

# R Internals

---

Version 2.6.2 (2008-02-08)

R Development Core Team

---

Permission is granted to make and distribute verbatim copies of this manual provided the copyright notice and this permission notice are preserved on all copies.

Permission is granted to copy and distribute modified versions of this manual under the conditions for verbatim copying, provided that the entire resulting derived work is distributed under the terms of a permission notice identical to this one.

Permission is granted to copy and distribute translations of this manual into another language, under the above conditions for modified versions, except that this permission notice may be stated in a translation approved by the R Development Core Team.

Copyright © 1999–2006 R Development Core Team

ISBN 3-900051-14-3

# Table of Contents

# 1 R Internal Structures

This chapter is the beginnings of documentation about R internal structures. It is written for the R core team and others studying the code in the ‘src/main’ directory.

It is a work-in-progress, first begun for R 2.4.0, and should be checked against the current version of the source code.

## 1.1 SEXPs

What R users think of as *variables* or *objects* are symbols which are bound to a value. The value can be thought of as either a SEXP (a pointer), or the structure it points to, a SEXPREC (and there are alternative forms used for vectors, namely VECSXP pointing to VECTOR\_SEXPREC structures). So the basic building blocks of R objects are often called *nodes*, meaning SEXPRECs or VECTOR\_SEXPRECs.

Note that the internal structure of the SEXPREC is not made available to R Extensions: rather SEXP is an opaque pointer, and the internals can only be accessed by the functions provided.

Both types of node structure have as their first three fields a 32-bit `sxpinfo` header and then three pointers (to the attributes and the previous and next node in a doubly-linked list), and then some further fields. On a 32-bit platform a node<sup>1</sup> occupies 28 bytes: on a 64-bit platform typically 56 bytes (depending on alignment constraints).

The first five bits of the `sxpinfo` header specify one of up to 32 SEXPTYPEs.

### 1.1.1 SEXPTYPEs

Currently SEXPTYPEs 0:10 and 13:25 are in use. Values 11 and 12 were used for internal factors and ordered factors and have since been withdrawn. Note that the SEXPTYPEs are stored in `saved` objects and that the ordering of the types is used, so the gap cannot easily be reused.

no	SEXPTYPE	Description
0	NILSXP	NULL
1	SYMSXP	symbols
2	LISTSXP	pairlists
3	CLOSXP	closures
4	ENVSXP	environments
5	PROMSXP	promises
6	LANGSXP	language objects
7	SPECIALSXP	special functions
8	BUILTINSXP	builtin functions
9	CHARSXP	internal character strings
10	LGLSXP	logical vectors
13	INTSXP	integer vectors
14	REALSXP	numeric vectors
15	CPLXSXP	complex vectors
16	STRSXP	character vectors
17	DOTSXP	dot-dot-dot object
18	ANYSXP	make “any” args work
19	VECSXP	list (generic vector)
20	EXPRSXP	expression vector
21	BCODESXP	byte code

---

<sup>1</sup> strictly, a SEXPREC node; VECTOR\_SEXPREC nodes are slightly smaller but followed by data in the node.

22	EXTPTRSXP	external pointer
23	WEAKREFSXP	weak reference
24	RAWSXP	raw vector
25	S4SXP	S4 classes not of simple type

Many of these will be familiar from R level: the atomic vector types are `LGLSXP`, `INTSXP`, `REALSXP`, `CPLXSP`, `STRSXP` and `RAWSXP`. Lists are `VECSXP` and names (also known as symbols) are `SYMSXP`. Pairlists (`LISTSXP`, the name going back to the origins of R as a Scheme-like language) are rarely seen at R level, but are for example used for argument lists. Character vectors are effectively lists all of whose elements are `CHARSXP`, a type that is rarely visible at R level.

Language objects (`LANGSXP`) are calls (including formulae and so on). Internally they are pairlists with first element a reference<sup>2</sup> to the function to be called with remaining elements the actual arguments for the call (and with the tags if present giving the specified argument names). Although this is not enforced, many places in the code assume that the pairlist is of length one or more, often without checking.

Expressions are of type `EXPRSXP`: they are a vector of (usually language) objects most often seen as the result of `parse()`.

The functions are of types `CLOSXP`, `SPECIALSXP` and `BUILTINSXP`: where `SEXPTYPE`s are stored in an integer these are sometimes lumped into a pseudo-type `FUNSPXP` with code 99. Functions defined via `function` are of type `CLOSXP` and have `formals`, `body` and `environment`.

The `SEXPTYPE` `S4SXP` was introduced in R 2.4.0 for S4 classes which were previously represented as empty lists, that is objects which do not consist solely of a simple type such as an atomic vector or function.

### 1.1.2 Rest of header

The `sxpinfo` header is defined as a 32-bit C structure by

```
struct sxpinfo_struct {
    SEXPTYPE type      : 5; /* discussed above */
    unsigned int obj   : 1; /* is this an object with a class attribute? */
    unsigned int named : 2; /* used to control copying */
    unsigned int gp    : 16; /* general purpose, see below */
    unsigned int mark  : 1; /* mark object as 'in use' in GC */
    unsigned int debug : 1;
    unsigned int trace : 1;
    unsigned int spare : 1; /* unused */
    unsigned int gcgen : 1; /* generation for GC */
    unsigned int gccls : 3; /* class of node for GC */
}; /* Tot: 32 */
```

The `debug` bit is used for closures and environments. For closures it is set by `debug()` and unset by `undebug()`, and indicates that evaluations of the function should be run under the browser. For environments it indicates whether the browsing is in single-step mode.

The `trace` bit is used for functions for `trace()` and for other objects when tracing duplications (see `tracemem`).

