

CAINTEGRATOR2 v.1.0

User's Guide



Center for Biomedical Informatics
and Information Technology

CREDITS AND RESOURCES

calIntegrator2 Development and Management Teams			
Development	Quality Assurance	Documentation	Project and Product Management
JP Marple ²	Quy Phung ⁴	JP Marple ²	Shine Jacob ⁴
Will Fitzhugh ²		Jill Hadfield ¹	Anand Basu ¹
Eric Tavela ²			Juli Klemm ¹
TJ Andrews ⁵			
Ngoc Nguyen ⁶			
Matt Reh fuss ²			
Huaitian Liu ⁴			
Yuri Kotliarov ¹			
Karen Ketchum ⁷			
Systems and Application Support		Training	
Cuong Nguyen ²			
Deanna Siemaszko ³			
¹ NCI Center for Biomedical Informatics and Information Technology (CBIIT)		² 5AM Solutions	³ Terrapin Systems
⁴ Enterprise Solutions And Consulting (ESAC)	⁵ ScenPro	⁶ Claris LLC	⁷ Science Application International Corporation (SAIC)

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

TABLE OF CONTENTS

Credits and Resources	i
Using the caIntegrator2 v.1.0 User's Guide	1
Introduction to the caIntegrator2 User's Guide	1
Organization of this Guide	1
User's Guide Text Conventions	2
Chapter 1	
Getting Started with caIntegrator2	5
About caIntegrator2	5
Registering as a New caIntegrator2 User	6
Logging In	8
Using the caIntegrator2 Workspace	8
caIntegrator2 Functions	9
Using Online Help	10
Logging Out	10
Application Support	11
Chapter 2	
Creating a New Study	13
Creating a Study – Overview	13
Configuring and Deploying a Study	14
Creating/Editing a Study	15
Adding Clinical Data	16
Adding/Editing Genomic Data	24
Adding Imaging Data	29
Deploying the Study	31
Managing a Study	31
Managing Platforms	32

Chapter 3

Searching a caIntegrator2 Study35

Search Overview	35
Searching a Study	36
Results Type Tab	41
Sorting Tab	42
Managing Queries	43
Saving a Query	43
Editing a Query	44
Exporting Query Results	44

Chapter 4

Viewing Query Results45

Query Results Overview	45
Browsing Query Results	46
Clinical and Imaging Data	46
Genomic Data	46
Expanding Imaging Data Results	50
Relationship of Patient to Study to Series to Images	54
Exporting Data	55

Chapter 5

Analyzing Studies57

Data Analysis Overview	57
Creating Kaplan-Meier Plots	58
K-M Plot for Annotations	58
K-M Plot for Gene Expression	60
K-M Plot for Queries	63
Creating Gene Expression Plots	65
Gene Expression Value Plot for Annotation	66
Gene Expression Value Plot for Genomic Queries	69
Gene Expression Value Plot for Clinical Queries	71
Understanding a Gene Expression Plot	75
Analyzing Data with GenePattern	78
GenePattern Modules	80
Comparative Marker Selection (CMS) Analysis	81
Principal Component Analysis (PCA)	83
GISTIC-Supported Analysis	86

Chapter 6

Administering User Accounts89

Administering caIntegrator2 User Accounts Using UPT	89
Steps for Creating User Access to caIntegrator2	90
Creating a New caIntegrator2 User	90
Creating a New User Group	92
Creating a New Protection Group	93
Assigning a User Group to a Protection Group	94
Adding a User to a User Group	97
Changing a User Password	99

Appendix A

Data Import Configurations101

Subject Clinical Data Configuration	101
Delimited-Text Annotation Import	101
Annotation Field Configuration	102
Sample Data Configuration	102
Genomic Data Configuration	103
Imaging Data Configuration	103

Index105

USING THE CAINTEGRATOR2 v.1.0 USER'S GUIDE

This chapter introduces you to the *calIntegrator2 v.1.0 User's Guide* and suggests ways you can maximize its use.

Topics in this chapter include:

- [Introduction to the calIntegrator2 User's Guide](#) on this page
- [Organization of this Guide](#) on this page
- [User's Guide Text Conventions](#) on page 2

Introduction to the calIntegrator2 User's Guide

The *calIntegrator2 v.1.0 User's Guide* is the companion documentation to the calIntegrator2 software application. The user's guide includes information and instructions for the end user about using calIntegrator2.

Organization of this Guide

The *calIntegrator2 v.1.0 User's Guide* contains the following chapters and appendices:

Using the calIntegrator2 User's Guide — This chapter introduces you to the *calIntegrator2 v.1.0 User's Guide* and suggests ways you can maximize its use.

Chapter 1 Getting Started in calIntegrator2 — This chapter introduces general calIntegrator2 procedures and how to obtain help to use calIntegrator2. .

Chapter 2 Creating a Study — This chapter describes the processes for creating and managing studies in calIntegrator2.

Chapter 3 Searching a calIntegrator2 Study — This chapter describes the processes for searching studies within calIntegrator2 using the search and browse tools.

Chapter 4 Viewing Search Results — This chapter describes search results that calIntegrator2 returns after queries.

Chapter 5 Analyzing Studies — This chapter describes how to use calIntegrator2 tools to analyze data in clinical or genomic studies that have been deployed in calIntegrator2.

Chapter 6 Administering User Accounts —This chapter describes the process for creating and managing user accounts in caIntegrator2.

Appendix A Exporting Data — This appendix describes how MAGE-TAB documents are parsed, validated and imported into caIntegrator2. It also provides examples of the types of MAGE-TAB documents that are expected by caIntegrator2.

Index—This section of the guide provides a complete index.

User's Guide Text Conventions

Table 2.1 illustrates how text conventions are represented in this guide. The various typefaces differentiate between regular text and menu commands, keyboard keys, toolbar buttons, dialog box options and text that you type.


Convention	Description	Example
Bold & Capitalized Command Capitalized command > Capitalized command	Indicates a Menu command Indicates Sequential Menu commands	Admin > Refresh
TEXT IN SMALL CAPS	Keyboard key that you press	Press ENTER
TEXT IN SMALL CAPS + TEXT IN SMALL CAPS	Keyboard keys that you press simultaneously	Press SHIFT + CTRL and then release both.
Monospace type	Used for filenames, directory names, commands, file listings, and anything that would appear in a Java program, such as methods, variables, and classes.	URL_definition ::= url_string
Icon	A toolbar button that you click	Click the Paste button () to paste the copied text.
Boldface type	Options that you select in dialog boxes or drop-down menus. Buttons or icons that you click.	In the Open dialog box, select the file and click the Open button.
<i>Italics</i>	Used to reference other documents, sections, figures, and tables.	<i>caCORE Software Development Kit 1.0 Programmer's Guide</i>
<i>Italic boldface monospace type</i>	Text that you type	In the New Subset text box, enter <i>Proprietary Proteins.</i>
Note:	Highlights a concept of particular interest	Note: This concept is used throughout the installation manual.

Table 2.1 caIntegrator2 User's Guide Text Conventions

Convention	Description	Example
Warning!	Highlights information of which you should be particularly aware.	Warning! Deleting an object will permanently delete it from the database.
{ }	Curly brackets are used for replaceable items.	Replace {root directory} with its proper value, such as c:\cabio

Table 2.1 caIntegrator2 User's Guide Text Conventions (Continued)

CHAPTER

1

GETTING STARTED WITH CAINTEGRATOR2

This chapter introduces general calIntegrator2 procedures and how to obtain help to use calIntegrator2.

Topics in this chapter include:

- [About calIntegrator2](#) on this page
- [Registering as a New calIntegrator2 User](#) on page 6
- [Logging In](#) on page 8
- [Using the calIntegrator2 Workspace](#) on page 8
- [Using Online Help](#) on page 10
- [Logging Out](#) on page 10
- [Application Support](#) on page 11

About calIntegrator2

NCI, Center for Biomedical informatics and Information Technology (CBIIT) is developing a novel translational informatics platform called calIntegrator that allows researchers and bioinformaticians to access and analyze clinical and experimental data across multiple clinical trials and studies. The calIntegrator framework provides a mechanism for integrating and aggregating biomedical research data and provides access to a variety of data types (e.g. Immunohistochemistry (IHC), microarray-based gene expression, SNPs, clinical trials data, etc.) in a cohesive fashion.

calIntegrator2 is a web based or locally installed portal that allows researchers and study managers to access the biomedical informatics infrastructure and data analysis tools established by calIntegrator from one common software platform. As a calIntegrator2 user, you can perform the following tasks:

- Integrate clinical data with genomic and/or imaging data
- Import data of various types in a predefined flat format, and create new studies with multiple study data
- Update an existing study to add new attributes or to add/modify data
- Perform analyses on study data.

Registering as a New caIntegrator2 User

To request a caIntegrator2 user account, you must register as a new user, completing the following steps:

1. Go to the CBIIT caIntegrator2 login page <http://caintegrator2.nci.nih.gov> or use the URL provided by your System Administrator for the caIntegrator2 instance at your institution.
2. Click the **Register Now** hypertext link, under the caIntegrator2 login section in the upper left of the page. This opens the account registration form (*Figure 1.1*).

Register

Security Information

Do you have an LDAP Account?: ☒ Yes ☐ No

Username*:

Password*:

Requested Role(s)*: ☐ Study Manager ☐ Study Investigator

Existing Studies to be Accessed:

Account Details

First Name*:

Last Name*:

Email*:

Organization*:

Address 1*:

Address 2:

City*:

State*:

Country*:

Postal Code*:

Phone*:

Fax:

Figure 1.1 New user account registration form

3. In the Register form, enter the appropriate information¹.
 - **Security Information**
 - **Do you have an LDAP account** [a user profile with your institution] at [NCICB or your institution]?

¹. Items with an asterisk or highlight are required.

If **Yes**, enter your username and case-sensitive password for the purposes of verifying that it is correct. After you submit your request, you can continue to use calIntegrator2 without an account to browse and search available experiments and download data while your account is verified and activated.

–**Username***

–**LDAP Password***

–**Requested role(s)*** – Select one or more of the roles. Roles are described in [Table 1.1](#).

If your LDAP profile is not validated, calIntegrator2 indicates that the LDAP credentials do not check out. You are asked to reenter them, but you can choose to answer no, and the System Administrator will manually ensure you don't get a duplicate LDAP account during provisioning. You can **Cancel** or talk with your System Administrator about the problem.

If you select **No** [you do not have an LDAP account], the text boxes for entering the LDAP account information disappear. You must indicate the role you would like to be assigned in calIntegrator2, and continue entering the appropriate information in the **Account Details** section.

Role	Description	Permissible 1.0 Actions
Study Manager	Creates, owns and manages studies	Create studies Assign annotations to studies Edit studies Search studies Perform analyses on study data
Study Investigator	Investigates and queries the study data	Query study data Save queries Analyze using K-M Plot Analyze using Gene Expression Plots Analyze using GenePattern

Table 1.1 calIntegrator2 role descriptions

◦ **Account Details**

— **First Name***

— **Last Name***

— **Email [address]***

— **Organization***

— **Address [Lines 1* and 2]**

— **City***

- **State***
- **Country***
- **Postal [or Zip] Code***
- **Phone***
- **Fax**

4. Click **Submit Registration Request** to execute the request, or click **Cancel** to abort the registration.

After registration is sent, the screen displays a confirmation message.

At this point, an email containing all of the information you specified in the new user request form is sent to the caIntegrator2 system administrator and an account request confirmation email is also sent to you, the prospective user, at your specified email address. In response, the caIntegrator2 system administrator uses UPT to create your user account and assign the requested roles (in predefined groups like Study Investigator). When your account is created, the system administrator sends you an email to alert you, after which you can login to caIntegrator2.

When your account is registered, the UserID and password you are assigned determines your access rights for the software.

Logging In

To log into caIntegrator2, follow these steps:

1. On the login page, enter your **username** and **password**.
2. Click the **Login** button. If your login is successful, the Welcome to Browse/Study page appears (*Figure 1.2*).

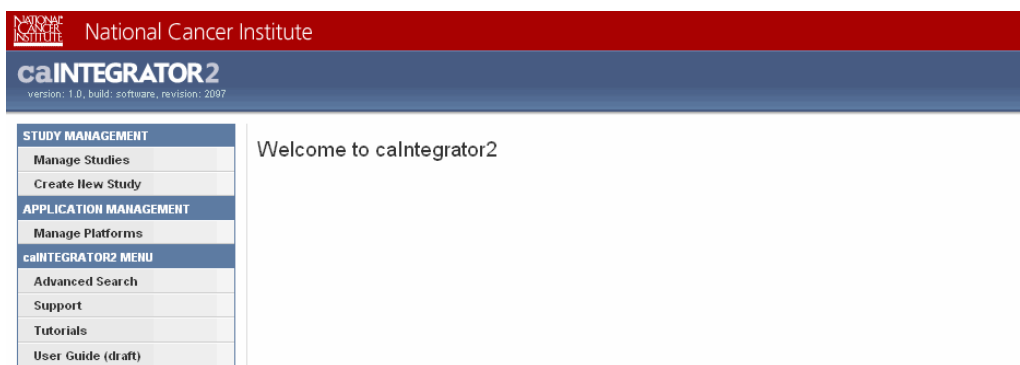


Figure 1.2 The caIntegrator2 workspace before any studies have been deployed

Using the caIntegrator2 Workspace

The caIntegrator2 workspace enables quick access to all caIntegrator2 functions and information. To access caIntegrator2 functions, use the options listed on the left sidebar of the workspace.

caIntegrator2 Functions

When you log into caIntegrator2, before any studies have been created the workspace opens with a Welcome page, as shown in ([Figure 1.2](#)). Once a study is created, its name is listed at the top of the left sidebar.

[Table 1.2](#) describes each caIntegrator2 option in the workspace ([Figure 1.2](#)).

Sidebar Option	Function
[Study Name]	When you log in, one study displays in the left sidebar by default. Any study that you select in the My Studies drop-down list in the upper right of the page replaces this default selection.
Home	Click this to return to the home page for the selected study.
Search [study name]	Click this option to open the Search [study name] page from which you can launch queries into your selected study. For more information, see Searching a caIntegrator2 Study .
Study Data	Click Queries > My Queries to open the list of previous queries you saved. Click any item in the list to open the saved query, which displays on the Criteria, Columns and Sorting tabs. From those tabs, you can modify criteria and/or launch the query again. For more information, see Saving a Query on page 43.
Analysis Tools	Click either of the listed options to open a page where you can launch an analysis of the data in the selected study. <ul style="list-style-type: none"> • Generate a K-M Plot. See Creating Kaplan-Meier Plots on page 58. • Generate a Gene Expression Plots. See Creating Gene Expression Plots on page 65. • Launch GenePattern Analysis. Analyzing Data with GenePattern on page 78.
Study Management	Click either of the listed options to manage the selected study through editing or deleting it or by creating a new study. <ul style="list-style-type: none"> • Click Manage Studies. See Managing a Study on page 31. • Click Create a New Study. See Configuring and Deploying a Study on page 14.
Application Management	Click Manage Platforms to identify, add or remove platforms that caIntegrator2 supports. For more information, see Managing Platforms on page 32.
caIntegrator2 Menu	<ul style="list-style-type: none"> • Click Support to view contact information for Application Support. • Click Tutorials to view a tutorial to help you get started using caIntegrator2. • Click User Guide to open the caIntegrator2 v.1.0 User's Guide in PDF format.

Table 1.2 caIntegrator2 tabs

In the **My Studies** drop-down list in the upper right of the page, select the study you want to use for your current session. (The list includes all studies to which you are subscribed.) As you do so, the following left sidebar contents change to reflect options relevant to your study selection:

- the logo for the selected study (if it exists)

- the name for the selected study
- the list of saved queries for that study

Using Online Help

The online help explains how to use all of the features.

To access online help, click the help icon at the top of each page to open a context-sensitive topic. Context-sensitive help displays information that corresponds to the page from which help was opened.

Help opens displaying the table of contents in the left panel.

Once you are in online help, several buttons and/or options help you locate topics of interest.


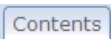

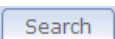




<i>Icon or Button</i>	<i>Description</i>
	Locates and highlights your current topic in the table of contents.
	Select a topic from the complete online help table of contents.
	Select a topic from the online help index.
	Perform word searches of Help by entering query text in the search text box.
	Create a list of your frequently-accessed topics.
 Related Topics 	Opens other closely related topics.
	Prints the current topic.
Topic Name > Topic Name	The breadcrumb trail shows the relative location of the current help topic relative to neighboring topics. Click a breadcrumb link to display that help topic.
Back Forward	Navigates through previously viewed topics.

Table 1.3 Online help tips

Logging Out

To log out of calIntegrator2, click the **logout** link in the upper right-hand corner of the page.

Application Support

For any general information about the application, application support or to report a bug, contact NCICB Application Support.

Email: ncicb@pop.nci.nih.gov	When submitting support requests via email, please include: <ul style="list-style-type: none">• Your contact information, including your telephone number.• The name of the application/tool you are using• The URL if it is a Web-based application• A description of the problem and steps to recreate it.• The text of any error messages you have received
Application Support URL	http://ncicb.nci.nih.gov/NCICB/support
Telephone: 301-451-4384 Toll free: 888-478-4423	Telephone support is available: Monday to Friday, 8 am – 8 pm Eastern Time, excluding government holidays.

CHAPTER 2

CREATING A NEW STUDY

This chapter describes the processes for creating and managing studies in calIntegrator2.

Topics in this chapter include:

- [*Creating a Study – Overview*](#) on this page
- [*Configuring and Deploying a Study*](#) on page 14
- [*Managing a Study*](#) on page 31

Creating a Study – Overview

You can create a calIntegrator2 study by importing clinical study data, genomics data and imaging data, using a combination of spreadsheet/files and existing caGrid applications as source data. Each instance of calIntegrator2 can support multiple studies. As the manager creating a study, it is important that you understand the study well and that the data you wish to aggregate has been submitted to the applications whose data can be integrated in calIntegrator2.

- **Clinical** – The clinical data should be available in CSV files, with a unique patient identifier in one column, one patient per row. Other relevant data can be supplied in other columns to be identified as annotations in the file from within calIntegrator2. You, as the study creator, must have access to the clinical data file, as the file does not come from a caBIG[®] repository.
- **Genomic** – To use calIntegrator2 to integrate array data, the data should be imported into caArray, either locally or the CBIIT installation, using that system's data file import functionality. You must also have a mapping file in CSV format. This file indicates correlations between array files and the clinical subjects in the clinical data files. A mapping file consists of two columns: one with the patient ID, and one with the sample ID.

- **Imaging** – Imaging data should have been submitted to the NBIA grid node as public data, either locally or as part of the CBIIT installation. Image annotations, which includes information about images provided by radiologists or other researchers can include such information as tumor size, tumor location, etc. It must be in CSV format, with unique image series IDs in one column and annotation IDs in the second column. You must also have an image mapping file in CSV format. This file indicates correlations between clinical subjects or images in NBIA and clinical subjects in the clinical data files. A mapping file consists of two columns: one with the patient ID, and one with the NBIA image series ID in the other column.

As you create the study, you define its structure in the process, identifying the data sources and mapping the data between different source data. After the study has been created and deployed, the study can then be used to perform analyses.

Configuring and Deploying a Study

Note: Only a user with a Study Manager role can create a study.

When you create a study, you must specify different data-types (clinical, array, image, etc), data sources (caGrid applications – caArray and NBIA) and map the data, (patient to sample, image series, etc.).

To create a new study, follow these steps:

1. In the Study Management section of the left sidebar, click **Create New Study**.
2. In the Create New Study dialog box that opens, provide a name and description for the study you are creating ([Figure 2.1](#)).

Create New Study

Figure 2.1 Create Study page

3. Click **Save**.

This opens an Edit Study page where you can add identify data files for your study.

Creating/Editing a Study

The Edit Study page displays the Name and Description that you entered for a new study, or for an existing study that you are editing ([Figure 2.2](#)).

Figure 2.2 Edit Study page

To continue creating a study or to modify a study, on the Edit Study page complete these steps:

1. Change the name and or description, if you so choose. Click **Save**.
Note: You can save the study at any point in the process of creating it. You can resume the definition and deployment process later.
2. If you choose to add a logo for the study, click the **Browse** button corresponding to **Logo File**. Navigate for the file, then click **Upload Now**. Once you save the study (or its edit), the logo displays in the center of the page ([Figure 2.3](#)). On the home page for the study, the logo displays in the upper left, above the sidebar.



Figure 2.3 Example of a logo added to the caIntegrator2 browser on the Edit Study page

To continue, you can add clinical data sources, genomic data or imaging data.

Adding Clinical Data

The Edit Study page opens after you save a new study or click to edit an existing study.

Note: To edit information for an existing study, follow the same basic directions in this section. Instead of entering new information, you will modify existing information.

To add or edit clinical metadata in this page, follow these steps:

1. On the Edit Study page, click the **Browse** button in the Clinical Data Sources section. Navigate to locate the file. Files must be in CSV file format.

In the Clinical Data Sources section, if a file has already been selected, its information displays in the varying fields.

2. Click **Add Clinical Data Source**. This opens the Define Fields for Clinical Data page ([Figure 2.4](#)).

Define Fields for Clinical Data

Assign annotation definitions to data fields and click **Done**.

Field Definition	Field Header from File	Data from File
Identifier Change Assignment	PATIENT_ID	151
DC_STUDY_ID Change Assignment	DC_STUDY_ID	B-NCI_U133A_1L.CHP
MICROARRAY Change Assignment	MICROARRAY	NCI_U133A_1L
SITE Change Assignment	SITE	MSKCC
IN_DC_STUDY Change Assignment	IN_DC_STUDY	1
GENDER Change Assignment	GENDER	Male
AGE_AT_DIAGNOSIS Change Assignment	AGE_AT_DIAGNOSIS	64
RACE Change Assignment	RACE	White(01)
ADJUVANT_CHEMO Change Assignment	ADJUVANT_CHEMO	Yes
ADJUVANT_HORMONE Change Assignment	ADJUVANT_HORMONE	No

Figure 2.4 Define Fields for Clinical Data page

The Field Header from File column on the Define Fields... page displays column headers taken from the source *CSV file. The page also displays data values in the file you have designated. You must map each column name to an existing

column name in the caIntegrator2 database or in caDSR. If it doesn't yet exist, you can create a custom column name ([Figure 2.5](#)).

	A	B	C	D	E	F	G	H
1	Pa	Age	Gender	Survival	Disease	Grade	Race	
2	ASP221	50-54	M		ASTROCYTOMA		WHITE	
3	ASP308	50-54	M		GBM		WHITE	
4	FPH113	20-24	M		UNKNOWN		WHITE	
5	FPH114	40-44	M		UNKNOWN		WHITE	
6	FPH118	55-59	M		GBM		WHITE	
7	FPH309	50-54	M		GBM		WHITE	
8	E09238	45-49	M	18-24M	GBM		WHITE	
9	E09239	25-29	M		UNKNOWN		WHITE	
10	E09262	35-39	M		ASTROCYTOMA		WHITE	
11	E09278	30-34	M		UNKNOWN		WHITE	
12	E09331	35-39	M		UNKNOWN		ASIAN NOS	
13	E09332	55-59	M		GBM		WHITE	
14	E09336	30-34	M		GBM		WHITE	
15	E09348	60-64	M		GBM		WHITE	
16	E09378	45-49	M		UNKNOWN		WHITE	
17	E09449	50-54	M		UNKNOWN		OTHER	
18	E09454	0-4	M		UNKNOWN		WHITE	
19	E09489	55-59	M		GBM		WHITE	
20	E09515	35-39	M		UNKNOWN		WHITE	
21	E09569	45-49	M		UNKNOWN		WHITE	
22	E09587	35-39	M		UNKNOWN		OTHER	
23	E09601	40-44	M		GBM		WHITE	
24	E09610	55-59	M		GBM		WHITE	
25	E09611	60-64	M		UNKNOWN		ASIAN NOS	
26	E09615	45-49	M		UNKNOWN		WHITE	
27	E09624	35-39	M		GBM		WHITE	
28	E09645	45-49	M		UNKNOWN		WHITE	
29	E09657	50-54	M		UNKNOWN		WHITE	
30	E09730	40-44	M		UNKNOWN		WHITE	

Figure 2.5 Example of a source CSV file whose data you are mapping in caIntegrator2

The MOST important steps in creating a new study in caIntegrator2:

- You MUST designate one column in the file as a unique “identifier” column type.
- You MUST review and define column annotation definitions for each column header in the file.

If caIntegrator2 “recognizes” the same column header in other files already in the system, a term, for example “age” or “survival”, which is the current definition appears in the **Field Definition** column above the blue **Change Assignment** link. If the area above the blue **Change Assignment** link is blank, no correlating term exists in the database; you must specify the field type, and then the term will populate the space.

