# CGHub: Kick-starting the Worldwide Genome Web

**Christopher Wilks[1]\*, Dan Maltbie[2]\*, Mark Diekhans[1], and David Haussler[1]**

[1] University of California, Santa Cruz
1156 High Street,
Santa Cruz, CA 95064, USA

[2] Annai Systems
475 Alberto Way, Suite 130
Los Gatos, California 95032, USA

E-mail: cwilks@soe.ucsc.edu, danm@annaisystems.com
\* Authors to whom correspondence should be addressed

**Abstract:**

The University of California, Santa Cruz (UCSC) is under contract with the National Cancer Institute (NCI) to construct and operate the Cancer Genomics Hub (CGHub), a nation-scale library and user portal for cancer genomics data. This contract covers growth of the library to 5 Petabytes. The NCI programs that feed into the library currently produce about 20 terabytes of data each month. We discuss the receiver-driven file transfer mechanism Annai GeneTorrent (GT) for use with the library. Annai GT uses multiple TCP streams from multiple computers at the library site to parallelize genome downloads. We review our performance experience with the new transfer mechanism and also explain additions to the transfer protocol to support the security required in handling patient cancer genomics data.

**Keywords:** cancer; genomics; bittorrent; genes; 10Gbps; protocol

## 1. Introduction and Background

Currently various cancer genomics projects of the National Cancer Institute (NCI) [8] are producing a comprehensive characterization of the genomic changes in dozens of types of human cancer; the projects include the Cancer Genome Atlas (TCGA) [13], the Therapeutically Applicable Research to Generate Effective Treatments (TARGET) [12] a project addressing childhood cancers, and the Cancer Genome Characterization Initiative (CGCI) [14]. These three projects represent the first genomics data sets in the world that are large enough to provide researchers with the statistical power to address the true complexity of cancer. Thus, the resulting data have an incalculable value in cancer research. The methods and knowledge gained from these and subsequent cancer genomics projects could ultimately lead to a revolution in cancer care [4]. A large-scale portal dedicated to cancer genomics sequence data with efficient retrieval of these large data by software clients was urgently needed. This paper outlines the background, design, and efficient transfer protocol of the Cancer Genomics Hub (CGHub).



**Figure 1.** Worldwide web maturity progression diagram

In the last thirty years, the Internet and the World Wide Web of information has undergone a substantial maturing process as shown in Figure 1. The "slow start" when transfers were limited to relatively non-user friendly approaches such as File Transfer Protocol (FTP) but ramped up through slightly better applications such as Gopher and then WAIS, and finally reached a point of true mass usability and efficiency in the form of the ubiquitous HTTP and newer interfaces ever progressing such as Web 2.0.

The genomics "Web" is also undergoing a maturing process, albeit several years behind the information internetworks we have just described. In Figure 2 we see the analogous process. However, we're still in the early stages where storage, access, and transmission methods are being worked out.
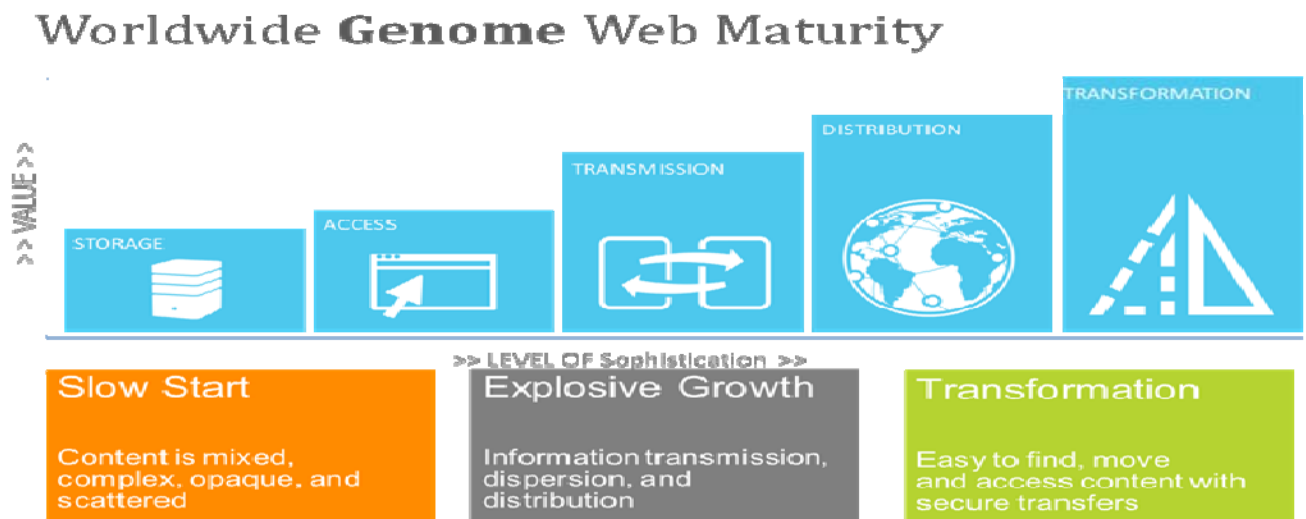


