

CABIO 4.2

Technical Guide



NATIONAL[®]
CANCER
INSTITUTE

Center for Biomedical Informatics
and Information Technology

CREDITS AND RESOURCES

caBIO Development and Management Teams		
Development	Technical Guide	Program Management
Hong Dang ³	Jill Hadfield ¹	Juli Klemm ^{1c}
Sharon Gaheen ²	Konrad Rokicki ²	Krishnakant Shanbhag ¹
Norval Johnson ⁴	Jim Sun ²	George Komatsoulis ¹
Sriram Kalyanasundaram ⁴	Lalitha Viswanath ²	
Ying Long ²	Sharon Gaheen ²	
Konrad Rokicki ²		
Jim Sun ²		
Lalitha Viswanath ²		
Ye Wu ²		
¹ National Cancer Institute Center for Biomedical Informatics and Information Technology (NCI CBIIT)		² Science Application International Corporation (SAIC)
³ Alpha Gamma Technologies		⁴ TerpSys

Contacts and Support	
NCICB Application Support	http://ncicbsupport.nci.nih.gov/sw/ Telephone: 301-451-4384 Toll free: 888-478-4423

LISTSERV Facilities Pertinent to caBIO		
LISTSERV	URL	Name
CABIO_USERS	https://list.nih.gov/archives/cabio_users.html	Users
CABIO_DEVELOPERS	https://list.nih.gov/archives/cabio_developers.html	Developers

Table 1:

GForge is a cross project collaboration site for NCICB caBIO developers located at <http://gforge.nci.nih.gov/projects/cabiodb/>.

Release Schedule

This guide has been updated for the caBIO 4.2 release. It may be updated between releases if errors or omissions are found. The current document refers to the 4.2 version of caBIO, released in November 2008 by the NCICB

TABLE OF CONTENTS

Credits and Resources	i
About This Guide	1
Purpose	1
Release Schedule	1
New Features in caBIO 4.2	2
Audience	2
Additional caBIO Documentation	2
Organization of This guide	3
Text Conventions Used	3
Chapter 1	
caCORE Overview	5
Architecture Overview	5
Components of caCORE	6
Enterprise Vocabulary Services (EVS)	6
Cancer Data Standards Repository (caDSR)	6
Cancer Bioinformatics Infrastructure Objects (caBIO)	6
Common Security Model (CSM)	6
Common Logging Module (CLM)	7
Chapter 2	
caBIO Overview	9
General Overview	9
caBIO Data Overview	9
Architecture Overview	10
Chapter 3	
caBIO Architecture	13
caBIO System Architecture	13
The Application Service Layer	14
The Data Source Delegation Layer	14

Object-Relational Mapping (ORM)	14
Non-ORM Layer	15
Client Technologies	15
Major caBIO Domain Packages	15
The caBIO System Packages	17

Chapter 4

Interacting with caBIO19

caBIO Graphical User Interfaces (GUIs)	19
caBIO Home Page	19
caBIO Portet	21
caBIO Service Interface Paradigm	24
Java API	25
Installation and Configuration	25
A Simple Example	27
Service Methods	28
Examples of Use	31
Utility Methods	41
Web Services API	42
Configuration	42
Java WS Client	42
Operations	43
Examples of Use	45
XML-HTTP API	47
Service Location and Syntax	47
Examples of Use	49
Working With Result Sets	49
Limitations	50
Grid Service API	50
Service Location	51
Examples of Use	51
Working With Result Sets	52
Limitations	52

Chapter 5

caBIO Utilities53

FreestyleLM	53
Freestyle LM Java Client	54
Freestyle LM Web Interface	55
Grid Id Resolver	56
Grid Id Resolver Java API	56

Genome Range Queries	56
Range Query Java API	56
Array Annotation API	58
SVG Manipulator	63
Manipulating Pathway Diagrams	63
Chapter 6	
Cancer Bioinformatics Infrastructure Objects	67
Introduction	67
caBIO API	67
Data Sources in the caBIO Database	81
caGrid Identifiers	87
Classes to be indexed for FreeStyleLM Search	88
<u>Appendix</u>	<u>93</u>
Appendix A	
References	93
Technical Manuals/Articles	93
Scientific Publications	94
caBIG® Material	95
caCORE Material	96
Modeling Concepts	96
Applications Currently Using caBIO	96
Index	97

ABOUT THIS GUIDE

This section introduces you to the caBIO 4.2 Technical Guide. Topics in this section include:

- *Purpose* on this page
- *Release Schedule* on this page
- *New Features in caBIO 4.2* on page 2
- *Audience* on page 2
- *Additional caBIO Documentation* on page 2
- *Organization of This guide* on page 3
- *Text Conventions Used* on page 3

Purpose

The caBIO 4.2 Technical Guide describes the Cancer Bioinformatics Infrastructure Objects (caBIO) model, the data represented by its objects, and its application programming interfaces.

This guide describes:

- the purpose, architecture and components of caBIO
- the APIs for accessing the caBIO system including Java, Web services, and XML-HTTP
- Additional utilities for accessing the data in caBIO
- the data accessible through the caBIO APIs Audience.

Release Schedule

This guide is updated for each caBIO release. It may be updated between releases if errors or omissions are found. The current document refers to the 4.2 version of caBIO, which was released in November 2008 by NCI CBIIT.

New Features in caBIO 4.2

The primary feature of the caBIO 4.2 release is an enhanced pathway object model supporting data from the [Pathway Interaction Database \(PID\)](#).

The Pathway Interaction Database is a highly-structured, curated collection of information about known biomolecular interactions and key cellular processes assembled into signaling pathways. It is a collaborative project between the US National Cancer Institute (NCI) and Nature Publishing Group (NPG), and is an open access online resource. The PID exposes data regarding pathways, complex components and molecular interactions from NCI-Nature-curated dataset, Reactome and BioCarta.

An additional feature implemented during this release cycle involves the inclusion of a caBIO “Portlet” in the caGrid Portal. A Portlet is a pluggable user interface component that operates within a web portal. A Portlet leveraging caBIO simple and templated based search facilities was implemented in this release providing the user community with access to the caBIO FreeStyleLM search feature as well as pre-defined query templates supporting:

- Genome Range Queries
- Microarray Annotations
- Genomic Annotations
- Pathways
- Cancer Gene Index

The caBIO Portlet is now available on the Community Tab of the caGrid portal.

Audience

The primary audiences of this guide are:

- The application developer who wants to learn about the architecture and use of the caBIO APIs and utilities. These sections of the guide assume familiarity with the Java programming language and/or other programming languages, database concepts, and the Internet. If you intend to use caCORE resources in software applications, it assumes that you have experience with building and using complex data systems.
- The scientist end-user who wants to understand the data provided in caBIO (*Chapter 2*) and/or use the Freestyle Lexical Mine search (described in *Chapter 5*).

Additional caBIO Documentation

- The caBIO 4.2 Release Notes contain a description of the end user tool enhancements and bug fixes included in this release.
- The caBIO 4.2 JavaDocs contain the current caBIO API specification

- The caCORE SDK 4.0 Developer's Guide contains detailed instructions on the use of the SDK and how it helps create a caCORE-like software system. See <http://ncicb.nci.nih.gov/NCICB/infrastructure/cacoresdk#Documentation>.
- For more caBIO-related references, see *Appendix A*.

Organization of This guide

This brief overview explains what you will find in each section of this guide.

- This chapter introduces this technical guide.
- Chapter 1 provides an overview of caCORE and caBIG®.
- Chapter 2 provides an overview of caBIO.
- Chapter 3 describes the architecture of caBIO
- Chapter 4 describes the various methods for accessing the data in caBIO.
- Chapter 5 describes additional utilities provided by the caBIO API, which are not generated by the caCORE SDK.
- Chapter 6 describes the Cancer Bioinformatics Infrastructure Objects (caBIO)
- Appendix A includes a listing of references to provide the caBIO users with more in depth information about relevant topics.
- Appendix B provides technical users with instructions how migrate systems that are using previous versions of caBIO.

Text Conventions Used

This section explains conventions used in this guide. The various typefaces represent interface components, keyboard shortcuts, toolbar buttons, dialog box options, and text that you type.

Convention	Description	Example
Bold	Highlights names of option buttons, check boxes, drop-down menus, menu commands, command buttons, or icons.	Click Search .
<u>URL</u>	Indicates a Web address.	http://domain.com
text in SMALL CAPS	Indicates a keyboard shortcut.	Press ENTER.
text in SMALL CAPS + text in SMALL CAPS	Indicates keys that are pressed simultaneously.	Press SHIFT + CTRL.
<i>Italics</i>	Highlights references to other documents, sections, figures, and tables.	See <i>Figure 4.5</i> .
<i>Italic boldface monospaced type</i>	Represents text that you type.	In the New Subset text box, enter <i>Proprietary Proteins</i> .

Convention	Description	Example
Note:	Highlights information of particular importance	Note: This concept is used throughout the document.
{ }	Surrounds replaceable items.	Replace {last name, first name} with the Principal Investigator's name.

CHAPTER 1

CACORE OVERVIEW

This chapter provides an overview of the NCICB caCORE infrastructure.

Topics in this chapter include:

- Architecture Overview on this page
- Components of caCORE on page 6

Architecture Overview

The NCI Center for Bioinformatics (NCICB) provides biomedical informatics support and integration capabilities to the cancer research community. NCICB has created a core infrastructure called Cancer Common Ontologic Representation Environment (caCORE), a data management framework designed for researchers who need to be able to navigate through a large number of data sources. By providing a common data management framework, caCORE helps streamline the informatics development throughout academic, government and private research labs and clinics. The components of caCORE support the semantic consistency, clarity, and comparability of biomedical research data and information. caCORE is open-source enterprise architecture for NCI-supported research information systems, built using formal techniques from the software engineering and computer science communities. The four characteristics of caCORE include:

- Model Driven Architecture (MDA)
- *n*-tier architecture with open Application Programming Interfaces (APIs)
- Use of controlled vocabularies, wherever possible
- Registered metadata

The use of MDA and *n*-tier architecture, both standard software engineering practices, allows for easy access to data, particularly by other applications. The use of controlled vocabularies and registered metadata, less common in conventional software practices, requires specialized tools, generally unavailable.

As a result, the NCICB (in cooperation with the NCI Office of Communications) has developed the Enterprise Vocabulary Services (EVS) system to supply controlled vocabularies, and the Cancer Data Standards Repository (caDSR) to provide a dynamic metadata registry. When all four development principles are addressed, the resulting system has several desirable properties. Systems with these properties are said to be “caCORE-like”.

1. The n -tier architecture with its open APIs frees the end user (whether human or machine) from needing to understand the implementation details of the underlying data system to retrieve information.
2. The maintainer of the resource can move the data or change implementation details (Relational Database Management System, and so forth) without affecting the ability of remote systems to access the data.

Most importantly, the system is ‘semantically interoperable’; that is, there exists runtime-retrievable information that can provide an explicit definition and complete data characteristics for each object and attribute that can be supplied by the data system.

Components of caCORE

The components that comprise caCORE are EVS, caDSR, caBIO, CSM, and CLM. Each is described briefly below.

Enterprise Vocabulary Services (EVS)

EVS provides controlled vocabulary resources that support the life sciences domain, implemented in a description logics framework. EVS vocabularies provide the semantic ‘raw material’ from which data elements, classes, and objects are constructed.

Cancer Data Standards Repository (caDSR)

The caDSR is a metadata registry, based upon the ISO/IEC 11179 standard, used to register the descriptive information needed to render cancer research data reusable and interoperable. The caBIO, EVS, and caDSR data classes are registered in the caDSR, as are the data elements on NCI-sponsored clinical trials case report forms.

Cancer Bioinformatics Infrastructure Objects (caBIO)

The caBIO model and architecture is the primary programmatic interface to caCORE. Each of the caBIO domain objects represents an entity found in biomedical research. Unified Modeling Language™ (UML) models of biomedical objects are implemented in Java as middleware connected to various cancer research databases to facilitate data integration and consistent representation. Examining the relationships between these objects can reveal biomedical knowledge that was previously buried in the various primary data sources.

Common Security Model (CSM)

CSM provides a flexible solution for application security and access control with three main functions:

- Authentication to validate and verify a user's credentials
- Authorization to grant or deny access to data, methods, and objects

- User Authorization Provisioning to allow an administrator to create and assign authorization roles and privileges.

Common Logging Module (CLM)

CLM provides a separate service under caCORE for Audit and Logging Capabilities. It also comes with a web based locator tool. It can be used by a client application directly, without the application using any other components like CSM.

In September of 2007, the NCI CBIIT Infrastructure and Product Management Team made the decision to separate the caCORE Components that had previously been bundled and released together. This decision was geared toward allowing each of the infrastructure product teams to be more responsive in addressing specific needs of the user community. With each component release there is a product specific Technical Guide. For more details on the other components, refer to the caCORE Overview page at http://ncicb.nci.nih.gov:80/NCICB/infrastructure/cacore_overview, which directs you to other product specific Technical Guides.

CHAPTER 2

CABIO OVERVIEW

This chapter provides an overview of the caBIO data and infrastructure.

Topics in this chapter include:

- General Overview on this page
- caBIO Data Overview on this page
- Architecture Overview on page 10

General Overview

The cancer Bioinformatics Infrastructure Objects (caBIO) model and architecture is a synthesis of software, vocabulary, and metadata models for cancer research. Each of the caBIO domain objects represents an entity found in biomedical research such as Gene, Chromosome, Nucleic Acid Sequence, SNP, Library, Clone, and Pathway. Examining the relationships between these objects can reveal biomedical knowledge that was previously buried in the various primary data sources. Given the dynamic nature of this information, the data in caBIO is updated on a monthly basis through a series of ETL (Extract, Transform, and Load) processes.

The caBIO data can be accessed via a variety of Application Programming Interfaces (APIs) as well as a web user interface. The system also provides a number of utilities to search for and manipulate data objects programmatically.

caBIO Data Overview

The caBIO system is primarily a repository of molecular biology data. As such, the entities that concern the central dogma of molecular biology are the core of this model. Data for these classes come from a variety of sources. Gene, Chromosome, Nucleic Acid Sequence, and Taxon come from the NCBI UniGene data resource, both directly and via the Cancer Gene Anatomy Project (CGAP), an NCI activity that maintains information about expression of genes in tumors based on EST sequencing frequency.

Cytogenetic location information, as well as EST and mRNA sequence annotation data come from the University of California Santa Cruz (UCSC) Genome Project. Protein data (including data about sequences) is retrieved from the Protein Information Resource (PIR) at Georgetown University.

Pathway data in caBIO comes from the Pathway Interaction Database (PID), a highly-structured, curated collection of information about known biomolecular interactions and key cellular processes assembled into signaling pathways. The PID is a collaborative effort between the National Cancer Institute (NCI) and Nature Publishing Group (NPG), and is an open access online resource. The caBIO domain model provides classes representing PID components and interactions.

Annotations from the most widely used human microarrays from Affymetrix (<http://www.affymetrix.com>), Illumina (<http://www.illumina.com>), and Agilent (<http://www.agilent.com>) are stored in caBIO, to allow the rich annotations provided by the manufacturers to be integrated with the other caBIO data sources and served through a caCORE interface. Data for these microarray annotations is distributed across many of the caBIO classes; the core classes are Microarray and ArrayReporter with subclasses SNPArrayReporter, ExpressionArrayReporter, and TranscriptArrayReporter.

Single Nucleotide Polymorphism data are primarily supplied by dbSNP and by Affymetrix (<http://www.affymetrix.com>). The core classes are SNP, PopulationFrequency, and GeneRelativeLocation.

The caBIO system provides information about libraries used as part of EST sequencing efforts. This set of classes (which includes Library, Tissue, and Protocol) are populated with data from the NCI CGAP initiative. This group of classes provides an instructive example about the necessity of grounding metadata in concept-based controlled terminology. The Protocol object here represents the method used to create a library from a tissue (NCI Thesaurus concept code C41111) rather than the concept of a clinical trial protocol (C25320).

Clinical trial data in the caBIO system are obtained from the Cancer Therapeutics Evaluation Program (CTEP) of the NCI. This part of the model was originally driven by use cases from the NCI CMAP project, and as such, the classes contain top level about ongoing clinical trials. In addition, the caBIO system provides some very basic information on the mechanism of action of drug compounds, variations found in disease and the genetic targets of drug compounds.

Architecture Overview

The caBIO system is an n-tier client-server architecture. The caBIO application server, located at NCI (<http://cabioapi.nci.nih.gov/cabio41>), provides a home page for web-based searching and data retrieval, as well as access to a wide variety of APIs for programmatic access.

The caBIO server interacts with a database backend on behalf of the user to serve requested objects to the client. Hibernate-based object-relational mapping (ORM) is used to populated caBIO domain objects. Only object-oriented APIs are exposed to clients so that the relational database backend is completely transparent to the user.

Many of the APIs rely on a query-by-example mechanism that takes a partially-populated object in, and returns other objects like it back to the user. This paradigm of interaction makes semantic interoperability much easier.

The following APIs are available for programmatic data access:

- Java API
- Web Services API
- REST API
- Grid Service (The current grid service is for caBIO 3.2 and 4.0.)

The system also provides several additional utilities for accessing the data in caBIO:

- Freestyle Lexical Mine (FreestyleLM) search
- Genome Range Queries
- Grid ID Resolver
- Array Annotation APIs
- XML Marshalling Utility for caBIO domain objects
- Pathway image manipulation utility

CHAPTER 3

CABIO ARCHITECTURE

This chapter describes the architecture of the caBIO system. It includes information about the major components, such as database object-relational mappings (ORM), client-server communication, and how the system connects to non-ORM sources, such as FreestyleLM. It also describes the layout of the caBIO system packages. Topics in this chapter include:

- caBIO System Architecture on this page
- Client Technologies on page 14
- Major caBIO Domain Packages on page 15
- The caBIO System Packages on page 16

caBIO System Architecture

The caBIO infrastructure exhibits an n -tiered architecture with client interfaces, server components, backend objects, data sources, and additional backend systems (*Figure 3.1*). This n -tiered system divides tasks or requests among different servers and data stores. This isolates the client from the details of where and how data is retrieved from different data stores. The system also performs common tasks such as logging and provides a level of security.

Clients (browsers, applications) receive information from backend objects. Java applications also communicate with backend objects via domain objects packaged within `cabio41-beans.jar`. Non-Java applications can communicate via SOAP (Simple Object Access Protocol). Back-end objects communicate directly with data sources, either relational databases (using Hibernate) or non-relational systems.

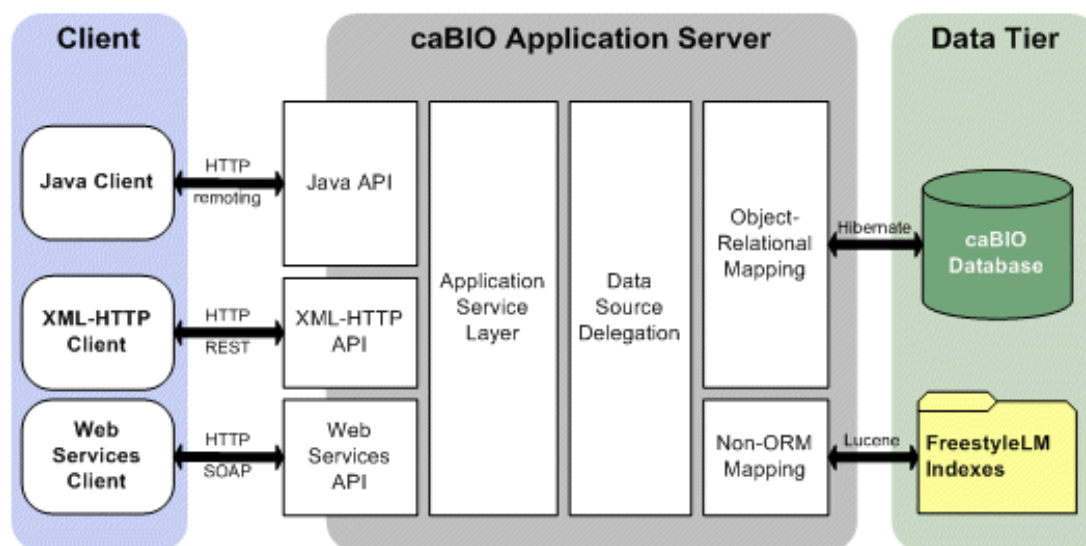


Figure 3.1 caBIO architecture

Most of the caBIO infrastructure is written in the Java programming language and leverages reusable, third-party components.

The infrastructure is composed of the following layers.

The Application Service Layer

The Application Service layer consolidates incoming requests from various interfaces and translates them to native query requests that are then passed to the data layers. This layer is also responsible for handling client authentication and access control using the Java API. (This feature is currently disabled for the caBIO system running at NCICB; all interfaces provide full, anonymous read-only access to all data.)

The Data Source Delegation Layer

This layer is responsible for conveying each query that it receives to the respective data source that can perform the query. The presence of this layer enables multiple data sources to be exposed by a single running instance of a caBIO server.

Object-Relational Mapping (ORM)

The ORM mapping is implemented using Hibernate. Hibernate is a high performance object/relational persistence and query service for Java. Hibernate provides the ability to develop persistent classes following common object-oriented (OO) design methodologies such as association, inheritance, polymorphism, and composition.

The Hibernate Query Language (HQL), designed as a "minimal" object-oriented extension to SQL, provides a bridge between the object and relational databases. Hibernate allows for real world modeling of biological entities without creating complete SQL-based queries to represent them.

Non-ORM Layer

The caBIO database is a relational database. A Lucene-based document index is generated on the backend to facilitate the Lucene full-text search feature. The caBIO API implements the non-ORM based Lucene search as well as the ORM based Hibernate Lucene Search.

Client Technologies

Applications using the Java programming language can access caBIO directly through the domain objects provided by the `cabio41-beans.jar` (see *caBIO System Architecture* on page 13.) The network details of the communication to the caCORE servers are abstracted away from the developer. Hence, developers need not deal with issues such as network and database communication, but can instead concentrate on the biological problem domain.

The caBIO system also allows non-Java applications to use SOAP clients to interface with caBIO web services. SOAP is a lightweight XML-based protocol for the exchange of information in a decentralized, distributed environment. It consists of an envelope that describes the message and a framework for message transport. caBIO uses the open source Apache Axis package to provide SOAP-based web services to users. This allows other languages, such as Python or Perl to communicate with caBIO objects in a straightforward manner.

The caBIO architecture includes a presentation layer that uses a J2SE application server (such as Tomcat or JBoss). The JSPs (Java Server Pages) are web pages with Java embedded in the HTML to incorporate dynamic content in the page. caBIO also employs Java Servlets, which are server-side Java programs that web servers can run to generate content in response to client requests. All logic implemented by the presentation layer uses Java Beans, which are reusable software components that work with Java. All caBIO domain objects can be transformed into XML, the eXtensible Markup Language, as a universal format for structured data on the Web.

