2008

# Introducing legacy program scripting to molecular biology toolkit (MBT)

Todd Newell

# Introducing Legacy Program Scripting to Molecular Biology Toolkit (MBT)

Graduate Project Defense for MS Computer Science Degree
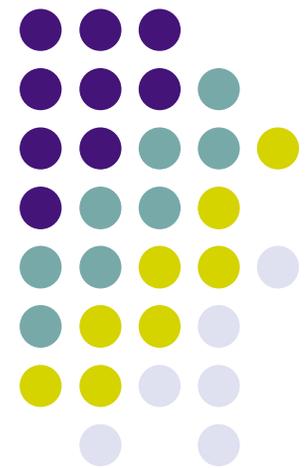
Todd A. Newell

October 28, 2005

# Table of Contents

# Abstract

- λ Ever-changing landscape of molecular visualization requires a common thread
- λ Introduction of powerful visualization language, MBT, gives users a comprehensive collection of libraries with the freedom to develop unique applications
- λ Lacks scripting capabilities, however, something that older, popular languages such as Chime and RasMol have
- λ MBT lacks handholds for widespread acceptance without scripting
- λ This project will use JavaCC to parse legacy commands from Chime and RasMol into single, simple commands usable by MBT – scripting
- λ Bridging to a user-centric and friendly method of use for end-users not entrenched in software development

# Introduction

- λ Description of Problem
    - λ Visualization software is always changing and it is tough for users to keep up
    - λ These users are often chemists, biologist, and physicists instead of computer scientists
    - λ Desire to change with the times comes at a price - users have to learn new application with potentially long learning curve
    - λ Methods of use should be transferable
        - λ Scripting
        - λ .pdb Files

# Chime, RasMol, and MBT

- λ Two popular, older languages have scripting abilities and are widely used
    - λ Chime: free but incomplete version supplied, not open source.
    - λ Chime: popular tool leads to many scripts being created
    - λ RasMol: older and unsupported
    - λ RasMol: free, open source, command list is subset of Chime's
    - λ RasMol: also large number of scripts existing for application
- λ New toolkit, MBT, incorporates many of the benefits from Chime and RasMol, however not all of them
    - λ Open source, free
    - λ Toolkit library that needs applications to render images and interact with molecule
    - λ Uses .pdb files to input compound data
        - λ Chemical information listed for each atom
        - λ .pdb files are the same wherever used

# Sample .pdb File

```
ATOM    1 OH  OSP   1     -0.332  0.440 -0.016 1.00 0.00
ATOM    2 1HH  OSP   1      0.596  0.456  0.228 1.00 0.00
ATOM    3 2HH  OSP   1     -0.596 -0.456 -0.228 1.00 0.00
TER     4     OSP   1
END
```
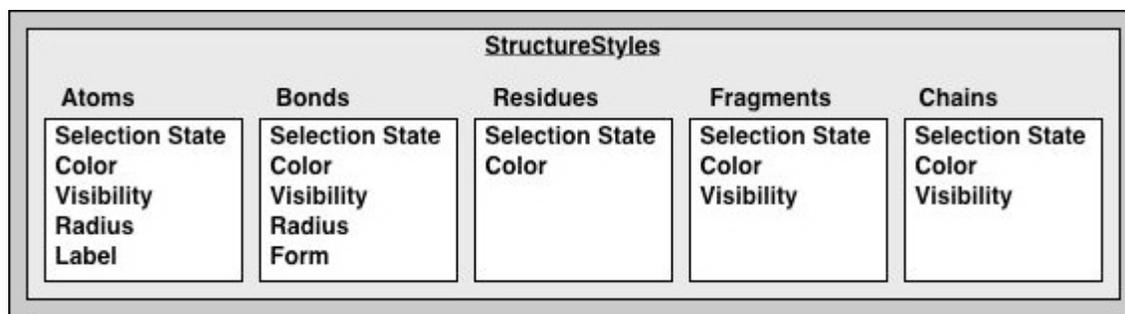
λ  Type of Component, Atom Serial Number, Atom Name, Coordinates of Atom

λ  Conforms to standards – identical each time used

λ  Read by the viewer to manufacture initial view of the molecule.

