

TCGA: A Community Resource Looking for a Broader Community

Kenna Shaw, Ph.D.
Director
The Cancer Genome Atlas

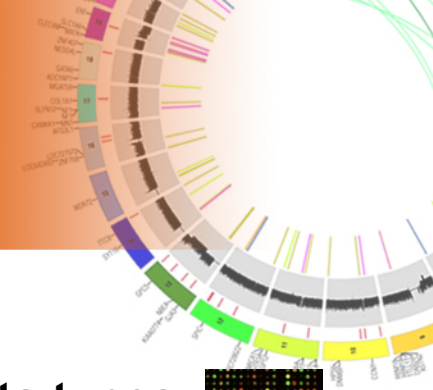
TCGA: Core Objectives



Launched in 2006 as a pilot and expanded in 2009, the goals of TCGA are to:

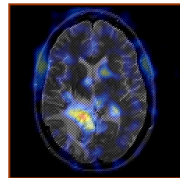
- Establish infrastructure for effective team science
- Develop a scalable “pipeline” beginning with highest quality samples
- Determine the feasibility of a large-scale, high throughput approach to identifying the molecular ‘parts-list’
- Evaluate using statistically-robust sample sets
- Make the data publicly and broadly available to the cancer community while protecting patient privacy

TCGA: “No Platform Left Behind”

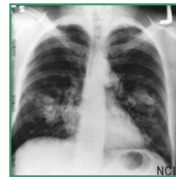


25* forms of cancer

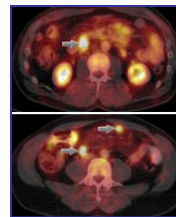
glioblastoma multiforme
(brain)



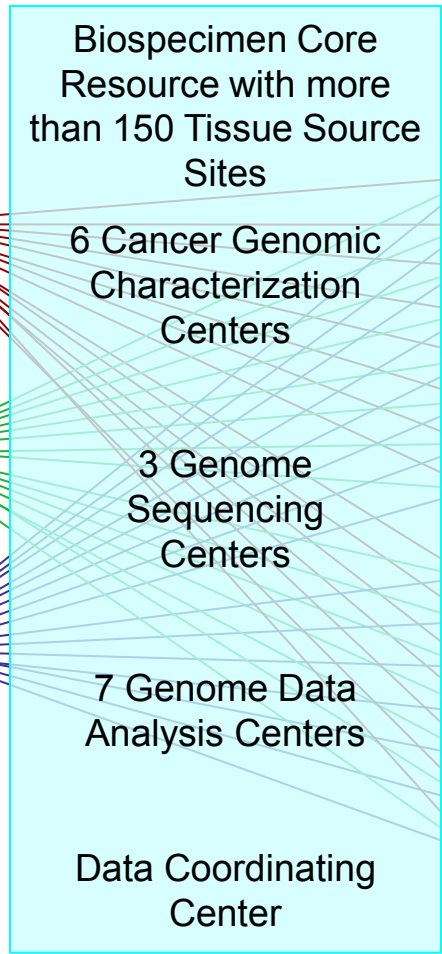
squamous carcinoma
(lung)



serous
cystadenocarcinoma
(ovarian)

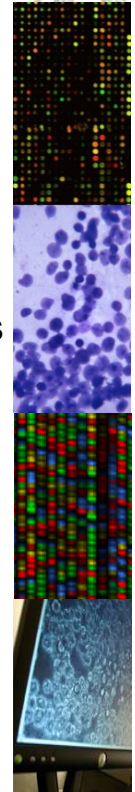


Etc. Etc. Etc.

