

Forth Recognizer -- Request For Discussion

Author: Matthias Trute

Contact: mtrute@web.de

Version: 3

Date: September, 4 2016

Status: Final with Bugfix

Change history

- 2014-10-03 Version 1 - initial version.
- 2015-05-17 Version 2 - extend rationale, added ' and []
- 2015-12-01 Version 3 - separate use cases, minor changes for nested recognizer stacks. New `POSTPONE` action

- **2015-09-27**
 - Move naming conventions to informal appendix (not mandatory but worth to be mentioned)
 - New `POSTPONE` action. That finalizes the `R:TABLE` API.
- **2015-10-17**
 - Recognizers *can/may/shall* but don't *have to* replace the existing interpreter (Euroforth 2015).
 - empty recognizer stacks are now allowed as a consequence from the above.
 - [R:FAIL necessity](#)
- **2015-10-31**
 - **‘Nested Recognizer Stacks’** and (fast) switching between them. Introducing `FORTH-RECOGNIZER` as entry for the interpreter and similar words.
 - fix stack effects at various places. More precious `POSTPONE` description.
- **2015-11-07**
 - finalize the building blocks.
 - Introduce the recognizer stack id as an additional parameter.
 - remove use cases from the proposal. Affects e.g. `interpret`, `tick`, `bracket-tick` and `POSTPONE`. These words went to the informal appendix if applicable.
- **2015-11-17**
 - expose the `R>INT/R>COMP/R>POST` words.
- **2015-11-28**
 - document structure revisited.
- **2015-12-01**
 - final version and public announcement
- **2016-09-01**
 - Fix a code bug for the float literal `POSTPONE` to comply with the words.

Background

I'm working on a Forth for 8-bit micro-controllers for almost 10 years now (amforth.sf.net). It is not only a useful tool for serious work but a nice playground to experiment with Forth too.

In 2011 my Forth got a floating point library. Since a micro-controller is (was) a resource constrained system it is not an option to include it permanently. It has to be a loadable module. Therefore I needed a way to keep the core system small but at the same time able to fully handle the new numbers. All but one problem were easy to fix. Only adding the number format to the Forth interpreter turned out to be serious one. I searched the net for ways to extend the Forth interpreter. What I found was having many hooks in the interpreter (U. Hoffman, Euroforth 2008) or a conditional re-compile of the sources with an autotool/configure like build system. Nothing really convinced me or my users. While googling I stumbled across the number parsing prefix discussion in `c.l.f` in 2007. The ideas sketched there looked promising so I stopped searching and started with them to invent my own solution.

I changed the Forth interpreter into a dumb tool, that delegates all data related work to modules, which can be changed at run-time. That made it possible to load the FP library into a running system that afterwards was able to deal with the new numbers like native ones. Surprisingly the new system had no disadvantages in speed or size compared the old one, something I consider very important on a micro-controller.

Shortly thereafter, Bernd Paysan got interested in what I did (we have regular IRC sessions on Forth topics) and started to implement recognizers in gforth. He suggested changes that further simplified my concept and made it more flexible.

By now we reached a point that justifies the public review. There are two very different Forth's available (both GPL'ed) that implement recognizers. A third implementation is in the proposal (public domain).

A recognizer written for one Forth works without modification for the other ones too. The words used to actually implement a recognizer (mostly string processing) need to be available of course. E.g. I wrote a recognizer for time stamp strings with gforth that converts the hh:mm:ss notation into a double cell number for the seconds since midnight. The code runs on amforth too. Gforth is a 64-bit system on the PC, amforth a 16-bit system on an 8-bit micro-controller (hence the double numbers). With that, something like

```
: test 01:00:01 d. ." seconds since midnight" ; ok
test 3601 seconds since midnight ok
01:01:00 01:00:01 d+ d. 7261 ok
```

is possible. Similarly strings: everything that starts with a " is a string until the closing " is reached. Further string handling get the addr/len without the enclosing ".

```
: test "A string" type ; ok
test A string ok
" Another string" type ok Another string
```

Another use case are name-spaces with word lists, without touching ORDER:

```
: test i2c.begin i2c.sendbyte i2c.end ;
```

where `begin/sendbyte/end` are words from the word-list identified with `i2c` (a constant with the wid). The recognizer splits the word at the first dot and uses the left sub-word to get the a word-list. In that word-list it searches with the remaining string and handles the result just like an ordinary dictionary search: interpret, compile (or not found).

Implementations for these examples are available in the respective Forth systems.

Problem

The Forth compiler can be extended easily. The Forth interpreter however has a fixed set of capabilities as outlined in section 3.4 of the standard text: Words from the dictionary and some numbers.

It's not easily possible to use the Forth text interpreter in an application or system extension context. The building blocks (`FIND`, `COMPILE`, `,`, `>NUMBER` etc) are available but there is a gap between them and what the Forth interpreter does. Applications need to use either additional system provided and system specific intermediate words (if available) or have to re-invent the wheel to use e.g. numbers with a sign or hex numbers with the \$ prefix.

Some Forth interpreters have ways to add new data types. That makes it possible to use a loadable library to implement new data types to be handled like the built-in ones. An example are the floating point numbers. They have their own parsing and data handling words including a stack of their own.

