# FOR CONSTRAINT-BASED DESIGN OF SYNTHETIC BIOLOGICAL SYSTEMS

Ernst Oberortner and Douglas Densmore

August 22, 2014

Boston University

Department of Electrical and Computer Engineering

Technical Report No. ECE-2014-NN

### BOSTON UNIVERSITY

## EUGENE — A DOMAIN-SPECIFIC LANGUAGE FOR CONSTRAINT-BASED DESIGN OF SYNTHETIC BIOLOGICAL SYSTEMS

Ernst Oberortner and Douglas Densmore



Boston University

Department of Electrical and Computer Engineering
8 Saint Mary's Street
Boston, MA 02215

www.bu.edu/ece

August 22, 2014

Technical Report No. ECE-2014-NN

Here comes a special preface page. For example, if the report is restricted, then a suitable note can be included. This page can also be used to indicate to whom the document is dedicated, etc.

#### Summary

Summary of the report. Maximum 1 page.

### Contents

1	Inti	roduction	
2	Lan	nguage Basics	
	2.1	Comments	
	2.2	Primitive Types	
	2.3	Primitive Values / Atoms	
	2.4	Identifiers	
	2.5	Declaration of Variables	
	2.6	Assignment Statements	
	2.7	Error Reporting	4
	2.8	Reserved Words	
3	Dec	clarative Language Features	ļ
	3.1	Declaration of Biological Facts	
		3.1.1 Properties of Biological Facts	
		3.1.2 Types	
		3.1.3 Instantiations	
		3.1.4 Composite Biological Facts	
		3.1.5 Relations	
	3.2	Design Templates	
	3.3	Rules and Constraints	
		3.3.1 Structural Constraints	
		3.3.2 Relational Expressions	
	3.4	Containers for Biological Facts and Rules	1
		3.4.1 Collections	1
		3.4.2 Arrays	1
4	Dat	ta Exchange Capabilities	1
	4.1	Inclusion of Eugene Scripts	1
	4.2	Data Import of Eugene Containers	1
	4.3	Data Import from the iGEM Partsregistry	1
	4.4	GenBank Data Import and Export	1
	4.5	SBOL Data Import and Export	1
	4.6	SBOL Visual Compliant Design Visualization	1

5	Imp	perative Language Features
	5.1	Conditional Branches
	5.2	Iterators
	5.3	Function Prototypes
	5.4	Loops
	5.5	Built-in Functions
		5.5.1 Design Verification
		5.5.2 Design Enumeration
6	Alg	orithms
	6.1	Enumeration of rule-compliant designs
		Verification of designs and rules

### List of Figures

### List of Tables

#### Introduction

#### Language characteristics

- $\bullet$  rule-based design paradigm  $\Rightarrow$  the specification of synthetic biological designs based on biological constraints and rules
- multi-paradigm language ⇒ combines declarative and imperative language features
- declarative features  $\Rightarrow$  the user specifies WHAT the design problem is based on rules and constraints; invokes built-in functions
- imperative features  $\Rightarrow HOW$  to (1) process the results returned by Eugene's declarative features and/or (2) build the WHAT question to ask Eugene
- interpreted language  $\Rightarrow$  a given Eugene script is being executed "on-the-fly". That is every statement is executed immediately. The Eugene parser is a one-phase parser, which requires that variables, biological facts, and functions must be declared before being used and/or referred to them.
- strongly typed language  $\Rightarrow$  three primitive types (num, txt, bool), two types of arrays (num and txt), language model for typing the biological facts (primitive and composite).

In this manual, we describe the syntax of the Eugene language using the Extended-Backus-Naur Form (EBNF).<sup>1</sup>

### Language Basics

Eugene's syntax is case sensitive. Language statements are separated by the colon operator (;).

