# Inter-connected Network of Knowledge – The NLB Journey

**Siang Hock KIA**
Technology & Innovation, National Library Board, Singapore
kia_siang_hock@nlb.gov.sg

**Yi Chin LIAU**
Technology & Innovation, National Library Board, Singapore
liau_yi_chin@nlb.gov.sg

**Ian ONG**
Technology & Innovation, National Library Board, Singapore
ian_ong@nlb.gov.sg

**Abstract:**

*Libraries and archives collect and provide access to digital resources of significant national, cultural and heritage values. Their collections are painstakingly curated, described and made accessible.*

*NLB adopted a holistic digital strategy that involves the digitisation of uniquely Singapore content, findability enhancement via search engine optimisation, and multi-screen delivery. The next leap is to link the digitised content into a network to enable contextual discovery.*

*NLB has successfully created high quality associations amongst its digital resources. Text analytics technologies were used to automatically sieve through millions of items, and cluster related resources together. The result is a massively inter-connected network of knowledge. The users are now able to view and explore related resources within and across collections and media formats. Information is no longer viewed in isolation, but seen as a part of the larger context.*

**Keywords:** contextual discovery, knowledge network, text analytics, digital strategy.

# 1 INTRODUCTION

The National Library Board of Singapore (NLB) oversees the National Library, the Public Libraries and the National Archives. The NLB's mission is to provide a trusted, accessible and globally-connected library and information service through the National Library and a comprehensive network of Public Libraries.

Through its innovative use of technology and collaboration with strategic partners, NLB ensures that library users have access to a rich array of information services and resources that are convenient, accessible and relevant.

A key function of NLB is to connect people to the country, by connecting the past to the present. As the key memory institution in Singapore, NLB collects and provides access to digital resources of significant national, cultural and heritage values. Akin to the corporate memories of enterprises, the NLB collection constitutes a major portion of the national memory of Singapore.

# 2 THE NLB DIGITAL STRATEGY

A seismic shift in the information seeking behaviour has taken shape. With the ever-expanding reach of the Internet, users are accustomed to quick and easy access to content, and are expecting an online experience that is rich and instant. This is a generation that is always connected, through multiple devices.

NLB saw great opportunities to ride on the confluence of technological advances and the connected lifestyle to bring relevant information to the users. Its digital strategy involves:

- digitising uniquely Singapore content
- making NLB content findable through the popular search engines
- delivering content on all devices through responsive web design

'Content is king' in the Internet era. NLB is in the privileged position to acquire and digitise uniquely Singapore content, through its statutory standing and through the extensive network of partners of content owners. Significant resources have been invested in the digitisation of valuable content, and the efforts continue unabated. More details on the NLB digital collections can be found in the next section of this paper.

Most users go to the popular search engines (e.g., Google) as their first port of call to search for information. Leveraging on this behaviour, NLB embarked on sustained search engine optimisation (SEO) efforts to ensure that its digitised resources are easily findable via the popular search engines. SEO and the user friendly content sites resulted in a wider reach and higher usage of these valuable digital resources. An oft-cited example is the increased usage of the Infopedia[1] service from 400 to 150,000 pageviews a month after the trial (Chellapandi, Chow & Tay, 2010). The average monthly pageviews in 2013 is 270,000.

---

[1] http://eresources.nlb.gov.sg/infopedia/

NLB has also adopted the Responsive Web Design framework since 2011, and implemented device-aware online viewing and streaming capabilities to ensure that the user experience on any device will be optimal.


## 3 THE NLB DIGITAL COLLECTIONS

Over the years, NLB has built up a huge collection of content to meet the diverse needs of its patrons. When the National Archives of Singapore (NAS) joined the NLB family in November 2012, it brought with it a huge and highly valuable collection of primary and unique materials on Singapore's history. Aggressive digitisation continues unabated to preserve as well as give wider access through more platforms. The digitised contents from NAS and NLB constitute a valuable repository of the memories of the nation.

