

# Ontario Library Research Cloud Status Report

7 November 2014

## Evolution

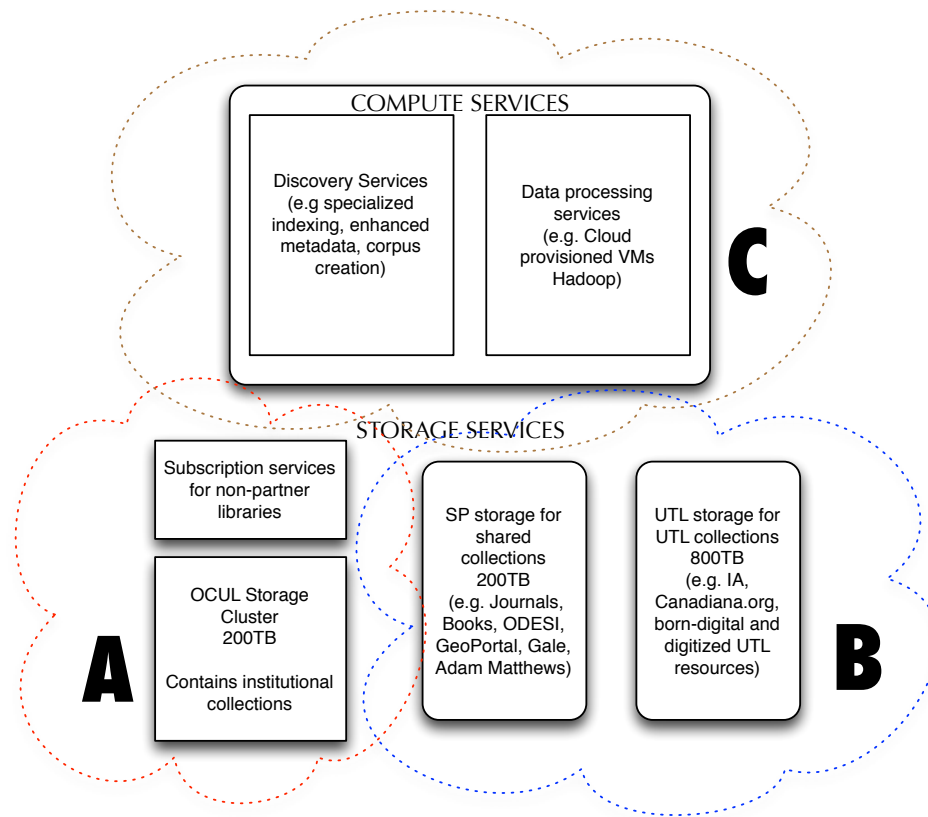
The Ontario Library Research Cloud (OLRC) began as a proposal to the OCUL directors in November 2012 for the development of an “OCUL Cloud Storage Network”. The proposal envisioned a distributed storage network of about 200TB, with nodes located at a small number of OCUL partner libraries. The network would provide enhanced storage capacity for the digital collections of those partners, of Scholars Portal, and of OCUL libraries participating in the initiative as cloud subscribers. Access to the cloud would be provided through web APIs, allowing for integration with common repository tools such as DSpace, Dataverse, Fedora and Archivematica. The OCUL Directors approved this proposal at their May 2013 meeting. The budget approved for the project was \$247,000, spread out over 18 months. The start date was set for November 2013.

During the summer of 2013, the scope of the project grew to incorporate a new element. At that time, the University of Toronto Library was developing plans for a new Canadian Text Archive Centre (CTAC). The CTAC would integrate high-capacity storage with advanced text-mining tools to allow researchers to explore digital texts “at scale”. To build the storage capacity required for CTAC, plans for the original OCUL Cloud Storage Network were extended so that, using the same cloud technologies, an 800TB storage cluster would be created for UTL alongside the original 200TB cloud envisioned for OCUL. Costs for this second cluster, estimated at about \$500,000, would be covered by UTL.

