

Dictation Evaluation Reddit Parser

Garcia, Benjamin Lembke, Logan MacMillan, Kyle Smith, Christopher
 Stelter, Andrew

September 21, 2018

Contents

Title

I	Language Tutorial	1
1	Introduction	1
2	Getting Started	1
3	Basic Usage	1
II	Refrence Guide	1
4	Syntax Notation	1
5	A Note About Language Extension	1
6	Lexical Conventions	2
6.1	White Space	2
6.2	Comment Syntax	2
6.3	Common Syntax Elements	2
7	Language Modes	2
7.1	Main Mode	2
7.2	Selection Mode	3
7.3	Criteria Mode	3
8	Data Types	3
8.1	Post	3
8.1.1	Retrieving Text Body	3
8.2	Source Modules	3
8.2.1	Loading Grammar	4
8.2.2	Performing Queries	4
8.2.3	Sorting	4
8.2.4	General Matching	4
8.3	Criterion	4
8.4	Selection	4
9	Statements and Expressions	4
9.0.1	Stop	5
9.1	Main Mode	5
9.1.1	Load/Unload	5
9.1.2	Recall	5
9.1.3	Clear	5
9.1.4	Read	6
9.1.5	Create	6
9.2	Criteria Mode	6
9.2.1	Recall	6
9.2.2	Save	6
9.2.3	Add/Remove Posts	7

Part I

Language Tutorial

1 Introduction

2 Getting Started

3 Basic Usage

Part II

Reference Guide

This section describes the DERP language, as outlined in the DERP Language white paper submitted to the course CSC 792 (Compilers) at the South Dakota School of Mines Technology in the Fall of 2018. This reference defines all of the main constructs of the language and gives some examples of their usages; it will not contain hard definitions of any of the specific language constructs defined in plugins to the language; however, it may use some as examples.

4 Syntax Notation

Syntax definitions in this reference will be separated from surrounding text by an empty line, indented, and written in a monospace font. The grammar definitions will follow Backus-Naur form, occasionally using Regular Expression syntax for brevity. Additionally, definitions will use the following conventions:

- Terminals will be unformatted; keywords will be expressed as a sequence of terminal characters.
- Nonterminals will be *italicized*
- Characters which could be considered either terminals or BNF/Regular Expression syntax will be **bold when they should be interpreted as part of BNF or Regular Expression syntax.**
- All terminals shown are lower-case text; however, DERP parsers should match them case-insensitive.
- Production definitions follow the format

nonterminal : *formula*

Where the *formula* is any syntactic rules for parsing the *nonterminal*. Furthermore, the rest of this manual will reference specific terminals, and nonterminals inline using the same monospace font and formatting.

5 A Note About Language Extension

Because one of the main features of DERP is that it supports loading and unloading of language extension modules, there are some guarantees we can make in this document which will only hold true for the base language with no loaded modules.

Modules are allowed to define additional parse rules for the following tokens of DERP:

- source

- field
- sortkey

Furthermore, modules are allowed to define any additional parse rules they would like, so long as those parse rules are used within the additional definitions for the above tokens and do not redefine other DERP parse rules. Fields and sources are defined further below in [section]. For example, a source module for interfacing with a Reddit API module might define the following rules:

```
source : subreddit string | reddit
field : title | upvotes | tags | nsfw
nsfw : nsfw | not safe for work
sortkey : popular
```

6 Lexical Conventions

6.1 White Space

White space in the DERP language is interpreted in one of three ways, depending upon its context. Any whitespace character (that is, any character in the set `[\r, \n, \t, ' ']`) can be used as a part of a string literal used to identify something to the interpreter. The whitespace character (a unix-style newline) following any sequence of non-whitespace text indicates the end of a DERP statement to be parsed. Any whitespace character in any other context is ignored.

6.2 Comment Syntax

DERP does not support in-code comments of any sort.

6.3 Common Syntax Elements

A number of common keywords and phrases exist in the DERP language and are used in multiple expressions. Some of them are just common definitions that are useful to have defined, and some are optional syntax elements. The optional syntax elements are injected into expressions to make them sound more like standard English speech when spoken, but they are often not required for the syntax to be valid. The common syntax elements and keywords are defined here.

```
article : a | an | the
string : "[a-zA-Z]+"
digit : [0-9]
number : digit + ([.,]digit+)?
```

7 Language Modes

The DERP language is defined in terms of different modes of operation. Switching to a new mode is achieved through the use of certain expressions which are specific to the current mode. When a mode is triggered, the subset of the language that is valid may change. The next sections give a brief description of each mode, as well as the statements used to switch between them.

