

Virtual Data Integration and

The Genboree Network

Aleksandar Milosavljevic
Bioinformatics Research Laboratory (BRL)
Baylor College of Medicine

6th Genboree Workshop on Epigenome
Informatics
March 4th, 2013

Two key integration problems

- Genboree Workbench: “Omic” **Toolset Integration**
- Genboree Network: Virtual **Data Integration**

- Genboree Workbench: “Omic” Toolset Integration
- Genboree Network: Virtual Data Integration

Epigenomic Toolset integrates Spark developed by the British Columbia Genome Center in Vancouver

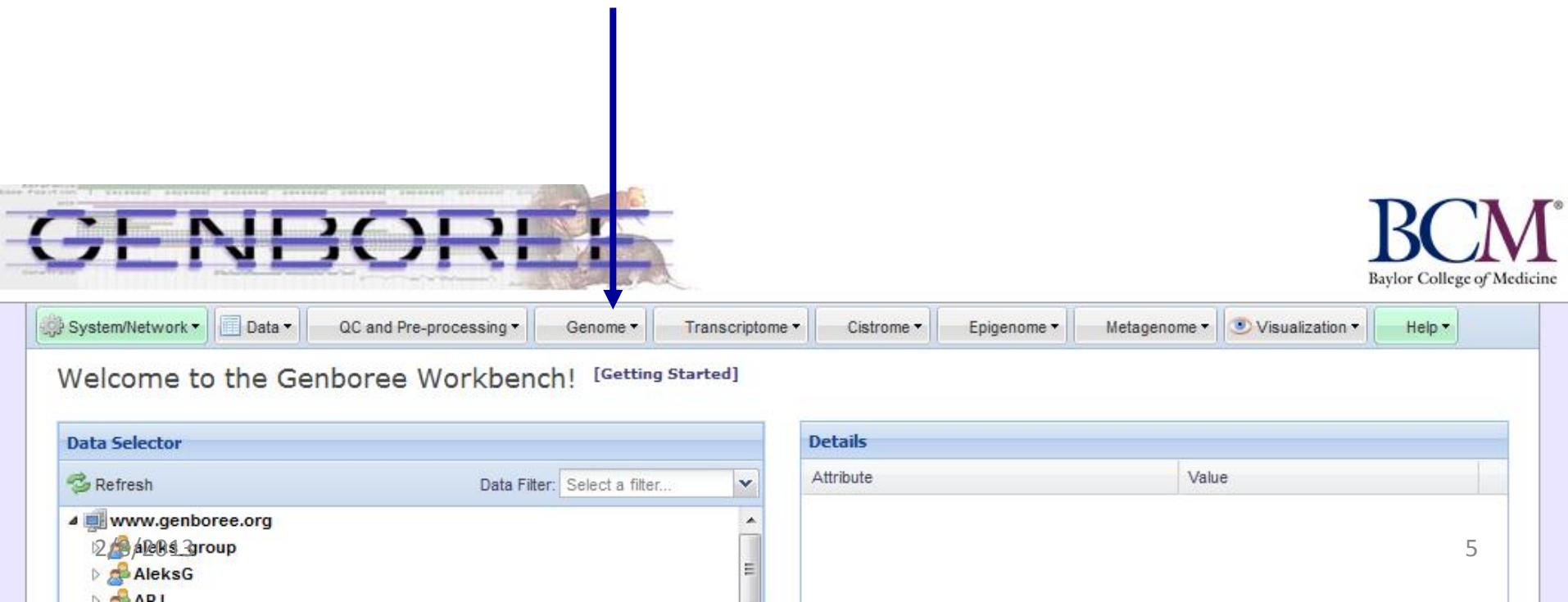
Spark:

Nielsen CB et al. Genome Res. (11):2262-9 2012

The screenshot shows the Genboree Workbench interface. At the top, there is a banner with the word "GENBOREE" and a small illustration of a person. Below the banner is a navigation bar with tabs: System/Network, Data, QC and Pre-processing, Genome, Transcriptome, Cistrome, Epigenome, Metagenome, Visualization, and Help. The "Epigenome" tab is currently selected. A large blue arrow points downwards from the "Epigenome" tab towards the main workspace. The workspace is divided into two main sections: "Data Selector" on the left and "Details" on the right. The "Data Selector" panel contains a "Refresh" button, a "Data Filter" dropdown set to "Select a filter...", and a tree view of data sources. The tree shows a root node "www.genboree.org" which has a child node "aleks_group". The "Details" panel has a header "Details" and a table with columns "Attribute" and "Value". The table is currently empty.

