

# ChatScript System Variables and Engine-defined Concepts

Copyright Bruce Wilcox, gowilcox@gmail.com www.brilligunderstanding.com  
Revision 7/18/2021 cs11.5

- Engine-defined Concepts
- System Variables
- Control over Input
- Interchange Variables

## Engine-defined concepts

In addition to concepts defined in script files, the system automatically defines a bunch of dictionary-based sets as well as dynamically computed concept members.

| set            | description                                |
|----------------|--|
| ~web_url       | word is a web url                          |
| ~email_url     | word is an email address                   |
| ~kindergarten  | word learned early in life                 |
| ~grade1_2      | word learned in these grades               |
| ~grade3_4      | word learned in these grades               |
| ~grade_5-6     | word learned in these grades.              |
|                | Unmarked words are learned even later      |
| ~utf8          | word has nonascii characters               |
| ~daynumber     | word could be a number of a day in a month |
| ~yearnumber    | word could be the number of a recent year  |
| ~dateinfo      | phrase is month day year of some kind      |
| ~kelvin        | temperature marker                         |
| ~celcius       | temperature marker                         |
| ~fahrenheit    | temperature marker                         |
| ~twitter_name  | twitter user name                          |
| ~hashtag_label | twitter topic reference                    |

## Interjections, “discourse acts”, and concept sets

Some words and phrases have interpretations based on whether they are at sentence start or not. E.g., *good day*, *mate* and *It is a good day* are different for *good day*.

Likewise *sure* and *I am sure* are different.

Words that have a different meaning at the start of a sentence are commonly called interjections.

In ChatScript these are defined by the `livedata/interjections.txt` file. In addition, the file augments this concept with “discourse acts”, phrases that are like an interjection. All interjections and discourse acts map to concept sets, which come thru as the user input instead of what they wrote.

For example *yes* and *sure* and *of course* are all treated as meaning the discourse act of agreement in the interjections file. So you don’t see *yes*, *I will go* coming out of the engine.

The interjections file will remap that to the sentence `~yes`, breaking off that into its own sentence, followed by *I will go* as a new sentence.

These generic interjections (which are open to author control via `interjections.txt`) are:

| interjection                   | description |
|--------------------------------|-------------|
| <code>~yes</code>              |             |
| <code>~no</code>               |             |
| <code>~emomaybe</code>         |             |
| <code>~emohello</code>         |             |
| <code>~emogoodbye</code>       |             |
| <code>~emohowzit</code>        |             |
| <code>~emothanks</code>        |             |
| <code>~emolaugh</code>         |             |
| <code>~emohappy</code>         |             |
| <code>~emosad</code>           |             |
| <code>~emosurprise</code>      |             |
| <code>~emomisunderstand</code> |             |
| <code>~emoskeptic</code>       |             |
| <code>~emoignorance</code>     |             |
| <code>~emobeg</code>           |             |
| <code>~emobored</code>         |             |
| <code>~emopain</code>          |             |
| <code>~emoangry</code>         |             |
| <code>~emocurse</code>         |             |
| <code>~emodisgust</code>       |             |
| <code>~emoprotest</code>       |             |

| interjection | description |
|--------------|-------------|
| ~emoapology  |             |
| ~emomutual   |             |

Because all interjections at the start of a sentence are broken off into their own sentence, this kind of pattern does not work:

**u:** (~yes \_\*)

You cannot capture the rest of the sentence here, because it will be part of the next sentence instead. This means interjections act somewhat differently from other concepts.

If you use a word in a pattern which may get remapped on input, the script compiler will issue a warning. Likely you should use the remapped name instead.

The following concepts are triggered by exactly repeating either the chatbot or oneself (to a repeat count of how often repeated). Repeats are within a recency window of about 20 volleys.

| concept       | description |
|---------------|-------------|
| ~repeatme     |             |
| ~repeatinput1 |             |
| ~repeatinput2 |             |
| ~repeatinput3 |             |
| ~repeatinput4 |             |
| ~repeatinput5 |             |
| ~repeatinput6 |             |

## POS (Part of Speech) Tags

Words will have pos-tags attached, specifying both generic and specific tag attributes, eg., ~noun, ~noun\_singular.

### Generic Specifics

| nouns                 | description |
|-----------------------|-------------|
| ~noun                 |             |
| ~noun_singular        |             |
| ~noun_plural          |             |
| ~noun_proper_singular |             |
| ~noun_proper_plural   |             |

| nouns                   | description |
|-------------------------|-------------|
| ~noun_gerund            |             |
| ~noun_number            |             |
| ~noun_infinitive        |             |
| ~noun_omitted_adjective |             |

| verbs                    | description |
|--------------------------|-------------|
| ~verb                    |             |
| ~verb_present            |             |
| ~verb_present_3ps        |             |
| ~verb_infinitive         |             |
| ~verb_present_participle |             |
| ~verb_past               |             |
| ~verb_past_participle    |             |
| ~aux_verb                |             |
| ~aux_verb_present        |             |
| ~aux_verb_past           |             |
| ~aux_verb_future         |             |
| ~aux_verb_tenses         |             |
| ~aux_be                  |             |
| ~aux_have                |             |
| ~aux_do                  |             |