**Figure 2.** World Wide Genome web maturity progression diagram

## 2. Methods

Referring to Figure 2, CGHub encompasses the storage, access, transmission, and distribution stages of development, and as such, is a critical component of the burgeoning Worldwide Genome Web. This paper demonstrates the approaches we've taken in building CGHub, focusing on the GeneTorrent file transfer protocol and the real world results in the everyday production use of the CGHub system.

The Cancer Genomics Hub focuses on cancer as its primary data type. However genomics data is certainly not limited to cancer data or even simply disease data. CGHub can be used as a model for the construction of other non-cancer, or even non-disease genomics data portals. In Figure 3 we show CGHub in the context of the genomic sequencing data flow.
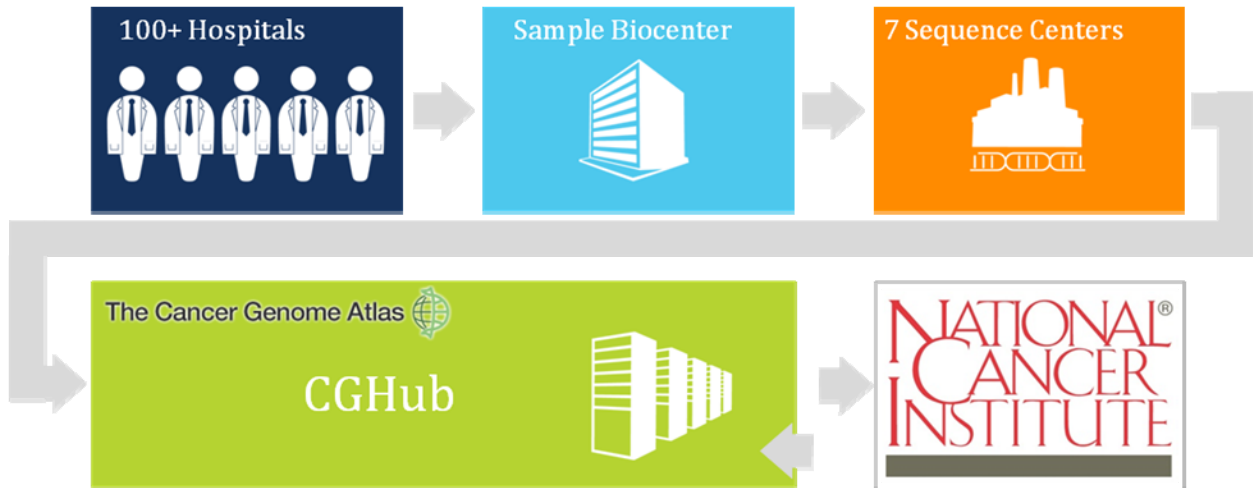
**Figure 3.** CGHub in the context of the TCGA data flow

