

# cPath Architectural Overview

Last Updated: March 1, 2004

<b>I. INTRODUCTION</b>	<b>2</b>
<b>II. TOOLS AND LIBRARIES</b>	<b>2</b>
<b>III. ARCHITECTURAL OVERVIEW</b>	<b>3</b>
<b>IV. CPATH DATABASE SCHEMA</b>	<b>4</b>
A. IMPORT TABLE	4
B. CORE CPATH TABLES	4
C. EXTERNAL DATABASE TABLES	5
<b>V. CPATH APPLICATION LOGIC</b>	<b>6</b>
A. PACKAGE STRUCTURE	6
B. FROM BROWSER REQUEST TO SYSTEM RESPONSE	7
C. STYLE SHEETS, JSP TEMPLATES AND JSP CUSTOM TAGS	9
D. WEB APPLICATION CONFIGURATION	10
<b>VI. DEVELOPMENT PRACTICES</b>	<b>10</b>
A. AUTOMATED BUILD PROCESS	10
B. CODING STANDARDS	11
C. UNIT TESTING	11
D. DOCUMENTATION	11

## I. Introduction

This document provides an architectural overview of the cPath database and web application. Information about the database schema, data access objects, servlet architecture, configuration, and unit testing is provided.

## II. Tools and Libraries

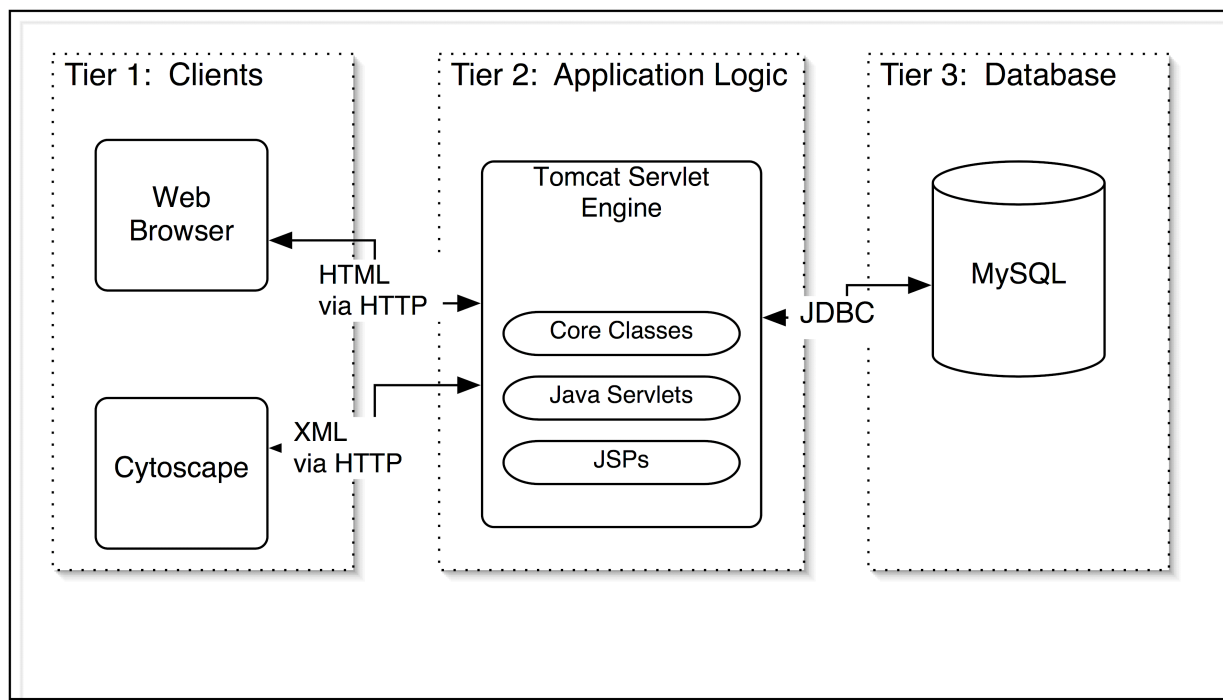
cPath is currently built with the following open source tools and libraries:

- **Apache Tomcat:** open source servlet/jsp engine. Information available at: <http://jakarta.apache.org/tomcat/>.
- **Ant:** used for the cPath build and deployment process. Information available at: <http://ant.apache.org/>.
- **Castor:** Java/XML Binding Framework. Information available at: <http://www.castor.org/>.
- **Checkstyle:** open source toolkit for checking coding conventions. cPath adheres to the Sun Java coding conventions. Information available at: <http://checkstyle.sourceforge.net/>.
- **CVS:** used for source control and revision tracking. Information available at: <http://www.cvshome.org/>.
- **Jakarta Database Connection Pool (DBCP):** provides a database connection pool for connecting to relational databases. Information available at: <http://jakarta.apache.org/commons/dbcp/index.html>
- **JDOM:** Java XML API. Information available at: <http://www.jdom.org>.
- **JUnit:** used for all cPath unit tests. Information available at: <http://www.junit.org>.
- **Log4j:** open source logging framework. Information available at: <http://logging.apache.org/log4j/docs/>.
- **Lucene:** open source text indexer; provides cPath full text search functionality. Information available at: <http://jakarta.apache.org/lucene/docs/index.html>.
- **Maven:** Maven is a “Java project management and project comprehension tool”. We currently use Maven to automatically generate the cPath Developer site every 24 hours.
- **MySQL:** Backend database is built in MySQL. Information available at: <http://www.mysql.com>.

- **Struts**: open source framework for building web applications. Information available at: <http://jakarta.apache.org/struts/>.
- **Xerces**: fast, validating XML Parser. Information available at: <http://xml.apache.org/xerces-j/>.

### III. Architectural Overview

cPath is a database driven web site that uses a traditional three-tier architecture (see Figure 1). The first tier consists of client applications, including web browsers and other software applications, such as Cytoscape. Web browser clients communicate with cPath by issuing http requests, and receiving HTML in response. Other client applications, such as Cytoscape communicate with cPath by issuing http requests, and receiving XML in response. The second tier consists of Java code running in the open source Tomcat servlet engine. This layer consists of core Java classes, Java Servlets, and Java Server Pages. The third tier consists of a backend database running mySQL. The mySQL database stores individual pieces of data, as well as XML document fragments, and we communicate with mySQL via JDBC.



**Figure 1:** cPath Three-Tier Architecture.

In the discussion that follows, we work backwards by first describing the database schema, and then moving our way to the application logic.

## IV. cPath Database Schema

All cPath data is stored in MySQL. The cPath database schema consists of eight relational tables. These tables are grouped into four sets:

- **import table:** used for importing new data into cPath.
- **core entity tables:** used to store core entities, such as interactors and interactions, and internal/external links.
- **external database tables:** used to store information about external databases, such as SWISS-PROT, NCBI, etc.
- **administrative tables:** used for logging and caching purposes.

Details on each set of tables is provided in the sections below.

### A. Import Table

The import table contains information about XML/text records, which are scheduled for import into cPath. For example, if an administrator wants to import a new PSI-MI file, the XML data is first loaded into the import table, where it is logged, and recorded. From here, the XML document is parsed and chopped into its constituent parts and loaded into the core cPath tables. The import table contains the following structure:

import	
IMPORT_ID	<i>Primary ID</i>
DESC	<i>Record Description</i>
DOC_BLOB	<i>Document Text/XML</i>
DOC_MD5	<i>MD5 Fingerprint</i>
STATUS	<i>Record Status</i>
CREATE_TIME	<i>Timestamp Created</i>
UPDATE_TIME	<i>Timestamp Updated</i>
EX_DB_ID	<i>Reference to External DB_ID</i>
LO_ID	<i>Reference to LinkedOut ID</i>

### B. Core cPath Tables

The cPath core consists of three tables: cpath, internal\_link, and external\_link. The cpath table contains core entities, such as interactors and interactions. Each entity record contains a short name, description and type. Type must be specified as either: PHYSICAL\_ENTITY or INTERACTION. Each record also contains an XML document fragment, written in PSI-MI format.<sup>1</sup> To obtain the full information for an entity, the cPath code must extract and parse the specified XML document fragment.

<sup>1</sup> Eventually, we will transition from PSI-MI to BioPax.

The `internal_link` table stores links between cPath records. For example, a cPath interaction record will specify bidirectional links between the interaction record and all its interactors. The `external_link` table records links to external databases, such as SWISS-PROT, NCBI, etc. Information about these external databases is provided in the next section.