The `named` field is set and accessed by the `SET_NAMED` and `NAMED` macros, and take values 0, 1 and 2. R has a ‘call by value’ illusion, so an assignment like

```
b <- a
```

---

<sup>2</sup> a pointer to a function or a symbol to look up the function by name, or a language object to be evaluated to give a function.

appears to make a copy of `a` and refer to it as `b`. However, if neither `a` nor `b` are subsequently altered there is no need to copy. What really happens is that a new symbol `b` is bound to the same value as `a` and the `named` field on the value object is set (in this case to 2). When an object is about to be altered, the `named` field is consulted. A value of 2 means that the object must be duplicated before being changed. (Note that this does not say that it is necessary to duplicate, only that it should be duplicated whether necessary or not.) A value of 0 means that it is known that no other `SEXP` shares data with this object, and so it may safely be altered. A value of 1 is used for situations like

```
dim(a) <- c(7, 2)
```

where in principle two copies of `a` exist for the duration of the computation as (in principle)

```
a <- 'dim<-'(a, c(7, 2))
```

but for no longer, and so some primitive functions can be optimized to avoid a copy in this case.

The `gp` bits are by definition ‘general purpose’. As of version 2.4.0 of R, bit 4 (i.e., the fifth bit) is turned on to mark S4 objects. Bits 0-3 and bits 14-15 have been used previously as described below (from detective work on the sources).

The bits can be accessed and set by the `LEVELS` and `SETLEVELS` macros, which names appear to date back to the internal factor and ordered types and are now used in only a few places in the code. The `gp` field is serialized/unserialized for the `SEXPTYPEs` other than `NILSXP`, `SYMSXP` and `ENVSXP`.

If we label the bits from 0, bits 14 and 15 of `gp` are used for ‘fancy bindings’. Bit 14 is used to lock a binding or an environment, and bit 15 is used to indicate an active binding. (For the definition of an ‘active binding’ see the header comments in file ‘`src/main/envir.c`’.) Bit 15 is used for an environment to indicate if it participates in the global cache.

Almost all other uses seem to be only of bits 0 and 1, although one reserves the first four bits.

The macros `ARGUSED` and `SET_ARGUSED` are used when matching actual and formal function arguments, and take the values 0, 1 and 2.

The macros `MISSING` and `SET_MISSING` are used for pairlists of arguments. Four bits are reserved, but only two are used (and exactly what for is not explained). It seems that bit 0 is used by `matchArgs` to mark missingness on the returned argument list, and bit 1 is used to mark the use of a default value for an argument copied to the evaluation frame of a closure.

Bit 0 is used by macros `DDVAL` and `SET_DDVAL`. This indicates that a `SYMSXP` is one of the symbols `..n` which are implicitly created when `...` is processed, and so indicates that it may need to be looked up in a `DOTSXP`.

Bit 0 is used for `PRSEEN`, a flag to indicate if a promise has already been seen during the evaluation of the promise (and so to avoid recursive loops).

Bit 0 is used for `HASHASH`, on the `PRINTNAME` of the `TAG` of the frame of an environment.

Bits 0 and 1 are used for weak references (to indicate ‘ready to finalize’, ‘finalize on exit’).

Bit 0 is used by the condition handling system (on a `VECSXP`) to indicate a calling handler.

As from R 2.5.0, bits 2 and 3 for a `CHARSXP` are used to note that it is known to be in Latin-1 and UTF-8 respectively. (These are not usually set if it is also known to be in ASCII, since code does not need to know the charset to handle ASCII strings.)

### 1.1.3 The ‘data’

A `SEXP` is a C structure containing the 32-bit header as described above, three pointers (to the attributes, previous and next node) and the node data, a union

```

union {
    struct primsxp_struct primsxp;
    struct symsxp_struct symsxp;
    struct listsxp_struct listsxp;
    struct envsxp_struct envsxp;
    struct closxp_struct closxp;
    struct promsxp_struct promsxp;
} u;

```

All of these alternatives apart from the first (an `int`) are three pointers, so the union occupies three words.

The vector types are `RAWSXP`, `CHARSXP`, `LGLSXP`, `INTSXP`, `REALSXP`, `CPLXSXP`, `STRSXP`, `VECSXP`, `EXPRSXP` and `WEAKREFSXP`. Remember that such types are a `VECTOR_SEXP`, which again consists of the header and the same three pointers, but followed by two integers giving the length and ‘true length’<sup>3</sup> of the vector, and then followed by the data (aligned as required: on most 32-bit systems with a 24-byte `VECTOR_SEXP` node the data can follow immediately after the node). The data are a block of memory of the appropriate length to store ‘true length’ elements (rounded up to a multiple of 8 bytes, with the 8-byte blocks being the ‘Vcells’ referred in the documentation for `gc()`).

The ‘data’ for the various types are given in the table below. A lot of this is interpretation, i.e. the types are not checked.

#### NILSXP

There is only one object of type `NILSXP`, `R_NilValue`, with no data.

**SYMSXP** Pointers to three nodes, the name, value and internal, accessed by `PRINTNAME` (a `CHARSXP`), `SYMVALUE` and `INTERNAL`. (If the symbol’s value is a `.Internal` function, the last is a pointer to the appropriate `SEXP`.) Many symbols have `SYMVALUE` `R_UnboundValue`.

**LISTSXP** Pointers to the `CAR`, `CDR` (usually a `LISTSXP` or `NULL`) and `TAG` (usually a `SYMSXP`).

**CLOSXP** Pointers to the formals (a pairlist), the body and the environment.

**ENVSXP** Pointers to the frame, enclosing environment and hash table (`NULL` or a `VECSXP`). A frame is a tagged pairlist with tag the symbol and `CAR` the bound value.

**PROMSXP** Pointers to the value, expression and environment (in which to evaluate the expression). Once an promise has been evaluated, the environment is set to `NULL`.

**LANGSXP** A special type of `LISTSXP` used for function calls. (The `CAR` references the function (perhaps via a symbol or language object), and the `CDR` the argument list with tags for named arguments.) R-level documentation references to ‘expressions’ / ‘language objects’ are mainly `LANGSXPs`, but can be symbols (`SYMSXP`s) or expression vectors (`EXPRSXP`s).

#### SPECIALSXP

##### BUILTINSXP

An integer giving the offset into the table of primitives/ `.Internals`.

**CHARSXP** `length`, `trueLength` followed by a block of bytes (allowing for the `nul` terminator).

##### LGLSXP

**INTSXP** `length`, `trueLength` followed by a block of `C` `ints` (which are 32 bits on all R platforms).

<sup>3</sup> This is almost unused. The only current use is for hash tables of environments (`VECSXP`s), where `length` is the size of the table and `trueLength` is the number of primary slots in use, and for the reference hash tables in serialization (`VECSXP`s), where `trueLength` is the number of slots in use.