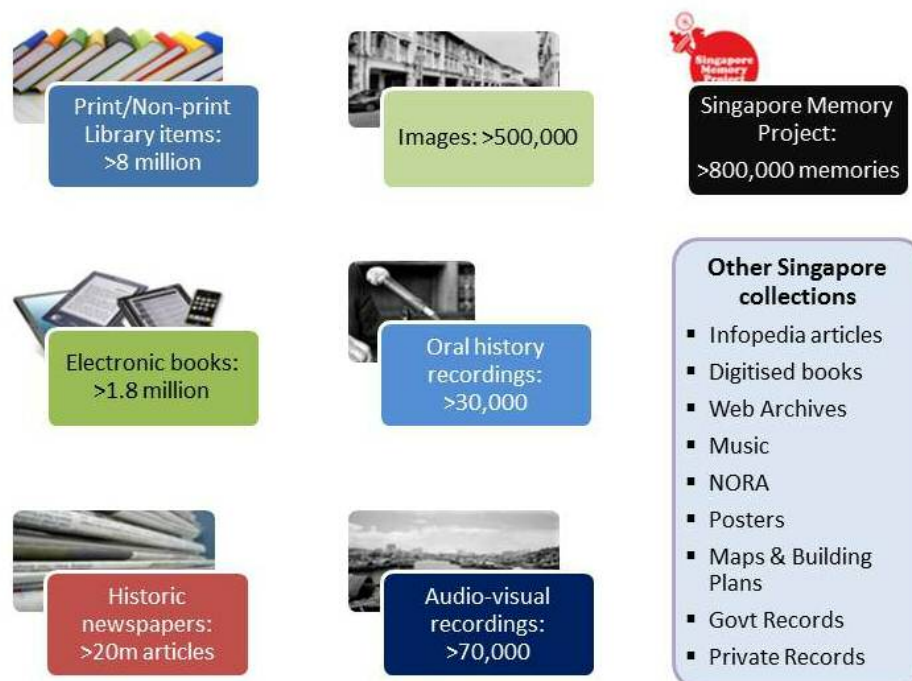Figure 1 shows the collections available for online and onsite access.



Figure 1: The NLB collections


## 4 CONTEXTUAL DISCOVERY

Riding on the success of the digital strategy, NLB looked for the next 'big thing'. Every year, NLB users collectively contribute to tens of millions of e-retrievals. We see every single one of these interactions as a golden opportunity to 'push' relevant content to the user. The success of the 'customers who bought this item also bought' recommendation feature at amazon.com is a clear testament of the power of pushing relevant recommendations.

However, to effectively connect people to knowledge, we need to connect knowledge to knowledge first. Figure 2 shows a subset of the related content within NLB's collections on the Cenotaph monument as described in the Cenotaph article highlighted in a red box.
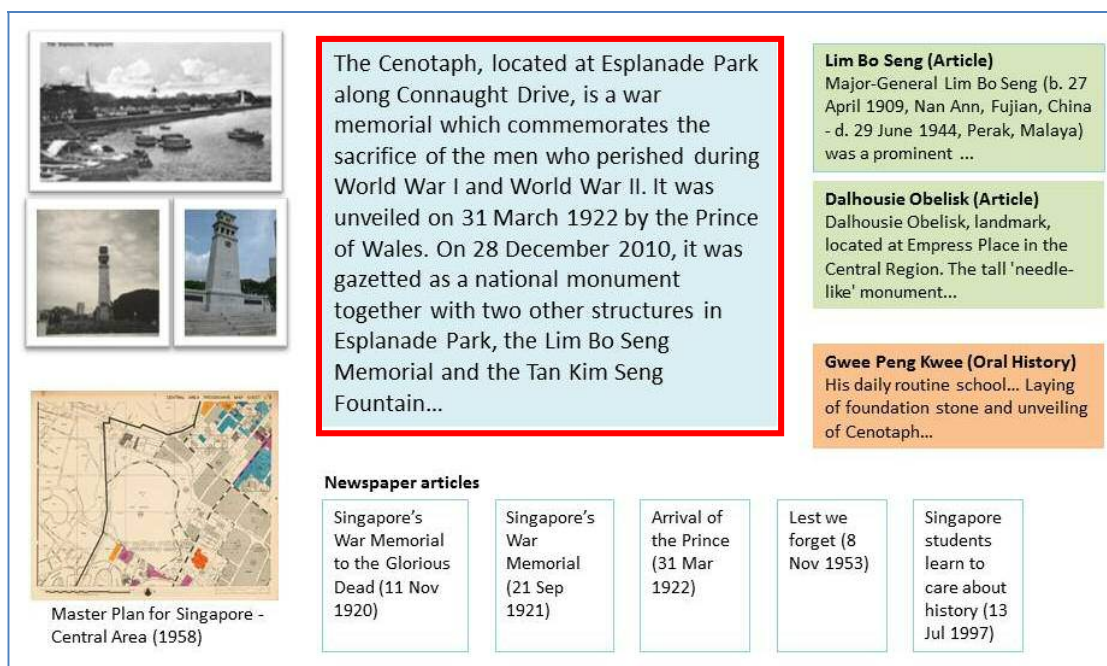
Figure 2: Contextual discovery through the recommendations of related content