To actually handle data in the Forth context, the processing actions need to be `STATE` aware. It would be nice if the Forth text interpreter, that maintains `STATE`, is able to do the data processing without exposing `STATE` to the data handling methods. For that the different methods need to be registered somehow.

Solution

The monolithic design of the Forth interpreter is factored into three major blocks: First the interpreter. It maintains `STATE` and organizes the work. Second the actual data parsing. It is called from the interpreter and analyses strings (usually sub-strings of `SOURCE`) if they match the criteria for a certain data type. These parsing words are grouped as a stack to achieve an order of invocation. The result of the parsing words is handed over by the interpreter to data specific handling methods. There are three different methods for each data type depending on `STATE` and to `POSTPONE` the data.

The combination of a parsing word and the set of data handling words to deal with the data is called a recognizer. There is no strict 1:1 relation between the parsing words and the data handling sets. A data handling set for e.g. single cell numbers can be used by different parsing words.

Whenever the Forth text interpreter is mentioned, the standard words `EVALUATE` (`CORE`), `'` (tick, `CORE`), `INCLUDE-FILE` (`FILE`), `INCLUDED``` (`FILE`), ```LOAD` (`BLOCK`) and `THRU` (`BLOCK`) are expected to act likewise. This proposal is not about to change these words, but to provide the tools to do so. As long as the standard feature set is used, a complete replacement with recognizers is possible.

The proposal is about the building blocks.

Proposal

XY. The optional Recognizer word set

XY.1 Introduction

A recognizer consists of two elements: a parsing word and information token(s) returned by the parsing words that identify the parsed data and provide methods to perform the various semantics of the data: interpret, compile and postpone. A parsing word can return different information tokens. A particular information token can be used by different parsing words.

There is a system provided information token called `R:FAIL`. It is used if no other token is applicable. This failure token is associated with the system error actions if used in step e) of the text interpreter (see Appendix). It is used to achieve the action d) of the section 3.4 text interpreter.

The parsing word of a recognizer has the stack effect (parsing word names start with `REC:` in this chapter)

```
REC:TABLE ( addr len -- i*x R:TABLE | R:FAIL )
```

"addr/len" is a string, if provided by the Forth text interpreter a sub-string inside `SOURCE`. The parsing word must not change the string content. It may change `>IN` however.

"i*x" is the result of the parse action of the string "addr/len". `R:TABLE` is the information token that the interpreter uses to execute the interpret, compile or postpone actions for the data "i*x".

All three actions are called with the "i*x" data as left by the parsing word and are generally expected to consume it. They can have additional stack effects, depending on what `R:TABLE:METHOD` actually does.

```
R:TABLE:METHOD ( ... i*x -- j*y )
```

The data items "i*x" don't have to be on the data stack, they can be at different places, if applicable. E.g. floating point numbers have a stack of their own. In this case, the data stack contains the `R:TABLE` information only.

The names `R:TABLE`, `REC:TABLE` and `R:TABLE:METHOD` don't have to actually exist, except the `R:FAIL` name.

A Forth system shall provide recognizers for integer numbers (both single and double precision) and the word look-up in the dictionary.

A recognizer stack is identified by its ID. The numeric value is system dependent and generally opaque. The elements of a recognizer stack are available with `GET/SET-RECOGNIZERS` only.

XY.2 Additional terms and notations

Information token: A single cell number. It identifies the data type and a method table to perform the data processing of the interpreter.

Recognizer: A combination of a text parsing word that returns information tokens together with parsed data if successful. The text parsing word is assumed to run in cooperation with `SOURCE` and `>IN`.

A recognizer stack is identified with its stack id. This is cell sized numeric value. It may but don't have to be an address accessible with `@` and `!` operations.

XY.3 Additional usage requirements

XY.3.1 Environment Queries

Obsolete.

XY.4 Additional documentation requirements

XY.4.1 System documentation

XY.4.1.1 Implementation-defined options

No additional options.

XY.4.1.2 Ambiguous conditions

- Change of the content of the parsed string during parsing.
- During `SET-RECOGNIZERS` the Recognizer stack size is exceeded. In this case, at least the following actions are possible * resize the stack, keeping its id * throw an exception (TBD) * execute an error action

XY.4.2 Program documentation

- No additional dependencies.

XY.5 Compliance and labeling

The phrase "Providing the Recognizer word set" shall be appended to the label of any Standard System that provides all of the Recognizer word set.

XY.6 Glossary

XY.6.1 Recognizer Words

DO-RECOGNIZER (addr len stack-id -- i*x R:TABLE | R:FAIL) RECOGNIZER

Apply the string at "addr/len" to the elements of the recognizer stack identified by `stack-id`. Terminate the iteration if either one recognizer returns a information token that is different from `R:FAIL` or the stack is exhausted. In this case return `R:FAIL`.

"i*x" is the result of the parsing word. It represents the data from the string. It may be on other locations than the data stack. In this case the stack diagram should be read accordingly.

FORTH-RECOGNIZER (-- stack-id) RECOGNIZER

A system VALUE with a recognizer stack id.

It is VALUE that can be changed using TO assigning a new recognizer stack id. This change has immediate effect.

The recognizer stack from this stack-id shall be used in all system level words like EVALUATE, LOAD etc.