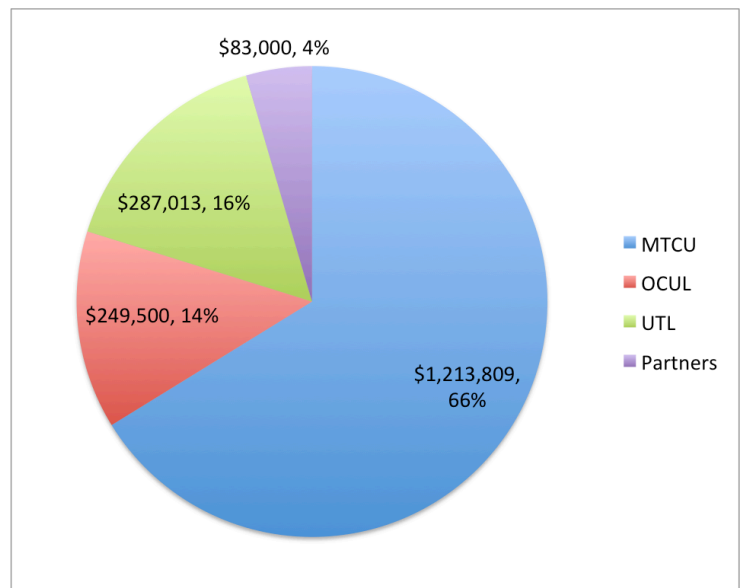
As staff began preparing for THE implementation of the combined OCUL/CTAC storage cloud, the Ministry of Training, Colleges and Universities (MTCU) announced a funding opportunity under its Productivity Innovation Fund (PIF) program. A proposal was submitted to the MTCU in September 2013 by 10 OCUL institutions, with the University of Toronto as the lead partner, to support the development of an Ontario Digital Library Research Cloud (ODLRC), subsequently renamed to the shorter Ontario Library Research Cloud (OLRC). The ODLRC proposal incorporated elements of the original OCUL Cloud project and CTAC, extending the capacity of the network to 1.2 PB and adding hardware and staffing to realize the text-mining ambitions of CTAC, making these services available to all OCUL libraries.

The architecture diagram for the OLRC, shown below, attempts to identify how elements of earlier proposals were brought together in the design of the OLRC. A common storage services layer, divided into two clusters (A and B), is shown connecting to a computational services layer (C), which provides web-based text mining support services. Just as cloud storage would be accessible through storage services APIs and tools, computational services of the OLRC would be available

through APIs so that these services could be incorporated into local delivery environments.



The extended scope of the project submitted to MTCU required extension of its duration from 18 months to two years (spread over 3 fiscal years). Likewise, the budget for the project increased to \$1.8M, with the MTCU grant accounting for \$1.2M of the revenue.



## Implementation

### RFP

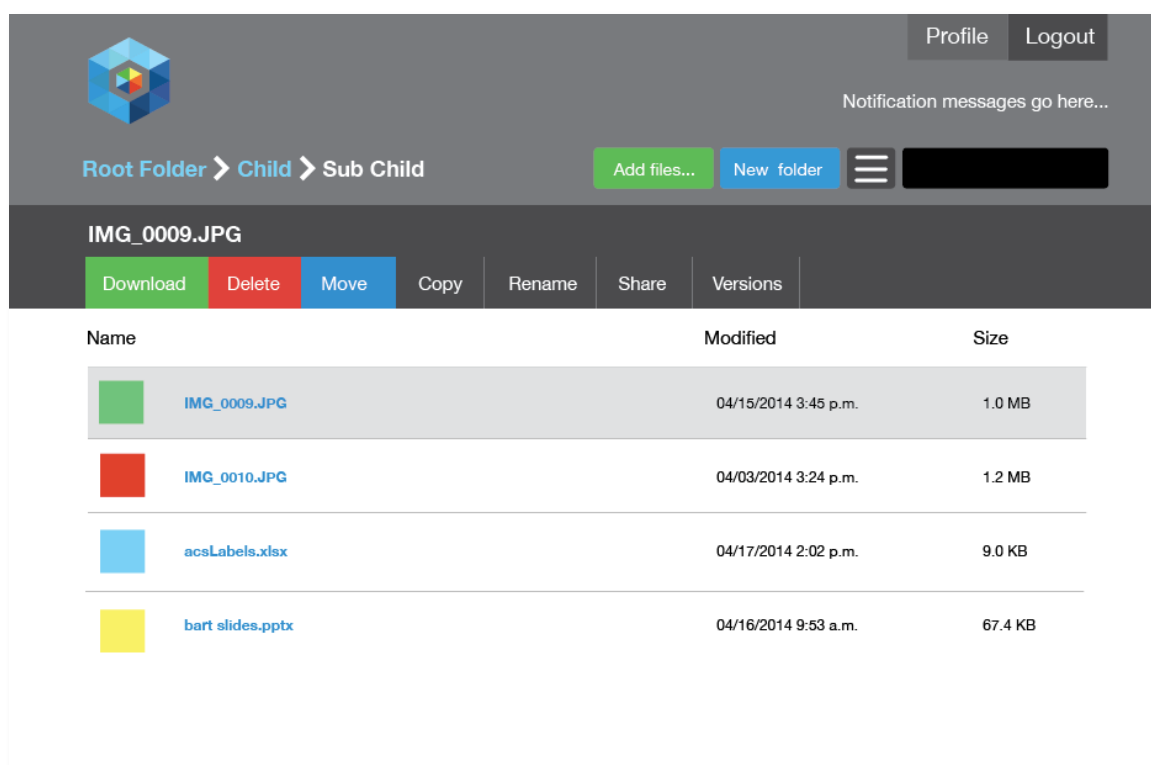
Following the award of the PIF grant late in 2013, staff worked quickly to develop a Request for Proposal (RFP) for the storage hardware. Under the terms of the grant, all MTCU funds had to be expended by the end of March 2014. Because the bulk of the MTCU funds were earmarked for hardware acquisition, it was imperative to move quickly to acquire that hardware. An RFP was issued in January 2014 with responses assessed in February and the RFP awarded in March to Dell Canada. Equipment was ordered and received by the end of March, and work on implementation of the project began in April 2014.