Genomic Toolset developed in collaboration with the Baylor Genome Center

Atlas2 genome resequencing:
Evani US et al. BMC Genomics 6:S19 2012



The screenshot shows the Genboree Workbench interface. At the top, there is a banner with the word "GENBOREE" and a small illustration of a person. Below the banner is a menu bar with the following items: System/Network, Data, QC and Pre-processing, Genome (highlighted with a blue arrow), Transcriptome, Cistrome, Epigenome, Metagenome, Visualization, and Help. The main area displays a "Welcome to the Genboree Workbench!" message and a "Getting Started" link. On the left, there is a "Data Selector" panel with a "Refresh" button and a "Data Filter" dropdown set to "Select a filter...". It lists several items: "www.genboree.org", "AleksG_group" (with a count of 12), "AleksG" (with a count of 1), and "API". On the right, there is a "Details" panel with a table header "Attribute" and "Value". The page number "5" is located in the bottom right corner.

Attribute	Value



Microbiome Toolset includes QIIME (U. Colorado)

The screenshot shows the Genboree Workbench interface. At the top, a red box highlights the text "Run at Bioinformatics Research Laboratory". In the top right corner, the Baylor College of Medicine (BCM) logo is visible. The main interface has a navigation bar with tabs: System/Network, Data, QC and Pre-processing, Genome, Transcriptome, Cistrome, Epigenome, Metagenome, Visualization, and Help. A "Data Selector" sidebar on the left lists various user profiles and projects, with a red box highlighting the "www.genboree.org" entry. The main workspace shows a "Welcome to the Genboree Workbench!" message and a "Getting Started" link. A dropdown menu under the "Data Initialization" tab is open, showing options like "Import 16S rRNA Sequences" and "Upload Metagenomic WGS Results". A red box highlights the "Data Analysis" option in this menu. Another red box highlights the "Data Analysis" tab in the main toolbar. A large red box encloses the "Data Analysis" section of the interface, which contains sub-options: "Taxonomic Classification (RDP)", "QIIME", "Alpha Diversity", "Machine Learning", and "WGS Sample Grid Viewer". Below this is an "Input Data" section with up, down, and remove buttons. The entire interface is framed by a red border.

The Genboree Microbiome Toolset and the analysis of 16S rRNA microbial sequences.

Riehle K, Coarfa C, Jackson A, Ma J, Tandon A, Paithankar S, Raghuraman S, Mistretta TA, Saulnier D, Raza S, Diaz MA, Shulman R, Aagaard K, Versalovic J, Milosavljevic A.

BMC Bioinformatics. 2012;13 Suppl 13:S11. doi: 10.1186/1471-2105-13-S13-S11. Epub 2012 Aug 24. PMID = 23320832

Other widely used tools are integrated within the Genboree Workench

- MACS
- TopHat
- CuffLiks
- CuffDiff.
-

- Genboree Workbench: “Omic” Toolset Integration
- Genboree Network: Virtual Data Integration

As the data volume explodes, *physical* data integration will become impossible

- **Global proliferation of epigenomic profiling** will make it impossible to physically integrate data for analysis
- **IRB-mandated local physical custody** of potentially identifiable human subject data.
- Network bandwidth increases are not keeping up with the **rapidly escalating bandwidth** of sequencing machines.