3. To indicate the unique identifier of choice, on the row showing the column header (PatientID in the figure, but other examples are subject identifier, sample identifier, etc), click **Change Assignment** in the **Field Definition** column.

Assigning An Identifier or Annotation

When you click **Change Assignment** on the Define Fields... page, the Assign Annotation Definition for Column dialog box opens ([Figure 2.6](#)). On this page you can change the column type and the field definition for the specific data field you selected.

Note: When you change an assignment, you must make sure the data types match--numeric, etc.

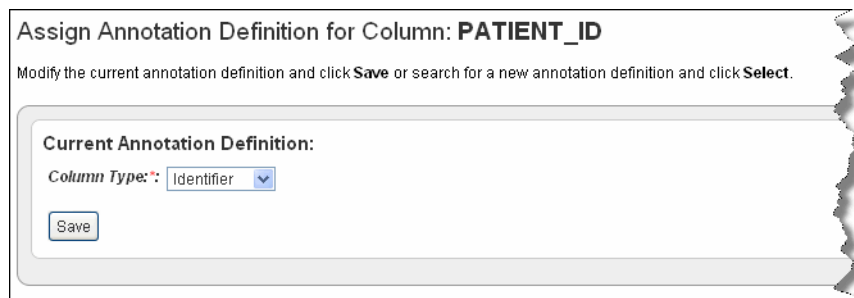


Figure 2.6 The Assign Annotation Definition dialog box

1. For the column (PatientID) that you choose to be the one and only Identifier column, in the **Column Type** drop-down list, select **Identifier**.
2. Click **Save** to save the identifier. This returns you to the Define Fields for Clinical Data page where the Identifier is noted in the Field Definition column.
3. After you have defined which field is the Identifier, you must ensure that ALL other fields also have a field definition assignment. For those fields without a Field Definition assignment or for those whose Field Definition you want to review, click **Change Assignment**.
4. In the Assign Annotation Definition for Column: [column header] dialog box, select **Annotation** in the drop-down list. The page extends, displaying additional fields for fleshing out the annotation description.

As you select the column type, you can work with column headers in one of four ways in this dialog box.

- You can accept existing default definitions (those that are inherent in the data file you selected). See Step 5.
 - You can create your own definitions manually. See Step 6.
 - You can search for and use definitions in other caIntegrator2 studies. See [Searching for Annotation Definitions](#) on page 20.
 - You can search for and use definitions found in caDSR. See [Searching for Annotation Definitions](#) on page 20.
5. If there is anything you want to change about an existing annotation definition of the field such as its name, or if you want to view or edit its definition, click the **Change Assignment** link on the Define Fields... page. The Assign Definition page opens, expanded now to include a Current Annotation Definition section above a section where you can still initiate a search for an annotation definition ([Figure 2.7](#)).

Note: If the column header you are working with already has a designated Field Definition, the Current Annotation Definition section of the Assign

Annotation Definition for Column page is already visible when you open this dialog box.

Assign Annotation Definition for Column: **MICROARRAY**

Modify the current annotation definition and click **Save** or search for a new annotation definition and click **Select**.

Current Annotation Definition:

Column Type: Annotation

Name: MICROARRAY

Definition: Created via selenium for DC Lung Full on 09/14/09 15:11:17.

Keywords: MICROARRAY

Data Type: string

Non-Permissible

Permissible

Permissible Values:

Add >

< Remove

New

Save

Search for an Annotation Definition:

Search

Search existing studies and caDSR for definitions.

Figure 2.7 Current Annotation Definition

- To enter a new name annotation, or any other information about the annotation definition, click the **New** button and enter the information described in [Table 2.1](#)

Annotation Field	Field Description
Name	Enter the name for the annotation.
Definition	Enter the term(s) that define the annotation.
Keywords	Insert keyword(s) that can be used to find the annotation in a search, separated by commas.
Data Type	Enter a string (default), numeric, or date

Table 2.1 Annotation fields for new definitions

Annotation Field	Field Description
Permissible/Non-permissible Values	<p>Note: The first time you load a file, before you assign annotation definitions (step 3 on page 17), these panels may be blank. If the column header for the data is already “recognizable” by calIntegrator2, the system makes a “guess” about the data type and assigns the values to the data type in the newly uploaded file. They will display in the Non-permissible values sections initially. Use the Add and Remove buttons to move the values shown from one list to the other, as appropriate.</p> <p>When you select or change annotation definitions by selecting matching definitions (described in Searching for Annotation Definitions on page 20), this may add (or change) the list of non-permissible values in this section.</p> <p>If you leave all values for a field in the Non-permissible panel, then when you do a study search, you can enter free text in the query criteria for this field.</p> <p>If there are items in the Permissible values list, then the values for this annotation are restricted to only those values. When you perform a study search, you will select from a list of these values when querying this field. If there are no items in the permissible values list then the field is considered free to contain any value.</p> <p>To edit a field's permissible values, you must change the annotation definition. You can do this even after a study has been deployed.</p>

Table 2.1 Annotation fields for new definitions

Searching for Annotation Definitions

An alternative to creating a new definition is to search for annotation definitions already present in calIntegrator2 studies or in caDSR.

1. Enter search keyword(s) in the **Search** text box on the Assign Annotation Definition page. Click **Search**. After a few moments, the search results display on the page (*Figure 2.8*).

Search for an Annotation Definition:

microarray Search existing studies and caDSR for definitions.

Matching Annotation Definitions from caintegrator2				
Name	Actions	CDE Public ID	Data Type	Definition
MICROARRAY	Select		string	Created via selenium for DC Lung Full on C
MICROARRAY	Select		string	Created via selenium for DC Lung Full on C

Matching Annotation Definitions from caDSR					
Name	Actions	CDE Public ID	Context	Status	Definition
Microarray Microarray Analysis Data float One Dimensional Array	Select View	2658378	caBIG	RELEASED	A microarray is a piece of glass or plastic on which different samples have been affix are usually DNA fragments but may also be antibodies, other proteins, or tissues. _Ana microarrays to profile the pattern of proteins). A collection or single item of factual infor drawn. _Generic value domain for a single dimensional array with floating numbers as i
Microarray Identifier java.lang.Long	Select View	2223905	caCORE	RELEASED	A microarray is a piece of glass or plastic on which different samples have been affix are usually DNA fragments but may also be antibodies, other proteins, or tissues. _One

Figure 2.8 Results for annotation definition search

2. To view the definitions corresponding to any of the “Matching Annotation Definitions”, which are those currently found in other caIntegrator2 studies, click the [term], such as “age”, hypertext link. The definition then appears in the Current Annotation Definition segment of the page just above.

In summary, when you click the link, that assigns the definition to the Define Fields for Clinical Data page, and it also closes the Annotation Definition page.

You can modify any portion of the definition, as described in [step 6](#) on page 19.

3. The matches from caDSR display some of the details of the search results. To view more details of a match, such as permissible values, click **View**, which opens caDSR to the term. If you click **Select**, the caDSR definition automatically replaces the annotation definition for this field with which you are working.

Caution: Take care before you add a caDSR definition that it says exactly what you want. caDSR definitions can have minor nuances that require specific and limited applications of their use.

4. Once you have settled on an appropriate field definition for the annotation, click **Save**. This returns you to the Define Field for Clinical Data page.

Note: If you have not clicked **Select** for alternate definitions in this dialog box, then click **Save** to return to the Define Field...dialog box without making any definition changes.

5. From the Define Fields for Clinical Data page, be sure and designate the annotations for each field in the file. Click **Save** on each page to save your entries or click **New** to clear the fields and start again. You will not be able to proceed until every Field Definition entry on the Fields for Clinical Data screen has a unique entry, one as an Identifier and the remainder as annotations.

The Data From File columns on the page display the column header *values* of the first three rows you designated as “annotations”.

6. Click **Done**. This saves the study by name and description, but does not deploy the study. See [Deploying the Study](#) on page 31.

Saving the study returns you to the Edit Study page where a “Not Loaded” status now appears for the file whose annotations (column headers) you have defined ([Figure 2.9](#)).

Study Overview

Study Name: test jbh

Study Description:

Status: Not Deployed

Status Description:

Study Logo:

Study Logo

Logo File: Browse...

JPEG/GIF, 200x72 maximum

Upload Now

Clinical Data Sources

Add New Edit Survival Values

Type	Description	Status	Action
DELIMITED_TEXT	dc_lung_clinical_data.csv	Not Loaded	

Figure 2.9 Example file whose annotations have been defined

7. Click the **Load Clinical** link in the Action section to load the data file you configured. At this point, the Status changes to “Loaded”.

Note: You can add as many files as are necessary for a study. Patients 1-20 in first file, 21-40 in second file, or many patients in first file and annotations in second file, etc. As long as IDs are defined correctly, it works.

8. Once you have assigned data types to every column header in the data file and have loaded the clinical data, click **Save and Deploy**. At that point, calIntegrator2 loads data from the file to the calIntegrator2 database.

Note: You can change assignments even after the study is deployed, using the Edit feature. For more information, see [Creating/Editing a Study](#) on page 15.

The Manage Studies page opens when the study is deployed. The **Deployed** status is indicated on the Manage Studies page as well as the Edit Study page. For more information, see [Managing a Study](#) on page 31.

You can continue to perform other tasks in calIntegrator2 while deployment is in process.

See also [Deploying the Study](#) on page 31.

Note: You can repeatedly upload additional or updated subject annotations, samples, image data, array data to the study at later intervals. These later imports do not remove any existing data; they instead insert any new subjects or update annotations for existing subjects.

Defining Survival Values

Survival value is the length of time a patient lived. If you plan to analyze your data in calIntegrator2 to create a Kaplan-Meier (K-M) Plot, then during the Annotation Definition process described above, you must make sure that you have defined at least three fields set to the “date” Data Type. These will be matched to the following three properties during Survival Value definition.

- **Survival Start Date**
- **Death Date**
- **Last Followup Date**

Note: Setting survival values is optional if you do not plan to use the K-M plot analysis feature or if you do not have this kind of data (survival values) in the file.

For some applications, such as REMBRANDT and I-SPY, survival values are pre-defined in the databases when you load the data. In calIntegrator2, however, you can review and define survival value ranges in a data set you are uploading to a study. To be able to do so, you need to understand the kind of data that can comprise the survival values.

To set up survival values, follow these steps:

1. On the Edit Study page, click **Edit Survival Values**. This opens the Survival Value Definitions dialog box ([Figure 2.10](#)).

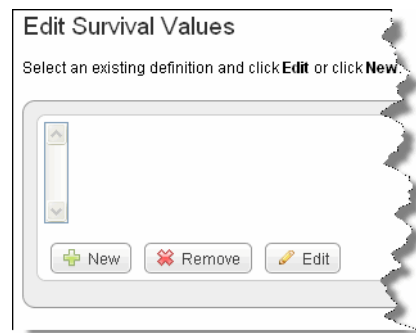


Figure 2.10 Survival Value Definition dialog box

2. Click **New** to enter new survival value definitions.

- OR -

Click **Edit** to edit existing survival value definitions.

- The dialog box extends, now displaying three drop-down lists that show column headers for date metadata in the spreadsheet you have uploaded. [Figure 2.11](#) displays survival value ranges that have already been added to a study.

Survival Value Definitions for 'Demo Study based on 101 Long-Full data'

Survival from enrollment

New

Remove

Edit

Survival Value Definition Properties for 'Survival from enrollment'

Name:	Survival from enrollment
Survival Start Date:	ENROLLMENT_DATE
Death Date:	DEATH_DATE
Last Followup Date:	LAST_CONTACT_DATE

Save

Figure 2.11 Survival Definitions example

In the drop-down lists, select the appropriate survival value definitions for each field listed. You might want to refer to the column headers in the data file itself. Dates covered by the definitions are already in the data set. You cannot enter specific dates.

- **Name** – Enter a unique name that adequately describes the survival values you are defining here. *Example:* Survival from Enrollment Date or Survival from Treatment Start. The name you enter displays later when you are selecting survivals to create the K-M plot.
- **Survival Start Date** – Select the column header for this data.
- **Death Date** – Select the column header for this data.
- **Last Followup Date** – Select the column header for this data.

See also [Creating Kaplan-Meier Plots](#) on page 58.

Adding/Editing Genomic Data

Note: Genomic data must be parsed and stored in caArray to be able to use it in caIntegrator2.

Once you have loaded clinical data and identified patient IDs, you can add either array genomic sample data from caArray, which caIntegrator2 maps by sample IDs to the patient IDs in the clinical data, covered in this section, or you can load imaging files from NBIA, also mapped by IDs to the patient data, covered in [Adding Imaging Data](#) on page 29. You can also edit genomic data information that you have already added to the study. Genomic sample data and imaging data are independent of each other, so neither is required before loading the other.

It is essential that you are well acquainted with the data you are working with--the clinical data, and the corresponding array data in caArray.

caIntegrator2 supports a limited number of array platforms. For more information, see [Managing Platforms](#) on page 32.

To add genomic data to your caIntegrator2 study, follow these steps:

1. On the Edit Study page where you have selected and added the clinical data, click the **Add** button under Genomic Data Sources. You can upload genomic data only from caArray.

This opens the Edit Genomic Data Source dialog box (Enter the appropriate information in the fields (*Figure 2.12*).

Edit Genomic Data Source

caArray Server Hostname:	array.nci.nih.gov
caArray server JNDI Port:	0
caArrayUsername:	
caArrayPassword:	
caArray Experiment Id:	
Vendor:	Affymetrix
Data Type:	Expression
Platform (only needed for Agilent):	

Cancel Save

Figure 2.12 Edit Genomic Source dialog box

- **caArray Host Name** – Enter the hostname for your local installation or for the CBIIT installation of caArray, array.nci.nih.gov. If you misspell it, you will receive an error message.
 - **caArray JNDI Port** – Enter the appropriate server port. See your administrator for more information. *Example:* For the CBIIT installation of caArray, enter **8080**.
 - **caArray Username** and **caArray Password** – If the data is private, you must enter your caArray account user name and password; you must have been given permissions in caArray for the experiment. If the data is public, you can leave these fields blank.
 - **caArray Experiment ID** – Enter the caArray Experiment ID which you know corresponds with the clinical data you uploaded. *Example:* Public experiment “beer-00196” on the CBIIT installation of caArray (array.nci.nih.gov). If you misspell your entry, you will receive an error message.
 - **Vendor** – Select either **Agilent** or **Affymetrix**
 - **Data Type** – Select **Expression** or **Copy Number**.
 - **Platform (needed only for Agilent)** – If appropriate, select the **Agilent** platform.
2. Click **Save**.

caIntegrator2 goes to caArray, validates the information you have entered here, finds the experiment and retrieves all the sample IDs in the experiment. Once

this finishes, the experiment information displays on the Edit Study page under the Genomic Data Sources section ([Figure 2.13](#)).

Host Name	Experiment Identifier	File Description	Data Type	Status	Action
nci-as-d227-v.nci.nih.gov	admin-00001	Mapping File(s): nci_sample_mapping.csv Control Sample Mapping File(s): jags_0034_control_samples.csv	Expression	Loaded	Edit Map Samples Delete

Figure 2.13 Genomic Data Sources section of the Edit Study page

3. If you want to redefine the caArray experiment information, you can edit it. Click the **Edit** link corresponding to the Experiment ID. The Edit Genomic Data Source dialog box reopens, allowing you to edit the information.

Note: At any point in the process of working within a study, you can create a gene list. For more information, see [Creating a Gene List](#) on page 47.

Mapping Genomic Data to Clinical Data

Because the goal of caIntegrator2 is to integrate data from clinical, genomic and imaging data sources, data from uploaded source files must be mapped to each other.

To map the samples from the caArray experiment to the patients in the clinical data you uploaded, follow these steps:

1. On the Edit Study page, click the **Map Samples** link. This opens the Edit Sample Mappings page ([Figure 2.14](#)).

Upload mapping files and click **Map Samples**.

Data Source

caArray Server Hostname: array-stage.nci.nih.gov

caArray server JNDI Port: 8080

caArray Username:

caArray Password:

caArray Experiment Id: rembr-00037

Subject to Sample Mapping File: [Browse...](#)

Control Sample Set Name*:

Control Samples File: [Browse...](#)

[Map Samples](#)

Control Sample Sets

Set Name	Sample Name

Sample Mappings

Unmapped Samples

Sample Name
E0907B_U133P2
GeneratedSample.ASTROCYTOMA_I_E07733_U133P2
GeneratedSample.ASTROCYTOMA_I_E09137_U133P2
GeneratedSample.ASTROCYTOMA_I_E09629_U133P2
GeneratedSample.ASTROCYTOMA_I_E09820_U133P2_2
GeneratedSample.ASTROCYTOMA_I_E09821_U133P2
GeneratedSample.ASTROCYTOMA_I_E09948_U133P2
GeneratedSample.ASTROCYTOMA_I_E09971_U133P2

Figure 2.14 Edit Sample Mappings page

When you first open this page, all of the samples in the caArray experiment you selected are listed as unmapped, because calIntegrator2 does not know how these sample names correlate to the patient data in the clinical file until you upload the subject to sample mapping file.

2. At the top of the page, click **Browse** to navigate for the CSV file that identifies the mapping information. Click the **Upload Mapping File** button.

The mapping file has only two columns (typically without headers)—one that shows the subject ID (designated in calIntegrator2 as the “Identifier”) and one that has “Sample name” field from the linked caArray experiment, with one subject per row (*Figure 2.15*). This provides calIntegrator2 with the information for mapping patients to caArray samples.

	A	B	C	D	E	F	G	H
1	E10216	GeneratedSample.UNKNOW	DISEASE_L_E10216_U133P2					
2	E10144	GeneratedSample.UNKNOW	DISEASE_L_E10144_U133P2					
3	E09212	GeneratedSample.UNKNOW	L_20070227_16-22-37-238_E09212_U133P2					
4	E09369	GeneratedSample.UNKNOW	L_20070227_16-22-37-238_E09369_U133P2					
5	E10162	GeneratedSample.UNKNOW	DISEASE_L_E10162_U133P2					
6	E10318	GeneratedSample.UNKNOW	DISEASE_L_E10318_U133P2					
7	E09264	GeneratedSample.OLIGO_L_20070227_11-27-27-881_E09264_U133P2						
8	E10252	GeneratedSample.UNKNOW	DISEASE_L_E10252B_U133P2					
9	E09624	GeneratedSample.GBM_L_20070226_13-30-40-39_JF0142_U133P2						

Figure 2.15 Example sample mapping file, in CSV format

Note: When you open the mapping file, make sure that the patient ID is used for mapping.

Unmapped samples continue to show at the top of the calIntegrator2 page. They were loaded from caArray, but they are not in the mapping file. These are not used for integration.

3. Scroll down the page to see samples that are mapped to the patients in the clinical data (*Figure 2.16*).

1338	GeneratedSample.OLIGO_L_20070227_11-49-51-876_JF0142_U133P2	
1338	GeneratedSample.UNKNOW	DISEASE_L_E10029_U133P2
1339	GeneratedSample.GBM_L_20070226_14-05-29-569_JF1356_U133P2	
1340	GeneratedSample.OLIGODENDROGLIOMA_L_JF0599_U133P2	
1342	GeneratedSample.GBM_L_20070226_13-30-40-39_JF0142_U133P2	
1345	GeneratedSample.GBM_L_20070226_14-31-20-427_JF1409_U133P2	
Samples Mapped to Subjects		
Sample ID	Sample Name	Subject Identifier
901	GeneratedSample.UNKNOW	DISEASE_L_E10216_U133P2
911	GeneratedSample.UNKNOW	DISEASE_L_E10144_U133P2
914	GeneratedSample.UNKNOW	L_20070227_16-22-37-238_E09212_U133P2
918	GeneratedSample.UNKNOW	L_20070227_16-22-37-238_E09369_U133P2
922	GeneratedSample.UNKNOW	DISEASE_L_E10162_U133P2
925	GeneratedSample.UNKNOW	DISEASE_L_E10318_U133P2
930	GeneratedSample.OLIGO_L_20070227_11-27-27-881_E09264_U133P2	
940	GeneratedSample.UNKNOW	DISEASE_L_E10252B_U133P2
954	GeneratedSample.GBM_L_20070226_13-30-40-39_JF0142_U133P2	
957	GeneratedSample.ASTROCYTOMA_L_E09137_U133P2	
958	GeneratedSample.UNKNOW	DISEASE_L_E09890_U133P2
968	GeneratedSample.UNKNOW	L_20070227_16-57-07-283_E09515_U133P2
1004	GeneratedSample.UNKNOW	L_20070227_17-26-09-910_E09722_U133P2

Figure 2.16 Example of samples mapped to patients' data

Uploading Control Samples

A Control Samples file is used to calculate fold change data, which compares “tumor” sample gene expression in the caArray experiment to the control samples to identify

those that exhibit up or down gene regulation. Control samples can be the “normal” samples, but that is not necessarily the case.

To upload the control samples, follow these steps:

1. On the Edit Sample Mappings page, click the **Map Samples** link.
2. Click **Browse** to navigate for the control samples file, and click the **Upload Control Samples** File button. The control sets display at the top of the page once they have been uploaded (*Figure 2.17*).

Set Name	Sample Name
Rembrandt controls	GeneratedSample.Normal_L_20070227_14-22-24-128_Normal_5_U133_P2
	GeneratedSample.Normal_L_20070227_14-01-17-731_HF0526_U133P2
	GeneratedSample.Normal_L_20070227_14-01-17-731_HF0131_U133P2
	GeneratedSample.Normal_L_20070227_14-01-17-731_HF0523_U133P2
	GeneratedSample.Normal_L_20070227_14-01-17-731_HF0120_U133P2
	GeneratedSample.Normal_L_20070227_14-01-17-731_HF0137_U133P2
	GeneratedSample.Normal_L_20070227_14-22-24-128_Normal_6_U133_P2
	GeneratedSample.Normal_L_20070227_14-01-17-731_HF0151_U133P2
	GeneratedSample.Normal_L_20070227_14-01-17-731_HF0398_U133P2

Figure 2.17 Example list of control samples

The control samples now display toward the bottom of the page.

3. This information will be used when performing other tasks in caIntegrator2, to be described in other sections.

Configuring Copy Number Data

You can add copy number data for a genomic data source by uploading the mapping file. This allows you to configure parameters to be used when segmentation data is being configured.

The name specified in the third column of the mapping file is specific for each array manufacturer as follows:

- Affymetrix – The third column of the mapping file must contain filenames that end in cnchp. The corresponding experiment in caArray must have these files and the extensions must match .cnchp.
- Agilent – The third column must name a file which contains level 2 copy number data. Level one copy number will not work. This file name is repeated for each line in the mapping file.

To add copy number data relating to the genomic data you are adding, follow these steps:

1. In the Genomic Data Sources section, for the data you have already added, click **Configure Copy Number Data** hypertext link.

Note: This link is available only if you have uploaded copy number data and you are configuring a Copy Number data type (as indicated by the Data Type column on the Edit Study page).

The Edit Copy Number page opens (*Figure 2.18*).

Figure 2.18 Edit Copy Number page

2. Browse for and enter appropriate information to identify the copy number mapping file. The fields are described in *Table 2.2*. An asterisk* indicates a required field..

Field	Description
Subject and Sample to Copy Number Mapping File	Browse for the appropriate CN mapping file
caDNACopy Service URL*	Control for selecting the URL which hosts the caDNACopy grid service
Change Point Significance Level	Significance levels for the test to accept change-points
Early Stopping Criteria	The sequential boundary used to stop and declare a change
Permutation Replicates	The number of permutations used for p-value computation
Random Number Seed	The segmentation procedure uses a permutation reference distribution. This should be used if you plan to reproduce the results.

Table 2.2 Fields for retrieving a copy number mapping file.

3. Click **Configure copy number data** for a genomic data source. On the screen upload a copy number mapping file (format: subject id, sample id, file name) and configure the parameters to be sent when computing segmentation data.

Adding Imaging Data

Once you have loaded clinical data and identified patient IDs, you can add either array genomic sample data from caArray which caIntegrator2 maps by sample IDs to the patient IDs in the clinical data, or you can load imaging files from NBIA, also mapped by IDs to the patient data, covered in this section. Genomic sample data and imaging data are independent of each other, so neither is required before loading the other.

It is essential that you are well acquainted with the data you are working with--the clinical data, and the corresponding imaging data in NBIA.

Any data in NBIA can be uploaded to calIntegrator2. Imaging data consist of images and or annotations for images.

To add imaging data to the study you are creating or are editing, follow these steps:

1. On the Edit Study page, click the **Add** button under Imaging Data Sources section. Imaging data can be NBIA images or image annotations, which are uploaded in spreadsheet format.

This opens the Edit Imaging Data Source dialog box. Enter the appropriate information in the fields (*Figure 2.19*). Asterisks indicate required fields..

Edit Imaging Data Source

Enter a NBIA Data Source and Image Mapping Data from a file and click **Save**.