### Building “Part A”

Following conclusion of the RFP, the implementation team switched focus to the design and implementation of “Part A” of the project -- the OCUL storage cluster. The design phase included work with GTANet and ORION to test different network configurations and measure the bandwidth required to operate the storage network over the provincial research network. Various disaster scenarios were emulated while network traffic was being monitored. From these tests, staff developed a set of “basic requirements” for potential hosting sites, including minimum network bandwidth and local support requirements. Conversations with the partner libraries over the summer helped to identify 5 to 6 candidate sites for hosting storage nodes.

The OCUL storage cluster is being assembled now, with nodes in place at Toronto, Ottawa, and Queens. Two additional sites, yet to be identified, will complete the network, which should become operational in “beta” mode for all partner libraries by January 2015. The network should move to production mode by March 2015, with subscription services available to non-partner sites by May 2015.

Access to the storage cloud is provided through web APIs – specifically the OpenStack Swift API. This protocol is already implemented in a number of client tools that can be used to copy content to and from the cloud (e.g. Cyberduck, ownCloud, ExpandDrive). The most basic mode of access to content in the cloud is through HTTP. Users will be able to link to content through simple URLs.



### *Simple web application for working with data in the cloud*

A Hackfest held at Ryerson in June 2014 brought together developers from many of the partner libraries to work with the Swift API and explore integration with common library repository tools. Likewise, Scholars Portal staff have been working with Artefactual to add support for the Swift API to its archival management tool, Archivematica. This integration opens up new opportunities for using Archivematica to manage very large archival collections, including video and research data. A release of Archivematica that includes Swift support should be available in two or three months. Scholars Portal staff are also actively looking into integration support for Swift with Dataverse and Islandora/Fedora. This work will be ongoing, both within the project and within the open source development communities supporting these repository tools.

### **Building “Part B”**

Necessary upgrades to the power capacity in the data centre at UTL have delayed implementation of the shared UTL/SP storage cluster (“Part B” of the architecture diagram above). While staff at UTL work with university facilities management on these upgrades, design work on the cluster continues. It is expected that implementation of the UTL/SP cluster will be completed by March 2015. This will bring over 1PB of storage online and will increase the capacity of UTL and SP to store large digital collections. The migration of selected content from the Pillar SAN will start in the summer of 2015.

## Building “Part C”

The development of computational services to support text mining will be the most challenging part of the OLRC project. Text mining is different for different disciplines. It encompasses procedures such as topic modeling and clustering, entity identification and enrichment, part of speech analysis, and visualization. Nor is it clear what model of computational support will be most useful for researchers. Some researchers will be skilled programmers, or will have access to graduate students who are skilled programmers. These researchers will be primarily interested in acquiring content for analysis in their own environments, including high performance computing clusters hosted by Compute Canada. Other researchers, however, won't have such resources and will need a more guided form of computational access to the collections, using tools such as n-gram viewers to help them explore texts online.

