# University of California Santa Cruz

# **CGHub User Guide**

January 11, 2011 Version 1.2

1	INTRODUCTION	3
	1.1 TECHNICAL SUPPORT	
	1.2 References	3
	1.3 DOCUMENT SCOPE	
2	SYSTEM OVERVIEW	4
3	SOFTWARE INSTALLATION	7
	3.1 CGQUERY SCRIPT	7
	3.1.1 Prerequistes	
	3.1.2 Download cgquery Software Package	
	3.1.3 Install cgquery Script	
	3.2 GENETORRENT CLIENT PROGRAM	8
	3.2.1 Prerequisites	8
	3.2.2 Download GeneTorrent Software Package	8
	3.2.3 Install GeneTorrent Software	
	3.2.4 Check Installation	9
4	USER AUTHENTICATION	10
	4.1 AUTHENTICATION	
5	SEQUENCE METADATA QUERIES	12
_	5.1 Query Attributes	
	5.2 USING THE CGQUERY SCRIPT	
	5.3 EXAMPLE CGQUERY COMMANDS	
	5.4 CGHUB REST API	21
	5.5 SAMPLE REST CALLS	22
6	SEQUENCE DOWNLOAD	26
	6.1 DOWNLOADING FILES USING GENETORRENT CLIENT	26
	6.2 QUERY AND DOWNLOAD	27
	6.3 OHERY AND DOWNLOAD USING COOLERY INTERACTIVE MODE	28

#### 1 Introduction

CGHub is a genome sequence storage repository that provides large capacity and fast access to sequence data generated as part of NCI research projects. The repository is accessed over the Internet, providing simple but secure sequence data search and transfer services. The CGHub currently contains sequence data for The Cancer Genome Atlas (TCGA) cancer research study sponsored by the NCI and NHGRI. NCI plans to add more studies (such as TARGET and CGCI) over time.

For the TCGA study, sequence data is considered personally identifying information (PII) because it comes from human clinical samples. In order to support research use guidelines and enforce participant privacy rights, only NIH authorized users of TCGA can access the genome sequence data via protected CGHub Web Services. The metadata that describes the sequence and related analysis, captured as XML-based metadata descriptors, is publicly available via the CGHub Web Services.

# 1.1 Technical Support

For technical or operational requests, contact: support@cghub.ucsc.edu

#### 1.2 References

CGHub Home Page: https://cghub.ucsc.edu/

TCGA Project Home Page at NCI: http://cancergenome.nih.gov/

TCGA Data Guide: <a href="https://wiki.nci.nih.gov/display/TCGA/TCGA+Data+Primer">https://wiki.nci.nih.gov/display/TCGA/TCGA+Data+Primer</a>

TCGA Code Tables:http://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm

dbGaP TCGA Study Page: http://www.ncbi.nlm.nih.gov/gap?term=phs000178

#### 1.3 Document Scope

This document describes how to search and download primary genome sequence data and related metadata from the CGHub repository.

Downloading sequence data from CGHub involves three user steps:

- Use your browser to access your eRA account and generate a CGHub authentication credential
- 2. Query the metadata to get the download URLs of the desired analysis objects that contain the sequence data of interest
- 3. Invoke the GeneTorrent Client with the selected URLs to download

These steps will be described in detail in the subsequent sections.

Once a generated authentication credential has been obtained it can be re-used for an unlimited number of download operations until the credential expires. The credential will need to be periodically renewed; a new credential can be obtained at any time. Keep this credential in a safe and secure place on the computer used for downloading.

Note that the authentication credential is not required to query the CGHub metadata; the authentication credential is required only for download. Query functionality is available to the public.

Queries can be performed with either the *cgquery* script or the CGHub Web Services API.

- The cgquery script provides a simple but powerful way to execute queries and file download. It is recommended that new users start with the cgquery script before attempting to use the CGHub Web Services API.
- The CGHub Web Services API returns an XML result set. This method is useful when scripting or embedding queries into your site-specific tools.

# 2 System Overview

CGHub serves as the central repository for sequencing centers to store and for analysis centers to retrieve genome sequence data and its associated metadata.

There are four main groups that work together to support the CGHub repository: Tissue Source Sites, Biospecimen Core Resource Centers, Genome Sequencing Resource Centers, and Data Coordinating Centers

A **Tissue Source Site (TSS)** collects samples (tissue, cell or blood) and clinical metadata, which are then sent to a <u>BCR</u>. A TSS is identified by its TSS ID.

A **Biospecimen Core Resource (BCR)** is a TCGA center where samples are carefully catalogued, processed, quality-checked and stored along with participant clinical information.

The work of the BCR includes the following important functions:

- Serving as the interface between the TCGA program and the different <u>Tissue Source Sites</u> that are collecting tumor and matched normal controls
- Ensuring and verifying that TCGA human subjects protections and guidelines are adhered to and that all regulations are followed at each Tissue Source Site
- Examining of all biospecimens to ensure they meet rigorous standards for each tumor type (including percent necrosis and percent tumor nuclei)

- Reviewing of pathology to ensure accurate diagnosis and inclusion in TCGA
- Collecting clinical information for each sample and applying standardized terminology, definitions and formats that are caBIG compliant
- Extracting and distributing DNA and RNA from samples to each of the genomic characterization and sequencing centers

A **Genome Sequencing Center (GSC)** is a TCGA center that uses high-throughput methods to sequence tumor and normal samples provided by the BCR.

The **Data Coordinating Center (DCC)** is the central provider of TCGA data. The DCC standardizes data formats and validates submitted data.

The work of the DCC includes the following important functions:

- Protecting participant privacy and confidentiality through secure access to research and clinical information that are classified as controlled access datasets
- Developing data standards and controlled vocabularies
- Establishing informatics pipelines for dataflow from production centers to a central repository
- Developing new analytical and visualization technologies for different audiences to facilitate data analysis
- Coordinating project level activities.

