# GARLIC - Genomic Analysis Result Library Interface Cell

Harvard Medical School
2011 Open Source Developer Retreat

- Brief Overview Of i2b2

- Assumptions – Requirements – Challenges

- Learning From Others

- Proposed Strategy

- Data Model

- Ontologies & Data Dictionaries

- Core Technology / Mongo / BioMart / Galaxy

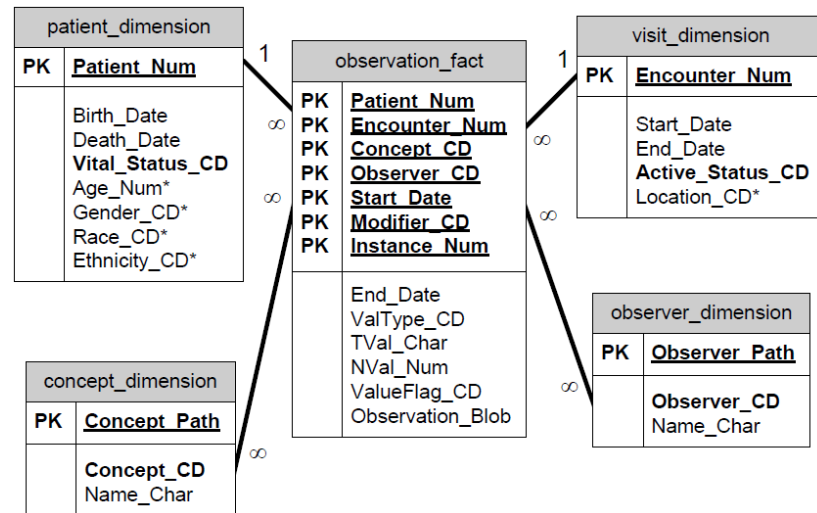- External System & Application Integration

# Brief Overview Of i2b2

> *….a scalable informatics framework that will enable clinical researchers to use existing clinical data for discovery research and, when combined with IRB-approved genomic data, facilitate the design of targeted therapies for individual patients with diseases having genetic origins.*

- Modular – "Cell"

- HTTP – XML Restful (& SOAP)

- Java 1.5
  - JBOSS 4 & AXIS

- Rich Client
  - Eclipse Frame Work

- Web Client
  - Yahoo UI Javascript

- Data Warehouse Star Schema
  - SQLServer & Oracle



https://www.i2b2.org/events/slides/i2b2_AMIA_Tutorial_20100310.pdf
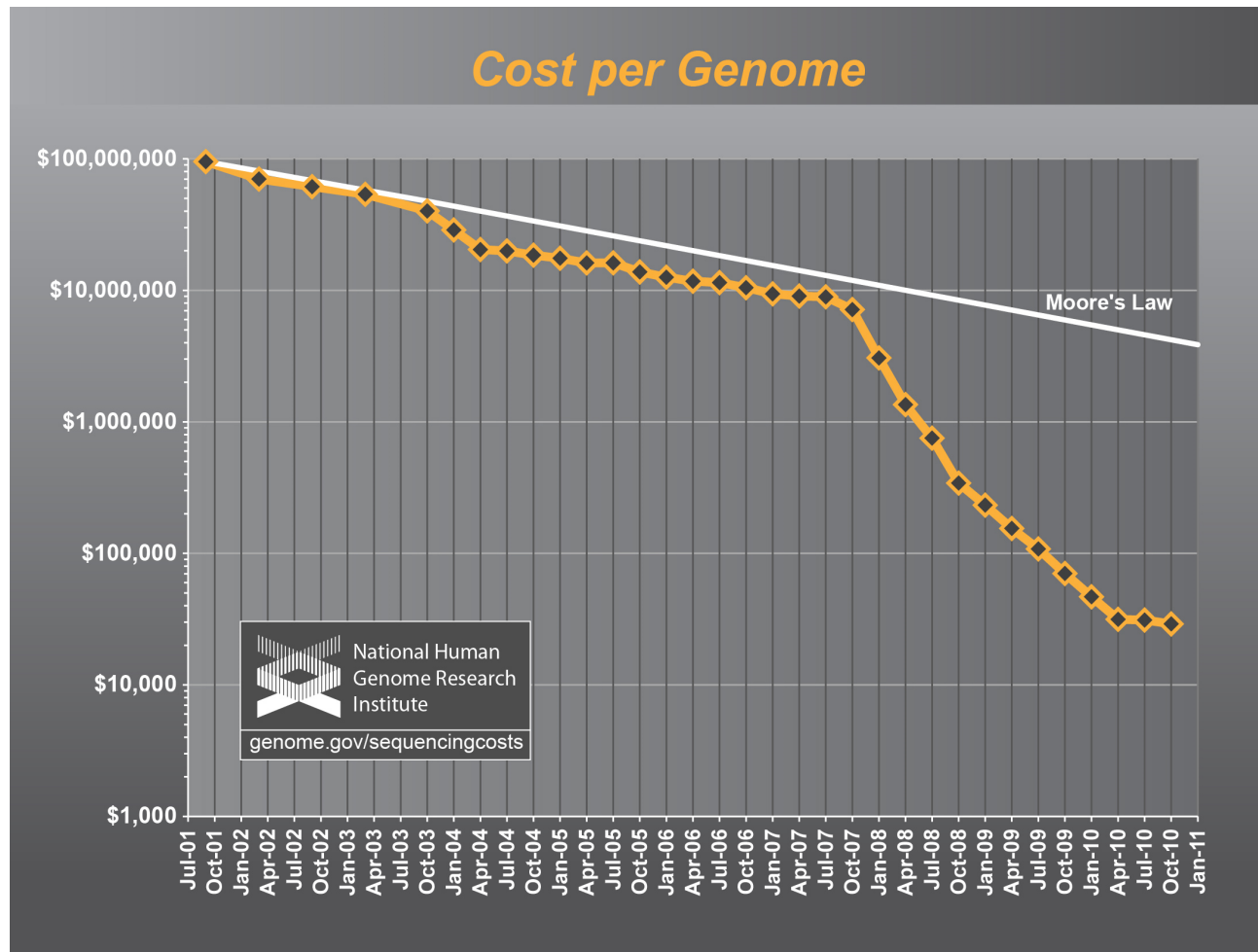
# Assumptions – Requirements - Challenges