The core cPath tables contain the following structure:

<b>cpath</b> Contains core cPath Entities	
CPATH_ID	<i>Primary ID</i>
NAME	<i>Entity Name</i>
DESC	<i>Entity Description</i>
TYPE	<i>Entity Type</i>
SPEC_TYPE	<i>Entity SubType</i>
NCBI_TAX_ID	<i>NCBI Taxonomy ID</i>
XML_CONTENT	<i>XML Document</i>
CREATE_TIME	<i>Timestamp Created</i>
UPDATE_TIME	<i>Timestamp Updated</i>

<b>external_link</b>	
EXTERNAL_LINK_ID	<i>Primary ID</i>
CPATH_ID	<i>Reference to CPATH_ID</i>
EXTERNAL_DB_ID	<i>Reference to External_DB_ID</i>
LINKED_TO_ID	<i>Linked to Identifier</i>
CREATE_TIME	<i>Timestamp Created</i>
UPDATE_TIME	<i>Timestamp Updated</i>

<b>internal_link</b>	
INTERNAL_LINK_ID	<i>Primary ID</i>
CPATH_ID_A	<i>Reference to First Entity</i>
CPATH_ID_B	<i>Reference to Second Entity</i>

### C. External Database Tables

The external database tables contain information about external databases and URL construction rules for connecting to those databases. The `external_db_cv` table contains controlled vocabulary terms, which match a specific database. For example, cPath recognizes: SWISS-PROT, SWP, and SWISSPROT, and all these CV terms point to a single SWISS-PROT record in the `external_database` table. The `external_database` table contains URLs for connecting to the external database. The token `%ID%` is replaced dynamically with the linked to ID. For example, here is a URL for connecting to SWISS-PROT:

<http://us.expasy.org/cgi-bin/niceprot.pl?%ID%>

The external database tables contain the following structure:

external_db	
EXTERNAL_DB_ID	Primary Key
NAME	Database Name
URL	URL Construction
DESC	Database Description
FIXED_CV_TERM	Foreign Key to CV_TERM
DBDB_ID	Link to Database of DBs
DBDB_URL	Link to Database of DBs
CREATE_TIME	Timestamp Created
UPDATE_TIME	Timestamp Updated

external_db_cv	
CV_ID	Primary Key
EXTERNAL_DB_ID	Foreign Key to external_db
CV_TERM	Term

## V. cPath Application Logic

### A. Package Structure

Application logic for cPath is written in Java, and is specified in the package: `org.mskcc.pathdb`. It contains the following package structure:

Package	Description
<b>action</b>	Struts Action Classes. Each action class corresponds to a user-initiated action.
<b>controller</b>	Classes for processing and validating web service API requests.
<b>form</b>	All Struts Action Forms. Each form corresponds to an HTML form.
<b>lucene</b>	Classes for connecting to the Lucene Full Text Search engine.
<b>model</b>	JavaBean objects which encapsulate cPath records, such as ImportRecord, CPathRecord, etc.
<b>servlet</b>	All cPath Servlets.
<b>sql</b>	All Database Access code.
<b>taglib</b>	All Custom JSP Tags.
<b>task</b>	Long-term Tasks, which require multi-threading

	execution.
<b>test</b>	All JUnit Unit Tests.
<b>tool</b>	All Command Line Utilities and Programs.
<b>util</b>	Misc. Utility classes, such as an XML validator, Cross-Site scripting filter, etc.
<b>xdebug</b>	Live Debugging/Diagnostics Facility.
<b>xmlrpc</b>	XML-RPC Services for uploading new data into cPath.

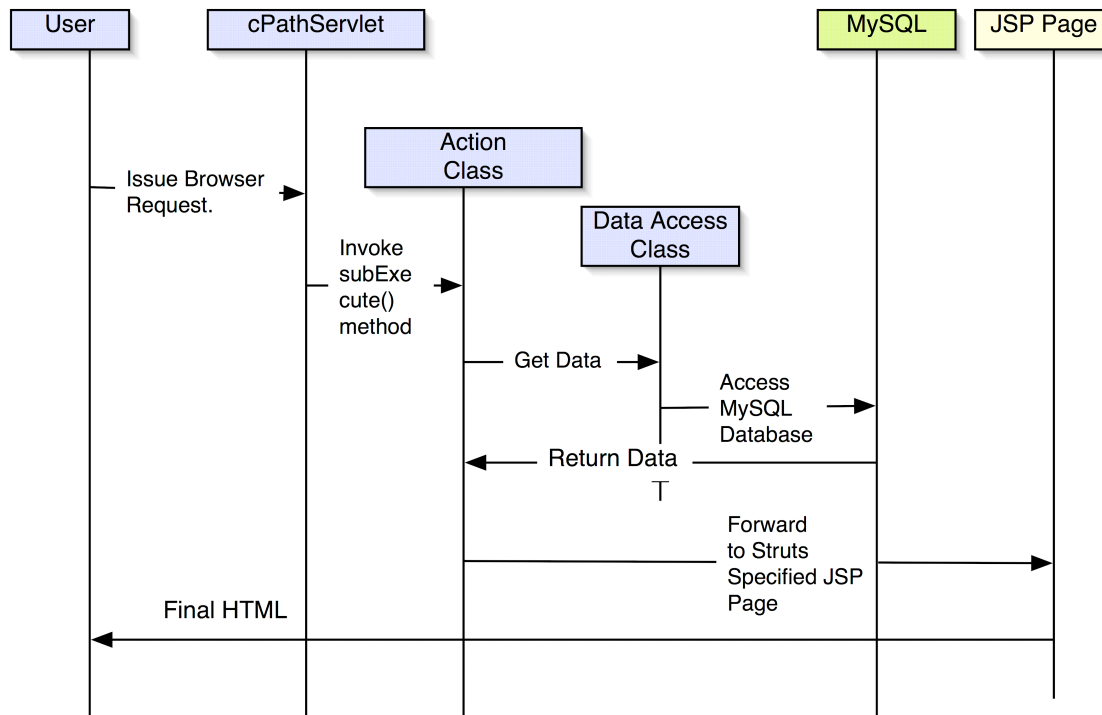
### ***B. From Browser Request to System Response***