Figure 2.19 Edit Image Data Source dialog box

- **NBIA Server Grid URL*** – Enter the URL for the grid connection to NBIA
- **NBIA Username and NBIA Password.** This information is not required, as currently all data in the NBIA grid is Public data.
- **Collection Name*** – Enter the name/source for the collection.
- **Current Mapping** – If a mapping file has already been uploaded to the study to map imaging data, the file name displays here.
- **Select Mapping File Type*** – Click to select the file type:
 - **Auto** – No file is required. Selecting this takes all clinical subject IDs and attempts to map them to the corresponding ID in the collection in NBIA. If the ID does not exist in NBIA, then no mapping is made for that ID.
 - **By Subject** – Requires a file to be uploaded. The “clinical to imaging mapping file” must be a two column mapping (CSV) from the calIntegrator2 clinical subject ID to the NBIA subject ID.
 - **By Image Series** – Requires a file to be uploaded. The clinical to imaging mapping file needs to be a two column mapping (CSV) from the calIntegrator2 clinical subject ID to the NBIA study instance UID.

- **Clinical to Imaging Mapping File** – Click **Browse** to navigate to the appropriate clinical to imaging mapping file. See **Select Mapping File Type*** field description.
- 2. Click **Add** to upload the data to caIntegrator2.

The imaging data information displays on the Edit Study page under the Imaging Data Sources section ([Figure 2.20](#)).

Host Name	Collection Name	File Description	Status	Action
imaging.nci.nih.gov	NCRI	Annotation File: ncri_image_annotations.csv Mapping File: ncri_image_mapping.csv	Loaded	Edit Edit Annotations Delete

Figure 2.20 Imaging Data Sources section of the Edit Study page

- 3. Once the data is uploaded, you must assign identifiers and annotations to the data in the same way you did with the clinical data. For more information, see [Assigning An Identifier or Annotation](#) on page 17 and [Searching for Annotation Definitions](#) on page 20.
- 4. To deploy the study, see [Deploying the Study](#).

Deploying the Study

When you are ready to deploy the study, click the **Deploy Study** button on the Edit Study page. caIntegrator2 retrieves the selected data from the data service(s) you defined and makes the study available to a study manager or to anyone else who may want to analyze the study's data. Using the Manage Studies feature, you can then configure and share data queries and data lists with all investigators who access the study.

Note that you can continue to work in caIntegrator2 while study is being deployed.

Managing a Study

Note: A user without management privileges has no access to this section of caIntegrator2.

Once you have started to create a study or have deployed it, you can update an existing study in the following ways:

- Add new attributes (annotations) and upload relevant data to an existing study.
- Delete a study
- Modify existing annotation definitions
- Reload subset of study data and re-deploy the study and perform new analyses
- Re-deploy the entire study with new set of data and mappings.

To update, edit or delete a study, follow these steps:

1. On the left sidebar, click **Manage Studies**. The Manage Studies page appears (*Figure 2.21*).

Manage Studies

View studies and click **Edit** to modify or click **Delete**.

Name	Description	Last Modified By	Status	Deployment Start Date	Deployment Finish Date	Action
Demo Rembrandt TCOA Agilent Copy Number Level 2	Mapping 4 samples to 2 identifiers.	manager	Deployed	2009/10/28 13:06:06	2009/10/28 17:10:56	Edit Delete
Demo Study based on DC Lung Full data.	DC Lung Full. Study created via selenium.	manager	Deployed	2009/10/27 17:30:24	2009/10/27 19:14:56	Edit Delete
Demo Study based on Rembrandt with NCR1 data.	Rembrandt with NCR1. Study created via selenium.	manager	Deployed	2009/10/29 13:21:52	2009/10/29 13:24:31	Edit Delete
Demo Study based on Rembrandt with no images data.	Rembrandt with no images. Study created via selenium.	manager	Deployed	2009/10/27 17:21:46	2009/10/27 20:09:24	Edit Delete
Demo Study based on Small Copy Number data.	Small Copy Number. Study created via selenium.	manager	Deployed	2009/10/27 17:37:34	2009/10/27 21:03:00	Edit Delete
Demo Study based on TCOA Agilent data	TCOA Agilent. Study created via selenium.	manager	Deployed	2009/10/27 17:32:33	2009/10/27 22:04:35	Edit Delete

Figure 2.21 Manage Studies page

All of the "in process" or "completed" studies display on this page, with associated metadata.

2. Click the **Edit** link corresponding to your study of choice to open the Edit Studies page.

On this page you can edit any details such as adding or deleting files, survival values, and so forth. For information about working in the Edit Study, see *Creating/Editing a Study* on page 15.

3. Click the **Delete** link to delete the corresponding study.

Managing Platforms

calIntegrator2 supports a limited number of array platforms, all of which originate from Agilent or Affymetrix. While they do not represent all of the platforms supported by caArray, calIntegrator2 must have array definitions loaded for the platforms it supports, and be able to properly load the data from caArray and parse it.

You can create a study without genomic data, but you cannot add genomic data to a calIntegrator2 study without a corresponding supported array platform.

On the Manage Platforms page, you can identify, add or remove supported platforms.

To manage platforms in calIntegrator2, follow these steps:

1. Click **Manage Platforms** on the left sidebar.

The Manage Platforms page that opens lists the platforms calIntegrator2 currently supports, those that the system can pull from caArray (*Figure 2.22*).

You can also add a new platform by entering information in the fields at the top of the page.

The screenshot shows the 'Manage Platforms' interface. At the top right, there is a '(draft)' status indicator. The form includes the following fields:

- Platform Type:** A dropdown menu currently set to 'Affymetrix: Gene Expression'.
- Platform Name (For NON-GEML.xml file):** An empty text input field.
- Annotation File:** An empty text input field with a 'Browse...' button to its right.
- Add Annotation File:** A button below the Annotation File field.
- Annotation File(s) Selected:** A large text area for listing selected files, with up and down arrow icons on the right.
- Create Platform:** A button at the bottom of the form.

Below the form is a table listing existing platforms:

Name	Vendor	Reporter List	Action
AgilentG4502A_07_01	AGILENT	AgilentG4502A_07_01, AgilentG4502A_07_01	Delete
GeneChip Human Mapping 100K Set	AFFYMETRIX	Mapping50K_Hind240, Mapping50K_Xba240	None
HG-U133A	AFFYMETRIX	HG-U133A, HG-U133A	Delete
HG-U133A	AFFYMETRIX	HG-U133A, HG-U133A	None

Figure 2.22 Manage Platforms page

- To add a platform, in the Platform Type field, select the appropriate platform type from the drop down list. Click **Browse** to navigate for the Affymetrix or Agilent file you want to add.
- Enter a **Platform Name** if the file is a NON-GEML.xml file.
Depending on what Platform Type is selected, there may be other parameters to provide here as well. Once all parameters have been provided, click **Create Platform**.
- Click the **Browse** button to browse for the appropriate annotation file. When you have located it, click **Add Annotation**. The system displays annotation files you select in the Associated File(s) Selected box.
- Click the **Add** button.

CHAPTER 3

SEARCHING A CAINTEGRATOR2 STUDY

This chapter describes the processes for searching studies within calIntegrator2.

Topics in this chapter include:

- [Search Overview](#) on this page
- [Searching a Study](#) on page 36
- [Managing Queries](#) on page 43

Search Overview

The search and browse functions in calIntegrator2 allow you to search for clinical data, genomic or imaging data that were uploaded into the application as part of a study. When gene expression and imaging data are uploaded into a calIntegrator2 study, mapping files that correlate the data in those files to patient IDs in the clinical data file must also be uploaded. When you launch a search, calIntegrator2 finds and integrates the clinical, genomic and imaging data based on the mapping files and the criteria that you define in the search query.

In a search query, you can specify criteria for just one of the data types, or configure complex search criteria that join two or three data types. The available criteria for the query were defined when the study was deployed.

The basic workflow for a study search follows these steps:

1. Select the study to be searched.
2. Select one data type:
 - **Clinical** – searches one or more uploaded CSV files for data identifiers or annotations (column headers) specified when the study was created
 - **Genomic** – Searches caArray experiments samples uploaded in the study for gene expression data by gene name or reporter ID.

- **Image Series** – Searches NBIA files uploaded in the study for image annotations or links to images, identified by subject identifiers or image series IDs.
3. Define criteria for the search in the selected data type and run the search.
 4. For a more complex search, select multiple criteria from more than one data type.
 5. Specify whether you want clinical/imaging annotations to display or genomic data to display.
 6. Review search results.
 7. Configure results column and sorting display settings. You can do this before or after you run a search. If you choose to do it after, you must re-run the search.
 8. Download annotation search results as a CSV file. The CSV file contains only the data you specified in the annotation and display configurations.
 9. Follows links to NBIA in the search results to view or download images located in the search.

Searching a Study

To initiate a search of all annotations and/or other data in a study, follow these steps:

1. In calIntegrator2, in the upper right hand corner, select the study you want to browse or perform a simple search.
2. On the left sidebar, under the first section that displays the study name, click **Search [Study Name]**. This opens a simple search query page with five tabs (*Figure 3.1*).

Search Demo Study based on Rembrandt with NCRI data.

The screenshot shows the 'Criteria' tab of the search interface. It features a dropdown menu with 'Clinical' selected and an 'Add' button. A message below the dropdown states 'No criteria added. Please select criteria from the pulldown box.' At the bottom, there are radio buttons for 'or' and 'and', and a 'Run Query' button.

Figure 3.1 Search page

3. On the Criteria tab, in the drop-down list, select the type of data you want to search. The listed options reflect the type of data that have been uploaded to the study.

Note: You can perform a search using one or more criteria you set in one of the data types, or you can define criteria in more than one data type per query, creating a more complex search.

- **Clinical**
- **Gene Expression**
- **Image Series**

4. Click **Add** to define annotation elements for the search.

Continue with:

Clinical and Image Series on page 37

Gene Expression on page 38

5. To add additional criteria for the search, repeat steps 3 and 4, as appropriate. You can set more than one data type or more than one criterion for a data type. The criteria become cumulative, thus refining the search.
6. Once you have configured the query criteria, select the Boolean **Or** or **And** search operator at the bottom of the page.
- **Or** finds a data subset with at least one of the search criteria
 - **And** finds a data subset with both/or all search criteria.
7. Click the **Remove** button to clear any data elements you have defined.
8. You can launch the search from this tab. Click the **Run Search** button. For information about the search results, see *Chapter 4 Viewing Query Results*. You may want to run the search first to see what kind of results you get before you configure the data display, described in step 9.
- or –
9. On the Results Type tab, you can specify the columns you want to display in the search results data. On the Sorting tab, you can specify how the data is to be sorted. For more information, see *Results Type Tab* on page 41 and *Sorting Tab* on page 42.

Note: As long as you are still in the current query session, you can return to the Criteria, Columns and Sorting tabs to add, modify or remove data and display criteria and re-run the search. If you configure another query without saving the first, the first query will be lost. If you save the query, your current search criteria are saved.

Clinical and Image Series

- If you select Clinical or Image Series data types, an additional drop-down list displays data elements that are annotation definitions specified when

the data was uploaded into the study ([Figure 3.2](#)). Select a search criterion from among the options. You can make only one selection at a time.

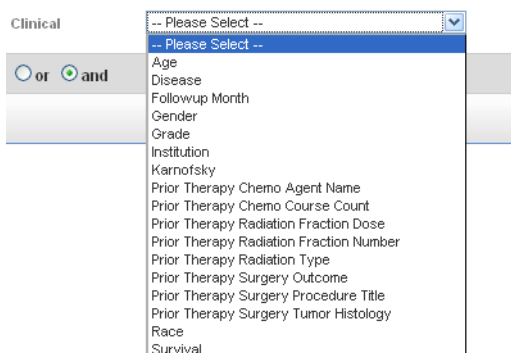


Figure 3.2 Additional clinical search criteria

- ° Each choice opens other fields relevant to the selection where you can further define your search query.
 - If permissible values were added when the annotation was defined, you must select among the values in a drop-list that displays on the right side of the page.
 - If no permissible values were defined as part of the annotation, you have the option to enter descriptive text in a text box on the right side of the page ([Figure 3.3](#)).

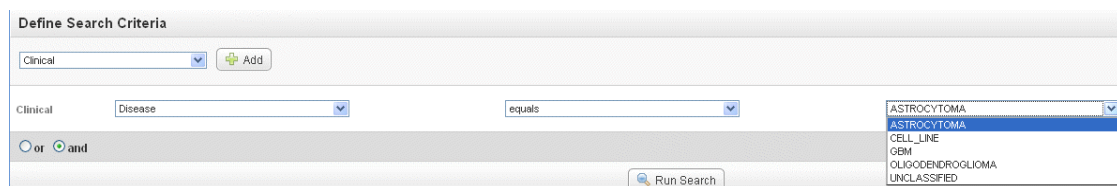


Figure 3.3 You may be able to further define search criteria when you select a specific clinical or imaging annotation element


Note: When working with image data, if only an Imaging Mapping file was uploaded when the study was created and not an Image Series Annotation file, you cannot enter image search criteria. The search results will, however, display a link that allows you to view the associated images in NBIA.

Continue with step 5 in [Searching a Study](#) on page 36.

Gene Expression

1. For the Gene Expression selection, select **Gene Name** or **Fold Change**.
2. **Gene Name** or **Fold Change** – Enter one or more gene symbols in the text box or click the icons to locate genes in the following databases. If you enter more than one gene in the text box, separate the entries by commas.

calIntegrator2 provides three methods whereby you can obtain gene names for a gene expression search:

- **caBio** – This link searches caBIO, then pulls identified genes into caIntegrator2 for analysis.
 - a. Click the **caBIO** icon ().
 - b. Enter **Search Terms**.
 - c. Select if you want to search in **Gene Keywords**, **Gene Symbols** or **Pathways** (from the drop-down list).
 - Selecting **Gene Keywords** searches only the Full Name field in caBio.
 - Selecting **Gene Symbols** searches only the Unigene and HUGO gene symbols in caBio.
 - Selecting **Pathways** searches only the pathway names in caBio. Note that searching in Pathways is a two step process. First, the initial Pathway search produces search results which are pathways. Second, from the pathway search results screen, you must select pathways of interest, then click **Search Pathways for Genes** to obtain a list of genes related to the selected pathways.
 - d. Select the **Any** or **All** choice to determine how your search terms will be matched. **Any** finds any match for any search term you entered. **All** finds only results that match all of the search terms.
 - e. Choose the **Taxon** from the drop-down list and click **Search**. The search results display (*Figure 3.4*).

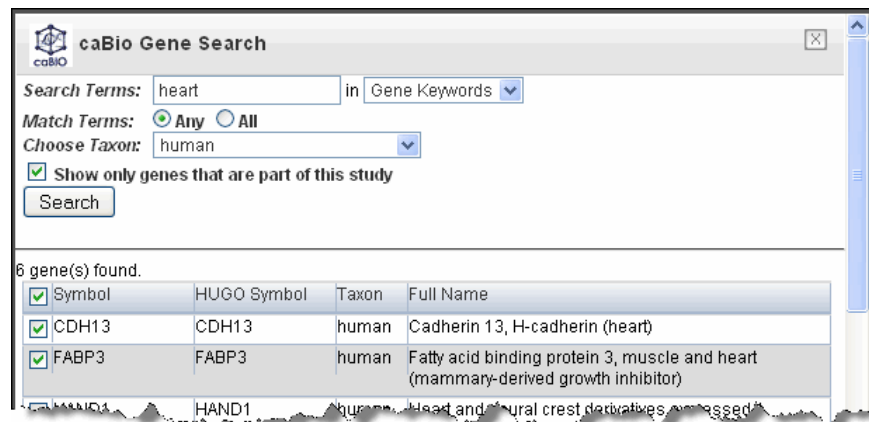



Figure 3.4 Example caBIO gene search results

- f. In the search results, use the check boxes to identify the genes whose symbols you want to use in the gene expression analysis.
- g. Click **Use Genes** at the bottom of the page. This pulls the checked genes into the Criteria tab (*Figure 3.5*).



Figure 3.5 Genes pulled in from caBIO display on the Criteria tab

- **Gene List** – This link locates gene lists saved in calIntegrator2.
 - a. Click the Genes List icon () to open a small dialog that lists prior-saved gene lists.
 - b. In the drop-down menu, select a gene list. In the list that appears, use the check boxes to identify the genes whose symbols you want to use in the gene expression analysis.
 - c. Click **Use Genes** at the bottom of the dialog. This pulls the checked genes into the Search Criteria tab.
- **CGAP** – Use this directory to identify genes. Before clicking this link you must enter gene symbols in the text box. This link does not pull anything into calIntegrator2 but does provide information about the gene(s) whose names you entered.

Additional fields display for the Fold Change selection.

The fold change option appears only if genomic control samples have been uploaded to the study. Fold change identifies genes with expression differences compared to control samples, as defined when the study was deployed in calIntegrator2. You can enter query values in greater/lesser-than-or-equal-to arguments.

3. Select or enter data for the Fold change fields shown in [Figure 3.6](#):

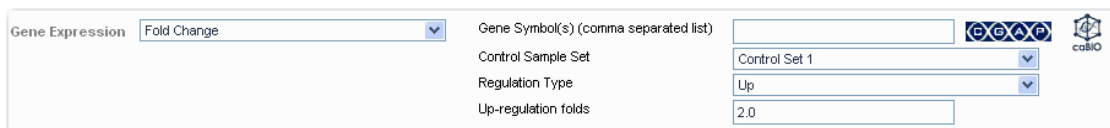


Figure 3.6 Fields for identifying fold change search criteria

- **Control Sample Set** – Select from the drop down list the name of the uploaded control sample set to serve as the fold change reference.
- **Regulation Type** – Select the term that describes the gene expression in comparison with the control samples: **Up** is increased expression; **Down** is decreased expression; **Up or Down** is increased or decreased; **Unchanged** means no change in expression.
- **Up-Regulation Folds** – The options here are dependent upon the Regulation Type you selected.
 - **Up** = Up Regulation Folds – Samples with a fold change greater than this value, when compared to the control samples, will be returned.
 - **Down** = Down Regulation Folds – Samples with a fold change less than this value, when compared to the control samples, will be returned.
 - **Up or Down** = Down Regulations Folds, Up Regulation Folds – Samples with a fold change either up or down, when compared to the control samples, will be returned.

- **Unchanged** = Samples with a fold change between the two specified values will be returned.

Continue with step 5 in [Searching a Study](#) on page 36.

Results Type Tab

You can specify columns for the way you want the search results to display either before or after you run the search. If you run the search directly from the Criteria tab before setting the results type/sorting features, by default only the Subject Identifiers display on the Search Results tab. You can then come back to the Results Type tab and [Sorting Tab](#) to expand the display options and re-run the search, having set the display parameters.

The selection you make on the Results Type tab determines whether caIntegrator2 displays search results for clinical or genomic data. It filters the search based on the criteria you set on the Criteria tab, whether it is clinical, gene expression or image series data type(s). In other words, if you select clinical criteria on the Criteria tab, but select Genomic on the Results Type tab, the data subset that displays on the Search Results tab is genomic data that is filtered by the clinical criteria you defined on the Criteria tab.

1. On the Results Type tab, select the **Clinical** or **Genomic** radio button to search clinical data ([Figure 3.7](#)).

Figure 3.7 Results Type tab

Clinical – Select the annotation elements that you want to display in the search results. All elements listed are column headers in the data uploaded to the study. You can make multiple selections on this list.

Note: For Clinical Annotations, the Patient or Subject Identifier display by default in the search results.

Results display as tabular data.

Genomic – Select the Reporter Type:

- **Gene Name** – Finds and summarizes at the gene level all reporters that match criteria for the gene you defined on the Criteria tab.
- **Reporter ID** – Finds all reporters that map to the gene(s) you identified on the Criteria tab

Results display in a gene expression data matrix.

Imaging – If imaging annotations have been added to the study, annotation elements also display on the lower right section of this page when you select **Clinical**. All elements listed are column headers in the image annotation data uploaded to the study. You can make multiple selections on this list.

Note: If you select even one Image Annotation on the Results Type tab, the Image Series IDs display by default in the search results. If you select no Image Annotations on the Results Type tab, however, even if you have selected image series criteria on the Criteria tab, no image series IDs display in the search results. The fact that images can be located, however, in NBIA is indicated by two image-related buttons at the bottom of the Query Results page. You can open the images in NBIA, but they will be at StudyInstance UID level. See [Relationship of Patient to Study to Series to Images](#) on page 54.

Results display as tabular data.

2. Use the **Select All** or **Unselect All** buttons to aid you in making your selections.

The column selection is saved as part of the query if you save it. See [Saving a Query](#) on page 43.

Sorting Tab

On the Sorting tab, you can set the sort order for data columns in the query results. You can also indicate whether column contents are sorted in ascending or descending order.

The columns that display on the Sorting tab are those criteria that you selected on the [Results Type Tab](#) for a Clinical Results type search.

1. Select the Sorting tab and indicate the left to right column order of the Search Results by changing one or more numbers in the Column Order column in this table (*Figure 3.8*).

Column	Column Order (L-R)	Row Order
Death Date	1	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Prior Therapy Radiation Fraction Number	2	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Prior Therapy Surgery Procedure Title	3	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Institution	4	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Gender	5	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Prior Therapy Surgery Tumor Histology	6	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Prior Therapy Radiation Type	7	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Last Followup Date	8	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Age	9	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Race	10	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Prior Therapy Surgery Outcome	11	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort
Karnofsky	12	<input type="radio"/> Ascending <input type="radio"/> Descending <input checked="" type="radio"/> No Sort

Figure 3.8 Sorting tab

2. In the **Row Order** column, indicate how you want columns sorted, **Ascending** or **Descending**, or leave the default, **No Sort**, if you choose.
3. Click **Run Query** at the bottom of the page to execute your sorting changes in the search results. When you do so, the change in column order is visible on the Query Results tab, as well as on the Sorting tab. For example, any column that you have indicated to be number 1 now appears in Query Results immediately after the Subject Identifier column and at the top of the Set Sort Order table on the Sorting tab.

Sorting parameters are saved as part of the query if you choose to save it using the Save Query feature. See [Saving a Query](#) on page 43.

4. If you click the **Reset** button before running the query from the Sorting tab, the original column settings are restored.

For information about the search results, see [Chapter 4 Viewing Query Results](#).

Managing Queries

When you create a search query in calIntegrator2, you can save the query for later use or edit it.

Exporting Query Results

Saving a Query


To save a query, follow these steps:

1. Click the **Save As** tab and enter a **Search Name** and **Search Description**, unique to the search. *Example: **Batch ID 6 and female***
2. Click **Save**.

Once the query is saved, it is listed by its name under the **Study Data > Queries > My Queries** in the left sidebar, whenever the study to which the query applies is selected. Click on the saved query in this list to either edit or re-run the query. Click on the query name to retrieve query results. If you hover over the Name text for the query, a popup displays the query description.

Editing a Query

To edit a query, follow these steps:

1. To edit a query, select it in the left sidebar under the **Study Data > Queries > My Queries**.
2. Click the **Edit** icon () corresponding to the study.
3. Change the query and display criteria on the Criteria, Columns and Sorting tabs.
4. On the Save As tab, check the appropriate options and click **Save As**. You can use the same name as the original query or modify the name as needed.

Exporting Query Results

After running a search, you can export the result set or a subset as a tab-delimited text file. For more information, see [Exporting Data](#) on page 55.

CHAPTER 4

VIEWING QUERY RESULTS

This chapter describes search results that calIntegrator2 returns after queries.

Topics in this chapter include the following:

- [Query Results Overview](#) on this page
- [Browsing Query Results](#) on page 46

Query Results Overview

After you launch a search of a calIntegrator2 study, the system automatically opens the Query Results tab showing the results of your search.

If you have not configured column and sort display parameters before launching the search, by default the tab shows only the subject identifiers and a column that allows you to select each row of the data subset.

To display and/or sort additional data, you must return to the Columns and/or Sorting tabs to set display parameters, then re-run the search. The new search results will display the additional information, with the columns and data sorted as you specified. See [Results Type Tab](#) on page 41.

calIntegrator2 paginates search results into pages of configurable size (default 20) with standard paginated navigation controls. To sort columns by ascending or descending parameters for on any displayed field, click on the underlined column header.

You can download search results as a CSV file. The file contains the annotations, columns and data sort configurations you specified in the search query. See [Exporting Query Results](#) on page 44.

Browsing Query Results

The query results that can display depend upon the criteria you established for the search. Follow the links below for more information about the category of data you searched.