# MBT

λ Requires additional application to run

λ MBT's libraries arranged in hierarchy which allows only a few classes to affect molecular components

λ StructureStyles.java have hooks into atoms, bonds, residues, fragments, and chains

| StructureStyles | | | | |
|---|---|---|---|---|
| **Atoms** | **Bonds** | **Residues** | **Fragments** | **Chains** |
| Selection State<br>Color<br>Visibility<br>Radius<br>Label | Selection State<br>Color<br>Visibility<br>Radius<br>Form | Selection State<br>Color | Selection State<br>Color<br>Visibility | Selection State<br>Color<br>Visibility |

# Goals

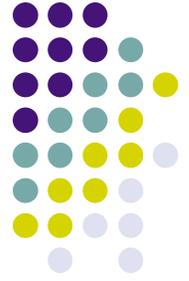λ Last component needed is a parser to evaluate input languages – JavaCC

λ Allows us to bridge from older languages to the new, via scripting

λ Goals

  λ "This project was designed to create a parsing language that accepts scripting commands from two ancestor programs, checking them for correctness, and eventually using the commands to affect the chemical components in the modified molecular visualization toolkit"

# Scripting

λ What is scripting?

  λ Classic method of execution – there is no way to change variables or run methods after compiling without a portal into the program

  λ Graphical programs have pull-down menus

  λ Scripting allows users to make changes without the knowledge of target applications

  λ Easy to remember, accessible commands that run more complex pieces of code

# Scripting Example

λ Scripting command to select all components in a molecule: "select all"

λ MBT method to select all components:
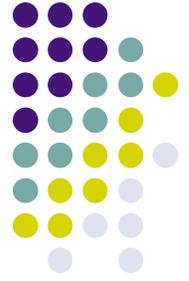
```
public void selectAll() {
    selection.clearAll();
    setAtomSelection( selection.getMin(),
            selection.getMax(), true );
    bondSelection.clearAll();
    setBondSelection( bondSelection.getMin(),
            bondSelection.getMax(), true ); }
```

# Importance of Scripting

λ   Chime and RasMol have scripting ability, while MBT does not

λ   It was the void we decided to fill – commands that were used to run Chime and RasMol could be made to execute MBT methods

λ   In this way, users without a comprehensive background in programming can jump to a toolkit like MBT by reusing familiar commands

λ   Gives the user a simple method of interaction once application has been compiled and executed

# Using MBT for Visualization

λ   MBT is a library and not a free-standing application

λ   Hierarchies are arranged to be intuitive

λ   Users can create applications to fit specific needs

λ   MBT can be used to render and display molecules either 2D or 3D – using Java3D

λ   Also allows non-graphical manipulation of molecules – loading the same .pdb files lets the user create a virtual representation and change parts of the components without ever seeing the molecule

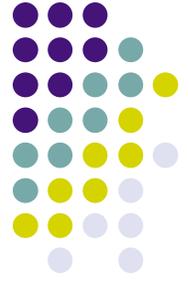λ   Scripting can act on either type of representation

# Architecture

- λ JavaCC
  - λ Designed to compile programming languages
  - λ Can analyze any left-decent grammar
  - λ Tokens are the lowest common denominator
  - λ Consume tokens, process them, and analyze them for correctness
  - λ Tokens act as reserve words that disallow their use elsewhere in the program
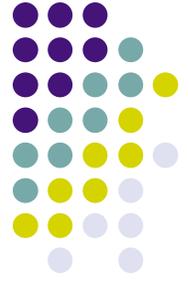
# Architecture - JavaCC