REALSXP	length, <code>truelength</code> followed by a block of C doubles
CPLXSXP	length, <code>truelength</code> followed by a block of C99 double complexes, or equivalent structures.
STRSXP	length, <code>truelength</code> followed by a block of pointers (SEXPs pointing to CHARSEXPs).
DOTSXP	A special type of LISTSXP for the value bound to a <code>...</code> symbol: a pairlist of promises.
ANYSXP	This is used as a place holder for any type: there are no actual objects of this type.
VECSXP	
EXPRSXP	length, <code>truelength</code> followed by a block of pointers. These are internally identical (and identical to STRSXP) but differ in the interpretations placed on the elements.
BCODESXP	For the future byte-code compiler.
EXTPTRSXP	
	Has three pointers, to the pointer, the protection value (an R object which if alive protects this object) and a tag (a SYMSXP?).
WEAKREFSXP	
	A WEAKREFSXP is a special VECSXP of length 4, with elements ‘key’, ‘value’, ‘finalizer’ and ‘next’. The ‘key’ is NULL, an environment or an external pointer, and the ‘finalizer’ is a function or NULL.
RAWSXP	length, <code>truelength</code> followed by a block of bytes.
S4SXP	two unused pointers and a tag.

### 1.1.4 Allocation classes

As we have seen, the field `gccls` in the header is three bits to label up to 8 classes of nodes. Non-vector nodes are of class 0, and ‘small’ vector nodes are of classes 1 to 6, with ‘large’ vector nodes being of class 7. The ‘small’ vector nodes are able to store vector data of up to 8, 16, 32, 48, 64 and 128 bytes: larger vectors are `malloc`-ed individually whereas the ‘small’ nodes are allocated from pages of about 2000 bytes.

## 1.2 Environments and variable lookup

What users think of as ‘variables’ are symbols which are bound to objects in ‘environments’. The word ‘environment’ is used ambiguously in R to mean *either* the frame of an `ENVSXP` (a pairlist of symbol-value pairs) *or* an `ENVSXP`, a frame plus an enclosure.

There are additional places that ‘variables’ can be looked up, called ‘user databases’ in comments in the code. These seem undocumented in the R sources, but apparently refer to the `RObjectTable` package at <http://www.omegahat.org/RObjectTables/>.

The base environment is special. There is an `ENVSXP` environment with enclosure the empty environment `R_EmptyEnv`, but the frame of that environment is not used. Rather its bindings are part of the global symbol table, being those symbols in the global symbol table whose values are not `R_UnboundValue`. When R is started the internal functions are installed (by C code) in the symbol table, with primitive functions having values and `.Internal` functions having what would be their values in the field accessed by the `INTERNAL` macro. Then `.Platform` and `.Machine` are computed and the base package is loaded into the base environment followed by the system profile.

The frames of environments (and the symbol table) are normally hashed for faster access (including insertion and deletion).

By default R maintains a (hashed) global cache of ‘variables’ (that is symbols and their bindings) which have been found, and this refers only to environments which have been marked to participate, which consists of the global environment (aka the user workspace), the base environment plus environments<sup>4</sup> which have been **attached**. When an environment is either **attached** or **detached**, the names of its symbols are flushed from the cache. The cache is used whenever searching for variables from the global environment (possibly as part of a recursive search).

### 1.2.1 Search paths

S has the notion of a ‘search path’: the lookup for a ‘variable’ leads (possibly through a series of frames) to the ‘session frame’ the ‘working directory’ and then along the search path. The search path is a series of databases (as returned by `search()`) which contain the system functions (but not necessarily at the end of the path, as by default the equivalent of packages are added at the end).

R has a variant on the S model. There is a search path (also returned by `search()`) which consists of the global environment (aka user workspace) followed by environments which have been attached and finally the base environment. Note that unlike S it is not possible to attach environments before the workspace nor after the base environment.

However, the notion of variable lookup is more general in R, hence the plural in the title of this subsection. Since environments have enclosures, from any environment there is a search path found by looking in the frame, then the frame of its enclosure and so on. Since loops are not allowed, this process will eventually terminate: until R 2.2.0 it always terminated at the base environment, but nowadays it can terminate at either the base environment or the empty environment. (It can be conceptually simpler to think of the search always terminating at the empty environment, but with an optimization to stop at the base environment.) So the ‘search path’ describes the chain of environments which is taken once the search reaches the global environment.

### 1.2.2 Name spaces

Name spaces are environments associated with packages (and once again the base package is special and will be considered separately). A package `pkg` with a name space defines two environments `namespace:pkg` and `package:pkg`: it is `package:pkg` that can be attached and form part of the search path.

The objects defined by the R code in the package are symbols with bindings in the `namespace:pkg` environment. The `package:pkg` environment is populated by selected symbols from the `namespace:pkg` environment (the exports). The enclosure of this environment is an environment populated with the explicit imports from other name spaces, and the enclosure of *that* environment is the base name space. (So the illusion of the imports being in the name space environment is created via the environment tree.) The enclosure of the base name space is the global environment, so the search from a package name space goes via the (explicit and implicit) imports to the standard ‘search path’.

The base name space environment `R_BaseNamespace` is another `ENVSNXP` that is special-cased. It is effectively the same thing as the base environment `R_BaseEnv` *except* that its enclosure is the global environment rather than the empty environment: the internal code diverts lookups in its frame to the global symbol table.

---

<sup>4</sup> Remember that attaching a list or a saved image actually creates and populates an environment and attaches that.

## 1.3 Attributes

As we have seen, every `SEXP` has a pointer to the attributes of the node (default `R_NilValue`). The attributes can be accessed/set by the macros/functions `ATTRIB` and `SET_ATTRIB`, but such direct access is normally<sup>5</sup> only used to check if the attributes are `NULL` or to reset them. Otherwise access goes through the functions `getAttrib` and `setAttrib` which impose restrictions on the attributes. One thing to watch is that if you copy attributes from one object to another you may (un)set the `"class"` attribute and so need to copy the object and S4 bits as well. There is a macro/function `DUPLICATE_ATTRIB` to automate this.