Auxilliary verbs are segmented into normal ones and special ones. Normal ones give their tense directly. Special ones give their root word. The tense of the be/have/do verbs can be had via `~properties()` and testing for verb tenses

| adjectives            | description |
|-----------------------|-------------|
| ~adjective            |             |
| ~adjective_normal     |             |
| ~adjective_number     |             |
| ~adjective_noun       |             |
| ~adjective_participle |             |

| adjectives in comparative form | description |
|--------------------------------|-------------|
| ~more_form~most_form           |             |
| ~adverb                        |             |
| ~adverb_normal                 |             |

| adverbs in comparative form                                      | description                       |
|--|-----------------------------------|
| ~more_form~most_form   |                                   |
| ~pronoun~pronoun_subject~pronoun_object                          |                                   |
| ~conjunction_bits~conjunction_coordinate~conjunction_subordinate |                                   |
| ~determiner_bits~determiner~pronoun_possessive~predeterminer     |                                   |
| ~possessive  | covers ' and 's at end of word    |
| ~to_infinitive   | “to” when used before a noun      |
| ~preposition~particle  | free-floating preposition tied to |
| ~comma   |                                   |
| ~quote   | covers ' and _“_ when not en      |
| ~paren   | covers opening and closing par    |
| ~foreign_word  | some unknown word                 |
| ~there_existential   | the word there used existential   |

In addition to normal generic kinds of pos tags, words which are serving a pos-tag role different from their putative word type are marked as members of the major tag they act as part of. E.g,

|                         | description   |
|-------------------------|---|
| ~noun_gerund            | verb used as a<br>~noun   |
| ~noun_infinitive        | verb used as a<br>~noun   |
| ~noun_omitted_adjective | an adjective<br>used as a<br>collective noun<br>(eg <i>the beautiful<br/>are kind</i> ) |
| ~adjectival_noun        | noun used as<br>adjective like<br>bank “bank<br>teller”                                 |
| ~adjective_participle   | verb participle<br>used as an<br>adjective  |

For ~noun\_gerund in *I like swimming* the verb gerund *swimming* is treated as a noun (hence called noun-gerund) but retains verb sense when matching keywords tagged with part-of-speech (i.e., it would match **swim~v** as well as **swim~n**).

Additionally, there is

|                                | description   |
|--------------------------------|---|
| <code>~number</code>           | is not a part of speech, but is comprise of<br><code>~noun_number</code><br>(a normal number value like <i>17</i> or <i>seventeen</i> ) |
| <code>~adjective_number</code> | also a normal numeral value and also<br><code>~placenumbers</code><br>like <i>first</i> .   |
| <code>~integer</code>          |   |
| <code>~float</code>            |   |
| <code>~positiveinteger</code>  |   |
| <code>~negativeinteger</code>  |   |
| <code>~modelnumber</code>      | not a true number, but a word with both alpha and numeric   |
| <code>~filename</code>         | looks like a filename with extension  |

To can be a preposition or it can be special. When used in the infinitive phrase To go, it is marked `~to_infinitive` and is followed by `~noun_infinitive`.

|                               | description   |
|-------------------------------|---|
| <code>~verb_infinitive</code> | refers to a match on the infinitive form of the verb ( <i>I hear John sing</i> or <i>I will sing</i> ). |

|                    | description   |
|--------------------|---|
| ~There_existential | refers to the use of where not involving location, meaning the existence of, as in There is no future.  |
| ~Particle          | refers to a preposition piece of a compound verb idiom which allows being separated from the verb. If you say <i>I will call off the meeting</i> , call_off is the composite verb and is a single token. But if you split it as in <i>I will call the meeting off</i> , then there are two tokens. The original form of the verb will be call and the canonical form of the verb will be call_off, while the free-standing off will be labeled ~particle. |
| ~verb_present      | will be used for normal present verbs not in third person singular like <i>I walk</i> and   |

|                   | description   |
|-------------------|---|
| ~verb_present_3ps | will be used for things like <i>he walks</i>  |
| ~possessive       | refers to 's and ' that indicate possession, while possessive pronouns get their own labeling<br>~pronoun_possessive. |
| ~pronoun_subject  | is a pronoun used as a subject (like <i>he</i> )  |
| ~pronoun_object   | refers to objective form like <i>him</i>  |

Individual words serve roles in the parse of a sentence, which are retrievable. These include:

|   | description  |
|---|--|
| ~mainsubject  |  |
| ~mainverb   |  |
| ~mainindirect   |  |
| ~maindirect   |  |
| ~subject2   |  |
| ~verb2  |  |
| ~indirectobject2  |  |
| ~object2  |  |
| ~subject_complement   | adjective object of sentence involving linking verb                    |
| ~object_complement  | 2ndary noun or infinitive verb filling modifying mainobject or object2 |
| ~conjunct_noun~conjunct_verb~conjunct_adjective~conjunct_adverb~conjunct_phrase |  |



|                              | description  |
|------------------------------|--|
| <b>~postnominalAdjective</b> | adjective<br>occurring<br>AFTER the<br>noun it<br>modified   |
| <b>~reflexive</b>            | reflexive<br>pronouns  |
| <b>~not</b>                  |  |
| <b>~address</b>              | noun used as<br>addressee of<br>sentence   |
| <b>~appositive</b>           | noun restating<br>and modifying<br>prior noun  |
| <b>~absolutephrase</b>       | special phrase<br>describing<br>whole sentence   |
| <b>~omittedtimeprep</b>      | modified time<br>word used as<br>phrase but<br>lacking<br>preposition<br>( <i>Next tuesday I<br/>will go</i> ) |
| <b>~phrase</b>               | a prepositional<br>phrase start<br>(except   |
| <b>~clause</b>               | a subordinate<br>clause start  |
| <b>~verbal</b>               | a verb phrase  |

and special concepts: | **~capacronym** | word is in all caps (and &) and is likely an acronym | **~emoji** | word starts and end with : and represents an emoji

## System Variables

The system has some predefined variables which you can generally test and use but not normally assign to. These all begin with % . Ones that are reasonable to set are written in bold underline. Boolean values are always 1 or null on returns. 1 or 0 if you are setting them.