Each time a web browser makes a request to cPath, the application logic layer handles the processing of the request, and the generation of dynamic HTML/XML pages. To aid in the processing of web requests, cPath utilizes the open source Struts Framework. The Struts framework provides a number of built-in advantages, including: clean separation of logic and presentation, form validation, centralization of request handling, and centralization of exception handling.

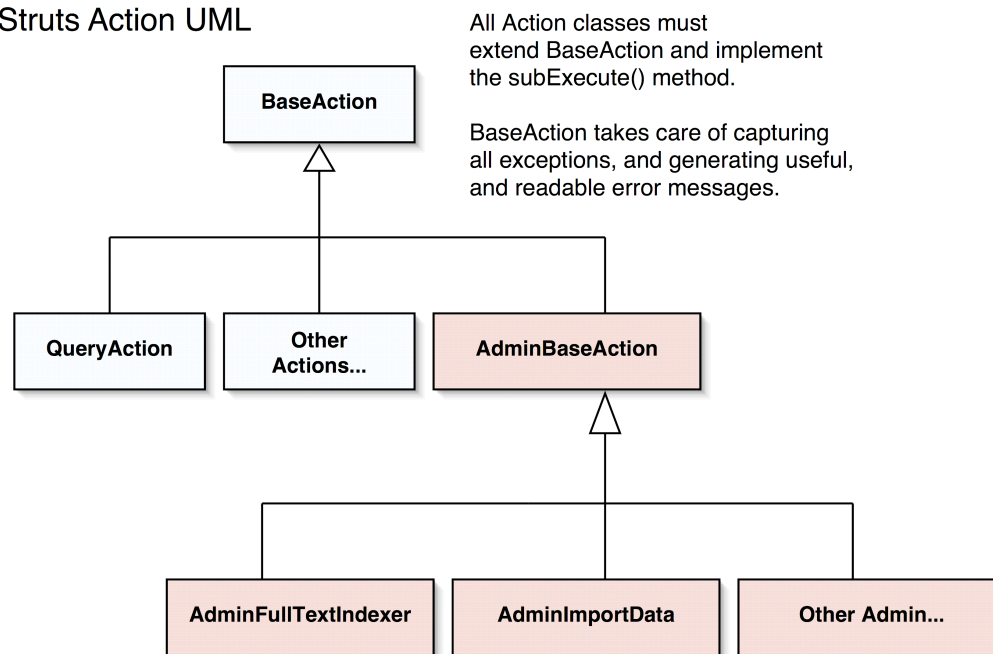
Flow through the application layer proceeds as follows (see Figure 2):

- All user requests go through the central **cPathServlet** class. This class provides a central spot for all cPath initialization and configuration.
- Based on the URL requested, Struts will check the Struts.xml configuration file, and determine which action class is invoked. Action classes represent user request actions, such as “run a full text search”, or “get all interactions for yeast”.
- All cPath action classes inherit from the **BaseAction** class, and all password protected administrator classes inherit from the **AdminBaseAction** class (see Figure 3). The base action classes provide a common framework for capturing exceptions, gathering real-time web diagnostics, and protecting specific resources via HTTP authentication.
- Some action classes will require access to the MySQL database. In that case, the action class will make calls to a DAO Object. A DAO, or Data Access Object encapsulates all database access code, and also takes care of handling connections to the database connection pool. There is one DAO object for each cPath table. For example, the import table has an ImportDao class.
- All final HTML construction is completed in Java Server Pages, and custom JSP tags.