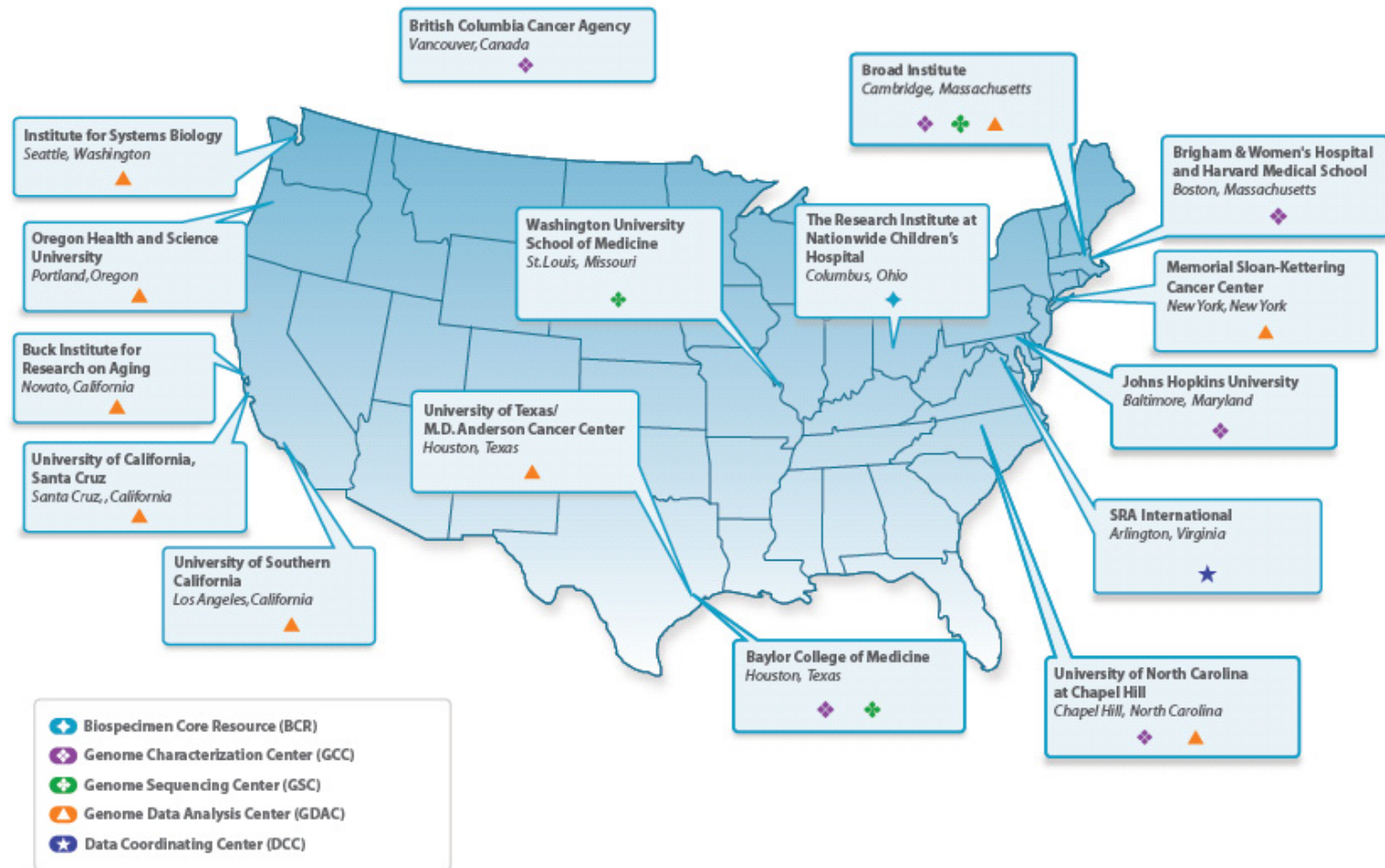
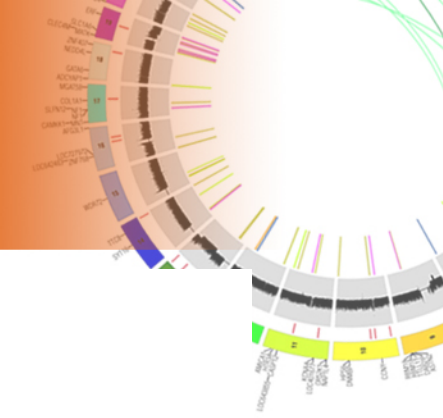