## Date & Time & Numbers

| variable                      | description  |
|-------------------------------|--|
| <code>%date</code>            | one or two digit day of the month  |
| <code>%day</code>             | Sunday, etc  |
| <code>%daynumber</code>       | 1-7 where 1 = Sunday   |
| <code>%fulltime</code>        | seconds representing the current time and date (Unix epoch time)   |
| <code>%fullmtime</code>       | Numeric full time/date in milliseconds (Unix epoch time)   |
| <code>%hour</code>            | 0-23   |
| <code>%timenumber</code>      | completely consistent full time info in numbers that you can do <code>_0 = ~burst(%timenumber)</code> to get <code>_0</code> =seconds (2digit) <code>_1</code> =minutes (2digit) <code>_2</code> =hours (2digit) <code>_3</code> =dayinweek(0-6 Sunday=0) <code>_4</code> =dateinmonth (1-31) <code>_5</code> =month(0-11 January=0) <code>_6</code> =year. You need to get it simultaneously if you want to do accurate things with current time, since retrieving <code>%hour</code> <code>%minute</code> separately allows time to change between calls |
| <code>%leapyear</code>        | boolean if current year is a leap year   |
| <code>%daylightsavings</code> | boolean if current within daylight savings   |
| <code>%minute</code>          | 0-59   |
| <code>%month</code>           | 1-12 (January = 1)   |
| <code>%monthname</code>       | January, etc   |
| <code>%second</code>          | 0-59   |
| <code>%volleytime</code>      | number of seconds of computation since volley input started  |
| <code>%time</code>            | hh:mm in military 24-hour time   |
| <code>%zulu</code>            | 2016-07-27T11:38:35.253Z   |
| <code>%week</code>            | 1-5 (week of the month)  |
| <code>%year</code>            | e.g., 2011   |
| <code>%rand</code>            | get a random number from 1 to 100 inclusive  |

Time and date information are normally local, relative to the system clock of the machine CS is running on. See `$cs_utcoffset` for adjusting time based on relationship to utc (e.g your server is in Virginia and you are in Colorado).

`%rand` is only pseudo-random. A specific username is assigned a seed based on their name. Thereafter the seed evolves by the dialog but it is repeatable when the same user starts over again. If you want truly random, use `%fullmtime` `%showmany` to get range 0 .. `$showmany-1`

## User Input

| variable      | description  |
|---------------|--|
| %bot          | current<br>bot<br>responding   |
| %revisedinput | Boolean<br>is<br>current<br>input<br>from<br>$\neg$ input<br>not<br>direct<br>from<br>user |
| %command      | Boolean<br>was the<br>user<br>input a<br>command   |
| %foreign      | Boolean<br>is bulk<br>of the<br>sen-<br>tence<br>com-<br>posed<br>of<br>foreign<br>words   |
| %impliedyou   | Boolean<br>was the<br>user<br>input<br>having<br>you as<br>implied<br>subject              |

| variable        | description   |
|-----------------|---|
| %impliedsubject | Boolean<br>was the<br>user<br>input<br>having<br>an<br>implied<br>subject<br>(not<br>you,<br>usually<br>I)    |
| %input          | the<br>count<br>of the<br>number<br>of<br>volleys<br>this<br>user<br>has<br>made<br>ever                      |
| %volley         | sae as<br>%input,<br>the<br>count<br>of the<br>number<br>of<br>volleys<br>this<br>user<br>has<br>made<br>ever |
| %ip             | ip<br>address<br>supplied   |
| %myip           | ip<br>address<br>of cs<br>server<br>responding  |

| variable       | description  |
|----------------|--|
| %language      | current<br>dictio-<br>nary<br>language   |
| %length        | the<br>length<br>in<br>tokens<br>of the<br>current<br>sentence   |
| %more          | Boolean<br>is there<br>another<br>sen-<br>tence<br>after<br>this   |
| %morequestion  | Boolean<br>is there<br>a ? or<br>ques-<br>tion<br>word in<br>the<br>pend-<br>ing<br>sentences  |
| %originalinput | all sen-<br>tences<br>user<br>passed<br>into<br>volley,<br>before<br>ad-<br>justed<br>in any<br>way<br>except<br>OOB<br>data is<br>stripped<br>off |

| variable          | description   |
|-------------------|---|
| %originalsentence | the<br>current<br>sen-<br>tence<br>after to-<br>keniza-<br>tion but<br>before<br>any<br>adjustments |
| %parsed           | Boolean<br>was<br>current<br>input<br>parsed<br>successfully  |
| %question         | Boolean<br>was the<br>user<br>input a<br>ques-<br>tion -<br>same as<br>? in a<br>pattern            |
| %quotation        | Boolean<br>is<br>current<br>input a<br>quotation  |
| %sentence         | Boolean<br>does it<br>seem<br>like a<br>sen-<br>tence<br>(sub-<br>ject/verb<br>or<br>command)       |

| variable                    | description  |
|-----------------------------|--|
| <code>%tableinput</code>    | current line being executed in a table expansion during script compilation   |
| <code>%tense</code>         | past , present, or future simple tense (present perfect is a past tense)     |
| <code>%user</code>          | user login name supplied   |
| <code>%userfirstline</code> | value of <code>%input</code> that is at the start of this conversation start |
| <code>%userinput</code>     | Boolean is the current input from the user (vs the chatbot)                  |

| variable          | description   |
|-------------------|---|
| <b>%voice</b>     | active<br>or<br>passive<br>on<br>current<br>input   |
| <b>%trace_on</b>  | Fake<br>empty<br>variable<br>used to<br>turn on<br>tracing<br>(see De-<br>bugging<br>commands)  |
| <b>%trace_off</b> | Fake<br>empty<br>variable<br>used to<br>turn off<br>tracing<br>(see De-<br>bugging<br>commands) |