The following describes the basic system workflow for creating and distributing genome sequencing data. This workflow is provided here to orient the reader to the end-to-end system operation and is for informational purposes only. More details can be found in the CGHub Submission Guide.

- A Tissue Source Site (TSS) submits a participant tissue sample to a Biospecimen Resource Center (BCR) where the sample is stored and prepared for sequencing.
  - 1. The analyte is shipped to a Genome Sequencing Center (GSC) with a unique identifier (UUID) for the sample aliquot.
  - Metadata describing the biospecimen is sent to TCGA Data Coordination Center (DCC), who coordinates project level activities.
- GSC receives and processes the analyte from the BCR
  - 1. The analyte is sequenced by a high throughput sequencer to produce raw read coverage

- 2. Additional processing, such as alignment, mapping, and expression counts are performed and captured in a BAM file.
- GSC submits the sequence information to CGHub
  - 1. A unique identifier (AnalysisUUID) for sequence data and metadata information for the single sample analysis is created by the GSC.
  - 2. This set of data (metadata and sequence data) is referred to as the Analysis Object.
  - 3. GSC performs validation of the sequence data and metadata prior to upload to ensure proper formatting.
  - 4. Once validated, the GSC will upload the Analysis Object to CGHub and initiate the workflow on CGHub to make the data available to downloaders.
- CGHub synchronizes with the TCGA DCC to determine status of sample.
  - 1. The DCC maintains the higher-level study attributes as well as clinical data associated with each sample.
  - 2. CGHub will query the DCC to verify that submitted data is associated with a valid sample that should be available to authorized researchers.
  - 3. CGHub will also query for redacted samples and suppress the sequence data from download.
- Researchers retrieve the CGHub catalog and perform queries across the public sequence metadata using selected metadata attributes such as cancer type, sequence type, source sequencing center, or date range.
- Researchers retrieve protected sequence data using their NIH authorized credentials, as long as the sequence data is in the "live" state.

#### 3 Software Installation

# 3.1 cgquery Script

The *cgquery* is a separate package from GeneTorrent and therefore must be downloaded and installed separately. Note: to permanently install this software you must be a System Administrator.

# 3.1.1 Prerequistes

Python 2.6.5 or later with SSL support.

# Check for Python Version

```
$ python -c 'import platform; print platform.python_version()'
2.6.5
```

# "2.6.5" or greater is required

# Check if Python supports SSL

Use the following command to check if the python installation supports SSL.

```
$ python -c 'import socket; print hasattr(socket, "ssl")'
True
```

# 3.1.2 Download cgquery Software Package

#### From a browser:

- go to GeneTorrent Software Downloads web page (<a href="https://cghub.ucsc.edu/downloads.html">https://cghub.ucsc.edu/downloads.html</a>)
- click on "cgquery-1.9.tar.gz" or latest version

#### From command line terminal window:

**Command:** wget http://cghub.ucsc.edu/cgquery-1.9.tar.gz

<sup>&</sup>quot;True" is required

# 3.1.3 Install cgquery Script

Install *cgquery* with the following command:

Command: sudo tar -C / -xzvf cgquery-1.9.tar.gz

For technical or operational requests, contact: <a href="mailto:support@cghub.ucsc.edu">support@cghub.ucsc.edu</a>

# 3.2 GeneTorrent Client Program

Note: to permanently install this software you must be a System Administrator.

Binary (ready-to-run) versions of the client are available for 64 bit versions of the following LINUX Systems:

- Centos6
- Centos 5.5
- SUSE11.4
- FedoraCore15
- Ubuntu11.4
- Ubuntu10.4

In addition, the source code for GeneTorrent can be found at **GeneTorrent-1.0.0.0-src.tgz** (https://cghub.ucsc.edu/GeneTorrent-1.0.0.0-src.tgz).

#### 3.2.1 Prerequisites

The GeneTorrent software requires certain software packages to be installed to function. Curl and OpenSSL are required.

# 3.2.2 Download GeneTorrent Software Package

To find the GeneTorrent software installation package, go to GeneTorrent Software Downloads (https://cghub.ucsc.edu/downloads.html). Select a download according to your operating system. This will download a file named: GeneTorrent-1.0.0.0-YOUR\_OS.x86\_64.tar.gz (where YOUR\_OS is the name of the operating system that you selected).

#### 3.2.3 Install GeneTorrent Software

Install GeneTorrent with the following command:

```
Command:sudo tar -C / -xzvf GeneTorrent-1.0.0.0.$YOUR OS.x86 64.tar.gz
```

Where \$YOUR\_OS is one of the operating systems listed in the OPERATING SYSTEMS SUPPORTED section.

```
]$ sudo tar -C / -xzvf GeneTorrent-1.0.0.0.Centos5.5.x86_64.tar.gz
[sudo] password for user:
usr/
usr/bin/
usr/bin/GTLoadBalancer
usr/bin/GeneTorrent
usr/bin/gtoinfo
usr/share/
usr/share/man/
usr/share/man/man1/
usr/share/man/man1/GeneTorrent.1
usr/share/GeneTorrent/
usr/share/GeneTorrent/
usr/share/GeneTorrent/GeneTorrent.openssl.conf
usr/share/GeneTorrent/dhparam.pem
```

In addition, there is an RPM package of GeneTorrent for Centos 6. This can installed using the following commands:

```
Command: gunzip GeneTorrent-1.0.0.0-1.x86_64.rpm.gz

Command: sudo yum --nogpg install GeneTorrent-1.0.0.0-

1.el6.CP.x86 64.rpm
```

Note: Remember to answer 'y' to the yum prompt "Is this ok [y/N]:"

#### 3.2.4 Check Installation

To verify that the installation succeeded and that the correct version of the software is installed type the following command.

```
[admin@centos6]$ GeneTorrent --version
GeneTorrent: Genomic Data Transfer Tool version 1.0.0.0
```

For technical or operational requests, contact: support@cghub.ucsc.edu

#### 4 User Authentication

In order to download CGHub data, users must have a valid NIH or eRA account and must have received authorization from the TCGA project Data Access Committee (DAC). Please refer to the dbGaP study page for more information on applying for access to controlled-access data:

http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\_id=phs000178.v5.p5.