Multiple data types

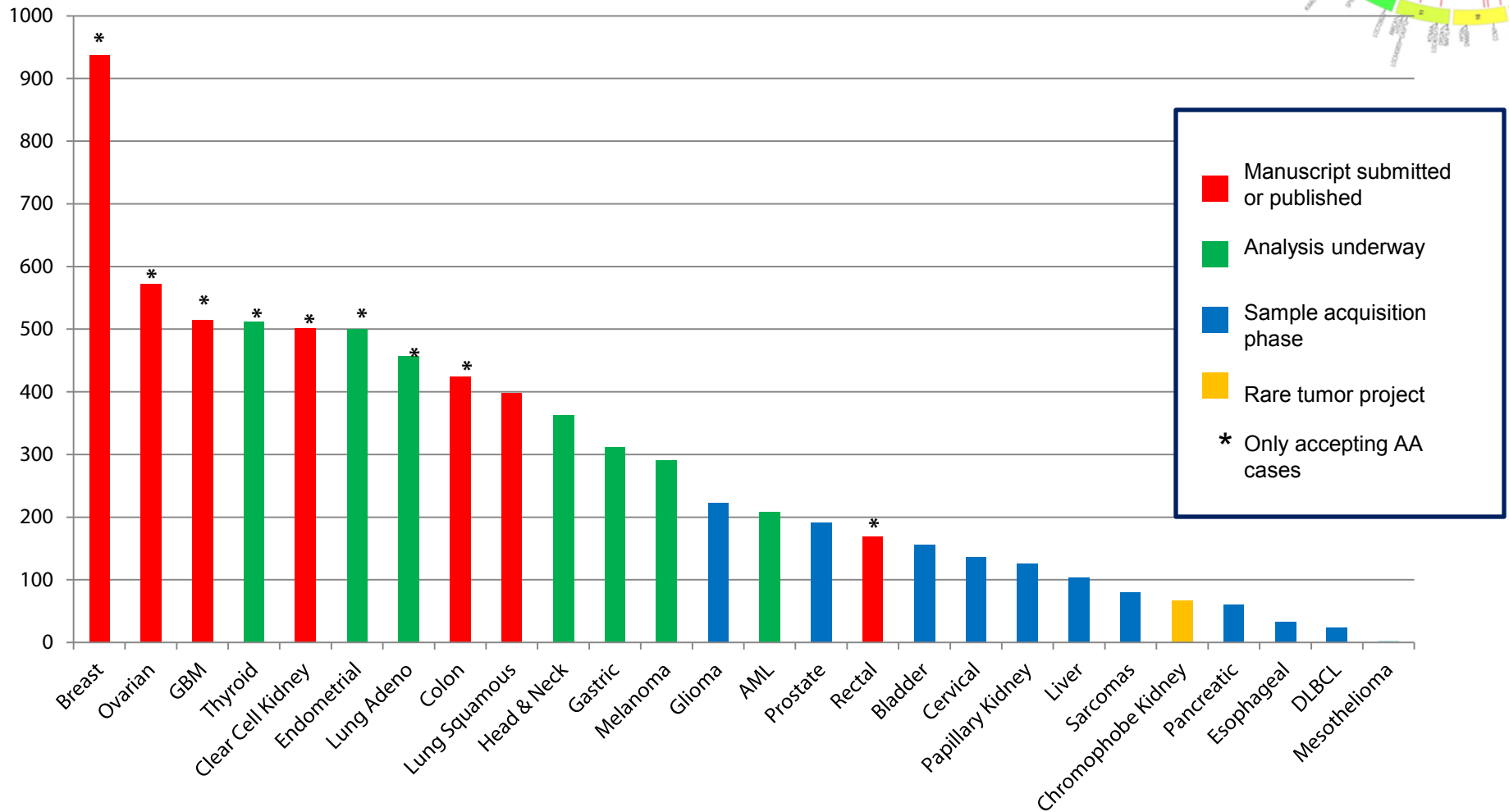
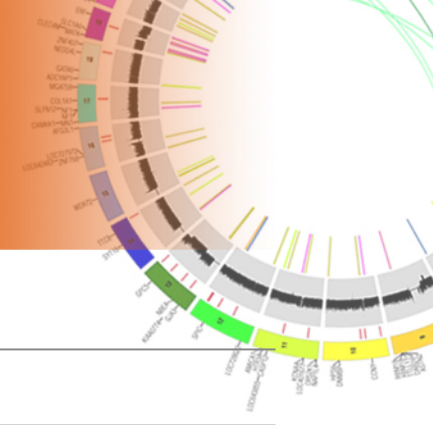
- Clinical diagnosis
- Treatment history
- Histologic diagnosis
- Pathologic report/images
- Tissue anatomic site
- Surgical history
- Gene expression/RNA sequence
- Chromosomal copy number
- Loss of heterozygosity
- Methylation patterns
- miRNA expression
- DNA sequence
- RPPA (protein)
- Subset for Mass Spec