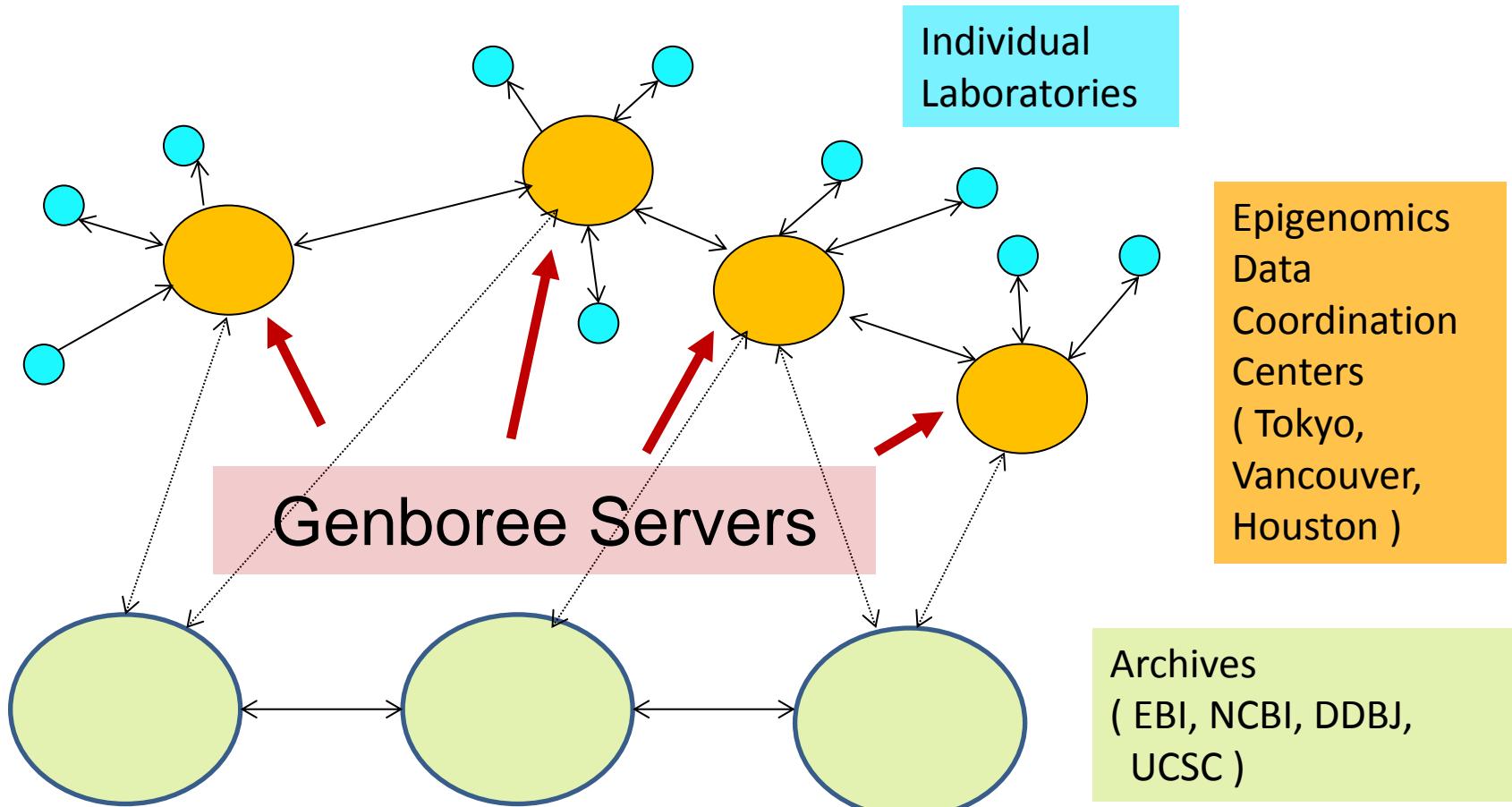


As the data volume explodes, *physical* data integration will become impossible

Solution: *virtual* data integration

Data Ecosystem Model proposed for IHEC