The code assumes that the attributes of a node are either `R_NilValue` or a pairlist of non-zero length (and this is checked by `SET_ATTRIB`). The attributes are named (via tags on the pairlist). The replacement function `attributes<-` ensures that `"dim"` precedes `"dimnames"` in the pairlist. Attribute `"dim"` is one of several that is treated specially: the values are checked, and any `"names"` and `"dimnames"` attributes are removed. Similarly, you cannot set `"dimnames"` without having set `"dim"`, and the value assigned must be a list of the correct length and with elements of the correct lengths (and all zero-length elements are replaced by `NULL`).

The other attributes which are given special treatment are `"names"`, `"class"`, `"tsp"`, `"comment"` and `"row.names"`. For pairlist-like objects the names are not stored as an attribute but (as symbols) as the tags: however the R interface makes them look like conventional attributes, and for one-dimensional arrays they are stored as the first element of the `"dimnames"` attribute. The C code ensures that the `"tsp"` attribute is an `REALSXP`, the frequency is positive and the implied length agrees with the number of rows of the object being assigned to. Classes and comments are restricted to character vectors, and assigning a zero-length comment or class removes the attribute. Setting or removing a `"class"` attribute sets the object bit appropriately. Integer row names are converted to and from the internal compact representation.

Care needs to be taken when adding attributes to objects of the types with non-standard copying semantics. There is only one object of type `NILSXP`, `R_NilValue`, and that should never have attributes (and this is enforced in `installAttrib`). For environments, external pointers and weak references, the attributes should be relevant to all uses of the object: it is for example reasonable to have a name for an environment, and also a `"path"` attribute for those environments populated from R code in a package.

When should attributes be preserved under operations on an object? Becker, Chambers & Wilks (1988, pp. 144–6) give some guidance. Scalar functions (those which operate element-by-element on a vector and whose output is similar to the input) should preserve attributes (except perhaps class, and if they do preserve class they need to preserve the `OBJECT` and S4 bits). Binary operations normally call `copyMostAttributes` to copy most attributes from the longer argument (and if they are of the same length from both, preferring the values on the first). Here ‘most’ means all except the `names`, `dim` and `dimnames` which are set appropriately by the code for the operator.

Subsetting (other than by an empty index) generally drops all attributes except `names`, `dim` and `dimnames` which are reset as appropriate. On the other hand, subassignment generally preserves such attributes even if the length is changed. Coercion drops all attributes. For example:

```
> x <- structure(1:8, names=letters[1:8], comm="a comment")
> x[]
 a b c d e f g h
 1 2 3 4 5 6 7 8
attr(,"comm")
```

---

<sup>5</sup> An exception is the internal code for `terms.formula` which directly manipulates the attributes.

```

[1] "a comment"
> x[1:3]
a b c
1 2 3
> x[3] <- 3
> x
a b c d e f g h
1 2 3 4 5 6 7 8
attr(,"comm")
[1] "a comment"
> x[9] <- 9
> x
a b c d e f g h
1 2 3 4 5 6 7 8 9
attr(,"comm")
[1] "a comment"

```

## 1.4 Contexts

*Contexts* are the internal mechanism used to keep track of where a computation has got to (and from where), so that control-flow constructs can work and reasonable information can be produced on error conditions, (such as *via* traceback) and otherwise (the `sys.xxx` functions).

Execution contexts are a stack of C structs:

```

typedef struct RCNTXT {
    struct RCNTXT *nextcontext; /* The next context up the chain */
    int callflag;              /* The context 'type' */
    JMP_BUF cjmpbuf;          /* C stack and register information */
    int cstacktop;            /* Top of the pointer protection stack */
    int evaldepth;           /* Evaluation depth at inception */
    SEXP promargs;           /* Promises supplied to closure */
    SEXP callfun;            /* The closure called */
    SEXP sysparent;          /* Environment the closure was called from */
    SEXP call;               /* The call that effected this context */
    SEXP cloenv;             /* The environment */
    SEXP conexit;           /* Interpreted on.exit code */
    void (*cend)(void *);    /* C on.exit thunk */
    void *cenddata;         /* Data for C on.exit thunk */
    char *vmax;             /* Top of the R_alloc stack */
    int intsusp;            /* Interrupts are suspended */
    SEXP handlerstack;      /* Condition handler stack */
    SEXP restartstack;      /* Stack of available restarts */
    struct RPRSTACK *prstack; /* Stack of pending promises */
} RCNTXT, *context;

```

plus additional fields for the future byte-code compiler. The 'types' are from

```

enum {
    CTXT_TOPLEVEL = 0, /* toplevel context */
    CTXT_NEXT      = 1, /* target for next */
    CTXT_BREAK     = 2, /* target for break */
    CTXT_LOOP      = 3, /* break or next target */
    CTXT_FUNCTION  = 4, /* function closure */

```

```

CTXT_CCODE      = 8, /* other functions that need error cleanup */
CTXT_RETURN     = 12, /* return() from a closure */
CTXT_BROWSER    = 16, /* return target on exit from browser */
CTXT_GENERIC    = 20, /* rather, running an S3 method */
CTXT_RESTART    = 32, /* a call to restart was made from a closure */
CTXT_BUILTIN    = 64 /* builtin internal function */
};

```

where the `CTXT_FUNCTION` bit is on wherever function closures are involved.

Contexts are created by a call to `begincontext` and ended by a call to `endcontext`: code can search up the stack for a particular type of context via `findcontext` (and jump there) or jump to a specific context via `R_JumpToContext`. `R_ToplevelContext` is the ‘idle’ state (normally the command prompt), and `R_GlobalContext` is the top of the stack.

Note that whilst all calls to closures set a context, those to special internal functions never do, and those to builtin internal functions have done so only recently (and prior to that only when profiling).

Dispatching from a S3 generic (via `UseMethod` or its internal equivalent) or calling `NextMethod` sets the context type to `CTXT_GENERIC`. This is used to set the `sysparent` of the method call to that of the `generic`, so the method appears to have been called in place of the generic rather than from the generic.