Communication between the client interfaces and the server components occurs over the Internet using the HTTP protocol. The server components are deployed in a web application container as a .war (Web archive) file that communicates with the back-end relational database management system that contains the actual data.

By using XSL/XSLT, the extensible stylesheet language for expressing stylesheets and XSL Transformations (XSLT), as a language for transforming XML documents, non-programmers can transform the information in the caBIO objects for use in documents or other systems.

Major caBIO Domain Packages

Table 3.1 shows the major caBIO domain packages from which you can access the Java interfaces and classes. All of the objects in the domain packages implement the

java.io.Serializable interface. To view the JavaDocs page for each package, go to <http://cabioapi.nci.nih.gov/cabio41/docs>.

Package	Description
caBIO	Contains the domain-specific classes in the caBIO object package such as Gene, Chromosome, etc. For a list of the domain objects and their descriptions, see <i>Chapter 6, Cancer Bioinformatics Infrastructure Objects</i> , on page 67. <code>gov.nih.nci.cabio.domain</code>
Common	Contains a single class, DatabaseCrossReference, which is used by the caBIO domain package objects to provide links to related data hosted by other sources. Further details about this class are available in the next section and in <i>Chapter 6, Cancer Bioinformatics Infrastructure Objects</i> , on page 67. <code>gov.nih.nci.common.domain</code>
Provenance	Contains classes that provide provenance (source) metadata about the origin of data in caBIO objects. Further details about the classes in this package are available in <i>Chapter 6, Cancer Bioinformatics Infrastructure Objects</i> , on page 67. <code>gov.nih.nci.common.provenance.domain</code>

Table 3.1 Major caBIO domain packages

The caBIO System Packages

Figure 3.2 displays a package for each application area in the caBIO API, described in Table 3.2

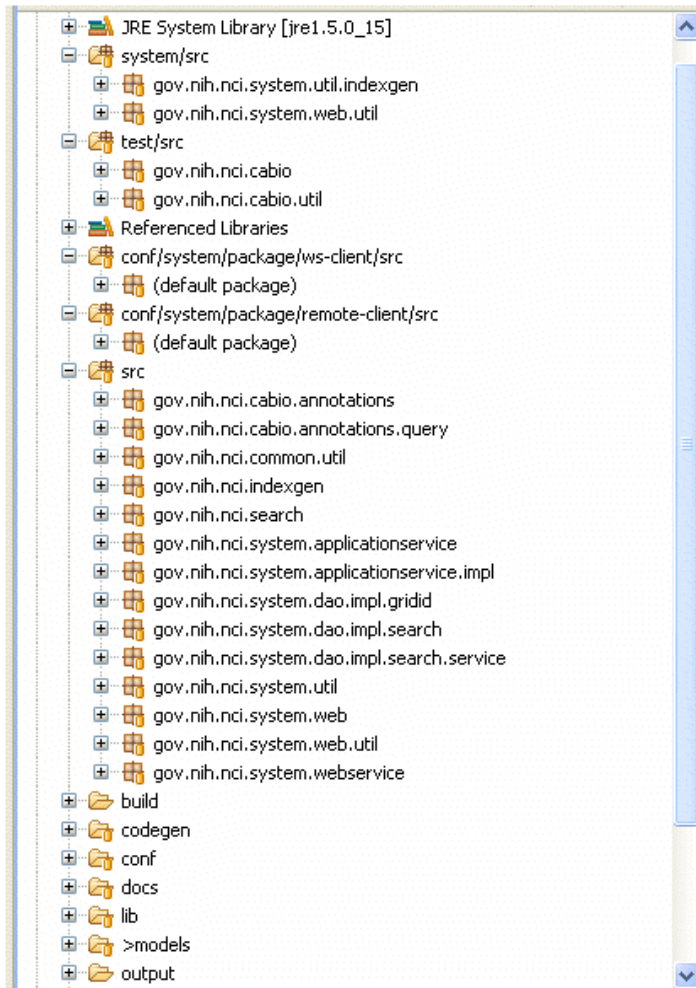


Figure 3.2 caBIO API package

In addition to domain packages, the caBIO client API includes the following framework packages:

Package	Description
<code>gov.nih.nci.common</code>	The common package contains utility classes for manipulating caBIO domain objects.
<code>gov.nih.nci.indexgen</code>	The indexgen package contains the classes that generate the Lucene indexes for the caBIO database.

Table 3.2 caBIO API framework client packages


Package	Description
<code>gov.nih.nci.search</code>	The search package contains domain objects for the FreestyleLM search component and Grid Ids. Clients can instantiate these objects and pass them to appropriate API methods to search the FreestyleLM indexes, Range Queries or the Grid Ids. For the range queries, it includes subclasses to search on a specific assembly of a genome as well as to search for locations relative to a genomic feature (defined by the subclass). It is possible to specify distances between, upstream and downstream of the feature.
<code>gov.nih.nci.system</code>	The system package has contains the client-side API, which communicates with the caBIO server, and helper classes for creating client-side proxies.
<code>gov.nih.nci.system.applicationService</code>	The ApplicationService package contains classes that extend the SDK-provided functionality to provide caBIO-specific functionalities such as searching for a caBIO domain objects by passing Grid Id, RangeQuery or FreestyleLM SearchQuery object.
<code>gov.nih.nci.system.client</code>	The client package consists of classes that are used to generate dynamic proxies. The Application Service proxy classes intercept all the calls to the actual ApplicationService and takes action to facilitate the lazy-loading mechanism for the caBIO API.
<code>gov.nih.nci.system.dao (Data Access)</code>	The data access package is the layer at which the query is executed. This layer has implementation for querying ORM and non-ORM based data stores and as with the applicationServices, provides functionality to perform Grid Id, Range Query and FreestyleLM searches
<code>gov.nih.nci.system.webservice</code>	The Web service package contains the Web service wrapper class that uses Apache's Axis and has been extended to provide caBIO-specific functionality like searching for domain objects using grid Ids.

Table 3.2 *caBIO API framework client packages (Continued)*

INTERACTING WITH caBIO

caBIO provides facilities enabling both end user interaction via Graphical User Interfaces (GUIs) and system interaction via the caBIO Service Interface Layer. This chapter describes the various methods for interacting with caBIO and provides examples of how to use each of the interface APIs.

Topics in this chapter include:

- *caBIO Graphical User Interfaces (GUIs)* on this page
- *caBIO Service Interface Paradigm* on page 24
- *Java API* on page 25 
- *Web Services API* on page 42g
- *XML-HTTP API* on page 47
- *Grid Service API* on page 50

caBIO Graphical User Interfaces (GUIs)

caBIO provides two primary graphical user interfaces (GUIs) for accessing data:

- caBIO Home Page <http://cabioapi.nci.nih.gov/cabio42/Home.action>
- caBIO Portlet

caBIO Home Page

The caBIO Home Page (*Figure 4.1*) is a web-based GUI that provides a Query-by-Example (QBE)-style interface for browsing the caBIO object model. The caBIO home page is auto-generated using the caCORE SDK and customized to support unique caBIO features like the FreeStyleLexical Mine (LM) search feature. The caBIO Home Page allows end users to traverse the caBIO object model and perform queries based

on search criteria defined by the attributes in each class caBIO Service Interface Paradigm. .

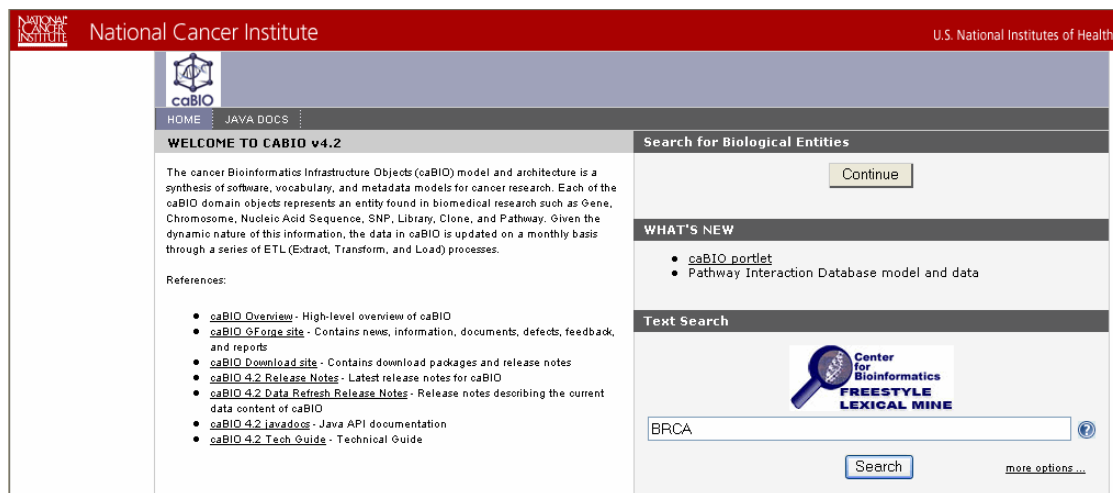


Figure 4.1 caBIO Home Page

The caBIO Home Page provides access to the FreeStyle Lexical Mine search, a “Google-like” search facility for querying caBIO data. Figure 4.2 illustrates the results of the FreeStyleLM search.

Class	sentence	celllineStatus
1.gov.nih.nci.cabio.domain.Evidence	We conclude that, despite their potential functional as sociation with the BRCA gene products, RAD51, RAD 52, and RAD54 are not themselves targeted by mutati ons in human breast cancer.	no
2.gov.nih.nci.cabio.domain.Evidence	We conclude that, despite their potential functional as sociation with the BRCA gene products, RAD51, RAD 52, and RAD54 are not themselves targeted by mutati ons in human breast cancer.	no
3.gov.nih.nci.cabio.domain.Evidence	The majority of BRCA1 and BRCA2 tumors showed a copy number reduction (relative to reference probe wit h ratios 4:2, 3:2, 4:3) at corresponding loci, suggestin g that a specific physical deletion of the wild-type BR CA gene allele has been followed by a duplication of t he remaining mutant allele via polyploidization.	no
4.gov.nih.nci.cabio.domain.Evidence	Familial breast cancer can be subdivided into three gr oups: those that can be accounted for by mutations i n the breast cancer susceptibility gene BRCA 1, famil ies harbouring mutations in BRCA 2 and families wher e neither BRCA 1 nor BRCA 2 appear to be involved.	no

Figure 4.2 FreestyleLM search results

Users can perform caBIO queries via the QBE interface illustrated in *Figure 4.3*. Users can traverse the caBIO object model and enter search criteria for each object attribute

Domain Class Browser

Please click on any of the tree nodes.

To view the search criteria for a class, expand a package listed below and select a class. To search for records, provide valid search criteria and click the Submit button. For any date attributes, please use the syntax: mm-dd-yyyy.

- Evidence
- EvidenceCode
- Exon
- ExonArrayReporter
- ExpressedSequenceTag
- ExpressionArrayReporter
- Gene**
- GeneAgentAssociation
- GeneAlias
- GeneCytogeneticLocation
- GeneDiseaseAssociation
- GeneFunctionAssociation
- GeneOntology
- GeneOntologyRelationship
- GenePhysicalLocation
- GeneRelationship

gov.nih.nci.cabio.domain.Gene

bigid:

clusterid:

fullName:

hugoSymbol:

id:

symbol:

Search Object:

Figure 4.3 QBE search interface

Search results for the QBE search are displayed in tabular form (*Figure 4.4*). Links are provided to allow users to retrieve additional information from associated objects.



National Cancer Institute



Criteria: gov.nih.nci.cabio.domain.Gene[@symbol=EGFR]

gov.nih.nci.cabio.domain.Gene								
bigid	clusterid	fullname	hugoSymbol	id	symbol	libraryCollection	geneFunctionAssociat	
hdi://2500.1.PMEUQUCL5.RJ3XIV2MP	488293	Epidermal growth factor receptor (erythroblastic leukemia viral (v-erb-b) oncogene homolog, avian)	EGFR	21065	EGFR	getLibraryCollection	getGeneFunctionAssociat	
hdi://2500.1.PMEUQUCL5.RJGUK7JP63	8534	Epidermal growth factor receptor	-	124410	Egfr	getLibraryCollection	getGeneFunctionAssociat	
hdi://2500.1.PMEUQUCL5.LY36OHKDDQ	439882	Epidermal growth factor receptor	-	178417	Egfr	getLibraryCollection	getGeneFunctionAssociat	

Figure 4.4 QBE search results

caBIO Portlet

A “portlet” is a pluggable user interface component that operates within a web portal. The caBIO Portlet is a pluggable GUI component that operates within the caGrid Portal as a community service. The caBIO Portlet provides a simple search for performing “Google-like” queries and a templated-based search allowing users to execute pre-defined queries.

The simple search (*Figure 4.5*) is based on the existing FreestyleLM search that allows a user to search by entering one or more keywords. The results of the simple search are formatted based on the caBIO object type, to show a relevant title and description

and class for each result. Pagination allows the user to retrieve further results. All functionality in the Simple Search Portlet is implemented using Ajax with the REST API.

The screenshot shows the National Cancer Institute caGrid Portal header. Below it is the caBIO Portlet interface. It includes a description of caBIO as a repository of biomedical research data. There are two tabs: 'Simple Search' (selected) and 'Templated Searches'. A search input field contains the text 'blood' and a 'Search' button is next to it. Below the search field are 'Advanced Options' including a 'Match' section with radio buttons for 'any' (selected) and 'all', an 'Exclude' text field, and two checkboxes: 'Domain-specific formatting' (checked) and 'Map to objects (slower)' (unchecked).

Figure 4.5 caBIO Portlet -- Simple search

Results from Simple Searches are displayed in “Google-like” form (Figure 4.6). Additional details can be obtained for each search result by selecting the search result link.

This screenshot shows the same caBIO Portlet interface as Figure 4.5, but with search results displayed below the 'Advanced Options' section. A purple bar indicates 'Results 1 - 5 of 3348 for blood'. The results are listed as follows:

- [Peripheral **blood** mobilized Hematopoietic Precursor Cell T...](#)
NHLBI, National Heart Lung and **Blood** Institute
[gov.nih.nci.cabio.domain.ClinicalTrialProtocol](#)
- [Hi-Dose Melphalan & TBI with Autologous BM & Peripheral B...](#)
NHLBI, National Heart Lung and **Blood** Institute
[gov.nih.nci.cabio.domain.ClinicalTrialProtocol](#)
- [In corresponding peripheral **blood** of HCC tissues, the pos...](#)
rare expression in corresponding peripheral **blood** of HCC tissues
[gov.nih.nci.cabio.domain.Evidence](#)
- [IDE 8139: NHLBI/Dr. A John Barrett IDE Studies with MoAb ...](#)
NHLBI, National Heart Lung and **Blood** Institute

Figure 4.6 caBIO Portlet -- Simple search results

Templated searches provide the user with pre-defined searches categorized by functionality (*Figure 4.7*).

National Cancer Institute

caGrid Portal

HOME SERVICES TOOLS **COMMUNITY** NEWS REGISTER

caBIO Portlet

caBIO is a repository of data useful in biomedical research, compiled from multiple primary sources. This portlet is an easy way to get started with caBIO. [Learn more](#) about caBIO and its various APIs.

Simple Search **Templated Searches**

Select category:

- ☒ Genome Range Queries
- ☒ Microarray Annotations
- ☒ Genomic Annotations

Gene Lookup by Symbol ▼ Gene Symbol

☒ Pathways

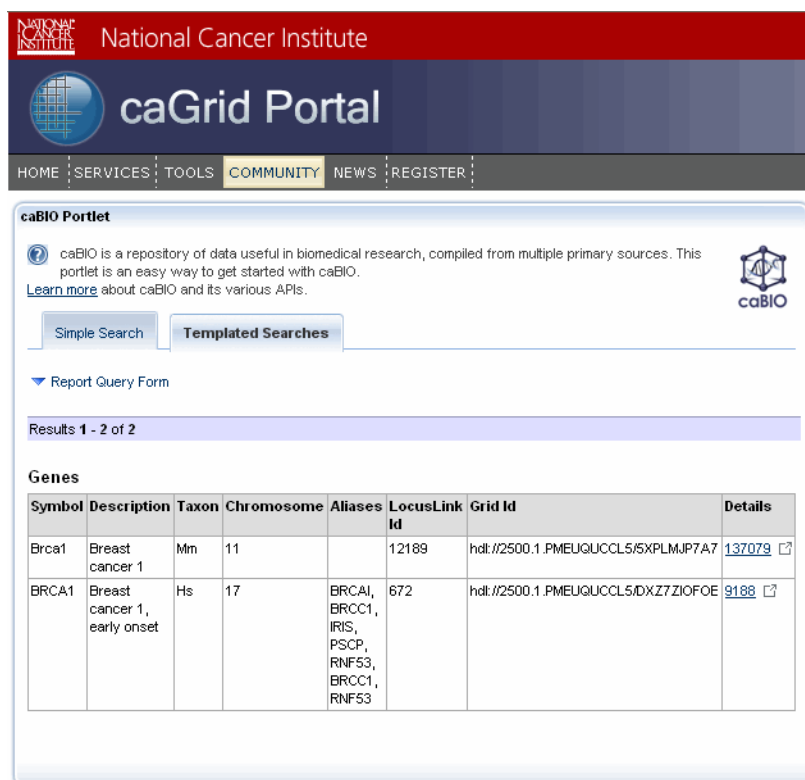
☒ Cancer Gene Index

Figure 4.7 caBIO Portlet -- Templated search

Templated searches supported by caBIO include:

- Genome Range Queries – Absolute range queries by Taxon or Chromosome and Grid Id range queries
- Microarray Annotations – Reporter lookup by Gene, SNP, or Reporter Name
- Genomic Annotations – Gene lookup by Gene symbol or HUGO symbol
- Pathways – Pathway search by Pathway name or Gene symbol
- Cancer Gene Index – Search for Diseases or Agents associated with a Gene

Results from Templated Searches are displayed in tabular form (Figure 4.8). Additional details can be obtained for each search result via the Details link.



The screenshot shows the caGrid Portal interface. At the top is the National Cancer Institute logo and the caGrid Portal header. Below the header is a navigation bar with links: HOME, SERVICES, TOOLS, COMMUNITY, NEWS, and REGISTER. The COMMUNITY tab is selected. The main content area displays the caBIO Portlet. It includes a description of caBIO as a repository of data useful in biomedical research, compiled from multiple primary sources. Below the description are two tabs: Simple Search and Templated Searches. The Templated Searches tab is active, showing a 'Report Query Form' link and a table of search results. The table has columns: Symbol, Description, Taxon, Chromosome, Aliases, LocusLink Id, Grid Id, and Details. Two results are shown: Brca1 (Breast cancer 1, Mm, 11, 12189, hdl://2500.1.PMEUGUCL5/5XPLMJP7A7, 137079) and BRCA1 (Breast cancer 1, early onset, Hs, 17, BRCA1, BRCC1, IRIS, PSCP, RNF53, BRCC1, RNF53, 672, hdl://2500.1.PMEUGUCL5/DXZ7ZIOFOE, 9188).

Symbol	Description	Taxon	Chromosome	Aliases	LocusLink Id	Grid Id	Details
Brca1	Breast cancer 1	Mm	11		12189	hdl://2500.1.PMEUGUCL5/5XPLMJP7A7	137079
BRCA1	Breast cancer 1, early onset	Hs	17	BRCA1, BRCC1, IRIS, PSCP, RNF53, BRCC1, RNF53	672	hdl://2500.1.PMEUGUCL5/DXZ7ZIOFOE	9188

Figure 4.8 Templated search results

The caBIO Portlet is available on the Community tab of the caGrid portal, <http://cagrid-portal.nci.nih.gov/web/guest/community>.

caBIO Service Interface Paradigm

The caBIO API is built on top of the caCORE SDK. The caCORE architecture includes a service layer that provides a single, common access paradigm to clients using any of the provided interfaces. As an object-oriented middleware layer designed for flexible data access, caBIO relies heavily on strongly-typed objects and an object-in/object-out mechanism. The methodology used for obtaining data from the caBIO system is often referred to as query-by-example, meaning that the inputs to the query methods are themselves domain objects that provide the criteria for the returned data. The major benefit of this approach is that it allows for run-time semantic interoperability when used as part of the caCORE infrastructure, which provides shared vocabularies and a metadata registry.

The basic order of operations required to access and use the caBIO system is as follows:

1. Ensure that the client application has knowledge of the objects in the domain space.
2. Formulate the query criteria using the domain objects.
3. Establish a connection to the server.

4. Submit the query objects and specify the desired class of objects to be returned.
5. Use and manipulate the result set as desired.

There are five primary application programming interfaces (APIs) native to the caBIO system. Each of the available interfaces described below uses this same paradigm to provide access to the caBIO domain model, but with minor changes relating primarily to the syntax and structure of the clients. The following sections describe each API, identify installation and configuration requirements, and provide code examples.

Java API

The Java API provides direct access to domain objects and all service methods. Because caBIO is natively built in Java, this API provides the fullest set of features and capabilities.

Installation and Configuration

The caBIO Java API uses the following on the client machine (*Table 4.1*).

Software	Version	Required
Java 2 Platform Standard Edition Software 5.0 Development Kit (JDK 5.0)	1.5.0_14 or higher	Yes
Apache Ant	1.6.5 or higher	Yes

Table 4.1 caBIO Java API client software

Accessing the caBIO system requires an Internet connection.

To use the Java API, download the client package provided on the NCICB web site.