The Cancer Genomics Hub uses a specially enhanced version of the BitTorrent (BT) [1] protocol to allow for more efficient use of high speed networks, such as Internet2, in a secure manner. This approach uses the Annai GeneTorrent (GT) protocol which is open-source and supports both download and upload to the CGHub server system (which utilizes the Annai Genome Network Operating System, or GNOS). One of the appealing elements of using BitTorrent is the fact that it is a highly vetted protocol as demonstrated by the usage statistics in Figure 4.
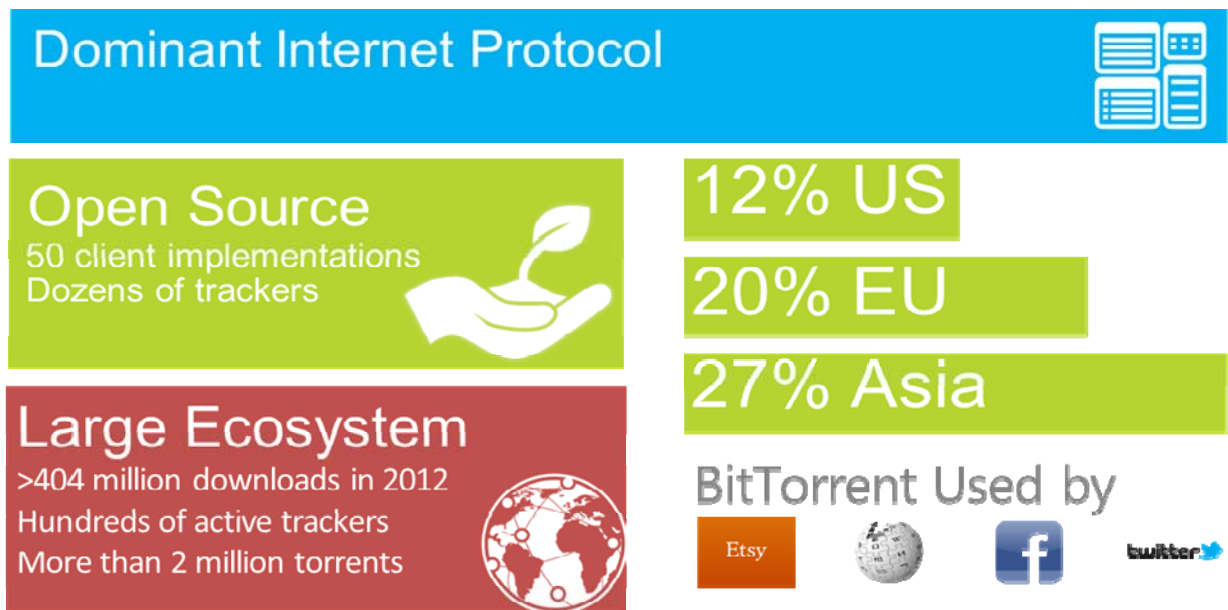


**Figure 4.** BitTorrent usage details [2][3][7][10][15][16]

The key elements of the GeneTorrent protocol for CGHub include leveraging the high degree of parallelism inherent in BT to support:

1) efficient encrypted data transfers

2) scalability of both server and client in a straight forward way

To demonstrate that 1) has been met we have tested encrypted data transfers of a single file with >3.2 Gbps rate on a 10Gpbs network. In support of 2) we have designed the system in such a way as to be simple to extend performance (up to a point) on either end of the connection. For the server, elevated aggregate traffic can be handled by simply adding additional COTS download servers to the current download server pool which is dynamically load balanced for every GT download. The servers are configured in a standard way with only the addition of the GT server software, minimizing administrative and operational overhead needed for scaling. In the case of the client, simply adding more CPUs/cores can make a difference up to a certain point if the client is not IO (disk subsystem) or network limited.