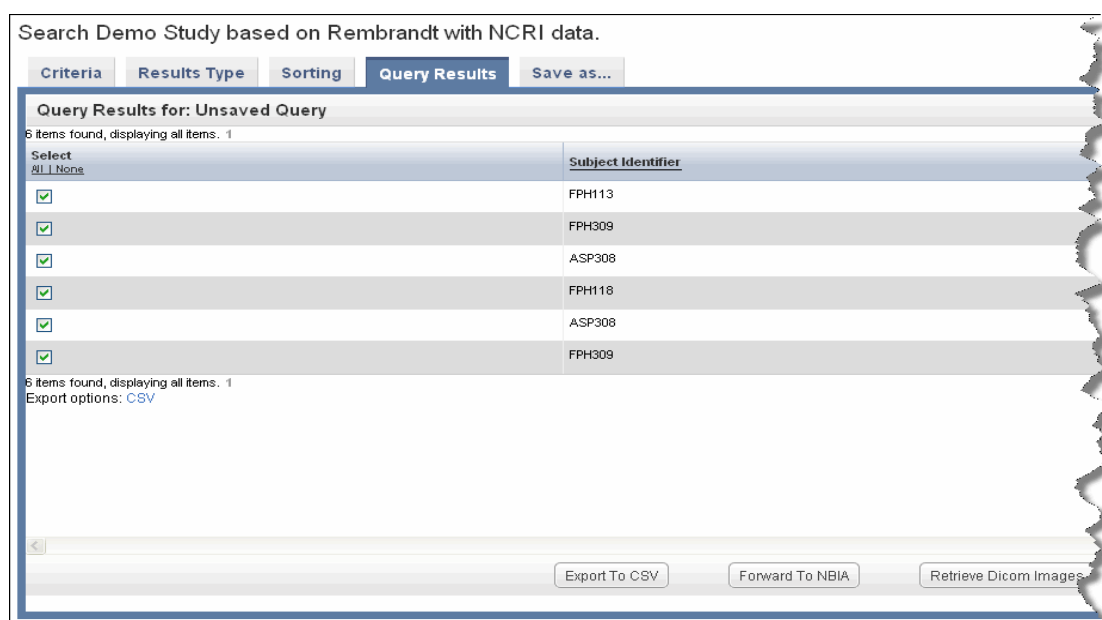
[Clinical and Imaging Data](#) on page 46

[Genomic Data](#) on page 46

[Expanding Imaging Data Results](#) on page 50

Clinical and Imaging Data

If you run the search before configuring column and sort display parameters, only the [subject] ID that meet the criteria and a column allowing you to select each row appear on the table ([Figure 4.1](#)). .



Search Demo Study based on Rembrandt with NCRI data.

Criteria Results Type Sorting **Query Results** Save as...

Query Results for: Unsaved Query

6 items found, displaying all items. 1

Select	Subject Identifier
<input checked="" type="checkbox"/>	FPH113
<input checked="" type="checkbox"/>	FPH309
<input checked="" type="checkbox"/>	ASP308
<input checked="" type="checkbox"/>	FPH118
<input checked="" type="checkbox"/>	ASP308
<input checked="" type="checkbox"/>	FPH309

6 items found, displaying all items. 1
Export options: [CSV](#)

Export To CSV Forward To NBIA Retrieve Dicom Images

Figure 4.1 Query Results page

You can add details for one or more subjects by configuring them on the Results Type tab. Annotations listed there are the column headers in the CSV file(s) that were uploaded to the study. For information about using the Results Type tab, see [Results Type Tab](#) on page 41.

Genomic Data

If you select the Genomic result type on the Results Type tab, genomic data search results display in a gene expression data matrix. Because the data was downloaded from caArray, the data permissions granted there still apply. In other words, if you have been given access to the data in caArray, you can see it in caIntegrator2.

In the matrix, samples in the experiment form the column headings. If the rows display Gene Name, the cells display the median gene expression value for each gene. If the

rows display Probe Set, the cells display the normalized signal-based value for a given reporter for a given sample (*Figure 4.2, Figure 4.3*).

Search Demo Study based on DC Lung Full data

Criteria

Results Type

Sorting

Query Results

Save as...

Query Results for: Unsaved Query

Subject ID	321	326	158	53	180	350	545	54	142	199	539	89	193	107	305	244	74	361
Sample ID	321	326	158	53	180	350	545	54	142	199	539	89	193	107	305	244	74	361
Gene Name																		
EGFR	151.48	395.43	179.63	70.39	91.81	175.26	97.61	80.96	111.84	115.25	309.05	81.36	136.61	84.75	145.52	208.45	66.68	14

Export options: CSV

Figure 4.2 Gene Name search EGFR, Reporter Type: Gene

Search Demo Study based on DC Lung Full data

Criteria

Results Type

Sorting

Query Results

Save as...

Query Results for: Unsaved Query

		Subject ID	21	362	369	332	480	135	93	345	255	486	333	115
		Sample ID	21	362	369	332	480	135	93	345	255	486	333	115
Gene Name	Reporter ID													
EGFR	201983_s_at		7071.93	957.15	1747.74	2614.96	875.01	646.42	503.99	4350.84	1297.73	2425.24	1275.78	365.28
EGFR	201984_s_at		2625.52	402.21	1225.3	1629.62	569.71	477.81	311.24	673.6	611.8	637.53	347.0	376.42
EGFR	210984_x_at		420.18	44.28	56.72	214.84	133.44	24.43	20.37	36.96	40.0	31.66	45.54	40.23
EGFR	211550_at		9.43	18.68	33.64	11.73	23.74	15.93	12.4	31.44	29.33	20.11	33.91	14.49
EGFR	211551_at		180.43	157.4	349.1	208.0	161.24	173.67	59.78	213.68	319.14	189.92	250.89	131.74
EGFR	211607_x_at		378.93	96.02	48.57	108.57	209.22	13.32	29.6	61.36	62.75	52.88	87.53	47.32

Export options: CSV

Figure 4.3 Gene Name search EGFR, Reporter Type: Reporter ID

- Genomic data does not display in tandem with clinical and imaging data; it only displays when you select the Genomic result type on the Results Type tab. Genomic data is however, filtered by clinical and imaging query criteria configured on the Criteria tab.
- Click the Export Options CSV link to download the CSV file whose data displays on the Search Results tab. When you do so, the CSV file opens automatically in MS Excel or similar applications for working with spreadsheets, showing the columns and sorting as you defined them in calIntegrator2 on the appropriate tabs.

Creating a Gene List


From any page in calIntegrator2, you can save a list of genes so you can use it for searches or analyses. To create a gene list, follow these steps:

1. Click the **Create New Gene List** link in the left sidebar. This opens the Manage Gene List page (*Figure 4.4*):

Figure 4.4 Manage Gene List page

2. Enter a name for the gene list.
3. Enter a description (optional).
4. For **Gene Symbol**, enter one or more gene symbols in the text box or click the icons to locate genes in the following databases. If you enter more than one gene in the text box, separate the entries by commas.

caIntegrator2 provides three methods whereby you can obtain gene names for creating a gene list:

- **caBio** – This link searches caBio, then pulls identified genes into caIntegrator2 for analysis.
 - a. Click the **caBio** icon ().
 - b. Enter **Search Terms**.
 - c. Select if you want to search in **Gene Keywords**, **Gene Symbols** or **Pathways** (from the drop-down list).
 - Selecting **Gene Keywords** searches only the Full Name field in caBio.
 - Selecting **Gene Symbols** searches only the Unigene and HUGO gene symbols in caBio.
 - Selecting **Pathways** searches only the pathway names in caBio. Note that searching in Pathways is a two step process. First, the initial Pathway search produces search results which are pathways. Second, from the pathway search results screen, you must select pathways of interest, then click **Search Pathways for Genes** to obtain a list of genes related to the selected pathways.

- d. Select the **Any** or **All** choice to determine how your search terms will be matched. **Any** finds any match for any search term you entered. **All** finds only results that match all of the search terms.
- e. Choose the **Taxon** from the drop-down list and click **Search**. The search results display in the same dialog box (*Figure 4.5*).

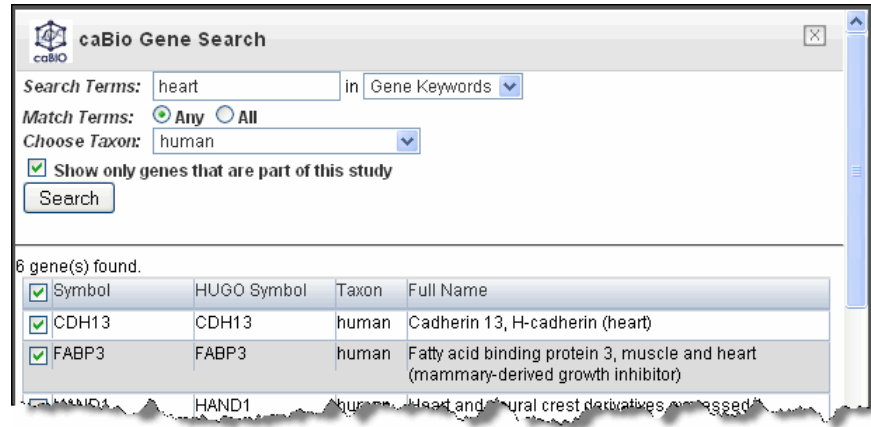



Figure 4.5 Example caBIO gene search criteria and results

- f. In the search results, use the check boxes to identify the genes whose symbols you want to include in the gene list.
 - g. Click **Use Genes** at the bottom of the page. This pulls the checked genes into the Gene Symbols field on the Gene List page.
- **Gene List** – This link locates genes lists saved in calIntegrator2.
 - a. Click the Gene List icon () to open a small dialog that lists prior-saved gene lists in calIntegrator2.
 - b. In the drop-down menu, select a gene list. In the list that appears, use the check boxes to identify the genes whose symbols you want to use in the plot analysis.
 - c. Click **Use Genes** at the bottom of the dialog. This pulls the checked genes into the Gene Symbols field on the Gene List page.
 - **CGAP** – Use this directory to identify genes. Before clicking this link you must enter gene symbols in the text box. This link does not pull anything into calIntegrator2 but does provide information about the gene(s) whose names you enter.
5. If you so choose, you can upload a gene list. For the Upload File field, click the **Browse** button to navigate to a .csv file made up of gene symbols. calIntegrator2 converts the comma-separated content to a gene list.

- Click **Create Gene List** at the bottom of the page. caIntegrator2 now opens the Edit Gene List page which shows the name and symbols of the newest gene list ([Figure 4.6](#)).

Figure 4.6 The Edit Gene List for reviewing, editing the name or deleting a gene list.

Editing a Gene List

To view a gene list in caIntegrator2, under Study Data in the left sidebar, click **Saved Lists > My Gene Lists**. The system displays gene lists that have been saved for the open study. Click on any of the list names to open an Edit Gene List dialog box.

- On this page you can review list metadata including the genes symbols included in the list.
- To rename the list in the **GeneList Name** text box, click the **Rename** button.
- To delete the study by clicking the **Delete** button.

See also [Creating a Gene List](#).

Expanding Imaging Data Results

In reviewing imaging search results, it is important to understand the hierarchy of submissions in NBIA. For more information, see [Relationship of Patient to Study to Series to Images](#) on page 54.

If you run the search before configuring column and sort display parameters, only the Subject Identifiers for the patients/images that meet the criteria and a column containing one check box per row display by default (*Figure 4.7*).

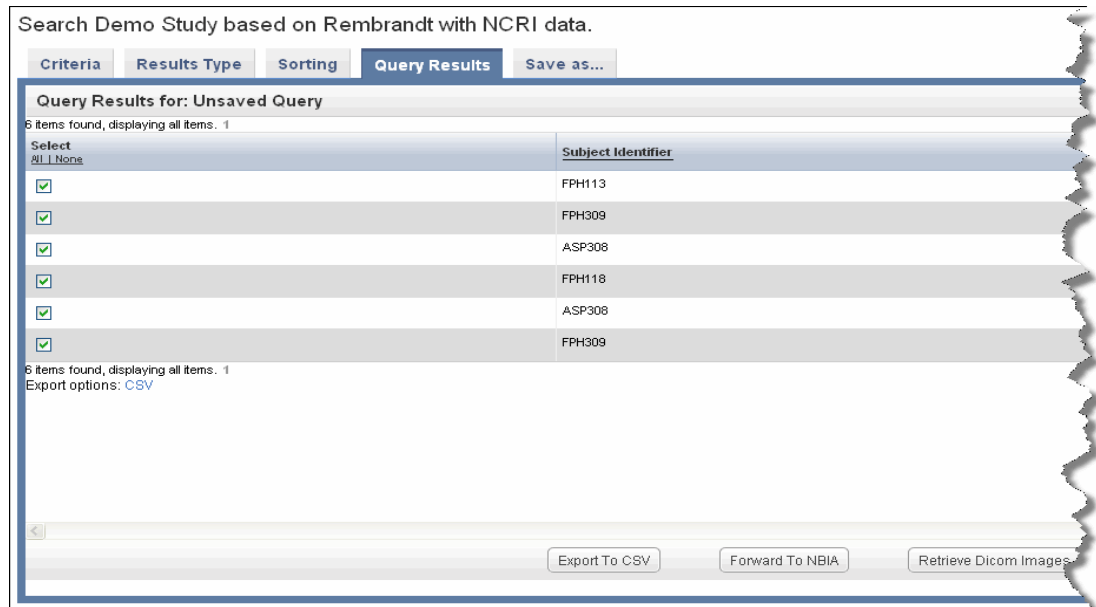


Figure 4.7 With imaging criteria only and no column definition, only Subject IDs display

If your annotation choice on the Columns page identifies annotations such as tumor size or tumor location, the search results display image series subsets that have those annotations. The check boxes work in conjunction with buttons at the bottom of the results page (*Figure 4.8*).

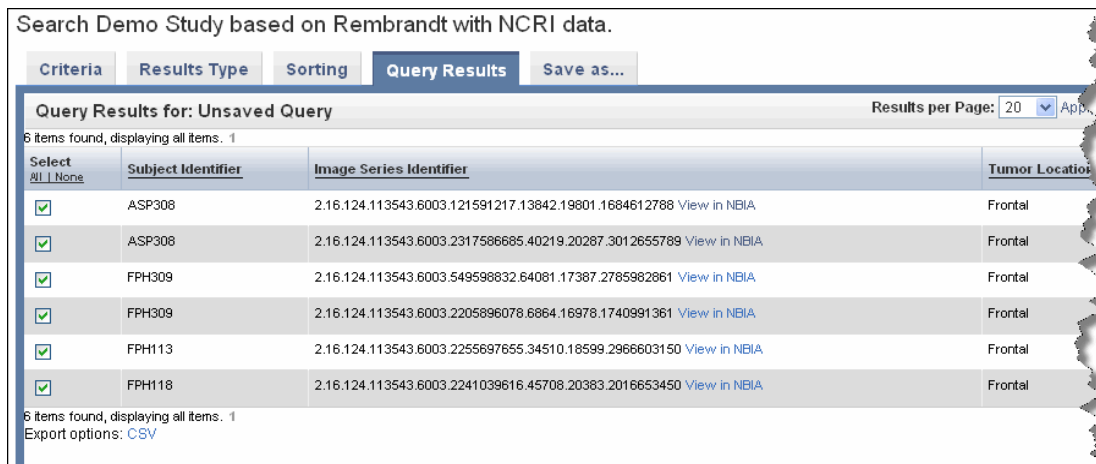


Figure 4.8 By expanding display parameters, you can view complete details for image search results

You can add more details for images by configuring image annotations on the Results Type tab. Annotations listed there are the column headers in the image series CSV file(s) that were uploaded to the study. Examples of image details include the following:

- All image details (name, size, etc.)

- The series that the image belongs to
- Image feature attributes
- The subject ID. Click the subject ID under Clinical Annotations on the Results Type tab to display this.

You can set display parameters for the results on the Columns and Sorting tabs. For more information, see [Results Type Tab](#) on page 41.

See also [caIntegrator2 and NBIA](#), [Retrieving Dicom Images](#) and [Example of Retrieving Images](#):

caIntegrator2 and NBIA

Images can be accessed in NBIA if you see buttons on the Search Results page. See the Imaging Note in [Results Type Tab](#) on page 41. You can click links on the Search Results tab to view or download image data.

- **View in NBIA** – This link corresponds to each Image Series listed in the results table. If you click the link, NBIA opens to the login page. After you log in, NBIA brings up the first image in the corresponding image series ([Figure 4.9](#)). You must log into NBIA to see the data. On the NBIA page that opens, you can opt to view the entire series containing this image, or you can display the image as a large JPEG-formatted image. You can also add the image to the NBIA basket. For more information, see the NBIA online help or user's guide accessible from NBIA.



Figure 4.9 An example of displaying the first image in image series

- **Forward to NBIA** – This button is linked to results you have selected by row. Click the button to open NBIA, where the image series you select are loaded in the NBIA image basket. In the event that the caIntegrator2 study was NOT configured with image annotation for an image series, caIntegrator2 sends NBIA a list of Study Instance UIDs, for which NBIA will add all corresponding image series to the basket. In the event that the caIntegrator2 study was configured with annotations for an image series, the system sends NBIA a list of

Image Series IDs, for which NBIA adds all corresponding image series to the basket.

Retrieving Dicom Images

On the Imaging data Search Results page, you can click the **Retrieve DICOM Images** button which is linked to results you have selected by row. caIntegrator2 retrieves the corresponding image(s) from NBIA through the grid. NBIA organizes the download file by patientID, StudyInstanceUID, and ImageSeriesUID, and compresses it into a zip file. When caIntegrator2 notifies you that the file is retrieved, the DICOM Retrieval page indicates whether the retrieved files are Study Instance UIDs or Image Series UIDs (Figure 4.10). For more information, see the note below.



Figure 4.10 DICOM Retrieval result

Click the **Download DICOM** link to download and save the file. caIntegrator2 unzips the file and displays the list of images in the file. To open the DICOM images, you must have a DICOM image viewer application installed on your computer. For more information, see <http://dicom.online.fr/fr/download.htm>.

In the search results, not all of the patients in the data subset may be mapped to image series IDs. If you select a mixture of patients that have image annotations as indicated by an image series ID and patients that do not have image annotations (no image series ID), when you click the **Retrieve DICOM Images** button, NBIA retrieves the images for the entire *NBIA study instance UID* that includes the image seriesIDs you checked.

If on the Search Results tab you select only patients that have image annotations as indicated by an image series ID, when you click the **Retrieve DICOM Images** button, NBIA retrieves images for the *NBIA image series* that were matched in the search. If the results are a mixture, but you select one specific row with a valid image annotation, caIntegrator2 aggregates to the *images series*. If results are a mixture and you select multiple rows, caIntegrator2 aggregates to the NBIA study in which multiple image series you have selected in the search results are found.

If your query does not have image annotations and all check boxes are selected, results will go up to image series UID and gives all image series in it. Search results may ultimately depend on how the study was created. For example, if no image series display in query results, it means they were not mapped. In the study. In that case, the results “move” up to Study Instance UIDs.

To best understand this, it is important to review the hierarchy of submissions in NBIA. For more information, see [Relationship of Patient to Study to Series to Images](#) on page 54.

Example of Retrieving Images:

You are searching a study that has image data and image annotation(s) for at least one image series.

1. Open a study that has imaging data associated with it that points to the production NBIA server.
2. Make a query that will have image series or patients who are associated to Image Studies and select a few of those patients in the check box.
3. Click the **Retrieve Dicom Images** button.

Note that it aggregates to the image study.

4. Now go back to Results Type tab, select all image annotations and run the query again.
5. Select an image series type column and click the **Retrieve Dicom Images** button.

calIntegrator2 now aggregates to the Image Series that were selected and not the Image Study.

6. Select a row that doesn't have image series data, and a row that does, and push the button.

This should aggregate to the study for the rows selected.

7. Click **Forward to NBIA**. You should see the same types of aggregation for these tests.

When the image Study is in the checked boxes (regardless of image series being there or not), the system aggregates up to the Image Study level.

Relationship of Patient to Study to Series to Images

This flowchart illustrates the relationship of patient to study to series and lastly to images.

Clinical trial > Patient (Subject) > Study > Series > Images

For example, the Study Instance UID is the set of images resulting from one patient office visit. When you upload a spreadsheet of an image series, the hierarchy of images in an image series might look like this:

Study Instance UID (one office visit):

Brain (image series)

- Brain image 1
- Brain image 2

- Brain image 3

Leg (image series)

- Leg image 1
- Leg image 2
- Leg image 3

You can add details for images by configuring image annotations on the Results Type tab. Annotations listed there are the column headers in the image series CSV file(s) that were uploaded to the study. Examples of image details include the following:

- All image details (name, size, etc.)
- The series that the image belongs to
- Image feature attributes
- The subject ID. Click the subject ID under Clinical Annotations on the Results Type tab to display this.

Exporting Data

You can choose to download tabular search results as a CSV file. Click the **Export .csv** link at the bottom of the page. You may need to scroll the page to see it. The file contains the annotations, columns and data sort configurations you specified in the search query.

Note: You will not see the Export option when genomic data displays as query results.

CHAPTER 5

ANALYZING STUDIES

This chapter describes how to use calIntegrator2 tools to analyze data in clinical or genomic studies that have been deployed in calIntegrator2.

Topics in this chapter include the following:

- [Data Analysis Overview](#) on this page
- [Creating Kaplan-Meier Plots](#) on page 58
- [Creating Gene Expression Plots](#) on page 65
- [Analyzing Data with GenePattern](#) on page 78

Data Analysis Overview

Once a study has been deployed, you can analyze the data using calIntegrator2 analysis tools.

You can verify that the study is in “Deployed” status by selecting the study name in the My Studies dropdown selector. After selecting the study name, click **Home** in the left sidebar of the calIntegrator2 Menu. A study summary should appear, including a status field. If the status is not deployed, or if the study summary does not appear, then the study is not deployed and available for analysis.

If the study is ready for analysis, you will see an **Analysis Tools** menu in the left sidebar with the following options:

- **K-M Plot:** This tool analyzes clinical data, generating a Kaplan-Meier (K-M) plot based on survival data sets. See [Creating Kaplan-Meier Plots](#) on page 58.
- **Gene Expression Plot:** This tool analyzes annotation, clinical or genomic data based on gene expression values. See [Creating Gene Expression Plots](#) on page 65.

- **GenePattern:** This feature provides an express link to GenePattern where you can perform analyses on selected calIntegrator2 studies, or it enables you to perform several GenePattern analyses on the grid. See [Analyzing Data with GenePattern](#) on page 78 .

After defining or running the analysis on selected data sets, analysis results display on the same page, allowing you to review the analysis method parameters you defined.

Creating Kaplan-Meier Plots

The Kaplan-Maier method analyzes comparative groups of patients or samples. In calIntegrator2, the K-M method compares survival statistics among comparative groups. You can configure the survival data in the application. For example, you might identify a group of patients with smoking history and compare survival rates with a group of non-smoking patients, or compare the survival data for two groups of patients with a specific disease type and based on Karnofsky scores . You could compare groups of patients with varying gene expression levels. You can also identify data sets using the query feature in the application, saving the queries, then configuring the K-M to compare groups identified by the queries.

The key is to first identify subsets of patients or samples that meet criteria you want to establish, thus filtering the data you want to compare. Next, generate a K-M plot based on their survival probability as a function of time. Survival differences are analyzed by the log-rank test.

Note: To perform a K-M plot analysis, survival data must have been identified for the study you want to analyze. For more information, see [Defining Survival Values](#) on page 23.

K-M Plot for Annotations

The groups identified for this K-M plot generation are based on clinical annotations.

1. Select the study whose data you want to analyze in the upper right portion of the calIntegrator2 page.
2. Under Analysis Tools on the left sidebar, select **K-M Plot**.
3. Select the **For Annotation** tab at the top of the page ([Figure 5.1](#)).

Kaplan-Meier Survival Plots (draft)

For Annotation For Gene Expression For Queries

Annotation Based Kaplan-Meier Survival Plots

	Annotation Type	Annotation	Values
1.) Patient Groups:	Select Annotation Type	Select Annotation	
Survival Value			
2.) Select Survival Measure:	Survival from enrollment		
Reset			

Figure 5.1 Fields for defining annotation data for a K-M plot

4. The groups to be compared in the K-M plot originate from one patient group. Varying data sets are based upon multiple values corresponding to the selected annotation. Define Patient Groups using these options:
 - **Annotation Type** – Select the annotation type that identifies the patient group. Selections are based on the data in the chosen study.
 - **Annotation** – Select an annotation. Fields are based on the annotation type you select. For example, if you choose **Subject**, then you could select **Gender** or **Radiation Type** or any field that would distinguish the patients into groups based upon their values.
 - **Values** – Using conventional selection techniques, select two or more values which will be the basis for the K-M plot. Permissible (available) values or “No Values” correspond to the selected annotation.
5. **Survival value** is the length of time the patient lived. For **Survival Value**, select the survival measure which is the unit of measurement for the survival value to be used for the plot.
6. Click the **Create Plot** button.

Note: The Create Plot button displays only after you have selected appropriate criteria.

calIntegrator2 generates the plot which then displays below the plot criteria ([Figure 5.2](#)).