GET-RECOGNIZERS (stack-id -- rec-n .. rec-1 n) RECOGNIZER

Return the execution tokens `rec-1 .. rec-n` of the parsing words in the recognizer stack identified with `stack-id`. `rec-1` identifies the recognizer that is called first and `rec-n` the word that is called last.

The recognizer stack is left unchanged.

R>COMP (R:TABLE -- XT-COMPILE) RECOGNIZER

Return the execution token for the compilation action from the recognizer information token.

R>INT (R:TABLE -- XT-INTERPRET) RECOGNIZER

Return the execution token for the interpretation action from the recognizer information token.

R>POST (R:TABLE -- XT-POSTPONE) RECOGNIZER

Return the execution token for the postpone action from the recognizer information token.

R:FAIL (-- R:FAIL) RECOGNIZER

An information token with two uses: First it is used to deliver the information that a specific recognizer could not deal with the string passed to it. Second it is a predefined information token whose elements are used when no recognizer from the recognizer stack could handle the passed string. These methods provide the system error actions.

The actual numeric value is system dependent.

RECOGNIZER (size -- stack-id) RECOGNIZER

Create a new recognizer stack with `size` elements.

RECOGNIZER: (XT-INTERPRET XT-COMPILE XT-POSTPONE "<spaces>name" --) RECOGNIZER

Skip leading space delimiters. Parse name delimited by a space. Create a recognizer information token "name" with the three execution tokens.

The words for XT-INTERPRET, XT-COMPILE and XT-POSTPONE are called with the parsed data that the associated parsing word of the recognizer returned. The information token itself is consumed by the caller.

Each of the words XT-INTERPRET, XT-COMPILE and XT-POSTPONE has the stack effect (... `i*x` -- `j*y`). The words to compile and postpone the data shall consume the data "`i*x`". If the data "`i*x`" is on different locations (e.g. floating point numbers), these words shall use that data.

SET-RECOGNIZERS (rec-n .. rec-1 n stack-id --) RECOGNIZER

Set the recognizer stack identified by `stack-id` to the recognizers identified by the execution tokens of their parsing words `rec-n .. rec-1`. `rec-1` will be the parsing word of the recognizer that is called first, `rec-n` will be the last one.

If the size of the existing recognizer stack is too small to hold all new elements, an ambiguous situation arises.

XY.7 Reference Implementation

The code has as little as possible dependencies.

```
: RECOGNIZER ( size -- stack-id )
  1+ ( size ) CELLS HERE SWAP ALLOT
  0 OVER ! \ empty stack
;
```

```

\ create the default recognizer stack
4 RECOGNIZER VALUE FORTH-RECOGNIZER

: SET-RECOGNIZERS ( rec-n .. rec-1 n stack-id -- )
  2DUP ! >R
  BEGIN
    DUP
  WHILE
    DUP CELLS R@ +
    ROT SWAP ! 1-
  REPEAT R> 2DROP
;

: GET-RECOGNIZERS ( stack-id -- rec-n .. rec-1 n )
  DUP @ DUP >R SWAP
  BEGIN
    CELL+ OVER
  WHILE
    DUP @ ROT 1- ROT
  REPEAT 2DROP
  R>
;

\ create a simple 3 element structure
: RECOGNIZER: ( XT-INTERPRET XT-COMPILE XT-POSTPONE "<spaces>name" -- )
  CREATE SWAP ROT , , ,
;

\ decode the data structure created by RECOGNIZER:
: R>POST ( R:TABLE -- XT-POSTPONE ) CELL+ CELL+ @ ;
: R>COMP ( R:TABLE -- XT-COMPILE ) CELL+ @ ;
: R>INT ( R:TABLE -- XT-INTERPRET ) @ ;

\ system failure recognizer
:NONAME -13 THROW ; DUP DUP RECOGNIZER: R:FAIL

: DO-RECOGNIZER ( addr len stack-id -- i*x R:TABLE | R:FAIL )
  DUP >R @
  BEGIN
    DUP
  WHILE
    DUP CELLS R@ + @
    2OVER 2>R SWAP 1- >R
    EXECUTE DUP R:FAIL <> IF
      2R> 2DROP 2R> 2DROP EXIT
    THEN
    DROP R> 2R> ROT
  REPEAT
  DROP 2DROP R> DROP R:FAIL
;

```

A.XY Informal Appendix

A.XY.1 POSTPONE

POSTPONE compiles the data returned by DO-RECOGNIZER (i*x) into the dictionary as literal(s) and appends the compilation action of the R:TABLE information token. Later at run-time the i*x data is read back and the compilation action is performed like it would have been called directly at compile time.

```
: POSTPONE ( "name" -- )
  PARSE-NAME FORTH-RECOGNIZER DO-RECOGNIZER DUP >R
  R>POST EXECUTE R> R>COMP COMPILE, ;
```

This implementation assumes a system that uses recognizers only.