To identify all these (and more) related resources would have required many searches to be performed with various search keywords, and sieving through hundreds if not thousands in the search result lists. This is very time-consuming. Given the current 'good enough' information seeking behaviour, it is unlikely that the typical users will be willing to spend the efforts needed to get to such a comprehensive set. It would be such a pity as they would be missing many treasures in the trove of digital resources available.

If we can 'push' the kind of relevant and related resources as illustrated in Figure 2 for each and every piece of digital resources in our collection, the users will have easy access to a comprehensive set of resources that provides a more complete 360 degree perspective of a topic.

Moreover, the Cenotaph article that the user is viewing (as a result of a search or browse action that the user just performed) provides a clear context of the user intent, which increases significantly the likelihood that the user will click through the recommendations, and discover many more resources in our collections. We call this 'contextual discovery', and believe strongly that it will be critical to the next generation of digital libraries.

This expands the NLB digital strategy, described in Section 2 above, by connecting the resources within and across collections, formats and languages to provide a rich contextual discovery experience.

## 5  USING TEXT ANALYTICS TO IDENTIFY RELATED CONTENT

The foundation to contextual discovery is the associations between the digital resources. Each association established a link between the resources, and together, a network of the resources is formed.

The value of the resources increases with the number of links, according to the 'network effect' phenomenon.  With a total collection size that goes into tens of millions, it will not be cost-effective to identify the associations manually.

Text analytics technologies have been identified to be suitable and scalable to perform this task (Lim & Chinnasamy, 2013).

## 5.1   The Proof-of-Concept

To test the feasibility of the use of text analytics to accurately identify related content, a proof-of-concept (PoC) was conducted.

The open source machine learning software Mahout[2] from the Apache Software Foundation[3] was selected for the PoC.  Newspaper articles for the year of 1989 from The Straits Times totalling over 50,000 were used.  Figure 3 shows the steps involved in the text analytics processing.
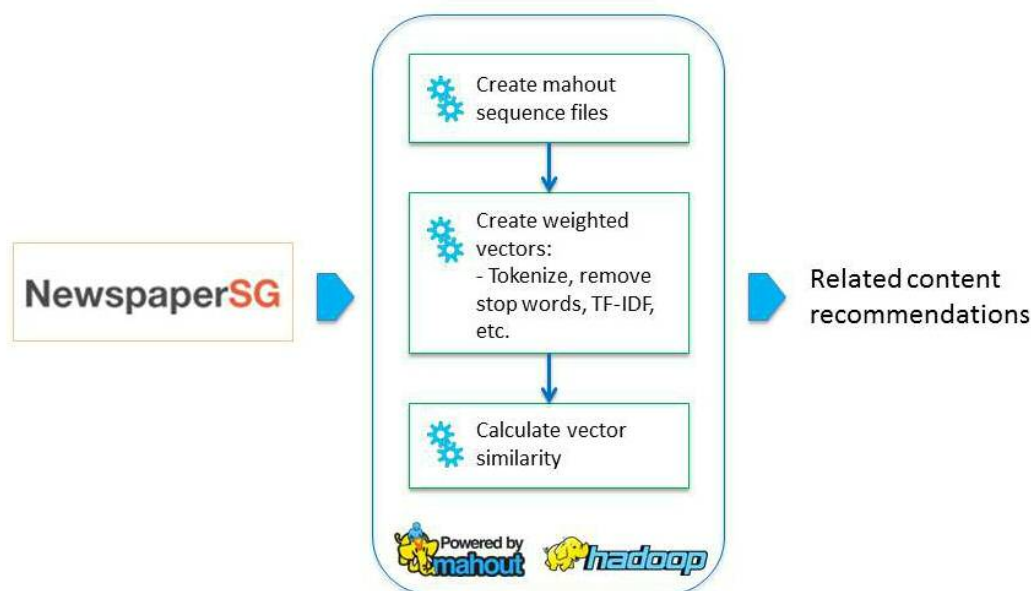


Figure 3: Text analytics PoC using Mahout

The text analytics processing generated a set of the most 'similar' newspaper articles for each of the 50,000 newspaper articles used in the PoC.  The similarity between 2 articles is measured by a value from 0 to 1, with 1 being a perfect match.