TCGA Research Network

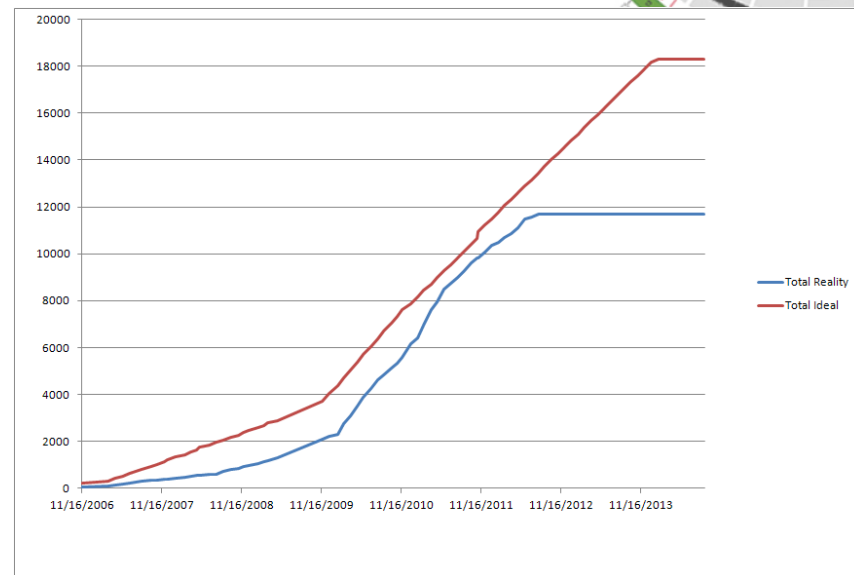


Tumor Project Progress



Rare Tumor Project (Initiated March 2012)

- Adrenocortical Carcinoma
- Adult ALL (B-cell and T-Cell)
- Anaplastic Thyroid
- Cholangiocarcinoma
- Chromophobe kidney
- High Risk MDS (del 5q- cases)
- Mesothelioma
- Paraganglioma/Pheochromocytoma
- Testicular Germ Cell
- Thymoma
- Uterine Carcinosarcoma
- Sarcomas
- Others??



TCGA: Platforms- Then and Now



Platform	Pilot	Expansion
SNP/CNV	Affy SNP 6.0 Agilent CGH Array Illumina 1M Duo	Affy SNP 6.0 Low Pass Sequencing*
Methylation	Infinium Array	Infinium Array
mRNA	Agilent 244K Array Affy Human Exon Array Affy U133 Array	RNAseq
miRNA	Agilent 8 x 15K Array	RNAseq
Mutation	600-1000 genes	DNaseq: 100% whole exomes 10% whole genomes
Proteomics	None	Reverse Phase Protein Arrays
Clinical Data	Minimum Enrollment & Follow-Up* H&E from Frozen Section Images Treatment Data	Minimum Enrollment & Follow-Up* Pathology Reports H&E from Frozen & Diagnostic Images

*- Not a core platform; Not all samples will have data file for this platform

More information on platforms and data available at: <http://tcga-data.nci.nih.gov/tcga/tcgaPlatformDesign.jsp>

TCGA: Data Availability



- 7,136 cases across 20+ tumor types
- 5865 with minimum clinical data set
- 3893 with at least 1 year follow-up; ~50% with treatment data
- 105,000 samples of RNA/DNA/protein shipped between 2006 & 6/2012
- All but 13,000 samples have data returned:
 - ~87% of all samples have Level 1 data publicly available
 - TCGA Program Office to complete accounting with Batelle/QMS by end of 2012