IHEC: International Human Epigenome Consortium



Virtual integration of reference epigenomes across three IHEC Data Centers



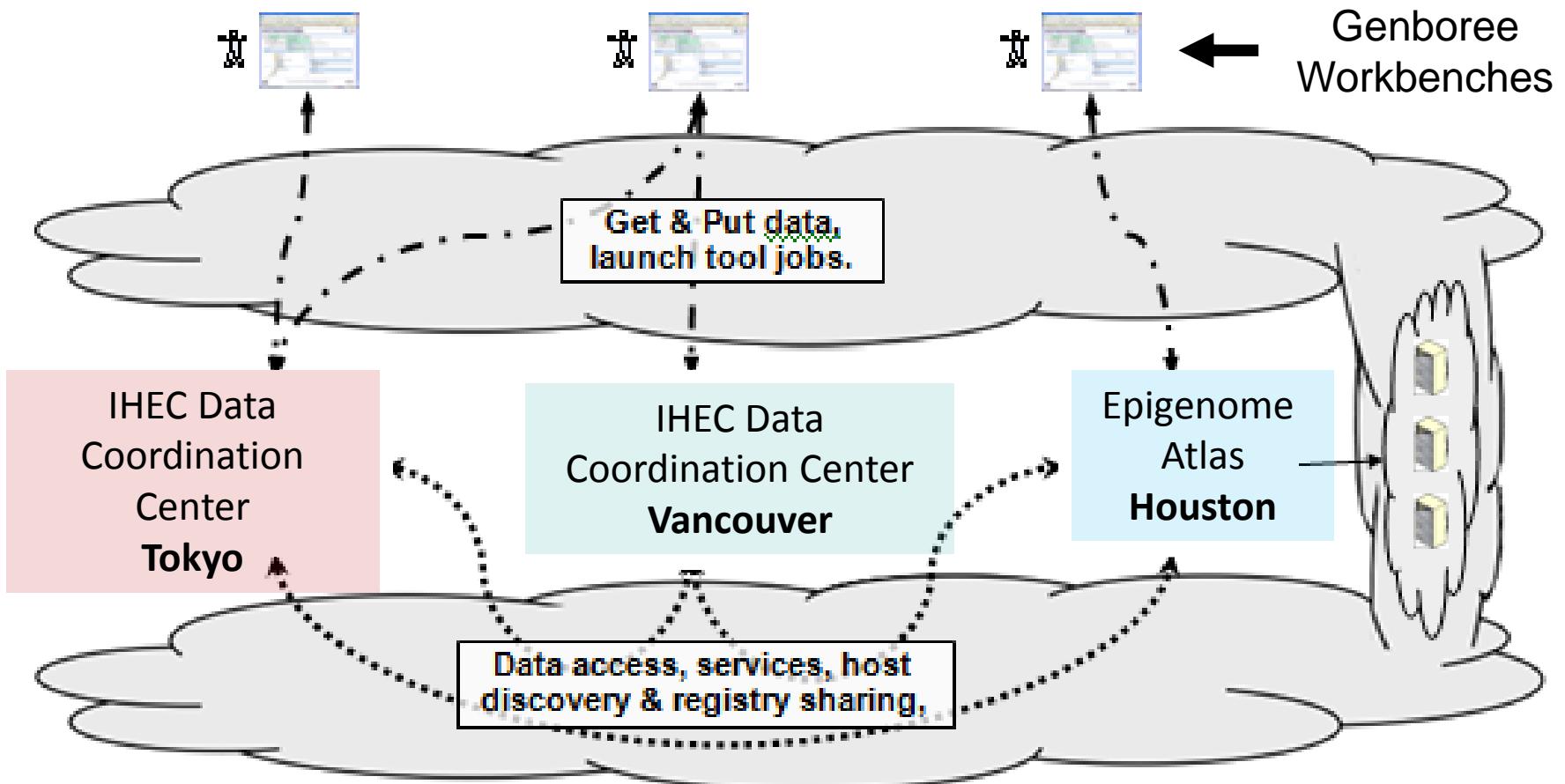
Vancouver Steven Jones