As previously mentioned this ease in scalability and high performance rates are primarily the result of the shared data serving capabilities inherent to BitTorrent and which GeneTorrent utilizes heavily. This software combined with a high performance shared filesystem such as the IBM General Purpose Filesystem (GPFS) [4] which CGHub uses as well, results in extremely parallelized data transfers, allowing for significant horizontal scaling (number of aggregate users) and the possibility of some vertical scaling as well (performance of a given download). GeneTorrent's basic parallel architecture is shown in Figure 5. And the built-out parallelism of CGHub is displayed in Figure 6.
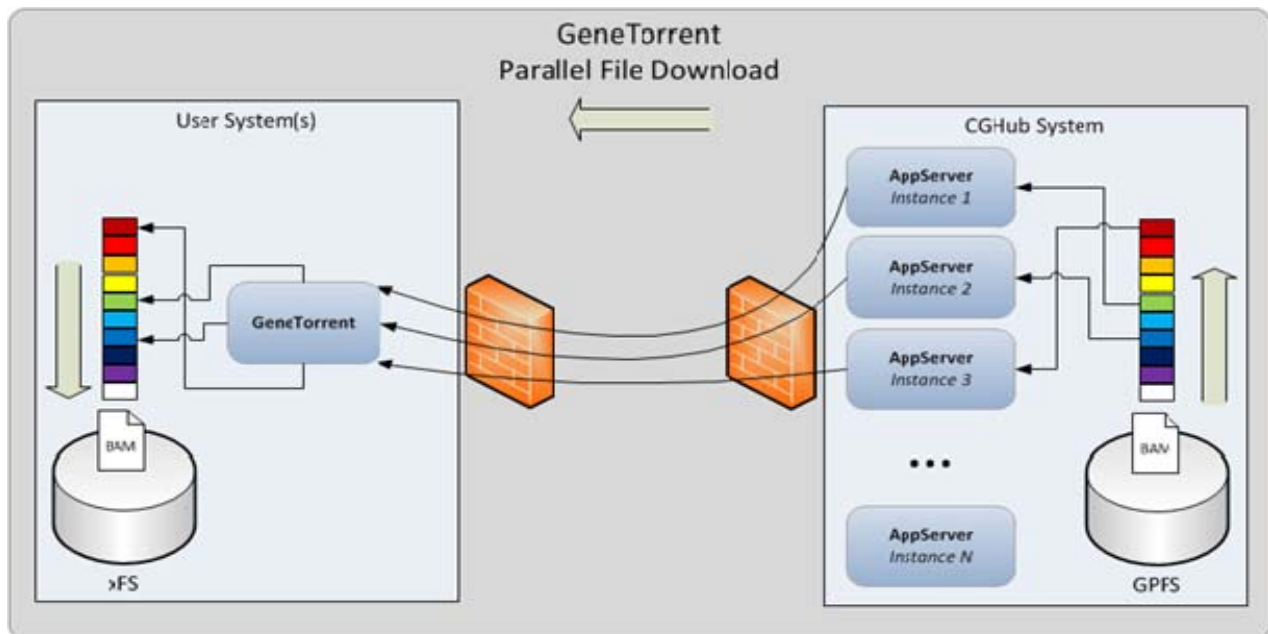


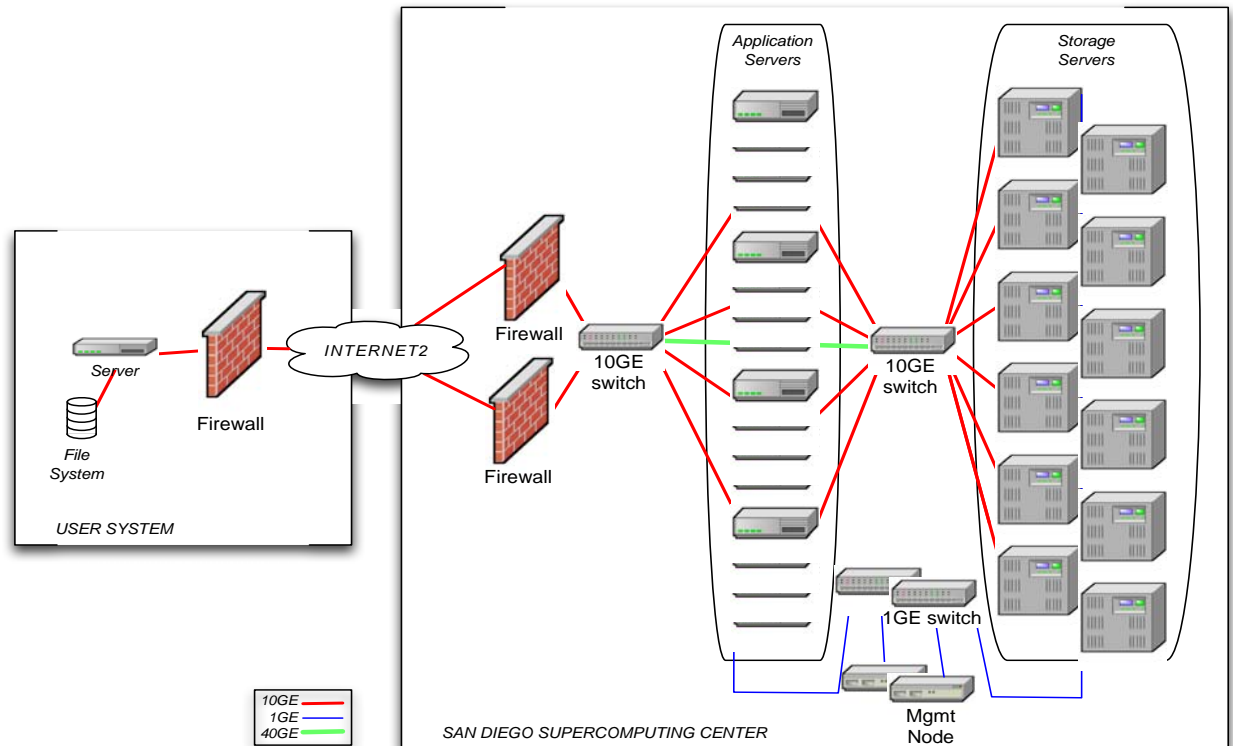**Figure 5.** Gene Torrent parallel download illustrated