To install and configure the Java API, follow these steps:

1. Open your browser and go to <http://ncicb.nci.nih.gov>. (*Figure 4.9*)

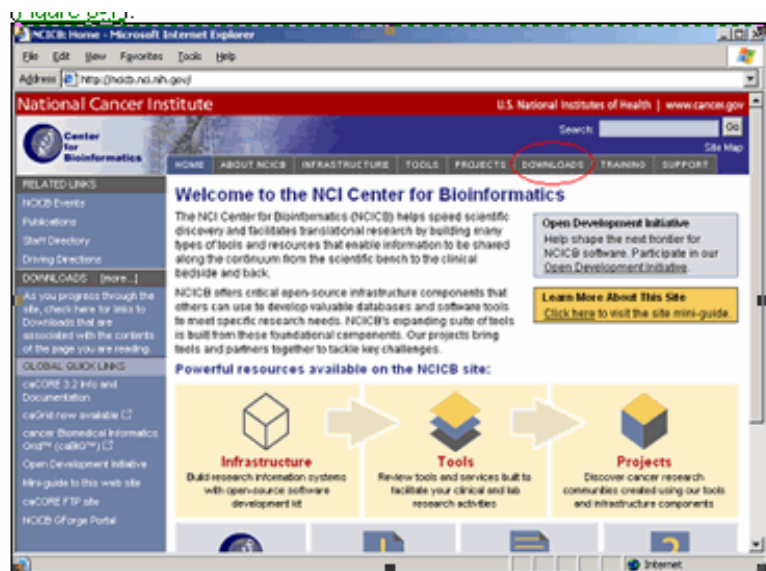


Figure 4.9 NCICB website Downloads option

- Click the **Download** link on the menu bar.
- Scroll down to the section titled **caBIO** and click the **Download** link.
- On the provided form, enter your name, email address and institution name and click to **Enter the Download Area**.
- Accept the license agreement.
- On the caBIO downloads page, download the caBIO Zip file from the **Primary Distribution** section.
- Extract the contents of the downloadable archive to a directory on your hard drive (for example, `c:\cabio` on Windows or `/usr/local/cabio` on Linux). The extracted directories and files include those in *Table 4.2*

<i>Directories and Files</i>	<i>Description</i>	<i>Component</i>
build.xml	Ant build file	Build file
Src directory	Contains demo code	Sample code
TestClient.java	Java API client sample	SDK sample code
TestGetXMLClient.java	REST client sample	
TestXMLClient.java	XML marshalling sample	
TechGuideExamples.java	Java API examples from this guide	caBIO sample code
TestQueries.java	Java API query examples	
TestFreestyleLM.java	FreestyleLM search examples	
TestCQL.java	CQL API examples	
TestSVG.java	SVG pathway image examples	
TestXML.java	XML utility sample	
lib directory	contains required jar files	Java Libraries
cabio41-beans.jar	domain objects	caBIO API
cabio41-client-framework.jar	caBIO API extensions to the caCORE SDK Framework	
Spring.jar	Spring framework	HTTP Remoting
hibernate*.jar	Hibernate	ORM layer
cglib-2.1.3.jar	CGLIB proxy library	
Log4j-1.2.13.jar	Log4j and Commons Logging	Logging Framework
commons-logging.jar		
commons-discovery-0.2.jar		

Table 4.2 Extracted directories and files in caBIO Java client package

Directories and Files	Description	Component
castor-1.0.2.jar	Castor serializer/deserializer	XML conversion
xercesImpl.jar	Apache Xerces XML parser	
*.xsd	XML schemas for objects	
handle.jar	Handle Service	Handle Service
handlerservice.jar	Custom handle code	
conf directory	Configuration files	
xml-mapping.xml	XML serialization mapping configuration files	XML mapping
unmarshaller-xml-mapping.xml		
application-config-client.xml	Spring configuration for client	Java API Spring configuration
Log4j.properties	Logging utilities configuration properties	Logging configuration

Table 4.2 Extracted directories and files in caBIO Java client package (Continued)

All of the jar files provided in the lib directory of the caBIO client package in addition to the files in the conf directory are required to use the Java API. These should be included in the Java classpath when building applications. The `build.xml` file that is included demonstrates how to do this when using Ant for command-line builds. If you are using an integrated development environment (IDE) such as Eclipse, refer to the tool's documentation for information on how to set the classpath.

A Simple Example

To run the simple example program after installing the caBIO client, open a command prompt or terminal window from the directory where you extracted the downloaded archive and enter `ant rundemo`. This will compile and run the `TechGuideExamples` class; successfully running this example indicates that you have properly installed and configured the caCORE client. The following is a short segment of code from the `TechGuideExamples` class along with an explanation of its functioning.

```

1  ApplicationService appService = ApplicationServiceProvider.getApplicationService();
2  Gene gene = new Gene();
3  gene.setSymbol("brca*"); // searching for all genes whose symbols start with brca
4  try {
5      List resultList = appService.search(Gene.class, gene);
6      for (Iterator resultsIterator = resultList.iterator(); resultsIterator.hasNext();)
7      {
8          Gene returnedGene = (Gene) resultsIterator.next();
9          System.out.println("Symbol: " + returnedGene.getSymbol() +
10              "\tName " + returnedGene.getFullName() +
11              "\tTaxon:" + returnedGene.getTaxon().getScientificName());
12      }
13  } catch (Exception e) {
14      e.printStackTrace();
15  }

```

This code snippet creates an instance of a class that implements the `ApplicationService` interface. This interface defines the service methods used to access data objects. A criterion object is then created that defines the attribute values for which to search. The search method of the `ApplicationService` implementation is called with

parameters that indicate the type of objects to return, `Gene.class`, and the criteria that returned objects must meet, defined by the gene object. The search method returns objects in a List collection, which is iterated through to print some basic information about the objects.

Although this is a simple example of the use of the Java API, a similar sequence can be followed with more complex criteria to perform sophisticated manipulation of the data provided by caBIO. Additional information and examples are provided in the sections that follow.

Service Methods

The methods that provide programmatic access to the caBIO server are located in the `gov.nih.nci.system.applicationservice` package. The `ApplicationServiceProvider` class uses the factory design pattern to return an implementation of the `ApplicationService` interface. The provider class determines whether there is a locally running instance of the caBIO system or whether it should use a remote instance. The returned `ApplicationService` implementation exposes the service methods that enable read/write operations on the domain objects.

Note: Although the infrastructure is capable of write operations, this functionality has been disabled for caBIO because it is primarily meant as a read-only data system.

The separation of the service methods from the domain classes is an important architectural decision that insulates the domain object space from the underlying service framework. As a result, new business methods can be added without needing to update any of the domain model or the associated metadata information from the object model. (This is critical for ensuring semantic interoperability over multiple iterations of architectural changes.)

Within the `ApplicationService` implementation, a variety of methods are provided allowing users to query data based on the specific needs and types of queries to be performed. In general, there are five types of searches:

- **Simple searches** are those that take one or more objects from the domain models as inputs and return a collection of objects from the data repositories that meet the criteria specified by the input objects.
- **Nested searches** also take domain objects as inputs but determine the type of objects in the result set by traversing a known path of associations from the domain model.
- **Detached criteria searches** use Hibernate detached criteria objects to provide a greater level of control over the results of a search (such as boolean operations, ranges of values, etc.)
- **HQL searches** provide the ability to use the Hibernate Query Language for the greatest flexibility in forming search criteria.
- **SDK Query Object criteria searches** were modeled similar to the Object representation of caBIG Query Language (CQL). The SDK Query Object criteria searches use syntax similar to the query by example (QBE) to specify the way results are to be retrieved. The system formulates the query based on

the navigation path specified in the query search criteria. The query mechanism allows the user to search for the objects in platform independent query syntax.

Method Signature	<code>List search(Class targetClass, Object obj)</code>
Search Type	Simple (One criteria object)
Description	Returns a List collection containing objects of type <code>targetClass</code> that conform to the criteria defined by <code>obj</code>
Example	<code>search(Gene.class, gene)</code>

Method Signature	<code>List search(Class targetClass, Object objList)</code>
Search Type	Simple (One criteria object)
Description	Returns a List collection containing objects of type <code>targetClass</code> that conform to the criteria defined by a collection of objects in <code>objList</code> . The returned objects must meet ANY criteria in <code>objList</code> (i.e. a logical OR is performed).
Example	<code>search(Gene.class, geneCollection)</code>

Method Signature	<code>List search(String path, Object obj)</code>
Search Type	Nested
Description	Returns a List collection containing objects conforming to the criteria defined by <code>obj</code> and whose resulting objects are of the type reached by traversing the node graph specified by <code>path</code>
Example	<code>search("gov.nih.nci.cabio.domain.Protein, gov.nih.nci.cabio.domain.Gene", nucleicAcidSequence)</code>

Method Signature	<code>List search(String path, List objList)</code>
Search Type	Nested
Description	Returns a List collection containing objects conforming to the criteria defined by the objects in <code>objList</code> and whose resulting objects are of the type reached by traversing the node graph specified by <code>path</code>
Example	<code>search("gov.nih.nci.cabio.domain.Protein, gov.nih.nci.cabio.domain.Gene", sequenceList)</code>

Method Signature	List query(DetachedCriteria detachedCriteria)
Search Type	Detached criteria
Description	Returns a List collection conforming to the criteria specified by detachedCriteria.
Example	query(criteria)

Method Signature	List query(Object criteria, Integer firstRow, String targetClassName)
Search Type	Detached criteria
Description	Identical to the previous query method, but allows for control over the size of the result set by specifying the row number of the first row.
Example	query(criteria, 101, targetClassName)

Method Signature	List query(HQLCriteria hqlCriteria)
Search Type	Hibernate Query Language
Description	Returns a List collection of objects that conform to the query in HQL syntax contained in hqlCriteria.
Example	query(hqlCriteria)

Method Signature	List query(CQLQuery cqlQuery)
Search Type	Common Query Language
Description	Returns a List collection of objects that conform to the query in SDK Query Object syntax contained in cqlQuery.
Example	query(cqlQuery)

In addition to the data access methods, the following helper methods are available via the ApplicationService class that provide flexibility in controlling queries and result sets.

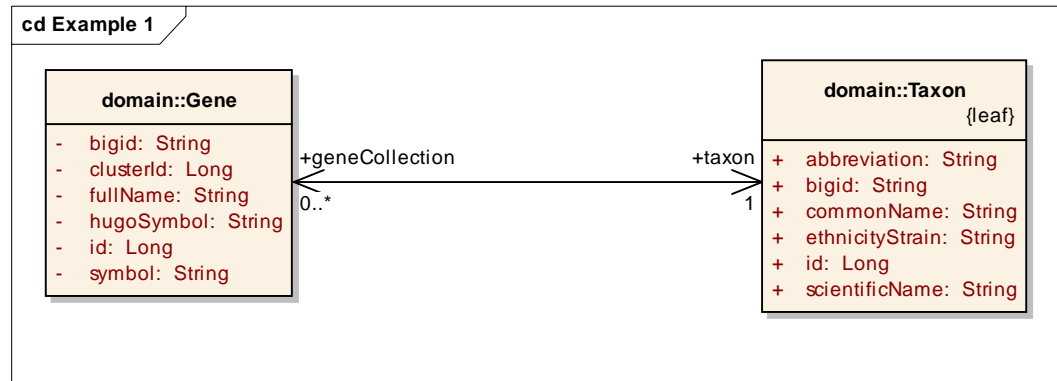
Method Signature	int getQueryRowCount(Object criteria, String targetClassName)
Description	Gets the number of objects in a given query. This is useful in determining the size of the data before querying for data
Example	getQueryRowCount(criteria, "gov.nih.nci.cabio.domain.Gene")

Examples of Use

This section includes a number of examples that demonstrate the use of the caBIO APIs. Included with each example is a brief description of the type of search being performed, a UML diagram depicting the domain objects used, and the example code accompanied by explanatory text.

Example One: Simple Search (Single Criteria Object)

In this example, a search is performed for all genes whose symbols start with 'brca'. The code iterates through the returned objects and prints out the symbol and name of each object along with the name of an associated object of type Taxon. The fetch of the associated Taxon object is done in the background and is completely transparent to the ser.



```

1  ApplicationService appService = ApplicationServiceProvider.getApplicationService();
2  Gene gene = new Gene();
3  gene.setSymbol("brca*"); // searching for all genes whose symbols start with 'brca'
4  try
5  {
6      List resultList = appService.search(Gene.class, gene);
7      for (Iterator resultsIterator = resultList.iterator(); resultsIterator.hasNext();)
8      {
9          Gene returnedGene = (Gene) resultsIterator.next();
10         System.out.println("Symbol: " + returnedGene.getSymbol() +
11             "\tTaxon:" + returnedGene.getTaxon().getScientificName()
12             "\tName " + returnedGene.getFullName());
13     }
14 } catch (Exception e) {
15     e.printStackTrace();
16 }
  
```

Lines	Description
1	Creates an instance of a class that implements the ApplicationService interface; this interface defines the service methods used to access data objects.
2-3	Creates a criterion object that defines the attribute values for which to query.
6	Calls the search method of the ApplicationService implementation and passes it the type of objects to return, Gene.class, and the criteria that returned objects must meet, defined by the gene object; the search method returns objects in a List collection.

<i>Lines</i>	<i>Description</i>
9	Casts an object from the result List and creates a variable reference to it of type Gene.
10	Prints the symbol attribute.
11	Prints the fullName attribute.
12	Fetches an associated object of type Taxon and prints its scientificName attribute.

Example Two: Simple Search (Criteria Object Collection)

This example demonstrates a search with multiple criteria objects that are passed to the search() method. The result set includes all objects of the specified type that match ANY of the criteria objects. In this case, the search returns all objects of type Gene that are associated with Taxon objects whose abbreviation attribute is either "hs" or "ms". In biological terms, this search returns all human and mouse genes.

```

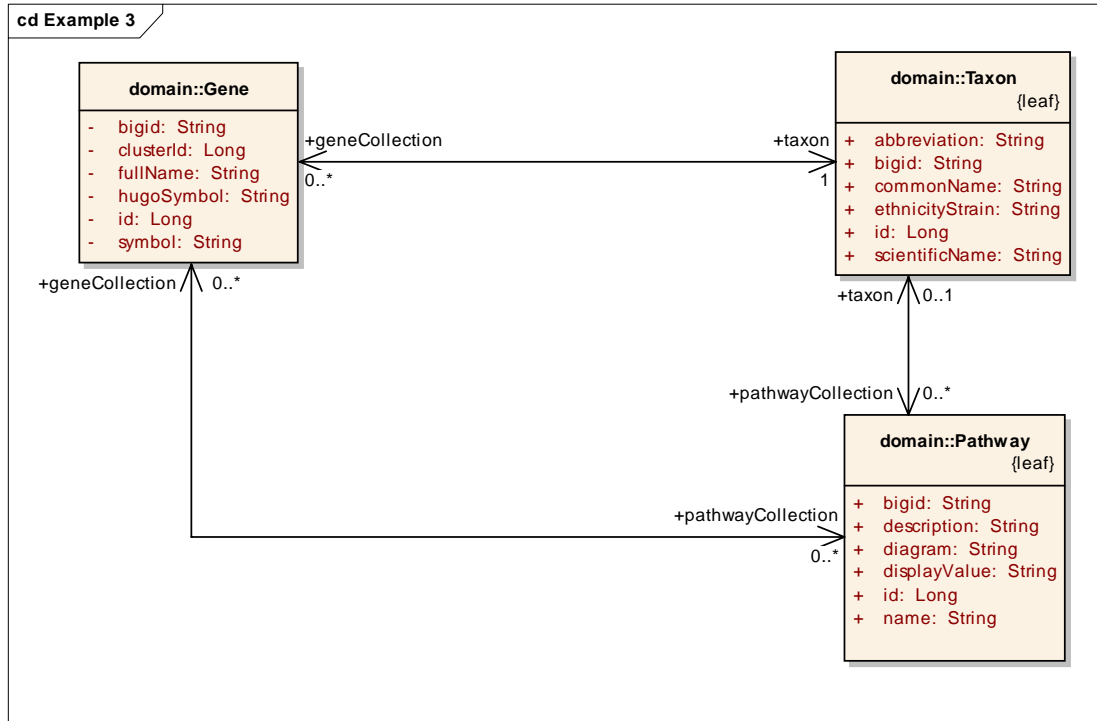
1  Taxon taxon1 = new Taxon();
2  taxon1.setAbbreviation("Hs");           // Homo sapiens
3  Taxon taxon2 = new Taxon();
4  taxon2.setAbbreviation("Mm");           // Mus musculus
5  List<Taxon> taxonList = new ArrayList<Taxon>();
6  taxonList.add(taxon1);
7  taxonList.add(taxon2);
8  try
9  {
10     List resultList = appService.search(Gene.class, taxonList);
11     System.out.println("Number of results: "+resultList.size());
12 }
13 catch (Exception e) {
14     e.printStackTrace();
15 }
```

<i>Lines</i>	<i>Description</i>
1-4	Creates two Taxon objects describing the search criteria.
5-7	Creates a List collection containing the two Taxon objects.
10	Searches for all Gene objects where the associated Taxon object matches ANY of the objects found in the taxonList collection.

Example Three: Simple Search (Compound Criteria Object)

In this example, the object that is passed to the search() method contains criteria values that are found in associated objects and collections of objects. This query returns those objects that match all of the criteria in the compound object. Note the distinction between this type of search and the previous example in which a collection of objects is passed into the search method. In the last example, the results will match ANY of the criteria objects. In this example, however, where a single compound object

is used, ALL criteria are matched. In biological terms, this search will return all pathways associated with the human Interleukin 5 gene.



```

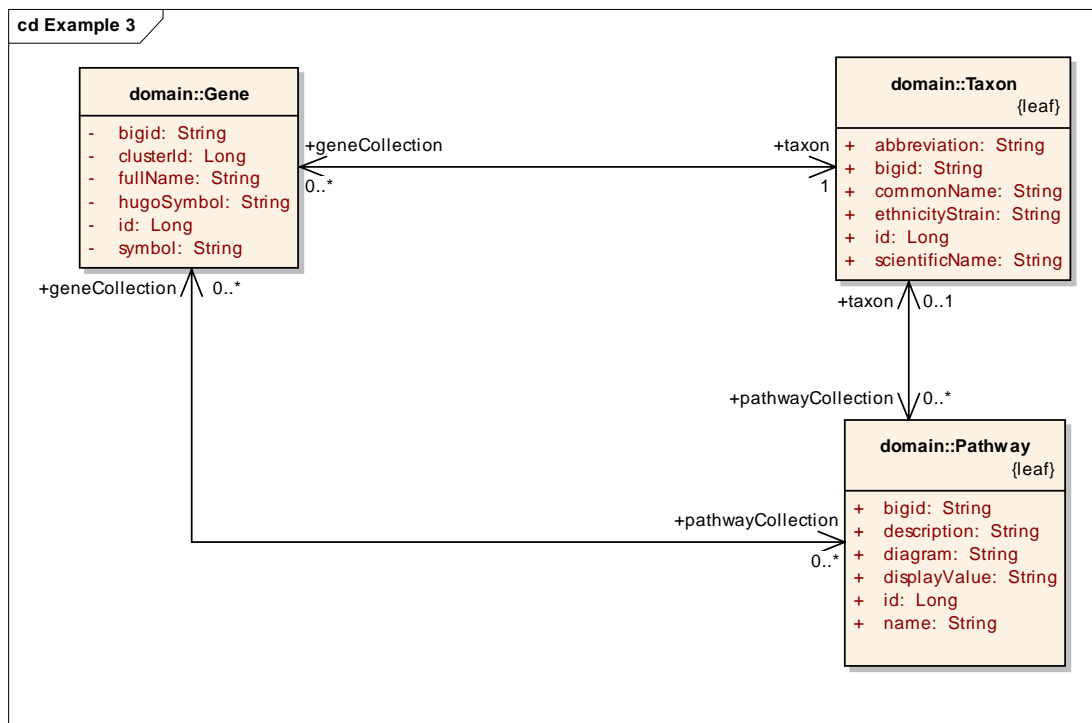
1  Taxon taxon = new Taxon();
2  taxon.setAbbreviation("hs");           // Homo sapiens
3  Gene gene = new Gene();
4  gene.setTaxon(taxon);
5  gene.setSymbol("IL5");                 // Interleukin 5
6  List<Gene> geneList = new ArrayList<Gene>();
7  geneList.add(gene);
8  Pathway pathway = new Pathway();
9  pathway.setGeneCollection(geneList);
10 try
11 {
12     List resultList = appService.search("gov.nih.nci.cabio.domain.Pathway", pathway);
13     for (Iterator resultsIterator = resultList.iterator(); resultsIterator.hasNext();)
14     {
15         Pathway returnedPathway = (Pathway)resultsIterator.next();
16         System.out.println("Name: " + returnedPathway.getName()
17             + "\tDisplayValue: " + returnedPathway.getDisplayValue());
18     }
19 } catch (Exception e) {
20     e.printStackTrace();
21 }
  
```

Lines	Description
1-5	Creates a Gene object and sets the symbol to "IL5" and the associated Taxon to an object whose abbreviation is set to "hs".
6-7	Because the Pathway and Gene classes are related by a many-to-many association, it is necessary to create a collection to contain the Gene object that will act as part of the compound criteria; multiple Gene objects could be added to this collection as needed.
8-9	Creates a Pathway object and sets the value of its geneCollection to the geneList object just created.

Lines	Description
12	Searches for all Pathway objects whose geneCollection contains objects that match the set criteria (i.e. the symbol is "IL5" and the associated Taxon objects' abbreviations are set to "hs").

Example Four: Nested Search

A nested search is one where a traversal of more than one class-class association is required to obtain a set of result objects given the criteria object. This example demonstrates one such search in which the criteria object passed to the search method is of type `Gene`, and the desired objects are of type `ProteinSequence`. Because there is no direct association between these two classes, the path of the traversal is passed to the search method enabling the query to be performed.



```

1  Gene gene = new Gene();
2  gene.setSymbol("TP53");           // Tumor protein p53 (Li-Fraumeni syndrome)
3  try
4  {
5      List resultList = appService.search(
6          "gov.nih.nci.cabio.domain.ProteinSequence,gov.nih.nci.cabio.domain.Protein"
7          gene);
8      for (Iterator resultsIterator = resultList.iterator(); resultsIterator.hasNext()
9      {
10         ProteinSequence returnedProtSeq = (ProteinSequence)resultsIterator.next();
11         System.out.println("Id: " + returnedProtSeq.getId() +
12             "Length: " + returnedProtSeq.getLength() );
13     }
14 } catch (Exception e) {
15     e.printStackTrace();
16 }

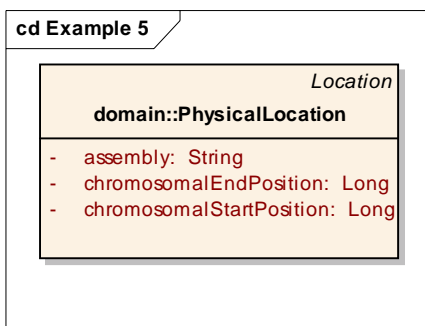
```

Lines	Description
1-2	Creates a Gene object and sets the symbol to "TP53".

<i>Lines</i>	<i>Description</i>
6	Defines search path as traversing from the criteria object of type Gene through Protein to ProteinSequence; note that the first element in the path is the desired class of objects to be returned, and that subsequent elements traverse back to the criteria object.
7	Sets the criteria object to the previously-created Gene.

Example Five: Detached Criteria Search

This example demonstrates the use of Hibernate detached criteria objects to formulate and perform more sophisticated searches. A detailed description of detached criteria is beyond the scope of this document; for more information, please consult the Hibernate documentation at http://www.hibernate.org/hib_docs/v3/api/org/hibernate/criterion/DetachedCriteria.html.



```

1 DetachedCriteria criteria = DetachedCriteria.forClass(PhysicalLocation.class);
2 try {
3     DetachedCriteria criteria = DetachedCriteria.forClass(PhysicalLocation.class);
4     criteria.add(Restrictions.eq("assembly", "reference"));
5     criteria.add(Restrictions.gt("chromosomalStartPosition", new Long(86851632)));
6     criteria.add(Restrictions.lt("chromosomalEndPosition", new Long(86861632)));
7     List resultList = appService.query(criteria);
8     System.out.println("\n Total # of records = " + resultList.size());
9 }
10 }
11 catch (Exception e) {
12     e.printStackTrace();
13 }

```

<i>Lines</i>	<i>Description</i>
1	Creates a DetachedCriteria object and sets the class on which the criteria will operate to PhysicalLocation.
5	Sets a restriction on the objects which states that the attribute assembly must be equal to "reference".
6	Sets a restriction on the objects that states that the attribute chromosomalStartPosition must be greater than ("gt") the value 86851632.
7	Sets a restriction on the objects that states that the attribute chromosomalEndPosition must be less than ("lt") the value 86861632.
8	Calls the query method of the ApplicationService implementation, passing the detached criteria object.