#### 4.1 Authentication

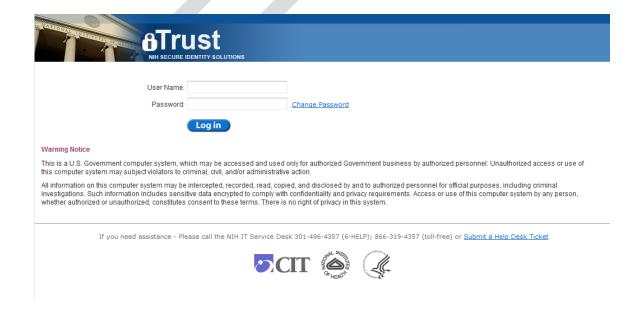
Steps 1-4 are only required to generate a new key. Download commands can be run repeatedly using the same key until it expires.

**Step 1:** Access URL for CGHub Authentication page.

The user begins to log in at the URL: https://cghub.ucsc.edu/secure.

**Step 2:** The CGHub system redirects the user to the NIH login page show below.

The user provides valid login credentials (Username and Password).



**Step 3:** NIH authenticates the user.

**Step 4:** CGHub generates and returns an authentication credential.

The user will be prompted to save the authentication credential in a local file.

This credential is valid for 365 days. You will need to repeat this process to generate a new credential whenever the credential expires or the user's account changes.

Pass the credential to the GeneTorrent download client as part of the sequence data download step.



# 5 Sequence Metadata Queries

Each sequence in the CGHub is described by metadata, which contains attributes such as the date, sample type, and disease abbreviation. (see section 5.1 for the list of searchable attributes).

Querying the sequence metadata is a powerful method of identifying sequences of interest.

cgquery is a command line tool that allows the user to construct and submit a querystring to search for sequences with one more sequence attributes. The results are formatted in either summary or complete view. (see section 5.2 cgquery)

cgquery uses public CGHub REST APIs. Users can call this interface directly through a browser or other tool such as wget or curl and develop custom GUIs or command line tools. (see section 5.4 for more detail)

# 5.1 Query Attributes

The querystring can include one or more of the attributes listed below using name=value pairs. Multiple queries can be combined using the '&' symbol. For example: "disease\_type=OV&sample\_type=01". Both attribute names and values are case sensitive. Values can use the "\*" wildcard. For example, "disease\_type=O\*" will match any values beginning with "O".

CGHub also stores the original submission XML documents, which users can perform free form text search against any strings within the analysis.xml, run.xml and experiment.xml files. xml\_text queries are case insensitive.

Name	Description	Example	Metadata Source
analysis_id	Globally unique ID (uuid) associated with a collection of sequencing data files and their metadata.  This is the primary identifier for CGHub data. It is equivalent to an SRA SRZxxxx accession #.	analysis_id=d29389cd -72F8-48B1-bd33- 7f525170c7b8	Generated by GSC prior to submission. Legacy data was assigned uuids when it was migrated from SRA into CGHub.
analysis_acc ession	Legacy SRZnnnnnn analysis accession # for data migrated from SRA.	analysis_accession=S RZ123456	analysis.xml <analysis accession=&gt;</analysis 
state	Determines if the submission is still being processed, is live or has been suppressed. Users can only download data with state=live.  Valid values: submitted,	state=live	Managed by CGHub.

	uploading, live, suppressed,		
	verifying_sample,		
	verifying_files, bad_data		
Study	The study for which the data	study=phs000178	analysis.xml
	was generated. This value is	, , , , , , , , ,	<study_ref< th=""></study_ref<>
	used to determine which		refname=>
	users can download the data.		remame=>
	users can download the data.		
	Valid valvas ab 200470		
	Valid values: phs000178		
aliquot_id	Globally unique ID for the	aliquot_id= a8C71dbb-	analysis.xml
	aliquot used by to generate	a809-4415-bde8-	<target refname=""></target>
	this analysis.	f24f70557077	
	This is a primary identifier for		
	each project's sample		
	database (e.g. TCGA's DCC		
	Portal).		
	For TCGA, this will be a		
	RFC4122 aliquot uuid		
	managed by the DCC.		A sal size at
sample_acce	Legacy sample accession #	sample_accession=S	Analysis.xml
ssion	for data migrated from SRA.	RS123456	<target< th=""></target<>
			accession=>
legacy_samp	This may be any project	legacy_sample_id=	TCGA DCC uuid to
le_id	specific ID.	TCGA-AB-2802-03C-	barcode mapping Web
		01W-0731-08	Service
	For TCGA, this will be the 7		
	part barcode retrieved from		
	the DCC.		
disease_abbr	The short disease name.	disease_abbr=OV	TCGA DCC uuidws
	The distributed harries	a.coa.co <u>_</u> a.co.	Web Service
	The set of legal TCGA values		<disease></disease>
	are defined by the DCC		<abbreviation></abbreviation>
	codes table at:		<abbic viation=""></abbic>
	http://tcga-		
	data.nci.nih.gov/datareports/c		
	odeTablesReport.htm		TOO 4 DOC
participant_i	The ID of the human	participant_id=	TCGA DCC uuidws
d	participant providing the	a85391a4-c6bc-4cd7-	Web Service
	sample. Previously	b53b-6e7cebf044b5	<participant href=""></participant>
	represented by the TCGA 3	5555-067 065104455	
	part barcode (e.g. TCGA-AB-		
	2802 <b>)</b> .		
	For TCGA this will be a		
	participant UUID provided by		
	the DCC.		
sample_id	The ID for the root sample		TCGA DCC uuidws
Campic_ia	from which the aliquot was	sample_id=	Web Service
		972a7b73-9323-4b02-	
	derived. Previously	90dd-bb420a18681f	<sample href=""></sample>
	represented by the TCGA 5		
	part barcode (e.g. TCGA-AB-		
	2802-03C-01W <b>)</b> .		