**Figure 6.** CGHub system layout emphasizing multi-layered parallelism

One of the primary challenges of sharing genomics data in a disease context is the private nature of individuals' genomics data. Organizations (including both research and industry) wishing to download are screened via a data access committee (DAC) residing at the NCI/NIH, and the data itself has to be kept confidential within the authorized organizations' control and at CGHub itself. The security implemented at CGHub is operative at various layers of the system. The data itself is kept private while at CGHub through the use of firewalls and access controls limited to administrative and developer personnel only (the director of operations for instance doesn't have a login to any of the production level machines). Network access is limited by firewall, and machine access is limited by a two-factor authentication approach utilizing a VPN sign-in as the first factor, and a machine login as the second. Finally, physical access is extremely restricted via a cage within the controlled data center at the San Diego Supercomputing Center. Data transfers both to and from CGHub entail the following process:

1) User authenticates (at least once every 365 days) via an InCommon [6] federated identity provider (IdP) at the National Institute of Health, which communicates the user's successful login status via the Shibboleth middleware system [11] to CGHub's service provider (SP) server.

2) An authorization token in the form of a cryptographic string is downloaded securely by the user from CGHub at the end of the process in step 1) to a secure location on the user's download system.

3) The authorization token from step 2 is securely passed over SSL to CGHub during the initiation of a download or upload session and is used to determine the user's ability to access their requested data (based on an authorization list of users for various data stored at CGHub)

4) If the user is authorized to access the data they've requested, a secure transfer will be negotiated with CGHub data servers using SSL. Data integrity is also guaranteed for any files fully downloaded from CGHub using GT thanks to the built-in SHA1 hash checking done by the BT protocol during the transfer.