λ JavaCC

  λ User defines the grammar which JavaCC program uses to analyze

  λ Rules are added as needed – making the grammar as flexible as needed

  λ Only requirement: grammar is non-ambiguous so no infinite loops

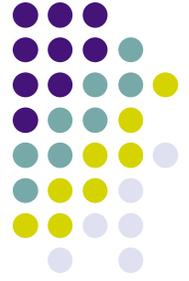# Architecture - JavaCC

λ JavaCC

  λ Tokens can be anything – regular expressions, keywords, or literals

  λ 2 distinct functions

    λ Token Manager – separates commands into tokens (as defined by grammar – whatever size)

    λ Parser – examines the rules to try to match tokens, in order, to the functions defined in the program

# Architecture - JavaCC

- λ JavaCC
  - λ LOOKAHEAD – used to notify the program of a way to avoid ambiguity
  - λ By specifying number of tokens ahead of consumed token to check, parsing path directed around spots with similar command structures
    - λ E.g.: LOOKAHEAD(2) color() – two or more functions begin with the same token < COLOR >.  Alert parser to check 2 tokens ahead of consumed token, instead of 1
  - λ Most of the program can operate in a processor efficient manner, only slowing down when meets with stumbling block

# Architecture – JavaCC / ChimeParser.jj

λ   PARSER_BEGIN( ChimeParser )
    /* IMPORTS – entire MBT library, each package, and any needed Java imports */
        public class ChimeParser {
                /* global variable declaration – must be static */
                private static Explorer explorer = new Explorer();
                public static void main( String args[] ) throws ParseException {
                        …
                        explorer.initialize();
                        try {
                                parser.cmdList();
                        } catch ( ParseException e ) {
                                System.out.println( e.getMessage() );
                                System.out.println
                                ( "ChimeParser: Encountered errors during parse.");
                        }
                }
        }
PARSER_END( ChimeParser )

# Architecture – JavaCC / ChimeParser.jj

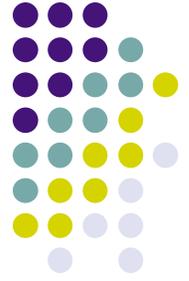λ   TOKEN : /* command list */ {
    < ANIM : "anim" | "animation" >
     | < DIRECTION : "direction" > }

λ   void anim() :
    {
            String animStr;
            Token animTok, lastTok;
    }
    {
            animTok = < ANIM >
            animStr = animParser()
            {
                        System.out.println( animTok.image + " "  + animStr );
            }
            [ lastTok = < SEMICOLON >
            {
                        insertSemiColon();
            } ]
    }

# Architecture - MBT

- λ MBT - packages
  - λ *edu.sdsc.mbt:* This package provides classes which define the core data storage containers for the MBT toolkit. **StructureMap**
  - λ *edu.sdsc.mbt.filters:* This package provides classes which enable filtering or subsetting of a structure's constituent components.
  - λ *edu.sdsc.mbt.gui:* This package provides classes which implement graphical user interface (GUI) component classes for MBT.
  - λ *edu.sdsc.mbt.io:* This package provides classes which enable molecular biology data sets (such as protein structures, sequence data, etc) to be loaded into an MBT application as a Structure object.
  - λ *edu.sdsc.mbt.util:* This package provides classes which implement extra functionality on top of the core MBT classes (that is, capabilities that are not absolutely required by, or would otherwise overcomplicate, the core classes).
  - λ *edu.sdsc.mbt.viewables:* This package provides classes which enable molecular biology data to be represented as visible/renderable viewable objects, plus, this package defines a top-level StructureDocument object to encapsulate the complete state of these objects and properties.
  - λ *edu.sdsc.mbt.viewers:* This package provides classes which enable data to be graphically visualized or to be processed by other means based upon coordinated events (that is, any MBT component that wishes to observe data and then respond and interact to changes in that data).
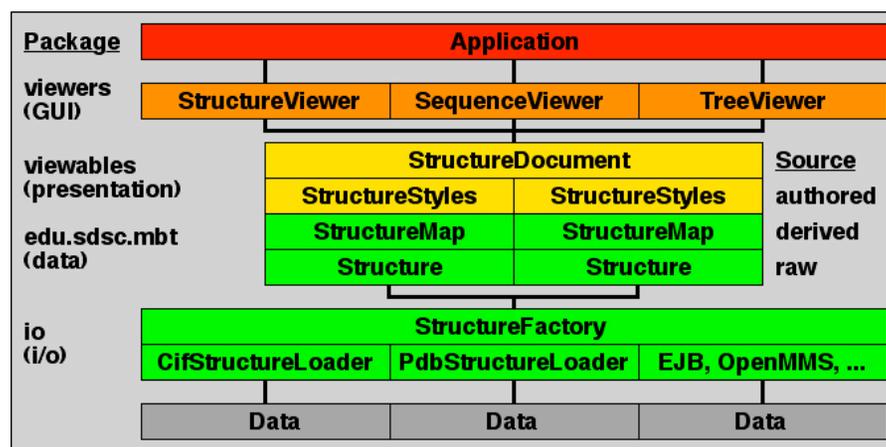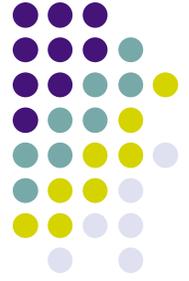
# Architecture - MBT