#### 2.1 Comments

Eugene supports two types of comments: single-line comments and multi-line comments.

```
// single line comment
/*
this is a
multi-line
comment
*/
```

#### 2.2 Primitive Types

Eugene supports five primitive types: which can be either a real number (**num**), a string of characters (**txt**), an array of floating point numbers (**num**[]), an array of strings (**txt**[]), or a boolean (**bool**).

```
primitive-type ::= num | txt | num[] | txt[] | bool
```

#### 2.3 Primitive Values / Atoms

Based on the supported primitive types, Eugene supports strings (i.e. character sequences enclosed in double quotas) and floating point numbers. In Eugene, identifiers are unique. That is, no two variables can have the same identifier. Identifiers are also case-sensitive.

```
string ::= " [^"] "
number ::= ( [0-9] )* (. ( [0-9] )* )?
```

#### 2.4 Identifiers

In Eugene, identifiers must be unique. An identifier is a sequence of characters that is constraint as follows:

```
id := [a-zA-Z] ( [a-zA-Z0-9_] )*
```

In a Eugene script, the identifiers must be defined before their usage and/or referring to them.

#### 2.5 Declaration of Variables

Variables have a unique identifier (see Section 2.4) and a primitive type (see Section 2.2). Depending on the variable's type, the variable can take only values of the specified type (see Section 2.3). Values can either be constants or expressions. Expressions are calculated during the execution of a Eugene script. Values can be assigned to variables either at the point of the variable's declaration or at later points of execution.

```
\langle variable\text{-}declaration \rangle ::= \langle primitive\text{-}type \rangle \langle id \rangle ;
\langle primitive\text{-}type \rangle ::= \underline{\mathbf{txt}} \mid \underline{\mathbf{txt}}[] \mid \underline{\mathbf{num}} \mid \underline{\mathbf{num}}[] \mid \underline{\mathbf{bool}}
\langle id \rangle ::= [a\text{-}zA\text{-}Z] ( [a\text{-}zA\text{-}Z0\text{-}9\_] )^*
```

#### **Examples:**

- declaration of a variable of type num with the identifier n: num n;
- declaration of a string variable with the identifier s: txt s;

#### 2.6 Assignment Statements

Assignment statements assign values to variables and are interpreted from right to left, i.e., the result of the expression specified on the right-hand-side (RHS) is assigned to the variable specified on the left-hand-side (LHS). The LHS and the RHS are separated by the equals operator (=).

```
\langle assignment \rangle ::= \langle lhs \rangle \equiv \langle expression \rangle ;
\langle lhs \rangle ::= \langle variable\text{-}declaration \rangle \mid \langle id \rangle
\langle expression \rangle ::= \langle multiplication \rangle \mid \langle division \rangle
\langle multiplication \rangle ::= \langle addition \rangle \mid \langle subtraction \rangle \ (\ \underline{*} \ \langle expression \rangle \ )?
\langle division \rangle ::= \langle addition \rangle \mid \langle subtraction \rangle \ (\ \underline{/} \ \langle expression \rangle \ )?
\langle addition \rangle ::= \langle atom \rangle \ (\ \underline{+} \ \langle expression \rangle \ )?
\langle subtraction \rangle ::= \langle atom \rangle \ (\ \underline{-} \ \langle expression \rangle \ )?
\langle atom \rangle ::= \langle id \rangle \mid \langle number \rangle \mid \langle string \rangle \mid (\ \langle expression \rangle \ )
```

#### 2.7 Error Reporting

In Eugene v2.0 we differentiate among various types of errors ranging from syntax errors, over unknown identifiers and incompatible types, to inconsistent rules. Any type of error is thrown as a Java IllegalArgumentException and reported to the output console. Also, the interpretation of a defective Eugene script stops after printing the error to the console. Error messages have the following format:

```
@Error
Line <line-no> Position <position-in-line>
<Error-Message>
```

#### 2.8 Reserved Words

Eugene has the following reserved words which can not be used as identifiers:

Property Part PartType \todo{Type} Device \todo{Template} Rule

The following list of reserved words can be specified in lower case too:

REPRESSES INDUCES DRIVES MORETHAN CONTAINS EXACTLY SAME\_COUNT WITH THEN BEFORE ALL\_BEFORE SOME\_BEFORE AFTER ALL\_AFTER SOME\_AFTER NEXTTO ALL\_NEXTTO SOME\_NEXTTO STARTSWITH ENDSWITH EQUALS ALL\_FORWARD ALL\_REVERSE SOME\_FORWARD SOME\_REVERSE FORWARD REVERSE SAME\_ORIENTATION ALTERNATE\_ORIENTATION AND OR NOT IF THEN ELSEIF ELSE FOR FUNCTION