The R `sys.frame` and `sys.call` work by counting calls to closures (type `CTXT_FUNCTION`) from either end of the context stack.

Note that the `sysparent` element of the structure is not the same thing as `sys.parent()`. Element `sysparent` is primarily used in managing changes of the function being evaluated, i.e. by `Recall` and method dispatch.

`CTXT_CCODE` contexts are currently used in `cat()`, `load()`, `scan()` and `write.table()` (to close the connection on error), by `PROTECT`, serialization (to recover from errors, e.g. free buffers) and within the error handling code (to raise the C stack limit and reset some variables).

## 1.5 Argument evaluation

As we have seen, functions in R come in three types, closures (`SEXPTYPE CLOSXP`), specials (`SPECIALSXP`) and builtins (`BUILTINSXP`). In this section we consider when (and if) the actual arguments of function calls are evaluated. The rules are different for the internal (special/builtin) and R-level functions (closures).

For a call to a closure, the actual and formal arguments are matched and a matched call (another `LANGSXP`) is constructed. This process first replaces the actual argument list by a list of promises to the values supplied. It then constructs a new environment which contains the names of the formal parameters matched to actual or default values: all the matched values are promises, the defaults as promises to be evaluated in the environment just created. That environment is then used for the evaluation of the body of the function, and promises will be forced (and hence actual or default arguments evaluated) when they are encountered. (Evaluating a promise sets `NAMED = 2` on its value, so if the argument was a symbol its binding is regarded as having multiple references during the evaluation of the closure call.)

If the closure is an S3 generic (that is, contains a call to `UseMethod`) the evaluation process is the same until the `UseMethod` call is encountered. At that point the argument on which to do dispatch (normally the first) will be evaluated if it has not been already. If a method has been found which is a closure, a new evaluation environment is created for it containing the matched arguments of the method plus any new variables defined so far during the evaluation of the body of the generic. (Note that this means changes to the values of the formal arguments in the

body of the generic are discarded when calling the method, but *actual* argument promises which have been forced retain the values found when they were forced. On the other hand, missing arguments have values which are promises to use the default supplied by the method and not the generic.) If the method found is a special or builtin it is called with the matched argument list of promises (possibly already forced) used for the generic.

The essential difference<sup>6</sup> between special and builtin functions is that the arguments of specials are not evaluated before the C code is called, and those of builtins are. In each case positional matching of arguments is used. Note that being a special/builtin is separate from being primitive or `.Internal: function` is a special primitive, `+` is a builtin primitive, `switch` is a special `.Internal` and `grep` is a builtin `.Internal`.

Many of the internal functions are internal generics, which for specials means that they do not evaluate their arguments on call, but the C code starts with a call to `DispatchOrEval`. The latter evaluates the first argument, and looks for a method based on its class. (If S4 dispatch is on, S4 methods are looked for first, even for S3 classes.) If it finds a method, it dispatches to that method with a call based on promises to evaluate the remaining arguments. If no method is found, the remaining arguments are evaluated before return to the internal generic.

The other way that internal functions can be generic is to be group generic. All such functions are builtins (so immediately evaluate all their arguments), and contain a call to the C function `DispatchGeneric`. There are some peculiarities over the number of arguments for the "Math" group generic, with some members allowing only one argument, some having two (with a default for the second) and `trunc` allows one or more but the default only accepts one.

### 1.5.1 Missingness

Actual arguments to (non-internal) R functions can be fewer than are required to match the formal arguments of the function. Having unmatched formal arguments will not matter if the argument is never used (by lazy evaluation), but when the argument is evaluated, either its default value is evaluated (within the evaluation environment of the function) or an error is thrown with a message along the lines of

```
argument "foobar" is missing, with no default
```

Internally missingness is handled by two mechanisms. The object `R_MissingArg` is used to indicate that a formal argument has no (default) value. When matching the actual arguments to the formal arguments, a new argument list is constructed from the formals all of whose values are `R_MissingArg` with the first `MISSING` bit set. Then whenever a formal argument is matched to an actual argument, the corresponding member of the new argument list has its value set to that of the matched actual argument, and if that is not `R_MissingArg` the missing bit is unset.

This new argument list is used to form the evaluation frame for the function, and if named arguments are subsequently given a new value (before they are evaluated) the missing bit is cleared.

Missingness of arguments can be interrogated via the `missing()` function. An argument is clearly missing if its missing bit is set or if the value is `R_MissingArg`. However, missingness can be passed on from function to function, for using a formal argument as an actual argument in a function call does not count as evaluation. So `missing()` has to examine the value (a promise) of a non-yet-evaluated formal argument to see if it might be missing, which might involve investigating a promise and so on . . . .

---

<sup>6</sup> There is currently one other difference: when profiling builtin functions are counted as function calls but specials are not.

## 1.5.2 Dot-dot-dot arguments

Dot-dot-dot arguments are convenient when writing functions, but complicate the internal code for argument evaluation.

The formals of a function with a `...` argument represent that as a single argument like any other argument, with tag the symbol `R_DotsSymbol`. When the actual arguments are matched to the formals, the value of the `...` argument is of `SEXPTYPE DOTSXP`, a pairlist of promises (as used for matched arguments) but distinguished by the `SEXPTYPE`.

Recall that the evaluation frame for a function initially contains the `name=value` pairs from the matched call, and hence this will be true for `...` as well. The value of `...` is a (special) pairlist whose elements are referred to by the special symbols `..1`, `..2`, `...` which have the `DDVAL` bit set: when one of these is encountered it is looked up (via `ddfndVar`) in the value of the `...` symbol in the evaluation frame.

Values of arguments matched to a `...` argument can be missing.

## 1.6 Autoprinting

Whether the returned value of a top-level R expression is printed is controlled by the global boolean variable `R_Visible`. This is set (to true or false) on entry to all primitive and internal functions based on the `eval` column of the table in `'names.c'`: the appropriate setting can be extracted by the macro `PRIMPRINT`.

The R primitive function `invisible` makes use of this mechanism: it just sets `R_Visible = FALSE` before entry and returns its argument.