Example Six: HQL Search

In this example, a search is performed for all genes whose symbols start with 'brca'. This is identical to Example One but uses a Hibernate Query Language (HQL) search string to form the query. For more information on HQL syntax, consult the Hibernate documentation at http://www.hibernate.org/hib_docs/v3/reference/en/html/queryhql.html

```

1 try {
2     String hqlString = "FROM Gene g WHERE g.symbol LIKE 'BRCA%'";
3     HQLCriteria hqlC = new HQLCriteria(hqlString);
4     List resultList = appService.query(hqlC);
5     System.out.println("\n Total # of records = " + resultList.size());
6 }
7 catch (Exception e) {
8     e.printStackTrace();
9 }

```

Example Seven: SDK Query Object Search

SDK Query Object Example One

```

ApplicationService appService = ApplicationServiceProvider.getApplicationService();

CQLQuery cqlQuery = new CQLQuery();
CQLObject target = new CQLObject(); // Create Gene Object
target.setName("gov.nih.nci.cabio.domain.Gene");

CQLAttribute attribute = new CQLAttribute();
attribute.setName("symbol");
attribute.setValue("%il%");
attribute.setPredicate(CQLPredicate.LIKE);
target.setAttribute(attribute);

cqlQuery.setTarget(target);

try {
    List resultList = appService.query(cqlQuery, gov.nih.nci.cabio.domain.Gene);
    //Iterate through retrieved list of objects
    for (Iterator i = resultList.iterator(); i.hasNext();) {
        Gene returnedGene = (Gene) i.next();
        System.out.println(
            "Symbol: " + returnedGene.getSymbol() + "\n" +
            "\tTaxon:" + returnedGene.getTaxon().getScientificName() + "\n" +
            "\tName " + returnedGene.getFullName() + "\n");
    }
}
catch (Exception e) {
    e.printStackTrace();
}

```

Lines	Description
2	Creates a new CQLQuery Object.
3-4	Creates a new search Target object (Gene).
5-8	Creates a new attribute for the search target object whose value should be compared with the value in the database for given predicate.
10	Sets the search target object in SDK Query Object Query.
12-13	Sets the search criteria object to the previously-created CQLQuery.

SDK Query Object Example Two

```

1      CQLQuery cqlQuery = new CQLQuery();
2      CQLObject target = new CQLObject();
3      target.setName("gov.nih.nci.cabio.domain.Gene");

4      CQLAttribute attribute = new CQLAttribute();
5      attribute.setName("symbol");
6      attribute.setValue("brca%");
7      attribute.setPredicate(CQLPredicate.LIKE);

8      //Create Taxon Object Query
9      CQLAssociation association1 = new CQLAssociation();
10     association1.setName("gov.nih.nci.cabio.domain.Taxon");
11     CQLAttribute attributel = new CQLAttribute();
12     attributel.setName("abbreviation");
13     attributel.setValue("%hs%");
14     attributel.setPredicate(CQLPredicate.LIKE);
15     association1.setAttributel(attributel);

16     //Create Taxon Object Query
17     CQLAssociation association2 = new CQLAssociation();
18     association2.setName("gov.nih.nci.cabio.domain.Taxon");
19     CQLAttribute attribute2 = new CQLAttribute();
20     attribute2.setName("abbreviation");
21     attribute2.setValue("%m%");
22     attribute2.setPredicate(CQLPredicate.LIKE);
23     association2.setAttributel(attribute2);

24     //Create Group(Collection) of Taxon Object Query
25     CQLGroup group = new CQLGroup();
26     group.addAssociation(association1);
27     group.addAssociation(association2);
28     group.setLogicOperator(CQLLogicalOperator.OR);
29     target.setAttributel(attribute);
30     target.setGroup(group);

31     cqlQuery.setTarget(target);

32     try {
33         List resultList = appService.query(cqlQuery);
34
35         //Iterate through retrieved list of objects
36         for (Iterator resultsIterator = resultList.iterator();
37 resultsIterator.hasNext();) {
38             Gene returnedGene = (Gene) resultsIterator.next();
39             System.out.println(
40                 "Symbol: " + returnedGene.getSymbol() + "\n" +
41                 "\tTaxon:" + returnedGene.getTaxon().getScientificName() +
42                 "\n" + "\tName " + returnedGene.getFullName() + "\n");
43         }
44     } catch (Exception e) {
45         e.printStackTrace();
46     }

```

Lines	Description
1	Creates a new CQLQuery Object.
2-3	Creates a new search Target object (Gene).
4-7	Sets attribute criteria for Target object
9-15, 17-23	Creates a new associated object (Taxon) whose attribute value should be compared with the value in the database for a given predicate.
25-30	Creates a group (collection) of objects and associates it to the search target object.
31	Sets the search target object in SDK Query Object Query.

Lines	Description
33	Search using the previously-created CQLQuery.

SDK Query Object Example Three

```

1      ApplicationService appService =
           ApplicationServiceProvider.getApplicationService();

2      CQLQuery cqlQuery = new CQLQuery();
3      CQLObject target = new CQLObject();
4      target.setName("gov.nih.nci.cabio.domain.Pathway");

5      //Create Gene Object Query
6      CQLAssociation association1 = new CQLAssociation();
7      association1.setName("gov.nih.nci.cabio.domain.Gene");
8      CQLAttribute attribute1 = new CQLAttribute();
9      attribute1.setName("symbol");
10     attribute1.setValue("IL5");
11     attribute1.setPredicate(CQLPredicate.EQUAL_TO);
12     association1.setAttribute(attribute1);

13     //Create Taxon Object Query
14     CQLAssociation association2 = new CQLAssociation();
15     association2.setName("gov.nih.nci.cabio.domain.Taxon");
16     CQLAttribute attribute2 = new CQLAttribute();
17     attribute2.setName("abbreviation");
18     attribute2.setValue("%hs%");
19     attribute2.setPredicate(CQLPredicate.LIKE);
20     association2.setAttribute(attribute2);

21     //Set Relationship between Gene and Taxon
22     association1.setAssociation(association2);

23     //Set Relationship between Pathway and Gene
24     //Role name is required as it can not be determined using
25     //reflection
26     association1.setTargetRoleName("geneCollection");
27     target.setAssociation(association1);

28     cqlQuery.setTarget(target);

29     try
30     {
31         List resultList = appService.query(cqlQuery);
32         //Iterate through retrieved list of objects
33         for (Iterator resultsIterator = resultList.iterator();
34             resultsIterator.hasNext();)
35         {
36             Pathway returnedPathway = (Pathway)resultsIterator.next();
37             System.out.println(returnedPathway.getDisplayValue());
38         }
39     } catch (Exception e) {
40         e.printStackTrace();
41     }

```

Lines	Description
2	Creates a new CQLQuery Object.
3-4	Creates a new search Target object (Pathway).
5-12	Creates a new associated object (Gene) whose attribute value should be compared with the value in the database for a given predicate.
13-20	Creates a new associated object (Taxon) whose attribute value should be compared with the value in the database for a given predicate.
22	Creates a relationship between two objects (Gene and Taxon).
26-27	Creates a relationship between a target object (Pathway) and its associated object (Gene). The role name is the name of the association used to retrieve the associated object. It is required in case it can not be determined by reflection.
28	Sets the search target object in a SDK Query Object Query.
31	Sets the search criteria object to the previously-created CQLQuery.

SDK Query Object Example Four

```

1      ApplicationService appService =
           ApplicationServiceProvider.getApplicationService();

2      CQLQuery cqlQuery = new CQLQuery();
3      CQLObject target = new CQLObject();
4      target.setName("gov.nih.nci.cabio.domain.ProteinSequence");

5      //Create Gene Object Query
6      CQLAssociation association1 = new CQLAssociation();
7      association1.setName("gov.nih.nci.cabio.domain.Gene");
8      CQLAttribute attributel = new CQLAttribute();
9      attributel.setName("symbol");
10     attributel.setValue("TP53");
11     attributel.setPredicate(CQLPredicate.EQUAL_TO);
12     association1.setAttribute(attributel);

13     //Create Protein Object Query
14     CQLAssociation association2 = new CQLAssociation();
15     association2.setName("gov.nih.nci.cabio.domain.Protein");

16     //Set Relationship between Gene and Protein
17     //Role name is required as it can not be determined using
18     //reflection
19     association1.setTargetRoleName("geneCollection");
20     association2.setAssociation(association1);

21     //Set Relationship between Protein and ProteinSequence
22     //Example of using relationship from target to source
24     association1.setSourceRoleName("proteinSequence");
25     target.setAssociation(association2);

26     cqlQuery.setTarget(target);

27     try
28     {
29         List resultList = appService.query(cqlQuery);
30         //Iterate through retrieved list of objects
31         for (Iterator resultsIterator = resultList.iterator();
32              resultsIterator.hasNext();)
33         {
34             ProteinSequence returnedProtSeq =
35                 (ProteinSequence)resultsIterator.next();
36             System.out.println("ID: " + returnedProtSeq.getId() +
37                               "Length: " + returnedProtSeq.getLength());
38         }
39     } catch (Exception e) {
40         e.printStackTrace();
41     }

```

Lines	Description
2	Creates a new CQLQuery Object.
3-4	Creates a new search Target object (ProteinSequence).
5-12	Creates a new associated object (Gene) whose attribute value should be compared with the value in the database for a given predicate.
14-15	Creates a new associated object (Protein).
19-20	Creates a relationship between a source object (Protein) and its associated object (Gene). The role name is the name of the association used to retrieve the associated object. It is required in case it can not be determined by reflection.
24-25	Creates a relationship between a target object (ProteinSequence) and its associated object (Protein).

<i>Lines</i>	<i>Description</i>
26	Sets the search target object in a SDK Query Object Query.
29	Sets the search criteria object to the previously-created CQLQuery.

Utility Methods

XML Utility

The caBIO Java Client provides a convenience class called `XMLUtility` in the `gov.nih.nci.common.util` package that provides a short cut to the caCORE SDK's utility for converting domain objects between native Java objects and XML serializations based on standard XML schemas. The XML schemas for all domain objects in caBIO, directly generated from the UML model, are included in the downloadable archive. Currently, the XML generated using the `XMLUtility` class includes only the object attributes; associated objects are not included.

The configuration used by the XML utility is stored in two files, `xml-mapping.xml` and `unmarshaller-xml-mapping.xml`. These files define the binding between class, attribute names and the corresponding XML element and attribute names.

A default marshaller and unmarshaller are automatically registered by the caBIO convenience class; Developers wishing to use their own should bypass the caBIO class and instantiate the SDK's `XMLUtility` in the `gov.nih.nci.system.client.util.xml` package.

In the following code, the XML utility is used to serialize an object and save it to a file stream. A new object is then instantiated from the file using the utility.

```

1 // Assume an object of type Gene called myGene
2 File myFile = new File("myGene.xml");
3 FileWriter myWriter = new FileWriter(myFile);
4 XMLUtility myUtil = new XMLUtility();
5 myUtil.toXML(myGene, myWriter);
6 Gene myNewGene = (Gene)myUtil.fromXML(myFile);
7 bool isSameGene = myNewGene.equals(myGene);    // true

```

<i>Lines</i>	<i>Description</i>
2-3	Creates a new file stream where the XML serialization will be saved.
4	Creates a new <code>XMLUtility</code> object; in this case, the default marshaller, and unmarshaller will be used.
5	Serializes the <code>myGene</code> object to XML using the mapping file and writes the output to the file stream <code>myGene.xml</code> .
6	Creates a new object called <code>myNewGene</code> by invoking the <code>fromXML()</code> method of the <code>XMLUtility</code> class and casting it to the proper type.
7	The newly created <code>Gene</code> object is equivalent to the old one.

For additional details, consult the caBIO JavaDocs at <http://cabioapi.nci.nih.gov/cabio41/docs/>.

Web Services API

The caBIO Web services API allows access to caBIO data from development environments where the Java API cannot be used, or where use of XML Web services is more desirable. This includes non-Java platforms and languages such as Perl, C/C++, .NET Framework (C#, VB.Net), Python, etc.

The Web services interface can be used in any language-specific application that provides a mechanism for consuming XML Web services based on the Simple Object Access Protocol (SOAP). In these environments, connecting to caBIO can be as simple as providing the endpoint URL. Some platforms and languages require additional client-side code to handle the implementation of the SOAP envelope and the resolution of SOAP types. A list of packages catering to different programming languages is available at <http://www.w3.org/TR/SOAP/> and at <http://www.soapware.org/>.

To maximize standards-based interoperability, the caBIO Web service conforms to the Web Services Interoperability Organization (WS-I) Basic Profile. The WS-I Basic Profile provides a set of non-proprietary specifications and implementation guidelines enabling interoperability between diverse systems. More information about WS-I compliance is available at <http://www.ws-i.org>.

On the server side, Apache Axis is used to provide SOAP-based inter-application communication. Axis provides the appropriate serialization and deserialization methods for the Java beans to achieve an application-independent interface. For more information about Axis, visit <http://ws.apache.org/axis/>.

Configuration

The caBIO WSDL file is located at <http://cabioapi.nci.nih.gov/cabio41/services/caBIOService?wsdl>.

In addition to describing the protocols, ports and operations exposed by the caBIO Web service, this file can be used by a number of IDEs and tools to generate stubs for caBIO objects. This allows code on different platforms to instantiate objects native to each for use as parameters and return values for the Web service methods. Consult the specific documentation for your platform for more information on how to use the WSDL file to generate class stubs.

The caBIO Web services interface has a single endpoint called caBIOService, which is located at <http://cabioapi.nci.nih.gov/cabio41/services/caBIOService>. Client applications should use this URL to invoke Web service methods.

Java WS Client

The Java WS client is a separate download from the Java API described earlier. For installation instructions, see the Java API on page 25. This distribution bundles the necessary Java libraries and example code that shows how to use the caBIO Web Service API.

Directories and Files	Description	Component
build.xml	Ant build file	Build
build.properties	Web Service configuration	Build
src directory	Contains demo code	Sample code

Directories and Files	Description	Component
TestClient.java	SDK WS client sample	
TestWS.java	Web services Java sample	
lib directory	contains required jar files	Java Libraries
cabio4l-beans.jar	domain objects	caBIO API
hibernate-search.jar	Hibernate Search Annotations	ORM layer
axis.jar	Apache Axis	Web Services
saa.jar	SOAP API for Java	
jaxrpc.jar	Java API for XML-based RPC	
wsdl4j-1.5.1.jar	WSDL for Java	
commons-logging.jar	Log4j and Commons Logging	Logging Framework
commons-discovery-0.2.jar		
conf directory	Configuration files	
undeploy.wsdd	XML serialization mapping configuration files	XML mapping
application-config-client.xml	Spring configuration for client	Java API Spring configuration
Log4j.properties	Logging utilities configuration properties	Logging configuration

All of the jar files provided in the lib directory of the caBIO client package in addition to the files in the conf directory are required to use the Java API. These should be included in the Java classpath when building applications. The build.xml file that is included demonstrates how to do this when using Ant for command-line builds. If you are using an integrated development environment (IDE) such as Eclipse, refer to the tool's documentation for information on how to set the classpath.

Operations

Through the caBIOService endpoint, developers have access to four operations:

Operation	queryObject
Input Schema	<pre> <complexType> <sequence> <element name="in0" type="xsd:string"/> <element name="in1" type="xsd:anyType"/> </sequence> </complexType> </pre>

Output Schema	<pre><sequence> <element name="queryReturn" type= "ArrayOf_xsd_anyType" /> </sequence></pre>
Description	Performs a search for objects conforming to the criteria defined by input parameter in1 and whose resulting objects are of the type reached by traversing the node graph specified by parameter in0; the result is a set of serialized objects (the type ArrayOf_xsd_anyType resolves to a sequence of xsd:anyType elements)

Operation	Query
Input Schema	<pre><complexType> <sequence> <element name="in0" type="xsd:string"/> <element name="in1" type="xsd:anyType"/> <element name="in2" type="xsd:int"/> <element name="in3" type="xsd:int"/> </sequence> </complexType></pre>
Output Schema	<pre><sequence> <element name="queryReturn" type="ArrayOf_xsd_anyType" /> </sequence></pre>
Description	Identical to the previous queryObject method, but allows for control over the result set by specifying the row number of the first row (in2) and the maximum number of objects to return (in3)

Developers should be aware of a significant behavioral decision that has been made regarding the Web services interface. When a query is performed with this interface, returned objects do not contain or refer to their associated objects. This means that a separate query invocation must be performed for each set of associated objects that need to be retrieved. One of the examples below demonstrates this functionality.

Examples of Use

Example One: Simple Search

The following code demonstrates a simple query written in the Java language that uses the Web services API. This example uses Apache Axis on the client side to handle the type mapping, SOAP encoding, and operation invocation.

```

1  Service service = new Service();
2  Call call = (Call) service.createCall();
3
4  QName qnGene = new QName("urn:ws.domain.cabio.nci.nih.gov", "Gene");
5  call.registerTypeMapping(
6      Gene.class,
7      qnGene,
8      new org.apache.axis.encoding.ser.BeanSerializerFactory(Gene.class, qnGene),
9      new org.apache.axis.encoding.ser.BeanDeserializerFactory(Gene.class, qnGene)
10 );
11
12 String url = "http://cabioapi.nci.nih.gov/cabio40/services/caBIOService";
13
14 call.setTargetEndpointAddress(new java.net.URL(url));
15 call.setOperationName(new QName("caCOREService", "queryObject"));
16 call.addParameter("arg1", org.apache.axis.encoding.XMLType.XSD_STRING, ParameterMode.IN);
17 call.addParameter("arg2", org.apache.axis.encoding.XMLType.XSD_ANYTYPE, ParameterMode.IN);
18 call.setReturnType(org.apache.axis.encoding.XMLType.SOAP_ARRAY);
19
20 gov.nih.nci.cabio.domain.ws.Gene gene = new gov.nih.nci.cabio.domain.ws.Gene();
21 gene.setSymbol("IL*");
22
23 try
24 {
25     Object[] resultList = (Object[])call.invoke(
26         new Object[]{"gov.nih.nci.cabio.domain.ws.Gene", gene });
27     System.out.println("Total objects found: " + resultList.length);
28     if (resultList.length > 0)
29     {
30         for(int resultIndex = 0; resultIndex < resultList.length; resultIndex++)
31         {
32             Gene returnedGene = (Gene)resultList[resultIndex];
33             System.out.println(
34                 "Symbol: " + returnedGene.getSymbol() + "\n" +
35                 "\tName " + returnedGene.getFullName() + "\n" +
36                 "");
37         }
38     }
39 } catch (Exception e) {
40     e.printStackTrace();
41 }

```

Lines	Description
1-2	Defines a new Web service call.
4-10	Maps a serialized object to its Java equivalent using the qualified name of the object from the WSDL file; in this case, the XML element Gene in the urn:ws.domain.cabio.nci.nih.gov namespace is mapped to the Java Gene class.
12	Defines the service endpoint.
14-18	Sets the call properties including the name of the operation to invoke, the input parameters that will be sent and the return type to expect.
20-21	Creates a Gene criteria object and sets its symbol attribute to "IL*"; note that the *.ws.* package is used.
25-26	Invokes the Web service operation using an array of two objects (target class name and criteria object) as input parameters and expecting an object array as its result.
32	Casts each object in the result array to type Gene.

Example Two: Searching Associations

This example is similar to the previous one but demonstrates how to search for associated elements by performing additional invocations of the query or queryObject operation.

```

1  try
2  {
3      Object[] resultList = (Object[])call.invoke(
4          new Object[]{"gov.nih.nci.cabio.domain.Gene", gene });
5      System.out.println("Total objects found: " + resultList.length);
6      if (resultList.length > 0)
7      {
8          for(int resultIndex = 0; resultIndex < resultList.length; resultIndex++)
9          {
10             Gene returnedGene = (Gene)resultList[resultIndex];
11             System.out.println(
12                 "Symbol: " + returnedGene.getSymbol() + "\n" +
13                 "\tName: " + returnedGene.getFullName() +
14                 "");
15             Object[] nestedResultList = (Object[])call.invoke(
16                 new Object[]{"gov.nih.nci.cabio.domain.ws.Taxon", gene });
17             if (nestedResultList.length > 0)
18             {
19                 Taxon returnedTaxon = (Taxon)nestedResultList[0];
20                 System.out.println("\tTaxon: " + returnedTaxon.getScientificName());
21             }
22         }
23     }
24 } catch (Exception e) {
25     e.printStackTrace();
26 }
```

Lines	Description
15-16	A second operation invocation requests objects of type Taxon based on the same gene criteria used for the original query.
19	Casts the objects resulting from the nested query as Taxon objects.

Limitations

By default, the queryObject operation limits the result set to 1000 objects, even if the size of the result set is larger. To retrieve the objects past the 1000th record, you must use the query operation and specify the first object index (parameter in2) to be greater than 1000.

Result sets serialized and returned by the Web services interface do not currently include associations to related objects. A consequence of this behavior is that nested criteria objects with one-to-many associations that are passed to the query queryObject operations will result in an exception being thrown.

The following code demonstrates a Web services invocation that would fail:

```

1 Gene gene = new Gene();
2 gene.setSymbol("IL*");
3 Pathway pathway = new Pathway();
4 pathway.setId(new Long(120));
5 List pathwayList = new ArrayList();
6 pathwayList.add(pathway);
7 gene.setPathwaycollection(pathwayList);
8 try
9 {
10     Object[] resultList = (Object[])call.invoke(
11         new Object[]{"gov.nih.nci.cabio.domain.Gene", gene });
12 } catch (Exception e) {
13     // Web Services Exception will be caught
14 }

```

Because the Web services invocation has an inherent timeout behavior, queries that take a long time to execute may not complete. If this is the case, use the query method to specify a smaller result count.

XML-HTTP API

The caCORE XML-HTTP API, based on the REST (Representational State Transfer) architectural style, provides a simple interface using the HTTP protocol. In addition to its ability to be invoked from most internet browsers, developers can use this interface to build applications that do not require any programming overhead other than an HTTP client. This is particularly useful for developing web applications using AJAX (asynchronous JavaScript and XML).