	For TCGA this will be a sample UUID provided by the DCC.		
analyte_code	The type of analyte (e.g. D=DNA, R=RNA, etc.).	analyte_code=D	TCGA DCC uuidws Web Service <analytetype></analytetype>
	The set of legal TCGA values are defined by the DCC at: http://tcga-		
	data.nci.nih.gov/datareports/c odeTablesReport.htm		
sample_type	The type of sample (e.g. Blood Derived Normal, Primary solid Tumor, etc.).	sample_type=01	TCGA DCC uuidws Web Service <sampletype></sampletype>
	The set of legal TCGA values are defined by the DCC codes table at : <a href="http://tcga-data.nci.nih.gov/datareports/c">http://tcga-data.nci.nih.gov/datareports/c</a>		
400 td	odeTablesReport.htm	to 31 04	TOOA DOO weiden
tss_id	The Tissue Source Site	tss_id=01	TCGA DCC uuidws Web Service
	The set of legal TCGA values are defined by the DCC codes table at :		<tissuesourcesite></tissuesourcesite>
	http://tcga- data.nci.nih.gov/datareports/c odeTablesReport.htm		
last_modified	Date the object was last modified.	1) "in the last month": last_modified=[NOW- 1MONTH TO NOW] 2) "all September 2011 records": last_modified=[2011- 08-31T23:59:59.99Z TO 2011-09- 30T23:59:59.99Z]	Created by CGHub
center_name	Short names defined by the project (e.g. BI, BCM).	center_name=BI	analysis.xml <analysis center_name=&gt;</analysis 
	The set of legal TCGA values are defined by the DCC codes table at : <a href="http://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm">http://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm</a>		
alias	Name of the analysis object as defined by the submitting center.	alias= NA12878h.HiSeq.WG S.bwa.cleaned.recal.b am.header	analysis.xml <analysis alias=""></analysis>
Title	Title of the analysis object as defined by the submitting center.	title= Reference alignment of M01005	analysis.xml <title>&lt;/th&gt;&lt;/tr&gt;&lt;tr&gt;&lt;th&gt;analysis_typ&lt;/th&gt;&lt;th&gt;The type of analysis. Uses&lt;/th&gt;&lt;th&gt;analysis_type=REFER&lt;/th&gt;&lt;th&gt;analysis.xml&lt;/th&gt;&lt;/tr&gt;&lt;tr&gt;&lt;th&gt;&lt;/th&gt;&lt;th&gt;&lt;/th&gt;&lt;th&gt;&lt;/th&gt;&lt;th&gt;&lt;/th&gt;&lt;/tr&gt;&lt;/tbody&gt;&lt;/table&gt;</title>

е	the values from SRA 1.3 schema. For TCGA, this will be REFERENCE_ALIGNMENT	ENCE_ALIGNMENT	<analysis_type></analysis_type>
library_strate gy	The sequencing technique used. Values are defined by the SRA 1.3 schema:  WGS WXS RNA-Seq WCS CLONE POOLCLONE AMPLICON CLONEEND FINISHING ChIP-Seq MNase-Seq DNase- Hypersensitivity Bisulfite-Seq EST FL-cDNA CTS MRE-Seq MeDIP-Seq MDB-Seq OTHER	library_startegy=WGS	experiment.xml <library_strateg Y&gt;</library_strateg 
platform	The machine used to generate the sequence data. Values are defined by the SRA 1.3 schema:  LS454 ILLUMINA HELICOS ABI_SOLID COMPLETE_GENO MICS PACBIO_SMRT ION_TORRENT	platform=ILLUMINA	experiment.xml <platform></platform>
filename	Name of a data file associated with an analysis.	filename=HG00099*	analysis.xml and/or experiment.xml <files></files>
xml_text	Free form search from the original submission XML documents.	xml_text= TCGA-AB- 2802-03C-01W	analysis.xml, run.xml and experiment.xml

# 5.2 Using the *cgquery* Script

The *cgquery* script calls the Web services APIs and parses the XML output in a human readable form. It takes the same query string arguments as the direct Web Services API.

```
Usage: cgquery [options] <querystring>
Options:
  --version
                      show program's version number and exit
                       show this help message and exit
  -h, --help
  -o OUTPUTXML, --output-xml=OUTPUTXML
                        file in which to store raw xml output
  -s SERVER, --server=SERVER
                       CGHub server location, including protocol and
port,
                       e.g. https://cghub-01.ucsc.edu:20000
  -a, --attributes query the /cghub/metadata/analysisAttributes
instead
                       of analysisObjects resource
                       enable interactive mode
  -i, --interactive
  -c CREDENTIAL, --credential=CREDENTIAL
                        file containing the GeneTorrent credential.
Only
                        required for interactive mode.
                        enable verbose output
  -v, --verbose
<querystring> should be the fully quoted query string, without the
question
mark separator, e.g. "disease abbr=COAD".
```

#### 5.3 Example *cgquery* Commands

The first query example searches for all sequences for a known participant.

