# Forth Recognizer -- Request For Discussion

**Author:**   Matthias Trute

**Contact:**   mtrute@web.de

**Version:**   4

**Date:**   23 July 2018

**Status:**   Final (Committee Supported Proposal)

# 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.

- **2018-07-23 Version 4 - Clarifications, Fixing typos, added test cases**

    - 2016-09-18 Added more test cases

    - 2016-09-25 Clarify that `>IN` is unchanged for an `REC-FAIL` (`RECTYPE-NULL`) result.

    - 2016-10-21 simpler reference implementation

    - 2016-11-05 first attempt to rename keywords and concept names

    - 2017-05-15 discussion of `LOCATE`

    - 2017-08-08 move example recognizers to discussion/rationale section.

    - 2017-09-12 renamed keywords in XY.6.1 as suggested by the Forth 200x committee

    - 2017-12-06 changed wording from "recognizer stack" to "recognizer sequence".

    - 2017-12-10 created Recognizer EXT section with recognizer sequence management words.

    - 2018-04-09 expanded EXT section with RECTYPE* words

    - 2018-05-11 add comments about `recognizable?`

    - 2018-07-23 finalized

# Background

I'm working on a Forth for 8-bit micro-controllers for more than 10 years now (amforth.sf.net). It is a useful tool for serious work and at the same time a nice playground for 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. 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 the running system to make it

work 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 that implement recognizers. A third implementation is in the proposal.

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 and at theforth.net.

# 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 number formats.

It's not possible to use the Forth text interpreter in an application or system extension context. Most interpreters in existing systems use a number of hooks to extent the interpreter. 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.

Furthermore applications need to use system provided and system specific words or have to re-invent the wheel to get numbers with a sign or hex numbers with the $ prefix. The building blocks (FIND, COMPILE,, >NUMBER etc) are available but there is a gap between them and what the Forth interpreter already does.

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. These 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 (sub-strings of `SOURCE`) if they match the criteria for a certain data type. These parsing words are grouped to achieve an order of invocation. The result of the parsing words is handed over to the interpreter with 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.

This proposal is about the building blocks.

# Proposal

## XY. The optional Recognizer word set

### XY.1 Introduction

The recognizer concept consists of two elements: parsing words that return data type information 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 data type information. A particular data type information can be used by different parsing words.

A system provided data type information is called `RECTYPE-NULL`. It is used if no other one is applicable. This 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.

A parsing word within the recognizer concept has the stack effect

```
REC-SOMETYPE ( addr len -- i*x RECTYPE-SOMETYPE | RECTYPE-NULL )
```

The parsing word must not change the string. Since it is called from the interpreter, it may access `SOURCE` and, if applicable, change `>IN`. If `>IN` is not used, any string may serve as input, otherwise "addr/len" is assumed to be a substring of the buffer `SOURCE`.

"i*x" is the result of the parse action of the string "addr/len". `RECTYPE-SOMETYPE` is the data type id 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 from the parsing word and are generally expected to consume it. They can have additional stack effects, depending on what `RECTYPE-SOMETYPE-METHOD` actually does.

```
RECTYPE-SOMETYPE-METHOD ( ... i*x -- j*y )
```

The data "i*x" doesn't have to be on the data stack, it 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 `RECTYPE-SOMETYPE` information only.

## XY.2 Additional terms and notations

**Data type id**

A cell sized number. It identifies the data type and a method set to perform the data processing in the text interpreter. The actual numeric value is system specific.

**Recognizer**

A string parsing word that returns a data type id together with the parsed data if successful. The string parsing word is assumed to run within the Forth interpreter and can access `SOURCE` and `>IN`.

**Recognizer Sequence**

An ordered set of recognizers. It is identified with a cell sized numeric id.

## XY.3 Additional usage requirements

### XY.3.1 Data type id

A data type id is a single cell value that identifies a certain data type. Append table the following table to table 3.1

| Symbol | Data type | Size on Stack |
|--------|-----------|---------------|
| dt | data type id | 1 cell |

## 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.

### 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

**FORTH-RECOGNIZER ( -- rec-seq-id ) RECOGNIZER**

A system VALUE with a recognizer sequence id.