Service Location and Syntax

The caBIO XML-HTTP interface uses the following URL syntax as described in *Table 4.3*.

```

http://{server}/{servlet}?query={returnClass}&{criteria}&
resultCounter={counter}&startIndex={index}&
pageSize={pageSize}&pageNumber={pageNumber}

```

Element	Meaning	Required	Example
server	Name of the web server on which the caBIO web application is deployed.	Yes	cabio.nci.nih.gov

Table 4.3 URL syntax used by the caBIO XML-HTTP interface

Element	Meaning	Required	Example
servlet	URI and the name of the servlet that will accept the HTTP GET requests	Yes	cabio41/server/ GetXML cabio41/server/ GetHTML
returnClasses	Class name indicating the type of objects that this query should return	Yes	query=Gene
criteria	Search request criteria describing the requested objects	Yes	Gene[@id=2]
counter	Number of top level objects returned by the search	No	resultCounter=500
index	Start index of the result set	No	startIndex=25
pageSize	Number of records to display on each "page"	No	pageSize=50
pageNumber	The number of the "page" for which to display results	No	pageNumber=3

Table 4.3 URL syntax used by the caBIO XML-HTTP interface

The caBIO architecture currently provides two servlets that accept incoming requests:

- **GetXML** returns results in an XML format that can be parsed and consumed by most programming languages and many document authoring and management tools
- **GetHTML** presents result using a simple HTML interface that can be viewed by most modern Internet browsers

Within the request string of the URL, the criteria element specifies the search criteria using XQuery-like syntax (Table 4.4). Within this syntax, square brackets ("[" and "]") represent attributes and associated roles of a class, the "at" symbol ("@") signals an attribute name/value pair, and a forward slash character ("/") specifies nested criteria. Criteria statements within XML-HTTP queries are generally of the following forms (although more complex statements can also be formed):

```
{ClassName}[@{attributeName}={value}] [@{attributeName}={value}]...
{ClassName}[@{attributeName}={value}]/
{ClassName}[@{attributeName}={value}]/...
```

Parameter	Meaning	Example
ClassName	The name of a class	Gene
attributeName	The name of an attribute of the return class or an associated class	symbol
Value	The value of an attribute	brca*

Table 4.4 Criteria statements within XML-HTTP queries

Examples of Use

The following examples demonstrate use of the XML-HTTP interface. In actual use, the queries shown here would either be submitted by a block of code or entered in the address bar of an Internet browser. Also note that the servlet name `GetXML` in each of the examples can be replaced with `GetHTML` to view with layout and markup in a browser.

Query	http://server/servlet/GetXML?query=Gene&Gene[@symbol=brca*]
Syntactic Meaning	Find all objects of type <code>Gene</code> whose symbol starts with 'brca'.
Biological Meaning	Find all BRCA genes.

Query	http://server/servlet/GetXML?query=Gene&Gene[@symbol=brca*]/Taxon[@scientificName=homo sapiens]
Syntactic Meaning	Find all objects of type <code>Gene</code> whose symbol starts with 'brca' and which have an associated <code>Taxon</code> object whose <code>scientificName</code> is equal to 'homo sapiens'.
Biological Meaning	Find all human BRCA genes.

Query	http://server/servlet/GetXML?query=Tissue&Tissue[@organ=eye][@histology=neoplasia]
Syntactic Meaning	Find all objects of type <code>Tissue</code> associated with attribute <code>organ</code> equal to 'eye' and <code>histology</code> equal to 'neoplasia'.
Biological Meaning	Find all tissues representing neoplasms of the eye.

Query	http://server/servlet/GetXML?query=Gene&Chromosome[@number=2]/Taxon[@scientificName=homo sapiens]
Syntactic Meaning	Find all objects of type <code>Gene</code> associated with <code>Chromosome</code> objects with number equal to 2 which themselves are related to <code>Taxon</code> objects with <code>scientificName</code> equal to 'homo sapiens' in biological terms
Biological Meaning	Find all human genes on chromosome number 2.

Working With Result Sets

Because HTTP is a stateless protocol, the caBIO server has no knowledge of the context of any incoming request. Consequently, each invocation of `GetXML` or `GetHTML` must contain all of the information necessary to retrieve the request, regardless of previous requests. Developers should consider this when working with the XML-HTTP interface.

Retrieving Related Results using XLinks

When using the GetXML servlet to retrieve results as XML, associations between objects are converted to XLinks within the XML. The link notation, shown below, allows the client to make a subsequent request to retrieve the associated objects.

```
<class name="gov.nih.nci.cabio.domain.Gene" recordNumber="1">
  ...
  <field name="taxon"
    xlink:type="simple"
    xlink:href=
"http://cabioapi.nci.nih.gov/cabio41/GetXML?query=Taxon&Gene[@id=5]">
    getTaxon
  </field>
  ...
</class>
```

Controlling the Number of Items Returned

The GetXML servlet provides a throttling mechanism to allow developers to define the number of results returned on any single request and where in the result set to start. For example, if a search request yields 500 results, specifying `resultCounter=450` returns only the last 50 records. Similarly, specifying `startIndex=50` returns only the first 50 records.

Paging Results

In addition to controlling the number of results to display, the GetXML servlet also provides a mechanism to support "paging". This concept, common to many web sites, allows results to be displayed over a number of pages, so that, for example, a request that yields 500 objects could be displayed over 10 pages of 50 objects each. When the paging feature is used, the GetXML servlet will include XLinks to each of the result pages in an XML `<page/>` element. The element data of the `<page/>` element is the number of the page, suitable for output as text or HTML when using an XSL stylesheet:

```
<page number="1"
  xlink:type="simple"
  xlink:href="http://cabioapi.nci.nih.gov/cabio41/GetXML?query=
{query}&pageNumber=4&resultCounter=1000&startIndex=0"> 4 </page>
```

Limitations

When specifying attribute values in the query string, use of the following characters generates an error: `[] / \ # & %`.

Grid Service API

The caBIO Grid Service API provides access to caBIO data through the caBIG[®] grid (caGrid).

Service Location

caBIO Grid Services can be accessed via the following URLs for each caBIO version:

V3.1	http://cabio-gridservice.nci.nih.gov:80/wsrf/services/cagrid/CaBIOSvc
V3.2	http://cabiogrid32.nci.nih.gov:80/wsrf/services/cagrid/CaBIO32GridSvc
V4.0	http://cabiogrid40.nci.nih.gov:80/wsrf/services/cagrid/CaBIO40GridSvc

These grid services can also be accessed through the caGrid Portal.

Examples of Use

caBIO grid services can be invoked using the caGrid Common Query Language (CQL). Below are examples of CQL queries for caBIO data.

Example One: Return all Genes

```
<CQLQuery xmlns="http://CQL.caBIG/1/gov.nih.nci.cagrid.CQLQuery">
  <Target name="gov.nih.nci.cabio.domain.Gene">
  </Target>
</CQLQuery>
```



Example Two: Return all Genes with a Symbol Beginning with "BRCA"

```
<CQLQuery xmlns="http://CQL.caBIG/1/gov.nih.nci.cagrid.CQLQuery">
  <Target name="gov.nih.nci.cabio.domain.Gene">
    <Attribute name="symbol" predicate="LIKE" value="BRCA%"/>
  </Target>
</CQLQuery>
```



Example Three: Return all Genes with a Symbol Beginning with "BRCA" and an Associated Taxon with a Scientific Name Equal to "Homo sapiens"

```
<CQLQuery xmlns="http://CQL.caBIG/1/gov.nih.nci.cagrid.CQLQuery">
  <Target name="gov.nih.nci.cabio.domain.Gene">
    <Group logicRelation="AND">
      <Attribute name="symbol" predicate="LIKE" value="BRCA%"/>
      <Association roleName="taxon"
name="gov.nih.nci.cabio.domain.Taxon">
        <Attribute name="scientificName" predicate="EQUAL_TO"
value="Homo sapiens"/>
      </Association>
    </Group>
  </Target>
</CQLQuery>
```



For examples on how to invoke caGrid Data Services via the Java programming language, Please refer to Chapter 6 of the caGrid 1.1 Programmer's Guide. Detailed instructions and examples for discovering and accessing caBIO grid services, are also available on the caGrid Wiki:

Working With Result Sets

Results from CQL queries are returned as XML. Below is an example of the CQL query results from a caBIO gene query:

```
<ns1:CQLQueryResults targetClassname="gov.nih.nci.cabio.domain.Gene"
  xmlns:ns1="http://CQL.caBIG/1/gov.nih.nci.cagrid.CQLResultSet">
- <ns1:ObjectResult>
  <ns2:Gene bigid="hdl://2500.1.PMEUQUCCCL5/HEGOZJZIY3" clusterId="2"
    fullName="N-acetyltransferase 2 (arylamine N-acetyltransferase)"
    symbol="NAT2" id="2"
    xmlns:ns2="gme://caCORE.caBIO/4.0/gov.nih.nci.cabio.domain" />
  </ns1:ObjectResult>
- <ns1:ObjectResult>
  <ns3:Gene bigid="hdl://2500.1.PMEUQUCCCL5/S344L2X73O" clusterId="4"
    fullName="Alcohol dehydrogenase 1B (class I), beta polypeptide"
    symbol="ADH1B" id="3"
    xmlns:ns3="gme://caCORE.caBIO/4.0/gov.nih.nci.cabio.domain" />
  </ns1:ObjectResult>
- <ns1:ObjectResult>
  <ns4:Gene bigid="hdl://2500.1.PMEUQUCCCL5/I62AZYNECK"
    clusterId="11" fullName="Carcinoembryonic antigen-related cell
    adhesion molecule 3" symbol="CEACAM3" id="4"
    xmlns:ns4="gme://caCORE.caBIO/4.0/gov.nih.nci.cabio.domain" />
  </ns1:ObjectResult>

  ...


- <ns1:ObjectResult>
  <ns1001:Gene bigid="hdl://2500.1.PMEUQUCCCL5/TWBNWHRUV4"
    clusterId="30213" fullName="Ceroid-lipofuscinosis, neuronal 5"
    symbol="CLN5" id="1544"
    xmlns:ns1001="gme://caCORE.caBIO/4.0/gov.nih.nci.cabio.domain"
    />
  </ns1:ObjectResult>
</ns1:CQLQueryResults>
```

Limitations


caBIO Grid services for a release are deployed after the caBIO API release has been deployed to Production to allow for time required to generate, test, and deploy new grid services for the release. To obtain the latest version of caBIO grid services, leverage the caBIO Grid Portal to discover caBIO services.

caBIO FreeStyleLM and Range Query search facilities are currently not available as grid services.

CHAPTER 5 CABIO UTILITIES

This chapter describes additional utilities provided by the caBIO API, which are not generated by the caCORE SDK. These customizations (this sentence fragment is in the Word doc--complete it?) 

Topics in this chapter include:

- *FreestyleLM* on this page
- *Grid Id Resolver* on page 56
- *Genome Range Queries* on page 56
- *Array Annotation API* on page 58
- *SVG Manipulator* on page 63 

FreestyleLM

All of the search methods discussed in the previous chapter allow for vertical searching in one caBIO domain object. In other words, one domain object must be selected to be searched. The FreestyleLM (Freestyle Lexical Mine) search component provides interfaces and APIs for conducting horizontal searches across caBIO domain objects.

FreestyleLM is built upon Hibernate and Apache Lucene full text search engines, providing full text search (Google™-like search) capabilities to the caBIO API. The classes and fields that are indexed are annotated in the caBIO domain model. Usually, these attributes are non-numeric fields with descriptions or identifiers. Thus, searching for “brca1” will return any object that makes mention of the “brca1” gene in its attributes, such as Genes, Proteins, Pathways, and so on.

The IndexGenerator utility within the caBIO API generates index files for these annotations. FreestyleLM provides the ability to query these indexes. The FreestyleLM search is facilitated by the classes listed below:

- `gov.nih.nci.search.SearchQuery`

- `gov.nih.nci.search.SearchResults`
- `gov.nih.nci.search.Sort`
- `gov.nih.nci.search.RangeFilter`

The search term is set as a keyword within the `SearchQuery` object. One or more words can be set as a keyword. Refer to the Lucene documentation for more information on the search syntax (<http://lucene.apache.org/java/docs/queryparsersyntax.html>).

Freestyle LM Java Client

The `CaBioApplicationService` provides a special data access method to perform the text based search.

Method Signature	<code>List search(SearchQuery searchQuery)</code>
Search Type	FreestyleLM Search
Description	Returns a List collection containing objects conforming to the criteria defined by the searchQuery.
Example	<code>Search(searchQuery);</code>

Example One: Full Text Search (Default)

In this example, a search is performed on the term 'brca*'. The caBIO FreestyleLM returns a set of `SearchResults` objects. The result set includes references to objects of different types that match the search term. In this case, the search results consist of information regarding Gene, Protein, Pathway, Tissue, GeneOntology and Library objects.

```

1 CaBioApplicationService appService =
  (CaBioApplicationService)ApplicationServiceProvider.getApplicationService();
2 try {
3     SearchQuery query = new SearchQuery();
4     query.setKeyword("brca*");
5     List results = appService.search(query);
6     for(int i=0; i<results.size(); i++) {
7         SearchResult result = (SearchResult)results.get(i);
8         System.out.println("Class: " + result.getClassName() + "\t" + result.getId());
9     }
10    System.out.println("Results: " + results.size());
11 } catch(Exception e) {
12     e.printStackTrace();
13 }
```

<i>Lines</i>	<i>Description</i>
1	Create a <code>CaBioApplicationService</code> instance.
3-4	Instantiate a <code>SearchQuery</code> object and set the keyword to 'brca*'
5	Call the <code>CaBioApplicationService</code> search method
6-9	Iterate through the results

Example Two: Hibernate Search

A search is performed on the term 'brca*' and the query type is specified as 'HIBERNATE_SEARCH'. (The default query type is `FULL_TEXT_SEARCH`). The caBIO FreestyleLM result set will include objects of different types that match the

search term. In this case, the search results consist of Gene, Protein, Pathway, Tissue, GeneOntology and Library objects.

```


1 CaBioApplicationService appService =
  (CaBioApplicationService)ApplicationServiceProvider.getApplicationService();
2 try {
3   SearchQuery query = new SearchQuery();
4   query.setKeyword("brca*");
5   query.setQueryType("HIBERNATE_SEARCH");
6   List results = appService.search(query);
7   for(int i=0; i<results.size(); i++){
8     Object result = (Object)results.get(i);
9     System.out.println("Class: " + result.getClass().getName());
10  }
11  System.out.println("Results: " + results.size());
12 } catch(Exception e){
13   e.printStackTrace();
14 }

1 CaBioApplicationService appService =
  (CaBioApplicationService)ApplicationServiceProvider.getApplicationService();
2 try {
3   SearchQuery query = new SearchQuery();
4   query.setKeyword("brca*");
5   Sort sorter = new Sort();
6   sorter.setSortByClassName(true);
7   query.setSort(sorter);
8   List results = appService.search(query);
9   for(int i=0; i<results.size(); i++){
10    SearchResult result = (SearchResult)results.get(i);
11    System.out.println("Class: " + result.getClassName() + "\t" + result.getId());
12  }
13  System.out.println("Results: " + results.size());
14 } catch(Exception e) {
15   e.printStackTrace();
16 }

```

Freestyle LM Web Interface

The caBIO API provides a simple web interface to perform text based queries. This user interface is built upon the FreestyleLM Search API (*Figure 5.1*). The following URL can be used to access the servlet: <http://cabioapi.nci.nih.gov/cabio41/search>



National Cancer Institute

Center for Bioinformatics
FREESTYLE LEXICAL MINE

Search

☐ match all words
☐ match any of the words
☐ exclude word
☒ Full Text Search
☐ Fuzzy Search
 Number of records per page 1000

Figure 5.1 Freestyle Lexical Mine interface

Grid Id Resolver

Each caBIO object can be uniquely identified (even across classes) using the Grid Id (bigId attribute). See *caGrid Identifiers* on page 87 for more details about Grid Identifiers. The API discussed here allows for simple retrieval of objects using these identifiers.

Grid Id Resolver Java API

The `CaBioApplicationService` provides special data access methods to perform Grid Id searches. The same methods are available in the web services interface.

Method Signature	<code>boolean exist(String bigId)</code>
Search Type	Grid Id Resolver
Description	Check if the object specified by the given Grid Id exists in caBIO.
Example	<pre>boolean exists = exist("hdl://2500.1.PMEUQUCCCL5/WX2PXQD7PB")</pre>

Method Signature	<code>boolean getDataObject(String bigId)</code>
Search Type	Grid Id Resolver
Description	Retrieve the domain object specified by the given Grid Id.
Example	<pre>Object obj = getDataObject("hdl://2500.1.PMEUQUCCCL5/WX2PXQD7PB")</pre>

Genome Range Queries

The Range Query API provides an alternative way to search for genomic features on a chromosome, based on the physical locations. The API is capable of performing arbitrary range searches on any given chromosome (absolute search), or range searches around a feature of interest (relative search).

The caBIO model contains a subclass of `PhysicalLocation` for every object type which has `PhysicalLocation` data. For example, there is a `GenePhysicalLocation` class for Gene location data. The Range Query API takes advantage of these classes by allowing the user to specify which subclass of `PhysicalLocation` to search for. If `SNPPhysicalLocation` is specified as the target class, then only `PhysicalLocations` pertaining to SNPs will be returned.

Range Query Java API

The `CaBioApplicationService` provides a special data access methods to perform range queries.

Method Signature	<code>List search(RangeQuery rangeQuery)</code>
Search Type	Genome Range Search

Description	Returns a List collection containing objects conforming to the criteria defined by the searchQuery.
Example	<code>search(rangeQuery);</code>

Method Signature	<code>List search(Class targetClass, RangeQuery rangeQuery)</code>
Search Type	Genome Range Search
Description	Returns a List collection containing objects of the type specified by the targetClass, conforming to the criteria defined by the given RangeQuery.
Example	<code>search(SNP.class, rangeQuery);</code>

Example One: Absolute Search

In this example, a search for SNPs on an arbitrary range of the human chromosome 1 is performed. The result set will only include SNP in the given range. By default, the reference assembly is searched. An alternative assembly may be specified with the `rangeQuery.setAssembly()` method.

```

1 CaBioApplicationService appService =
  (CaBioApplicationService)ApplicationServiceProvider.getApplicationService();

2 try {
3     Taxon taxon = new Taxon();
4     taxon.setAbbreviation("Hs");
5     Chromosome chromosome = new Chromosome();
6     chromosome.setNumber("1");
7     chromosome.setTaxon(taxon);
8     chromosome = (Chromosome)appService.search(
        Chromosome.class, chromosome).iterator().next();

9     AbsoluteRangeQuery query = new AbsoluteRangeQuery();
10    query.setChromosome(chromosome);
11    query.setStart(new Long(109433764));
12    query.setEnd(new Long(109434982));

13    List results = appService.search(SNPPhysicalLocation.class, query);
14    for(int i=0; i<results.size(); i++) {
15        SNPPhysicalLocation l = (SNPPhysicalLocation)results.get(i);
16        System.out.println("SNP: " + l.getSNP().getDBSNPID());
17    }

18 } catch(Exception e) {
19     e.printStackTrace();
20 }
```

Lines	Description
1	Create a CaBioApplicationService instance.
3-8	Retrieve the chromosome of interest
9-12	Create and configure the AbsoluteRangeQuery
13	Execute the RangeQuery with SNPPhysicalLocation as the target
14-17	Iterate through the results

Example Two: Relative Search

In this example, a search for SNPs on an arbitrary range of the human chromosome 1 is performed. The result set will only include SNP in the given range. By default, the

reference assembly is searched. An alternative assembly may be specified with the `rangeQuery.setAssembly()` method.

```

1 CaBioApplicationService appService =
    (CaBioApplicationService)ApplicationServiceProvider.getApplicationService();

2 try {
3     SNP snp = new SNP();
4     snp.setDBSNPID("rs10873500");
5     snp = (SNP)appService.search(SNP.class, snp).get(0);

6     FeatureRangeQuery query = new FeatureRangeQuery();
7     query.setFeature(snp);
8     query.setUpstreamDistance(new Long(1000));
9     query.setDownstreamDistance(new Long(1000));
10    List<ArrayReporterPhysicalLocation> results = appService.search(
11        ArrayReporterPhysicalLocation.class, query);

12    for(ArrayReporterPhysicalLocation l : results) {
13        System.out.println("Reporter: " + l.getArrayReporter().getName());
14    }
15 } catch(Exception e) {
16     e.printStackTrace();
17 }

```

<i>Lines</i>	<i>Description</i>
1	Create a CaBioApplicationService instance.
3-5	Retrieve the SNP of interest
6-9	Create and configure the FeatureRangeQuery
10-11	Execute the RangeQuery with ArrayReporterPhysicalLocation as the target
14-17	Iterate through the results

Array Annotation API

This client-side convenience API is intended for users of the Microarray annotations in caBIO. The API simplifies access to bulk microarray annotations such as reporters, genes and SNPs. The `ArrayAnnotationService` class provides the following Java methods:

Method Signature	Map<String, Collection<String>> getGenesForReporters(String arrayPlatform, Collection<String> reporterIds)
Description	Given an array and list of reporter names, return a mapping of those reporter names to HUGO gene symbols.
Example	getGenesForReporters("HG-U133_Plus_2", reporterIds)

Method Signature	Map<String, Collection<String>> getReportersForGenes(String arrayPlatform, Collection<String> geneSymbols)
Description	Given an array and a list of HUGO gene symbols, return a mapping of those gene symbols to reporters on the specified array.
Example	getReportersForGenes("HG-U133_Plus_2", geneSymbols)

Method Signature	<code>Collection<ArrayReporter> getReportersForPlatform(String arrayPlatform)</code>
Description	Returns all ArrayReporters on a given array platform.
Example	<code>getReportersForPlatform("HG-U133_Plus_2")</code>

Method Signature	<code>Collection<ExpressionArrayReporter> getExpressionReporterAnnotations(String arrayPlatform, Collection<String> reporterIds)</code>
Description	Given an array and list of reporter names, return the ExpressionArrayReporter objects corresponding to the reporter names, populated with the following annotations: <ul style="list-style-type: none"> • gene • physicalLocationCollection.chromosome • cytogeneticLocationCollection • cytogeneticLocationCollection.cytoband
Example	<code>getExpressionReporterAnnotations("HG-U133_Plus_2", reporterIds)</code>

Method Signature	<code>Collection<ExonArrayReporter> getExonReporterAnnotations(String arrayPlatform, Collection<String> reporterIds)</code>
Description	Given an array and list of reporter names, return the ExonArrayReporter objects corresponding to the reporter names, populated with the following annotations: <ul style="list-style-type: none"> • geneCollection • physicalLocationCollection • physicalLocationCollection.chromosome • cytogeneticLocationCollection • cytogeneticLocationCollection.cytoband
Example	<code>getExonReporterAnnotations("HuEx-1_0-st-v2 ", reporterIds);</code>