Take the article 'Goldsmith group to re-launch takeover plan' in Figure 4.   The top recommendations from the text analytics software are clearly about the takeover event at the British American Tobacco (BAT).  When the articles are organised in a chronological order (following the red arrows), the event unfolds.

---

[2] http://mahout.apache.org
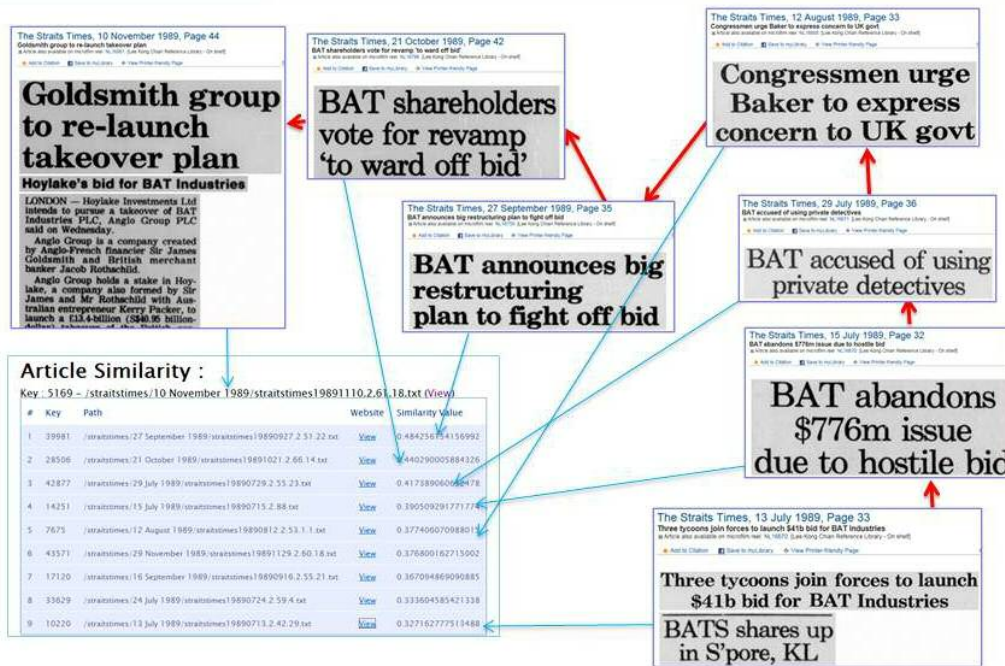[3] http://www.apache.org

Figure 4: Similar newspaper articles identified via text analytics

Algorithms implemented in Mahout, such as Term Frequency/Inverse Document Frequency (TF/IDF), a statistics that reflects the importance of a word in a corpus, and Cosine Similarity for text-based vector similarity are well-established in the domain of information retrieval.

The PoC confirmed that the use of text analytics algorithms can automatically, efficiently and accurately identify similar articles.

## 5.2 The Deployment

With the success of the PoC, NLB proceeded to implement the concept for its collections. We decided to start with the Infopedia collection given the smaller size of the collection (of around 1,800 articles) and the well written and fairly long text available. The text analytics processing was completed within an hour.

The recommendations were checked rigorously by the librarians on the staging Infopedia site. We went 'live' in June 2013 after the go-ahead was given (see Figure 5).

Thereafter, we extended the feature to other NLB digital collections, including PictureSG[4] and Singapore Memory Portal[5].