Consultation with faculty in Ontario (and beyond) will be critical to defining the computational services offered by the OLRC. Likewise, consultation with OCUL libraries will be necessary to understand their plans for supporting researchers engaged in text mining. Will those libraries be interested in an “out of the box” portal or in APIs that will allow them to integrate text mining services with local support infrastructure?

An academic advisory group will be established to guide the work of the OLRC in this phase of the project. The goal will be to identify a few key tools that can be built within the scope of the project and then work with researchers in testing these tools in the context of current research projects. Consultations should start early in 2015 and continue through to May 2015. Development and testing of tools will continue to the end of 2015.

## Governance

With funding for the project coming largely from MTCU, the governance of the project needs to conform to the agreement signed between MTCU and the partners. According to that agreement, the lead partner, the University of Toronto, is responsible to the ministry for financial reporting for the duration of the project. The contract also requires that a memorandum of understanding between the partners should be put in place to govern the relationship between those partners. The MTCU is silent on what that MOU should cover and on what kind of governance bodies should be put in place.

Ten OCUL libraries signed on as partners to the MTCU grant. Partner directors were asked to nominate individuals from their institutions to serve on an Administrative

Steering Committee, which would be assigned the task of providing oversight of the project. Even though OCUL was not (and could not have been) a named partner in the MTCU proposal, the Executive Director for OCUL was also asked to join the Administrative Steering Committee because of the significant financial investment OCUL as a whole has made in the project and because the results of the project are meant to benefit all OCUL members.

Project implementation is the responsibility of a joint working group composed of staff from Scholars Portal and UTL/ITS. This project team reports regularly to the Administrative Steering Committee, via email and through meetings hosted by teleconference. The Administrative Steering Committee has provided input into the selection of hosting sites, has reviewed and approved a Communications Plan for the project, and is finalizing a Memorandum of Understanding, which includes identification of responsibilities of hosting parties to the network as a whole. This Committee will remain in existence for the duration of the project, after which governance of the production services will transition to established models.

Much of the project involves the implementation of core infrastructure to support Scholars Portal services. The service agreement between OCUL and the University of Toronto Library provides a framework for managing shared infrastructure. The distributed cloud, however, is a consortial infrastructure service that, while operated largely by SP staff, will be under the joint management of the partners and run for the benefit of all OCUL members. Decisions about subscription pricing, sustainability, expansion of the network to new partners – all of these will have to be addressed cooperatively within the context of governance models used to provide oversight of other OCUL/SP services. It seems likely, for instance, that OCUL-SP will play an important role in oversight of this part of the OLRC once it moves to production status.

## Future

The storage cloud will enable a wide range of services for those libraries that have enough local technical expertise to unlock its value. For example, a school with a local repository tool will be able to connect that tool to the OLRC and be able to manage very large digital objects in the cloud. Another school implementing an archival management service for cultural or corporate records will be able to connect tools such as Archivematica and ATOM to the cloud and take advantage of the replication of content across the partner storage nodes.

But there are some schools in OCUL without local technical expertise to support such services. The OLRC opens up the opportunity for Scholars Portal to offer to OCUL members a range of digital curation and management tools in a hosted environment. The Duracloud/Archivematica partnership recently announced by the DuraSpace group is an example of the kind of value-added service that Scholars Portal might provide. In the Duraspace project, digital archivists can upload digital

content for archival processing to a hosted instance of Archivematica, with fully processed archival packages stored through Duracloud in various cloud storage providers. A similar service envisioned by Scholars Portal, tentatively called PERMAFROST, has been proposed to the OCUL Directors as an opt-in service for OCUL libraries (and possibly other cultural institutions in Ontario) looking for this kind of central hosting support to be able to advance local digital preservation activities.

Likewise, the OLRC may play an important role in Project ARC, championed by the Canadian Association of Research Libraries and Research Data Canada. Project ARC proposes a national technological and services infrastructure for research data management and archiving. A stack of services running on Compute Canada infrastructure would support the ingest and archiving of research data from various sources. With appropriate interfaces in place, the combination of Dataverse, Archivematica and the OLRC could be a very effective archiving service within the larger Project ARC framework for Ontario researchers who want to deposit and share their data outputs.