TCGA: The Pipeline for Comprehensive Characterization

Tissue Sample



Pathology QC

DNA & RNA Isolation, QC

Sequencing

Expression, CNA & LOH, Epigenetics

Data Storage at DCC & CGHub

GDAC

Integrative Analysis

Comprehensive Characterization of a Cancer Genome

3 months – 2 years

~90d

SNP 6.0 ~45d

Methylation ~60d

miRNAseq ~105d

mRNAseq ~120d

DNAseq Exome ~180d

ARTICLE

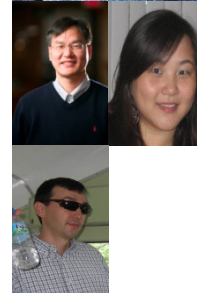
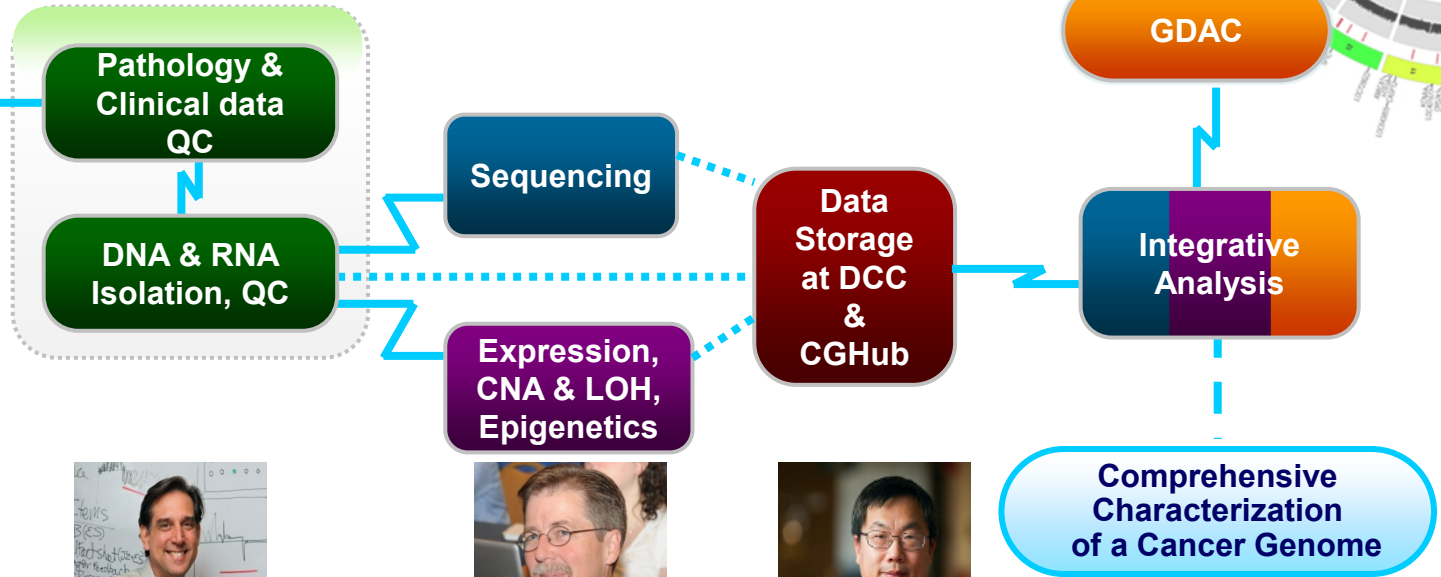
Comprehensive molecular characterization of human colon and rectal cancer



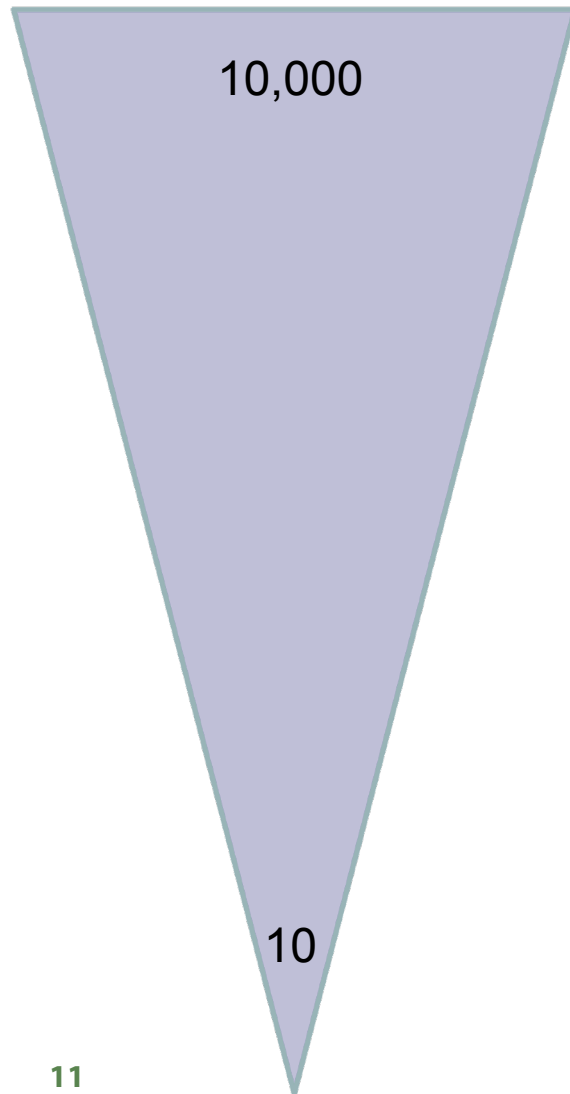
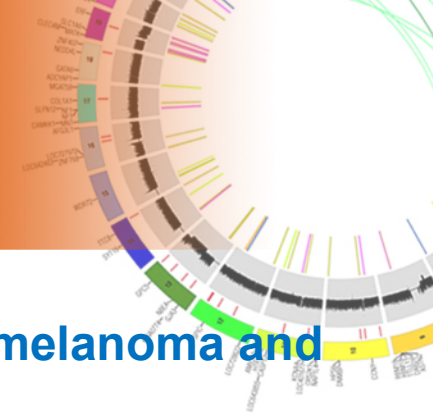
~12-24 months