---

[4] http://eresources.nlb.gov.sg/pictures/
[5] http://www.singaporememory.sg/

Figure 5: Infopedia recommendations via text analytics

At the onset of selecting the text analytics software for the PoC, one of the criteria is the ability to scale to handle large data sets that go into millions. Most of the algorithms supported by Mahout have been developed to operate in an Apache Hadoop[6] cluster. Hadoop is a popular open source software framework for distributed storage and large scale processing of data sets on a cluster of commodity hardware.

As we progress to process collections that are much larger, we started to implement a Hadoop cluster. The current Hadoop cluster comprises 13 virtual servers implemented in the NLB private cloud. This allows NLB to scale and re-configure the cluster quickly depending on the resource requirement of the analytics to be performed. This is critical as analytics algorithms are generally very resource intensive.

## 5.3  Handling large data sets

The 13-nodes Hadoop cluster proved to be unable to process collections that go beyond hundreds of thousands items. While technically feasible, it is not practical to continue to add nodes and storage to the Hadoop cluster.

It is in fact not necessary to compare every pair of resources as the bulk of the records in a huge data set will bear little similarity with one another. A better approach would be to break the large data set into smaller clusters of related content that can then be efficiently processed.

The good news is that Mahout also comes with the necessary clustering algorithms. With this, we went ahead to process the key English newspapers (6.7 million articles) and Chinese

newspapers (2.3 million articles) in NewspaperSG[7] and NAS Archives Online[8] collections (1 million records).

Table 1 shows the top 50 terms within example clusters for the English newspaper articles. The clustering certainly worked well. Note that stemming was done during the clustering process.

| Size of cluster | Top 50 terms |
|---|---|
| 52,678 | exhibit, art, artist, paint, museum, singapor, work, displai, galleri, open, mr, year, organis, centr, chines, cultur, held, on, nation, world, pictur, photograph, collect, hall, colour, includ, first, time, featur, sculptur, design, dai, peopl, piec, two, societi, part, visitor, fair, intern, painter, life, school, prize, old, show, road, book, trade, today |
| 86,881 | olymp, athlet, game, sport, medal, event, team, gold, championship, record, world, singapor, metr, swim, won, year, champion, win, nation, women, time, coach, asian, meet, competit, train, swimmer, two, second, race, compet, first, amateur, bronz, intern, associ, best, finish, yesterdai, silver, old, on, relai, men, set, track, medallist, mark, south, run |
| 142,289 | school, student, educ, teacher, univers, secondari, singapor, children, year, primari, studi, pupil, parent, mr, teach, colleg, cours, ministri, english, languag, chines, on, institut, examin, time, learn, princip, train, programm, work, graduat, help, nation, class, two, govern, girl, scienc, boi, first, level, malai, centr, make, dr, academ, dai, organis, scholarship, junior |
| 125,629 | polic, arrest, offic, suspect, two, men, yesterdai, man, investig, report, found, mr, gang, road, raid, on, station, crime, detain, year, night, arm, robberi, car, believ, peopl, forc, charg, singapor, robber, hous, todai, told, stolen, seiz, spokesman, old, held, four, escap, murder, reuter, chines, detect, member, street, made, drug, dai, home |

Table 1: Top 50 stemmed terms of example clusters from English newspapers

## 6 BENEFITS

With the use of text analytics, many high quality recommendations can be identified. Before this was implemented, associations were manually identified for some of the collections. However, due to the manual efforts required, the number of recommendations is necessarily limited. Now, the users have access to over a billion quality associations, while the librarians need not spend time identifying them manually.

While NLB is progressively introducing and enhancing the contextual discovery capability to more of its collections, the evidence so far has indicated a significantly higher usage of the digital collections after the contextual discovery implementations.

A good case in point is the cross-collection recommendations to related pictures at the Infopedia articles (red box in Figure 6).