It is `VALUE` that can be changed with `TO` to assign a new recognizer set. This change has immediate effect.

This recognizer set shall be used in all system level words like `EVALUATE, LOAD` etc.

**RECOGNIZE ( addr len rec-seq-id -- i*x RECTYPE-DATATYPE | RECTYPE-NULL ) RECOGNIZER**

Apply the string at "addr/len" to the elements of the recognizer set identified by `rec-seq-id`. Terminate the iteration if either a parsing word returns a data type id that is different from `RECTYPE-NULL` or the set is exhausted. In this case return `RECTYPE-NULL`.

"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.

**RECTYPE>COMP ( RECTYPE-DATATYPE -- XT-COMPILE ) RECOGNIZER**

Return the execution token for the compilation action from the recognizer date type id.

**RECTYPE>INT ( RECTYPE-DATATYPE -- XT-INTERPRET ) RECOGNIZER**

Return the execution token for the interpretation action from the recognizer data type id.

**RECTYPE>POST ( RECTYPE-DATATYPE -- XT-POSTPONE ) RECOGNIZER**

Return the execution token for the postpone action from the recognizer data type id.

**RECTYPE-NULL ( -- RECTYPE-NULL ) RECOGNIZER**

The null data type id. It is to be used if no other data type id is applicable but one is needed. Its associated methods perform system specific error actions. The actual numeric value is system dependent.

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

Skip leading space delimiters. Parse name delimited by a space. Create a data type id under the name `name` and associate the three execution tokens.

The words for XT-INTERPRET, XT-COMPILE and XT-POSTPONE are called with the parsed data `i*x` that e.g. `RECOGNIZE` has returned.

The word behind XT-INTERPRET shall have the stack effect ( `... i*x -- j*y` ). The words behind XT-COMPILE and XT-POSTPONE shall consume `i*x`.

### YZ.6.2 Recognizer Extension Words

A Forth system that uses recognizers in the core has words for numbers and dictionary look-ups. They shall be named as shown in the table:

| Name | Stack effect |
|------|--------------|
| REC-NUM | ( addr len -- n RECTYPE-NUM \| d RECTYPE-DNUM \| RECTYPE-NULL ) |
| REC-FLOAT | ( addr len -- RECTYPE-FLOAT \| RECTYPE-NULL ) (F: -- f \| ) |
| REC-FIND | ( addr len -- XT +/-1 RECTYPE-XT \| RECTYPE-NULL ) |
| REC-NT | ( addr len -- NT RECTYPE-NT \| RECTYPE-NULL ) |

The recognizer type names, if available, shall be as shown in the table below:

| Name | Stack items | Comment |
|------|-------------|---------|
| RECTYPE-NUM | ( -- n RECTYPE-NUM) | single cell number |
| RECTYPE-DNUM | ( -- d RECTYPE-DNUM) | double cell number |
| RECTYPE-FLOAT | ( -- RECTYPE-FLOAT)(F: -- f ) | floating point number , |
| RECTYPE-XT | ( -- XT +/-1 RECTYPE-XT) | word from the dictionary matching FIND |
| RECTYPE-NT | ( -- NT RECTYPE-NT) | word from the dictionary with name token NT |

The following words deal with changing and creating recognizer sequences.

**GET-RECOGNIZER ( rec-seq-id -- rec-n .. rec-1 n ) RECOGNIZER EXT**

Copy the recognizer sequence `rec-1 .. rec-n` to the data stack. The element `rec-1` is the first in the sequence.

The source is unchanged.

**SET-RECOGNIZER ( rec-n .. rec-1 n rec-seq-id -- ) RECOGNIZER EXT**

Replace the recognizer sequence identified by `rec-seq-id` with a new set of `n` recognizers `rec-x`.

If the capacity of the destination sequence is too small to hold all new elements, an ambiguous situation arises.

**NEW-RECOGNIZER-SEQUENCE ( size .. rec-seq-id ) RECOGNIZER EXT**

Create a new, empty recognizer sequence with at least `size` elements.

## *XY.7 Reference Implementation*

Basic recognizer sequence module. It is implemented as a separate stack.