## 3. Results

To demonstrate the working nature and successful deployment of the CGHub system, we have included a recent set of basic download statistics in Figure 7[1]. Also, in Figure 8 we display a week's worth of download traffic captured at the two main CGHub firewalls, showing multi-Gbps aggregate outbound bandwidth usage.
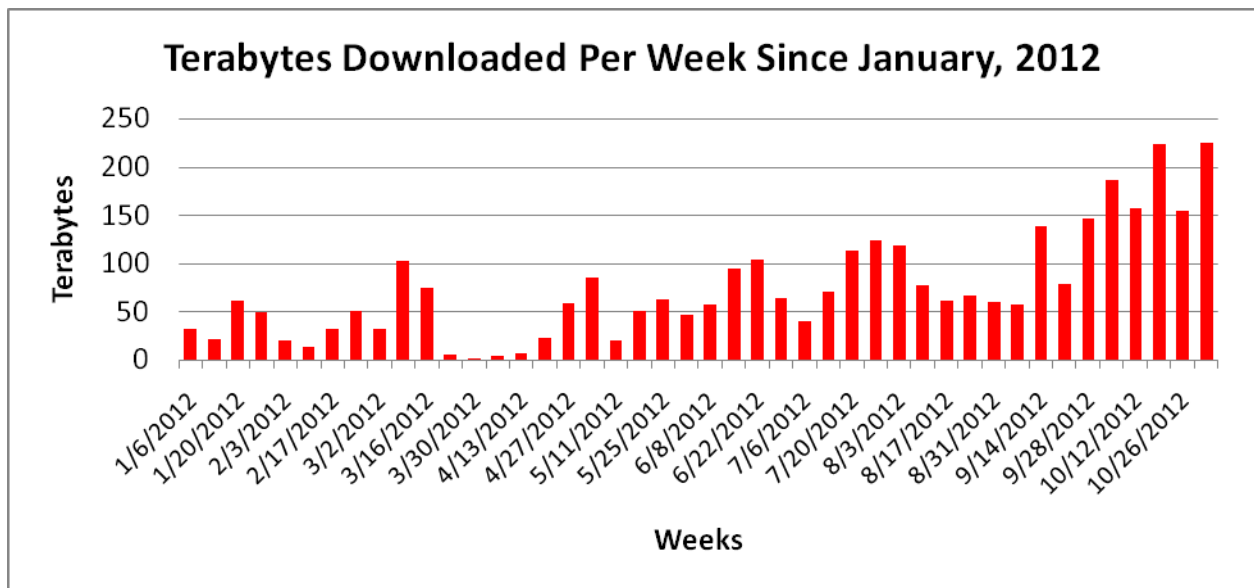


**Figure 7.** Approximate amount of data downloaded between 1/2012 and 11/2012

More than 4 petabytes in over 300,000 files have been downloaded by ~100 users during the period of a little more than a year (between 1/2012 and 1/2013).

---

[1] Download numbers are derived from initiation of a download, download completions are currently not captured, therefore numbers are approximate
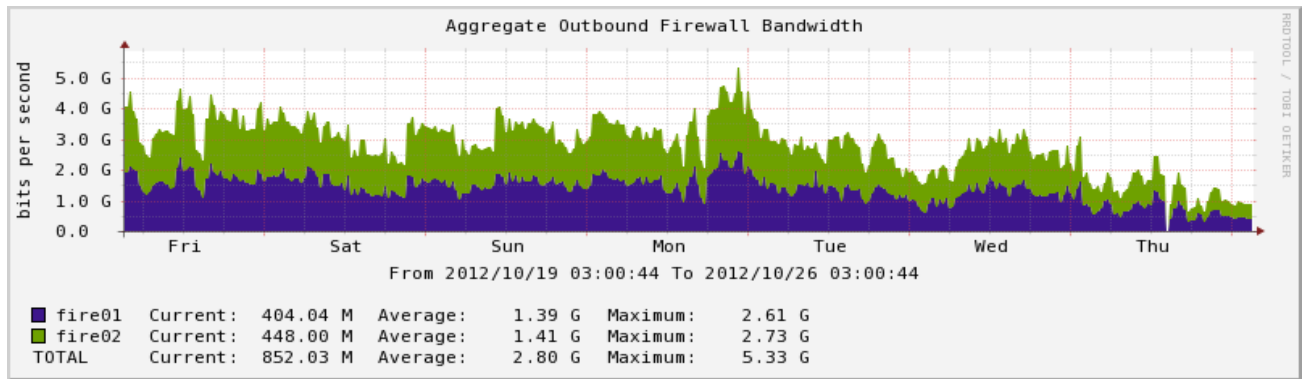
**Figure 8.** Aggregate outbound (download) bandwidth for a one week period at CGHub across the two main firewalls in the system

In addition to traffic graphs, we present a plot in Figure 9 of both the peak and average speeds we have seen when doing performance testing across a wide area network with a round trip time (RTT) of ~15 ms. These data were gathered running GT with the *–null-storage* option which discards data rather than writing to disk on the client side. The plot also shows both SSL enabled transfers and ones where encryption is disabled. The client download machine we used was equipped with 12 Intel Xeon X5690 CPU cores running at 3.47GHz with hyperthreading enabled. The client machine was also using a 10Gbps interface with a jumbo clean 10Gbps path to the CGHub servers over the Corporation for Education Network Initiatives in California (CENIC)'s network.