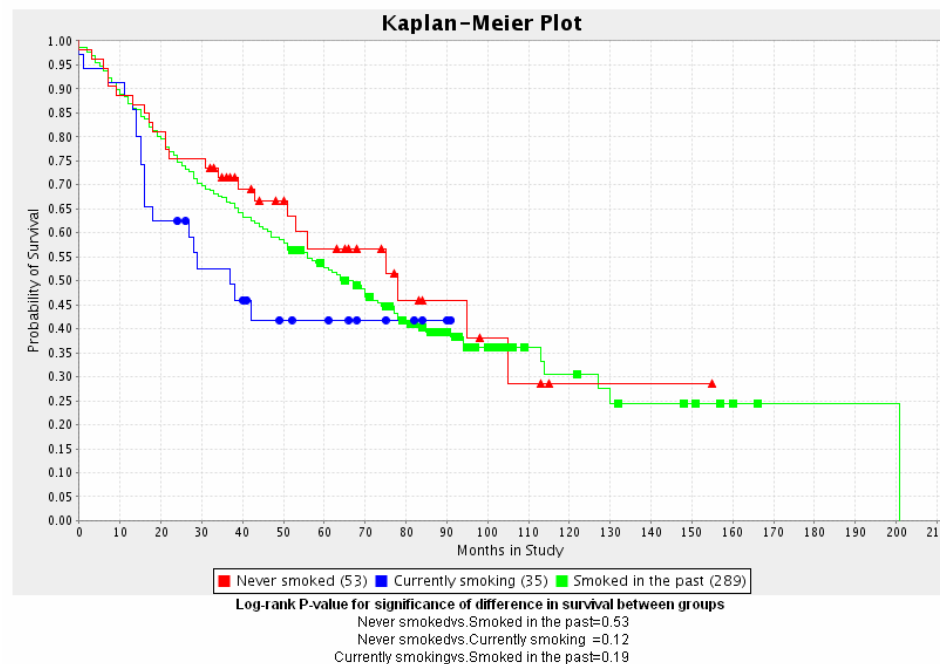


Figure 5.2 A K-M plot generated for groups based on clinical annotations

- The number of subjects for each group appears embedded in the legend of the graph below the plot.
- calIntegrator2 generates a P-value for the selected groups; it displays at the bottom of the page. A low P-value generally has more significance than a high P-value.

K-M Plot for Gene Expression

calIntegrator2 allows you to compare expression levels for one given gene at a time. The relative expression level is referred to as “fold change” and the numeric value for a given sample and reporter combination is the ratio of the expression value for that particular reporter for the given sample to a reference value calculated for that reporter across all control samples. The reference value is calculated by taking the mean of the \log_2 of the expression values for all control samples for the reporter in question. The \log_2 mean value (n) is then converted back to a comparable expression signal by returning 2 to the exponent n .


To create a K-M plot illustrating gene expression values, follow these steps:

1. Select the study whose data you want to analyze in the upper right portion of the calIntegrator page. You must select a study with gene expression data.
2. Under Analysis Tools on the left sidebar, select **K-M Plot**.
3. Select the **For Gene Expression** tab ([Figure 5.3](#)).

Figure 5.3 Fields for defining gene expression data for a K-M plot

4. For **Gene Symbol**, enter one or more gene symbols in the text box or click the icons to locate genes in the following databases. If you enter more than one gene in the text box, separate the entries by commas.

calIntegrator2 provides three methods whereby you can obtain gene names for calculating a KM plot:

- **caBio** – This link searches caBIO, then pulls identified genes into calIntegrator2 for analysis.
 - a. Click the **caBIO** icon ().

- b. Enter **Search Terms**.
- c. Select if you want to search in **Gene Keywords**, **Gene Symbols** or **Pathways** (from the drop-down list).
 - Selecting **Gene Keywords** searches only the Full Name field in caBio.
 - Selecting **Gene Symbols** searches only the Unigene and HUGO gene symbols in caBio.
 - Selecting **Pathways** searches only the pathway names in caBio. Note that searching in Pathways is a two step process. First, the initial Pathway search produces search results which are pathways. Second, from the pathway search results screen, you must select pathways of interest, then click **Search Pathways for Genes** to obtain a list of genes related to the selected pathways.
- d. Select the **Any** or **All** choice to determine how your search terms will be matched. **Any** finds any match for any search term you entered. **All** finds only results that match all of the search terms.
- e. Choose the **Taxon** from the drop-down list and click **Search**. The search results display in the same dialog box (*Figure 5.4*).

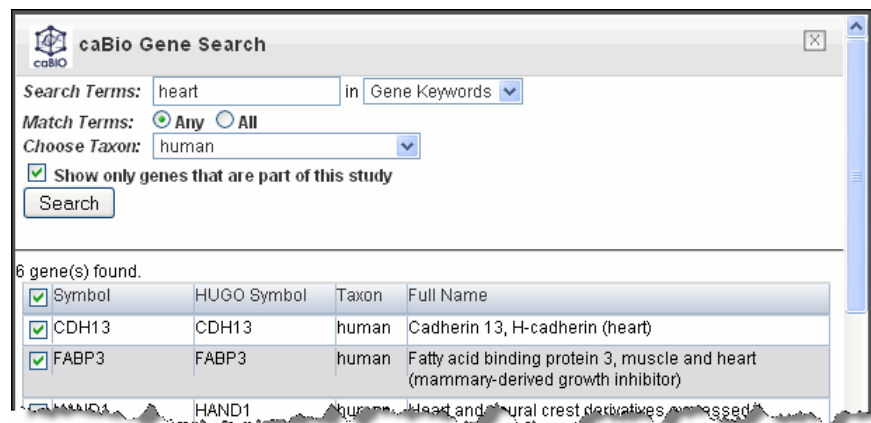



Figure 5.4 Example caBIO gene search criteria and results

- f. In the search results, use the check boxes to identify the genes whose symbols you want to use in the plot calculation.
- g. Click **Use Genes** at the bottom of the page. This pulls the checked genes into the For Gene Expression tab (*Figure 5.5*).



Figure 5.5 Genes pulled in from caBIO display on the selected tab

- **Gene List** – This link locates gene lists saved in caIntegrator2.

- a. Click the Gene List icon () to open a small dialog that lists prior-saved gene lists in calIntegrator2. See [Creating a Gene List](#) on page 47.
 - b. In the drop-down menu, select a gene list. In the list that appears, use the check boxes to identify the genes whose symbols you want to use in the plot analysis.
 - c. Click **Use Genes** at the bottom of the dialog. This pulls the checked genes into the For Gene Expression tab.
- **CGAP** – Use this directory to identify genes. Before clicking this link you must enter gene symbols in the text box. This link does not pull anything into calIntegrator2 but does provide information about the gene(s) whose names you enter.
5. **Over-expressed/Under-expressed** – Define the over- and under-expression criteria, expressed in terms of fold-change. Fold change is the ratio of the measured gene expression value for an experimental sample to the expression value for the control sample.
 6. **Survival value** – The length of time the patient lived. For **Survival Value**, select the survival measure which is the unit of measurement for the survival value to be used for the plot.
 7. **Control Sample Sets** – One or more are created by the study manager when a study is deployed. Select the **Control Sample Set** you would like to use to calculate fold-change.
 8. Click the **Create Plot** button. calIntegrator2 generates the plot which then displays below the plot criteria ([Figure 5.6](#)).

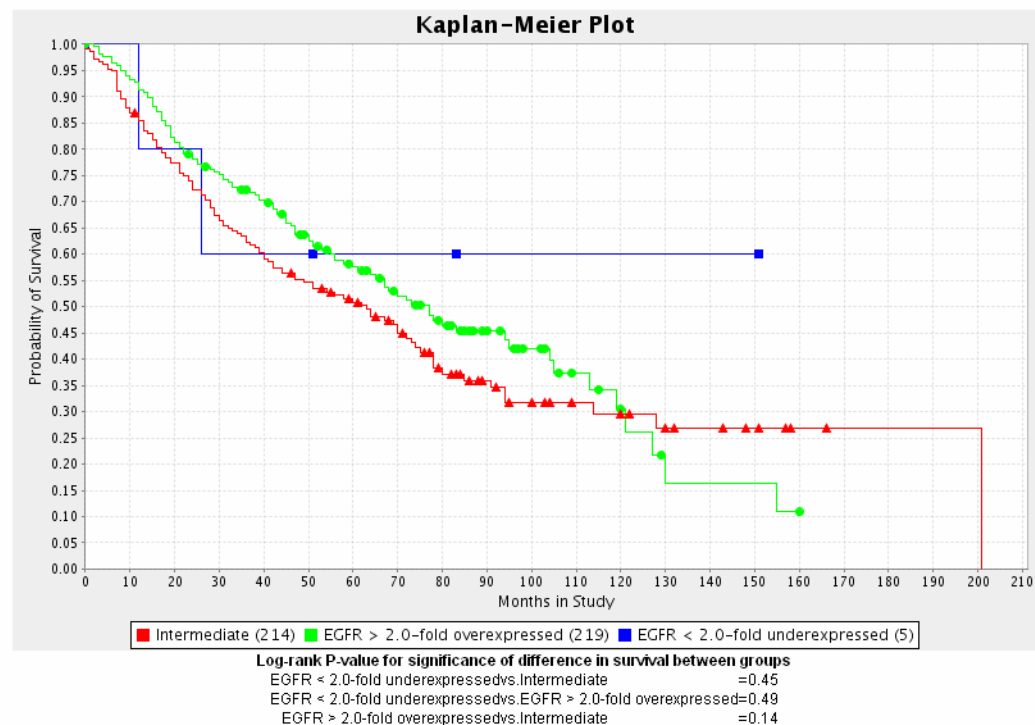


Figure 5.6 K-M plot generated from gene expression data.

- The number of subjects for each group appears embedded in the legend of the graph below the plot. Note the appearance of an intermediate group (red entries), which is a group with gene expression values that are not up-regulated nor down-regulated.
- In queries that include a fold change criterion and that are configured to return genomic data, raw expression values are replaced with calculated fold change values.
- A P-value is also generated for the selected groups; it displays at the bottom of the page. A low P-value generally has more significance than a high P-value.

K-M Plot for Queries

You can identify data sets using the query feature in the application. You can manipulate the queries to find the groups you want to compare, save the queries, then configure the K-M to compare the query groups. This is one method of limiting the data considered in the K-M plot calculation.

1. Select the study whose data you want to analyze in the upper right portion of the calIntegrator page. You must select a study for which the queries you will identify for the K-M plot have been saved.
2. Under Analysis Tools on the left sidebar, select **K-M Plot**.

3. Select the **For Queries** tab ([Figure 5.7](#)).

Kaplan-Meier Survival Plots (draft)

For Annotation For Gene Expression **For Queries**

Query Based Kaplan-Meier Survival Plots

1.) Select Queries:

All Available Queries

- gender female
- equal to or > 60
- equal to or > 70
- Never smoke
- equal to or > 40
- equal to or > 30

Add >

< Remove

Selected Queries

↓

↑

2.) ☐ Exclusive Subjects in Queries (Subjects in upper queries are removed from subsequent queries)

3.) ☐ Add additional group containing all other subjects not found in selected queries.

4.) Select Survival Value: Survival from enrollment

Reset Create Plot

Figure 5.7 Fields for defining K-M plot parameters based on saved queries in caIntegrator2

4. **Queries** – Select **Queries** whose data you want to analyze from the **All Available Queries** panel and move them to the **Selected Queries** panel using the **Add >>** button.

Note: Genomic queries do not appear in the lists; they cannot be selected for this type of K-M plot.

5. **Exclusive Subject in Queries** – Check the box if you want to exclude any subjects that appear in both (or all) queries selected for the plot, thus eliminating overlap.
6. **Add Additional Group...all other subjects** – Check the box to create an additional group of all other subjects that are not in selected query groups.
7. **Survival value** – The length of time the patient lived. Select the survival measure which is the unit of measurement for the survival value to be used for the plot.
8. Click the **Create Plot** button. caIntegrator2 generates the plot which then displays below the plot criteria ([Figure 5.8](#)).

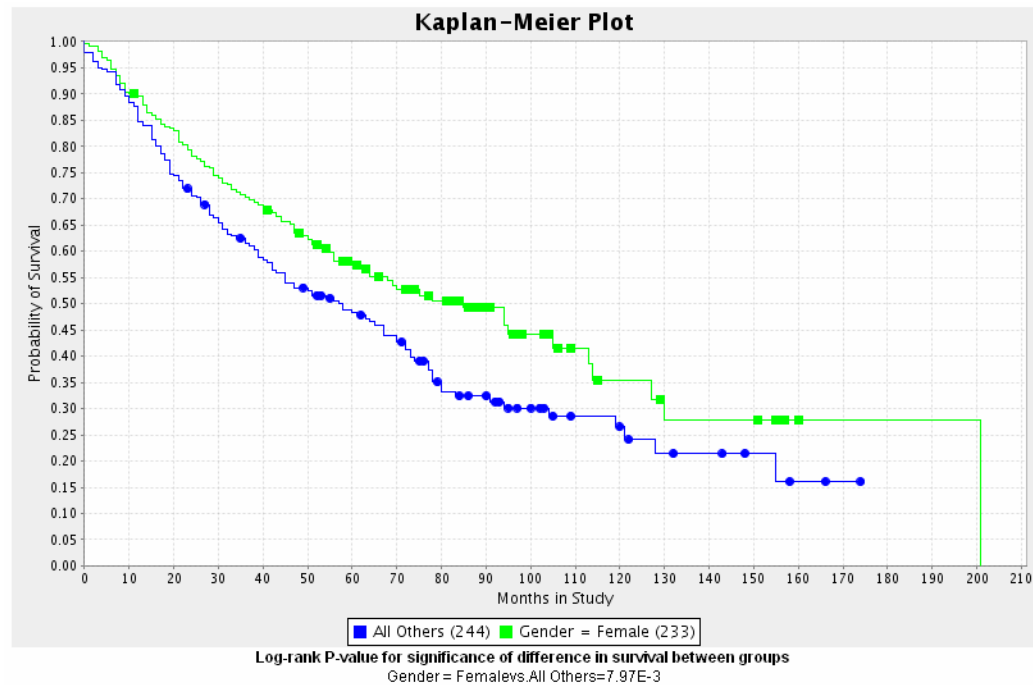


Figure 5.8 K-M Plot comparing statistics between subjects in two queries

- The number of subjects for each group appears embedded in the legend of the graph below the plot.
- A P-value is also generated for the selected groups; it displays at the bottom of the page. A low P-value generally has more significance than a high P-value.

Creating Gene Expression Plots

Gene expression plots compare signal values from reporters or genes. This statistical tool allows you to compare values for multiple genes at a time, but it does not require only two sets of data to be compared. It also allows you to compare expression levels for selected genes against expression levels for a set of control samples designated at the time of study definition.

caIntegrator2 provides three ways to generate meaningful gene expression plots, indicated by tabs on the page. The tabs are independent of each other and allow you to select the genes, reporters and sample groups to be analyzed on the plot.

- [Gene Expression Value Plot for Annotation](#) – You can locate genes in the caBio directories or caIntegrator2 Gene Lists. You can learn more about the genes in the CGAP directory. You can define criteria for the plot using clinical and image annotations.
- [Gene Expression Value Plot for Genomic Queries](#) – You can select data based on saved genomic queries.

- **Gene Expression Value Plot for Clinical Queries** – You can select data based on saved clinical queries. You can locate genes in the caBio directories or calIntegrator2 Gene Lists.

See also [Understanding a Gene Expression Plot](#) on page 75.

Gene Expression Value Plot for Annotation

To generate a gene expression plot, follow these steps:

1. Select the study whose data you want to analyze in the upper right portion of the calIntegrator2 page. (You must select a study which has genomic data.)
2. Under Analysis Tools on the left sidebar, select **Gene Expression Plot**. This opens a page with three tabs
3. Select the **For Annotation** tab ([Figure 5.9](#)).

The screenshot shows the 'Gene Expression Value Plots' interface with the 'For Annotation' tab selected. The interface includes a header with three tabs: 'For Annotation', 'For Genomic Queries', and 'For Clinical Queries'. Below the tabs is a section titled 'Annotation Based Gene Expression Plots'. It contains five numbered steps for configuration:


- 1.) Gene Symbol(s) (comma separated list): A text box containing 'DLEC1,PLUNC,KLF2,MYCL' and icons for caBio and CEXAP.
- 2.) Select Reporter Type: Radio buttons for 'Reporter Id' (selected) and 'Gene'.
- 3.) Sample Groups: A table with columns 'Annotation Type', 'Annotation', and 'Values'. The 'Annotation Type' column has a dropdown menu with 'Select Annotation Type'. The 'Annotation' column has a dropdown menu with 'Select Annotation'. The 'Values' column is empty.
- 4.) ☐ Add additional group containing all other subjects not found in selected queries.
- 5.) ☐ Add additional group containing all control samples for this study. A dropdown menu with 'Control Set 1'.

A 'Reset' button is located at the bottom right of the form.

Figure 5.9 Gene expression value tab for configuring gene expression annotation value plot

4. **Gene Symbol** – Enter one or more gene symbols in the text box or click the icons to locate genes in the following databases. If you enter more than one gene in the text box, separate the entries by commas.

calIntegrator2 provides three methods whereby you can obtain gene names for calculating a gene expression plot:

- **caBio** – This link searches caBio, then pulls identified genes into calIntegrator2 for analysis.
 - a. Click the **caBio** icon ().
 - b. Enter **Search Terms**.
 - c. Select if you want to search in **Gene Keywords**, **Gene Symbols** or **Pathways** (from the drop-down list).
 - Selecting **Gene Keywords** searches only the Full Name field in caBio.

- Selecting **Gene Symbols** searches only the Unigene and HUGO gene symbols in caBio.
- Selecting **Pathways** searches only the pathway names in caBio. Note that searching in Pathways is a two step process. First, the initial Pathway search produces search results which are pathways. Second, from the pathway search results screen, you must select pathways of interest, then click **Search Pathways for Genes** to obtain a list of genes related to the selected pathways.
- d. Select the **Any** or **All** choice to determine how your search terms will be matched. **Any** finds any match for any search term you entered. **All** finds only results that match all of the search terms.
- e. Choose the **Taxon** from the drop-down list and click **Search**. The search results display in the same dialog box (*Figure 5.4*).

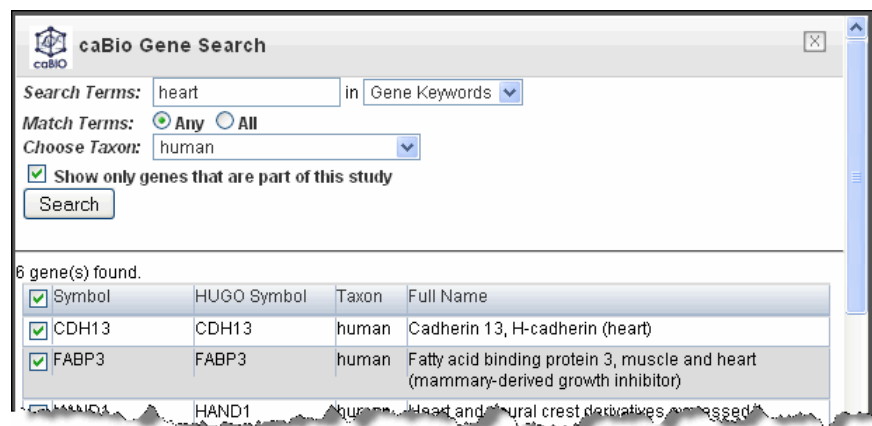



Figure 5.10 Example caBIO gene search results

- f. In the search results, use the check boxes to identify the genes whose symbols you want to use in the gene expression plot calculation.
- g. Click **Use Genes** at the bottom of the page. This pulls the checked genes into the For Annotation tab (*Figure 5.5*).



Figure 5.11 Genes pulled in from caBIO display on the tab

- **Gene List** – This link locates gene lists saved in caIntegrator2.
 - a. Click the Gene List icon () to open a small dialog that lists prior-saved gene lists in caIntegrator2.
 - b. In the drop-down menu, select a gene list. In the list that appears, use the check boxes to identify the genes whose symbols you want to use in the plot analysis.

- c. Click **Use Genes** at the bottom of the dialog. This pulls the checked genes into the For Annotation tab.
- **CGAP** – Use this directory to identify genes. Before clicking this link you must enter gene symbols in the text box. This link does not pull anything into calIntegrator2 but does provide information about the gene(s) whose names you enter.
5. **Reporter Type** – Select the radio button that describes the reporter type:
 - **Reporter ID** – Summarizes expression levels for all reporters you specify.
 - **Gene Name** – Summarizes expression levels at the gene level.
6. **Sample Groups** – Choose among the following options:
 - **Annotation Type** – Select the annotation type. Selections are based on the data in the chosen study
 - **Annotation** – Select an annotation. Fields are based on the annotation type you select. For example, if you choose Subject, then you could select Gender or Radiation Type or any field that would distinguish the patients into groups based upon study values.
 - **Values** – Using conventional selection techniques, select one or more values which will be the basis for the plot. Permissible (available) values or “No Values” correspond to the selected annotation.
7. **Add Additional Group...** – Define as follows:
 - **...all other subjects** – Check the box to create an additional group of all other subjects that are not in selected query groups.
 - **...control group** – Check the box to display an additional group of control samples for this study.
8. Click the **Create Plot** button. calIntegrator2 generates the plot which then displays below the plot criteria in bar graph format ([Figure 5.6](#)).

Legends below the plot indicate the plot input. By default, the plot shows the mean of the data. [Figure 5.12](#) displays a plot with gene expression median calculation summaries.

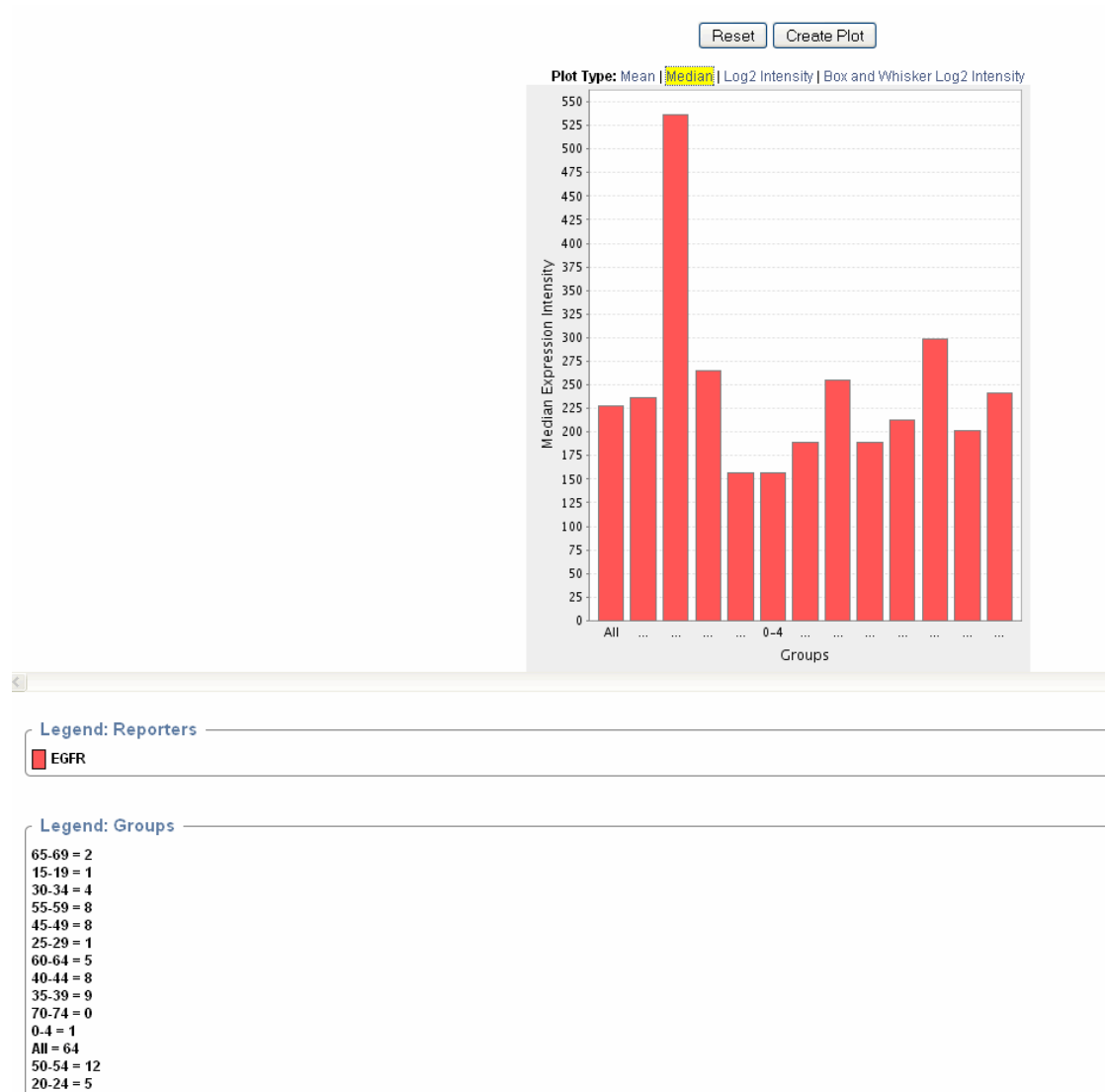


Figure 5.12 Gene expression plot based on selected annotations

- You can recalculate the data display by clicking the **Plot Type** above the graph. See [Understanding a Gene Expression Plot](#) on page 75.
- You can modify the plot parameters and click the **Reset** button to recalculate the plot.