**%inputsize** | gives how many characters were passed in input

**%inputlimited** | 1 if too many characters were given (relative to fullinputlimit)

## Chatbot Output

| variable               | description   |
|------------------------|---|
| <b>%inputrejoinder</b> | integer of<br>any pending<br>rejoinder for<br>input or null<br>if none<br>pending |



| variable                | description   |
|-------------------------|---|
| <b>%lastoutput</b>      | the text of the last generated response for the current volley - always null across volleys                                       |
| <b>%lastquest</b>       | Boolean did last output end in a ?  |
| <b>%outputrejoinder</b> | Boolean if system set a rejoinder for its current output or 0   |
| <b>%response</b>        | number of committed responses that have been generated for this sentence (see Advanced User-Advanced Output: Committed Responses) |

## System variables

| variable    | description                                |
|-------------|--|
| <b>%all</b> | Boolean is the :all flag on? (:all to set) |

| variable           | description   |
|--------------------|---|
| %document          | Boolean<br>is :docu-<br>ment<br>running                                     |
| %fact              | Numeric<br>value<br>most<br>recent<br>fact id                               |
| %freetext          | kb of<br>avail-<br>able<br>text<br>space                                    |
| %freedict          | number<br>of<br>unused<br>dictio-<br>nary<br>words                          |
| %freefact          | number<br>of<br>unused<br>facts   |
| %maxmatchvariables | highest<br>number<br>of<br>match<br>vari-<br>ables,<br>cur-<br>rently<br>20 |
| %maxfactsets       | highest<br>number<br>of<br>@fact-<br>sets,<br>cur-<br>rently<br>20          |

| variable           | description  |
|--------------------|--|
| <b>%host</b>       | name of the current host machine or “local”                              |
| <b>%regression</b> | Boolean is the regression flag on  |
| <b>%server</b>     | Boolean is the system running in server mode                             |
| <b>%rule</b>       | get a tag to the current executing rule. Can be used in place of a label |

| variable     | description  |
|--------------|--|
| %topic       | name of the current “real” topic . if control is currently in a topic or called from a topic which is not system or nostay, then that is the topic. Otherwise the most recent pending topic is found |
| %actualtopic | literally the current topic being processed (system or not)  |

| variable             | description   |
|----------------------|---|
| <b>%trace</b>        | Numeric value of the trace flag (:trace to set)   |
| <b>%httpresponse</b> | return code of most recent ^jsonopen call   |
| <b>%pid</b>          | Linux process id or 0 for other systems   |
| <b>%restart</b>      | You can set and retrieve a value here across a system restart.                          |
| <b>%timeout</b>      | Boolean tells if a timeout has happened, based on the time-limit command line parameter |

| variable   | description   |
|--|---|
| %lastcurltime  | Time  |
| %analy-<br>sis:<br>Name<br>Look<br>up:<br>Host/proxy<br>con-<br>nect:<br>App(SSL)<br>con-<br>nect:<br>Pre-<br>trans-<br>fer:<br>Total<br>Transfer: |   |
| %crosstalk4k   | 4k<br>buffer<br>in<br>server<br>visible<br>be-<br>tween<br>users to<br>pass<br>data<br>back<br>and<br>forth |
| %crosstalk1k   | 1k<br>buffer<br>in<br>server<br>visible<br>be-<br>tween<br>users to<br>pass<br>data<br>back<br>and<br>forth |

| variable           | description  |
|--------------------|--|
| <b>%logging</b>    | bit<br>status<br>of<br>server-<br>Log,<br>user-<br>Log,<br>and<br>host<br>name -<br>0=off<br>1=file<br>2=<br>stdout<br>4=stderr<br>8=prelog) |
| <b>%forkcount</b>  | number<br>of forks<br>re-<br>quested<br>in linux<br>evserver<br>environment  |
| <b>%servertype</b> | parent<br>or fork<br>in linux<br>evserver<br>environ-<br>ment,<br>server<br>or null<br>otherwise   |

| variable         | description   |
|------------------|---|
| <b>%dbparams</b> | copy of the server params given to db used as file-server (pg or mysql or mssql or mongo) |
| <b>%botid</b>    | bot id number in use  |

## Build data

| variable        | description                         |
|-----------------|-------------------------------------|
| <b>%dict</b>    | date/time the dictionary was built  |
| <b>%engine</b>  | date/time the engine was compiled   |
| <b>%os</b>      | os involved (linux windows mac ios) |
| <b>%script</b>  | date/time build1 was compiled       |
| <b>%version</b> | engine version number               |

You actually can assign to any of them. This will override them and make them return what you tell them to and is a particularly BAD thing to do if this is running on a server since it affects all users (unless you reset the variable at the end of the volley. Assigning a period to a variable resets it).

Typically one does this as a temporary assignment in a **#!** comment line to set up conditions for testing using **:verify**.

Making them return a new value is NOT the same thing as making the engine have a different value. Unless the variable is marked as settable, setting a value affects only the value returned by a future call to the system variable. It does not change engine values the variable is meant to reflect.



## Control Over Input

The system can do a number of standard processing on user input, including spell correction, proper-name merging, expanding contractions etc. This is managed by setting the user variable `$cs_token`.

The default `$cs_token` that comes with Harry is:

```
$cs_token = #DO_INTERJECTION_SPLITTING |  
            #DO_SUBSTITUTE_SYSTEM |  
            #DO_NUMBER_MERGE |  
            #DO_PROPERNAME_MERGE |  
            #DO_SPELLCHECK |  
            #DO_PARSE
```

The `#` signals a named constant from the `dictionarySystem.h` file. One can set the following:

These enable various LIVEDATA files to perform substitutions on input:

| flag                               | description   |
|------------------------------------|---|
| <code>#DO_ESSENTIALS</code>        | perform LIVEDATA/systemessentials which mostly strips off trailing punctuation and sets corresponding flags instead |
| <code>#DO_SUBSTITUTES</code>       | perform LIVEDATA/substitutes  |
| <code>#DO_CONTRACTIONS</code>      | perform LIVEDATA/contractions, expanding contractions   |
| <code>#DO_INTERJECTIONS</code>     | perform LIVEDATA/interjections, changing phrases to interjections   |
| <code>#DO_BRITISH</code>           | perform LIVEDATA/british, respelling brit words to American   |
| <code>#DO_SPELLING</code>          | performs the LIVEDATA/spelling file (manual spell correction)   |
| <code>#DO_TEXTING</code>           | performs the LIVEDATA/texting file (expand texting notation)  |
| <code>#DO_SUBSTITUTE_SYSTEM</code> | do all LIVEDATA file expansions   |