Method Signature	<code>Collection<SNPArrayReporter> getSNPReporterAnnotations(String arrayPlatform, Collection<String> reporterIds)</code>
Description	Given an array and list of reporter names, return the SNPArrayReporter objects corresponding to the reporter names, populated with the following annotations: SNP <ul style="list-style-type: none"> • physicalLocationCollection • physicalLocationCollection.chromosome • cytogeneticLocationCollection • cytogeneticLocationCollection.cytoband
Example	<code>getSNPReporterAnnotations("Mapping250K_Sty", reporterIds)</code>

Method Signature	<code>Collection<Gene> getGeneAnnotations(Collection<String> geneSymbols)</code>
Description	Returns a collection of the requested genes, populated with the following annotations: <ul style="list-style-type: none"> • <code>cytogeneticLocationCollection</code> • <code>cytogeneticLocationCollection.startCytoband</code> • <code>databaseCrossReferenceCollection</code> • <code>chromosome</code>
Example	<code>getGeneAnnotations(geneSymbols)</code>

Method Signature	<code>List<CytobandPhysicalLocation> getCytobandPositions(String chromosomeNumber)</code>
Description	Returns all the cytobands on the specified chromosome (found on the reference assembly).
Example	<code>getCytobandPositions("1")</code>

Method Signature	<code>List<CytobandPhysicalLocation> getCytobandPositions(String chromosomeNumber, String assembly)</code>
Description	Returns all the cytobands on the specified chromosome of the given assembly.
Example	<code>getCytobandPositions("1", "Celera")</code>

Method Signature	<code>Microarray getMicroarray(String platformName)</code>
Description	Returns the Microarray with the given name.
Example	<code>getMicroarray("HG-U133_Plus_2")</code>

Method Signature	<code>getGenesForSymbol(String hugoSymbol)</code>
Description	Returns the genes with the given HUGO symbol.
Example	<code>getGenesForSymbol("IL5")</code>

Method Signature	<code>Collection<GeneAlias> getAliasesForGene(String symbol)</code>
Description	Returns the aliases associated with the given gene symbol.
Example	<code>getAliasesForGene("IL5")</code>

Method Signature	<code>HashMap<String, Collection<String>> getReportersForSnps(String arrayPlatform, List<String> dbSnpsIds)</code>
Description	Returns all reporter names for the specified SNP on the given array.
Example	<code>getReportersForSnps("Mapping250K_Sty", dbSnpsIds)</code>

Method Signature	<code>Collection<SNP> getSNPAnnotations(List<String> dbSnpsIds)</code>
Description	Returns the specified SNP preloaded with annotations. The associations which are preloaded: <ul style="list-style-type: none"> • <code>physicalLocationCollection</code> • <code>physicalLocationCollection.chromosome</code> • <code>cytogeneticLocationCollection</code> • <code>cytogeneticLocationCollection.startCytoband</code>
Example	<code>getSNPAnnotations(dbSnpsIds)</code>

Method Signature	<code>Collection<Gene> getGenesForDbSnpId(String dbSnpId)</code>
Description	Returns the genes associated with the given dBSNP id by array annotations (relative locations).
Example	<code>getGenesForDbSnpId("rs200577")</code>

Method Signature	<code>Collection<SNP> getSnpNearGene(String symbol, Long kbUpstream, Long kbDownstream, String assembly)</code>
Description	Returns all the SNPs within the specified range of the given gene on the specified assembly.
Example	<code>getSnpNearGene("A1CF", 0L, 0L, "Celera")</code>

Method Signature	<code>Collection<SNP> getSnpNearGene(String symbol, Long kbUpstream, Long kbDownstream)</code>
Description	Returns all the SNPs within the specified range of the given gene, on the reference assembly.
Example	<code>getSnpNearGene("A1CF", 1000L, 1000L)</code>

Example One: Retrieve Gene Symbols for Reporters on the HGU-133 Array

In this example, a search for SNPs on an arbitrary range of the human chromosome 1 is performed. The result set will only include SNP in the given range. By default, the

reference assembly is searched. An alternative assembly may be specified with the `rangeQuery.setAssembly()` method.

```

1 CaBioApplicationService appService =
  (CaBioApplicationService)ApplicationServiceProvider.getApplicationService();

2 ArrayAnnotationService aaService = new ArrayAnnotationServiceImpl(appService);

3 List<String> list = new ArrayList<String>();
4 list.add("1554206_at");
5 list.add("1555490_s_at");
6 list.add("1556516_at");

7 Map<String, Collection<String>> results =
8   aaService.getGenesForReporters("HG-U133_Plus_2", list);

9 for(String reporter : results.keySet()) {
10   System.out.println(reporter+" : "+results.get(reporter).iterator().next());
11 }

```

<i>Lines</i>	<i>Description</i>
1	Create a CaBioApplicationService instance.
2	Create the ArrayAnnotationService instance.
3-6	Create a list of reporter names.
7-8	Execute the convenience method
9-11	Iterate through the results

Example Two: Retrieve SNPs with Pre-loaded Annotations

In this example, a search for SNPs on an arbitrary range of the human chromosome 1 is performed. The result set will only include SNP in the given range. By default, the reference assembly is searched. An alternative assembly may be specified with the `caBIO 4.2 Technical Guide` method.

```

1 CaBioApplicationService appService =
  (CaBioApplicationService)ApplicationServiceProvider.getApplicationService();

2 ArrayAnnotationService aaService = new ArrayAnnotationServiceImpl(appService);

3 List<String> list = new ArrayList<String>();
4 list.add("rs16893403");
5 list.add("rs1598492");
6 list.add("rs6768736");

7 Collection<SNP> results = aaService.getSNPAnnotations(list);

8 for(SNP snp : results) {
9   PhysicalLocation pl = snp.getPhysicalLocationCollection().iterator().next();
10   System.out.println(snp.getDBSNPID()+" on chromosome
11   "+pl.getChromosome().getNumber());

```

<i>Lines</i>	<i>Description</i>
1	Create a CaBioApplicationService instance.
2	Create the ArrayAnnotationService instance.
3-6	Create a list of SNP ids
7	Execute the convenience method to retrieve the SNP annotations

<i>Lines</i>	<i>Description</i>
8	Iterate through the results
9-10	Follow preloaded associations. None of these method calls produce additional queries.

SVG Manipulator

In the `gov.nih.nci.common.util` package, the caBIO API includes a class (`SVGManipulator`) that provides useful services to manipulate Scalable Vector Graphics (SVG) diagrams retrieved from the caBIO Pathway domain object.

Manipulating Pathway Diagrams

BioCarta and its Proteomic Pathway Project (P3) provide detailed graphical renderings of pathway information. NCI's CMAP web site captures pathway information from BioCarta, and transforms the downloaded image data into Scalable Vector Graphics (SVG) representations that support interactive manipulation of the online images. The `SVGManipulator` utility class provides the capability to do the following:

- Change the display colors for each gene contained in an SVG diagram.
- Modify the URL linking a gene in the SVG diagram to external gene information. The default gene URL links to the CMAP website.
- Disable all genes or nodes within an SVG diagram.
- Retrieve a gene's color.
- Reset a gene or node to its original state.
- Retrieve/set SVG diagram attributes via getter/setter methods

Figure 5.2 shows an example of how to use SVGManipulator class to modify content of an SVG diagram replace with original image (resolution here is horrible).

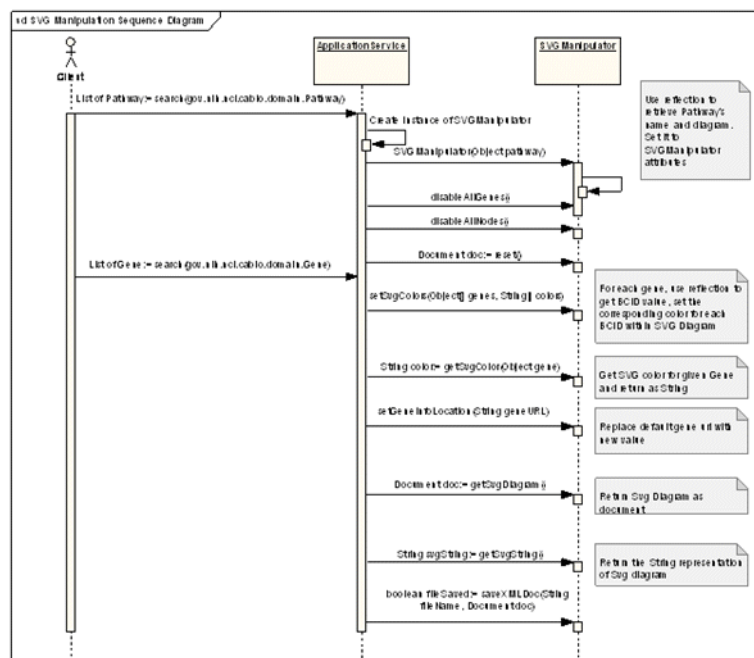


Figure 5.2 Sequence diagram using caBIO Pathway and Gene objects

SVG Diagram Manipulation Utility Example

The following test client demonstrates how to use the SVGManipulation class to change the appearance of a pathway diagram associated with a given pathway.

```

1  CaBioApplicationService appService =
    (CaBioApplicationService)ApplicationServiceProvider.getApplicationService();

2  Pathway pw = new Pathway();
3  pw.setName("h_freePathway");
4  List results = appService.search(Pathway.class, pw);
5  Pathway returnedPw = (Pathway) results.get(0);
6  SVGManipulator svgM = new SVGManipulator(returnedPw);

7  Gene[] genes = new Gene[2];

8  Gene gene1 = new Gene();
9  gene1.setClusterId(new Long(443914)); // SOD1
10 List resultList1 = appService.search(Gene.class, gene1);
11 genes[0] = (Gene) resultList1.get(0);

12 Gene gene2 = new Gene();
13 gene2.setClusterId(new Long(241570)); // TNF
14 List resultList2 = appService.search(Gene.class, gene2);
15 genes[1] = (Gene) resultList2.get(0);

16 String[] colors = new String[2];
17 colors[0] = "255,0,0"; // SOD1 = red
18 colors[1] = "0,0,255"; // TNF = blue
19 svgM.setSvgColors(genes, colors);

20 svgM.saveXMLDoc("out.svg", svgM.getSvgDiagram());
  
```


<i>Lines</i>	<i>Description</i>
2-6	Retrieve "h_freePathway" and create a new SVManipulator instance to operate on it.
7-15	Retrieve genes of interest in the pathway
16-18	Create a parallel list of colors for the genes
19	Set the colors for the genes
20	Save the modified SVG image



CANCER BIOINFORMATICS INFRASTRUCTURE OBJECTS

This chapter describes the Cancer Bioinformatics Infrastructure Objects (caBIO) model and its application programming interfaces.

Topics in this chapter include:

- Introduction on this page
- caBIO API on this page
- Data Sources in the caBIO Database on page 78
- caGrid Identifiers on page 85

Introduction

The Cancer Bioinformatics Infrastructure Objects (caBIO) model and architecture was the first of several model-driven information systems that make up caCORE and continues to be an on-going effort to model the genomic domain. The caBIO objects simulate the behavior of actual genomic components such as genes, chromosomes, sequences, libraries, clones, ontologies, etc. They provide access to a variety of genomic data sources including GenBank, Unigene, LocusLink, Homologene, Ensemble, UCSC Genome Sequencing Center's Annotations, NCI's CTEP (Cancer Therapy Evaluation Program) and NCICB's CGAP (Cancer Genome Anatomy Project) data repositories. The full list of data sources is listed starting on [page 81](#).

caBIO API

Most of the domain objects defined in the caBIO API are objects that specialize in bioinformatics applications. The caBIO domain objects are implemented as Java beans in the `gov.nih.nci.cabio.domain`, `gov.nih.nci.common.domain`, `gov.nih.nci.cabio.pathway`, and

`gov.nih.nci.common.provenance.domain` packages and include those classes that correspond to biological entities and bioinformatics concepts.

The caBIO UML model is published as an EA (Enterprise Architect) in the Files section of the caBIO GForge site (https://gforge.nci.nih.gov/frs/?group_id=51). Table 6.1 and Table 6.2 list and describe the classes in these packages. Detailed descriptions about each class and its methods are present in the caBIO 4.2 JavaDocs ([Add a link to the Javadocs?](#)).

All caBIO Domain Objects in the following table that are not hyperlinked (lacking blue underlined formatting) have “broken” (inactive) links in the MS Word source document (even though the formatting in the Word doc appears to be hyperlinked).

caBIO Domain Object	Description
Agent	<p>An object representing a therapeutic agent (drug, intervention therapy) used in a clinical trial protocol. It provides the NSC Number, name and source of the agent.</p> <p>An Agent provides access to the associated protocols used in clinical trials (<code>ClinicalTrialProtocol</code>) and associated genes targeted (<code>Targets</code>) by that agent via the <code>getClinicalTrialProtocolCollection</code> and <code>getTargetCollection</code> functionalities respectively.</p> <p>While the data in the Agent object is largely static, it has been updated to include Agents studied as part of the Cancer Gene Data Curation project at NCI.</p>
Anomaly	<p>An irregularity in either the expression of a gene or its structure (i.e., a mutation). It provides access to the associated histopathological information (<code>Histology</code>) as well as the disease (<code>DiseaseOntology</code>) and organ ontology (<code>OrganOntology</code>).</p> <p>A user can also retrieve the anomalies (<code>Anomaly</code>) associated with a given gene target using the <code>getAnomalyCollection</code> functionality.</p> <p>The data in Anomaly is currently static.</p>

Table 6.1 caBIO domain objects and descriptions

caBIO Domain Object	Description
ArrayReporter	<p>A superclass for different kinds of microarray reporters such as <code>ExonArrayReporter</code>, <code>SNPArrayReporter</code>, <code>TranscriptArrayReporter</code> and <code>ExpressionArrayReporter</code>.</p> <p>These objects get populated with annotations from Affymetrix, Illumina and Agilent Arrays such as the Affymetrix HG-U133, HG-U133 Plus2, HT-HG-U133 A, HT-HG-U133 A2.0 Arrays, HuMapping 250kNsp, HuMapping 250kSty, HuMapping 50kHind240, HuMapping Xba240 and Exon Arrays, Agilent 44K and Acgh 244k arrays and Illumina HumanHap Arrays and are synchronized monthly with the latest releases of the respective datasets.</p> <p>In case of SNP Arrays such as Affymetrix HuMapping and Illumina arrays, the location of a SNP relative to a Marker is available through <code>MarkerRelativeLocation</code> and relative through a Gene is available through <code>GeneRelativeLocation</code>.</p>
ArrayReporterCytogenetic Location	<p>Provides cytogenetic starts and stops for array reporters in <code>ArrayReporter</code>, where available. It exposes cytogenetic starts and stops for Affymetrix HG-U133, HG-U133 Plus2, HT-HG-U133A, HG-HG-U133A 2.0, HuMapping 250k Nsp, HuMapping 250k Sty, HuMapping 50k Hind240, HuMapping 50k Xba 240, Agilent 44K and Agilent CGH 244K arrays and gets populated with data provided by the manufacturer.</p>
ArrayReporterPhysical Location	<p>Provides chromosomal starts and stops for array reporters in <code>ArrayReporter</code>. It exposes chromosomal starts and stops for Affymetrix HG-U133, HG-U133 Plus2, HT-HG-U133A, HG-HG-U133A 2.0, HuMapping 250k Nsp, HuMapping 250k Sty, HuMapping 50k Hind240, HuMapping 50k Xba 240, Exon Arrays, Agilent 44K and Agilent aCGH 244K arrays and Illumina Human Hap Arrays and gets populated with data provided by the manufacturer.</p>
BiochemicalReaction	<p>A type of pathway interaction representing an interaction among physical entities, and treated as a single event. It is populated with data from PID.</p>
Chromosome	<p>An object representing a specific chromosome in an organism; It provides access to all known genes on the chromosome, as well as to the <code>ChromosomalLocation</code> and <code>CytogeneticLocation</code> of features such as SNP, ESTs, mRNAs, Markers, ArrayReporters, etc through the <code>getLocationCollection</code> functionality.</p> <p>This ability to associate multiple features with a chromosome aids in performing range-query searches, described in greater detail elsewhere.</p> <p>The data in the chromosome table is static and as with data in caBIO, encompasses only human and mouse chromosomes.</p>

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
ClinicalTrialProtocol	<p>The protocols and administrative information about the trial such as lead organization, participants, current phase of the clinical trial, and associated NIH Administrative Id.</p> <p>A <code>ClinicalTrialProtocol</code> object provides the user access to the associated agents, histopathology and disease ontologies through the <code>getAgentCollection</code>, <code>getDiseaseOntologyCollection</code> and <code>getHistopathologyCollection</code> functionalities respectively. This object gets populated with data from NCI's CTEP Program which gets updated on a monthly basis.</p>
Clone	<p>An object used to hold information pertaining to I.M.A.G.E. clones associated with human and mouse genes from Unigene's repository. A user can access the associated sequences (<code>NucleicAcidSequence</code>) using the <code>getNucleicAcidSequenceCollection</code> functionality.</p> <p>A clone is always associated with a clone-library further details of which are available through the associated Library Object.</p> <p>The orientation of the clone along with the associated sequence identifier is available through the <code>CloneRelativeLocation</code> object that provides further information regarding the clone insert read.</p> <p>The Clone object gets populated with clone data from Unigene's human and mouse annotations and that can be associated to existing libraries in caBIO's Library Object.</p> <p>As indicated below, the Library object in turn gets populated with clone-library data from CGAP who integrate information from ORESTES, dbEST, Unigene and other such resources.</p>
CloneRelativeLocation	<p>Provides the orientation (5'/3') of a clone insert read, along with links to the associated sequence information (<code>NucleicAcidSequence</code>).</p> <p>It gets populated with relative location of clones in the Clone Object as available in Unigene.</p>
ComplexComponent	<p>A pathway participant (subclass of <code>PhysicalParticipant</code>) which is part of a complex (<code>ComplexEntity</code>). It is populated with data from PID.</p>
ComplexEntity	<p>A pathway participant (subclass of <code>PhysicalParticipant</code>) which is part of a complex (<code>ComplexEntity</code>). It is populated with data from PID.</p> <p>A type of pathway entity (subclass of <code>PhysicalEntity</code>) which is composed of pathway participants (<code>ComplexComponents</code>). It is populated with data from PID.</p>

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
Condition	<p>A pathway participant (subclass of <code>PhysicalParticipant</code>) which is part of a complex (<code>ComplexEntity</code>). It is populated with data from PID.</p> <p>A type of pathway entity (subclass of <code>PhysicalEntity</code>) which is composed of pathway participants (<code>ComplexComponents</code>). It is populated with data from PID.</p> <p>A kind of pathway participant (subclass of <code>Participant</code>) representing a prerequisite condition for an interaction to occur. It is populated with data from PID.</p>
Cytoband	<p>The <code>Cytoband</code> Object provides information such as the name and location of cytobands in human and mouse available from UCSC's Genome Sequencing Center for every new build of the Human and Mouse Genome.</p> <p>It also gets populated with human cytoband information available in Unigene, Affymetrix and Agilent array probeset to facilitate retrieval of cytobands that identify the start and stop of a SNP, Gene and <code>ArrayReporter</code> in <code>SNPCytogeneticLocation</code>, <code>GeneCytogeneticLocation</code> and <code>ArrayReporterCytogeneticLocation</code> respectively.</p>
<code>CytobandPhysicalLocation</code>	An abstract class that provides the chromosomal starts and stops of Cytobands. It gets populated with data from UCSC's Genome Sequence Center.
CytogeneticLocation	<p>An abstract superclass that provides cytogenetic starts and stops for different features such as SNP, Gene and <code>ArrayReporters</code> through their respective subclasses, <code>SNPCytogeneticLocation</code>, <code>GeneCytogeneticLocation</code> and <code>ArrayReporterCytogeneticLocation</code> respectively.</p> <p>Each of these subclasses gets populated with data from different sources as indicated below.</p>
<code>DatabaseCrossReference</code>	<p>An object that cross-references caBIO Genes, SNPs, and <code>Proteins</code> with identifiers from different databanks.</p> <p>For example, it associates caBIO genes (Unigene) to Entrez and OMIM, caBIO SNPs (dbSNP) to TSC and caBIO Proteins (Uniprot) to corresponding identifiers in Ensembl and Enzyme Commission.</p> <p>It uses references to Entrez-Gene and LocusLink in Unigene's human and mouse data to associate Unigene, and uses Entrez-based datasets to associate caBIO Genes with OMIM.</p> <p>It uses references to Ensembl and Enzyme Commission Ids from Affymetrix SNP array-annotations to associate caBIO Proteins with Ensembl, RefSeq and EC identifiers.</p>

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
DiseaseOntology	<p>An object representing various oncological disease ontologies. It provides access to the clinical trials associated with these diseases (<code>ClinicalTrialProtocol</code>) and the associated histopathological (<code>Histopathology</code>). Gene-Disease associations made through the Cancer Gene Index Project are available through the <code>getGeneFunctionAssociation</code> functionality.</p> <p>Apart from existing static data, this object gets populated with data from CGDC Data.</p>
DiseaseOntologyRelationship	<p>Provides a hierarchical relationship among various diseases and associated ontologies. Thus it allows a user to obtain <code>DiseaseOntologys</code> that share a parent/child relationship. The data in this table is static.</p>
EntityAccession	The accession number or native database id of an entity in a pathway. It is populated with data from PID.
EntityName	The name of an entity in a pathway. It is populated with data from PID.
Evidence	<p>Provides details of the scientific publication used to make the Gene-Agent or Gene-Disease association through text-mining, as part of NCI's CGDC Project.</p> <p>It provides details such as the Pubmed identifier, the sentence(s) from the publication highlighting the association, the status of the curation process, the NCI-provided Evidence Code assigned to this Evidence (that summarizes the role of the Gene in the disease or the interaction of the Agent with the Gene). It provides access to further details about the assigned EvidenceCode through the <code>getEvidenceCode</code> functionality.</p>
EvidenceCode	<p>Provides access to NCI-curated Evidence Codes to represent status of Gene-Disease and Gene-Agent associations made through the CGDC Project.</p> <p>It gets populated with data from CGDC and provides a link to the associated Evidence (PubMed publications used to form the Gene-Agent or Gene-Disease associations through text-mining)</p>
Exon	<p>An <code>Exon</code> represents an expressed part of a gene. In caBIO, an exon constitutes part of a <code>Transcript</code> that is used as a probe in Affymetrix Exon Arrays.</p> <p>It gets populated with data from Affymetrix Human Exon Arrays. Transcripts associated with this Exon, per Affymetrix Exon Array annotations are available through the <code>Transcript</code> object, while the probesets themselves are available through <code>ArrayReporter's ExonArrayReporter</code> object.</p>