---

[7] http://eresources.nlb.gov.sg/newspapers/
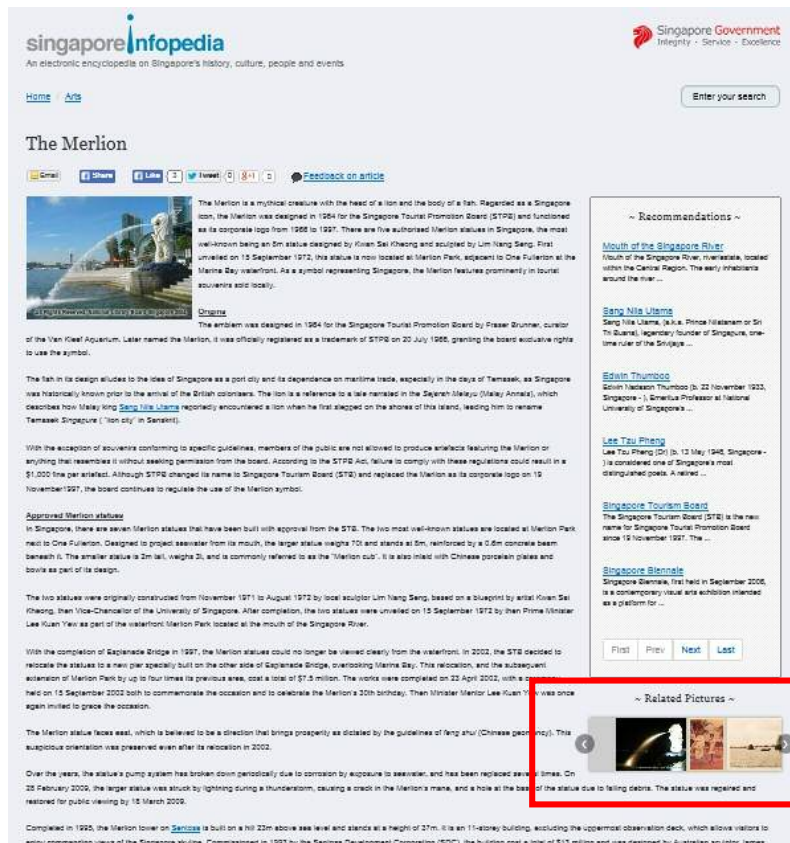[8] http://www.nas.gov.sg/archivesonline

Figure 6: Cross-collection recommendations

The recommendations within the PictureSG collection was implemented in September 2013, while the cross-collection recommendations from Infopedia to PictureSG was launched in November 2013. Table 2 shows some of the web statistics for PictureSG before and after the implementations.

| Month | Total referrals | # from Infopedia | % | Page views | Page views per visit |
|---|---|---|---|---|---|
| Apr-13 to Aug-13 | - | - | - | 37,841 (average) | 3.64 (average) |
| Sep-13 | 58,376 | 84 | 0.14 | 65,036 | 4.00 |
| Oct-13 | 21,088 | 93 | 0.44 | 41,969 | 3.34 |
| Nov-13 | 30,684 | 1,421 | 4.63 | 62,018 | 5.27 |
| Dec-13 | 27,628 | 1,554 | 5.62 | 55,605 | 5.07 |
| Jan-14 | 28,149 | 2,424 | 8.61 | 70,313 | 5.73 |
| Feb-14 | 29,871 | 3,181 | 10.65 | 84,341 | 6.41 |
| Mar-14 | 31,832 | 3,560 | 11.18 | 96,626 | 6.53 |

Table 2: PictureSG web statistics

Within 4 months, over 10% of the traffic came through Infopedia!! The usage and the page views per visit for PictureSG have also increased significantly.

This is certainly very encouraging, and we are very confident that the implementation of the contextual discovery capability across all our collections will provide a significant boost to the already high usage of the collections.

It is also interesting to note that although the average lengths of the textual information available differs significantly between the Infopedia and PictureSG collections, it did not prevent the use of text analytics to identify the associations. The average amount of textual data used for the text analytics for the Infopedia collection was 1,054 words, while that for PictureSG was 47. Similarly, the presence of OCR errors in the NewspaperSG newspaper articles has not been an issue too. There is therefore a fair amount of resilience built into the text analytics algorithms to deal with them.

The text analytics approach provides a quantitative measure of the degree of similarity between two resources. This similarity value is based on established information retrieval models and best practices. It provides us the mean to identify the best recommendations, and also the cut-off point for our recommendations.

Libraries and archives play a key role in cultivating an informed society in a knowledge intensive economy like Singapore. Contextual discovery, with the ease it provides to users to explore and deepen their understanding and insights, enables NLB to connect people to knowledge and to the nation. It makes knowledge come alive, sparks imagination and creates possibilities.

# 7 LESSIONS LEARNT

The implementation of text analytics to enable the contextual discovery of millions of digital resources is not a walk in the park, more like a walk in the dark.