The contents of the files above are pairs of tokens per line. Left is the word to replace and right is the replacement. When multiple words are involved, the left side uses underscores to represent this and the right side uses `+`. If the right side is missing, it means just delete. | `#DO_INTERJECTION_SPLITTING` | break off leading interjections into own sentence | `#DO_NUMBER_MERGE` | merge multiple word numbers into one (*four and twenty*) | `#DO_PROPERNAME_MERGE` | merge multiple proper name into one (*George Harrison*) | `#DO_DATE_MERGE` | merge month day and/or year sequences (*January 2, 1993*) | `#JSON_DIRECT_FROM_OOB` | asking the tokenizer to directly process OOB

data. See `^jsonparse` in JSON manual. | `#NO_FIX_UTF` | do not adjust inputs with html or utf8 encodings to simple ascii.

| `#TOKENIZE_BY_CHARACTER` | Every non-whitespace character becomes its own token and canonical form. (good for Japanese)

If any of the above items affect the input (except `TOKENIZE_BY_CHARACTER`), they will be echoed as values into `%tokenFlags` so you can detect they happened. The next changes do not echo into `%tokenFlags` and relate to grammar of input:

| flag                               | description  |
|------------------------------------|--|
| <code>DO_POSTAG</code>             | allow pos-tagging (labels like <code>~noun ~verb</code> become marked)   |
| <code>DO_PARSE</code>              | allow parser (labels for word roles like <code>~main_subject</code> )  |
| <code>DO_CONDITIONAL_POSTAG</code> | perform pos-tagging only if all words are known. Avoids wasting time on foreign sentences in particular  |
| <code>NO_CONDITIONAL_IDIOM</code>  | will not perform substitutions in the dictionary which are considered conditional idioms   |
| <code>NO_ERASE</code>              | where a substitution would delete a word entirely as junk, don't   |
| <code>DO_SPLIT_UNDERSCORES</code>  | happens after all other input tokenization and adjustments except number merge, and separates words that have been conjoined either because the dictionary has them ( <i>credit_card</i> ) or because they were merged by proper name merging, or by substitution. The result is only words without underscores (excluding number words like <i>five_thousand_and_four</i> ) |
| <code>MARK_LOWER</code>            | if a word is considered a proper name in CS and is marked as an upper case word, this will force it to perform any markings for its lower case form as well. Sometimes users type stuff in upper case that really should be lower  |

Normally the system tries to outguess the user, who cannot be trusted to use correct punctuation or casing or spelling. These block that:

| flag                       | description   |
|----------------------------|---|
| STRICT_CASING              | for 1st word of a sentence, assume user uses correct casing on words  |
| NO_INFERENCE_QUESTION_MARK | system will not try to set the QUESTION-MARK flag if the user didn't input a ? and the structure of the input looks like a question |
| DO_SPELL_CHECK             | internal spell checking   |

| flag            | description   |
|-----------------|---|
| ONLY_LOWER_CASE | input (except "I") to be lower case, refuse to recognize upper-case forms of anything |
| NO_IMPERATIVE   | don't   |
| NO_WITHIN       | match fragments within a composite word   |
| NO_SENTENCE_END | break input into sentences  |

Normally the tokenizer breaks apart some kinds of sentences into two. These prevent that:

| flag     | description                          |
|----------|--------------------------------------|
| NO_COLON | break apart a sentence after a colon |

| flag             | description  |
|------------------|--|
| NO_SEMICOLON_END | break apart a sentence after a semi-colon  |
| UNTOUCHED_INPUT  | this alone, will tokenize only on spaces, leaving everything but spacing untouched |

| flag        | description   |
|-------------|---|
| LEAVE_QUOTE | <p>If the output is found within " " it will become a single token exactly as it is seen.</p> <p>W/o Leave_Quote, it is converted into a word without quotes and using underscores instead of spaces.</p> <p>So “My Fair Lady” becomes My_Fair_Lady, which would match a movie title if you had one, unlike <i>My Fair Lady</i> becoming the resulting token and unrecognized</p> |

| flag         | description   |
|--------------|---|
| SPLIT_QUOTES | If output is found within " " the quotes will be removed. |

Note

you can change `$cs_token` on the fly and force input to be reanalyzed via `^retry(SENTENCE)`. I do this when I detect the user is trying to give his name, and many foreign names might be spell-corrected into something wrong and the user is unlikely to misspell his own name.

Just remember to reset `$cs_token` back to normal after you are done. Here is one such way, assuming `$stdtoken` is set to your normal tokenflags in your bot definition outputmacro:

```
#! my name is Rogr
s: (name is _*)

    if ($cs_token == $stdtoken)
    {
        $cs_token = #DO_INTERJECTION_SPLITTING |
                    #DO_SUBSTITUTE_SYSTEM | #DO_NUMBER_MERGE |
                    #DO_PARSE
        retry(SENTENCE)
    }
    _0 is the name.
    $cs_token = $stdtoken
```

If you type *my name is Rogr* into a topic with this, the original input is spell-corrected to *my name is Roger*, but this will change the `$cs_token` over to one without spell correction and redo the sentence, which will now come back with *my name is Rogr* and be echoed correctly, and `$cs_token` reset.

That's assuming nothing else would run differently and trap the response elsewhere. If you were worried about that, it would be possible for the script to save where it is using `^getrule(tag)` and modify your control script to return immediate control to here after input processing if you had changed `$cs_token`.