### Declarative Language Features

Eugene v2.0 supports the specification of synthetic biology designs based on biological Facts and Rules. Additionally, Eugene v2.0 supports the invocation of user-defined Functions for the verification and enumeration of rule-compliant synthetic biology designs.

```
\langle script \rangle ::= \langle declaration \rangle +
\langle declaration \rangle ::= \langle fact \rangle \mid \langle rule \rangle \mid \langle function \rangle
```

#### 3.1 Declaration of Biological Facts

Facts can either be specified manually or imported using data exchange standards. Here, we define the Eugene language for the manual specification of facts. We describe Eugene's data exchange facilities in Section 4. Every specification of a fact must be closed with the semi-colon character (;).

```
\langle fact \rangle ::= (\langle property \rangle \mid \langle type \rangle \mid \langle instantiation \rangle \mid \langle relation \rangle);
```

In the following, we define the syntax for the specification of properties, types, instantiations, and relations in more detail.

#### 3.1.1 Properties of Biological Facts

Properties represent attributes and characteristics of genetic components, either primitive components or composite components. Each property must have an identifier (see Section 2.4) and a primitive type (see Section 2.2).

```
\langle \mathit{property} \rangle ::= \, \underline{\mathbf{Property}} \,\, \mathtt{id} \,\, \underline{\big(} \,\, \langle \mathit{primitive-type} \rangle \,\, \underline{\big)}
```

#### Example:

#### **3.1.2** Types

Types are abstract representations of the common structure of genetic elements, such as parts. Eugene v2.0 only supports the definition of types of parts. Types define the properties and structure of genetic elements. For example, a type Promoter can defines the common characteristics (i.e. properties) of genetic elements, such as DNA sequence, list of operators, or strength.

```
\langle type \rangle ::= \langle part-type \rangle
\langle part-type \rangle ::= \underline{\mathbf{PartType}} \ \mathtt{id} \ \underline{(} \ \langle list-of-ids \rangle \ \underline{)}
\langle list-of-ids \rangle ::= \mathtt{id} \ (, \ \langle list-of-id \rangle )?
```

For the specification of part types, the identifiers in the list-of-ids> must refer to the identifiers of specified properties (see Section 3.1.1). As defined in Section 2.4, properties must have been defined before specifying a type.

#### 3.1.3 Instantiations

In Eugene, types define the common properties and structure of genetic elements. Hence, Eugene also provides facilities to instantiate types with concrete genetic facts. That is, instances of genetic types have a unique name and property values. For examples, a concrete instance of the type Promoter can have a specific DNA sequence (e.g. ATCG).

```
\begin{split} &\langle instance \rangle ::= \text{ id id } \underline{\big(} \ \langle list\text{-}of\text{-}property\text{-}values \rangle \ \underline{\big)} \\ &\langle list\text{-}of\text{-}property\text{-}values \rangle ::= (\ \langle dot\text{-}notation \rangle \ | \ \langle list\text{-}notation \rangle \ ) \\ &\langle dot\text{-}notation \rangle ::= \underline{\cdot} \text{ id } \underline{\big(} \text{ primitive-value } \underline{\big)} \ (\ \underline{\cdot} \ \langle dot\text{-}notation \rangle \ ) \\ &\langle list\text{-}notation \rangle ::= \langle value \rangle \ (\ , \ \langle list\text{-}notation \rangle \ ) \end{split}
```

The first identifier of an instance definition refers to the type of the instance (see Section 3.1.2). The second identifier is the name of the instance. In Eugene, we support two notations of specifying the property values of instances, the dot (.) notation and the list notation. Using the list notation necessitates to specify the property values in the same order as specified in the type specification. The dot notation offers the possibility to specify the property name of the property to that the specified value must be assigned.

Eugene does check the types of the property values with the type of the property. For example, it is not possible to assign a num value to a property of type txt.