- Initial Focus On Variant Data
  **Sequence** & Chip based

- Balancing the Needs Of Multiple Consumer / Producer Types
  Clinical Researcher, Bioinformatician, Biostatistician, Programmer...
  Variant Data From Research
  Variant Data From Clinical Practice

- Extensibility & Scalability

  – Store 2M Variants per Patient

  – Store 11 Facts Per Variant

  – Store 5000 Patients

  – Vertical Vs Horizontal scaling

- Minimize Impact On Exist i2b2 Core Software Components

Currently in the demo i2b2 project, the OBSERVATION_FACT table requires around 0.18KB per fact (row)

2M Variants x 11 Facts x 5000 patients => ~ 19TB

Utilizing OBSERVATION_BLOB field could offset this overhead but reduces data accessibility

# Predicted Demand

# The challenge……

*Distill and present the information derived from an NGS variant analysis dataset within i2b2 to enable researches to exploit the genomic knowledge it contains…*

# Learning From Others
## Clinical Researcher Perspective

i2b2
Genomic Analysis Result Library Interface Cell

Informatics for Integrating Biology & the Bedside

- The Cancer Genome Atlas
- 23andMe
- Navigenics….

# i2b2
## Genomic Analysis Result Library Interface Cell

# Learning From Others
## Domain Expert Perspective

**Integrative Genomics Viewer**

**BioDAS**

**REACTOME**

| Home | About | Content | Documentation | Tools | Download | Contact Us | Outreach |
|------|-------|---------|---------------|-------|----------|-----------|----------|

| About Reactome | Reactome Milestone |
|----------------|--------------------|

Search

Pathway Browser

Pathway Analysis

Species Comparison

Expression Analysis

If you would prefer to use our old website, click here.

REACTOME is an open-source, open access, manually curated and peer-reviewed pathway database. Pathway annotations are authored by expert biologists, in collaboration with Reactome editorial staff and cross-referenced to many bioinformatics databases. These include NCBI Entrez Gene, Ensembl and UniProt databases, the UCSC and HapMap Genome Browsers, the KEGG Compound and ChEBI small molecule databases, PubMed, and Gene Ontology. ... [more]

Reactome has achieved its milestone of curating reactions and pathways involving at least 5000 distinct human proteins... [more]

**Download**
The following links allow you to download Reactome data in

| Tutorial | Pathway of the M |
|----------|------------------|

**SequenceVariantAnalyzer**
http://www.svaproject.org/
Center for Human Genome Variation
Duke University School of Medicine