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

: SET-STACK ( item-n .. item-1 n stack-id -- )
  2DUP ! CELL+ SWAP CELLS BOUNDS
  ?DO I ! CELL +LOOP ;

: GET-STACK ( stack-id -- item-n .. item-1 n )
   DUP @ >R R@ CELLS + R@ BEGIN
     ?DUP
   WHILE
     1- OVER @ ROT CELL - ROT
   REPEAT
   DROP R> ;
```

The recognizer sequence uses the stack module. Hence the stack-id becomes the rec-seq-id.

```
: NEW-RECOGNIZER-SEQUENCE STACK ;
: SET-RECOGNIZER SET-STACK ;
: GET-RECOGNIZER GET-STACK ;

\ create the default recognizer sequence
4 NEW-RECOGNIZER-SEQUENCE VALUE FORTH-RECOGNIZER

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

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

\ the null token
:NONAME -1 ABORT" FAILED" ; DUP DUP RECTYPE: RECTYPE-NULL

\ depends on the stack implementation
: RECOGNIZE   ( addr len rec-seq-id -- i*x RECTYPE-SOMETYPE | RECTYPE-NULL )
    DUP >R @
    BEGIN
      DUP
```

```
    WHILE
      DUP CELLS R@ + @
      2OVER 2>R SWAP 1- >R
      EXECUTE DUP RECTYPE-NULL <> IF
        2R> 2DROP 2R> 2DROP EXIT
      THEN
      DROP R> 2R> ROT
    REPEAT
    DROP 2DROP R> DROP RECTYPE-NULL
;
```

# A.XY Informal Appendix

## A.XY.1 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 sequence 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 RECTYPE-NULL and no further data.

 Continue with the next element in the recognizer set until either all are used or the information token returned from the parsing word is not the system provided failure token RECTYPE-NULL.

 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 RECOGNIZE
      STATE @ IF RECTYPE>COMP ELSE RECTYPE>INT THEN
      EXECUTE
      ?STACK  \ simple housekeeping
  REPEAT 2DROP
;
```

## A.XY.2 POSTPONE

POSTPONE compiles the data returned by RECOGNIZE (i*x) into the dictionary as literal(s) and appends the compilation action of the RECTYPE-TOKEN data type id. 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 RECOGNIZE DUP >R
  RECTYPE>POST EXECUTE R> RECTYPE>COMP COMPILE, ;
```

This implementation assumes a system that uses recognizers only.

### A.XY.3 Test Cases

The test cases assume a stack to implement the recognizer set.

```
T{ 4 NEW-RECOGNIZER-SEQUENCE constant RS -> }T

T{ :NONAME 1 ;  :NONAME 2 ;  :NONAME 3  ; RECTYPE: rectype-1 -> }T
T{ :NONAME 10 ; :NONAME 20 ; :NONAME 30 ; RECTYPE: rectype-2 -> }T

T{ : rec-1 NIP 1 = IF rectype-1 ELSE RECTYPE-NULL THEN ; -> }T
T{ : rec-2 NIP 2 = IF rectype-2 ELSE RECTYPE-NULL THEN ; -> }T

T{ rectype-1 RECTYPE>INT  EXECUTE -> 1 }T
T{ rectype-1 RECTYPE>COMP EXECUTE -> 2 }T
T{ rectype-1 RECTYPE>POST EXECUTE -> 3 }T

\ testing RECOGNIZE
T{         0 RS SET-RECOGNIZER -> }T
T{ S" 1"     RS RECOGNIZE   -> RECTYPE-NULL }T
T{ ' rec-1 1 RS SET-STACK -> }T
T{ S" 1"     RS RECOGNIZE   -> rectype-1 }T
T{ S" 10"    RS RECOGNIZE   -> RECTYPE-NULL }T
T{ ' rec-2 ' rec-1 2 RS SET-STACK -> }T
T{ S" 10"    RS RECOGNIZE   -> rectype-2 }T
```

The dictionary lookup has the following test cases

```
T{ S" DUP" REC-FIND  -> ' DUP -1 RECTYPE-XT }T
T{ S" UNKOWN WORD" REC-FIND -> RECTYPE-NULL }T
```

The number recognizer has the following checks

```
VARIABLE OLD-BASE BASE @ OLD-BASE !