For most functions the intention will be that the setting of `R_Visible` when they are entered is the setting used when they return, but there need to be exceptions. The R functions `identify`, `options`, `system` and `writeBin` determine whether the result should be visible from the arguments or user action. Other functions themselves dispatch functions which may change the visibility flag: examples<sup>7</sup> are `.Internal`, `do.call`, `eval`, `eval.with.vis`<sup>8</sup>, `if`, `NextMethod`, `Recall`, `recordGraphics`, `standardGeneric`, `switch` and `UseMethod`.

'Special' primitive and internal functions evaluate their arguments internally *after* `R_Visible` has been set, and evaluation of the arguments (e.g. an assignment as in PR#9263) can change the value of the flag. Prior to R 2.5.0, known instances of such functions reset the flag after the internal evaluation of arguments: examples include `[`, `[[`, `$`, `c`, `cbind`, `dump`, `rbind` and `unlist`, as well as the language constructs (which are primitives) `for`, `while` and `repeat`.

The `R_Visible` flag can also get altered during the evaluation of a function, with comments in the code about `warning`, `writeChar` and graphics functions calling `GText` (PR#7397). (Since the C-level function `eval` sets `R_Visible`, this could apply to any function calling it. Since it is called when evaluating promises, even object lookup can change `R_Visible`.) From R 2.1.0 internal functions that were marked to set `R_Visible = FALSE` enforced this when the function returned. As from R 2.5.0 both internal and primitive functions force the documented setting of `R_Visible` on return, unless the C code is allowed to change it (the exceptions above are indicated by `PRIMPRINT` having value 2).

The actual autoprinting is done by `PrintValueEnv` in `'print.c'`. If the object to be printed has the `S4` bit set and `S4` methods dispatch is on, `show` is called to print the object. Otherwise, if the object bit is set (so the object has a `"class"` attribute), `print` is called to dispatch methods: for objects without a class the internal code of `print.default` is called.

<sup>7</sup> the other current example is left brace, which is implemented as a primitive.

<sup>8</sup> a `.Internal`-only function used in `source`, `withVisible` and a few other places.

## 1.7 The write barrier and the garbage collector

R has since version 1.2.0 had a generational garbage collector, and bit `gcgen` in the `sxpinfo` header is used in the implementation of this. This is used in conjunction with the `mark` bit to identify two previous generations.

There are three levels of collections. Level 0 collects only the youngest generation, level 1 collects the two youngest generations and level 2 collects all generations. After 20 level-0 collections the next collection is at level 1, and after 5 level-1 collections at level 2. Further, if a level- $n$  collection fails to provide 20% free space (for each of nodes and the vector heap), the next collection will be at level  $n+1$ . (The R-level function `gc()` performs a level-2 collection.)

A generational collector needs to efficiently ‘age’ the objects, especially list-like objects (including `STRSXPs`). This is done by ensuring that the elements of a list are regarded as at least as old as the list *when they are assigned*. This is handled by the functions `SET_VECTOR_ELT` and `SET_STRING_ELT`, which is why they are functions and not macros. Ensuring the integrity of such operations is termed the *write barrier* and is done by making the `SEXP` opaque and only providing access via functions (which cannot be used as lvalues in assignments in C).

All code in R extensions is by default behind the write barrier. The only way to obtain direct access to the internals of the `SEXPRECs` is to define ‘`USE_RINTERNALS`’ before including ‘`Rinternals.h`’, which is normally defined in ‘`Defn.h`’. To enable a check on the way that the access is used, R can be compiled with flag ‘`--enable-strict-barrier`’ which ensures that ‘`Defn.h`’ does not define ‘`USE_RINTERNALS`’ and hence that `SEXP` is opaque in most of R itself. (There are some necessary exceptions: foremost ‘`memory.c`’ where the accessor functions are defined and also ‘`size.c`’ which needs access to the sizes of the internal structures.)

For background papers see <http://www.stat.uiowa.edu/~luke/R/barrier.html> and <http://www.stat.uiowa.edu/~luke/R/gengcnotes.html>.

## 1.8 Serialization Formats

Serialized versions of R objects are used by `load/save` and also at a lower level by `.saveRDS/.readRDS` and `serialize/unserialize`. These differ in what they serialize to (a file, a connection, a raw vector) and whether they are intended to serialize a single object or a collection of objects (typically a workspace). `save` writes a header indicating the format at the beginning of the file (a single LF-terminated line) which the lower-level versions do not.

R has used the same serialization format since R 1.4.0 in December 2001. Reading of earlier formats is still supported via `load`, but they are not described here. (Files of most of these formats can still be found in ‘`data`’ directories of packages.) The current serialization format is called ‘version 2’, and has been expanded in back-compatible ways since R 1.4.0, for example to support additional `SEXPTYPEs`.

`save()` works by first creating a tagged pairlist of objects to be saved, and then saving that single object preceded by a single-line header (typically `RDX2\n` for a binary save). `load()` reads the header line, unserializes a single object (a pairlist or a vector list) and assigns the elements of the list in the appropriate environment.

Serialization in R needs to take into account that objects may contain references to environments, which then have enclosing environments and so on. (Environments recognized as package or name space environments are saved by name.) Further, there are ‘reference objects’ which are not duplicated on copy and should remain shared on unserialization. These are weak references, external pointers and environments other than those associated with packages, name spaces and the global environment. These are handled via a hash table, and references after the first are written out as a reference marker indexed by the table entry.

Serialization first writes a header indicating the format (normally ‘X\n’ for an XDR format binary save, but ‘A\n’, ASCII, and ‘B\n’, native word-order binary<sup>9</sup>, can also occur) and the version number of the format and of two R versions (as integers). (Unserialization interprets the two versions as the version of R which wrote the file followed by the minimal version of R needed to read the format.) Serialization then writes out the object recursively using function `WriteItem` in file ‘src/main/serialize.c’.

Some objects are written as if they were `SEXPTYPE`s: such pseudo-`SEXPTYPE`s cover `R_NilValue`, `R_EmptyEnv`, `R_BaseEnv`, `R_GlobalEnv`, `R_UnboundValue`, `R_MissingArg` and `R_BaseNamespace`.