The OLRC project has focused on realizing the potential of cloud computing for storage functions. But cloud computing embraces the entire scope of data center design and management. Managing compute and network resources is not only about managing hardware; it is more and more about managing software to carve out resources from hardware efficiently and flexibly. As OCUL moves forward to explore radical opportunities for collaboration, shared systems and flexible network designs running over the ORION research network will be critical for the realization of those opportunities. The OLRC is a first step along a path that will enable these new models of collaboration by utilizing the network not merely to connect pockets of IT infrastructure at the various universities but to enable new kinds of services that use the ORION network to tie IT infrastructure into a coordinated shared resource.

Appendix A: Budget Report

OLRC 3-Year Budget Submitted to MTCU in September 2013 Fiscal Year Runs from April to March	2013-14			2014-15			2015-16		
	Budget (as of March 31, 2014)	Expended (as of March 31, 2014)	Variance	Budget (as of Sept 2014)	Expended (as of Sept 2014)	Variance	Budget	Expended	Variance
<b>HARDWARE</b>									
Storage Nodes	\$610,433	\$704,022							
Server racks	\$23,970								
UPS and power management	\$82,232	\$90,713							
Network switches	\$58,908								
Proxy nodes	\$23,830								
Authentication nodes	\$23,830	\$121,435							
Test and integration cluster	\$88,219								
Data processing servers	\$214,468	\$216,347							
<b>SUBTOTAL</b>	<b>\$1,125,890</b>	<b>\$1,132,517</b>	<b>(\$6,627)</b>						
<b>IMPLEMENTATION COSTS</b>									
Cloud administration software	\$60,000	\$0		\$60,000			\$60,000		
Software integration services	\$20,000				\$22,378				
Training and travel	\$30,000	\$24,403		\$30,000					
Communication and promotion	\$7,500	\$3,589		\$7,500			\$10,000		
Scholars engagement event	\$5,000	\$0		\$5,000					
<b>SUBTOTAL</b>	<b>\$122,500</b>	<b>\$27,992</b>	<b>\$94,508</b>	<b>\$102,500</b>			<b>\$70,000</b>		
<b>STAFFING</b>									
1 Systems Support Specialist (ITS)	\$41,460	\$36,823		\$84,577	\$42,288				
1 Systems Support Specialist (SP)	\$41,460	\$41,459		\$84,577	\$42,288				
1 Programmer				\$79,386			\$80,973		
<b>SUBTOTAL</b>	<b>\$82,919</b>	<b>\$78,282</b>	<b>\$4,637</b>	<b>\$248,540</b>	<b>\$106,954</b>		<b>\$80,973</b>		
<b>TOTAL</b>	<b>\$1,331,309</b>	<b>\$1,238,791</b>	<b>\$92,518</b>	<b>\$351,040</b>	<b>\$106,954</b>		<b>\$150,973</b>		
<b>REVENUE</b>									
MTCU		\$1,213,809			\$0			\$0	
OCUL		\$50,000			\$129,500			\$70,000	
UTL		\$37,500			\$168,540			\$80,973	
Partners		\$30,000			\$53,000 (not collected yet)			\$0	
<b>TOTAL REVENUE</b>		<b>\$1,331,309</b>			<b>\$351,040</b>			<b>\$150,973</b>	
Carry Over					\$92,518				



## Appendix B: Hardware Purchased

### HARDWARE DETAILS (2013-14)

STORAGE	Unit Cost	Quantity	Extended Price
R720xd	\$13,710	19	\$269,373
MD1200	\$5,457	77	\$434,650
<b>Total</b>			<b>\$704,022</b>

RACKS/UPSs	Unit Cost	Quantity	Extended Price
Racks	\$877	18	\$16,329
UPS	\$2,536	26	\$68,197
KVM	\$2,134	2	\$4,413
Keyboard and Console	\$858	2	\$1,775
<b>Total</b>			<b>\$90,713</b>

SWITCHES and PROXY NODES	Unit Cost	Quantity	Extended Price
Switches	\$4,815	13	\$64,732
Switches	\$5,758	2	\$11,909
Proxy / Auth Nodes (Dell R620)	\$10,829	4	\$44,795
<b>Total</b>			<b>\$121,435</b>