TCGA: The Pipeline for Comprehensive Characterization

Tissue Sample

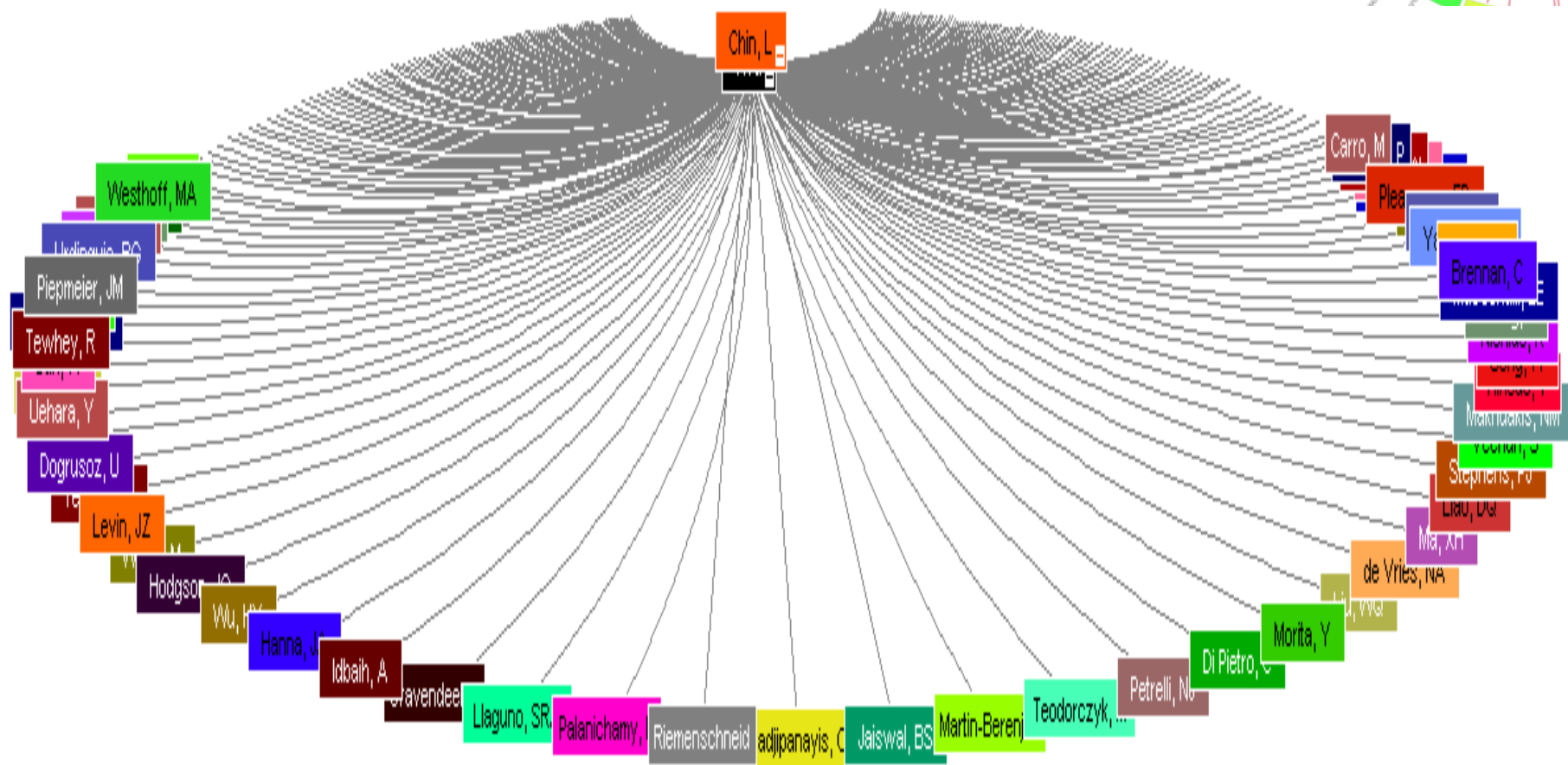


Sample Criteria Limit 'Askable' Questions



- **Primary, adult** tumors (except for melanoma and triplets)
- **Malignant** (no *in situ* cases)
- Snap **frozen**, <60min from clamp to LN2
- ~ **50 mg** (**biopsies starting to be feasible**)
- Pathology review of tissue sent to TCGA
- No more than 20% necrosis ; **≥ 60%*** tumor cells
- **No prior treatment**
- **Matched source of germline**: Blood (buffy coat/white cells)/saliva or skin for liquid tumors
- **Clinical annotation**; but not pre-analytic variables
- **IRB approval for use in TCGA; proactive consenting for genomic studies**
- **MTA w/out retention of IP**

End Goal: Making an Exhaustible Resource Inexhaustible



Source: ISI Web of Knowledge™, www.thomsonscientific.com

Where to find TCGA Sequence Data



- Moved from Short Read Archive (SRA at NCBI) to UCSC
- Open for downloads as of January 2012

UPGRADE NOTICE: The CGHub system has been upgraded to use GeneTorrent 2.1.2.0. All CGHub users: please update your systems to GeneTorrent version 2.1.2.0, available [here](#) now. Please email us at: support@cghub.ucsc.edu if you have any concerns or questions regarding this process.

Cancer Genomics Hub BETA (initial limited release)
The Cancer Genomics Hub (CGHub) is a secure repository for storing, cataloging, and accessing cancer genome sequences, alignments, and mutation information from The Cancer Genome Atlas (TCGA) consortium project and other cancer related projects. The CGHub mission is to facilitate the work of scientific researchers. CGHub is designed to be a fully automated resource, appearing to the user as an extension of their home institute computing resource. User