T{ S" 1234"    REC-NUM -> 1234 RECTYPE-NUM }T
T{ S" 1234."   REC-NUM -> 1234. RECTYPE-DNUM }T
T{ S" %-10010110" REC-NUM -> -150 RECTYPE-NUM }T
T{ S" %10010110"  REC-NUM ->  150 RECTYPE-NUM }T
T{ S" 'Z'"     REC-NUM -> char Z RECTYPE-NUM }T
T{ S" ABCXYZ" REC-NUM -> RECTYPE-NULL }T

\ check whether BASE is unchanged
T{ BASE @ OLD-BASE @ = -> -1 }T
```

# 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.
- FICL, a system close to Forth, has parse-steps since approx 2001.

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, and `RECTYPE-<name>` for the data type word created with `RECTYPE:` for the data type "name".

# 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-1234 S" 1234" ; -> }T
T{ : D-1234 S" 1234." ; -> }T
T{ : S-UNKNOWN S" unknown word" ; -> }T
T{ : S-DUP  S" DUP" ; -> }T

T{ S-1234 FORTH-RECOGNIZER RECOGNIZE -> 1234  RECTYPE-NUM   }T
T{ D-1234 FORTH-RECOGNIZER RECOGNIZE -> 1234. RECTYPE-DNUM  }T
T{ S-DUP  FORTH-RECOGNIZER RECOGNIZE -> ' DUP -1 RECTYPE-XT }T
T{ S-UNKNOWN FORTH-RECOGNIZER RECOGNIZE  -> RECTYPE-NULL }T
```

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

```
T{ S-DUP     REC-FIND -> ' DUP -1 RECTYPE-XT }T
T{ S-UNKNOWN REC-FIND -> RECTYPE-NULL }T
T{ S-1234    REC-FIND -> RECTYPE-NULL }T
T{ D-1234    REC-FIND -> RECTYPE-NULL }T

T{ S-UNKNOWN REC-NUM -> RECTYPE-NULL }T
T{ S-1234    REC-NUM -> 1234 RECTYPE-NUM }T
T{ D-1234    REC-NUM -> 1234. RECTYPE-DNUM }T
T{ S-DUP     REC-NUM -> RECTYPE-NULL }T
```

Floating point numbers are handled likewise

```
T{ : S-1234e5 S" 1234e5" ; -> }T
T{ S-1234e5 REC-FLOAT -> 1234e5 RECTYPE-FLOAT }
T{ S-1234e5 FORTH-RECOGNIZER RECOGNIZE -> 1234e5 RECTYPE-FLOAT }T
```

# Discussion / Rationale

This section reflects the years of discussion. Some parts of it may seem look odd. It may be wise to consult the earlier versions of this RFD. Despite they still apply.

## 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 a wrapper for FIND to use 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 ;

:NONAME ( i*x XT +/-1 -- j*y )  \ INTERPRET
  DROP EXECUTE ;
:NONAME ( XT +/-1 -- )          \ COMPILE
  0> IF COMPILE, ELSE EXECUTE THEN ;
:NONAME ( XT +/-1 -- )          \ POSTPONE
  POSTPONE 2LITERAL ;
RECTYPE: RECTYPE-XT

: REC-FIND ( addr len -- XT +/-1 RECTYPE-XT | RECTYPE-NULL )
    find-word ( addr len -- XT +/-1 | 0 )
    ?DUP IF RECTYPE-XT ELSE RECTYPE-NULL 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 (yet) part of the Forth standard, thus the postpone action here prints the number and throws an exception to enforce an error handling.

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

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

## Data Type Id's

Earlier revisions of this RFD called them information tokens. In fact they describe a data type (e.g. a number) and they point to a sequence of action for the actions of the Forth core system: interpret, compile and postpone. These tokens are not related to the process they created. They are only related to the data they are describing. Esp. the RECTYPE-NULL is in fact not a failure but a null information, just like the 0 (zero) is the logical FALSE information.

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 (partially modified with new keywords)

> 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:token` and then appends the `r:token` as literal and the compilation part of the `r:token`.