A.XY.2 Example Recognizer

The first example looks up the dictionary for the word and returns the execution token and the header flags if found. The data processing is the usual interpret/compile action. The Compile actions checks for immediacy and act accordingly. A portable postpone action is not possible. Amforth and gforth do it in a system specific way.

```
\ find-word is close to FIND but takes addr/len as input
256 BUFFER: find-word-buf \ counted string
: place ( c-addr1 u c-addr2 ) 2DUP C! CHAR+ SWAP MOVE ;
: find-word ( addr len -- xt +/-1 | 0 )
  find-word-buf place find-word-buf
  FIND DUP 0= IF NIP THEN ;

: immediate? ( flags -- true/false ) 0> ;
:NONAME ( i*x XT flags -- j*y ) \ INTERPRET
  DROP EXECUTE ;
:NONAME ( XT flags -- ) \ COMPILE
  immediate? IF COMPILE, ELSE EXECUTE THEN ;
:NONAME POSTPONE 2LITERAL ; ( XT flag -- )
RECOGNIZER: R:WORD

: REC:WORD ( addr len -- XT flags R:WORD | R:FAIL )
  find-word ( addr len -- XT flags | 0 )
  ?DUP IF R:WORD ELSE R:FAIL THEN
;
;
```

The second example deals with floating point numbers. The interpret action is a do-nothing since there is nothing that has to be done in addition to what the parsing word already did. The compile action takes the floating point number from the FP stack and compiles it to the dictionary. Postponing numbers is not defined, thus the postpone action here is printing the number and throwing an exception.

```
:NONAME ; ( -- ) ( F: f -- f ) \ INTERPRET
:NONAME POSTPONE FLITERAL ; ( -- ) ( F: f -- ) \ COMPILE
:NONAME FS. -48 THROW ; ( -- ) ( F: f -- ) \ POSTPONE
RECOGNIZER: R:FLOAT

: REC:FLOAT ( addr len -- R:FLOAT | R:FAIL ) ( F: -- f | )
  >FLOAT IF R:FLOAT ELSE R:FAIL THEN ;
```

A.XY.3 Text Interpreter

The Forth text interpreter can be changed into a generic tool that is capable to deal with any data type. It maintains `STATE` and calls the data processing methods according to it. The example is a full replacement if all necessary recognizers are available.

The algorithm of the Forth text interpreter as described in section 3.4 is modified. All subsections of 3.4 apply unchanged. Change the steps b) and c) from section 3.4 to make them optional, they can be performed with recognizers. Replace the step d) with the following steps d) to f)

- d. For each element of the recognizer stack provided by `FORTH-RECOGNIZER`, starting with the top element, call its parsing method with the sub-string "name" from step a).

Every parsing method returns an information token and the parsed data from the analyzed sub-string if successful. Otherwise it returns the system provided failure token `R:FAIL` and no further data.

Continue with the next element in the recognizer stack until either all are used or the information token returned from the parsing word is not the system provided failure token `R:FAIL`.

- e. Use the information token and do one of the following

1. if interpreting execute the interpret method associated with the information token.
2. if compiling execute the compile method associated with the information token.

- f. Continue with a)

```

: INTERPRET
  BEGIN
    PARSE-NAME DUP
  WHILE
    FORTH-RECOGNIZER DO-RECOGNIZER
    STATE @ IF R>COMP ELSE R>INT THEN \ get the right XT from R:*
    EXECUTE \ do the action.
    ?STACK \ simple housekeeping
  REPEAT 2DROP
;

```

A.XY.4 Naming Conventions

A Forth system that uses recognizers in the core has words for numbers and dictionary look-ups. These recognizers are useful for other data formats and use cases as well. They shall be named identically as shown in the table:

Name	Stack items	Comment
<code>R:NUM</code>	(-- n R:NUM)	single cell numbers, based on <code>>NUMBER</code>
<code>R:DNUM</code>	(-- d R:DNUM)	double cell numbers, based on <code>>NUMBER</code>
<code>R:FLOAT</code>	(-- R:FLOAT) (F: -- f)	floating point numbers, based on <code>>FLOAT</code>
<code>R:WORD</code>	(-- XT flags R:WORD)	words from the dictionary, <code>FIND</code>
<code>R:NAME</code>	(-- NT R:NAME)	words from the dictionary, with name tokens <code>NT</code>

The matching parsing words should be available as

Name	Stack effect
<code>REC:NUM</code>	(addr len -- n R:NUM d R:DNUM R:FAIL)
<code>REC:FLOAT</code>	(addr len -- R:FLOAT R:FAIL) (F: -- f)

REC:WORD	(addr len -- XT flags R:WORD R:FAIL)
REC:NAME	(addr len -- NT R:NAME R:FAIL)

Experience

First ideas to dynamically extend the Forth text interpreter were published in 2005 at comp.lang.forth by Josh Fuller and J Thomas: [Additional Recognizers?](#)

A specific solution to deal with number prefixes was roughly sketched by Anton Ertl at comp.lang.forth in 2007 with <https://groups.google.com/forum/#!msg/comp.lang.forth/r7Vp3w1xNus/Wre1BaKeCvcJ>

There are a number of specific solutions that can at least partly be seen as recognizers in various Forth's:

- prefix-detection in ciforth
- W32Forth uses its "chain" concept to achieve similar effects.
- various commercial Forth's seem to have ways to extent the interpreter.

A first generic recognizer concept was implemented in amforth version 4.3 (May 2011). The design presented in this RFD is implemented with version 5.3 (May 2014). gforth has recognizers since 2012, the ones described here since June 2014.

Existing recognizers cover a wide range of data formats like floating point numbers and strings. Others mimic the back-tick syntax used in many Unix shells to execute OS sub-process. A recognizer is used to implement OO notations.