For all `SEXPTYPE`s except `NILSXP`, `SYMSXP` and `ENVSXP` serialization starts with an integer with the `SEXPTYPE` in bits 0:7<sup>10</sup> followed by the object bit, two bits indicating if there are any attributes and if there is a tag (for the pairlist types), an unused bit and then the `gp` field<sup>11</sup> in bits 12:27. Pairlist-like objects write their attributes (if any), tag (if any), `CAR` and then `CDR` (using tail recursion): other objects write their attributes after themselves. Atomic vector objects write their length followed by the data: generic vector-list objects write the length followed by a call to `WriteItem` for each element. The code for `CHARSXP`s special-cases `NA_STRING` and writes it as length `-1` with no data.

Environments are treated in several ways: as we have seen, some are written as specific pseudo-`SEXPTYPE`s. Package and name space environments are written with pseudo-`SEXPTYPE`s followed by the name. ‘Normal’ environments are written out as `ENVSXP`s with an integer indicating if the environment is locked followed by the enclosure, frame, ‘tag’ (the hash table) and attributes.

In the ‘XDR’ format integers and doubles are written in bigendian order: however the format is not fully XDR as defined in RFC 1832 as byte quantities (such as the contents of `CHARSXP` and `RAWSXP` types) are written as-is and not padded to a multiple of four bytes.

The ‘ASCII’ format writes 7-bit characters. Integers are formatted with `%d` (except that `NA_integer_` is written as `NA`), doubles formatted with `%.16g` (plus `NA`, `Inf` and `-Inf`) and bytes with `%02x`. Strings are written using standard escapes (e.g. `\t` and `\013` for non-printing and non-ASCII bytes).

## 1.9 Encodings for `CHARSXP`s

Character data in R are stored in the sexptype `CHARSXP`. Until R 2.1.0 it was assumed that the data were in the platform’s native 8-bit encoding, and furthermore it was quite often assumed that the encoding was ISO Latin-1 or a superset (such as Windows’ CP1252 or Latin-9).

As from R 2.1.0 there was support for other encodings, in particular UTF-8 and the multi-byte encodings used on Windows for CJK languages. However, there was no way of indicating which encoding had been used, even if this was known (and e.g. `scan` would not know the encoding of the file it was reading). This led to packages with data in French encoded in Latin-1 in `.rda` files which could not be read in other locales (and they would be able to be displayed in a French UTF-8 locale, if not in most Japanese locales).

R 2.5.0 introduced a limited means to indicate the encoding of a `CHARSXP` via two of the ‘general purpose’ bits which are used to declare the encoding to be either Latin-1 or UTF-8. (Note that it is possible for a character vector to contain elements in different encodings.) Both printing and plotting notice the declaration and convert the string to the current locale (possibly using `<xx>` to display in hexadecimal bytes that are not valid in the current locale). Many (but

<sup>9</sup> there is no R-level interface to this format

<sup>10</sup> only 0:4 will currently be used for `SEXPTYPE`s but values 241:255 are used for pseudo-`SEXPTYPE`s.

<sup>11</sup> Currently the only relevant bits are 0:1, 4, 14:15.

not all) of the character manipulation functions will either preserve the declaration or re-encode the character string.

Eventually strings that refer to the OS such as file names will need to be passed through a wide-character interface on some OSes (e.g. Windows), but currently they are just recoded to the current locale.

When are character strings declared to be of known encoding? One way is to do so directly via **Encoding**. The parser declares the encoding if this is known, either via the **encoding** argument to **parse** or from the locale within which parsing is being done at the R command line. Functions **scan**, **read.table**, **readLines** and **source** have an **encoding** argument, but do not assume anything about files from the current locale. Also, **iconv** marks character strings it converts to Latin-1 or UTF-8.

It is not necessary to declare the encoding of ASCII strings as they will work in any locale, but the overhead in doing so is small since they will never be passed to **iconv** for translation.

The rationale behind considering only UTF-8 and Latin-1 is that most systems are capable of producing UTF-8 strings and this is the nearest we have to a universal format. For those that do not (for example those lacking a powerful enough **iconv**), it is likely that they work in Latin-1, the old R assumption.

## 1.10 Warnings and errors

Each of **warning** and **stop** have two C-level equivalents, **warningcall**, **error** and **errorcall**. The relationship between the pairs is similar: **warning** tries to fathom out a suitable call, and then calls **warningcall** with that call as the first argument if it succeeds, and with **call = R\_NilValue** it does not. When **warningcall** is called, it includes the deparsed call in its printout unless **call = R\_NilValue**.

**warning** and **error** look at the context stack. If the topmost context is not of type **CTXT\_BUILTIN**, it is used to provide the call, otherwise the next context provides the call. This means that when these function are called from a primitive or **.Internal**, the imputed call will not be to primitive/**.Internal** but to the function calling the primitive/**.Internal**. This is exactly what one wants for a **.Internal**, as this will give the call to the closure wrapper. (Further, for a **.Internal**, the call is the argument to **.Internal**, and so may not correspond to any R function.) However, it is unlikely to be what is needed for a primitive.

The upshot is that that **warningcall** and **errorcall** should normally be used for code called from a primitive, and **warning** and **error** should be used for code called from a **.Internal** (and necessarily from **.Call**, **.C** and so on, where the call is not passed down). However, there are two complications. One is that code might be called from either a primitive or a **.Internal**, in which case probably **warningcall** is more appropriate. The other involves replacement functions, where the call will be of the form (from R < 2.6.0)

```
> length(x) <- y ~ x
Error in "length<-"('*tmp*', value = y ~ x) : invalid value
```

which is unpalatable to the end user. For replacement functions there will be a suitable context at the top of the stack, so **warning** should be used. (The results for **.Internal** replacement functions such as **substr<-** are not ideal.)

## 1.11 S4 objects

[This section is currently a preliminary draft and should not be taken as definitive. The description assumes that **R\_NO\_METHODS\_TABLES** has not been set.]

### 1.11.1 Representation of S4 objects

[The internal representation of objects from S4 classes changed in R 2.4.0. It is possible that objects from earlier representations still exist, but there is no guarantee that they will be handled correctly. An attempt is made to detect old-style S4 objects and warn when binary objects are loaded or a workspace is restored.]