7.1 Main Mode

As the name implies, this is the mode that a DERP interpreter should begin in. In this mode, statements that change the loaded plugins, output the instructions for a saved query, read a saved query, clear the currently saved result, or begin query/criteria creation are all valid. See section [section] for more details about Main Mode statements.

Use a *create_expression* with the selection keyword to switch from Main Mode to Selection Mode

create_expression : create article? new? selection

Use a *create_expression* with the criteria keyword to switch from Main Mode to Criteria Mode

create_expression : create article? new? criteria

To end your DERP program or close the DERP interpreter use a *stop_expression* in Main Mode

stop_expression : stop | exit

7.2 Selection Mode

Selection mode is used to create a new selection statement which can be executed at a later time from Main Mode. While in Selection Mode, statements which build the query, read the query statements, save the query, or exit Selection Mode are all valid. See section [section] for more details about Selection Mode statements.

To exit Selection Mode (and return to Main Mode), use a *stop_expression*.

7.3 Criteria Mode

Criteria mode is very similar to selection mode; it is used to create criteria rather than selections. The difference between a criterion and a selection is that a criterion cannot be specific to any source; it is a filter that can be applied to a selection, but it is not a valid selection on its own. See section [section] for more details about Criteria Mode statements.

To exit Criteria Mode (and return to Main Mode), use a *stop_expression*.

8 Data Types

8.1 Post

Posts are the data type that DERP is built around. A Post is nothing more than some body of text and any number of fields containing information about the post. Fields of a post are key-value pairs mapping some name for the field to its value.

While defining the interface of a post is the task of a DERP interpreter writer, some of the functionalities it will likely be required to have are described here.

8.1.1 Retrieving Text Body

Posts should provide some method by which the DERP interpreter can acquire the entire text in the body of the post. It is not required that this information be stored directly in the post object itself, however – Retrieving this data can be delayed until it is requested by the DERP interpreter.

8.2 Source Modules

The first data type that someone writing DERP code will encounter is the Source Module. Source Modules are not created directly by the programmer, but are loaded and unloaded to modify the language. A source consists of, at minimum, new syntax for defining how to reference the source, and a public code interface that the DERP interpreter can use to get results from online.

Source modules are allowed to define new syntax rules for the following DERP nonterminals:

- source
- field
- sortkey

When input is found to match the rule a module provided for source, it indicates that the source module will be able to convert that input to a set of posts.

Source modules are allowed to populate an arbitrary set of key-value fields in the posts they create. Defining new syntax for *field* extends the DERP parser to be able to recognize the keys that the source module may populate in posts it retrieves.

By default, every field is also a sort key; however, it may be desirable to sort posts according to some ordering other than one of their fields. Defining new syntax for *sortkey* extends the language to recognize these new sorting methods.

While defining the interface of a source is the task of a DERP interpreter writer, some of the functionalities it will likely be required to have are described here.

8.2.1 Loading Grammar

As stated above, source modules are allowed to define new syntax and use that syntax to overload parts of the DERP language; therefore, any interface to the source module will likely provide a way for the DERP interpreter to get those changes to the grammar when the module is loaded.

It is anticipated that the grammar definition provided by a module may involve extra definitions of the fields added. These extra definitions can include information such as the data types associated with those fields. Without such a definition, DERP interpreters will not be able to perform type-checking on *qualifiers* before attempting to executing them.

8.2.2 Performing Queries

The most important part of a source module is its ability to take the input that was parsed according to source syntax and produce posts. Because source modules will likely deal with very large databases of information that posts can be retrieved from, this interface will likely need to provide functionalities to start retrieving posts from a specific point in the database and to retrieve a finite number of results from there.

8.2.3 Sorting

If the source module defines some non-standard *sortkey* it will likely also need to provide some method by which the interpreter can sort things when that *sortkey* is specified.

8.2.4 General Matching

One of the special qualifiers defined in [section], the *on_exp* qualifier, is used as a general-purpose match of posts against some text. If the interpreter does not implement some sort of generic matching routine, it will likely be the responsibility of source modules to determine if a post satisfies this generalized match.

8.3 Criterion

A criterion is a user-defined filter for specifying a subset of the Posts returned from executing a query. Criteria are created by entering Criteria Mode, making a number of statements to define the criterion, and then saving the criterion with a *save_expression*.

8.4 Selection

A selection follows many of the same rules as a Criterion, with the important difference that it can reference a specific source using the syntax rules defined in some source module. Because Selections contain one or more sources to get Posts from, they can be executed with a read statement to return a set of posts.