In Eugene v2.0, we do not support the assignment of property values to undefined properties. Each property must be defined in the type specification.

#### 3.1.4 Composite Biological Facts

#### 3.1.5 Relations

### 3.2 Design Templates

#### 3.3 Rules and Constraints

In Eugene, rules are conditions which can only take the values **true** or **false**.

#### 3.3.1 Structural Constraints

integrate table from miniEugene paper

#### 3.3.2 Relational Expressions

Relational expressions are rules of the following form:

```
\langle rel-exp \rangle ::= \langle exp \rangle \langle rel-operator \rangle \langle exp \rangle
\langle exp \rangle ::= \langle exp-operand \rangle (\langle exp-operator \rangle \langle exp-operand \rangle)^*
\langle exp-operand \rangle ::= \langle property \rangle \mid \langle constant \rangle
\langle property \rangle ::= \text{id} (\langle index \rangle)^* (\cdot \text{id} (\langle index \rangle)^*)^* (\cdot \langle function \rangle)^*
\langle constant \rangle ::= \text{number} \mid \text{string}
\langle exp-operator \rangle ::= + \mid -\mid *\mid /
\langle index \rangle ::= [\text{number}]
\langle function \rangle ::= \text{size} \mid \text{length} \mid ...
\langle rel-operator \rangle ::= < \mid <= \mid == \mid != \mid > \mid > \mid > =
```

#### **Examples:**

- Select all promoters whose strength is greater than 5: Promoter.strength > 5
- Select all pairs of repressors and promoters that have the following relation: Repressor.repress == Promoter.name
- Select all pairs of promoters and RBSs that have following relation: Promoter.sequence.size + RBS.sequence.size <= 500
- Select all devices on that the following equation holds **true**. Device[0].sequence.size <= Device[1].sequence.size
- find and provide more examples

#### 3.4 Containers for Biological Facts and Rules

Eugene v2.0 provides two types of containers for biological facts and rules:

- Collections are unordered groups of biological data.
- Arrays are ordered groups of biological data.

Eugene v2.0 also supports nested container types. That is, collections can contain collections and arrays, as well as arrays can contain other arrays and collections.

#### 3.4.1 Collections

#### **3.4.2** Arrays

### Data Exchange Capabilities

Eugene v2.0 supports the following types of data exchange facilities:

- INCLUDE of Eugene scripts
- IMPORT of Eugene containers
- IMPORT of biological data from the iGEM partsregistry<sup>1</sup>
- IMPORT and EXPORT of biological data from Genbank files<sup>2</sup>
- IMPORT and EXPORT of biological data compliant with the Synthetic Biology Open Language (SBOL) standard<sup>3</sup>
- VISUALIZATION of designs compliant with the SBOL Visual graphical notation<sup>4</sup> using Pigeon<sup>2</sup>

#### 4.1 Inclusion of Eugene Scripts

Eugene v2.0 supports to include Eugene scripts into a Eugene script via the include keyword. The syntax is defined as follows:

```
\langle include\text{-}statement \rangle ::= \underline{include} string
```

The string argument must specify the relative location of the included Eugene script to the actual Eugene script. If the specified Eugene script does not exist or is not a well-formed Eugene script, then an error is reported.

#### Examples:

• Select all promoters whose strength is greater than 5:

<sup>1</sup>http://parts.igem.org
2http://www.ncbi.nlm.nih.gov/genbank/
3http://www.sbolstandard.org/
4http://www.sbolstandard.org/visual

**Note!** In Eugene v2.0 we do not check for cyclic inclusions of Eugene script. For example, if the Eugene script A includes Eugene script B and Eugene script B includes Eugene script A, then an endless loop is the result.

- 4.2 Data Import of Eugene Containers
- 4.3 Data Import from the iGEM Partsregistry
- 4.4 GenBank Data Import and Export
- 4.5 SBOL Data Import and Export
- 4.6 SBOL Visual Compliant Design Visualization

### Imperative Language Features

Imperative languages features serve for the specification of the control-flow of design synthesis. In Eugene, imperative language features enable to

- control the execution of statements if certain conditions are satisfied (Section 5.1),
- iterate over the elements of a Eugene container, which are either collections or arrays (Section 5.2),
- group recurring blocks of statements into reusable functions (Section 5.3), and
- execute blocks of statements multiple times while a certain condition is satisfied (Section 5.4).