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
ExonArrayReporter	<p>One of the two 'transcript-array' objects in caBIO, the other being <code>ExpressionArrayReporter</code>.</p> <p>An object representing Affymetrix Exon Arrays providing details such as the manufacturer's PSR Id, the probe count and the strand information for each probeset. A user can retrieve the associated Exons, Transcripts and Genes using the <code>getTranscriptCollection</code> and <code>getGeneCollection</code> methods respectively.</p>
ExpressedSequenceTag	<p>A subclass of <code>NucleicAcidSequence</code>. It represents short nucleotide sequences (typically 400-600 bases) from one or both ends of randomly selected cDNA clones and gets populated with data from Unigene's human and mouse datasets.</p> <p>It provides access to the same functionalities such as available through the parent <code>NucleicAcidSequence</code> class, such as clones, relative locations of these Clones, the chromosomal start and stop (<code>NucleicAcidPhysicalLocation</code>), corresponding Genes and hence any associated <code>ExpressionArray</code> reporters.</p>
ExpressionArrayReporter	<p>Gene objects representing mRNA-based arrays in caBIO. It provides information on probesets in the Affymetrix Human U133, HG U133 Plus2, HT HG U133 A and HT HG U133 A2.0 Arrays, and <code>NucleicAcidSequence</code> respectively.</p> <p>A user can retrieve the associated Genes and protein domains (<code>ProteinDomain</code>) through the <code>getGeneCollection</code> and <code>getProteinDomainCollection</code> methods respectively. The sequence information for the representative Mrna sequence of this array-probe can be obtained from <code>NucleicAcidSequence</code> using <code>getNucleicAcidSequenceCollection</code> method.</p> <p>The mRNA-transcripts used to form the array probe itself is encapsulated in the Transcript object, along with such transcripts from Affymetrix Human Exon Arrays.</p>

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
Gene	<p>Gene objects are the effective portal to most of the genomic information provided by the caBIO data services; It provides access to the following objects:</p> <ul style="list-style-type: none"> a. associated proteins (Protein) b. gene-targets (Target) c. organs for which histopathological information of tumor-samples is available that implicates the gene (OrganOntology) d. histopathologies (Histopathology) associated with c) e. curated gene-ontologies (GeneOntology) f. BioCarta pathways (Pathway) involving the gene g. aliases (GeneAlias) from HUGO and Entrez h. clones (Clone) from Unigene i. Mrna/EST-sequences (NucleicAcidSequence) j. SNPs located on the Gene (GeneRelativeLocation) as per Affymetrix array annotation data files. k. chromosomal start and stop (GenePhysicalLocation) from MapView l. Agents it interacts with, per NCI's CGDC Text Mining Project (GeneAgentAssociation) m. Tumors it is implicated in (GeneDiseaseAssociation) as per NCI's CGDC Text Mining Project. n. Corresponding Identifiers from OMIM and Entrez in DatabaseCrossReference o. UniSTS Markers corresponding to this gene in the Marker object p. Symbols from HUGO <p>The Gene object itself is populated with human and mouse data from Unigene. Other associated objects (and hence their associations to Gene) get populated from different sources.</p>
GeneAlias	<p>An alternative name for a gene; provides descriptive information about the gene (as it is known by this alias), as well as access to the primary Gene it refers to. It gets populated with data from HUGO, Entrez and CGAP</p>
GeneAgentAssociation	<p>Agents interacting with a Gene as per results of text mining scientific publications from PubMed in the CGDC Project. It provides an association to details of the corresponding publication through the Evidence object.</p> <p>It is a subclass of the GeneFunctionAssociation class.</p>
GeneCytogeneticLocation	<p>Start and stop of the Gene in terms of Cytobands. It provides details of the cytoband through the associated Cytoband object.</p> <p>It gets populated with data from Unigene.</p>

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
GeneDiseaseAssociation	DiseaseOntologies associated with a Gene as per results of text mining scientific publications from PubMed in the CGDC Project. It provides an association to details of the corresponding publication through the Evidence object. It is a subclass of the GeneFunctionAssociation class.
GeneFunctionAssociation	An abstract superclass of the GeneDiseaseAssociation and GeneAgentAssociation classes.
GenePhysicalLocation	An object representing chromosomal starts and stops of genes. It gets populated with data from NCBI's MapView application
GeneOntology	An object providing information on a Gene's position in the Gene Ontology Consortium's controlled vocabularies. It provides access to Genes corresponding to an ontological term and to the relevant parent/child ontological relationship, if any via the corresponding GeneOntologyRelationship. It gets populated with GO data from CGAP and is kept synchronized with the parent GO Consortium datasets.
GeneOntology Relationship	Specifies the parent/child relationship, if any, for a given Gene Ontology and allows the user to retrieve such associated GeneOntologys. It gets populated with data from CGAP, as with the GeneOntology object above.
GeneRegulation	A type of pathway interaction (subclass of Interaction) that regulates gene expression. It is populated with data from PID.
GeneRelativeLocation	An object providing location of SNPs on Genes such as whether the SNP is located in the 3' UTR, 5' UTR, coding region of the Gene. It gets populated with data from Affymetrix Human Mapping 250K Nsp and Sty, 50K Xba240 and Hind240 and Illumina 550K genotyping arrays and provides a link to the associated Gene.
Histopathology	Represents anatomical changes in a diseased tissue sample associated with an expression experiment; captures the relationship between an organ and disease. A user can also retrieve any associated Pathways.
Homologous Association	Represents degrees of homology (as a percentage) amongst various genes in human and mouse, along with access to further details about the Genes themselves. It is populated with human and mouse data from CGAP.
Input	A pathway participant (subclass of PhysicalParticipant) that is consumed by an interaction. It is populated with data from PID.
Interaction	A biological event which is part of a larger process (Pathway). It is populated with data from PID.

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
Library	<p>An object representing a CGAP library; provides access to information about: the tissue sample and its method of preparation, the library protocol that was used, the clones contained in the library, and the sequence information derived from the library.</p> <p>CGAP in turn amalgamates data from a variety of different Clone Libraries such as dbEST, ORESTES, Unigene, etc.</p>
Location	<p>Super class for <Feature>PhysicalLocation and <Feature>CytogenicLocation objects. Provides start and stop locations in terms of chromosomes and cytobands for</p> <ol style="list-style-type: none"> Genes : GenePhysicalLocation and GeneCytogeneticLocation SNPs : SNPPhysicalLocation and SNPCytogeneticLocation Markers : MarkerPhysicalLocation Exons : ExonPhysicalLocation Transcripts : TranscriptPhysicalLocation ArrayReporters : ArrayReporterPhysicalLocation and ArrayReporterCytogeneticLocation NucleicAcidSequence (ESTs and MRNAs) : NucleicAcidPhysicalLocation Cytobands : CytobandPhysicalLocation
Macroprocess	A type of pathway interaction (subclass of Interaction) annotated with controlled vocabulary terms from the Gene Ontology's Biological Process vocabulary. It is populated with data from PID.
Marker	<p>An object representing different kinds of Markers such as 5'/3' Primers, Microsatellites and SNPs.</p> <p>It gets populated with human and mouse primers from UniSTS, microsatellite markers from Affymetrix HuMapping Arrays and SNPs from TSC.</p> <p>The location of a SNP relative to a marker is available through MarkerRelativeLocation and is populated with data from Affymetrix HuMapping arrays.</p> <p>Chromosomal start and stop of a Marker is available through MarkerPhysicalLocation</p>
MarkerAlias	An object representing different aliases or names for a Marker. It gets populated with data from UniSTS
MarkerPhysicalLocation	Class that provides the chromosomal start and stop of the Marker. It gets populated with data from NCBI's MapView.
MarkerRelativeLocation	<p>An object providing location of SNP from Affymetrix HuMapping Arrays relative to a Microsatellite Marker from these same arrays.</p> <p>The associated probeset information is available through ArrayReporter</p>

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
Microarray	An object representing various microarrays in caBIO such as the Affymetrix Human Mapping, HG-U133, HG-U133 Plus2, HT-HG-U133, HT-HG-U133 A 2.0 and Illumina HumanHap 550K 'SNP'/genotyping arrays, Affymetrix 'Exon' Arrays and 'oligo' arrays such as the Agilent Human Genome 44k and aCGH 244k arrays.
NegativeControl	A type of pathway participant (subclass of <code>PhysicalParticipant</code>) that exercises negative control over an interaction, but is not consumed by the interaction. It is populated with data from PID.
NucleicAcidPhysical Location	Provides chromosomal start and stop of an mRNA or an EST, as provided by UCSC's Genome Sequencing Center.
NucleicAcidSequence	<p>An object representation of a gene sequence; provides access to the clones from which it was derived, the ASCII representation of the residues it contains, and any <code>ExpressionArrayReporters</code> for which it is a representative sequence.</p> <p>It is populated with data pertaining to human and mouse from Unigene.</p> <p>Its chromosomal start and stop is available through <code>NucleicAcidPhysicalLocation</code></p>
OrganOntology	<p>A representation of an organ whose name occurs in a controlled vocabulary; provides access to any associated histopathological information, and anomalies. Because it is "ontolog-able", It provides access to any associated ontologies that it shares a parent/child relationship with, through the linked <code>OrganOntologyRelationship</code>.</p> <p>Data in this table is static.</p>
OrganOntologyRelationship	<p>An object that describes parent/child relationships, if any, amongst <code>OrganOntologys</code>.</p> <p>Data in this table is static.</p>
Output	A type of pathway participant (subclass of <code>PhysicalParticipant</code>) which is produced by an interaction. It is populated with data from PID.
Participant	A participant in a pathway interaction or part of a complex. May be a <code>Condition</code> or a <code>PhysicalEntity</code> . It is populated with data from PID.

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
Pathway	<p>An object representation of a molecular/cellular pathway in human and mouse compiled by BioCarta. It provides access to further details about the associated Genes and any available histopathological information for genes involved in clinical trials through <code>getGeneCollection</code> and <code>getHistopathologyCollection</code> methods.</p> <p>It gets populated with data from CGAP and is composed of BioCarta Pathways. It provides links to genes associated with every pathway as well as histopathological information pertaining to these Genes through the associated Histopathology object.</p>
PathwayReference	A kind of pathway interaction (subclass of <code>Interaction</code>) which stands for another pathway. It is populated with data from PID.
PhysicalEntity	A reference to a molecular species (protein, complex, RNA, or small molecule) that may participate in a pathway. It is populated with data from PID.
PhysicalEntityFamily	A set of related molecular species (<code>PhysicalEntity</code> objects) constituting a family. It is populated with data from PID.
PhysicalLocation	<p>Provides chromosomal start and stop positions for SNPs, ESTs, MRNAs, Markers and Cytobands. Further details for the SNP can be obtained from the SNP object. The sequence associated with the ESTs and MRNAs can be obtained from the linked <code>NucleicAcidSequence</code> object.</p> <p>The SNP Data is obtained from NCBI's dbSNP, while the latter are obtained from UCSC's Genome Sequencing Center for every Human and Mouse Genome build.</p>
PhysicalParticipant	A condition or physical entity that participates in an interaction or is a component of a complex (<code>ComplexEntity</code>). It is populated with data from PID.
PopulationFrequency	Represents the major and minor alleles of a SNP and their respective frequencies in different populations.
PositiveControl	A pathway participant (subclass of <code>PhysicalParticipant</code>) that exercises a positive (promoting) control over an interaction but is not consumed by the interaction. It is populated with data from PID.
Protein	An object representation of a protein; provides access to the encoding Gene via its GenBank ID, the Taxon in which this instance of the protein occurs, and references to homologous proteins in other species. It is populated with data from Uniprot-Swissprot.
ProteinAlias	An alternate name for a protein. This data is obtained from Uniprot-Swissprot.

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
ProteinDomain	An object representing protein domains from Interpro. This object is populated with Interpro protein domains associated with probesets from Affymetrix HG-U133, HG-U133 Plus2, HT-HG-U133 and HT-HG-U133A 2.0 Array annotations.
ProteinEntity	A kind of physical entity (subclass of <code>PhysicalEntity</code>), consisting of a sequence of amino acids. This object serves as a reference to one or more Proteins, for use within a Pathway. It is populated with data from PID.
ProteinSequence	The sequence of a protein from Swissprot.
ProteinSubunit	A cleaved product of a protein, for use within a Pathway. It is populated with data from PID.
Protocol	An object representation of the protocol used in assembling a clone library. It is populated with data from NCI's CTEP (Cancer Therapy Evaluation Program)
ProtocolAssociation	An association class relating <code>ClinicalTrialProtocols</code> to Diseases.
RelativeLocation	An abstract super class for <code>MarkerRelativeLocation</code> and <code>GeneRelativeLocation</code> that provides the relative locations of SNPs with regards to Genes and Markers. These classes get populated with data from Affymetrix HuMapping arrays and Illumina HapMap SNP Arrays.
RNAEntity	A kind of physical entity (subclass of <code>PhysicalEntity</code>), which is an RNA macromolecule. This object serves as a reference to one or more <code>NucleicAcidSequence</code> objects, for use within a Pathway. It is populated with data from PID.
SmallMoleculeEntity	A kind of physical entity (subclass of <code>PhysicalEntity</code>), which is not a macromolecule (i.e. coded by DNA) and is not a complex. This object serves as a reference to one or more <code>Agent</code> objects, for use within a Pathway. It is populated with data from PID.
SNP	<p>An object representing a Single Nucleotide Polymorphism; It provides the two most common substitutions at that position and the offset of the SNP in the parent sequence.</p> <p>A reference to related SNPs from The SNP Consortium (TSC) is available through <code>DatabaseCrossReference</code>.</p> <p>It provides a link to the associated Population Frequencies. The location of a SNP is available through the <code>PhysicalLocation</code> and <code>Location</code> objects as well as through the <code>SNPPhysicalLocation</code> object. Location of SNPs in terms of cytogenetic starts and stops is available through the <code>SNPCytogeneticLocation</code> object.</p>
SNPArrayReporter	An object representing probeset information for genotyping arrays such as Affymetrix Human Mapping Arrays and Illumina HumanHap 550k Arrays; It provides details on the corresponding SNP through the linked SNP object.

Table 6.1 caBIO domain objects and descriptions (Continued)

caBIO Domain Object	Description
SNPCytogeneticLocation	Provides start and stops for SNPs in the SNP object in terms of cytobands. It gets populated with data from Affymetrix HuMapping array annotations.
SNPPhysicalLocation	Provides chromosomal start and stop for SNPs in the SNP object. Gets populated with data from NCBI's dbSNP
Target	A gene thought to be at the root of a disease etiology, and which is targeted for therapeutic intervention in a clinical trial. Data in this object is static. Each target links to the associated Gene object through which further details of the Target are available.
Taxon	An object representing the various names (scientific, common, abbreviated, etc.) for human and mouse genomes. It can be used to retrieve the associated Genes, Chromosomes, Pathways if any, Proteins, or Tissues. Data in this object is static.
Tissue	A group of similar cells united to perform a specific function. Data in this object is static
Transcript	An object representation of a collection of Exons, that constitutes a probe in an Affymetrix Exon Array. It can be used to retrieve the associated Exons, ExonArrayReporters and Genes. It gets populated with data from Affymetrix HuExon Arrays
TranscriptArrayReporter	An abstract superclass for ExpressionArrayReporter and ExonReporter.
TranscriptPhysicalLocation	Chromosomal Start and Stop of a Transcript obtained from Affymetrix HuExon Arrays
Vocabulary	A collection of terms associated with a given Anomaly. Data in this object is static

Table 6.1 caBIO domain objects and descriptions (Continued)

Provenance Objects	Description
Provenance	An object modeling the provenance of various entities such as SNP, Gene, NucleicAcidSequence and Protein data in caBIO. It exposes information such as the source of the data (dbSNP, Unigene, LocusLink, Uniprot, etc) along with the class name modeling that object in caBIO, evidence code, etc.
InternetSource	An object modeling details of the primary data source providers for SNP, Gene, etc in caBIO
PublicationSource	An object modeling details of publications associated with caBIO data such as the author, date of publication, etc.

Table 6.2 Provenance objects and descriptions

Provenance Objects	Description
ResearchInstitutionSource	An object modeling details of the research institution such as name, address, contact information, etc from where caBIO data is obtained.
Source	An object modeling the primary data source for objects like SNP, Gene, Protein, etc in caBIO
SourceReference	An object providing web-enabled access to further details regarding Genes, Proteins, SNPs, and other caBIO objects from their primary data sources.
URLSourceReference	An object providing web-enabled access to obtain further details regarding Genes, Proteins, SNPs, and other caBIO objects from their primary data sources.

Table 6.2 Provenance objects and descriptions (Continued)

Data Sources in the caBIO Database

This section describes the internal and external data sources for caBIO and how the information these sources provide can be accessed via caBIO objects.

The caBIO application programming interfaces were developed primarily in response to the need for programmatic access to the information at several NCI web sites, including:

- Cancer Genome Anatomy Project (CGAP)
- Cancer Therapy Evaluation Program (CTEP)
- CGAP Genetic Annotation Initiative (GAI)
- Cancer Molecular Analysis Project (CMAP)
- Affymetrix, Agilent and Illumina Microarrays
- Genome Sequencing Center at University of California, Santa Cruz (UCSC)
- SNP Annotations from NCBI
- Protein Data from Uniprot-Swissprot
- SNP Annotations from SNP Consortium
- NCI's Cancer Gene Index Data Project
- NCBI's MapViewer data
- UniSTS
- HUGO
- and others

While this information is, in theory, available from multiple public sites, the number of links to traverse and the extent of collation that would have to be performed is daunting. The CGAP, CMAP, and GAI web sites have distilled this information from both internal

and public databases, and the caBIO data warehouses have optimized it for access with respect to the types of queries defined in the APIs.

While the caBIO data are extracted from many sources that include information from a wide variety of species, we emphasize that only genomic data pertaining to human and mouse are available from caBIO.

caBIO provides access to curated data from both internal (NCI) and external sites.

Table 6.3 contains data source internal to NCI and Table 6.4 contains data sources from sites outside of NCI.

NCI Data Source	Description
CGAP	<p>CGAP (Cancer Genome Anatomy Project) provides a collection of gene expression profiles of normal, pre-cancer, and cancer cells taken from various tissues. The CGAP interface allows users to browse these profiles by various search criteria, including histology type, tissue type, library protocol, and sample preparation methods. The goal at NCI is to exploit such expression profile information for the advancement of improved detection, diagnosis, and treatment for the cancer patient. Researchers have access to all CGAP data and biological resources for human and mouse, including ESTs, gene expression patterns, SNPs, cluster assemblies, and cytogenetic information.</p> <p>The CGAP web site provides a powerful set of interactive data-mining tools to explore these data, and the caBIO project was initially conceived as a programmatic interface to these tools and data. Accordingly, most of the data that are available from CGAP can also be accessed through the caBIO objects. Exceptions are those data sets having proprietary restrictions, such as the Mittleman Chromosome Aberration database.</p> <p>CGAP also provides access to lists of sequence-verified human and mouse cDNA IMAGE clones supplied by Invitrogen.</p>
CMAP	<p>CMAP (Cancer Molecular Analysis Project) is powered by caCORE. The goal of CMAP is to enable researchers to identify and evaluate molecular targets in cancer.</p> <p>The CMAP Profile Query tool finds genes with the highest or lowest expression levels (using SAGE and microarray data) for a given tissue and histology. Selecting a gene from the resulting table then leads to a Gene Info page. This page provides information about cytogenetic location, chromosome aberrations, protein similarities, curated and computed orthologs, and sequence-verified as well as full-length MGC clones, along with links to various other databases.</p>

Table 6.3 NCI-related data sources in the caBIO database

NCI Data Source	Description
CTEP	<p>CTEP (Cancer Therapy Evaluation Program) funds an extensive national program of basic and clinical research to evaluate new anti-cancer agents, with a particular emphasis on translational research to elucidate molecular targets and drug mechanisms. In response to this emergent need for translational research, there has been a groundswell of translational support tools defining controlled vocabularies and registered terminologies to enhance electronic data exchange in areas that have heretofore been relatively non-computational. The caCORE trials data are updated with new CTEP data on a quarterly basis, and many of the objects are designed to support translational research.</p> <p>For example, a caCORE Target object represents a molecule of special diagnostic or therapeutic interest for cancer research, and an Anomaly object is an observed deviation in the structure or expression of a Target. An Agent is a drug or other intervention that is effective in the presence of one or more specific Targets. The ClinicalTrialProtocol object organizes administrative information pertaining to that protocol.</p> <p>Data from CTEP are used to populate the Protocols, ProtocolAgents and ProtocolDiseases objects.</p>
GAI	<p>GAI (CGAP Genetic Annotation Initiative) is an NCI research program to explore and apply technology for identification and characterization of genetic variation in genes important in cancer. The GAI utilizes data-mining to identify "candidate" variation sites from publicly available DNA sequences, as well as laboratory methods to search for variations in cancer-related genes. All GAI candidate, validated, and confirmed genetic variants are available directly from the GAI web site, and all validated SNPs have been submitted to the NCBI dbSNP database as well that are in turn available through the SNP .</p>
HomoloGene	<p>HomoloGene is an NCBI resource for curated and calculated gene homologs. The caBIO data sources capture only the calculated homologs stored by HomoloGene. These calculated homologs are the result of nucleotide sequence comparisons performed between each pair of organisms represented in UniGene clusters.</p> <p>caBIO provides this information via HomologousAssociation and updates this data on a monthly basis.</p>
PID	<p>The Pathway Interaction Database is a highly-structured, curated collection of information about known biomolecular interactions and key cellular processes assembled into signaling pathways. It is a collaborative project between the US National Cancer Institute (NCI) and Nature Publishing Group (NPG), and is an open access online resource.</p>

Table 6.3 NCI-related data sources in the caBIO database (Continued)

NCI Data Source	Description
LocusLink/Entrez Gene	<p>LocusLink contains curated sequence and descriptive information associated with a gene. Each entry includes information about the gene's nomenclature, aliases, sequence accession numbers, phenotypes, UniGene cluster IDs, OMIM IDs, gene homologies, associated diseases, map locations, and a list of related terms in the Gene Ontology Consortium's ontology. Sequence accessions include a subset of GenBank accessions for a locus, as well as the NCBI Reference Sequence.</p> <p>The LocusLink Identifier corresponding to a caBIO Gene is available in DatabaseCrossReference.</p>
MapView	<p>Entrez Genomes presents a unified graphical view of maps (genetic and physical) and sequence data for a selected organism.</p> <p>The Entrez Map Viewer, is a software component of Entrez Genomes which provides an organism's complete genome, integrated maps (when available) for each chromosome, and sequence data for a region of interest.</p> <p>Data from MapView is used to populate GenePhysicalLocation and MarkerPhysicalLocation objects.</p>
dbSNP	<p>In collaboration with the National Human Genome Research Institute, the NCBI has established the dbSNP database to serve as a central repository for both single base nucleotide substitutions and short deletion and insertion polymorphisms. Once discovered, these polymorphisms could be used by additional laboratories, using the sequence information around the polymorphism and the specific experimental conditions. (Note that dbSNP takes the looser 'variation' definition for SNPs, so there is no requirement or assumption about minimum allele frequency.) The data from dbSNP is updated approximately every 3-4 months.</p> <p>Relevant information is available through SNP, Provenance, Source, URLSourceReference, SourceReference, PhysicalLocation, Location, SNPPPhysicalLocation, SNPCytogeneticLocation, GeneRelativeLocation and MarkerRelativeLocation objects.</p>
UniSTS	<p>Database at the National Center for Biotechnology Information providing "a unified, non-redundant view of sequence tagged sites (STS's). UniSTS integrates marker and mapping data from a variety of public resources." Used in eGenome as a source of STS's, and to collect and manage element names. Data from UniSTS is used to populate Marker and MarkerAlias objects.</p>