Most of the small words that constitute a recognizer don't need a name actually since only their execution tokens are used. For the major words a naming convention is suggested: REC:<name> for the parsing word of the recognizer "name", and R:<name> for the information token word created with RECOGNIZER: for the data type "name".

There is no REC:FAIL that would be the companion of the system provided R:FAIL. It's simply

```
: REC:FAIL ( addr len -- R:FAIL )
  2DROP R:FAIL ;
```

That way, REC:FAIL can be seen as the parsing word of the recognizer that is always present as the last one in the recognizer stack and that cannot be deleted.

Test cases

The hardest and ultimate test case is to use the interpreter with recognizers enabled. Some parts can be tested separately, however.

```
T{ s" unknown word" FORTH-RECOGNIZER DO-RECOGNIZER -> R:FAIL }T
T{ s" 1234" FORTH-RECOGNIZER DO-RECOGNIZER -> 1234 R:NUM }T
T{ s" 1234." FORTH-RECOGNIZER DO-RECOGNIZER -> 1234. R:DNUM }T
T{ s" DO-RECOGNIZER" FORTH-RECOGNIZER DO-RECOGNIZER -> ' DO-RECOGNIZER -1 R:WORD }T
```

The system provided recognizers, if available, work as follows:

```
T{ s" 1234" REC:NUM -> 1234 R:NUM }T
T{ s" 1234." REC:NUM -> 1234. R:DNUM }T
T{ s" DO-RECOGNIZER" REC:WORD -> ' DO-RECOGNIZER -1 R:WORD }T
```

Discussion / Rationale

GET/SET-RECOGNIZERS

These commands can create a deep data stack usage. They are modeled after the well established `GET-/SET-ORDER` and `N>R/NR>` word pairs.

An alternative solution are words following `>R` and `R>`. Likewise a `>RECOGNIZER` would put the new item on the top of the recognizer stack. Since this element is processed first in `DO-RECOGNIZER`, this action prepends to the recognizer stack. Having the recognizer loop acting the other way (bottom up) is confusing and therefore not an option too. Furthermore I expect that most changes to the recognizer stack take place at the *end* (bottom) of it appending a new recognizer. There is no commonly agreed way to change a stack at its bottom. Even more difficult is an insert or remove operation of a recognizer in the middle. Again the standard data stack words are the simplest way to do it. Since the recognizer stack is smaller than the data stack and stack changes are expected to happen seldom the proposed solution is considered the simplest solution.

POSTPONE

Adding the `POSTPONE` method has been seen as overly complex. At least with the current standard text it is necessary however. One reason is that `POSTPONE` has a lot of special cases which cannot be implemented without system knowledge. The postpone method carries this information for all data types. Recent discussions indicate that this may be solved cleanly in a future version of Forth, until this discussion is finished, a separate postpone action is the only way to implement what recognizers can achieve.

Bernd Paysan wrote in `clf`

Concerning the postpone action and `'` and `[']` using recognizers: IMHO, there's not much point in generating a super-efficient postpone, but you can use `'` and `[']` together with literals, if the postpone method is modified to *only* contain the work to save the `i*x` part of the recognizer output into the dictionary. The remaining action of postpone is generic. So `POSTPONE` executes the literal-append part of the `r:table` and then appends the `r:table` as literal and the compilation part of the `r:table`.

`'` and `[']` can check if the literal-append part is empty (a noop), and if not, create a quotation that contains that literal, and appends the `_r>int` part of the table. I.e. `['] 3` becomes something like `[: 3 noop ;]`, with an easy opportunity to optimize away the noop.

This is not mandatory, but I'd like to implement it that way. And that means the postpone part has to be changed to the essential core (the handling of the recognizer-specific `i*x`), and the rest is done by `POSTPONE`.

That means `r:num` is defined as

```
: lit, postpone literal ;
' noop ' lit, ' lit, recognizer: r:num
```

and `POSTPONE` is defined as

```
: POSTPONE ( "name" -- )
  PARSE-NAME FORTH-RECOGNIZER DO-RECOGNIZER >R
  R@ _r>post EXECUTE r> _r>comp COMPILE, ;
```

following your reference implementations.

This also makes the simple two-part table easier to implement, as *only* the compilation part (perform literal part+append interpretation part) needs to be generated.

From 6.1.2033 POSTPONE: "Append the compilation semantics of name to the current definition." This `postpone` does exactly this.

The suggested ' ' is part of the implementation and can be left to the system provider.

Multiword Parsing

The RFD suggests that the input stream is split into white-space delimited words as part of the general text interpreter. The parse actions of the recognizers get these words only.

A recognizer that deals with "sentences" (multiple words) needs more. It has to communicate back, where it finished its work so that subsequent parse action start at the right point. There are a few possibilities

- The input for recognizer comes from within `SOURCE` and is managed with `>IN`. That is the designated environment for recognizers. Systems are free to make a copy of the word before calling the parsing words from the recognizer. A multi-word recognizer nevertheless needs access the `SOURCE` buffer and changes `>IN` accordingly. It must not change the content of the string however.
- The input comes from an arbitrary string. `SOURCE` and `>IN` are not used. The word `DO-RECOGNIZER` has to tell now, how far it went in addition to the actual results. The standard already has a word that works that way: `>NUMBER (ud1 addr len -- ud2 addr' len')`. A similar `DO-RECOGNIZER` would have the stack effect `(addr len -- i*x addr' len' R:TABLE | R:FAIL)`.