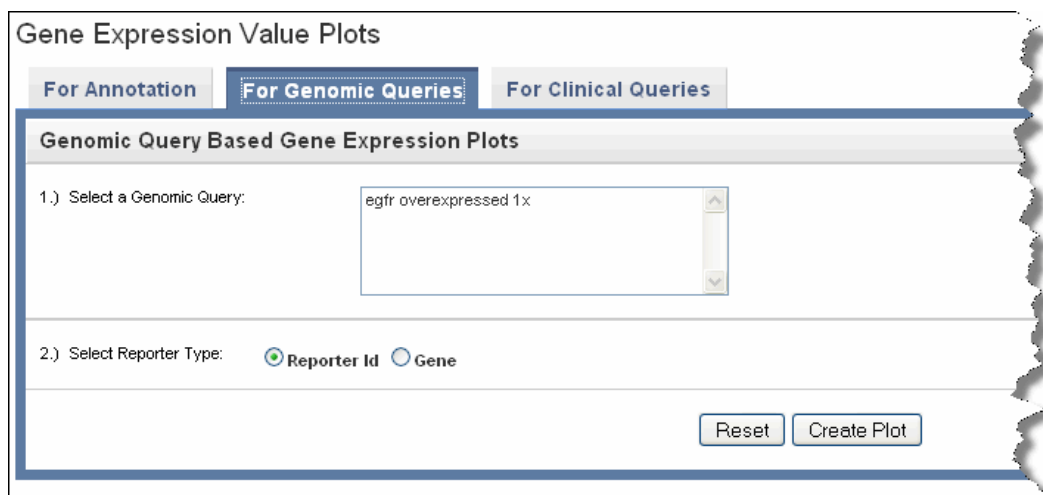
Gene Expression Value Plot for Genomic Queries

Data to be analyzed on this tab must have been saved as a genomic query. For more information, see [Saving a Query](#) on page 43.

To generate a gene expression plot using a genomic query, follow these steps:

- Select the study whose data you want to analyze in the upper right portion of the calIntegrator page. (You must select a study which has genomic data.)

2. Under Analysis Tools on the left sidebar, select **Gene Expression Plot**.
3. Select the **For Genomic Queries** tab (*Figure 5.13*).



Gene Expression Value Plots

For Annotation For Genomic Queries For Clinical Queries

Genomic Query Based Gene Expression Plots

1.) Select a Genomic Query: egfr overexpressed 1x

2.) Select Reporter Type: ☒ Reporter Id ☐ Gene

Reset Create Plot

Figure 5.13 Gene expression value tab for configuring gene expression genomic queries plot

4. **Genomic Query** – Click on the genomic query upon which the plot is to be based.
5. **Reporter Type** – Select the radio button that describes the reporter type:
 - **Reporter ID** – Summarizes expression levels for all reporters you specify.
 - **Gene Name** – Summarizes expression levels at the gene level..

- Click the **Create Plot** button. caIntegrator2 generates the plot which then displays below the plot criteria. Legends below the plot indicate the plot input ([Figure 5.14](#)).



Figure 5.14 A gene expression plot (Mean) based on a genomic query.

- You can recalculate the data display by clicking the **Plot Type** above the graph. See [Understanding a Gene Expression Plot](#) on page 75.
- You can modify the plot parameters and click the **Reset** button to recalculate the plot.

Gene Expression Value Plot for Clinical Queries

Data to be analyzed on this tab must have been saved as a clinical query, but it must have genomic data identified in the query. For more information, see [Adding/Editing Genomic Data](#) on page 24. For the genomic data, you must identify genes whose expression values are used to calculate the plot.

To generate a gene expression plot using a clinical query, follow these steps:



- Select the study whose data you want to analyze in the upper right portion of the caIntegrator page. You must select a study saved as a clinical study, but which has genomic data.
- Under Analysis Tools on the left sidebar, select **Gene Expression Plot**.

3. Select the **For Clinical Queries** tab (*Figure 5.15*).

Gene Expression Value Plots

For Annotation For Genomic Queries **For Clinical Queries**

Clinical Query Based Gene Expression Plots

1.) Gene Symbol(s) (comma separated list):  

2.) Select Reporter Type: ☒ Reporter Id ☐ Gene

3.) Select Queries:

All Available Queries

gender is M
Pathological T Stage = all
gender is F

Add >

< Remove

Selected Queries

▼ ▲

4.) ☐ Exclusive Subjects in Queries (Subjects in upper queries are removed from subsequent queries)

5.) ☐ Add additional group containing all other subjects not found in selected queries.


6.) ☐ Add additional group containing all control samples for this study:

Reset Create Plot

Figure 5.15 Gene expression value tab for configuring gene expression clinical queries plot

4. **Gene Symbol** – Enter one or more gene symbols in the text box or click the icons to locate genes in the following databases. If you enter more than one gene in the text box, separate the entries by commas.

caIntegrator2 provides three methods whereby you can obtain gene names for calculating a gene expression plot:

- **caBio** – This link searches caBIO, then pulls identified genes into caIntegrator2 for analysis.
 - a. Click the caBIO icon ().
 - b. Enter **Search Terms**.
 - c. Select if you want to search in **Gene Keywords**, **Gene Symbols** or **Pathways** (from the drop-down list).
 - Selecting **Gene Keywords** searches only the Full Name field in caBio.
 - Selecting **Gene Symbols** searches only the Unigene and HUGO gene symbols in caBio.
 - Selecting **Pathways** searches only the pathway names in caBio. Note that searching in Pathways is a two step process. First, the initial Pathway search produces search results which are pathways. Second, from the pathway search results screen, you must select pathways of interest, then click **Search Pathways for Genes** to obtain a list of genes related to the selected pathways.
 - d. Select the **Any** or **All** choice to determine how your search terms will be matched. **Any** finds any match for any search term you entered. **All** finds only results that match all of the search terms.

- e. Choose the **Taxon** from the drop-down list and click **Search**. The search results display in the same dialog box (*Figure 5.4*).

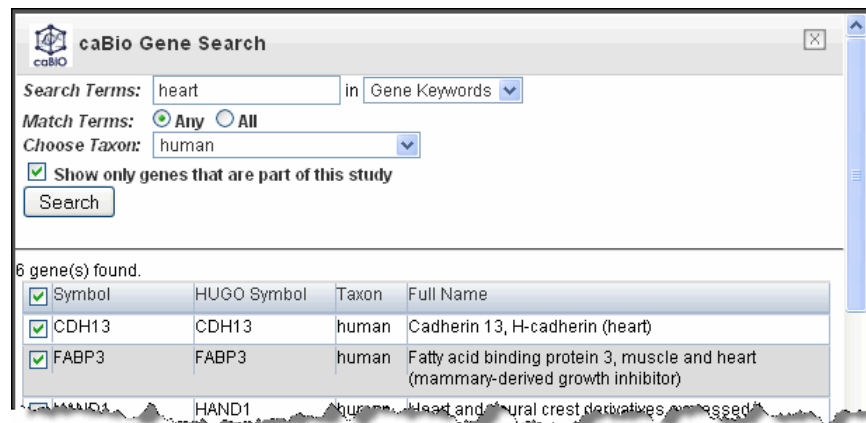



Figure 5.16 Example caBIO gene search results

- f. In the search results, use the check boxes to identify the genes whose symbols you want to use in the plot calculation.
- g. Click **Use Genes** at the bottom of the page. This pulls the checked genes into the tab (*Figure 5.5*).



Figure 5.17 Genes pulled in from caBIO display on the tab

- **Gene List** – This link locates gene lists saved in caIntegrator2.
 - a. Click the Gene List icon () to open a small dialog that lists prior-saved gene lists in caIntegrator2.
 - b. In the drop-down menu, select a gene list. In the list that appears, use the check boxes to identify the genes whose symbols you want to use in the gene expression analysis.
 - c. Click **Use Genes** at the bottom of the dialog. This pulls the checked genes into the tab.
 - **CGAP** – Use this directory to identify genes. Before clicking this link you must enter gene symbols in the text box. This link does not pull anything into caIntegrator2 but does provide information about the gene(s) whose names you enter.
5. For **Reporter Type**, select the radio button that describes the reporter type:
 - **Reporter ID** – Summarizes expression levels for all reporters you specify.
 - **Gene Name** – Summarizes expression levels at the gene level.
 6. For **Sample Groups**, choose among the following options:

- **Annotation Type** – Select the annotation type. Selections are based on the data in the chosen study
 - **Annotation** – Select an annotation. Fields are based on the annotation type you select. For example, if you choose Subject, then you could select Gender or Radiation Type or any field that would distinguish the patients into groups based upon study values.
 - **Values** – Using conventional selection techniques, select one or more values which will be the basis for the plot. Permissible (available) values or “No Values” correspond to the selected annotation.
7. For the **Add Additional Group...** options, define as follows:
- **...all other subjects** – Check the box to create an additional group of all other subjects that are not in selected query groups.
 - **...control group** – Check the box to display an additional group of control samples for this study.
8. Click the **Create Plot** button. caIntegrator2 generates the plot which then displays below the plot criteria in bar graph format (*Figure 5.6*).

By default, caIntegrator2 displays the mean of the data below the plot criteria. Legends below the plot indicate the plot input.

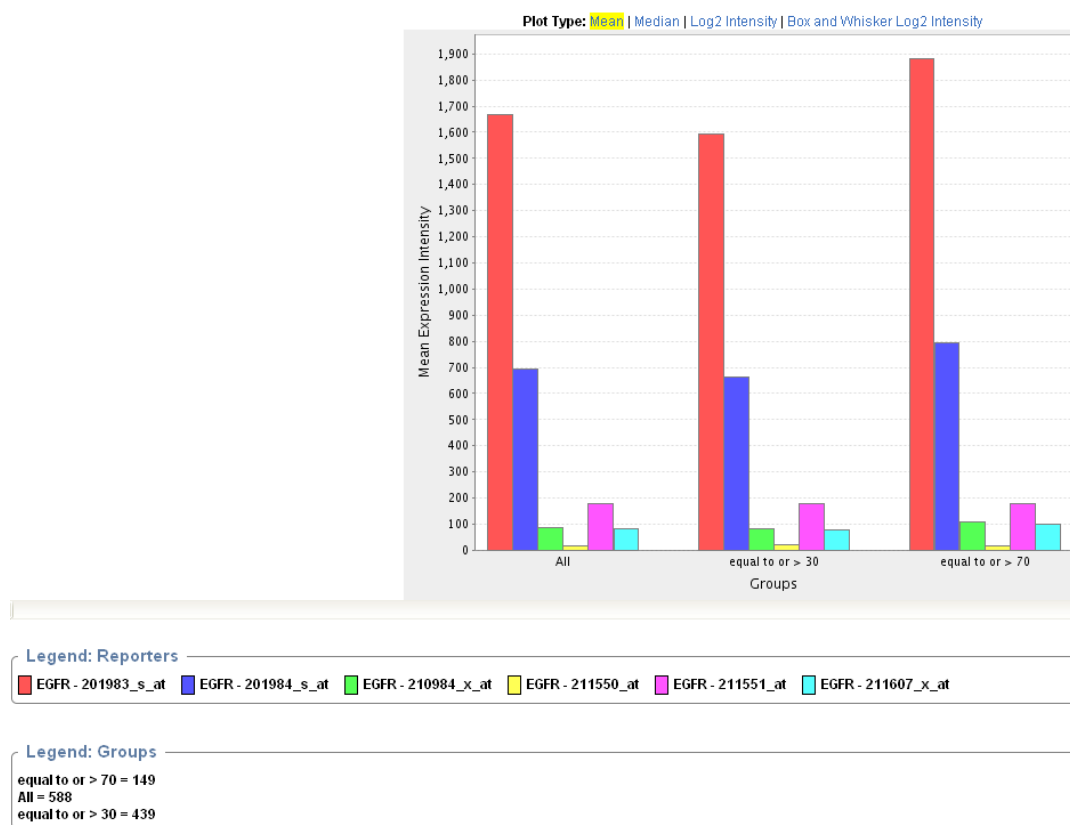


Figure 5.18 Gene expression plot based on clinical queries gene expression values

- You can recalculate the data display by clicking the **Plot Type** above the graph. See [Understanding a Gene Expression Plot](#) on page 75.
- You can modify the plot parameters and click the **Reset** button to recalculate the plot.

Understanding a Gene Expression Plot

Above the plot, you can select various plot types. When you do so, the plot is recalculated. Although all of the plots in this section appear similar, note the differences in calculation results and legends between the Y axis on each of the plots.

When you perform a Gene Expression simple search, by default the **Mean** Gene Expression Plot ([Figure 5.19](#)) appears.

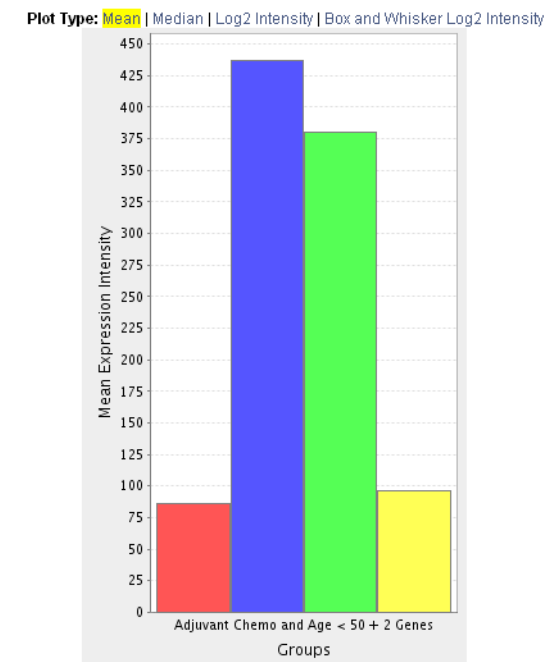


Figure 5.19 Gene expression plot calculating the mean

The **Mean** Gene Expression Plot ([Figure 5.19](#)) displays mean expression intensity (Geometric mean) versus Groups.

The **Median** Gene Expression Plot ([Figure 5.20](#)) displays the median expression intensity versus Groups..

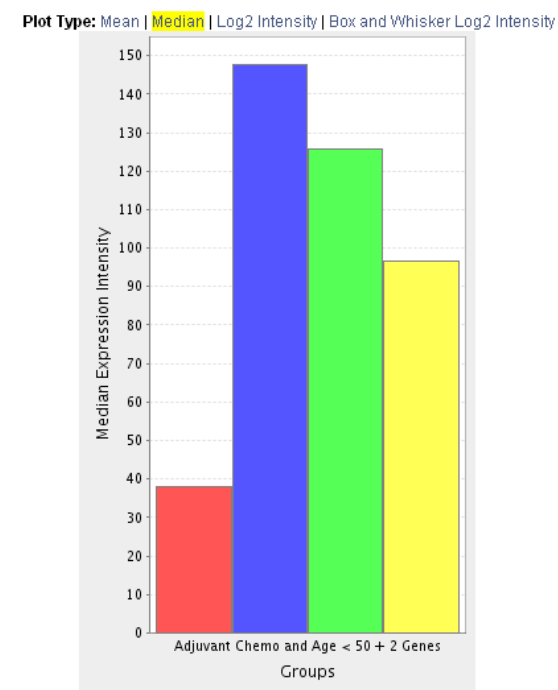


Figure 5.20 Gene expression plot calculating the median

The **Log2 Intensity** Gene Expression Plot ([Figure 5.21](#)) displays average expression intensities for the gene of interest based on Affymetrix GeneChip arrays (U133 Plus 2.0 arrays).

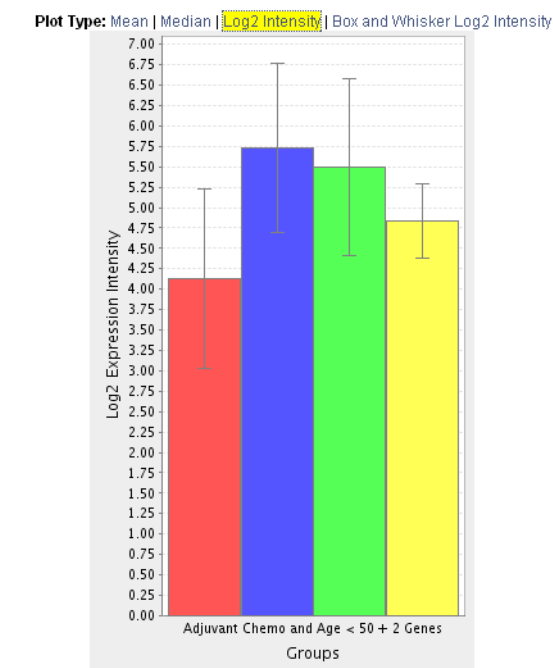


Figure 5.21 Gene expression plot displaying log2 intensity values

The box and whisker log₂ expression intensity plot displays a box plot ([Figure 5.22](#), [Figure 5.23](#)). Example uses of box and whisker plots include the following:

- Indicate whether a distribution is skewed and whether there are potential unusual observations (outliers) in the data set.
- Perform a large number of observations.
- Compare two or more data sets.
- Compare distributions because the center, spread, and overall range are immediately apparent.

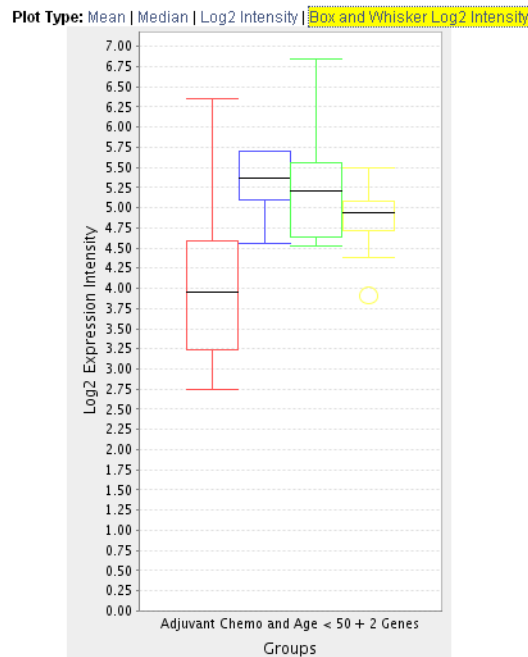


Figure 5.22 Box and whisker plot based on the same data set as represented in [Figure 5.19](#), [Figure 5.20](#), [Figure 5.21](#)

In descriptive statistics, a box plot or boxplot, also known as a box-and-whisker diagram or plot, is a convenient way of graphically depicting groups of numerical data through their five-number summaries (the smallest observation excluding outliers, lower quartile [Q1], median [Q2], upper quartile [Q3], and largest observation excluding outliers).

The box is defined by Q1 and Q3 with a line in the middle for Q2. The interquartile range, or IQR, is defined as Q3-Q1. The lines above and below the box, or 'whiskers', are at the largest and smallest non-outliers. Outliers are defined as values that are

more than $1.5 \times \text{IQR}$ greater than Q3 and less than $1.5 \times \text{IQR}$ than Q1. Outliers, if present, are shown as open circles (*Figure 5.23*).

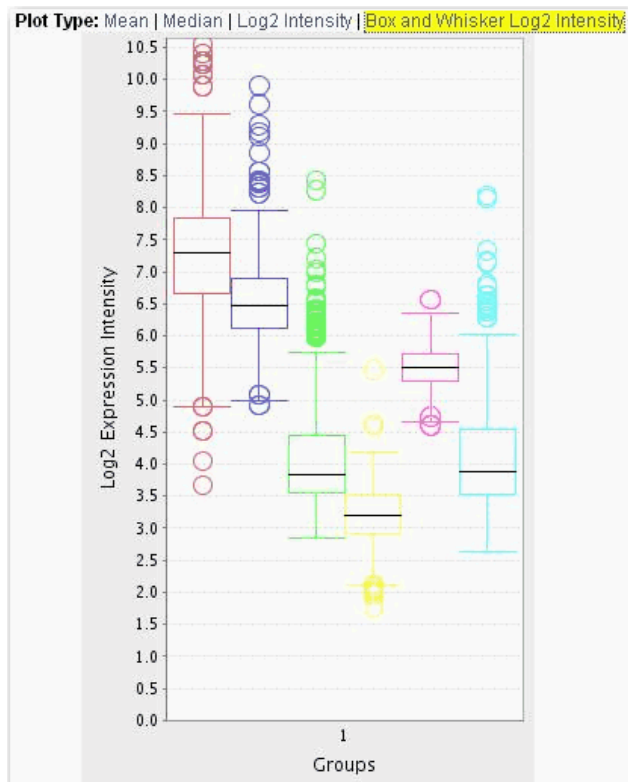


Figure 5.23 Box and whisker plot showing outliers

Boxplots can be useful to display differences between populations without making any assumptions of the underlying statistical distribution: they are non-parametric. The spacings between the different parts of the box help indicate the degree of dispersion (spread) and skewness in the data.

Analyzing Data with GenePattern

GenePattern is an application developed at the Broad Institute that enables researchers to access various methods to analyze genomic data. calIntegrator2 provides an express link to GenePattern where you can analyze data in any calIntegrator2 study.

Information is included in this section for connecting to GenePattern from calIntegrator2. Specifics for launching GenePattern tools from calIntegrator2 are included as well, but you may want to refer to additional GenePattern documentation, available at this website: http://www.broadinstitute.org/cancer/software/genepattern/tutorial/gp_concepts.html.

You have two options for using GenePattern from calIntegrator2:

- Option 1 – Use the web-interface of any available GenePattern instances.
 - a. To use the public instance from Broad, first register for an account at <http://genepattern.broad.mit.edu/gp/pages/login.jsf>
 - b. In calIntegrator2, enter the URL for connecting: <http://genepattern.broad.mit.edu/gp/services/>, then enter your userId and password.
- Option 2 – Use GenePattern on the grid.

The GenePattern feature in calIntegrator2 currently supports three analyses on the grid: Comparative Marker Selection (CMS), Principal Component Analysis (PCA) and GISTIC-supported analysis.

Tip: If you are using the web interface to access GenePattern (option #1 listed above), then you can run other GenePattern tools in addition to CMS, PCA and GISTIC.

1. Select the study whose data you want to analyze in the upper right portion of the calIntegrator2 page.
2. Click **GenePattern Analysis** in the left sidebar of calIntegrator2. This opens the GenePattern Analysis Status page (*Figure 5.24*).

GenePattern Analysis Status

Gene Pattern Modules		New Analysis Job		
Job Name	Job Type	Status	Creation Date	Status
JP - CMS - 2	Comparative Marker Selection	Completed - Download	2009/08/26 21:38:03	2009/08
CMS 1	Comparative Marker Selection	Completed - Download	2009/08/26 11:43:44	2009/08
PCA1	Principal Component Analysis	Completed - Download	2009/08/26 11:38:41	2009/08

Figure 5.24 GenePattern Analysis Status page

3. Select from the drop-down list the type of GenePattern analysis you want to run on the data.
 - **GenePattern Modules** – This option launches a session within GenePattern from which you can launch analyses. See *GenePattern Modules* on page 80.
 - **Comparative Marker Selection (Grid Service)**. This option enables you to run this GenePattern analysis on the grid. See *Comparative Marker Selection (CMS) Analysis* on page 81.
 - **Principal Component Analysis (Grid Service)**. This option enables you to run this GenePattern analysis on the grid. See *Principal Component Analysis (PCA)* on page 83.
 - **GISTIC (Grid Service)**. This option enables you to run this GenePattern analysis on the grid. See *GISTIC-Supported Analysis* on page 86.

- Click the **New Analysis Job** button to open a corresponding page where you can configure the analysis parameters.

GenePattern Modules

Note: To launch the analyses described in this section, you must have a registered GenePattern account. For more information, see <http://genepattern.broad.mit.edu/gp/pages/login.jsf>.

To configure the link for accessing GenePattern from caIntegrator2, open the appropriate page as described in *Analyzing Data with GenePattern* on page 78.

- Select the study whose data you want to analyze in the upper right portion of the caIntegrator2 page.
- Click **GenePattern Analysis** in the left sidebar of caIntegrator2. This opens the GenePattern Analysis Status page.
- Make sure **GenePattern Modules** is selected in the drop down list. Click **New Analysis Job**.
- In the GenePattern Analysis dialog box (*Figure 5.25*), specify connection information, as described *Table 5.1*.