Table 6.3 NCI-related data sources in the caBIO database (Continued)

NCI Data Source	Description
UniGene	<p>UniGene provides a nonredundant partitioning of the genetic sequences contained in GenBank into gene clusters. Each cluster has a unique UniGene ID and a list of the mRNA and EST sequences that are included in that cluster. Related information stored with the cluster includes tissue types in which the gene has been expressed, mapping information, and the associated LocusLink, OMIM, and HomoloGene IDs, thus providing access to related information in those NCBI databases as well through DatabaseCrossReference.</p> <p>Because the information in UniGene is centered around genes, access to Unigene is provided via the caBIO Gene objects. Specifically, the method <code>getClusterId</code> associated with a Gene object can be used to fetch the gene's UniGene ID. The corresponding IDs to cross-reference these genes into the NCBI OMIM and LocusLink databanks, Enzyme Commission, Ensembl and RefSeq databases can be obtained from DatabaseCrossReference using the <code>getDatabaseCrossReferenceCollection</code> method. While there is no explicit caBIO object corresponding to a Unigene cluster, all of the information associated with the cluster is available directly via the caBIO Gene object's methods.</p> <p>The sequences are available via <code>NucleicAcidSequence</code> and the associated aliases are available in the <code>GeneAlias</code>. Corresponding clone library and its location information is exposed via <code>Clone</code> and <code>CloneRelativeLocation</code>.</p> <p>Appropriate 'provenance' information is available in <code>Provenance</code>, <code>Source</code>, <code>URLSourceReference</code> and <code>SourceReference</code> objects.</p> <p>Markers associated with a Gene are available through the Marker Object whereas its chromosomal and cytogenetic start-stops are available through <code>GenePhysicalLocation</code> and <code>GeneCytogeneticLocation</code> objects.</p> <p>Associated Histopathological and Pathway information are available through <code>Histopathology</code> and <code>Pathway</code> respectively.</p> <p>Unigene data is updated on a monthly basis.</p>
Cancer Gene Data Curation Project	<p>The Cancer Gene Data Curation Project is an attempt to create a database of associations between genes and diseases and genes and drug compounds derived from the biomedical literature. The project involves a mixture of automatic text mining, semi-automatic verification, and manual validation/scoring of results.</p> <p>Data from this project is exposed through caBIO's <code>Evidence</code>, <code>EvidenceCode</code>, <code>GeneDiseaseAssociation</code> and <code>GeneAgentAssociation</code> objects</p>

Table 6.3 NCI-related data sources in the caBIO database (Continued)

External data sources in the caBIO database:

External Data source	Description
Affymetrix	<p>Affymetrix provides the majority of Microarray data for caBIO.</p> <p>The data provides information on allele frequencies of the SNP in different populations, and is represented by the PopulationFrequency object.</p> <p>The probeset information is available through ArrayReporter, ExpressionArrayReporter and SNPArrayReporter objects, amongst others.</p> <p>Details regarding the arrays exposed in caBIO are available through Microarray</p> <p>The GeneRelativeLocation object provides the location (intron, upstream, downstream etc.) of a SNP referenced in SNP Array datasets with respect to its associated genes. The validation status for a SNP comes from NCBI. The SNP Consortium (TSC) Ltd. a non-profit foundation provides the TSC id's for SNPs in DatabaseCrossReference.</p>
Agilent	<p>Data from Agilent's Whole Human Genome 44K Arrays and aCGH 244K Arrays are exposed through ArrayReporter, MicroArray and ExpressionArrayReporter. Associated Interpro Protein Domains and Unigene Genes are available through ProteinDomain and Gene objects.</p>
BioCarta	<p>BioCarta and its Proteomic Pathway Project (P3) provide detailed graphical renderings of pathway information concerning adhesion, apoptosis, cell activation, cell signaling, cell cycle regulation, cytokines/chemokines, developmental biology, hematopoiesis, immunology, metabolism, and neuroscience. NCI's CMAP web site captures pathway information from BioCarta, and transforms the downloaded image data into Scalable Vector Graphics (SVG) representations that support interactive manipulation of the online images. The CMAP web site displays BioCarta pathways selected by the user and provides options for highlighting anomalies, which include under- or over expressed genes as well as mutations.</p> <p>The pathway information is available via the Pathway object in caBIO.</p> <p>caBIO provides a class for manipulating SVG diagrams, which is described in SVG Manipulator on page 63.</p>
Illumina	<p>One of the SNP-Arrays in caBIO. Data from Illumina is used to populate MicroArray and ExpressionArrayReporter, amongst others.</p>
UniProt PIR	<p>Universal Protein Resource (UniProt) is a complete annotated protein sequence database and is a central repository of protein sequence and function created by joining the information contained in Swiss-Prot, TrEMBL, and PIR. The UniProt Knowledge base provides access to extensive curated protein information, including the amino acid sequence, protein name or description, taxonomic data and protein aliases.</p> <p>caBIO exposes information from the Swissprot databanks through ProteinSequence and ProteinAlias objects.</p> <p>Mappings to RefSeq Ids are available through DatabaseCrossReference.</p> <p>Provenance-related information is available through Provenance, Source, SourceReference and URLSourceReference objects.</p> <p>Protein Domain information from Interpro is exposed through ProteinDomain object</p>

Table 6.4 External data sources in the caBIO database

External Data source	Description
Gene Ontology Consortium	<p>The Gene Ontology Consortium provides a controlled vocabulary for the description of molecular functions, biological processes, and cellular components of gene products. The terms provided by the consortium define the recognized attributes of gene products and facilitate uniform queries across collaborating databases.</p> <p>In general, each gene is associated with one or more biological processes, and each of these processes may in turn be associated with many genes. In addition, the GO ontologies define many parent/child relationships among terms. For example, a branch of the ontology tree under biological process contains the term "cell cycle control", which in turn bifurcates into the "child" terms cell cycle arrest, cell cycle checkpoint, control of mitosis, etc.</p> <p>caBIO does not extract ontology terms directly from the Gene Ontology Consortium but, instead, extracts those terms stored with the LocusLink entry for that gene.</p> <p>This information is available via GoClosure, GoGenes, GeneOntology and GeneOntologyRelationship objects.</p>
SNP Consortium	<p>The SNP Consortium Ltd. is a non-profit foundation organized for providing public genomic data. Its mission is to develop up to 300,000 SNPs distributed evenly throughout the human genome and to make the information related to these SNPs available to the public without intellectual property restrictions.</p> <p>The TSC Ids corresponding to a SNP from dbSNP are available through DatabaseCrossReference.</p>
UCSC Genome	<p>UCSC (University of California, Santa Cruz Distributed Annotation System) provides the data for the Chromosomal start and end positions of mRNA, EST and Cytoband sequences.</p> <p>This data is used to populate the PhysicalLocation, Location, Cytoband, CytogeneticLocation and Location objects with the locations of ESTs, MRNAs and Cytobands respectively.</p> <p>The ESTs and MRNAs in PhysicalLocation and Location are linked with the corresponding sequence identifiers from NucleicAcidSequence</p>

Table 6.4 External data sources in the caBIO database (Continued)

caGrid Identifiers

The functionality provided by caGrid's Identifier Services Framework and the integration in caBIO is to have "identifiers" for each caBIO domain object. The identifier is essentially a forever globally unique name for the caBIO data-object such that it can be unambiguously used to refer to the data from different application contexts.

This identifier, obtained from the caGrid Identifier Service Framework, is essentially a string and a forever globally unique name for the caBIO domain object. Furthermore, the identifier can be (globally) resolved to the associated caBIO domain object.

In order to abstract the identifier's object properties, the data service implementations and the resolution mechanisms, the identifier's value must be treated as a "meaningless" opaque string by the consumer applications. Any leaking of design and implementation choices for the identifier framework in the applications is undesirable

from an architecture point of view as it makes the implementations brittle and susceptible to future changes. Of course resolution information will have to be embedded in identifier name, but this should only be meaningful for resolution service related components that are layered below the application.

Classes to be indexed for FreeStyleLM Search

The FreestyleLM (Freestyle Lexical Mine) search component, provides and interfaces and APIs for conducting horizontal searches across caBIO domain objects and is described in greater detail in Chapter 5.

FreeStyleLM is build upon Hibernate and Apache Lucene full text search engines, providing full text search (Google™-like search) capabilities to the caBIO API. The classes and fields that are indexed are annotated in the caBIO domain model. Usually, these attributes are non-numeric fields with descriptions or identifiers. Thus, searching for “brca1” will return any object that makes mention of the “brca1” gene in its attributes, such as Genes, Proteins, Pathways, and so on.

The following classes are indexed as part of FreeStyleLM Indexing application:

Class	Attribute
Agent	Id Comment EVSId Name NSCNumber Source
Anomaly	Id Description
Chromosome	Id Number
ClinicalTrialProtocol	Id CurrentStatus DocumentNumber LeadOrganizationId LeadOrganizationName NIHAdminCode ParticipationType PDQIdentifier Phase Title
Condition	Name
Cytoband	Id Name

Table 6.5 Classes indexed in FreeStyleLM

Class	Attribute
DiseaseOntology	Id Name EVSId
EntityAccession	Accession Database
EntityName	Name
EvidenceCode	Id EvidenceCode
ExpressionArrayReporter	Id SequenceSource SequenceType TargetDescription
Gene	Id ClusterId FullName Symbol
GeneAlias	Id Name Type
GeneOntology	Id Name
GeneFunctionAssociation	Id Role
Histopathology	Id AgeOfOnset Comments GrossDescription MicroscopicDescription
Library	Id CloneProducer CloneVector CloneVectorType Description Keyword LabHost Name UniGenId
MacroProcess	Name

Table 6.5 Classes indexed in FreeStyleLM (Continued)

Class	Attribute
Microarray	Id Manufacturer Name Type
OrganOntology	Id Name
PhysicalProperties	ActivityState Location PostTranslationalMod
ProteinAlias	Id Name
ProtocolAssociation	Id CTEPName DiseaseCategory DiseaseSubCategory
Protein	Id Keywords Name PrimaryAccession SecondaryAccession UniProtCode
Protocol	Id Description Name Type
ProteinDomain	Id AccessionNumber Description Source
Pathway	Id Description DisplayValue Name

Table 6.5 Classes indexed in FreeStyleLM (Continued)

Class	Attribute
Tissue	Id CellLine CellType Description DevelopmentalStage Histology Name Organ Supplier Type
Taxon	Id Abbreviation CommonName EthnicityStrain ScientificName
Target	Id Name Type
Vocabulary	Id CoreTerm GeneralTerm

Table 6.5 Classes indexed in FreeStyleLM (Continued)



APPENDIX

A

REFERENCES

Technical Manuals/Articles

1. National Cancer Institute. caCORE SDK 4.0 Developer's Guide <http://ncicb.nci.nih.gov/NCICB/infrastructure/cacoresdk#Documentation>
2. Java Bean Specification: <http://java.sun.com/products/javabeans/docs/spec.html>
3. Foundations of Object-Relational Mapping: <http://www.chimu.com/publications/objectRelational/>
4. Object-Relational Mapping articles and products: <http://www.service-architecture.com/object-relational-mapping/>
5. Hibernate Reference Documentation: <http://www.hibernate.org/5.html>
6. Basic O/R Mapping: http://www.hibernate.org/hib_docs/v3/reference/en/html/mapping.html
7. Java Programming: <http://java.sun.com/learning/new2java/index.html>
8. Jalopy User Manual: <http://jalopy.sourceforge.net/existing/manual.html>
9. Javadoc tool: <http://java.sun.com/j2se/javadoc/>
10. JUnit: <http://junit.sourceforge.net/>
11. Extensible Markup Language: <http://www.w3.org/TR/REC-xml/>
12. XML Metadata Interchange: <http://www.omg.org/technology/documents/formal/xmi.htm>
13. Ehcache: <http://ehcache.sourceforge.net/documentation/>

Scientific Publications

1. Ansher SS and Scharf R (2001). The Cancer Therapy Evaluation Program (CTEP) at the National Cancer Institute: industry collaborations in new agent development. *Ann N Y Acad Sci* 949:333-40.
2. Boon K, Osorio EC, Greenhut SF, Schaefer CF, Shoemaker J, Polyak K, Morin PJ, Buetow KH, Strausberg RL, De Souza SJ, and Riggins GJ (2002). An anatomy of normal and malignant gene expression. *Proc Natl Acad Sci U S A* 2002 Jul 15.
3. Buetow KH, Klausner RD, Fine H, Kaplan R, Singer DS, and Strausberg RL (2002). Cancer Molecular Analysis Project: Weaving a rich cancer research tapestry. *Cancer Cell* 1(4):315-8.
4. Boguski & Schuler (1995). ESTablishing a human transcript map. *Nature Genetics* 10: 369-71.
5. Clifford R, Edmonson M, Hu Y, Nguyen C, Scherpbier T, and Buetow KH (2000). Expression-based genetic/physical maps of single-nucleotide polymorphisms identified by the Cancer Genome Anatomy Project. *Genome Res* 10(8):1259-65.
6. Covitz P.A., Hartel F., Schaefer C., De Coronado S., Sahni H., Gustafson S., Buetow K. H. (2003). caCORE: A common infrastructure for cancer informatics. *Bioinformatics*. 19: 2404-2412.
7. Dowell RD, Jokerst RM, Day A, Eddy SR, Stein L. The Distributed Annotation System. *BMC Bioinformatics* 2(1):7.
8. The Gene Ontology Consortium. (2000). Gene ontology: tool for the unification of biology. *Nature Genetics* 25:25-9.
9. The Gene Ontology Consortium. (2001). Creating the gene ontology resource: design and implementation. *Genome Res* 11:1425-33.
10. Hartel FW and de Coronado S (2002). Information standards within NCI. In: *Cancer Informatics: Essential Technologies for Clinical Trials*. Silva JS, Ball MJ, Chute CG, Douglas JV, Langlotz C, Niland J and Scherlis W, eds. Springer-Verlag.
11. Pruitt KD, Katz KS, Sicotte H, and Maglott DR (2000). Introducing RefSeq and LocusLink: curated human genome resources at the NCBI. *Trends Genet* 16(1):44-7.
12. Pruitt KD, and Maglott DR (2001). RefSeq and LocusLink: NCBI gene-centered resources. *Nucleic Acids Res* 29(1):137-40.
13. Schuler et al. (1996). A gene map of the human genome. *Science* 274: 540-6.
14. Schuler (1997). Pieces of the puzzle: expressed sequence tags and the catalog of human genes. *J Mol Med* 75(10):694-8.
15. Strausberg RL (1999). The Cancer Genome Anatomy Project: building a new information and technology platform for cancer research. In: *Molecular*

- Pathology of Early Cancer (Srivastava S, Henson DE, Gazdar A, eds. IOS Press, 365-70.
16. Strausberg RL (2001). The Cancer Genome Anatomy Project: new resources for reading the molecular signatures of cancer. *J Pathol* 195:31-40.
 17. Zhang, Schwartz, Wagner, and Miller (2000). A Greedy algorithm for aligning DNA sequences. *J Comp Biol* 7(1-2):203-14.
 18. Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K (2001). dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* 1;29(1):308-11.
 19. Karolchik D, Baertsch R, Diekhans M, Furey TS, Hinrichs A, Lu YT, Roskin KM, Schwartz M, Sugnet CW, Thomas DJ, Weber RJ, Haussler D and Kent WJ(2003) The UCSC Genome Browser Database. *Nucl. Acids Res* 31(1), 51-54
 20. Lennon G, Auffray C, Polymeropoulos M, Soares MB(1996) The I.M.A.G.E. Consortium: an integrated molecular analysis of genomes and their expression. *Genomics.* 33(1):151-2
 21. Wheeler DL, Barrett T, Benson DA, Bryant SH, Canese K, Chetvernin V, Church DM, DiCuccio M, Edgar R, Federhen S, Geer LY, Kapustin Y, Khovayko O, Landsman D, Lipman DJ, Madden TL, Maglott DR, Ostell J, Miller V, Pruitt KD, Schuler GD, Sequeira E, Sherry ST, Sirotkin K, Souvorov A, Starchenko G, Tatusov RL, Tatusova TA, Wagner L, Yaschenko E.(2007) Epub(2006) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.* 35(Database issue):D5-12
 22. Pontius JU, Wagner L, Schuler GD (2003) UniGene: a unified view of the transcriptome In: *The NCBI Handbook*. Bethesda (MD): National Center for Biotechnology Information.
 23. Schuler GD (1997) Pieces of the puzzle: expressed sequence tags and the catalog of human genes. *J Mol Med* 75:694-698
 24. Schuler GD, et al (1996) A gene map of the human genome. *Science* 274:540-546.
 25. Boguski MS, Schuler GD (1995) ESTablishing a human transcript map. *Nature Genetics* 10: 369-371.
 26. Thorisson GA, Stein LD (2003) The SNP Consortium website: past, present and future. *Nucleic Acids Res.* 31(1):124-7.

caBIG[®] Material

1. aBIG[®]: <http://cabig.nci.nih.gov/>
1. caBIG[®] Compatibility Guidelines: http://cabig.nci.nih.gov/guidelines_documentation

caCORE Material

1. caCORE: http://ncicb.nci.nih.gov/NCICB/infrastructure/cacore_overview
2. caBIO: http://ncicb.nci.nih.gov/NCICB/infrastructure/cacore_overview/caBIO
3. caDSR: http://ncicb.nci.nih.gov/NCICB/infrastructure/cacore_overview/cadsr
4. EVS: http://ncicb.nci.nih.gov/NCICB/infrastructure/cacore_overview/vocabulary
5. CSM: http://ncicb.nci.nih.gov/NCICB/infrastructure/cacore_overview/csm

Modeling Concepts

1. Enterprise Architect Online Manual: <http://www.sparxsystems.com.au/EAUserGuide/index.html>
2. OMG Model Driven Architecture (MDA) Guide Version 1.0.1: <http://www.omg.org/docs/omg/03-06-01.pdf>
3. Object Management Group: <http://www.omg.org/>

Applications Currently Using caBIO

1. BIO Browser: <http://www.jonnywray.com/java/index.html>
2. caPathway: <http://cgap.nci.nih.gov/Pathways>

INDEX

A

- Affymetrix 86
- Agilent 86
- annotations, source of 10
- Apache Axis 15, 42
- Application Service layer 14
- ApplicationService package 18
- architecture, caBIO system 10, 13
- Array Annotation API
 - description 58
 - use examples 61, 62

B

- BioCarta 86

C

- caBIG Query Language 28
- caBIO
 - API description 67
 - architecture overview 10
 - data sources 81
 - description 6
 - domain packages 15
 - framework packages 17
 - grid services API, access to caBIO
 - data 50
 - GUIs 19
 - home page 19
 - model described 67
 - package 16, 17
 - queries with QBE interface 21
 - service interface paradigm 24
 - system architecture description 13
 - utilities 63
- caBIO grid service API See also grid services API
- caBIO Portlet
 - description 21
 - search 21
 - simple search 21

- templated search 23

- caCORE
 - CLM 7
 - components 6
 - description 5
- caDSR, description 6
- caGrid, identifier service 87
- Cancer Gene Data Curation Project 85
- CGAP 82
- client package 18
- CLM, description 7
- CMAP 82
- Common package 16, 17
- CSM, description 6
- CTEP 83

D

- data access package 18
- Data Source Delegation layer 14
- data sources, caBIO 81
- dbSNP 84
- domain packages, caBIO 15

E

- Eclipse 43
- EVS, description 6
- examples
 - grid services API, use examples
 - Java API, simple example 27
 - Java API, use examples 31, 32, 34, 35, 36
 - web services API, use examples 45, 46
 - XML-HTTP API, use examples 49
 - XML utility 41

F

- framework packages 17, 18

FreeStyle Lexical Mine See FreeStyle LM

FreeStyleLM

- classes for search 88
- description 53
- java client 54
- search 20
- use examples 54, 55
- web interface 55

G

GAI 83

Gene Ontology Consortium 87

Genome range queries

- description 56
- Java API 56
- use examples 57

GetXML servlet

- controlling items returned 50
- paging results 50

Grid ID Resolver

- description 56
- Java API 56

grid services API

- access to caBIO data 50
- limitations 52
- service locations 51
- use examples
- working with result sets 52

GUI

- caBIO 19
- caBIO Portlet 21

H

Hibernate 13, 28, 35

Hibernate Query Language 14

HomoloGene 83

I

Illumina 86

integrated development environment 43

J

Java API

- configuration 25
- description 25
- installation 25
- search types 28
- service methods 28
- simple example 27
- use examples 31

Java Bean 15

Java Server Page 15

Java Servlet 15

JBoss 15

L

LocusLink 84

M

Model driven architecture 5

N

Non-ORM layer 15

N-tier architecture 5

O

Object Relational Mapping 14

P

pathway data, source of 10

PIR 86

Proteomic Pathway Project 86

Provenance package 16, 18

purpose of guide 1

Q

QBE interface search 21

R

references

- applications using caBIO 96
- caBIG material 95
- caCORE material 96
- modeling concepts 96
- scientific publications 94
- technical manuals 93

REST 47

S

Scalar Vector GraphicSee SVG 10
search

- caBIO Portlet 21
- caBIO Portlet, simple 21
- caBIO Portlet, templated 23
- FreeStyleLM 20

search package 18

service interface paradigm 24

SNP Consortium 87

SNP data, source of 10

SOAP 13, 15, 42

SVG 86
 manipulation utility, description 63
 manipulation utility, pathway
 diagrams 63
SVG Diagram Manipulation
 use examples 64
System architecture
 client technologies 15
system package 18

T

Tomcat 15
TrEMBL 86

U

UCSC 87
Unigene 85
UniProt 86
UniSTS 84
Utility Methods
 SVG Manipulation Utility 63
utility methods, XML Utility 41

W

war 15
Web service package 18
Web Services API
 endpoint URL 42
 operations 42
 WSDL file 42
web services API
 caBIOSService endpoint
 operations 43
 configuration 42
 description 42
 Java client 42
 limitations 47
 use examples 45

X

XLinks, retrieving related results
 using 50
XML-HTTP API
 description 47
 limitations 50
 results sets 49
 service location 47
 syntax 47
 use examples 49
XML utility 41
XSLT 15