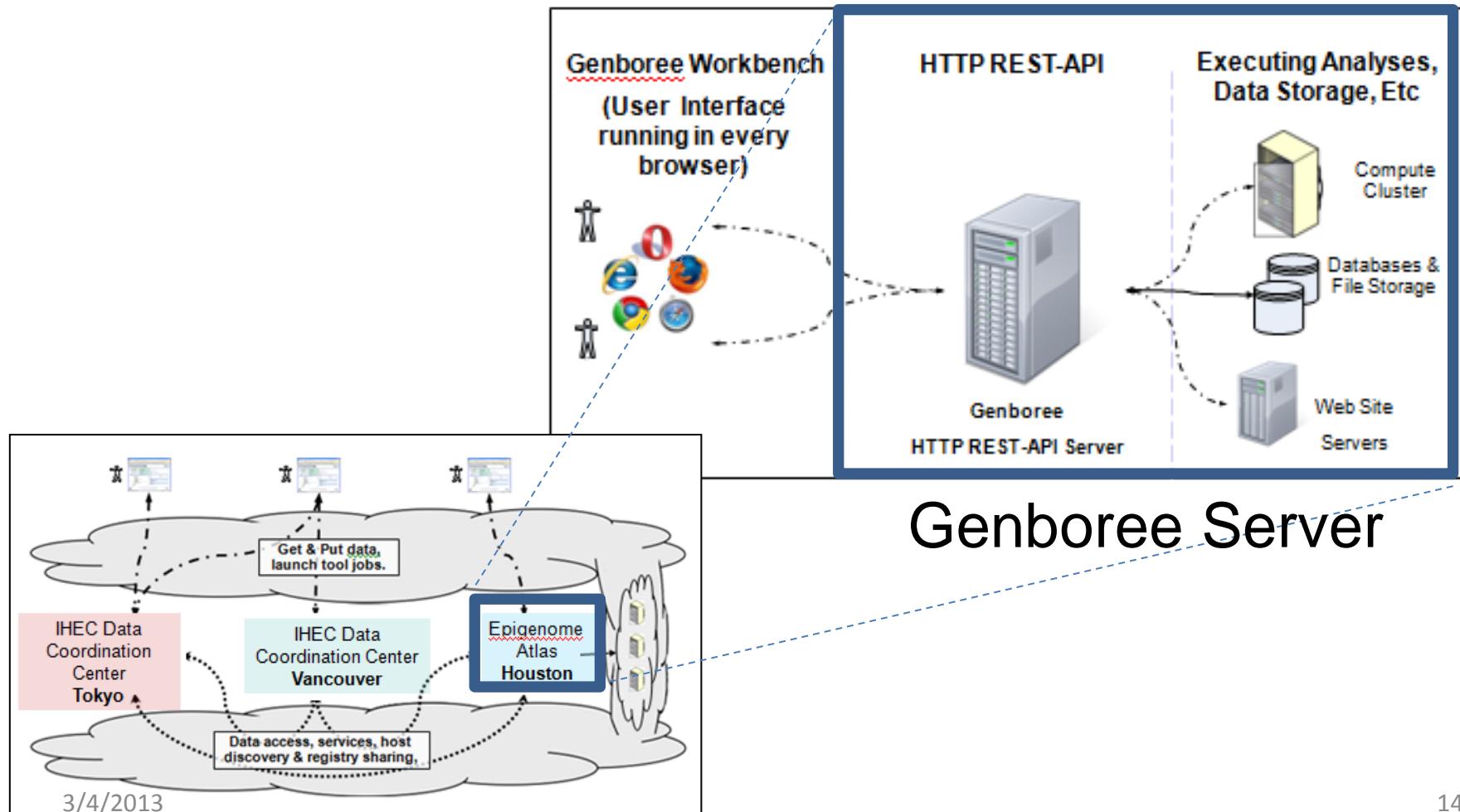
Tokyo, Toutai Mituyama

A pilot project in virtual data integration within IHEC using Genboree

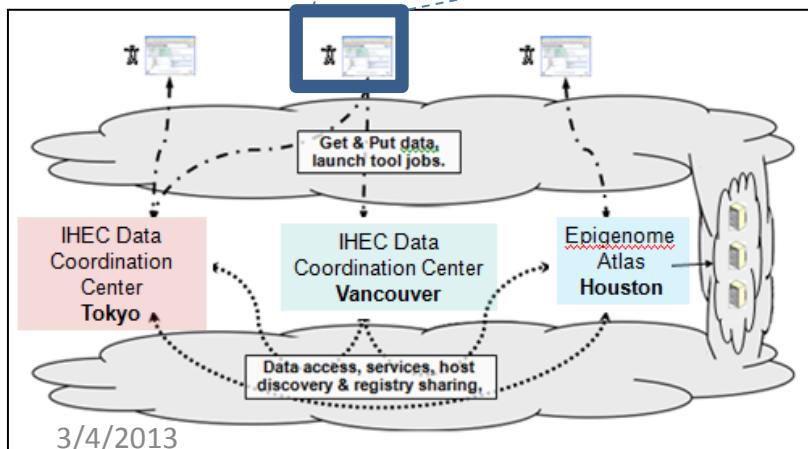