## **%tokenflags**

These are the values that %tokenflags may have after analysis of a sentence...

```
#define PRESENT 0x0000000000002000ULL
#define PAST 0x0000000000004000ULL // basic tense- both present perfect
and past perfect map to it #define FUTURE 0x0000000000008000ULL
#define PRESENT_PERFECT 0x0000000000010000ULL // distinguish PAST
PERFECT from PAST PRESENT_PERFECT #define CONTINUOUS
0x0000000000002000ULL
#define PERFECT 0x0000000000004000ULL
#define PASSIVE 0x0000000000008000ULL
```

**define IMPLIED\_\_SUBJECT**

**define QUESTIONMARK**

**define EXCLAMATIONMARK**

**define PERIODMARK**

**define USERINPUT**

**define COMMANDMARK**

**define IMPLIED\_\_YOU**

**FOREIGN\_\_TOKENS**

**FAULTY\_\_PARSE**

**QUOTATION**

**NOT\_\_SENTENCE**

One or more of these will be set if input was changed do to use of these files



```
#DO_ESSENTIALS
#DO_SUBSTITUTES
#DO_CONTRACTIONS
#DO_INTERJECTIONS
#DO_BRITISH
#DO_SPELLING
#DO_TEXTING
#DO_NOISE
#DO_PRIVATE
#DO_NUMBER_MERGE
#DO_PROPERNAME_MERGE
#DO_SPELLCHECK
#DO_INTERJECTION_SPLITTING
```

## Private Substitutions

While in general, substitutions are defined in the LIVEDATA folder, you can define private substitutions for your specific bot using the scripting language. You can say

```
replace: xxx yyyyy
```

which defines a substitution just like a livedata substitution file. It actually creates a substitution file called `private0.txt` or `private1.txt` in your TOPIC folder.

Even then, those substitutions will not be enacted unless you explicitly add to the `$cs_token` value `#DO_PRIVATE`, eg

```
$cs_token = #DO_INTERJECTION_SPLITTING |
            #DO_SUBSTITUTE_SYSTEM |
            #DO_NUMBER_MERGE |
            #DO_PROPERNAME_MERGE |
            #DO_SPELLCHECK |
            #DO_PARSE |
            #DO_PRIVATE
```

The left side of the substitution pair is case insensitive (matches either case on input) and can be placed in double-quotes (which converts spaces to underscores internally).

The right side of the substitution pair is case sensitive and can be placed in double-quotes (which converts spaces to plus signs internally).

Note: if you privately define a substitution that leads to a known interjection, it will be treated as an interjection, marked as `DO_INTERJECTIONS` rather than `DO_PRIVATE`. Interjections do not perform an actual substitution, does not replace the words on the left with the interjection concept name on the right.

Instead interjections merely mark the phrase as being a member of that concept, leaving the actual words unchanged.

Similarly while canonical values of words can be defined in `LIVEDATA/SYSTEM/canonical.txt`, you can define private canonical values for your bots by using the scripting language. You can say:

```
canon: oh 0
canon: faster fast
```

which defines new canonical values for things and creates a file `canon0.txt` or `canon1.txt` in your TOPIC folder.

You can optionally add `MORE_FORM` or `MOST_FORM` as a 3rd argument, to set those flags for adjectives and adverbs.

If you want to set a canonical pair from a table during compilation, you can use a function to do the same thing (but only 1 pair at a time).

```
^canon(word canonicalform)
```

## Numeric Substitutions

A special kind of private substitution (equally applicable in regular substitution files) is the numeric substitution.

```
replace: ?_km kilometers
```

The `?_` matches a digit number followed immediately by `km`, like `1.2km` and will separate the number and replace the units with the given replacement. The input can be singular or have an 's' like `10.5dollars`. And it can be with or without abbreviation periods, like `10kps` or `10k.p.s`

## Apostrophe Substitutions replace

```
replace: 'xxx yyy
```

allows you to split during tokenization any word followed by 'xxx into two words, original sans 'xxx and yyy. eg

```
replace: 've have
```

gives “companies’ve =>”companies have“.

## Replacing to a word with + in it

Normally `replace: x y+z` will generate 2 words, `y` and `z`. If you need a plus in your word, you can escape your 2nd word:

```
replace: "black and decker" \BLACK+DECKER
```

## Interchange Variables

The following variables can be defined in a script and the engine will react to their contents.

| interchange variable | description                            |
|----------------------|--|
| <b>\$cs_token</b>    | described<br>exten-<br>sively<br>above |

| interchange variable       | description   |
|----------------------------|---|
| <code>\$cs_response</code> | controls<br>auto-<br>matic<br>han-<br>dling of<br>outputs<br>to user.<br>By<br>default<br>it<br>consists<br>of<br><code>\$cs_response</code><br>=<br><code>#Response_upperstart</code><br> <br><code>#response_removespacebeforecomma</code><br> <br><code>#response_alterunderscores</code><br> <br><code>#response_removetilde</code><br>If you<br>want<br>none of<br>theses,<br>use<br><code>\$cs_response</code><br>= 0 (all<br>flags<br>turned<br>off).<br>See<br><code>^print</code><br>for<br>expla-<br>nation<br>of flags.<br><code>#response_noconvertspecial</code><br>- leave<br>escaped<br>n r and<br>t alone<br>in<br>output<br>and<br><code>^log</code> ,<br>36 <code>#response_upperstart</code><br>- makes<br>the first<br>letter of<br>an<br>output<br>sen-<br>tence |

| interchange variable          | description   |
|-------------------------------|---|
| <code>\$cs_jsontimeout</code> | seconds before JsonOpen declares a timeout failure. If unspecified the default is 300   |
| <code>\$cs_crashmsg</code>    | in server mode, what to say if the server crashes and we return a message to the user. By default the message is <i>Hey, sorry. I forgot what I was thinking about.</i> |
| <code>\$cs_abstract</code>    | used with :abstract   |