λ  MBT
- λ  Classes StructureStyles and StructureMap used the most when developing ChimeParser.jj
- λ  Spans the hierarchy – can make direct calls to the basic components as well as accessing collections
- λ  StructureStyles – contains most of the methods to make changes such as color, selection, and visibility
  - λ  Method calls directly from Explorer.java : structureStyles.setBondSelection ( index, index, true );
- λ  StructureMaps – edu.sdsc.mbt package and arranges the components after a Structure object is created
  - λ  Accessed whenever need to isolate a component, like the atoms, and check each one in turn:
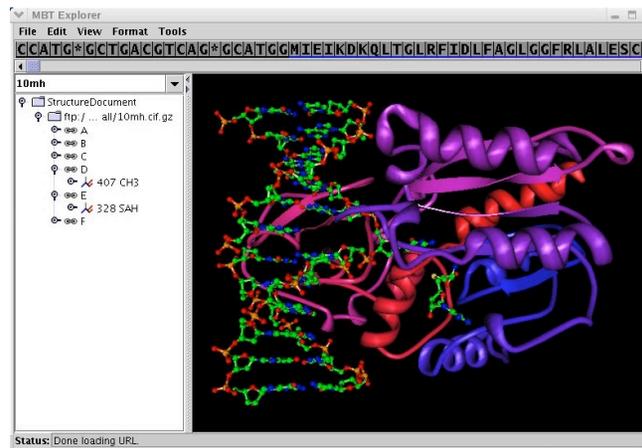
```
for( int index = 0; index < bondCount; index++ ) {
        Bond bond = structureMap.getBond( index );
}
```
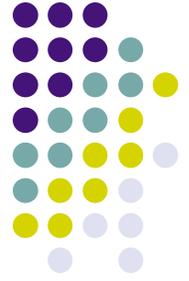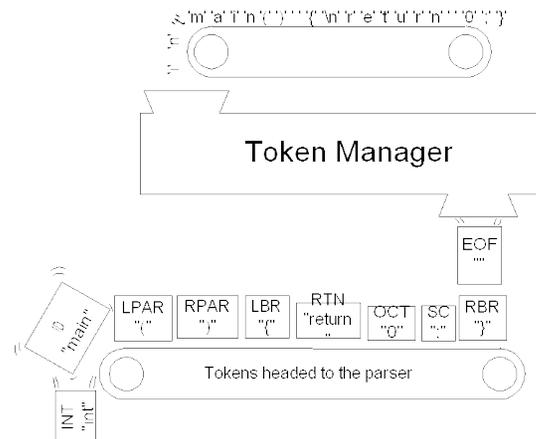
# Architecture - MBT