Besides user-defined functions, Eugene also provides various built-in functions which can be invoked within a Eugene script (see 5.5).

#### 5.1 Conditional Branches

- 5.2 Iterators
- 5.3 Function Prototypes
- 5.4 Loops
- 5.5 Built-in Functions
- 5.5.1 Design Verification

Built-in functions:

assert

note

#### 5.5.2 Design Enumeration

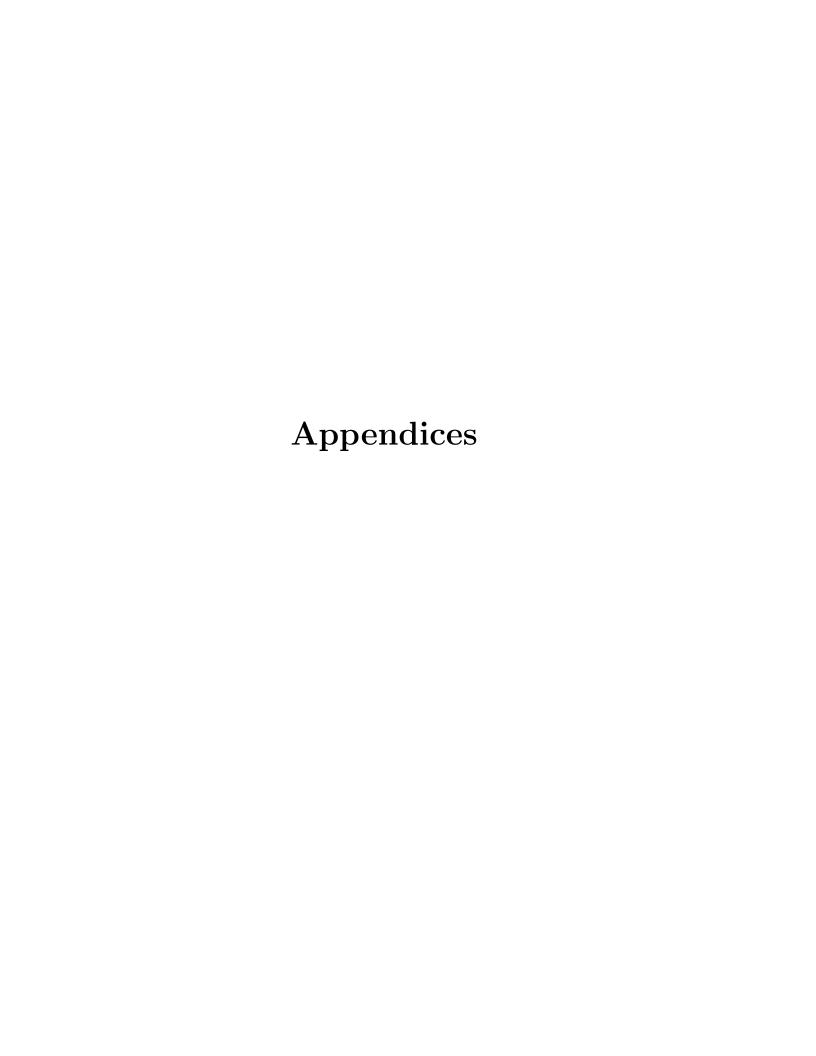
Built-in functions:

product permute

If a Eugene script contains a function prototype with an equivalent signature as one built-in function, then the built-in function is overwritten by the user-prototyped function.

### Algorithms

- 6.1 Enumeration of rule-compliant designs
- 6.2 Verification of designs and rules



# Appendix A<br/>Railroad Diagrams

### **Bibliography**

- [1] Standard, E. S. S. http://www.cl.cam.ac.uk/~mgk25/iso-14977.pdf
- [2] Bhatia, S.; Densmore, D. ACS Synthetic Biology 2013, 2, 348–350.