DATA PROCESSING SERVERS	Unit Cost	Quantity	Extended Price
R720xd	\$13,710	5	\$70,888
R720xd + SSD	\$19,270	2	\$39,855
R720xd + 10K RPM	\$19,484	3	\$60,446
MD1200	\$5,459	8	\$45,158
<b>Total</b>			<b>\$216,347</b>

\$1,132,518

## Appendix C. Project Timeline

November 2012	Initial proposal for an OCUL Cloud Storage Network presented to the Directors Proposal calls for implementation of a 200TB cloud with participation of 3-4 libraries as hosting partners
May 2013	Second draft of proposal is presented and approved Project will start in November 2013 and run for 18 months Project budget is \$247,000 over 2 fiscal years; Call for partners issued
June 2013	Discussions with UTL lead to expansion of plans for the network to add 800TB to hold collections planned to be part of the Canadian Text Archive Centre This adds about \$500,000 to the project budget to be covered by UTL
September 2013	A revised proposal is prepared for the Productivity Innovation Fund of the MTCU, with Toronto as the lead institution and nine other OCUL institutions signing on as partners Ottawa; Carleton; Queens; Guelph; Laurier; McMaster; Waterloo; Windsor; York This proposal extends the project timeline to 3 years, increases the proposed capacity of the cloud to 1.2TB, and adds computational support for text mining as a goal Total budget for the project is now \$1.8M over 3 years, with \$1.2M coming from MTCU and the balance from UTL, the other partner institutions, and OCUL, including \$79,500 from the New Initiatives Fund
November 2013	An Administrative Steering Group is established to provide partner oversight of the project, composed of a director-appointed representative from each partner library and the ED from OCUL MTCU grant is awarded
December 2013	Staff at UTL and SP develop RFI for acquiring storage hardware
January 2013	RFP is issued
February 2013	RFP Responses reviewed
March 2013	Contract awarded to Dell Canada; Storage and computational hardware purchased for the project and received by the end of March
April 2014	Administrative Team meets again and begins developing Memorandum of Understanding defining governance for the project, as required by MTCU grant
May 2014	Network design and planning begins; partnership with GTANet established to assess possible capacity issues between partner nodes; pilot involves Toronto, York and Ryerson
June 2014	Administrative Team meets again Final report to MTCU prepared by lead partner (Toronto) Hackfest event held at Ryerson University; goal is to introduce programmers to the cloud APIs and explore integration opportunities
July 2014	Basic requirements for hosting sites are established – network capacity, local support, etc... Partner sites are asked to review requirements and identify local capacity Provincial research network, ORION, takes lead role in network planning and testing, working with staff at UTL and SP
September 2014	Administrative Team meets and is presented with findings of the GTANet pilot Initial set of hosting sites is identified, including Ottawa, Queens, Toronto
October 2014	Storage nodes installed at first three sites Decision pending on location of the final two storage nodes

## Appendix D: Testing Results from GTANet Pilot

Network graphs are taken from the GTANet Statistics site:

<http://www.gtanel.ca/stats-new/>

**IMPORTANT NOTE:** On the graphs, 'Inbound' and 'Outbound' are reversed from what you would expect, because these are the stats from the network, not the node. So, 'Inbound' traffic refers to traffic leaving the node and flowing **into** the GTANet router, and vice versa. Also note that there are some gaps in GTANet's stats, which results in stats only being available toward the end of Round 2. SP staff have stats gathered using iftop for this period.