Since many standard words are already grouped around `SOURCE` and `>IN` it seems to be overkill to maximize the flexibility. That's why option 1 is preferred. Furthermore it leads to simpler code and easier integration into existing systems. There is no dependency on `SOURCE` and `>IN` for the single-word recognizer use case.

Another aspect with multiword recognizers is that it is possible that the closing syntactic element of the multi-word sentence is not within the current input string One or more `REFILL` may be necessary to get it. Since that may be troublesome in the long run, the closing element shall be in the same input string as the opening one.

Keep the Interpreter

The Forth 2015 meeting in Bath as well as (earlier) Andrew Haley added the wish / requirement to keep the current interpreter and make recognizers an truly optional part. Changed in the proposal to make the Forth 2012 interpreter steps to search the dictionary (step b) and convert numbers (step c) optional. That way the current interpreter can work without changes and at the same time the hard coded steps b) and c) from section 3.4 could be replaced with recognizers. The recognizer steps are added as step d) to f) It should be clear that the example implementation of the interpreter is not mandatory.

Nevertheless the full power of the concept cannot be achieved with such a two-class interpreter. For that, one need to be able to replace the standard actions `FIND` and number recognition too.

As a related change the words `R>COMP`, `R>INT` and `R>POST` became part of the proposal since they are needed to write an interpreter and similar words portably.

Switching Recognizer Stacks

The Forth 2015 meeting wishes a possibility to switch between prepared recognizer stacks. To achieve this, the words `set-recognizers`, `do-recognizer` and `get-recognizers` are changed to have an additional parameter `stack-id` that identifies the recognizer stack to be used. The elements of the recognizer stack may not be accessible with the normal fetch and store operation, the numeric value of the `stack-id` is implementation defined. The stack may have a limited size too resulting in an error condition if the maximum size is exceeded.

The new word `FORTH-RECOGNIZER` is introduced to have a global (drift) anchor to provide a common starting point to be used by various words like `EVALUATE` from whom a consistent behavior is expected. It is a `VALUE` to switch the whole stack at once.

Nesting Recognizer Stacks

An extension of the [Switching Recognizer Stacks](#).

Example is a number recognizer. Instead of checking for both single and double numbers, only one type is checked. All number checks are collected in the `recstack:numbers` recognizer stack.

```
: REC:SNUM ( addr len -- n R:NUM | R:FAIL )
...
;
: REC:DNUM ( addr len -- d R:DNUM | R:FAIL )
...
;

2 RECOGNIZER CONSTANT recstack:numbers

' REC:SNUM ' REC:DNUM 2 recstack:numbers SET-RECOGNIZERS

: REC:NUM ( addr len -- n R:NUM | d R:DNUM | R:FAIL )
  recstack:numbers DO-RECOGNIZER
;

' REC:NUM ' REC:WORD 2 FORTH-RECOGNIZER SET-RECOGNIZERS
```

Flags, R:FAIL or Exceptions

The `R:FAIL` word has two purposes. One is to deliver a boolean information whether a parsing word could deal with a word. The other task is the method table of for the interpreter to actually handle the parsed data, this time by generating a proper error message and to leave the interpreter. While the method table simplifies the interpreter loop, the flag information seems to be odd. On the other hand a comparison of the returned `R:*` token with the constant `R:FAIL` can be easily optimized.

A completely different approach is using exceptions to deliver the flag information from `DO-RECOGNIZER` to its callers. Using them requires the exception word set, which may not be present on all systems. In addition, an exception is a somewhat elaborate error handling tool and usually means that something unexpected has happened. Matching a string to a sequence of patterns means that exceptions are used in a normal flow of compare operations.

That `R:FAIL` is used in two ways is an optimization. The flag information can be carried with the equation `R:* R:FAIL <>` as well.

That there is no final `REC:FAIL` in the recognizer stack is an optimization too. Earlier versions of the recognizer concept did have such a bottom element. It turned out that it caused a lot of trouble. If it got deleted, the interpreter loop did not recognize this as an error and crashed without further notice. To circumvent this situation, the current recognizer stack depth is needed. Adding a check for an empty recognizer stack was more code. The second argument against is that adding a recognizer to the recognizer becomes more complex since there is a bottom element, that has to be kept, essentially making appending a recognizer always an insert-in-the-middle action.

R:FAIL necessity

Compare different implementations. Esp the dual use as a flag and a information is discussed with code examples. Exception are not an option as already discussed.

Parse `REC:*` actions

For simplicity the float recognizer. Others work likewise.

With `R:FAIL`

```
: REC:FLOAT ( addr len -- R:FLOAT | R:FAIL ) ( F: -- f )
>FLOAT IF R:FLOAT ELSE R:FAIL THEN ;
```

Without R:FAIL

```
: REC:FLOAT ( addr len -- ( R:FLOAT | 0 ) ( F: -- f )
>FLOAT IF R:FLOAT ELSE 0 THEN ;
```

almost the same.