> `'` 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 `RECTYPE>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 `RECTYPE-NUM` is defined as

```
: lit, postpone literal ;
' noop ' lit, ' lit, RECTYPE: RECTYPE-NUM
```

> and `POSTPONE` is defined as

```
: POSTPONE ( "name" -- )
  PARSE-NAME FORTH-RECOGNIZER  RECOGNIZE  >R
  R@ RECTYPE>POST EXECUTE R> RECTYPE>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.

# Multi-word 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 `RECOGNIZE` 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 similiar `RECOGNIZE` would have the stack effect ( addr len -- i*x addr' len' RECTYPE-TOKEN | RECTYPE-NULL).

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 multi-word 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.

The Forth interpreter makes sure that `>IN` points to the first character after the addr/len string that is passed as input to the parsing words.

## Keep the Interpreter

The Euro-Forth 2015 meeting 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 `RECTYPE>COMP`, `RECTYPE>INT` and `RECTYPE>POST` became part of the proposal since they are needed to write an interpreter and similar words portably.

## Switching Recognizer Sets

On the Euro-Forth 2015 meeting the wish to switch between prepared recognizer sets came up. To achieve this, the word `RECOGNIZE` is changed to have an additional parameter `rec-set-id` that identifies the recognizer set to be used. The elements of the recognizer sequence may not be accessible with the normal fetch and store operation, the numeric value of the `rec-set-id` is implementation defined. The sequence data 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 sequence at once.

## Nesting Recognizer Sets

An extension of the Switching Recognizer Sets.

Example is a number recognizer. Instead of checking for all number formats from the Forth 2012 spec in one recognizer, every variant is handled by an individual one. All number checks are collected in the `rec-numbers` recognizer set.

```
\ 'y'
: REC-CHAR ( addr len -- n RECTYPE-NUM | RECTYPE-NULL )
  ....
;
\ single cell numbers
```

```
: REC-SNUM ( addr len -- n RECTYPE-NUM | RECTYPE-NULL )
  ...
;
\ double cell numbers
: REC-DNUM ( addr len -- d RECTYPE-DNUM | RECTYPE-NULL )
  ...
;

3 STACK CONSTANT rec-numbers

' REC-CHAR ' REC-SNUM ' REC-DNUM 3 rec-numbers SET-STACK

: REC-NUM ( addr len -- n RECTYPE-NUM | d RECTYPE-DNUM | RECTYPE-NULL )
  rec-numbers RECOGNIZE
;

' REC-NUM ' REC-FIND 2 FORTH-RECOGNIZER SET-STACK
```

## Flags, `RECTYPE-NULL` or Exceptions

The `RECTYPE-NULL` 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 `RECTYPE-*` token with the constant `RECTYPE-NULL` can be easily optimized.

A completely different approach is using exceptions to deliver the flag information from `RECOGNIZE` 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 than something unexpected has happened. Matching a string to a sequence of patterns means that exceptions are used in a normal sequence of compare operations. The third argument against exceptions is that if used for recognizers, they mandate too much implementation details on system providers which is not considered useful.

`RECTYPE-NULL` is used in two ways is an optimization. The flag information can be carried with the equation `RECTYPE-* RECTYPE-NULL <>` as well.

## No `REC-FAIL`

There is no final `REC-FAIL` in the recognizer set. Earlier versions of the recognizer concept did have such a bottom element. 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 sequence size is needed. Adding a check for an empty recognizer sequence is more code. The second argument against is that adding a recognizer to the recognizer becomes more complex since the bottom element has to be kept, essentially making appending a recognizer always an insert-in-the-middle action.

### RECTYPE-sometype

Every recognizer returns the data and a id, called `RECTYPE-sometype`. This id is used to identify a data type and it provides all information necessary to handle the data inside the interpreter. Each data item is used in three different actions: interpret, compile and postpone. The interpret and compile action are used depending on `STATE`. The postpone action serializes the data and adds the data specific compile action to be executed later.

This design follows the name tokens and `TRAVERSE-WORDLIST` from the Programming Tools wordset.

# recognizable?

A common question was "what if I want to check only a given string whether it's recognizable or not". Esp the result of the parsing was of no interest. The `recognize` word returns both the datatype information and the data itself. A simple solution to get only the datatype information is using exceptions