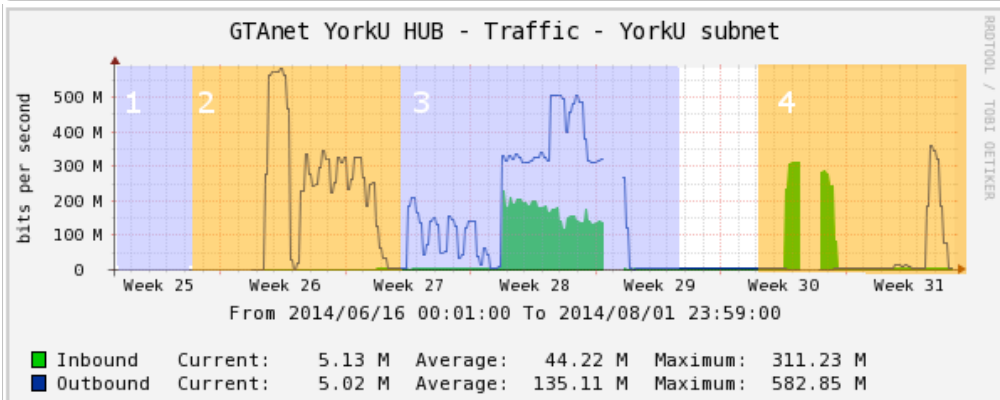
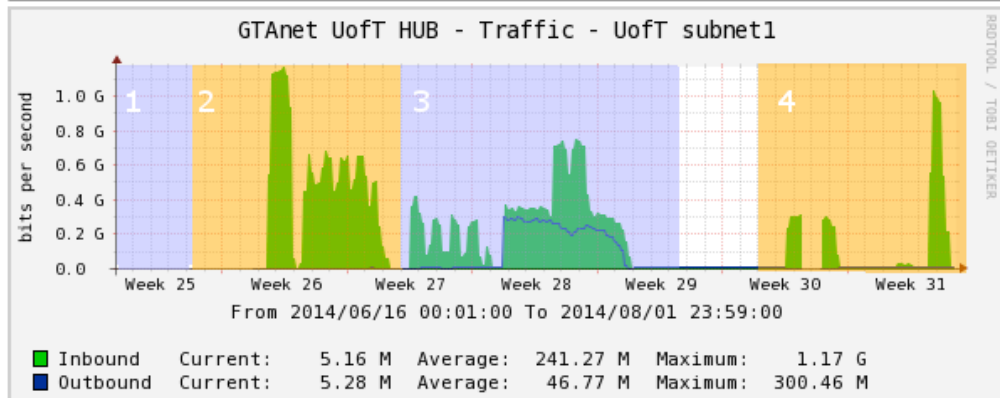
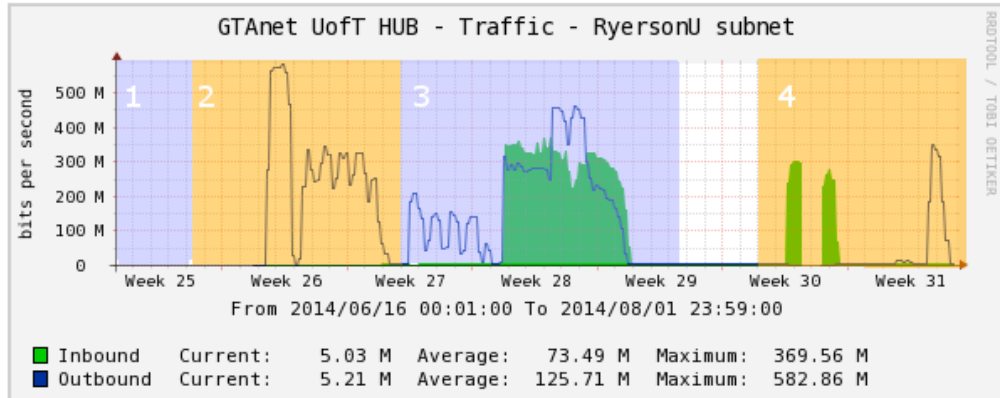
It is worth noting that these scenarios represent extreme changes to the storage network. Outside of initial setup conditions, upgrades and other planned events will occur in a much smoother fashion, which should ameliorate the need for such high utilization. That said, they may be fairly representative of possible outage scenarios, in which the network might lose a node and have it reappear.

Round 1 - Three storage nodes were set up, one each at UofT, Ryerson, and York. Additionally, a proxy node (which controls routing of traffic to the storage nodes) was set up at UofT. Once operational, baseline network traffic between the nodes was negligible (5-10Mb/s).

Round 2 - 35TB of data was loaded into the cluster. At 1Gb network speeds, replication and distribution of the data across the cluster took a couple of weeks. Network utilization for the storage nodes averaged 350 Mb/s, which is not bad, but the proxy node's network connection was fully utilized the entire time.