**NCBI** | Home | PubMed | GenBank | BLAST

**Homo sapiens polycystic kidney disease 1 (autosomal dominant) (PKD1), RefSeqGene on chromosome 16**
gi|209417925|ref|NG_008617.1|

Link To This Page

NG_008617.1 (56,511 bases)

Sequence | Load Accession | Set Origin | Views & Tools | Markers

1 | 2 K | 4 K | 6 K | 8 K | 10 K | 12 K | 14 K | 16 K | 18 K | 20 K | 22 K | 24 K | 26 K | 28 K | 30 K | 32 K | 34 K | 36 K | 38 K | 40 K | 42 K

5,172 : 5,272 (101 bases shown, positive strand)

Sequence | Flip Strands | Tools

5,180 | 5,190 | 5,200 | 5,210 | 5,220 | 5,230 | 5,240

Sequence NG_008617.1: Homo sapiens po
I dominant) (PKD1), RefSeqGene on chromosome 16

- Genes

- Alignments
NM_001009944.2
NM_000296.3

**Bioconductor**
OPEN SOURCE SOFTWARE FOR BIOINFOR

**BioPerl**

# Proposed Strategy

1. **Store Summarized Genomic Annotation Information Within the current OBSERVATION_FACT table**

   - Genomic Landmark-Centric e.g. Gene

2. **Separate Storage Detailed Genomic Annotation & Result Information**

   - Object Based (MongoDB?)
   - BioMart (http://www.biomart.org/) ?

3. **Store Genomic Datasets (BAM, PED etc…) Within A Secure File System – Indexed within i2b2 Data Mart**

4. **Analytics "Workflow Engine"**

   - Galaxy (http://galaxy.psu.edu)

# Component Diagram

**i2b2 Hive Core**

| PM-Cell (Authentication) | CRC-Cell (Summary Annotations) |
| R | Perl | cURL |

**Interface**

| I2b2 Hive | Domain Power Users |

| I2b2 Web Service API | Genomic Report API |

**Report Engine**

| Report Request Broker | Genomic Data Importer ( PED, GFF3 ... ) | Genomic Data Exporter ( PED, GFF3, BED, WIG ... ) |

Export Gene Level Results to CRC

I2b2-Galaxy Adaptor

**Galaxy**
- Raw Data Storage
- 'Canned' Workflow Reports

**Other Resources**

**Domain Experts**

Low Level Data Access

| Genomic Features | Data Analysis Metadata |
| Experiment Metadata | gridFS - BAM Files |

• Data Persistence
  MongoDB ?
  BioMart ?

Annovar(?)

SafeGenes (HMS) ?

# Data Model

- i2b2 Encounter = 'genomic analysis'
    - Single Patient
    - Single Sample
    - Single Assay / Platform

- Analysis Results – Annotation (GFF3/GVF)
    - Utilize The GFF3/GVF File Format For Describing Annotations
    - Summary stored within OBSERVATION_FACT Table
    - Leverage the MODIFIER_CD concept (> i2b2 1.6 RC2 )
    - Detailed Annotations Stored Within A Genomic Centric DB (MongoDB, Cassandra, CouchDB, BioMART )

- Analysis Results – Data Sets (PED, BAM, FASTA, CEL)
    - Stored Within A File System
    - Indexed Within i2b2 Observation Fact table
    - Ultimately, create a workflow that can create the GFF annotations based on the uploaded Data Sets (Nice To Have)

# Ontologies & Data Dictionaries

- ## Sequence Ontology Feature Annotation (SOFA)
  *Describe the genomic features detected*

  - http://www.sequenceontology.org/

- ## Experimental Factor Ontology (EFO)
  *Broadly describe the type of genomic assay performed*

- ## Software Ontology (SWO)
  *Track data file formats*

- ## UMLS (SNOMEDCT / HL7 Ver3)
  *Sample Type & Pathology*

- ## Human Genome Organisation (HUGO)
  *Gene Symbols*

- ## Annotation Pipelines (Annovar, Safegenes, PharmGKB)

  - Disease Gene relationship
  - Drug Gene relationship

# Core Element

- Java 1.6

- JBOSS 5

- JERSEY Restful Web Service Layer

- Oracle & SQL Server

- jQuery ( YUI plug-in development)