The file being downloaded is more than 100 Gigabytes in size. The main point of the plot is to show the impact of splitting the transfer up into more chunks (adjusting the client number of child processes dedicated to the single file download). Even when over subscribing the 24 usable threads on the client machine (with hyperthreading enabled) by running 30 child GT processes, GT is still able to increase its performance over those runs with a lower number of child processes. Thus more parallelization via increased utilization of the client's processors in this case leads to more efficient use of the high bandwidth network link for a single download.
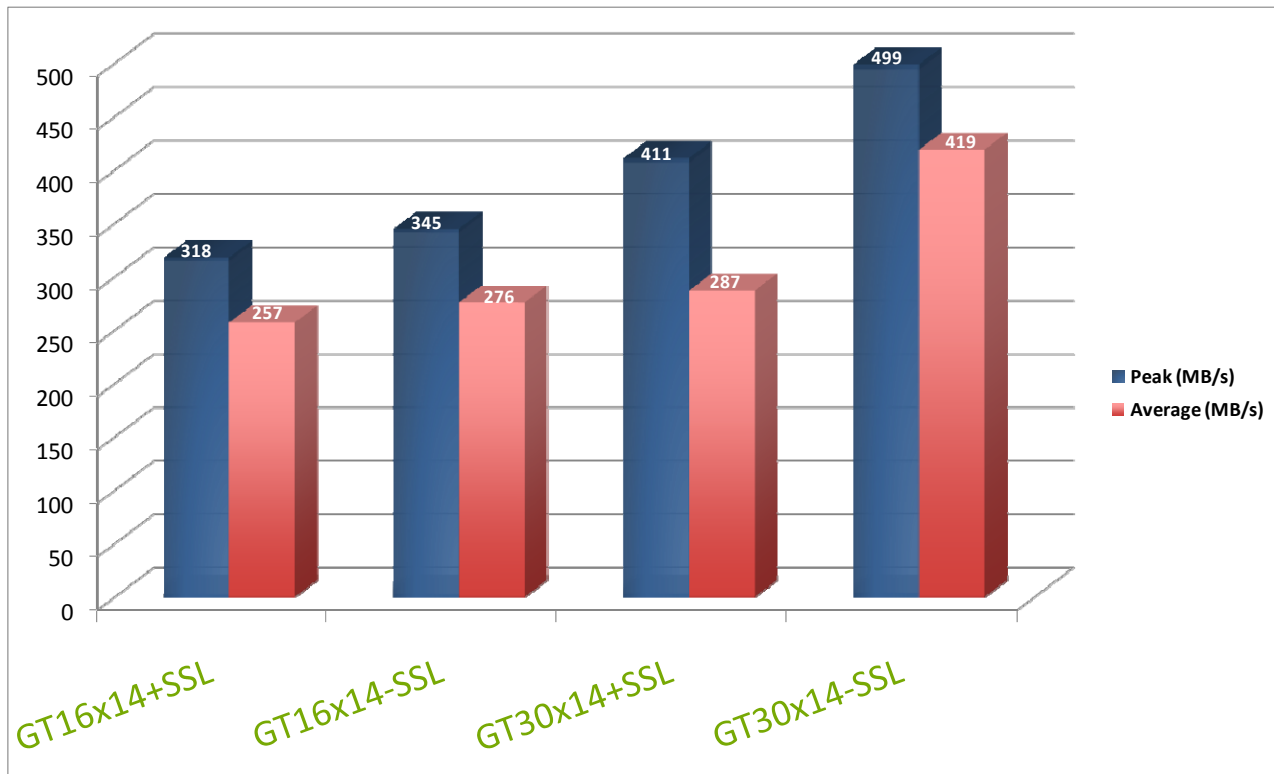
**Figure 9.** Relative performance of GeneTorrent showing both peak and average rates for four distinct download instances of the same file using different parameters on the same client machine. The first number in the x-axis labels is the total client processes transferring data in parallel (e.g. 16); the second number is the number of download machines at CGHub serving this transfer (14); the +SSL or –SSL denotes whether encryption was turned on or off for that download. All downloads were run with the GT *–null-storage* command line option so that no data was written to the client IO subsystem. All data was read from the GPFS filesystem on the CGHub servers with caching enabled.