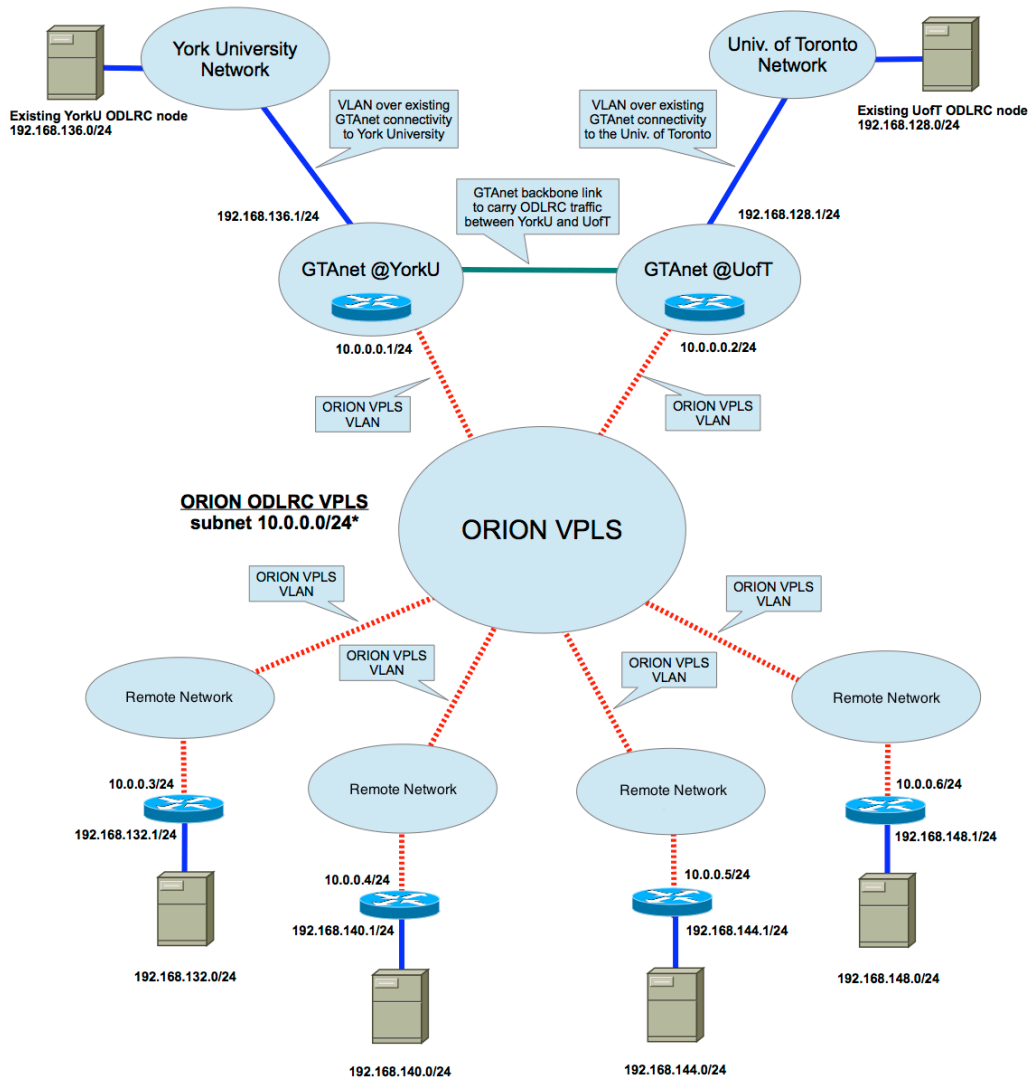
Round 3 - Once loading of the 35TB was completed, 12 drives were added to each storage node, which led to a rebalancing of the content within the cluster. Network utilization for the storage nodes hovered between 500-850 Mb/s, while the proxy node did not come under significant load. Attempts to load additional content at this time led to complete saturation of the 1 Gb network connections between the storage nodes.

Round 4 - This test added an additional storage node at the UofT site, which led to another rebalancing of the loaded content. It took approximately 2.5 days to import ~20TB of content to the new node, during which time, the 1Gb connection to the new node was fully saturated. Network utilization of the other nodes was around 350-500 Mb/s.



## Appendix E. Proposed Network Configuration for OLRC

### **The Ontario Digital Library Research Cloud (ODLRC)** Network Proposal – August 19, 2014



- ODLRC node equipment and an ODLRC router/switch will be installed at each site
- ORION will provision a VPLS network to link the remote ODLRC nodes, subnet 10.0.0.0/24 will be used over this network
- sites would pass the ORION VPLS VLAN through their networks to the ODLRC router/switch
- ODLRC node equipment would use a private /24 subnet at each site