GenePattern Analysis

GenePattern Server URL*:	<input type="text"/>
GenePattern Username*:	<input type="text"/>
GenePattern Password:	<input type="password"/>
<input type="button" value="Connect"/>	

Figure 5.25 Dialog box for configuring the link to GenePattern

Fields	Description
Server URL	Enter any GenePattern publicly available URL, such as http://genepattern.broad.mit.edu/gp/services/Analysis .
GenePattern Username	Enter your GenePattern user name.
GenePattern Password	Enter your GenePattern password.

Table 5.1 Fields for selecting GenePattern configurations

If you choose to access GenePattern in this way, you can continue to use GenePattern tools from within that application. See GenePattern user documentation for more information.

Tip: If you run these analysis within GenePattern itself, you may be able to view results in the GenePattern visualization module. If you run them on the grid from caIntegrator2, your results will be available only in spreadsheet and XML format.

You can run GenePattern analyses for Comparative Marker Selection, Principal Component Analysis and GISTIC-based analysis on the grid if you choose.

Comparative Marker Selection (CMS) Analysis

The Comparative Marker Selection (CMS) module implements several methods to look for expression values that correlate with the differences between classes of samples. Given two classes of samples, CMS finds expression values that correlate with the difference between those two classes. If there are more than two classes, CMS can perform one-vs-all or all-pairs comparisons, depending on which option is chosen.

For more information, see the GenePattern website: http://www.broad.mit.edu/cgi-bin/cancer/software/genepattern/modules/gp_modules.cgi.

To perform a CMS analysis, follow these steps:

1. Select the study whose data you want to analyze in the upper right portion of the calIntegrator2 page. You must select a study saved as a clinical study, but which has genomic data.
2. Click **GenePattern Analysis** in the left sidebar of calIntegrator2. This opens the GenePattern Analysis Status page.
3. In the GenePattern Analysis Status page, select **Comparative Marker Selection (Grid Service)** from the drop down list and click **New Analysis Job**. This opens the Comparative Marker Selection Analysis page (*Figure 5.26*).

Comparative Marker Selection Analysis

Job Name*:

Preprocess Server*: Default Broad service - <http://node255.broad.mit.edu:6060/wsrf/services/cagrid/PreprocessDatasetMAGEService> ▼

Comparative Server*: Default Broad service - <http://node255.broadinstitute.org:11010/wsrf/services/cagrid/ComparativeMarkerSelMAGESvc> ▼

Clinical Queries*: Must select two clinical queries, which are used to group the samples into two separate classifications to run against ComparativeMarkerSelection. The queries selected here have been previously saved by the user. Selected queries will result in the processing of only those samples which are mapped to patients in the saved query result.

All Available Queries

Selected Queries

Filter flag: ☐

Preprocessing Flag*: no-disc-or-norm ▼

Min Change*:

Min Delta*:

Threshold*:

Ceiling*:

Max Sigma Binning*:

Probability Threshold*:

Num Exclude*:

Log Base Two: ☐

Number Of Columns Above Threshold*:

Test Direction*: two-sided ▼

Test Statistic*: T-test ▼

Min Std*:

Number Of Permutations*:

Complete: ☐

Balanced: ☐

Random Seed*:

Smooth Pvalues: ☐

Phenotype Test*: one-versus-all ▼

Figure 5.26 Comparative Marker Selection analysis parameters

4. Select or define CMS analysis parameters, described in [Table 5.2](#). An asterisk indicates required fields. The default settings are valid; they should provide valid results.

CMS Parameter	Description
Job Name*	Assign a unique name to the analysis you are configuring.
Preprocess Server*	A server which hosts the grid-enabled data GenePattern PreProcess Dataset module. Select one from the list and calIntegrator2 will use the selected server for this portion of the processing.
Comparative Server*	A server which hosts the grid-enabled data GenePattern Comparative Marker Selection module. Select one from the list and calIntegrator2 will use the selected server for this portion of the processing.
Clinical Queries*	All clinical queries with appropriate data for the analysis are listed. Select and move 2 or more queries from the All Available Queries panel to the Selected Queries panel. Note: If a query has a genomic component (e.g. gene criteria), it does not display in the queries field.
Filter Flag	Variation filter and thresholding flag
Preprocessing Flag*	Discretization and normalization flag
Min Change*	Minimum fold change for filter
Min Delta*	Minimum delta for filter
Threshold*	Value for threshold
Ceiling*	Value for ceiling
Max Sigma Binning*	Maximum sigma for binning
Probability Threshold*	Value for uniform probability threshold filter
Num Exclude*	Number of experiments to exclude (max & min) before applying variation filter
Log Base Two	Whether to take the log base two after thresholding
Number of Columns Above Threshold*	Remove row if n columns no \geq than the given threshold
Test Direction*	The test to perform (up-regulated for class0; up-regulated for class1, two sided). By default, Comparative Marker Selection performs the two-sided test.
Test Statistic*	Select the statistic to use.
Min Std*	The minimum standard deviation if test statistic includes the min std option. Used only if test statistic includes the min std option.

Table 5.2 Comparative Marker Selection analysis options

CMS Parameter	Description
Number of Permutations*	<p>The number of permutations to perform. (Use 0 to calculate asymptotic P-values.) The number of permutations you specify depends on the number of hypotheses being tested and the significance level that you want to achieve (3). The greater the number of permutations, the more accurate the P-value.</p> <p>Complete – Perform all possible permutations. By default, complete is set to No and Number of Permutations determines the number of permutations performed. If you have a small number of samples, you might want to perform all possible permutations.</p> <p>Balanced – Perform balanced permutations</p>
Random Seed*	The seed for the random number generator.
Smooth Pvalues	Whether to smooth P-values by using the Laplace's Rule of Succession. By default, Smooth Pvalues is set to Yes , which means P-values are always less than 1.0 and greater than 0.0.
Phenotype Test*	<p>Tests to perform when class membership has more than 2 classes: one versus-all, all pairs.</p> <p>Note: The P-values obtained from the one-versus-all comparison are not fully corrected for multiple hypothesis testing.</p>

Table 5.2 Comparative Marker Selection analysis options

- When you have completed the form, click **Perform Analysis**.

calIntegrator2 takes you to the JobStatus/Launch page where you will see the job and its status in the Status column of the list ([Figure 5.27](#)).

GenePattern Analysis Status

(draft)

Gene Pattern Modules


Job Name	Job Type	Status	Creation Date	Status Update Date
Well-diff vs adjuvant chemo	Comparative Marker Selection	 Processing Locally	2009/08/14 11:48:35	2009/08/14 11:48:35
Filter out non-interesting genes	Gene Pattern	Completed - View 122444	2009/08/14 10:16:29	2009/08/14 10:19:47

Figure 5.27 The progress of a GenePattern analysis that has been launched displays in the status column of page

- When the job is complete, the system displays a completion date on the GenePattern Analysis status page. Click the **Download** link. This downloads zipped result files to your local work station. The number of files and their file type will vary according to the processing. The results format is compatible with GenePattern visualizers and can be uploaded within GenePattern.

Principal Component Analysis (PCA)

Principal Component Analysis is typically used to transform a collection of correlated variables into a smaller number of uncorrelated variables, or components. Those components are typically sorted so that the first one captures most of the underlying variability and each succeeding component captures as much of the remaining variability as possible.

You can configure GenePattern grid parameters for preprocessing the dataset in addition to PCA module parameters. For more information, see the GenePattern website: http://www.broad.mit.edu/cgi-bin/cancer/software/genepattern/modules/gp_modules.cgi.

To perform a PCA analysis, follow these steps:

1. Select the study whose data you want to analyze in the upper right portion of the calIntegrator page. You must select a study with gene expression data.
2. Click **GenePattern Analysis** in the left sidebar of calIntegrator2. This opens the GenePattern Analysis Status page.
3. In the GenePattern Analysis Status page, select **Principal Component Analysis (Grid Service)** from the drop down list and click **New Analysis Job**. This opens the Principal Component Analysis page ([Figure 5.28](#)).

Principal Component Analysis

(draft)

This form submits a job which analyzes samples using the GenePattern Principal Component Analysis module.

Job Name - Please enter a job name.
Principal Component Analysis Server - Select a PCA grid service from the dropdown.
Clinical Queries - Select saved Clinical queries to specify which samples will be processed.
Enable Preprocess Dataset - (Optional) Check this to display and configure preprocessing parameters.

* Job Name:

* Principal Component Analysis Server: Default Broad service - <http://node255.broad.mit.edu:6060/awstf/services/cagrid/PCA>

* Clinical Queries: Clinical Queries enable the user to specify which samples will be processed using PCA. The queries selected here have been previously saved by the user. Selected queries will result in the processing of only those samples which are mapped to subjects in the saved query result. If multiple queries are selected, all of the sample from each saved query are processed PLUS the results set will be classified according to those queries. (One class per selected query.)

All Available Queries

- gender female
- Never smoke

Selected Queries

Add >

< Remove

Enable Preprocess Dataset: ☐
(check to display preprocess parameters)

Perform Analysis

Figure 5.28 Principal Component Analysis parameters

4. Select or define PCA analysis parameters, described in [Table 5.3](#). You must enter a job name and select a clinical query, but you can accept the other default settings..

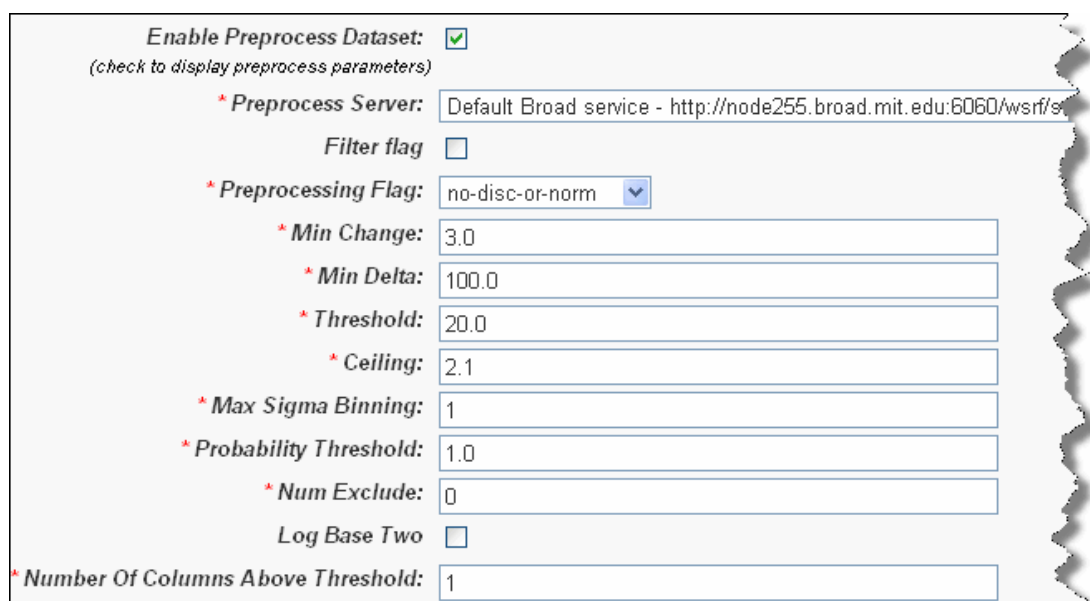
PCA Parameters	Description
Job Name*	Assign a unique name to the analysis you are configuring.
Principal Component Analysis Server*	A server which hosts the grid-enabled data GenePattern Principal Component Analysis module. Select one from the list and calIntegrator2 will use the selected server for this portion of the processing.
Clinical Queries*	All clinical queries display in this list. Select one or more of these queries to define which samples are analyzed using PCA. If you select more than one query, then the union of the samples returned by the multiple queries is analyzed.

Table 5.3 PCA analysis options

PCA Parameters	Description
Cluster By*	Selecting rows looks for principal components across all expression values, and selecting columns looks for principal components across all samples.

Table 5.3 PCA analysis options

5. If you want to preprocess the data set, click **Enable the Preprocess Dataset**. This opens an additional set of parameters (Figure 5.29), discussed in Table 5.4. The preprocessing is executed prior to running the PCA.



Enable Preprocess Dataset: ☒
 (check to display preprocess parameters)

* Preprocess Server: Default Broad service - http://node255.broad.mit.edu:6060/wsrf/s

Filter flag ☐

* Preprocessing Flag: no-disc-or-norm

* Min Change: 3.0

* Min Delta: 100.0

* Threshold: 20.0

* Ceiling: 2.1

* Max Sigma Binning: 1

* Probability Threshold: 1.0

* Num Exclude: 0

Log Base Two ☐

* Number Of Columns Above Threshold: 1

Figure 5.29 Parameters for pre-processing parameters for PCA

PCA Preprocessing Parameters	Description
Preprocess Server*	A server which hosts the grid-enabled data GenePattern PreProcess Dataset module. Select one from the list and caIntegrator2 will use the selected server for this portion of the processing.
Filter Flag	Variation filter and thresholding flag
Preprocessing Flag	Discretization and normalization flag
Min Change	Minimum fold change for filter
Min Delta	Minimum delta for filter
Threshold	Value for threshold
Ceiling	Value for ceiling
Max Sigma Binning	Maximum sigma for binning
Probability Threshold	Value for uniform probability threshold filter

Table 5.4 Parameters for preprocessing data sets for PCA

PCA Preprocessing Parameters	Description
Num Exclude	Number of experiments to exclude (max & min) before applying variation filter
Log Base Two	Whether to take the log base two after thresholding
Number of Columns Above Threshold	Remove row if n columns no \geq than the given threshold

Table 5.4 Parameters for preprocessing data sets for PCA

- When you have completed the form, click **Perform Analysis**.
- When the job is complete, the system displays a completion date on the GenePattern Analysis status page. Click the **Download** link. This downloads zipped result files to your local work station. The number of files and their file type will vary according to the processing. The results format is compatible with GenePattern visualizers and can be uploaded within GenePattern.

GISTIC-Supported Analysis

The GISTIC Module is a GenePattern tool that identifies regions of the genome that are significantly amplified or deleted across a set of samples. For more information, see http://www.broad.mit.edu/cgi-bin/cancer/software/genepattern/modules/gp_modules.cgi.

To perform a GISTIC-supported analysis, follow these steps:

- Select the study whose data you want to analyze in the upper right portion of the calIntegrator2 page. You must select a study with copy number (either Affymetrix SNP or Agilent Copy Number) data.
- Click **GenePattern Analysis** in the left sidebar of calIntegrator2. This opens the GenePattern Analysis Status page.

3. In the GenePattern Analysis Status page, select **GISTIC (Grid Service)** from the drop down list and click **New Analysis Job**. This opens the GISTIC Analysis page (*Figure 5.30*).

GISTIC Analysis

This form submits a job which analyzes samples using the GenePattern GISTIC module.

Job Name - Please enter a job name.
GenePattern Server URL / GISTIC Server - Select whether to use the GISTIC web service or grid service and provide or select the service address. If the web service is selected, authentication information is also required.
Clinical Queries - (Optional) Select a saved Clinical query to specify which samples will be processed.
Exclude Sample Control Set - (Optional) Select a Control Sample Set to be excluded from the Clinical Query.

* Job Name:

* GenePattern Server URL:

* GenePattern Username:

GenePattern Password:

* GISTIC Server:

For the Clinical query parameter below, choose either "All Samples" or a clinical query. If "All Samples" is selected, then all samples will be used. If a clinical query is selected, only those samples which map to the subjects in the clinical query results will be used. The clinical queries in this list have been previously saved by the user. Control samples can be excluded from this processing by selecting a control set name in the Exclude Sample Control Set dropdown.

Clinical query:

* Exclude Sample Control Set:

* Amplifications Threshold:

* Deletions Threshold:

* Join Segment Size:

* QV Thresh:

* Remove X:

cnv File:

Figure 5.30 GISTIC analysis criteria

4. Select or define GISTIC analysis parameters, as described in *Table 5.2*. You must indicate a Job Name, but you can accept the other defaults settings, which are valid and should produce valid results.

GISTIC Parameters	Description
Job Name*	Assign a unique name to the analysis you are configuring.
GISTIC Server*	A server which hosts the grid-enabled data GISTIC-based analysis module. Select one from the list and caIntegrator2 will use the selected server for this portion of the processing.
Refgene File*	Enter or select the cytoband file to use in the analysis. Allowed values: {Human Hg18, Human Hg17, Human Hg16}. Default = Human Hg16.
Clinical Query	All clinical queries display in this list as well as an option to select all non-control samples. Select a clinical query if you wish to run GISTIC on a subset of the data and select all non-control samples if wish to include all samples.
Amplifications Threshold*	Threshold for copy number amplifications. Regions with a log2 ratio above this value are considered amplified. Default = 0.1.
Deletions Threshold*	Threshold for copy number deletions. Regions with a log2 ratio below the negative of this value are considered deletions. Default = 0.1.

Table 5.5 GISTIC analysis parameters

GISTIC Parameters	Description
Join Segment Size*	Smallest number of markers to allow in segments from the segmented data. Segments that contain fewer than this number of markers are joined to the neighboring segment that is closest in copy number. Default = 4.
QV Thresh[hold]*	Threshold for q-values. Regions with q-values below this number are considered significant. Default = 0.25.
Remove X*	Flag indicating whether to remove data from the X-chromosome before analysis. Allowed values = {1,0}. Default = 1(yes).
cnv File	<p>This selection is optional.</p> <p>Browse for the file. There are two options for the cnv file.</p> <p>Option #1 enables you to identify CNVs by marker name. Permissible file format is described as follows:</p> <p>A two column, tab-delimited file with an optional header row. The marker names given in this file must match the marker names given in the markers_file. The CNV identifiers are for user use and can be arbitrary. The column headers are:</p> <ol style="list-style-type: none"> 1. Marker Name 2. CNV Identifier <p>Option #2 enables you to identify CNVs by genomic location. Permissible file format is described as follows:</p> <p>A 6 column, tab-delimited file with an optional header row. The 'CNV Identifier', 'Narrow Region Start' and 'Narrow Region End' are for user use and can be arbitrary. The column headers are:</p> <ol style="list-style-type: none"> 1. CNV Identifier 2. Chromosome 3. Narrow Region Start 4. Narrow Region End 5. Wide Region Start 6. Wide Region End

Table 5.5 GISTIC analysis parameters

5. When you have completed the form, click **Perform Analysis**.
6. When the job is complete, the system displays a completion date on the GenePattern Analysis status page. Click the **Download** link. This downloads zipped result files to your local work station. The number of files and their file type will vary according to the processing. The results format is compatible with GenePattern visualizers and can be uploaded within GenePattern.

CHAPTER 6

ADMINISTERING USER ACCOUNTS

This chapter describes the process for creating and managing user accounts in calIntegrator2. It also discusses the processes for managing ownership and access to studies in calIntegrator2.

Note: The options for performing user management tasks are visible in calIntegrator2 on the left sidebar of the browser only if you have these Admin privileges.

Administering calIntegrator2 User Accounts Using UPT

Note: If you are interested in registering an account in calIntegrator2, see *Registering as a New calIntegrator2 User* on page 6.

In calIntegrator2, all tasks related to creating and managing user accounts can be performed only by a calIntegrator2 administrator using the CBIIT User Provisioning Tool (UPT) v. 4.2. The following sections discuss the use of the UPT for performing these tasks. For further information about UPT, see Chapter 3 of the CSM 4.2 Programmer's Guide located here: https://gforge.nci.nih.gov/docman/view.php/12/18945/caCORE_CSM_v42_ProgrammersGuide.pdf

The UPT is a separately installed application which serves as the user management interface for all National Cancer Institute CBIIT Life Sciences Distribution (LSD) applications, including calIntegrator2. The UPT application is the central point for all user management functionality within calIntegrator2. You can use UPT to add new users and to apply user group assignments to the calIntegrator2 database directly. The UPT groups can refer to predefined groups such as Study Manager or Study Investigator, which determine what roles the user has.

The following terms are used both in this chapter and in the UPT to define user-related roles:

- **User** – a person who is accessing calIntegrator2. The user has an associated account and user id.

- **User Group** – a group of users, typically grouped by organization and role, for example, “Columbia University Study Managers”
- **Protection Group** – a group of studies given a secure status and typically grouped by organization, for example, “Columbia University Protected Studies”.

Steps for Creating User Access to caIntegrator2

The following steps summarize the process for establishing user access to caIntegrator2:

1. A potential user requests a user account in caIntegrator2. See [Registering as a New caIntegrator2 User](#) on page 6.
1. You, as a caIntegrator2 administrator, check if the **User** already exists in caIntegrator2. If not, create the new user. See [Creating a New caIntegrator2 User](#) on page 90.
2. Check if the requestor's **User Group** already exists in caIntegrator2. If not, create a new **User Group**. See [Creating a New User Group](#) on page 92.
3. Check if the **Protection Group** (e.g. “Columbia University Protected Studies”), containing the studies to which this user wants access currently exists. If not, create a new **Protection Group**. See [Creating a New Protection Group](#) on page 93.

Note: If the Protection Group already exists, contact the Organizational Contact person to confirm that it is OK to give this person access to this Protection Group.

4. Give the requestor's **User Group** access to the **Protection Group**. See [Assigning a User Group to a Protection Group](#) on page 94.
5. Add the **User** to the **User Group**. See [Adding a User to a User Group](#) on page 97

Creating a New caIntegrator2 User

To create a new User in caIntegrator2, follow these steps:

1. Login to UPT as a caIntegrator2 Admin.
2. First, search to see if the user already exists. Click the **User** menu option.
3. On the User page that opens, click **Select an Existing User**.

- Use the form and search for the user. If you define no criteria, UPT returns a list of all caIntegrator2 users currently in the system (*Figure 6.1*).

The screenshot shows the UPT interface with the following components:

- Header:** Common Security Module User Provisioning Tool. Login ID: boalt, Application: caIntegrator2, Role: Admin.
- Navigation Bar:** HOME, USER, PROTECTION ELEMENT, PRIVILEGE, GROUP, PROTECTION GROUP, ROLE, INSTANCE LEVEL, LOG OUT.
- Section Header:** User
- SEARCH RESULTS Table:**

Select	User Login Name	User First Name	User Last Name	User Organization	User Department	User Email Id
<input type="radio"/>	admin	UPT	Administrator			
<input type="radio"/>	cai2admin	cai2	Admin			
<input type="radio"/>	gumanager	Georgetown	Study Manager			
<input type="radio"/>	investigator	Research	Investigator			
<input type="radio"/>	manager	Study	Manager			
<input type="radio"/>	manager2	Study	Manager2			
<input type="radio"/>	manager3	Study	Manager3			
<input type="radio"/>	manager4	Study	Manager4			
<input type="radio"/>	manager5	Study	Manager5			
<input type="radio"/>	nblamanager	NBLA	Study Manager			
<input type="radio"/>	tcgaprivate	TCGA	Manager			
- Buttons:** View Details, Back

Figure 6.1 A list of current caIntegrator2 users displays in UPT after a user search

- If the user does not already exist (is not listed in the search results), then create a new user. To do so, select the **User** menu option again, then click **Create a New User**.

This opens the page for creating a new caIntegrator2 user (*Figure 6.2*).

Common Security Module User Provisioning Tool

Login ID : boalt
Application : caIntegrator2
Role : Admin

HOME USER PROTECTION ELEMENT PRIVILEGE GROUP PROTECTION GROUP ROLE INSTANCE LEVEL LOG OUT

Enter the details to add a new User. The **User Login Name** uniquely identifies the User and is a required field. The **User First Name** and **User Last Name** identifies the User. The **User Organization**, **User Department** and **User Title** provides his work details. The **User Phone Number** and **User Email Id** provides the contact details for the User. The **User Password** can be entered if the same schema is also going to be used for Authentication. The **User Start Date** and **User End Date** determine the period for which the User is a valid User.

* indicates a required field

ENTER THE NEW USER DETAILS	
*	User Login Name <input type="text"/>
*	User First Name <input type="text"/>
*	User Last Name <input type="text"/>
	User Organization <input type="text"/>
	User Department <input type="text"/>
	User Title <input type="text"/>
	User Phone Number <input type="text"/>
	User Password <input type="text"/>
	Confirm Password <input type="text"/>
	User Email Id <input type="text"/>
	User Start Date <input type="text"/> (MM/DD/YYYY)
	User End Date <input type="text"/> (MM/DD/YYYY)

Add Reset Back

Figure 6.2 UPT page for creating new user details

6. Enter details for the following required fields:

- **User Login Name**
- **User First Name**
- **User Last Name**
- **User Password**

Caution: If the requestor is an LDAP user, then the User Login Name must match the LDAP login id AND the User Password field must be left blank. If the requestor is not an LDAP user, then provide a password.