**Start with a Proof-of-concept (PoC)**

NLB decided to perform a PoC to determine the feasibility of the concept. This has allowed us to start small and quickly. The PoC results also helped us to communicate the possibilities to the stakeholders.

As importantly, the PoC has enabled us to develop a better understanding of capabilities of the technologies and tools, and build up internal competency in this domain. This is especially critical in the long run given the shortage of actual hands-on knowledge in this domain locally and in the region.

**Established and affordable software in text analytics are available**

We selected the Mahout as the software for the PoC as it was the most established open source machine learning and data mining software that met our requirements. That the algorithms were implemented on top of the Apache Hadoop framework provided us comfort in terms of scalability.

There were good documentation, including several well written books, on Mahout, which helped us tremendously in our PoC, and subsequent deployments.

**Do not simply add more Hadoop nodes when hitting performance issues**

With the Hadoop platform, it is tempting to simply add more compute and storage resources whenever we hit a performance issue. After we expanded the Hadoop cluster to 13 nodes, we decided that we should stop adding more nodes, but instead look for alternative solutions. As a result, we started exploring the clustering capabilities of Mahout, conducted more PoCs, and adopted the use of clustering to handle large data sets.

**Learning curve**

There is a substantial skill gap in text analytics, Mahout and Hadoop internally, in Singapore and in the region. Actual implementations of projects of this nature are few and far between. Adoption of Hadoop is still at a nascent stage in Asia Pacific.

The PoCs enabled NLB to explore, experiment and pick up critical skills needed. The documentation, books, online resources and forums, and online courses at Coursera[9] are good sources that help the team with the implementation.

## 8 NEXT STEPS

NLB is currently extending the contextual discovery capability with the use of text analytics to all the NLB digital collections. This is expected to be completed by 2014.

At the same time, we are exploring several possibilities to bring contextual discovery to the next level:

- *Cross-institution*. With NAS coming under NLB, we now have an expanded collection of digital resources on Singapore. We will be working on recommendations across NAS and the National Library.

- *Cross-language*. There are 4 official languages in Singapore: English, Chinese, Malay and Tamil. We are considering the use of machine translation to provide cross-language recommendations.

- *Content analytics*. Describing the digital content is a very labour intensive process. There is currently a large portion of the digital content that has not been described. As a result, they remain 'unsearcheable' and text analytics cannot be applied to make them discoverable.

  We are exploring the use of content analytics software to extract information from the 'raw' content. One example is the use of voice-to-text technologies for oral history and audio visual recordings. Another possibility is to perform feature extraction on images to identify similar images.

---

[9] https://www.coursera.org/

# 9 CONCLUSION

The contextual discovery capabilities described in this paper will be a critical component of the next generation of digital libraries. This is particularly important in an era of data deluge and shorter attention span.

Memory institutions such as libraries, archives and museums play a key role in connecting the people to their country. With contextual discovery, they can now bring together the full spectrum of resources within its repositories regardless of institutions, formats and language, and present them in a way that encourages a deeper understanding of the precious memories of the nation.

It is also what knowledge seekers have been dreaming of. They can spend less time gathering the pieces, and focus on digesting and analysing the 'dossier' of relevant information to derive new insights.

The core of this rich contextual discovery experience is a massively connected network of knowledge resources. Text analytics has been shown to be a good and scalable solution to automatically identify highly accurate connections. The results so far have been very promising.

The journey has begun. Over time, with more data, better technologies and associations, the 'knowledge network' will form the bedrock of next generation knowledge management. This promotes higher level of engagement and interactivity with the knowledge network and opens up more opportunities for innovative ideas to be conceived. It will power an informed society in the new knowledge economy.

**References**

Sharmini Chellapandi, Chow Wun Han, Tay Chiew Boon, (2010) "The National Library of Singapore experience: harnessing technology to deliver content and broaden access", Interlending & Document Supply, Vol. 38 Iss: 1, pp.40 – 48

LIM, Chee Kiam and CHINNASAMY, Balakumar (2013) Connecting library content using data mining and text analytics on structured and unstructured data. Paper presented at: IFLA WLIC 2013 — Singapore — "Future Libraries: Infinite Possibilities" in Session 152 - Reference and Information Services.