scripts employ the CGHub Application Programming Interface (API) to retrieve the CGHub catalog (index) of files, select files for download based on metadata attributes such as cancer type, sequence type, source sequencing center or data range, then initiate download, and finally confirm success.

BAM files and metadata are available for download using the GeneTorrent client compiled for a variety of Linux operating systems on Intel x86_64 platforms. CGQuery is a Python program for querying the metadata using the Web Services API that is compatible with a wide range of systems and is available on the [downloads](#) page. Please view the [CGHub User Guide](#) and the [Install README](#) for installation instructions and troubleshooting information.

All researchers using CGHub must meet the access and use criteria established by the NIH to ensure the privacy, security, and integrity of participant data (see directions below).

Important notice concerning TCGA metadata: The migration of the TCGA data to CGHub has uncovered inconsistent and missing XML metadata files. We have a project underway to correct these problems, however it is a time consuming process. Please check your search and download results carefully see they are what is expected. Please report any inconsistencies that are encountered to support@cghub.ucsc.edu.

CGHub Links

- CGHub Download User's Guide
- CGHub Download Quick Start Guide
- CGHub Download Troubleshooting
- GeneTorrent Technical Manual Page
- GeneTorrent Software Downloads
- CGHub Submissions/Upload Guide

E-Mail CGHub Technical Support

Subscribe to the CGHub Updates Email List

This low volume list will include announcements regarding updates, maintenance windows, downtime, and other important CGHub issues.