```
: (recognizable?) ( addr len -- rectype-data )
   [: forth-recognizer recognize throw ;]
   catch nip nip ;
: recognizable? ( addr len -- flag )
   (recognizable?) RECTYPE-NULL <> ;
```

The code assumes that the numeric value of any rectype-data item is never zero.

## LOCATE

`LOCATE` is an interactive tool found in many Forth systems to display information about an item `<something>` that follows immediately in the input line. `LOCATE` is non-standard and may thus has different meanings and implementations. It usually depends on carnal knowledge of the system.

With recognizers the fear came up, that a `LOCATE` may not work any longer due to complex syntactic schemas that are not easy to handle.

### *Existing Practice*

Common usage is `LOCATE word` giving a brief information where the source code of the definition can be found or directly displaying this information.

Only words in wordlists are subject to be `LOCATE`'d. Numbers and other literal-like data are not expected to work and produce various error messages.

The actions taken during `LOCATE` can be customized in many ways, defers, macros and substitutions are used.

Gforth (file *locate.fs*): `LOCATE word` opens a file called *TAGS*, searches there for `word`, constructs a command line from the information found to invoke the vi editor and executes it. If something unexpected happens exceptions are thrown at various stages.

Swiftforth has a header field for `LOCATE` information, VFXlin keeps somehow track of the file names during compilations. Both systems use them to display display the data and/or execute command lines.

### *Possible implementations*

The first approach assumes that the information `LOCATE` uses are tied to the item itself. E.g. a header element in the wordlist entry. The systems that go that way have words that make header information available starting from the execution tokens (XT) or the name token (NT). This information is part of the usual `RECOGNIZE` step.

E.g. `LOCATE FOO` may display "UNKNOWN" assuming `FOO` is not defined anywhere. `LOCATE IF` may display "XT address 1" (for an immediate `IF`) a user supplied recognizer (for simplicity the name token lookup) may display "NT address". A `->` recognizer that implements the `TO` operation can display the "TO address" information from the (hypothetical) `RECTYPE-TO` data token.

System specific knowledge in a `RECTYPE>STRING ( RECTYPE -- addr len)` transforms the RECTYPE-XT into something human readable. This is similar to the `NAME>STRING`. The recognizer sequence that is used to identify the data may be the same as the text interpreter is supposed to use (`FORTH-RECOGNIZER`). That way the `LOCATE` can be implemented as

```
: LOCATE
  DEPTH N>R \ save current data
```

```
   PARSE-NAME FORTH-RECOGNIZER RECOGNIZE DISPLAY-RECTYPE-DATA
   NR> DROP  \ restore previously saved data
;
```

with `DISPLAY-RECTYPE-DATA` to show the data actually is something like

```
: DISPLAY-RECTYPE-DATA RECTYPE>STRING TYPE DEPTH 0 ?DO . LOOP ;
```

This `DISPLAY-RECTYPE-DATA` can be expanded to work with any system provided recognizer data and may have a hook for user supplied ones.

The second version of `LOCATE` is a recognizer itself. This is illustrated for the `TAGS` file based `LOCATE` as in gforth. The recognizer returns a new data type id, called `RECTYPE-TAGS` This data type id does not need support compiling and postponing actions. The `LOCATE` command uses the interpret action only. The parsing action may be located in a recognizer sequence of its own or may be added temporarily to the standard set.

```
:NONAME ( addr len -- ) TYPE ; :NONAME 2DROP ; DUP RECTYPE: RECTYPE-TAGS
: REC-TAGS ( addr len -- addr' len' RECTYPE-TAGS | RECTYPE-NULL )
   \ open TAGS file, search for addr/len and create a new
   \ string with data from the TAGS file at addr' len' if found
;
1 REC-STACK LOCATE-RECOGNIZER
' REC-TAGS 1 LOCATE-RECOGNIZER SET-STACK
: LOCATE PARSE-NAME LOCATE-RECOGNIZER RECOGNIZE RECTYPE>INT EXECUTE ;
```

With the `LOCATE-RECOGNIZER` as a separate set, user supplied data type id's can be added to the LOCATE sequence easily. Moreover any non-locate-able strings (literals) are handled automatically without interfering with other data locations (floating point stack) due to the standard `RECTYPE-NULL` action.

## RECTYPE-NULL **necessity**

Comparing the different implementations. Esp the dual use as a flag and a token is discussed with code examples.

Exceptions are not an option as already discussed.

### *Parse* `REC-*` *actions*

For simplicity the recognizer for floating point numbers.

With `RECTYPE-NULL`

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

Without `RECTYPE-NULL`

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

Conclusion: almost the same.

RECOGNIZE

with `RECTYPE-NULL`