## 4. Conclusion

There are significant efforts underway at the NIH and some of its various institutes including the NCI to understand the genetic basis of cancer through the sequencing of massive numbers of patients' tumors. This research effort can only be fully realized through the use of a resource that supports an efficient transfer protocol while maintaining the necessary security that surrounds patient data. We have described the Cancer Genomics Hub utilizing the Annai GeneTorrent protocol as a solution to this need. Further, we have demonstrated its use over the last year to download more than 4 petabytes of cancer genomics data from the TCGA project queried by ~100 users. Ongoing efforts include the expansion of the repository to encompass the additional hundreds of terabytes of TARGET and CGCI cancer genomics data. One impact of

this expansion is that it will extend CGHub to incorporate data specifically aimed at addressing childhood cancer (TARGET).  Additionally, we presented the relative performance of GT to demonstrate its ability to scale on the client to reach greater transfer rates.  Based on this evidence, it is clear that CGHub is achieving its intended purpose of bringing cancer genomics data within easier reach of the researchers at their home institutions, while mitigating the problems of big data management, not the least of which is large data transfer over wide area networks such as Internet2.

**Acknowledgments**

The authors wish to thank the following:

**References and Notes**

1. BitTorrent (BT). http://www.bittorrent.com/ (accessed on 29-01-2013)

2. Chao Zhang, Prithula Dhungel, Di Wu, and Keith W. Ross. Unraveling the BitTorrent Ecosystem. IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, VOL. 22, NO. 7, JULY 2011.

3. Code as Craft. Turbocharging Solr Index Replication with BitTorrent. http://codeascraft.etsy.com/2012/01/23/solr-bittorrent-index-replication/ (accessed on 31-01-2013)

4. Hudson TJ, Anderson W, Artez A, Barker AD, Bell C, Bernabe RR, Bhan MK, Calvo F, Eerola I, Gerhard DS, et al. 2010. International network of cancer genome projects. Nature 464: 993–998.

5. IBM General Parallel File System (GPFS), 2012. http://www-03.ibm.com/systems/software/gpfs/ (accessed on 29-01-2013)

6. InCommon. http://www.incommon.org/ (accessed on 29-01-2013)

7. Musicmetric Digital Music Index. http://www.musicmetric.com/dmi/ (accessed on 30-01-2013)

8. National Cancer Institute (NCI). http://www.cancer.gov/ (accessed on 29-01-2013)

9. National Institutes of Health (NIH). http://www.nih.gov/ (accessed on 29-01-2013)

10. Sandvine Global Internet Phenomena Report: 1H 2012. http://www.sandvine.com/downloads/documents/Phenomena_1H_2012/Sandvine_Global_Internet_Phenomena_Report_1H_2012.pdf (accessed on 30-01-2013)

11. Shibboleth. http://shibboleth.net/ (accessed on 29-01-2013)

12. Therapeutically Applicable Research to Generate Effective Treatments (TARGET). http://target.cancer.gov/ (accessed on 29-01-2013)

13. The Cancer Genome Atlas (TCGA). http://cancergenome.nih.gov (accessed on 29-01-2013)

14. The Cancer Genome Characterization Initiative (CGCI). http://cgap.nci.nih.gov/cgci.html. (accessed on 29-01-2013)

15. TorrentFreak. Facebook Uses BitTorrent, and They Love It. http://torrentfreak.com/facebook-uses-bittorrent-and-they-love-it-100625/ (accessed on 31-01-2013)

16.       Wikimedia Foundation. Video Labs: P2P Next Community CDN for Video Distribution. http://blog.wikimedia.org/2010/09/27/video-labs-p2p-next-community-cdn-for-video-distribution/ (accessed on 31-01-2013)