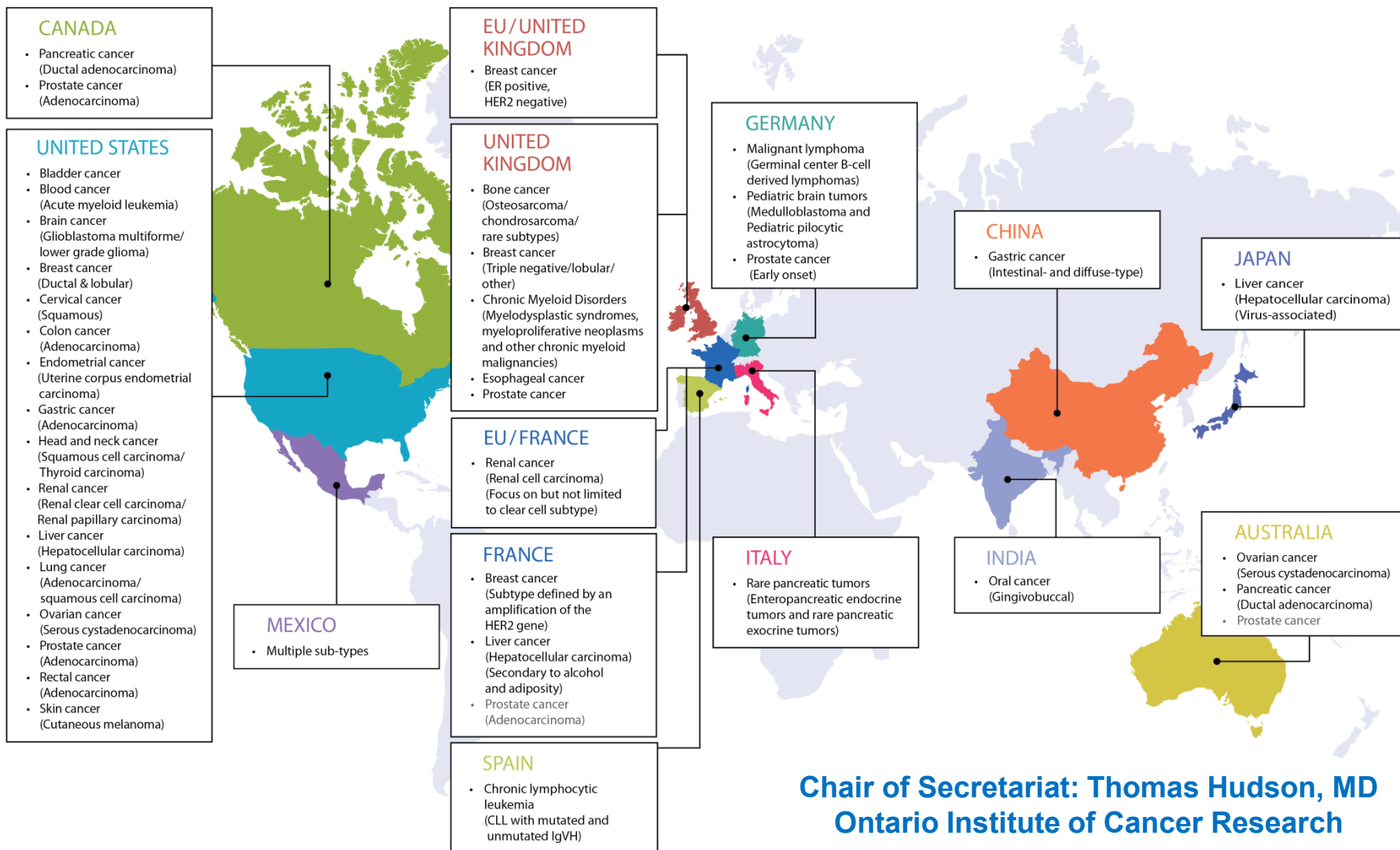
CGHub Authorized Access

- 2 petabytes now, 5Pb (5×10^{15}) total by 2014
- General Parallel File System, Dual RAID 6 subsystems, Redundant I/O paths
- Currently holds 10,000 files; expected to double in next 3 months
- Co-location opportunities in same data center for groups who want to compute on the data
- User support: support@cghub.ucsc.edu

<https://cghub.ucsc.edu>



March 2011 : International Cancer Genome Consortium Projects



**Chair of Secretariat: Thomas Hudson, MD
Ontario Institute of Cancer Research**

Acknowledgements



NCI Program Office

Kenna Shaw Liming Yang Roy Tarnuzzer Zhining Wang Emma Spaulding Margi Sheth John Demchok Julia Zhang Martin Ferguson Greg Eley



Center for Cancer Genomics
Stephen Chanock & Lou Staudt



NHGRI

Brad Ozenberger Heidi Sofia Lindsay Lund Mark Guyer