λ Explorer.java
  λ The viewer chosen to render the molecule
  λ Not fully implemented, but was designed to make basic changes to the molecule
  λ Other viewers developed to access specific type of component or chemical interaction (i.e.: hydrogen bonds or visualizing single nucleotide polymorphisms (SNP) )
  λ Additional methods added to support ChimeParser.jj
    λ hbonds( boolean selected)
    λ labels( boolean selected )
    λ selectAtoms( boolean selected)
    λ setBonds( boolean selected) / setAllBonds( boolean selected)
    λ setAtomColorByRGB( float[] colors ) / setBondColorByRGB( float[] colors )
    λ updateAtomColors( StructureStyles ss, boolean selected ) / updateBondColors( StructureStyles ss, boolean selected )

# Execution

- λ ChimeParser.jj – first point of execution, responsible for parsing input commands
  - λ Running: % java ChimeParser
  - λ Await input command to begin token manager
  - λ i.e.: color list {listname} [color] {list of colors} broken down to defined tokens: <color>, <list>, <allowable string regular expression>, [<color>], (<valid color>)+

# Execution

- λ ChimeParser.jj
  - λ Parser was then applied to token stream to check for correctness
  - λ Initializes an Explorer object when executes the main method
  - λ Drops into the command functions, checking each for possible entrance
    - λ Once finds matching first token, enters into that command's parsing tree (left-decent), comparing to operating grammar

# Execution

λ Program execution
- λ % javacc ChimeParser.jj
  - λ Creates 7 .java files:
    - ♣ ChimeParser.java – transformation of grammar definitions to Java code
    - ♣ ChimeParserConstants.java
    - ♣ ChimeParserTokenManager.java
    - ♣ ParseException.java
    - ♣ SimpleCharacterStream.java – because using default options with ChimeParser.jj
    - ♣ Token.java
    - ♣ TokenMgrError.java
- λ % javac *.java – compiles all .java in directory (including Explorer.java )
- λ % java ChimeParser ['<'] [<command filename>]
  - λ '<': force commands from file to be taken individually from Standard.in
  - λ <command filename>: file user wishes to parse
  - λ with no "load" command, commands will echo back to user if correct
- λ % <commands user wishes to parse>

# Future Work

- Implement full Chime and RasMol libraries
  - Improve the depth of the chemistry
  - Cross-department work to redesign effect of commands on MBT
- Expanding supported input languages
  - PyMol – amazing graphics, interface is lacking
  - Use ChimeParser.jj as backend for any visualization language that conforms to grammar rules (non-ambiguous)
- Artificial Intelligence
  - Protein active site anticipation and extrapolation
  - Study protein interactions with hypothetical molecules
  - Drug research: i.e.: G-protein coupled receptors
- Add scripting ability directly to MBT
  - Eliminate ChimeParser.jj and JavaCC – tokenless?
- Additional implementation techniques
  - Web-based
  - Databases

# References

λ   1. MDL Chime
http://www.mdl.com

λ   2. RasMol Home Page
http://www.umass.edu/microbio/rasmol/

λ   3. Norvell, Theodore S. (2004). *The JavaCC Faq*. Retrieved from Memorial University of Newfoundland Web site: http://www.engr.mun.ca/~theo/JavaCC-FAQ/javacc-faq.pdf

λ   4. JavaCC Tutorial
http://www.cs.rit.edu/~jmg/javacc/

λ   5. The Molecular Biology Toolkit. J.L. Moreland, A.Gramada, O.V. Buzko and P.E. Bourne 2005 The Molecular Biology Toolkit (MBT): A Modular Platform for Developing Molecular Visualization Applications. BMC Bioinformatics, 6:21. Funded by Grant NIGMS 1P01GM63208. Web site: http://mbt.sdsc.edu/

λ   6. Kenekin, T. (2005). New Bull's-Eyes For Drugs, *Scientific American*. October 2005, pp. 50-57.