-penetration-worldwide/> [accessed 20 May 2019]

[2] 'Explore the Archives', in *Qatar Digital Library* <*https://www.qdl.qa/en/search/site/?f[0]=document_source:archive_source*> [accessed 20 May 2019]

[3] 'File 4/49 II Foreign Office Account.', British Library: India Office Records and Private Papers, IOR/R/15/2/1158, in *Qatar Digital Library* <https://www.qdl.qa/archive/81055/vdc_100000000241.0x00031d> [accessed 20 May 2019]

[4] 'Automated image analysis with IIIF', in *Cogapp Blog* <https://blog.cogapp.com/automated-image-analysis-with-iiif-6594ff5b2b32> [accessed 20 May 2019]

[5] The archive, in *The Arabian Gulf Digital Archive* <https://www.agda.ae/en/search> [accessed 20 May 2019]

[6] IIIF <https://iiif.io>

[7] Roddis, Tristan and Farquhar, Adam. "From at risk to open access: The Endangered Archives of the world." *MW18: MW 2018*. Published February 4, 2018. Consulted May 21, 2019. <https://mw18.mwconf.org/paper/from-at-risk-to-open-access-the-endangered-archives-of-the-world/>

[8] Universal Viewer <https://www.universalviewer.io>

[9] Gaskiya ta fi Kwabo [1939-1958], in *Endangered Archives Programme* <https://eap.bl.uk/archive-file/EAP485-2-1> [accessed 20 May 2019]

[10] Mirador Viewer <http://projectmirador.org/>

[11] al-Qānūn fī al-ṭibb الﻗﺎﻧون ﻓﻲ اﻟطب Avicenna اﺑن ﺳﯾﻧﺎ [1r] (12/564), British Library: Oriental Manuscripts, Or 4946, in *Qatar Digital Library* <https://www.qdl.qa/archive/81055/vdc_100052704274.0x00000c> [accessed 20 May 2019]

[12] 'Hold the Front Page', in *Cogapp Blog* <*https://blog.cogapp.com/hold-the-front-page-125f140147c2*> [accessed 20 May 2019]

[13] ALTO – Technical Metadata for Layout and Text Objects
<https://www.loc.gov/standards/alto/techcenter/elementSet/index.html> [accessed 20 May 2019]

[14]  'Ralph Fitch, England's Pioneer to India and Burma. His companions and contemporaries. With his remarkable narrative told in his own words', British Library: Printed Collections, T 36804, in *Qatar Digital Library*
<https://www.qdl.qa/en/archive/81055/vdc_100023516428.0x000001/search/elephant>
[accessed 20 May 2019]

[15] Storiiies – Experiments in Digital Storytelling <https://storiiies.cogapp.com> [accessed 20 May 2019]

[16] Slow looking <https://slowlooking.cogapp.com> [accessed 20 May 2019]

[17] Transcriptinator <https://labs.cogapp.com/transcriptinator/> [accessed 20 May 2019]

[18] Trade Card Explorer <https://labs.cogapp.com/tc/> [accessed 20 May 2019]