```
: RECOGNIZE    ( addr len rec-set-id -- i*x RECTYPE-sometype | RECTYPE-NULL )
    DUP >R @
    BEGIN
      DUP
    WHILE
      DUP CELLS R@ + @
      2OVER 2>R SWAP 1- >R
      EXECUTE DUP RECTYPE-NULL <> IF
        2R> 2DROP 2R> 2DROP EXIT
      THEN DROP R> 2R> ROT
    REPEAT
    DROP 2DROP R> DROP RECTYPE-NULL
;
```

Without `RECTYPE-NULL`

```
: RECOGNIZE    ( addr len rec-set-id -- i*x RECTYPE-sometype | 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 RECTYPE-NULL
;
```

again, almost the same.

POSTPONE

with `RECTYPE-NULL`

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

without `RECTYPE-NULL`

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

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

### *Interpreter*

With `RECTYPE-NULL`

```
: INTERPRET
  BEGIN
      PARSE-NAME DUP
  WHILE
      FORTH-RECOGNIZER RECOGNIZE
      STATE @ IF RECTYPE>COMP ELSE RECTYPE>INT THEN
      EXECUTE \ do the action.
      ?STACK  \ simple housekeeping
  REPEAT 2DROP
;
```

Without `RECTYPE-NULL`

```
: INTERPRET
  BEGIN
      PARSE-NAME DUP
  WHILE
      FORTH-RECOGNIZER RECOGNIZE
      ?DUP IF \ we got an RECTYPE-*
        STATE @ IF RECTYPE>COMP ELSE RECTYPE>INT THEN
        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 extending the text interpreter specification too.

### *Conclusion*

`RECTYPE-NULL` 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: `RECTYPE-NULL` vs `0`.

# Use Cases

These use cases are purely informative.

## 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 `REC-NT ( addr len -- nt RECTYPE-NT | RECTYPE-NULL)`. 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
RECTYPE: RECTYPE-NT
```

The actual `REC-NT` 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-nt)    ( addr len wid -- nt RECTYPE-NT | RECTYPE-NULL )
  search-name ?DUP IF RECTYPE-NT ELSE RECTYPE-NULL THEN
;

\ checks only the standard wordlist
: REC-NT ( addr len -- nt RECTYPE-NT | RECTYPE-NULL )
  FORTH-WORDLIST (rec-nt)
;
```

# 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 set.

The words dealing with the order stack (`ALSO`, `PREVIOUS`, `FORTH`, `ONLY` etc) may be extended/changed to handle the recognizer sequence too/instead. On the other hand, `ALSO` is essentially `DUP` on a different stack. `ONLY` and `FORTH` set a predefined stack content.

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) RECTYPE>INT EXECUTE ;
: (interpret-c) RECTYPE>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 RECOGNIZE  \ 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
RECTYPE: RECTYPE-notfound

: rec-notfound ( addr len -- )
   RECTYPE-notfound
;
```

With that recognizer put at the end (bottom) of the recognizer set, 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 sequence 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 RECTYPE-XT to work.

```
: executable? ( RECTYPE-TOKEN -- f )
        RECTYPE>INT \ get the interpretation action for the given token
   RECTYPE-XT RECTYPE>INT \ get the system specific interpret action
   =
;

: ' ( "<spaces>name" -- XT )
  PARSE-NAME FORTH-RECOGNIZER RECOGNIZE
  executable? 0= IF
    \ call the system specific error action "invalid tick"
    -13 THROW
  THEN
  DROP \ remove the immediate flag
  \ the XT from the RECTYPE-XT data 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 RECOGNIZE action (REC-* words) and the literal-like and final-action set replaces the proposed 3 method set in RECTYPE-*. 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 RECTYPE: RECTYPE-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 behavior. 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 behavior 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 behavior.

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.