| interchange variable        | description  |
|-----------------------------|--|
| <code>\$cs_looplimit</code> | loop()<br>defaults<br>to 1000<br>itera-<br>tions<br>before<br>stop-<br>ping.<br>You can<br>change<br>this<br>default<br>with<br>this |

| interchange variable    | description  |
|-------------------------|--|
| <code>\$cs_trace</code> | if this variable is defined, then whenever the user's volley is finished, the value of this variable is set to that of <code>:trace</code> and <code>:trace</code> is cleared to 0, but when the user is read back in, the <code>:trace</code> is set to this value. For a server, this means you can perform tracing on a user w/o making all user transactions dump trace data |

| interchange variable               | description   |
|------------------------------------|---|
| <code>\$cs_control_pre</code>      | name of<br>topic<br>(flag it<br>SYS-<br>TEM)<br>to run<br>in<br>gambit<br>mode<br>on pre-<br>pass,<br>set by<br>author.<br>Runs<br>before<br>any sen-<br>tences<br>of the<br>input<br>volley<br>are ana-<br>lyzed.<br>Good<br>for<br>setting<br>up<br>initial<br>values |
| <code>\$cs_usermessagelimit</code> | max<br>number<br>of mes-<br>sage<br>pairs<br>(user<br>input &<br>bot<br>output)<br>saved<br>in topic<br>file  |



| interchange variable          | description   |
|-------------------------------|---|
| <code>\$cs_externaltag</code> | name of<br>a topic<br>to use<br>to<br>replace<br>existing<br>internal<br>English<br>pos-<br>parser.<br>See<br>bottom<br>of<br>ChatScript<br>PosParser<br>manual<br>for<br>details |

| interchange variable      | description  |
|---------------------------|--|
| <code>\$cs_prepass</code> | <p>name of a topic (mark it SYS-TEM) to run in responder mode on main volleys, which runs before <code>\$cs_control_main</code> and after all of the above and parsing is done. Used to amend preparation data coming from the engine. You can use it to add your own spin on input processing before going to your main control. I use it to, for example, label commands</p> |

| interchange variable           | description   |
|--------------------------------|---|
| <code>\$cs_control_main</code> | name of<br>topic<br>(flag it<br>SYS-<br>TEM)<br>to run<br>in re-<br>sponder<br>mode<br>on<br>main<br>volleys,<br>set by<br>author |
| <code>\$cs_control_post</code> | name of<br>topic<br>(flag it<br>SYS-<br>TEM)<br>to run<br>in<br>gambit<br>mode<br>on post-<br>pass,<br>set by<br>author           |
| <code>\$botprompt</code>       | message<br>for<br>console<br>window<br>to label<br>bot<br>output  |
| <code>\$userprompt</code>      | message<br>for<br>console<br>window<br>to label<br>user<br>input<br>line  |

| interchange variable       | description  |
|----------------------------|--|
| <code>\$cs_crashmsg</code> | message<br>to use if<br>a crash<br>occurs.<br>see also<br><code>\$cs_crash</code>                              |
| <code>\$cs_crash</code>    | topic to<br>execute<br>in<br>gambit<br>mode if<br>a crash<br>occurs.<br>see also<br><code>\$cs_crashmsg</code> |
| <code>\$cs_language</code> | if<br>spanish,<br>will<br>adjust<br>spell<br>check-<br>ing for<br>spanish<br>colloquial                        |

| interchange variable    | description   |
|-------------------------|---|
| <code>\$cs_token</code> | bits<br>control-<br>ling<br>how the<br>tok-<br>enizer<br>works.<br>By<br>default<br>when<br>null,<br>you get<br>all bits<br>as-<br>sumed<br>on. The<br>possible<br>values<br>are in<br>src/dictionarySystem.h<br>(hunt<br>for<br>\$token)<br>and you<br>put a #<br>in front<br>of them<br>to gen-<br>erate<br>that<br>named<br>nu-<br>meric<br>constant |

| interchange variable       | description   |
|----------------------------|---|
| <code>\$cs_abstract</code> | topic used by :abstract to display facts if you want them displayed   |
| <code>\$cs_prepass</code>  | topic used between parsing and running user control script. Useful to supplement parsing, setting the question value, and revising input idioms |

| interchange variable   | description  |
|------------------------|--|
| \$cs_wildcardseparator | when a match variable covers multiple words, what should separate them by default it's a space, but underscore is handy too. Initial system character is space, creating fidelity with what was typed. Useful if _ can be recognized in input (web addresses). Changing to _ is consistent with multi-word representation and keyword recogni- |

| interchange variable      | description   |
|---------------------------|---|
| <b>\$cs_userfactlimit</b> | how many of the most recent permanent facts created by the script in response to user inputs are kept for each user. Std default is 100. * means all. |
| <b>\$cs_outputchoice</b>  | for regression: forces specific one of a [] [] output choice block - base 0   |
| <b>\$cs_response</b>      | controls some characteristics of how responses are formatted  |



| interchange variable        | description                                 |
|-----------------------------|---|
| <code>\$cs_randIndex</code> | the<br>random<br>seed for<br>this<br>volley |

| interchange variable | description  |
|----------------------|--|
| \$cs_utcoffset       | <p>if defined, then %time returns current utc time + time-zone offset. The offset is usually a simple number, meaning hours, and can have + or - in front of it. It can also be a normal time reference like 02:30 which means plus 2 hours and 30 minutes beyond utc, or - 01:30:20 which means 1 hour, 30 minutes, and 20 seconds before utc (as if anyone would</p> |