The querystring is "participant\_id=6d72de06-232a-4983-a06c-eba6d82cb3f1"

```
: /cghub/metadata/analysisObject
   REST Resource
                            : participant id=6d72de06-232a-4983-
   QueryString
a06c-eba6d82cb3f1
  Output File
                            : None
   Results Returned
______
   Result 1
      analysis_id : e29aa109-d508-4621-9a92-9f7ff7e0018f analysis_data_uri :
https://cghub.ucsc.edu/cghub/data/analysis/download/e29aa109-d508-4621-
9a92-9f7ff7e0018f
       analysis attribute uri :
https://cghub.ucsc.edu/cghub/metadata/analysisAttributes/e29aa109-d508-
4621-9a92-9f7ff7e0018f
      last_modified
center_name
: 2011-06-13T07:00:00Z
: BI
      state
                            : live
      aliquot_id
                            : 087484e8-dc3e-461a-be5f-4217b7c39732
      Files
          filename
                           : C509.TCGA-55-1594-11A-01D-1040-
01.2.bam
         filesize
                          : 23140869502
        checksum
                         : d39c62bdb7abf0213b837fe738f81821
```

The guery returns 7 results. The output has been truncated after the first result.

The next example does the same query but adds the –a option to show all attributes associated with the search results.

```
]$ ./cqquery -a "participant id=6d72de06-232a-4983-a06c-eba6d82cb3f1"
______
  Script Version
                    : https://cghub.ucsc.edu
: /cghub/metadata/analysisAttributes
: participant_id=6d72de06-232a-4983-
   CGHub Server
  REST Resource
  QueryString
a06c-eba6d82cb3f1
  Output File
                        : None
  Results Returned
______
  Result 1
     https://cghub.ucsc.edu/cghub/data/analysis/download/e29aa109-d508-4621-
9a92-9f7ff7e0018f
                      : 2011-06-13T07:00:00Z
     last_modified
center_name
                        : BI
     state study
                        : live
                       : phs000178
                      : WXS
: ILLUMINA
     library_strategy
      platform
```

```
      sample_accession
      : SRS127193

      legacy_sample_id
      : TCGA-55-1594-11A-01D-1040-01

      disease_abbr
      : LUAD

      analyte_code
      : D

      sample_type
      : 11

      tss_id
      : 55

      participant_id
      : 6d72de06-232a-4983-a06c-eba6d82cb3f1

      sample_id
      : 99dc64d5-e61d-4e34-9972-f17d292aec3d

      aliquot_id
      : 087484e8-dc3e-461a-be5f-4217b7c39732

      analysis_xml
      : 134749 bytes of XML

      run_xml
      : 95798 bytes of XML

      experiment_xml
      : 19466 bytes of XML

      Files
      : C509.TCGA-55-1594-11A-01D-1040-

      01.2.bam
      : 23140869502

      checksum
      : d39c62bdb7abf0213b837fe738f81821
```

The query again returns 7 results. The output has been truncated after the first result.

The next query example searches for all sequences that have been modified in the last 5 months.

# The querystring is "last modified=[NOW-5MONTH+TO+NOW]"

```
]$ ./cgquery "last modified=[NOW-5MONTH+TO+NOW]"
______
                 : 1.9
: https://cghub.ucsc.edu
: /cghub/metadata/analysisObject
: last_modified=[NOW-5MONTH+TO+NOW]
: None
   Script Version
  CGHub Server
REST Resource
QueryString
Output File
   Results Returned
                           : 1832
_____
   Result 1
      analysis_id : a03207f9-5fba-4e0f-a9e9-981b13e5c193 analysis_data_uri :
https://cqhub.ucsc.edu/cqhub/data/analysis/download/a03207f9-5fba-4e0f-
a9e9-981b13e5c193
      analysis attribute uri :
https://cghub.ucsc.edu/cghub/metadata/analysisAttributes/a03207f9-5fba-
4e0f-a9e9-981b13e5c193
      state aliquot_id
                           : live
                           : 0031d433-d703-4b2c-9fdf-2920008eb457
      Files
          filename : UNCID_339878.TCGA-A6-2681-01A-01R-
1410-07.110309 UNC3-
RDR300156 00080 FC 62J42AAXX.4.trimmed.annotated.translated to genomic.
```

```
filesize : 1906592634
checksum : 6703e2b337ebac2a4e12cd09f0d1de77
```

This query returned 1832 results. The output has been truncated after the first result.

In order to refine the set of sequences of interest the next example searches for all sequences that have been modified in the last 5 months and have a sample type of 11 (solid tissue normal). The querystring is "last\_modified=[NOW-5MONTH+TO+NOW] &sample\_type=11"

```
]$ ./cgquery "last modified=[NOW-5MONTH+TO+NOW]&sample type=11"
______
   Script Version
                            : 1.9
   CGHub Server
  CGHub Server : https://cghub.ucsc.edu

REST Resource : /cghub/metadata/analysisObject

QueryString : last_modified=[NOW-
                           : https://cghub.ucsc.edu
5MONTH+TO+NOW]&sample type=11
                           : None
   Output File
   Results Returned
                            : 48
Result 1
      analysis_id :
analysis_data_uri :
     analysis_id
                           : 187cc211-9de3-44aa-91f1-74e35a0bf98b
https://cghub.ucsc.edu/cghub/data/analysis/download/187cc211-9de3-44aa-
91f1-74e35a0bf98b
       analysis attribute uri
https://cghub.ucsc.edu/cghub/metadata/analysisAttributes/187cc211-9de3-
44aa-91f1-74e35a0bf98b
      last modified
                           : 2011-08-12T07:00:00Z
      center name
                            : UNC-LCCC
                            : live
      state
      aliquot id
                            : 02222fea-f3da-4328-8188-8ef24d1f55e4
      Files
         filename
                           : UNCID 368056.TCGA-BH-A18U-11A-23R-
A12D-07.110714 UNC13-
SN749 0082 ADODGMABXX.7.trimmed.annotated.translated to genomic.bam
          filesize
                           : 3085973081
          checksum
                         : 6cd3773bb591a255b93223e02daa022b
```

This search returns 48 results. The output has been truncated after the first result.