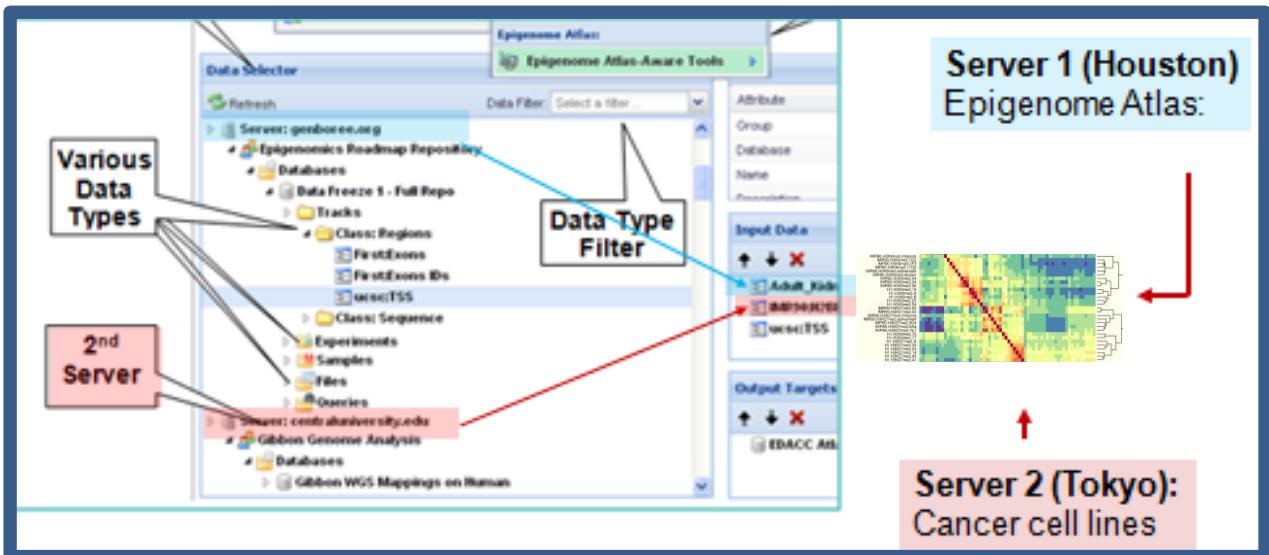
IHEC: International Human Epigenome Consortium



“Programmable Web”: REpresentational State Transfer (REST) Application Programming Interfaces (APIs)