## Struts Sequence Diagram

**Figure 2:** Data flow within the Application Logic Layer.

## Struts Action UML

**Figure 3:** UML Diagram of the Struts Action Classes. Each Action class corresponds to a user-requested action.



### C. Style Sheets, JSP Templates and JSP Custom Tags

cPath is built using open source CSS style sheets created at <http://style.tigris.org>. Tigris provides a cross-browser compatible set of style sheets and open source icons. This enables cPath to have a very professional look, with minimal effort. It also enables us to easily provide a “printer-friendly” version of all cPath pages.

All presentation is done via Java Server Pages (JSPs) and custom JSP tags. The WEB-INF/jsp/global directory contains a number of reusable JSP components which are used throughout the site. These include the following:

JSP Component	Description
<b>header.jsp</b>	Global header component, used on all cPath pages.
<b>footer.jsp</b>	Global footer component, used on all cPath pages.
<b>tabs.jsp</b>	Global tabs component, seen at the top of all cPath pages.
<b>navbar.jsp</b>	Global navigation bar component, seen in left column of all cPath pages.
<b>xdebug.jsp</b>	Global component used to output real-time web diagnostics.

The custom JSP tags provide reusable tags, which can be inserted into any cPath page. All cPath custom tags are defined in: WEB-INF/tag-lib/cbio-taglib. A summary of all custom tags is provided below:

JSP Custom Tag	Description	Example Usage
<b>DiagnosticsTable</b>	Outputs internal tests used in the Admin cPath self-test page.	<code>&lt;cbio:diagnosticsTable /&gt;</code>
<b>ErrorMessage</b>	Centralizes all display of error messages to the user.	<code>&lt;cbio:errorMessage throwable=" "%= exception %&gt;" /&gt;</code>
<b>HtmlTable</b>	Base Class for all Custom JSP Tags which create HTML tables.	Not used directly.
<b>InteractionTable</b>	Displays a list of matching interactions.	<code>&lt;cbio:interactionTable interactions=" "%= interactions %&gt;" protocolRequest=" "%= protocolRequest %&gt;" /&gt;</code>
<b>InteractorTable</b>	Displays all information about a specific interactor.	<code>&lt;cbio:interactorTable /&gt;</code>
<b>SearchResultsTable</b>	Displays results from a	<code>&lt;cbio:searchResultsTable</code>

	full-text search.	protocolRequest= "<%= protocolRequest %>"/>
<b>TaskTable</b>	Displays all admin tasks (in progress and completed tasks are shown.)	<cbio:taskTable/>

## D. Web Application Configuration

cPath configuration is centralized in three files.

File	Description
<b>WEB-INF/web.xml</b>	The main configuration file for the cPath web application. Here, you can modify the Struts debugging level, the user/name password for connecting to cPath, and the admin user/name password for accessing administrative functionality.
<b>WEB-INF/struts-config.xml</b>	The main configuration file for the Struts Framework. This file specifies all the action URLs, HTML form objects, and JSP pages.
<b>config/cpathResources.properties</b>	Defines a number of String constants used throughout the cPath site. All user error messages are defined here.

## VI. Development Practices

### A. Automated Build Process

The cPath build process is automated with Ant. A number of Ant targets are defined:

Target	Description
<b>Main Development Targets</b>	
boot	Performs a complete and clean build.
check	Run CheckStyle on all source code.
clean	Deletes all build files and starts fresh.
compile	Compile all Java sources.

<b>Test Targets</b>	
boot	Loads Bootstrap data into MySQL. Required for running JUnit Tests.
test	Runs JUnit Test Suite.
<b>Local Tomcat Development Targets</b>	
install	Installs cPath to local development servlet container.
list	Lists installed applications on servlet container.
reload	Reloads cPath application on servlet container.
remove	Removes cpath application on servlet container.
<b>Deployment / WAR Targets</b>	
deploy_war	Deploy WAR to cBIO Production server.
war	Creates binary WAR distribution.

## ***B. Coding Standards***

cPath adheres to the Sun Java coding standards. These are enforced via the ant checkstyle target. To see if your code adheres to the standards, run “ant check”.

## ***C. Unit Testing***

cPath uses the JUnit testing framework. All unit tests are organized into suites, and all suites are organized into `org.mskcc.pathdb.test.AllTest`. To run all cPath tests, run “ant test”. To create a new test, create the new test, and add it to the correct package suite.

## ***D. Documentation***

cPath uses the open source Maven toolkit for automatically generating the cPath developer site every 24 hours. All documentation is stored in cpath/xdocs. Add new documents to this directory.