DO-RECOGNIZER

with R:FAIL

```
: DO-RECOGNIZER ( addr len stack-id -- i*x R:TABLE | R:FAIL )
  DUP >R @
  BEGIN
  DUP
  WHILE
  DUP CELLS R@ + @
  2OVER 2>R SWAP 1- >R
  EXECUTE DUP R:FAIL <> IF
  2R> 2DROP 2R> 2DROP EXIT
  THEN DROP R> 2R> ROT
  REPEAT
  DROP 2DROP R> DROP R:FAIL
;
```

Without R:FAIL

```
: DO-RECOGNIZER ( addr len stack-id -- i*x R:TABLE | 0 )
  DUP >R @
  BEGIN
  DUP
  WHILE
  DUP CELLS R@ + @
  2OVER 2>R SWAP 1- >R
  EXECUTE DUP IF
  2R> DROP 2R> 2DROP EXIT
  THEN DROP R> 2R> ROT
  REPEAT
  DROP 2DROP R> DROP R:FAIL
;
```

again almost the same.

POSTPONE

with R:FAIL

```
: POSTPONE ( "name" -- )
  PARSE-NAME FORTH-RECOGNIZER DO-RECOGNIZER DUP >R
  R>POST EXECUTE R> R>COMP COMPILE, ;
```

without R:FAIL

```

: POSTPONE ( "name" -- )
  PARSE-NAME FORTH-RECOGNIZER DO-RECOGNIZER
  ?DUP IF
    DUP >R
    R>POST EXECUTE R> R>COMP COMPILE,
  ELSE
    NOT-RECOGNIZED
  THEN ;

```

special casing "not-recognized" and slightly more complex due to NOT-RECOGNIZED.

Interpreter

With R:FAIL

```

: INTERPRET
  BEGIN
    PARSE-NAME DUP
  WHILE
    FORTH-RECOGNIZER DO-RECOGNIZER
    STATE @ IF R>COMP ELSE R>INT THEN \ get the right XT from R:*
    EXECUTE \ do the action.
    ?STACK \ simple housekeeping
  REPEAT 2DROP
;

```

Without R:FAIL

```

: INTERPRET
  BEGIN
    PARSE-NAME DUP
  WHILE
    FORTH-RECOGNIZER DO-RECOGNIZER
    ?DUP IF \ we got an R:* table
      STATE @ IF R>COMP ELSE R>INT THEN \ get the right XT from R:*
      EXECUTE \ do the action.
    ELSE
      \ no recognizer did the job
      NOT-RECOGNIZED
    THEN
      ?STACK \ simple housekeeping
  REPEAT 2DROP
;

```

Like POSTPONE Special casing the "not-found" condition and slightly more complex due to NOT-RECOGNIZED.

Adapting the special case "not recognized" requires to extent the text interpreter specification too.

Conclusion: R:FAIL is essential since it simplifies both the concept and the implementation by not special casing any result. Furthermore the code for the recognizers is easier to read and understand: R:FAIL vs 0. "Notation matters".

Use Cases

Name Tokens

Name Tokens (nt) are part of the Forth 2012 Programming Tools word set. This section is just a use case description deploying an optional word set.

The words found in the dictionary with FIND return the execution token and the immediate flag. Using the Programming Tools word set, the dictionary look-up can be made based on TRAVERSE-WORDLIST with a recognizer called e.g. REC:NAME (addr len -- nt R:NAME|R:FAIL). The major difference to FIND is that all header information is available to handle the token:

```
:NONAME NAME>INTERPRET EXECUTE ; ( nt -- ) \ interpret
:NONAME NAME>COMPILE EXECUTE ; ( nt -- ) \ compile
:NONAME POSTPONE LITERAL ; ( nt -- ) \ postpone
RECOGNIZER: R:NAME
```

The actual REC:NAME is slightly more complex and usually benefits from system knowledge.

```
\ the analogon to search-wordlist
: search-name ( addr len wid -- nt | 0 )
  >R 0 \ used as flag inside the following quotation
  [: ( addr len flag nt -- addr len false true | nt false )
    >R DROP 2DUP R@ NAME>STRING COMPARE
    IF R> DROP 0 -1 ELSE 2DROP R> 0 THEN
  ;] R> TRAVERSE-WORDLIST ( -- addr len false | nt )
  DUP 0= IF NIP NIP THEN
;

\ a single wordlist is checked
: (rec:name) ( addr len wid -- nt R:NAME | R:FAIL )
  search-name
  ?DUP IF R:NAME ELSE R:FAIL THEN
;

\ checks only the standard word-list
: rec:name ( addr len -- nt R:NAME | R:FAIL )
  FORTH-WORDLIST (rec:name)
;
```

Search Order Word Set

A large part of the Search Order word set is close to what recognizers do while dictionary searches. The order stack can be seen as a subset of the recognizer stack.

The words dealing with the order stack (ALSO, PREVIOUS, FORTH, ONLY etc) may be extended/changed to handle the recognizer stack too/instead. On the other hand, ALSO is essentially DUP on a different stack. ONLY and FORTH set a predefined stack content. With the GET/SET-RECOGNIZERS words all changes can be prepared on the data stack with the usual data stack words.

A complete redesign of the Search Order word set affects many programs, worth an own RFD. The common tools to actually implement both recognizer and search order word sets may be useful for themselves.

Completely unrelated is SET/GET-CURRENT. Recognizers don't deal with the places, new words are put into. Possible changes here are not considered part of the recognizer word set proposal.

Stateless interpreter

An implementation of the interpreter without an explicit `STATE`. For legacy applications a `STATE` variable is maintained but not used.