- **User Organization**
- **User Department**

7. Click **Add** to confirm the new user.

Creating a New User Group

You can assign a user group to a protection group. The advantage of working with a user group is that you do not have to assign roles to each user individually. You can assign users to a user group to which you assign a role, and then assign that user group to the protection group, or you can assign a role collectively to a protection group after it is created.

To create a new User Group in calIntegrator2, follow these steps:

1. Login to UPT as calIntegrator2 Admin.
2. First search for an existing group that the user wishes to join. Click the **Group** menu option.
3. On the Group page that opens, click **Select an Existing Group**.
4. Use the form and search for the group. If you define no criteria, UPT returns a list of all calIntegrator2 groups currently in the system
5. If a user group does not already exist, then create a new user group. Click the **Group** menu option, then click **Create a new Group**.
6. On the form that opens (*Figure 6.3*), enter a unique **Group Name** and a description, if appropriate. Click **Add**.

Note: The recommended naming convention for a new User Group is *[insert organization name] Study [insert role]s* Example: "Columbia University Study Managers".

Common Security Module
User Provisioning Tool

Login ID : boalt
Application : calIntegrator2
Role : Admin

HOME USER PROTECTION ELEMENT PRIVILEGE **GROUP** PROTECTION GROUP ROLE INSTANCE LEVEL LOG OUT

Enter the details to add a new Group. The **Group Name** uniquely identifies the Group and is a required field. The **Group Description** is a brief summary about the Group.

* indicates a required field

ENTER THE NEW GROUP DETAILS

* **Group Name**

Group Description

Add Reset Back

Figure 6.3 UPT page for creating a new group

Creating a New Protection Group

: If you prefer that a study or group of studies have limited access, you can assign a user to a particular protection group and assign roles which allow the users in the protection group study access. : A protection group provides security or limited access for studies listed there.

To create a new Protection Group in calIntegrator2, follow these steps:

1. Login to UPT as calIntegrator2 Admin.
2. Click the **Protection Group** menu option.

- On the page that opens, click **Create a New Protection Group**. The page opens for defining PG Group details (*Figure 6.4*).

Common Security Module
User Provisioning Tool

Login ID : boalt
Application : calintegrator2
Role : Admin

HOME USER PROTECTION ELEMENT PRIVILEGE GROUP PROTECTION GROUP ROLE INSTANCE LEVEL LOG OUT

Enter the details to add a new Protection Group. The **Protection Group Name** uniquely identifies the Protection Group and is a required field. The **Protection Group Description** is a brief summary about the Protection Group. The **Protection Group Large Count Flag** is used to indicate if the Protection Group has a large number of associated Protection Elements.

* indicates a required field

ENTER THE NEW PROTECTION GROUP DETAILS

* **Protection Group Name**

Protection Group Description

Protection Group Large Count Flag ☐ Yes ☒ No

Add Reset Back

Figure 6.4 UPT page for creating a new protection group

- Enter a unique **Protection Group Name** and Description, if appropriate. Click **Add**.

Note: The recommended naming convention is *[insert organization name here] Protected Studies*. Example: "Columbia University Protected Studies".

Assigning a User Group to a Protection Group

To give a User Group access to a Protection Group (a group of protected studies), follow these steps:

- Login to UPT as calIntegrator2 Admin.
- Find the user group that you want to assign. Click the **Group** menu option and click **Select an Existing Group**. In the page that opens, click **Search**. If you define no criteria, UPT returns a list of all calIntegrator2 groups currently in the system (*Figure 6.5*).

Group

SEARCH RESULTS		
Select	Group Name	Group Description
<input type="radio"/>	Study Managers Group 3	Study Managers who can create/modify any Group 3 studies.
<input type="radio"/>	Study Managers Group 4	Study Managers who can create/modify any Group 4 studies.
<input type="radio"/>	Study Managers Group 5	Study Managers who can create/modify any Group 5 studies.
<input type="radio"/>	NCI Study Investigators	Study investigators for the NCI studies.
<input type="radio"/>	NCI Study Managers	Study Managers who can create/modify any NCI studies.
<input type="radio"/>	Platform Manager Group	The platform manager group.
<input type="radio"/>	TCGA Study Managers	Study Managers who can create/modify any TCGA studies.

View Details Back

Figure 6.5 UPT page showing Group search results

3. Select the radio button next to the group name you want to assign to the Protection Group. Click **View Details**. This opens the Group Details page (Figure 6.6).

Common Security Module
User Provisioning Tool

Login ID : boalt
Application : calintegrator2
Role : Admin

HOME USER PROTECTION ELEMENT PRIVILEGE GROUP PROTECTION GROUP ROLE INSTANCE LEVEL LOG OUT

Update the details of the displayed Group. The **Group Name** uniquely identifies the Group and is a required field. The **Group Description** is a brief summary about the Group. The **Update Date** indicates the date when this Group's Details were last updated.

GROUP DETAILS	
* Group Name	NCI Study Managers
Group Description	Study Managers who can create/modify any NCI studies.
Group Update Date	09/24/2009 (MM/DD/YYYY)

Update Delete Back

Associated Users Associated PE & Privileges Associated PG & Roles Assign PG & Roles

Figure 6.6 UPT page showing details for a selected group

4. Below the group details, click **Associated PG & Roles**. The page that opens displays any PG to which the user group is already assigned (Figure 6.7).

Group, Protection Group and Roles

SELECTED GROUP

Group Name
NCI Study Managers

Select the **Protection Group** association which to be removed for the selected **Group** or whose **Roles** Association needs to be updated.

SEARCH RESULTS		
Select	Associated Protection Group Name	Associated Role Name
<input type="radio"/>	NCI Protected Studies	STUDY_MANAGER_ROLE

Remove PG & Roles Associated Roles Back

Figure 6.7 UPT page that shows any PGs to which the select user group is assigned

5. Below the group name, examine if the Protection Group of your choice is already listed there. If so, this means your user group is already assigned to the protection group of choice, and you can skip the remainder of the steps in this section. If the Protection Group is not listed there, then click **Back**.

- Back on the User Group details page, click **Assign PG & Roles**. This opens the Group, Protection Group and Roles Association page ([Figure 6.8](#)).

Group, Protection Group and Roles Association

SELECTED GROUP	
Group Name	NCI Study Managers

Select a single **Protection Group** to associate with the selected **Group**.

AVAILABLE PROTECTION GROUPS
Platforms
Protected Studies for Group 3
Protected Studies for Group 4
Protected Studies for Group 5
TCGA Protected Studies

ASSIGNED PROTECTION GROUP

Select **Roles** which are to be associated with the selected **Group**.

AVAILABLE ROLES
PLATFORM_MANAGER_ROLE
STUDY_INVESTIGATOR_ROLE
STUDY_MANAGER_ROLE

ASSIGNED ROLES

Figure 6.8 UPT page for assigning user group to a protection group and selected roles

- From the list of Available Protection Groups, highlight your PG of choice and click **Assign**.

Now you can assign a role to the user. The caIntegrator2 Roles are defined in [Table 6.1](#):

Role Name	Role Definition
STUDY_MANAGER_ROLE	Assigning this role allows the user to modify existing studies, create new studies, and deploy existing studies.

Table 6.1 Names and definitions for caIntegrator2 roles

Role Name	Role Definition
STUDY_INVESTIGATOR_ROLE	Assigning this role allows the user to search the study, save queries about the study and perform analyses.
PLATFORM_MANAGER_ROLE	Assigning this role allows the user to create and delete array platforms for the entire calIntegrator2 installation. Caution: Array platforms are shared by all users and studies in the calIntegrator2 installation. A user with this role can affect the platforms that are used by by all users and studies in the calIntegrator2 installation.

Table 6.1 Names and definitions for calIntegrator2 roles

8. If this user group is a group of study managers, then select STUDY_MANAGER_ROLE. If this user group is a group of study investigators, then select STUDY_INVESTIGATOR_ROLE. Click **Assign**.
9. Click **Update Association** at the bottom of the page. This completes the assigning of the user group to the protection group you chose.

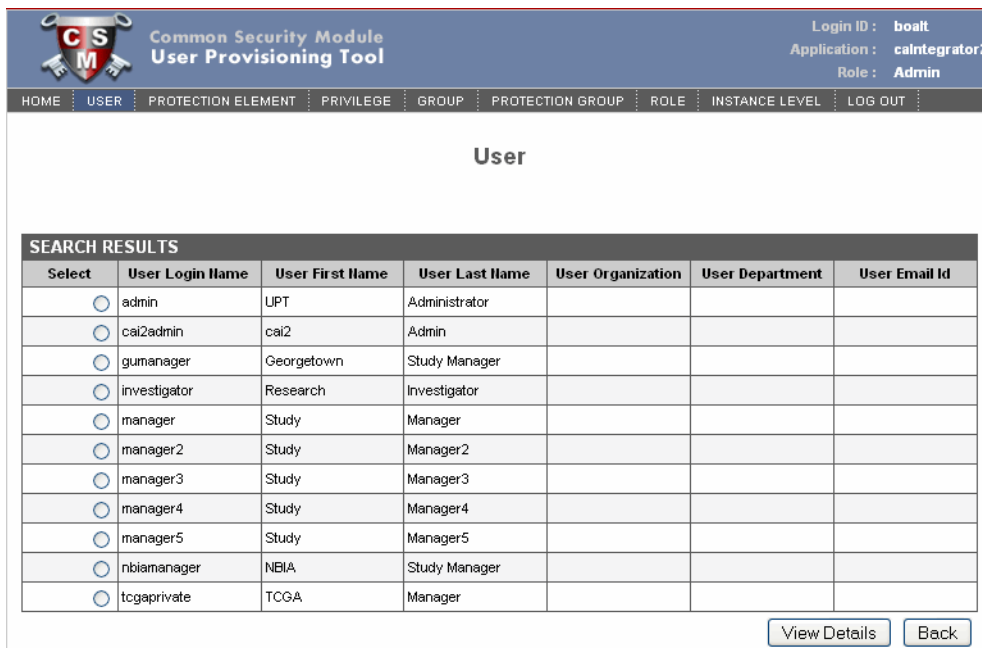
Note: If a **User** has the STUDY_MANAGER_ROLE role for more than one **Protection Group**, then any study that the **User** creates will be assign to each of those **Protection Groups**.

Adding a User to a User Group

To add a user to an existing user group, follow these steps:

1. Login to UPT as calIntegrator2 Admin.
2. Find the user that you want to assign to a user group. Click the **User** menu option, then click **Select an Existing User**.

- Enter the name of the user you are looking for and click **Search**. If you define no criteria, UPT returns a list of all calIntegrator2 users currently in the system (Figure 6.10).



The screenshot shows the 'Common Security Module User Provisioning Tool' interface. The top navigation bar includes links: HOME, USER, PROTECTION ELEMENT, PRIVILEGE, GROUP, PROTECTION GROUP, ROLE, INSTANCE LEVEL, and LOG OUT. The user's session information is displayed as: Login ID: boalt, Application: calIntegrator2, Role: Admin.

The main section is titled 'User' and displays 'SEARCH RESULTS' in a table. The table has the following columns: Select, User Login Name, User First Name, User Last Name, User Organization, User Department, and User Email Id.

Select	User Login Name	User First Name	User Last Name	User Organization	User Department	User Email Id
<input type="radio"/>	admin	UPT	Administrator			
<input type="radio"/>	cal2admin	cal2	Admin			
<input type="radio"/>	gumanager	Georgetown	Study Manager			
<input type="radio"/>	investigator	Research	Investigator			
<input type="radio"/>	manager	Study	Manager			
<input type="radio"/>	manager2	Study	Manager2			
<input type="radio"/>	manager3	Study	Manager3			
<input type="radio"/>	manager4	Study	Manager4			
<input type="radio"/>	manager5	Study	Manager5			
<input type="radio"/>	nbiamanager	NBIA	Study Manager			
<input type="radio"/>	tcgaprivate	TCGA	Manager			

At the bottom right of the table, there are two buttons: 'View Details' and 'Back'.

Figure 6.9 UPT page showing a list of calIntegrator2 users

- Select the radio button next to the name and click **View Details** (Figure 6.10).



The screenshot shows the 'Common Security Module User Provisioning Tool' interface for editing a user. The top navigation bar is the same as in Figure 6.9. The user's session information is: Login ID: boalt, Application: calIntegrator2, Role: Admin.

The main section is titled 'USER DETAILS' and contains a form for updating user information. The form includes the following fields:

- User Login Name: manager
- User First Name: Study
- User Last Name: Manager
- User Organization: (empty)
- User Department: (empty)
- User Title: (empty)
- User Phone Number: (empty)
- User Password: (masked with dots)
- Confirm Password: (masked with dots)
- User Email Id: (empty)
- User Start Date: (empty) (MM/DD/YYYY)
- User End Date: (empty) (MM/DD/YYYY)
- User Update Date: 09/24/2009 (MM/DD/YYYY)

At the bottom right of the form, there are three buttons: 'Update', 'Delete', and 'Back'. Below these buttons, there are four links: 'Associated Groups', 'Associated PE & Privileges', 'Associated PG & Roles', and 'Assign PG & Roles'.

Figure 6.10 UPT page showing details for a selected user

- Click the **Associated Groups** button at the bottom of the page. This opens the page where you can assign a user to a group (*Figure 6.11*).

User and Groups Association

SELECTED USER	
User Login Name	manager

Assign or Deassign multiple **Groups** for the selected **User**. To remove the complete association Deassign all the **Groups**.

AVAILABLE GROUPS
Study Managers Group 3
Study Managers Group 4
Study Managers Group 5
NCI Study Investigators
TCGA Study Managers

ASSIGNED GROUPS
Platform Manager Group
NCI Study Managers

Figure 6.11 UPT page for assigning a user to user groups

- Select the group(s) that you want the user to be in and click **Assign**.
- At the bottom of the page click **Update Association**. This completes the assigning of the user to the user group. Now the user will have access to any studies to which the user group has been given access.

Note: You can add a user to more than one user group. For example, a user could be assigned to “Columbia University Study Managers” as well as to “Columbia University Study Investigators”.

Changing a User Password

To change a password for a User, follow these steps:

- Confirm if the User is an LDAP user or not. If the User is an LDAP user, then this person must change their password using the NCI password change utility. Skip the rest of these steps.
- If the User is not an LDAP user, then continue with the rest of these steps.
- Login to UPT as calIntegrator2 Admin.
- Find the User that you want to change. Click the **User** menu option, then **Select an Existing User**.
- Enter the name of the user you are looking for and click **Search**. If you define not criteria, UPT returns a list of all calIntegrator users.
- Select the radio button next to the name and click **View Details**

- Replace the **User Password** and **Confirm Password** fields with the new password (Figure 6.12).



Common Security Module
User Provisioning Tool

Login ID : boalt
Application : calIntegrator2
Role : Admin

HOME USER PROTECTION ELEMENT PRIVILEGE GROUP PROTECTION GROUP ROLE INSTANCE LEVEL LOG OUT

Update the details of the displayed User. The **User Login Name** uniquely identifies the User and is a required field. The **User First Name** and **User Last Name** identifies the User. The **User Organization**, **User Department** and **User Title** provides his work details. The **User Phone Number** and **User Email Id** provides the contact details for the User. The **User Password** can be entered if the same schema is also going to be used for Authentication. The **User Start Date** and **User End Date** determine the period for which the User is a valid User. The **Update Date** indicates the date when this User's Details were last updated.

USER DETAILS	
* User Login Name	manager
* User First Name	Study
* User Last Name	Manager
User Organization	
User Department	
User Title	
User Phone Number	
User Password	••••••
Confirm Password	••••••
User Email Id	
User Start Date	(MM/DD/YYYY)
User End Date	(MM/DD/YYYY)
User Update Date	09/24/2009 (MM/DD/YYYY)

Update Delete Back

Associated Groups Associated PE & Privileges Associated PG & Roles Assign PG & Roles

Figure 6.12 UPT page where you can edit user details, such as a password

- At the bottom of the page click **Update**.

DATA IMPORT CONFIGURATIONS

This appendix describes configurations for importing data into a study.

Topics in this appendix include the following:

- [Subject Clinical Data Configuration](#) on this page
- [Delimited-Text Annotation Import](#) on page 101
- [Annotation Field Configuration](#) on page 102
- [Sample Data Configuration](#) on page 102
- [Genomic Data Configuration](#) on page 103
- [Imaging Data Configuration](#) on page 103

Subject Clinical Data Configuration

The following clinical data configuration information is collected:

- Clinical Data Source (delimited text)
- Protocol Id (of study to import)
- For delimited text, see [Delimited-Text Annotation Import](#). For subject annotation files, one field must be identified as the subject identifier.
- See [Annotation Field Configuration](#) for details on specification of visibility and browse configuration.

Delimited-Text Annotation Import

Delimited-text annotation files must be in standard comma-separated value format. The file must include a header line that specifies the name for each field. Each row of data must contain the same number of values as the header row. The file must include a column that will be designated as the identifier (e.g. subject identifier, sample identifier,

etc.) for each row. Optionally a file may include a single column that will be designated as a time-point indicator. Each row must contain a unique combination of identifier and time-point indicator of a unique identifier if no time-point is included. An example of the content of a file including a time-point is shown below.

```
"patientId", "timepoint", "bloodPressure", "weight"  
"1234", "T1", "120/80", "180"  
"1234", "T2", "125/80", "190"  
"5678", "T1", "120/85", "200"
```

After upload of the file, the Study Manager must indicate for each field:

- Field type (identifier, timepoint indicator, text, integer, float or boolean)
- After specification of these types, the file will be validated to ensure that the values are valid for the types selected and that the file conforms to the requirements given above.

Annotation Field Configuration

For each annotation field (regardless of the source), the Study Manager must specify the following information:

- Annotation semantics: each annotation field (whether associated with a subject, image series, image or sample) must either:
 - be associated with an existing annotation definition known to the system,
 - be associated to an existing CDE in caDSR or
 - have sufficient semantic metadata recorded so that the field may be submitted for registration as a CDE in caDSR.
- Field authorization: Each field must be either declared publicly visible or restricted to a list of groups. The default will be the visibility settings given at the study level.
- Whether the field is to be included in the results list for a given entity type (i.e. Subject, Sample, ImageSeries or Array Data) when browsing data (See Use Case: Browse Study Data).
- Whether the field is to be included in simple single-input searches when browsing data (See Use Case: Browse Study Data).

Sample Data Configuration

Sample data may be uploaded from either caArray 2 or from delimited-text import. Samples imported from caArray 2 may have annotation updated by use of the delimited-text import functionality if sample annotation is required. Import from caArray 2 requires specification of the following information:

- caArray server hostname
- caArray server JNDI port
- caArray username

- caArray password
- Either the experiment identifier (to import all samples in the experiment) or a file containing a comma-separated list of samples in the format “experiment identifier”, “sample name”.
- Mapping of samples to subjects. This may be specified by a comma-separated list in the format “subject identifier”, “sample identifier” or by a regular-expression based mapping formula.

When samples are imported via delimited-text import, the time-point is associated to the sample itself. This means that each sample may be associated with only one time-point (i.e. multiple time-points for the same sample are invalid).

Genomic Data Configuration

All genomic data (i.e. array data) is imported from caArray 2. First the Study Manager must specify sufficient information to map study samples to caArray 2 samples. If all samples were imported directly from caArray 2 as described in Special Requirement: Sample Data Configuration, no further information is required for this step. If samples were imported via delimited-text, the Study Manager must specify

- caArray server hostname
- caArray server JNDI port
- caArray username
- caArray password
- A mapping of calIntegrator2 sample identifiers to caArray 2 samples, specified as a comma-separated list in the format “calIntegrator2 sample identifier”, “caArray 2 experiment identifier”, “caArray 2 sample name”.

The system will enable the Study Manager to navigate easily to the selected caArray 2 instance.

Next, the system will indicate the available platforms and array data types available for the study samples. The Study Manager will indicate which platforms and data types to import and for each platform/data type combination will indicate:

- Whether to import the data
- The visibility of the data; either public or restricted to a set of groups. Low-level genotyping data (raw data and normalized) will always have restricted visibility.

Imaging Data Configuration

The following imaging data configuration information is collected:

- NBIA grid server hostname (defaults to NCICB instance)
- NBIA grid server port (defaults to NCICB instance port)
- Protocol Id

- Mapping of NBIA Patients to subjects imported from clinical data source. This may be specified by a comma-separated list in the format “subject identifier”, “NBIA patient identifier” or by a regular-expression based mapping formula.
- Which annotation fields to import from NBIA.
- The system will enable the Study Manager to navigate easily to the selected caArray 2 instance.

Additional annotation for either images or image series may be imported using the delimited-text import functionality.

INDEX

A

- account, requesting new user 6
- adding
 - clinical data 16
 - genomic data 24
 - imaging data 29
 - new user to user group 97
- annotation
 - assigning identifier 17
 - configuring field 102
 - importing delimited text 101
 - K-M plot for 58
 - searching for definitions 20
- Application Support i, 11
- assigning, annotation identifier 17

B

- box and whisker plot
 - interpretation 77
 - uses for 77

C

- caBio search 39, 48, 60, 66, 72
- caIntegrator2
 - logging in 8
 - logging out 10
 - requesting user account 6
 - using workspace 8
 - workspace 8
- caIntegrator2 User's Guide
 - introduction 1
 - organization 1
 - text conventions 2
- clinical study
 - adding data 16
 - data for 13
- columns, defining display 41, 42
- Comparative Marker Selection (CMS)
 - data analysis 81

configuring

- annotation fields for import 102
- copy number data 28
- genomic data for import 103
- imaging data for import 103
- sample data for import 102
- subject clinical data for import 101

control samples, uploading 27

copy number data, configuring 28

creating

- gene list 47
- K-M plot 58
- new user 92
- protection group 93
- study 14, 15
- user account 89

D

- data analysis, overview 57
- defining survival values 23
- delimited text annotation import 101
- DICOM, retrieving images 53

E

- editing
 - gene list 50
- editing a query 44
- exporting
 - data 55
 - query results 44

F

- fold change
 - control samples file 27
 - search 40

G

- gene expression, K-M plot for 60
- gene expression plot
 - description 57, 65, 75

- for clinical queries 71
- for genomic queries 69
- interpreting 75
- plot display, box & whisker 77
- plot display, log2 intensity 76
- plot display, mean 75
- plot display, median 76

gene list

- creating 47
- deleting 50
- editing 50
- search 40, 49, 62, 67, 73

GenePattern

- analyses, description 78
- analyses, in caIntegrator2 79
- analyses, modules 80
- CMS analysis 81
- GISTIC analysis 86
- PCA 83
- plot description 58

genomic data

- adding copy number data to 28
- adding to study 24
- configuring for import 103
- for study 13
- mapping to clinical data 26

GISTIC-based data analysis 86

H

hierarchy of objects, NBIA 54

I

imaging data

- adding to study 29
- configuring for import 103
- for study 14

importing delimited text annotations 101

K

Kaplan_Meier plot see K-M plot

K-M plot

- creating 58
- description 57, 58
- for annotations 58
- for gene expression 60
- for queries 63

L

logout link 10

M

managing

platforms 32

queries 43

study 31

user accounts 89

mapping genomic to clinical data 26

N

NBIA

forwarding imaging results to 52

viewing imaging results in 52

NCICB Application Support i, 11

O

objects in NBIA, hierarchy of 54

overview, chapters in guide 1

P

password, changing 99

patient

relationship to study, series, images 54

platforms, managing 32

plot

gene expression, description 57, 65

gene expression description 75

GenePattern description 58

K-M description 57

Principal Component Analysis

data analysis 83

protection group

creating 93

definition 90

Q

query

editing 44

exporting results 44

K-M plot for 63

managing 43

saving 43

query See also searching

R

registering new user 6

Results Type tab 41

S

sample data, configuring for import 102

saving query 43

searching

annotation definitions 20

caBIO 39, 48, 60, 66, 72

- fold change 40
- gene list 40, 49, 62, 67, 73
- overview 35
- study 36
- search results
 - browsing 46
 - exporting data 55
 - forwarding imaging results to NBIA 52
 - genomic data 46
 - imaging data 50
 - overview 45
 - retrieving DICOM images 53
 - viewing imaging data in NBIA 52
- Sorting tab 42
- study
 - adding clinical data 16
 - adding genomic data 24
 - adding imaging data 29
 - clinical data, description 13
 - configuring copy number data 28
 - creating 13, 14, 15
 - deploying 31
 - editing 31
 - genomic data, description 13
 - imaging data, description 14
 - managing 31
 - mapping genomic data to clinical 26
 - relationship patient, study, series, images 54
 - searching 36
 - uploading control samples to 27
- subject clinical data, configuring for import 101
- survival values, defining 23

T

- Technical Support i
- text conventions in user guide 2

U

UPT

- adding new user to user group 97
- assigning user group to protection group 94
- creating new user 90
- creating protection group 93
- creating user group 92
- description 89
- summary of steps 90

user

- adding to user group 97
- changing password 99
- creating new 90
- creating user access, summary 90
- definition 89

- user's manual conventions 2

- user account, new 89