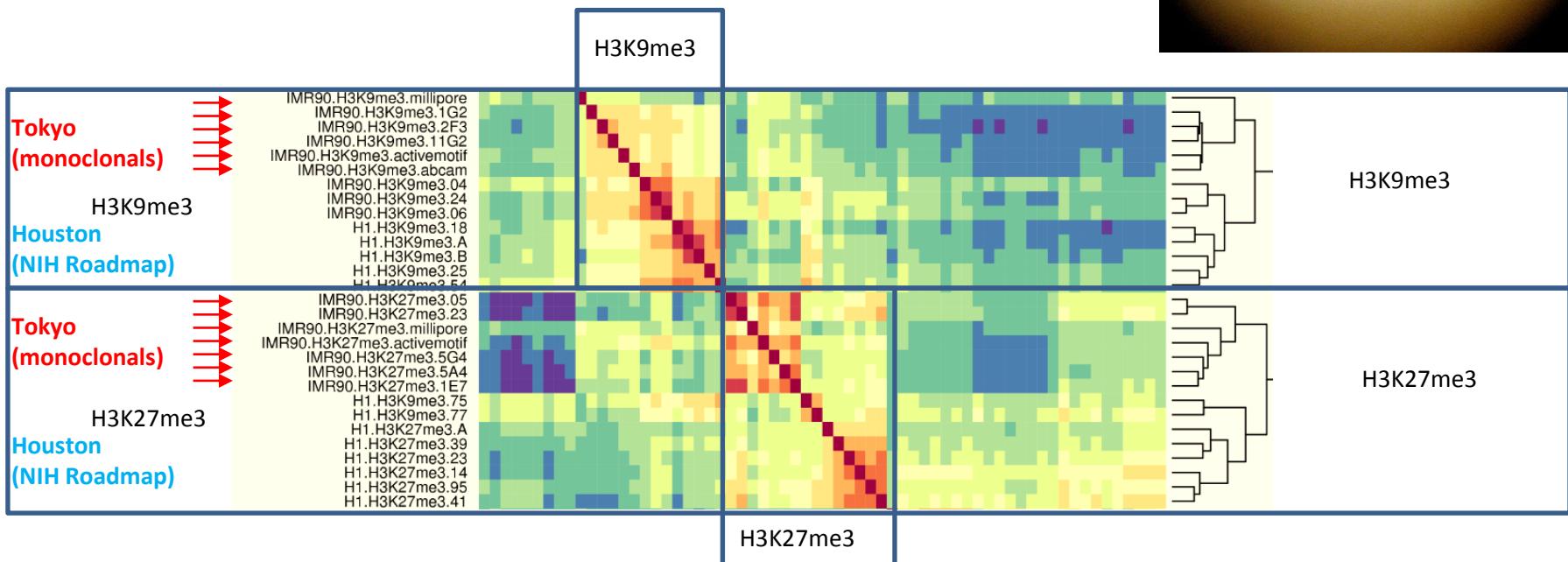
Drag-and Drop comparison of epigenomes across different Genboree servers



Genboree Workbench
(30+ Tools 50+ Utilities)

Virtual integration across IHEC DCCs

Only relevant portions of H3K9me3 and H3K27me3 signals (averages over promoters) were transferred between Tokyo and Houston.



Collaboration with

Drs. **T. Ushijima**, National Cancer Center, Tokyo, **T. Mituyama**, AIST, Tokyo, **Y. Kanai**, National Cancer Center, Tokyo, **H. Aburatani**, Univ. of Tokyo, **K. Shirahige**, Univ. of Tokyo, **Y. Wada**, Univ. of Tokyo, **H. Kimura**, Osaka Univ., **Y. Suzuki**, Univ. of Tokyo

Virtual integration principles

- Data integrated *just in time* for analysis. No data hoarding!
- Seamless integration *by “apps” accessing multiple sites* on behalf of the user.
- Data integrated *at the highest level of processing (most compressed)*.
- *Only the data relevant* for particular analysis (e.g., epigenomic marks over enhancers or promoters, not the whole genome).

Requirements for virtual data integration

- **Data standards**
(NIH Roadmap, IHEC)
- **Metadata standards**
(NIH Roadmap, IHEC)
- **Web-based interoperability**
(Genboree REST APIs)

Epigenomic data and metadata standards

www.ihec-epigenomes.org



IHEC
International Human Epigenome Consortium

[Home](#) [Areas of focus](#) [Standard Operating Procedures](#) [Tools / Useful Information](#) [Policies and Guidelines](#) [IHEC Structure](#) [Outreach and Training](#) [Feedback](#) [Intranet](#)

Links

The EPGENOME Network of Excellence
The focal point for the European epigenetics research community

Centre for Epigenetics
Centre of Excellence funded by The Danish National Research Foundation.