The code depends on `DEFER` and `IS` from `CORE EXT`. Similar code can be found in `gforth` and `win32forth`.

```
\ legacy state support
VARIABLE STATE
: on ( addr -- ) -1 SWAP ! ;
: off ( addr -- ) 0 SWAP ! ;

\ the two modes of the interpreter
: (interpret-i) R>INT EXECUTE ;
: (interpret-c) R>COMP EXECUTE ;
DEFER (interpret) ' (interpret-i) IS (interpret)

\ switch interpreter modes
: ] STATE on [' (interpret-c) IS (interpret) ;
: [ STATE off [' (interpret-i) IS (interpret) ; IMMEDIATE

: interpret
  BEGIN
    PARSE-NAME DUP \ get something
  WHILE
    FORTH-RECOGNIZER DO-RECOGNIZER \ analyze it
    (interpret) \ act on it
    ?stack \ simple housekeeping
  REPEAT 2DROP
;
```

Not-Found Hooks

Many systems have a not-found hook that is called if a word is not found and is not a number. This hook is usually a deferred word. With recognizers it can be implemented as follows:

```
: throw-13 -13 THROW ;

DEFER interpret-notfound ( addr u -- )
  ' throw-13 IS interpret-notfound
DEFER compiler-notfound ( addr u -- )
  ' throw-13 IS compiler-notfound
DEFER postpone-notfound ( addr u -- )
  ' throw-13 IS postpone-notfound

' interpret-notfound
' compiler-notfound
' postpone-notfound
RECOGNIZER: rec:notfound
```

With that recognizer places at the end (bottom) of the recognizer stack the final action if a word could not be handled is a set of words that can be changed independently. These hooks are most useful for existing code that uses the `not-found` deferred word API. (Idea and basic code structure taken from `gforth`).

' and [']

' (tick) and its companion ['] (bracket-tick) are affected too. It is common practice that the sequence ' foo execute does the same as calling foo directly (in interpret mode). Now consider special recognizer that searches an otherwise hidden word-list (think of name spaces). Words from it may be interpreted and compiled without problems, but could not be found with '. Therefore it is desirable to use the recognizer stack here too. The difficulty here is to decide whether a recognized item is an executable "tick-able" word. E.g. numbers and compile-only words are not.

Implementation requires system specific knowledge. The following code depends on R:WORD to work.

```
: executable? ( R:TABLE -- f )
    R>INT \ get the interpretation action for the given token
    R:WORD R>INT \ get the system specific interpret action
    =
;

: ' ( "<spaces>name" -- XT )
    PARSE-NAME FORTH-RECOGNIZER DO-RECOGNIZER
    executable? 0= IF
        \ call the system specific error action "invalid tick"
        -13 THROW
    THEN
    DROP \ remove the immediate flag
    \ the XT from the R:WORD result set is left
;
```

Older Remarks

2-Method API

Anton Ertl suggested an alternative implementation of recognizers. Basically all text data is converted into a literal at parse time. Later the interpreter decides whether to execute or compile the literal data, depending on STATE. POSTPONE is a combination of storing the literal data together with their compile time action.

```
interpretation: conv final-action
compilation:    conv literal-like postpone final-action
postpone:
    conv literal-like postpone literal-like postpone final-action
```

The conv-action is what is done inside the DO-RECOGNIZER action (REC:* words) and the literal-like and final-action set replaces the proposed 3 method set in R:*. It is not yet clear whether this approach covers the same range of possibilities as the proposed one or may solve the tick-problem mentioned above. Another side effect is that postponing literals like numbers becomes possible without further notice.

For simple use cases (literals) it's possible to automatically convert this approach into the 3-method API (Anton Ertl and Bernd Paysan):

```
: rec-methods { : literal-xt final-xt -- interpret-xt compile-xt postpone-xt : }
    final-xt
    :noname literal-xt compile, final-xt ]] literal compile, ; [[ dup >r
    :noname literal-xt compile, r> compile, postpone ;
;
```

With that command, the standard number recognizer can be rewritten as

```
\ numbers
:NONAME ; \ final-action do nothing
' LITERAL \ literal-action
rec-methods RECOGNIZER: R:NUM
```

Anton Ertl writes in comp.lang.forth:

If you define recognizers through these components, you don't need to specify the three components, in particular not a POSTPONE action; and yet POSTPONEing literals works as does any other POSTPONEing of recognizers. With that, one might leave it up to systems whether they support POSTPONEing recognizers or not.

Disadvantage: Does not combine with doing the dictionary look-up as a recognizer for immediate words:

If you make the immediate word a parse-time action with a noop for literal-like and noop for run-time, it works correctly for interpretation and compilation, but not for POSTPONE. And anything else is even further from the desired behaviour. One could enhance this scheme to support immediate words correctly, but I don't see a clean way to do that.

So there seems to be a choice:

1. Compose the behaviour of recognizers of these components, but do not treat the dictionary as a recognizer.
2. Treat the dictionary as a recognizer, but build recognizers from interpretation, compilation, and postponeing behaviour.

A complete reference implementation does not exist, many aspects were published at comp.lang.forth by Jenny Brien.

Acknowledgments

The following people did major or minor contributions, in no particular order.

- Bernd Paysan
- Jenny Brien
- Andrew Haley
- Alex McDonald
- Anton Ertl