# Visual Example

# Accessing Detailed Results

- Simple 'Entrez' like API for Bioinformaticians

## MongoDB

- Document Oriented database written in C++

- JSON formatted Documents

- Uses Javascript

- Supports sharding

- Use MapReduce

- Language specific drivers (C, C++, Java, Perl, Python, Ruby, Scala…)


- Document ~ Row in RDBMS

- Collection ~ Table in RDBMS

# MongoVUE



## Screen content

MongoVUE

File  Server  Database  Collection  Index  JavaScript  Tools  Window  Help

Connect  View  Find  Update  Remove  Map Reduce  JavaScript

Database Explorer

- local
  - 10gen
    - Collections
      - analysis
      - gff
        - Indexes
          - _id_
    - Stored JavaScript
    - GridFS
    - Users

gff

Enter Find Criteria

{Find}  { "start": 556655 }  $Where
{Fields}  {Sort}

Tree View | Table View | Text View | Explain

Json Text of resulting Documents:

```
/* 0 */
{
  "_id": 17,
  "seqid": "chr1",
  "strand": "+",
  "score": ".",
  "type": "SNV",
  "Variant_seq": "C,T",
  "Genotype": "heterozygous",
  "frame": ".",
  "Reference_seq": "T",
  "source": "CSHL",
  "Variant_reads": "75,4",
  "start": 556655,
  "ID": "BJW-1126442",
  "end": 556655,
  "Total_reads": "79"
}
```

Learn Shell

```
[ 10:21:16 AM ]
db.gff.find({ "start" : 556655 });
db.gff.find({ "start" : 556655 }).explain();

[ 10:18:22 AM ]
db.gff.find({ "_id" : 17 });
db.gff.find({ "_id" : 17 }).explain();

[ 10:18:08 AM ]
db.gff.find({ "_id" : 20 });
db.gff.find({ "_id" : 20 }).explain();

[ 10:16:37 AM ]
db.gff.find({ "_id" : 1 });
```

## Callout boxes

```
/* 0 */
{
  "_id": 17,
  "seqid": "chr1",
  "strand": "+",
  "score": ".",
  "type": "SNV",
  "Variant_seq": "C,T",
  "Genotype": "heterozygous",
  "frame": ".",
  "Reference_seq": "T",
  "source": "CSHL",
  "Variant_reads": "75,4",
  "start": 556655,
  "ID": "BJW-1126442",
  "end": 556655,
  "Total_reads": "79"
}
```

```
db.gff.update({ "_id" : 17 }, { "$set" : {"flanking" : "AGGGGTCGT" }});

db.gff.find({ "start" : 556655 });

db.gff.find({ "_id" : 17 });
```
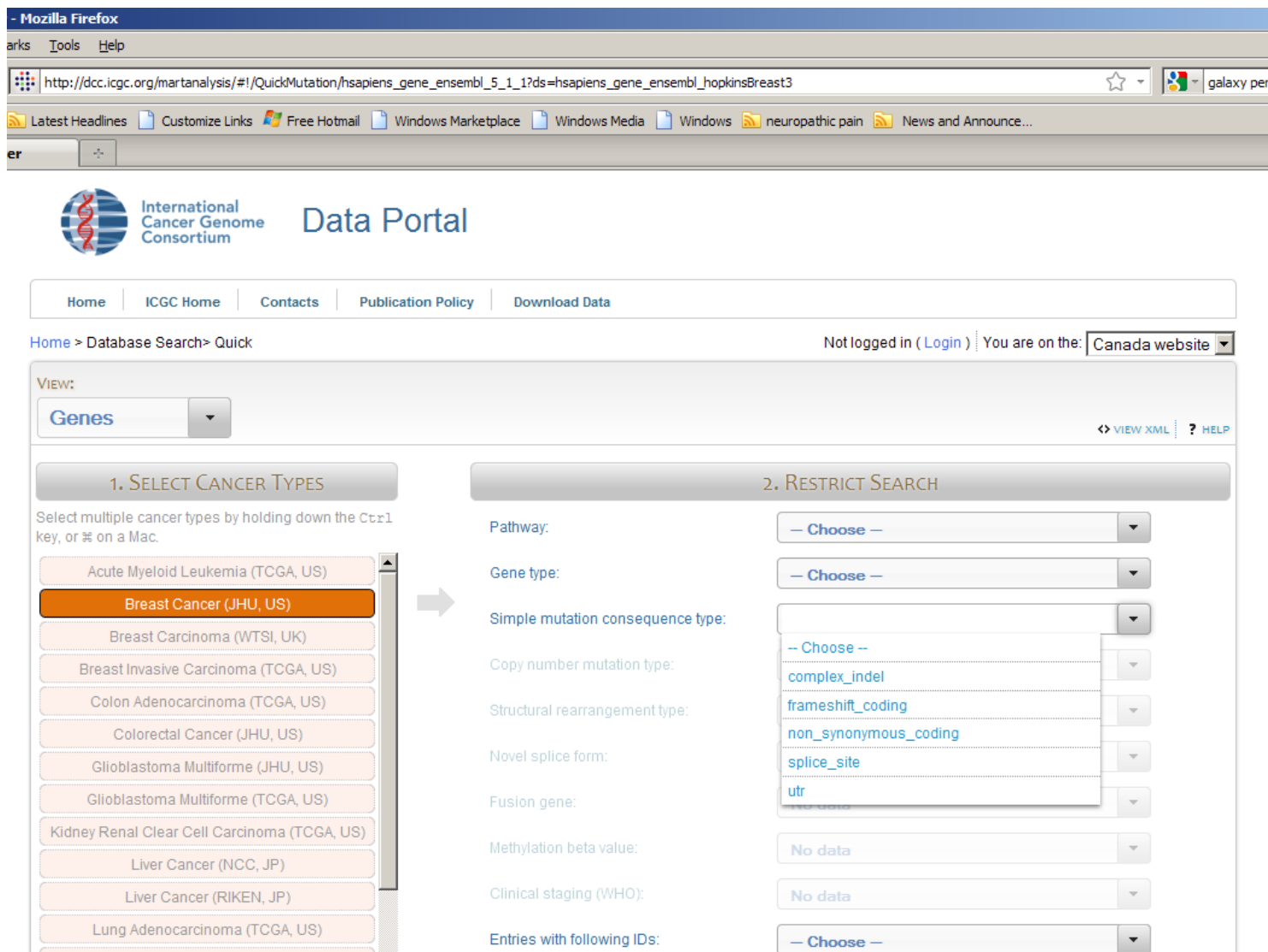
**db.<COLLECTION>.<FUNCTION>**