The final example shows how to search using the wildcard character "\*". Suppose you have a partial aliquot id for a sample and you want all attributes associated with the sequence. The querystring is "aliquot\_id=\*c90-a762-4df7-8b44-b4facd\*"

```
]$ ./cgquery -a "aliquot id=*c90-a762-4df7-8b44-b4facd*"
______
   Script Version
                                      : 1.9
                                    : https://cghub.ucsc.edu
: /cghub/metadata/analysisAttributes
: aliquot_id=*c90-a762-4df7-8b44-
    CGHub Server
    REST Resource
    QueryString
b4facd*
   Output File
                                      : None
   Results Returned
    Result 1
         analysis_id : 8ca52fbf-60ab-477f-90cd-0704d6b5b42d analysis_data_uri :
https://cghub.ucsc.edu/cghub/data/analysis/download/8ca52fbf-60ab-477f-
90cd-0704d6b5b42d
        04d6b5b42d
last_modified : 2011-08-12T07:00:00Z
center_name : UNC-LCCC
state : live
study : phs000178
library_strategy : RNA-Seq
platform : ILLUMINA
. SRS136282
                                  : SRS136282
: TCGA-BP-4342-01A-01R-1289-07
: KIRC
         sample accession
         legacy sample id
         disease_abbr
         css_id
participant_id
sample_id
aliquot_id
ana1
                                     : R
                                 : 01
: BP
                                     : c0e26587-3191-4396-a20b-da9e7de213f4
                             : CUe2658/-3191-4396-a200-ua9e/ue21314

: 04d3168b-ae52-400c-94fa-8ec43eface20

: 005a0c90-a762-4df7-8b44-b4facd06e9b1

: 2167 bytes of XML

: 663 bytes of XML

: 2432 bytes of XML
         aliquot_id
analysis_xml
run_xml
experiment_xml
         Files
                               : UNCID_374694.TCGA-BP-4342-01A-01R-
             filename
1289-07.110511 UNC11-
SN627 0085 BC01AFABXX.5.trimmed.annotated.translated to genomic.bam
             filesize : 7978810705
checksum : 09e2d4d12be
                                       : 09e2d4d12bec19701b46b2b558a0b084
______
   Parse Summary
        rarse Errors : 0
Parse Warnings : 0
        Parse Errors
```

#### 5.4 CGHub REST API

The table below describes the query APIs available from CGHub.

	Resource URI	Purpose
Analysis	/cghub/metadata/analysisObject?	Returns minimal fields to
Object	{querystring}	identify the desired object as
Query		well as a URI to download the
		data files.
Analysis	/cghub/metadata/analysisAttributes?	Returns the full Analysis
Attribute	{querystring}	object attributes including the
Query		sample information and
		submission XML objects.

Details of the {querystring} are described in section 5.1 **Error! Reference source ot found.**.

The CGHub REST API returns query results in XML format. The results will be saved to a file if the –O <filename> option is used.

The **analysisObject** query returns minimal attributes to locate files for download. This format can be used to identify a list of sequences to pass as input to GeneTorrent to perform the transfer. It will return the following attributes:

- analysis\_id
- last\_modified
- center\_name
- state
- aliquot\_id
- files
  - o filename
  - o filesize
  - o checksum
- analysis data URI (used to request a download of this object)
- analysis attributes URI (used to get the detailed attributes for this object)

The **analysisAttributes** query returns all of the analysisObject attributes as well as the following additional attributes:

• study

- sample accession
- legacy\_sample\_id
- disease\_abbr
- analyte\_code
- sample\_type
- tss\_id
- participant\_id
- sample\_id
- XML documents from the original submission
  - o analysis.xml
  - o run.xml
  - experiment.xml

# 5.5 Sample REST Calls

The queries will return XML formatted data, including a URL to download the associated data files for each returned analysis object.

For example enter the following query into a Web browser:

```
https://cghub.ucsc.edu/cghub/metadata/analysisObject?aliquot_id=c0cfafbc-6d07-4ed5-bfdc-f5c3bf8437f6
```

#### OR execute the command line:

```
wget --no-check-certificate -0 xmlout
https://cghub.ucsc.edu/cghub/metadata/analysisObject?aliquot_id=
c0cfafbc-6d07-4ed5-bfdc-f5c3bf8437f6
```

#### The returned XML will look like:

The results return the same attribute fields that a *cgquery* search does except that they are in XML format.

The same query can be done to return the full attribute set by substituting "analysisAttribute" for "analysisObject" in the query URI. This query saves the results to the file "xmlout".

```
] $ wget --no-check-certificate -0 xmlout
https://cghub.ucsc.edu/cghub/metadata/analysisAttributes?aliquot id=c0c
fafbc-6d07-4ed5-bfdc-f5c3bf8437f6
--2012-01-09 12:44:48--
https://cghub.ucsc.edu/cghub/metadata/analysisAttributes?aliquot id=c0c
fafbc-6d07-4ed5-bfdc-f5c3bf8437f6
Resolving cghub.ucsc.edu... 192.35.223.5
Connecting to cghub.ucsc.edu|192.35.223.5|:443... connected.
WARNING: certificate common name "*.cghub.ucsc.edu" doesn't match
requested host name "cghub.ucsc.edu".
HTTP request sent, awaiting response... 200 OK
Length: 3072 (3.0K) [text/xml]
Saving to: "xmlout"
2012-01-09 12:44:54 (6.74 MB/s) - "xmlout" saved [3072/3072]
```