ROADMAP

IHEC Recommendations for Epigenomic Analysis (DRAFT)

Download

Data and Metadata Models Developed by the NIH Roadmap Epigenomics Project
(Under Consideration by the International Human Epigenome Consortium (IHEC) as a Basis for IHEC Recommendation)

These assay standards developed/in use by the NIH Roadmap Reference Epigenome Mapping Centers are being considered for potential use by IHEC.

RNA-seq Standards
ChIP-seq Standards
MethylC-seq Standards

Data and Metadata Models Developed by the NIH Roadmap Epigenomics Project
(Under Consideration by the International Human Epigenome Consortium (IHEC) as a Basis for IHEC Recommendation)

R. Alan Harris¹, Robert A. Waterland^{1,2}, Ryan Lister³, Martin Hirst⁴, Aleksandar Milosavljevic¹
¹Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas
²USDA/ARS Children's Nutrition Research Center, Department of Pediatrics, Baylor College of Medicine, Houston, Texas, USA.
³Genomic Analysis Laboratory, The Salk Institute for Biological Studies, La Jolla, California, USA.
⁴Genome Sciences Centre, BC Cancer Agency, Vancouver, British Columbia, Canada.

Genboree Installations

Other locations that use Genboree



Installation in-progress,
Targeted completion 2Q13



Installation completed,
Dec 2012



Cloud: Elastic Computing Plus Dedicated Hosting

The screenshot shows the Genboree Workbench interface. At the top, there are four tabs: System, Data, Analysis, and Visual. Below the tabs, a header says "Welcome to the Genboree W". A bulleted list provides instructions for using the Data Selector tree:

- The **Data Selector** tree on the left shows all available items.
- Drag items to be used as tool *inputs* or *outputs*.
- Drag items to be used as *output destinations*.
- Tools which can be run on your selected items.
- Unsure about what kinds of items a particular tool takes?
 - Just click the tool button when it is highlighted.

The main area is titled "Data Selector" and contains a tree view of various projects and groups. Projects listed include shortTags, SOLID-SV-CC-JR, Spark Access, Targeted Atlases, TCGA, TCGA-Reporting, testAcgh, testEDACC, Tumor Sequencing Project (TSP), Universal Probes, weilie_group, Yue, yxb4544_group, zfranco_group, and zuozhouc_group. Two specific entries are highlighted with red boxes:

- www.brain-research-lab.org
- www.microbiome-center.org

The screenshot shows the "Brain Research Lab #1" project page. The title is "Brain Research Lab #1" and the logo is "GENBOREE hosted site". A message says "You are currently not logged in." Below it, there's information about the 3rd Epigenome Informatics Workshop and links for FAQs and Use Cases. A section for "Support Site" mentions a Community Support Site where registered users can access forums and make feature requests. A login form is present with fields for "Login Name" and "Password", and buttons for "Login", "Forgot your password?", and "Guest/Public View". A link to "Register here!" is also provided. The bottom of the page has a link to "Genboree access at www.genboree.org".

The screenshot shows the "Microbiome Center #1" project page. The title is "Microbiome Center #1" and the logo is "GENBOREE hosted site". A message says "You are currently not logged in." Below it, there's information about the 3rd Epigenome Informatics Workshop and links for FAQs and Use Cases. A section for "Support Site" mentions a Community Support Site where registered users can access forums and make feature requests. A login form is present with fields for "Login Name" and "Password", and buttons for "Login", "Forgot your password?", and "Guest/Public View". A link to "Register here!" is also provided. The bottom of the page has a link to "Genboree access at www.genboree.org".

A combination of dedicated hosting and elastic cloud computing accessible via the Genboree Workbench