# BioMart

- Developed By
    - European Bioinformatics Institute (http://www.ebi.ac.uk)
    - Cold Spring Harbor Laboratory (http://www.cshl.edu)
    - Ontario Cancer Research Institute (http://www.oicr.on.ca)

- 3 Tier architecture
    - One ore more RDBMS (Oracle, MySQL, Postgres) w/ BioMart compliant schemas & multiple datasets per schema
    - APIs (Perl & Java)
    - Query interfaces
        - MartView, a web browser interface, based on the Perl API.
        - MartService, a web services interface, based on the Perl API.
        - MartURLAccess, a URL based access to MartView, based on Perl API.
        - MartExplorer, a standalone GUI tool, based on the Java API.
        - MartShell, a command-line tool, also based on the Java API.

- Distributed Annotation System ( http://www.ensembl.org/info/docs/das/index.html)
    *DAS is a specification of a protocol for requesting and returning annotation data for genomic regions. DAS allows sequence annotation to be stored in a decentralised manner, by multiple third-party annotators, and integrated on an as-needed basis by client-side software.*

# BioMart – 0.8 RC5

## Galaxy

- ## Developed By

  - Penn State Uni. (http://galaxy.psu.edu)

  - Python based analysis bioinformatic environment

*Galaxy is a framework for integrating computational tools. It allows nearly any tool that can be run from the command line to be wrapped in a structured well defined interface.*

*On top of these tools, Galaxy provides an accessible environment for interactive analysis that transparently tracks the details of analyses, a workflow system for convenient reuse, data management, sharing, publishing, and more.*



https://bitbucket.org/galaxy/galaxy-central/wiki/ISMB2010_GalaxyTutorial_3_RunningYourOwn

# Integration

- ## Genome Browsers

  - UCSC

  - Ensembl

  - NCBI

- ## Applications

  - **GALAXY** (http://galaxy.psu.edu/)
    Penn State University

  - **WGAViewer** (http://people.genome.duke.edu/~dg48/WGAViewer/whatis.php)
    Duke University School Of Medicine

  - **SequenceVariantAnalyzer** (http://www.svaproject.org/)
    Duke University School Of Medicine

  - **Cytoscape**
    w/ Reactome Plugin

  - **JalView**
    Sequence alignment visualization