# Note that the results contain the raw XML in the <analysys\_xml> section.

```
<?xml version="1.0" encoding="utf-8" standalone="ves"?>
<ResultSet date="2012-01-09 12:44:54">
      <Query>aliquot id:c0cfafbc-6d07-4ed5-bfdc-f5c3bf8437f6</Query>
      <Hits>1</Hits>
      <Result id="1">
            <analysis id>8d0b3af6-0f78-4282-bde5-
c5238fa3d4c2</analysis id>
            <state>live</state>
            <last modified>2011-08-10T07:00:00Z</last modified>
            <upload date></upload date>
            <center name>BCCAGSC</center name>
            <study>phs000178</study>
            <aliquot id>c0cfafbc-6d07-4ed5-bfdc-
f5c3bf8437f6</aliquot id>
            <files>
                  <file>
                        <filename>TCGA-CZ-5466-11A-01R-1502-
13 mirna.bam</filename>
                        <filesize>132976422</filesize>
                        <checksum
type="MD5">6ede057c98171b986a36e77a9bf7efc0</checksum>
                  </file>
            </files>
            <sample accession>SRS238705</sample accession>
            <legacy sample id>TCGA-CZ-5466-11A-01R-1502-
13</legacy sample id>
            <disease abbr>KIRC</disease abbr>
            <tss id>CZ</tss id>
            <participant id>5722df9f-5631-476d-a11b-
b3c1e9a40fbf</participant id>
            <sample id>f072dc5f-6b9a-4935-8ceb-f9f3fd7ebe73/sample id>
            <analyte code>R</analyte code>
            <sample type>11</sample type>
            library strategy></library strategy>
            <platform></platform>
            <analysis xml><ANALYSIS SET</pre>
xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance">
  <ANALYSIS center name="BCCAGSC" analysis date="2011-07-27T18:28:48Z"
analysis center="BCCAGSC" alias="M02084" accession="SRZ017243">
    <TITLE>Reference alignment of M02084</TITLE>
    <STUDY REF accession="SRP000677"/>
    <DESCRIPTION/>
    <ANALYSIS TYPE>
      <REFERENCE ALIGNMENT>
        <ASSEMBLY>
          <STANDARD short name="GRCh37-lite"/>
        </ASSEMBLY>
        <PROCESSING>
          <PIPELINE>
            <PIPE SECTION section name="Alignment">
              <STEP INDEX>1</STEP INDEX>
              <PREV STEP INDEX>NIL</prev STEP INDEX>
              <PROGRAM>BWA</PROGRAM>
              <VERSION>0.5.7</VERSION>
              <NOTES/>
```

```
</PIPE SECTION>
          </PIPELINE>
          <DIRECTIVES>
<alignment includes unaligned reads>true</alignment includes unaligned</pre>
reads>
<alignment marks duplicate reads>true</alignment marks duplicate reads>
<alignment includes failed reads>true</alignment includes failed reads>
          </PROCESSING>
     </REFERENCE ALIGNMENT>
    </ANALYSIS TYPE>
    <TARGETS>
      <TARGET accession="SRS238705" refcenter="TCGA"
sra object type="SAMPLE" refname="TCGA-CZ-5466-11A-01R-1502-13"/>
   </TARGETS>
    <DATA BLOCK>
      <FILES>
        <FILE checksum="6ede057c98171b986a36e77a9bf7efc0"</pre>
filetype="bam" filename="TCGA-CZ-5466-11A-01R-1502-13 mirna.bam"
checksum method="MD5"/>
     </FILES>
    </DATA BLOCK>
    <ANALYSIS ATTRIBUTES>
     <ANALYSIS ATTRIBUTE>
        <TAG/>
       <VALUE/>
      </ANALYSIS ATTRIBUTE>
    </ANALYSIS ATTRIBUTES>
  </ANALYSIS>
</ANALYSIS SET>
</analysis xml>
            <run xml></run xml>
            <experiment xml></experiment xml>
      <analysis data uri>https://cghub.ucsc.edu/cghub/data/analysis/dow
nload/8d0b3af6-0f78-4282-bde5-c5238fa3d4c2</analysis data uri>
      </Result>
</ResultSet>
```

# 6 Sequence Download

To download sequence data files from the CGHub system securely, efficiently and reliably, users are required to use the GeneTorrent download client, an open-source application available at no cost to users.

6.1 Downloading Files using GeneTorrent Client

The GeneTorrent client runs as a command line tool on all systems

The download mode command line arguments are:

GeneTorrent -d [ URI | UUID | .xml | .gto ] -c credentials [ -p path ]

-d content-specifier, --download content-specifier

The content-specifier should be one of the following:

- A fully-qualified URI pointing to an analysis object at the GeneTorrent Executive.
- A UUID denoting an analysis object at the GeneTorrent Executive. In this case, GeneTorrent will construct a URI based on the default server, currently https://cghub.ucsc.edu.
- An XML file, which will be parsed to obtain a list of URIs.
- A .gto file directly, which eliminates all calls to the GeneTorrent Executive and causes a download to begin. In this case you do not need to supply access credentials using the -c option.

Note that multiple -d options can be specified on the same command line, in which case multiple analysis objects will be downloaded.

-c credential-file, --credentialFile credential-file

A credential file is required for download mode, unless you are directly passing a .gto file. Pass the full or relative path to the file containing the access credentials (security token) previously received from the User Authentication, (see section 4 for details).

-p path, --path path

Path to save data files in the gto file(s). UUID is part of the gto and will always be added to path, so data files will be found at path/UUID. The current directory will be used by default.

#### --maxChildren max-children

The maximum number of parallel children that should be spawned to perform the download. By default up to 8 children are used, but this number may be adjusted. In general you should have fewer children than you have CPU cores on your machine.