9 Statements and Expressions

All statements in DERP are a single line of text beginning with some expression and ending with the newline character. There is no provision for splitting a statement across multiple lines. Most expressions listed in

this section are valid only in the mode they are listed under, but there is one special expression that is always valid: *stop_expression*. (A version of the *recall_expression* is also valid at all times, but the syntax differs depending on the mode, so each version will be listed under the appropriate mode heading below)

statement : *expression* "

Rather than list every possible expression type here, we use the convention throughout this document that any nonterminal ending with *_expression* is a valid production of *expression*.¹

9.0.1 Stop

The Stop expression can be used in any context to switch modes back to the previous mode. In Selection or Criteria mode, this means returning to the Main Mode. In Main Mode, the Stop expression is used to end the DERP program. If DERP code is being run directly through an interactive interpreter, this should end the interpreter loop.

The stop expression is simply one of the keywords stop or exit

stop_expression : stop | exit

9.1 Main Mode

The main mode allows users to load and unload modules and to delete, execute, and display the definitions of existing criteria and selections.

9.1.1 Load/Unload

As mentioned above, source modules are imported code. These imports are done with the *load_expression*, and sources modules can be unloaded with the *unload_expression*. While a source is loaded, its parse rules will be used to parse any statements in addition to the standard DERP language rules.

Loading source modules is done with the expression

unload_expression : unload *string*

Where the string given should be the same string name that was used to load the module originally. It is not valid to use the unload keyword with a name that was not previously used to load a source module.

9.1.2 Recall

The recall expression can be used from Main Mode to retrieve the lines of DERP code that make up a criteria or selection. It follows the syntax

recall_expression : recall *string*

The string following the keyword recall is the name that was given in a preceding save expression from within Criteria Mode or Selection Mode.

9.1.3 Clear

The clear expression is used to delete a created Criteria or Selection. After a name is used with the *clear_expression*, it will not be valid in any future *read_expression*, *recall_expression*, or *match_expression* unless a *create_expression* is used to assign the string a new selection or criterion

clear_expression : clear *string*

¹It is intentional that on_exp and with_exp do not match this convention.

9.1.4 Read

Read expressions are used to execute a selection and output the resulting Posts. Optionally, a read expression can have a sort predicate, which defines the order that the resulting Posts should be outputted.

Read expressions follow the syntax

read_expression : read *string* *sort_predicate*?
sort_predicate : (sorted | ordered) by *sortkey* *order*?
order : ascending | descending

The string given immediately after the read keyword should correspond to a previously created selection. Any field defined by a source module is to be a valid *sortkey*, however source modules may define additional syntax for *sortkey* to allow sorting by more context-specific means. Any posts returned from the selection for which the *sortkey* is not relevant (for example, if multiple sources are used in the selection, and the *sortkey* is specific to one of the sources) should be sorted to the end of the list in any order.

The DERP language does not define what constitutes ‘outputting posts’ after they are found using a query. This leaves the interpreter open to any number of formats for conveying post information to the user. It should be noted that DERP selection semantics guarantee that determining which posts can be done on a group of posts, and after they are processed they do not need to be re-processed if another group of posts is retrieved. This allows interpreters to work with stream sources that can provide a semi-infinite very large result as multiple smaller pieces.

9.1.5 Create

The create expression is used exclusively to trigger either selection mode or criteria mode, depending upon whether it is used with the criteria or selection keyword.

Create expressions take the form

create_expression : create *article*? new? (selection | criteria)

9.2 Criteria Mode

After a *create_expression* using the keyword **criteria** is read, all following statements will be interpreted using the criteria creation syntax defined here. Criteria mode ends when the stop statement is found.

Valid expressions in criteria mode can be used to add things to the filter or remove them, read back the criteria so far, and save the criteria with a name.

9.2.1 Recall

The *recall_expression* can be used in criteria mode to repeat back the add and remove expressions used for the criterion so far.

recall_expression : recall

While this is not a particularly useful functionality for users who are writing DERP code in a text file to be parsed, it exists in anticipation of DERP interpreters which take input and produce output in some format other than text on a screen. For example, it would be useful in an implementation of the DERP interpreter which takes voice commands.

9.2.2 Save

The *save_expression* is used to save the set of *add_expressions* and *remove_expressions* that have been executed since criteria mode was last entered. This set of commands is named according to the second part (the *string*) of the expression. Saved criteria become available for use immediately. Using saved criteria is covered with the **matching** qualifier

save_expression : save as *string*

9.2.3 Add/Remove Posts

The most important expressions in criteria mode are the *add_expression* and *remove_expression*. They build a sequence of executable statements that are used to filter the list of elements. Both expressions have similar form

```
add_expression :  add posts selector  
remove_expression :  remove posts selector
```

Both expressions use the selector, which is a series of one or more qualifiers strung together to indicate what items from a list of posts should be accepted by the criteria.

selector