| interchange variable           | description  |
|--------------------------------|--|
| <code>\$\$db_error</code>      | error<br>mes-<br>sage<br>from a<br>post-<br>gres<br>failure<br>\$\$find-<br>text_start<br>- ^find-<br>text<br>return<br>the end<br>nor-<br>mally,<br>this is<br>where it<br>puts<br>the<br>start |
| <code>\$\$tcpopen_error</code> | error<br>mes-<br>sage<br>from a<br>tcpopen<br>error  |
| <code>\$\$document</code>      | name of<br>the doc-<br>ument<br>being<br>read in<br>docu-<br>ment<br>mode  |
| <code>\$cs_randindex</code>    | current<br>value of<br>the<br>random<br>genera-<br>tor<br>value  |

| interchange variable           | description  |
|--------------------------------|--|
| <code>\$cs_bot</code>          | name of<br>the bot<br>cur-<br>rently<br>in use   |
| <code>\$cs_login</code>        | login<br>name of<br>the user   |
| <code>\$\$csmatch_start</code> | start of<br>found<br>words<br>from<br>^match   |
| <code>\$\$csmatch_end</code>   | end of<br>found<br>words<br>from<br>^match   |
| <code>\$cs_fullfloat</code>    | if<br>defined,<br>causes<br>the<br>system<br>to gen-<br>erate<br>full<br>float<br>64-bit<br>preci-<br>sion on<br>out-<br>puts,<br>other-<br>wise<br>you get<br>2 digit<br>preci-<br>sion by<br>default |

| interchange variable | description  |
|----------------------|--|
| \$cs_botid           | when non-zero creates facts and functions restricted by this bit-mask so facts and functions created by other masks cannot be seen. allows you to separate facts and functions per bot in a multi-bot environment. During compilation if this is set by a bot: command, then functions created and facts created by tables |

| interchange variable               | description  |
|------------------------------------|--|
| <code>\$cs_numbers</code>          | if defined, causes the system to output numbers in a different language style: french, indian. All other values are english. |
| <code>\$cs_topicretrylimit</code>  | if defined changes how many times you can pass back <code>RETRY_TOPIC</code> before it fails (current limit is 30)           |
| <code>\$\$topic_retry_limit</code> | if exceeded, exceeded topic retry limit is encountered   |

| interchange variable               | description   |
|------------------------------------|---|
| <code>\$cs_topicretrylimit</code>  | if defined changes how many times you can pass back RETRY_TOPIC before it fails (current limit is 30) |
| <code>\$cs_userhistorylimit</code> | if not null, indicates how many volleys back are tracked as what was said by both parties             |

| interchange variable               | description   |
|------------------------------------|---|
| <code>\$cs_saveusedJson</code>     | if not null, the only JSON facts CS will write into the user's topic files that are referred to (directly or indirectly) from user variables being saved. (see below) |
| <code>\$cs_proxycredentials</code> | See ^JSONOPEN in JSON manual  |
| <code>\$cs_proxyserver</code>      | See ^JSONOPEN in JSON manual  |
| <code>\$cs_proxymethod</code>      | See ^JSONOPEN in JSON manual  |



| interchange variable                               | description   |
|--|---|
| <code>\$cs_addresponse</code>                      | provides a function name hook onto the output q to the user. See below.   |
| <code>%trace_on</code> and <code>%trace_off</code> | Pseudo system variable used by the <code>^test-pattern</code> and <code>^testout-put</code> call to let code request a trace be returned. |
| <code>\$cs_indentlevel</code>                      | controls indenting when tracing in <code>^test-pattern</code> . 3 is a good number usually  |

| interchange variable                                     | description  |
|--|--|
| <code>\$cs_tracetestoutput</code>                        | after  |
| set to 1 to  | this   |
| force tracing in   | many   |
| ^testoutput  | sen-   |
| \$cs_sentences_limit                                     | tences   |
|  | in   |
|  | volley,  |
|  | cs   |
|  | ignores  |
|  | the rest   |
|  | (default   |
|  | 50)  |
| <code>\$cs_outputlimit</code>                            | set as a   |
| Generating more  | json   |
| output than this   | struc-   |
| will report a bug  | ture of  |
| into   | move   |
| LOGS/bugs.txt  | its  |
| cs_summary' After volley prints to                       | terminal milliseconds of time used in preparation, run |
| After volley   | a  |
| prints to  | mongo  |
| terminal   | query  |
| milliseconds of  |  |
| time used  |  |
| cs_inputlimit' Restrict user input size (excluding goob) | 'cs_new_user   |
| set to 1, treat  |  |
| user as always   |  |
| new (don't try to  |  |
| read topic file)   |  |
|  |  |
| \$cs_mongoqueryparams'                                   |  |

Note for %trace\_on and %trace\_off - you can use the command line parameter **blockapitrace** to prevent tracing in any code you accidentally leave in place.

`$cs_saveusedJson` exists as a kind of garbage collection. Nowadays most facts will come from JSON data either from a website or created in script. But keeping on top of deleting obsolete JSON may be overlooked. When this variable is non-null, ChatScript will automatically destroy any JSON fact that cannot trace a JSON fact path back to some user variable. Variables that have as values the name of a JSON object or array automatically protect all JSON facts underneath. JSON references merely within some text string will not protect anything, nor will references from some other non-JSON fact.

`$cs_inputlimit=x:y` for excessively long user input (excluding oob portion), the input will be truncated by keeping the first x characters and the last y characters.

`$cs_crash` - This topic can generate an appropriate dummy output and CS completes that volley but does not save an updated user file. The NEXT volley coming in will force cs to completely reload itself before processing. Making a dummy output hopefully means the same fatal input will not be sent back into CS to crash it again (due to external retry when no answer is received from CS). E.g.,

```
topic: ~crashtopic system ()
      t: Huh?
```

`$cs_addresponse` names a function of 2 arguments that will be called when CS wants put text into the output queue of the user. The first argument will be what CS wants to output. The second is the rule tag that generated this output. If the function returns a failure code, the message will be aborted and not put into the queue. If the function returns a text value (not null) then that message will replace what was intended to go to the user.