#### 6.2 Query and Download

#### **Step 1:** Query for the desired files

```
% cgquery "center name=CGHubTest"
______
   Script Version
                    : https://cghub.ucsc.edu
: /cghub/metadata/analysisObject
: center_name=CGHubTest
: None
                            : 1.9
   CGHub Server
   REST Resource
QueryString
                             : None
   Output File
   Result 1
      analysis id
                             : FC31132C-CD53-44D5-A9F2-0F4AEE52965C
       analysis_data_uri
https://cghub.ucsc.edu/cghub/data/analysis/download/FC31132C-CD53-44D5-
A9F2-0F4AEE52965C
       analysis attribute uri :
https://cqhub.ucsc.edu/cqhub/metadata/analysisAttributes/FC31132C-CD53-
44D5-A9F2-0F4AEE52965C
                        : 0011-09-18 12:00:00.0
                            : CGHubTest
: live
       last modified
       center_name
       state
       aliquot_id
                             : 1643378B-CBAF-43E7-A9A0-157EB1FF3F02
      study
                             : CGHUB
      Files
          filename
HG00114.mapped.ILLUMINA.bwa.GBR.low coverage.20101123.bam
          filesize
                     : 11643625129
   Result 2
      https://cghub.ucsc.edu/cghub/data/analysis/download/0344FD0B-3435-4405-
B284-C462849DF935
       analysis attribute uri :
https://cghub.ucsc.edu/cghub/metadata/analysisAttributes/0344FD0B-3435-
4405-B284-C462849DF935
      last_modified : 0011-09-17 12:00:00.0 center_name : CGHubTest
       state
                            : verifying sample
       aliquot_id
                            : 38220E07-7242-4274-8155-24E45972A08F
       study
                             : CGHUB
       Files
```

```
filename :
HG00099.unmapped.SOLID.bfast.GBR.low_coverage.20101123.bam
filesize : 24032208912
```

**Step 2:** Invoke GeneTorrent on the command line with the desired file URI(s) and user credentials. Multiple –d options can be provided to download multiple files.

```
GeneTorrent -v -c <credential> -d
https://cghub.ucsc.edu/cghub/data/analysis/download/FC31132C-CD53-44D5-
A9F2-0F4AEE52965C
```

The GeneTorrent client connects to the CGHub system, bi-directionally authenticates using the credential, verifies that the user is authorized to access the file, and starts the SSL-secured file transfer.

GeneTorrent, by default, runs silently to completion. If invoked with the '-v' (verbose) option, the user will see:

```
Welcome to GeneTorrent version x.y.z, download mode.
Status: checking (r) downloaded: ( )
uploaded: ( )
Status: downloading downloaded: 7.34MB (718.kB/s)
uploaded: 9.88kB (21.7kB/s)
Status: downloading downloaded: 14.4MB (1.33MB/s)
uploaded: 19.5kB (39.3kB/s)
```

# **Step 3:** Verify that the download was successful.

When GeneTorrent finishes transferring the file, the user will find the BAM file in the current working directory, or in the location specified by the '-p' (path) option.

The file length and the MD5 checksum of the file can be verified using standard tools provided by the operating system.

# 6.3 Query and Download using cgquery Interactive Mode

The *cgquery* –i option will prompt the user to download one or more of the returned files using GeneTorrent. The interactive command requires the user credentials.

```
Output File
                                : None
   Result 1
                                : FC31132C-CD53-44D5-A9F2-0F4AEE52965C
       analysis_id
analysis_data_uri
        analysis id
https://cghub.ucsc.edu/cghub/data/analysis/download/FC31132C-CD53-44D5-
A9F2-0F4AEE52965C
       analysis attribute uri
https://cghub.ucsc.edu/cghub/metadata/analysisAttributes/FC31132C-CD53-
44D5-A9F2-OF4AEE52965C
       last_modified center_name
                               : 2011-09-18T19:00:00Z
                               : CGHubTest
                               : live
       state
       aliquot id
                               : 1643378B-CBAF-43E7-A9A0-157EB1FF3F02
                               : CGHUB
       study
       Files
           filename
HG00114.mapped.ILLUMINA.bwa.GBR.low coverage.20101123.bam
                               : 11643625129
           filesize
   Result 2
                              : 0344FD0B-3435-4405-B284-C462849DF935
       analysis id
       analysis_id
analysis_data_uri
https://cghub.ucsc.edu/cghub/data/analysis/download/0344FD0B-3435-4405-
B284-C462849DF935
        analysis attribute uri
https://cghub.ucsc.edu/cghub/metadata/analysisAttributes/0344FD0B-3435-
4405-B284-C462849DF935
                             : 2011-09-17T19:00:00Z
       last modified
       center name
                               : CGHubTest
       state
                               : live
                               : 38220E07-7242-4274-8155-24E45972A08F
       aliquot id
                                : CGHUB
       study
        Files
           filename
HG00099.unmapped.SOLID.bfast.GBR.low coverage.20101123.bam
           filesize : 24032208912
Enter index of URI to download (0 for all) or 'q' to quit
    [ 0] : All URIS
    [ 1] : https://cghub-
test.ucsc.edu/cghub/data/analysis/download/FC31132C-CD53-44D5-A9F2-
0F4AEE52965C
https://cghub.ucsc.edu/cghub/data/analysis/download/0344FD0B-3435-4405-
B284-C462849DF935
Index> 2
Processing URI
https://cqhub.ucsc.edu/cqhub/data/analysis/download/0344FD0B-3435-4405-
B284-C462849DF935
Executing command '/usr/bin/GeneTorrent -v -c mykey.pem -d
https://cghub.ucsc.edu/cghub/data/analysis/download/0344FD0B-3435-4405-
```

```
B284-C462849DF935'

Welcome to GeneTorrent version x.y.z, download mode.
Status: checking (r) downloaded: ( )
uploaded: ( )
Status: downloading downloaded: 7.34MB (718.kB/s)
uploaded: 9.88kB (21.7kB/s)
Status: downloading downloaded: 14.4MB (1.33MB/s)
uploaded: 19.5kB (39.3kB/s)
```