S4 objects can be of any `SEXPTYPE`. They are either an object of a simple type (such as an atomic vector or function) with S4 class information or of type `S4SXP`. In all cases, the ‘S4 bit’ (bit 4 of the ‘general purpose’ field) is set, and can be tested by the macro/function `IS_S4_OBJECT`.

S4 objects are created via `new()`<sup>12</sup> and thence via the C function `R_do_new_object`. This duplicates the prototype of the class, adds a class attribute and sets the S4 bit. All S4 class attributes should be character vectors of length one with an attribute giving (as a character string) the name of the package (or `.GlobalEnv`) containing the class definition. Since S4 objects have a class attribute, the `OBJECT` bit is set.

It is currently unclear what should happen if the class attribute is removed from an S4 object, or if this should be allowed.

### 1.11.2 S4 classes

S4 classes are stored as R objects in the environment in which they are created, with names `._C__classname`: as such they are not listed by default by `ls`.

The objects are S4 objects of class `"classRepresentation"` which is defined in the `methods` package.

Since these are just objects, they are subject to the normal scoping rules and can be imported and exported from name spaces like other objects. The directives `importClassesFrom` and `exportClasses` are merely convenient ways to refer to class objects without needing to know their internal ‘metaname’ (although `exportClasses` does a little sanity checking via `isClass`).

### 1.11.3 S4 methods

Details of methods are stored in S4 objects of class `"MethodsList"`. They have a non-syntactic name of the form `._M__generic:package` for all methods defined in the current environment for the named generic derived from a specific package (which might be `.GlobalEnv`).

There is also environment `._T__generic:package` which has names the signatures of the methods defined, and values the corresponding method functions. This is often referred to as a ‘methods table’.

When a package without a name space is attached these objects become visible on the search path. `library` calls `methods:::cacheMetaData` to update the internal tables.

During an R session there is an environment associated with each non-primitive generic containing objects `.AllMTable`, `.Generic`, `.Methods`, `.MTable`, `.SigArgs` and `.SigLength`. `.MTable` and `AllMTable` are merged methods tables containing all the methods defined directly and via inheritance respectively. `.Methods` is a merged methods list.

Exporting methods from a name space is more complicated than exporting a class. Note first that you do not export a method, but rather the directive `exportMethods` will export all the methods defined in the name space for a specified generic: the code also adds to the list of generics any that are exported directly. For generics which are listed via `exportMethods` or

<sup>12</sup> This can also create non-S4 objects, as in `new("integer")`.

exported themselves, the corresponding "MethodsList" and environment are exported and so will appear (as hidden objects) in the package environment.

Methods for primitives which are internally S4 generic (see below) are always exported, whether mentioned in the 'NAMESPACE' file or not.

Methods can be imported either via the directive `importMethodsFrom` or via importing a namespace by `import`. Also, if a generic is imported via `importFrom`, its methods are also imported. In all cases the generic will be imported if it is in the namespace, so `importMethodsFrom` is most appropriate for methods defined on generics in other packages. Since methods for a generic could be imported from several different packages, the methods tables are merged.

When a package with a name space is attached `methods:::cacheMetaData` is called to update the internal tables: only the visible methods will be cached.

### 1.11.4 Mechanics of S4 dispatch

This subsection does not discuss how S4 methods are chosen: see <http://developer.r-project.org/howMethodsWork.pdf>.

For all but primitive functions, setting a method on an existing function that is not itself S4 generic creates a new object in the current environment which is a call to `standardGeneric` with the old definition as the default method. Such S4 generics can also be created *via* a call to `setGeneric`<sup>13</sup> and are standard closures in the R language, with environment the environment within which they are created. With the advent of name spaces this is somewhat problematic: if `myfn` was previously in a package with a name space there will be two functions called `myfn` on the search paths, and which will be called depends on which search path is in use. This is starkest for functions in the base name space, where the original will be found ahead of the newly created function from any other package with a name space.

Primitive functions are treated quite differently, for efficiency reasons: this results in different semantics. `setGeneric` is disallowed for primitive functions. The `methods` namespace contains a list `.BasicFunsList` named by primitive functions: the entries are either `FALSE` or a standard S4 generic showing the effective definition. When `setMethod` (or `setReplaceMethod`) is called, it either fails (if the list entry is `FALSE`) or a method is set on the effective generic given in the list.

Actual dispatch of S4 methods for almost all primitives piggy-backs on the S3 dispatch mechanism, so S4 methods can only be dispatched for primitives which are internally S3 generic. When a primitive that is internally S3 generic is called with a first argument which is an S4 object and S4 dispatch is on (that is, the `methods` name space is loaded), `DispatchOrEval` calls `R_possible_dispatch` (defined in 'src/main/objects.c'). (Members of the S3 group generics, which includes all the generic operators, are treated slightly differently: the first two arguments are checked and `DispatchGroup` is called.) `R_possible_dispatch` first checks an internal table to see if any S4 methods are set for that generic (and S4 dispatch is currently enabled for that generic), and if so proceeds to S4 dispatch using methods stored in another internal table. All primitives are in the base name space, and this mechanism means that S4 methods can be set for (some) primitives and will always be used, in contrast to setting methods on non-primitives.

The exception is `%*%`, which is S4 generic but not S3 generic as its C code contains a direct call to `R_possible_dispatch`.

The primitive `as.double` is special, as `as.numeric` and `as.real` are copies of it. The `methods` package code partly refers to generics by name and partly by function, and was modified in R 2.6.0 to map `as.double` and `as.real` to `as.numeric` (since that is the name used by packages exporting methods for it).

<sup>13</sup> although this is not recommended as it is less future-proof.

Some elements of the language are implemented as primitives, for example `}. This includes the subset and subassignment ‘functions’ and they are S4 generic, again piggybacking on S3 dispatch.`

`.BasicFunsList` is generated when `methods` is installed, by computing all primitives, initially disallowing methods on all and then setting generics for members of `.GenericArgsEnv`, the S4 group generics and a short exceptions list in `‘BasicFunsList.R’`: this currently contains the subsetting and subassignment operators and an override for `c`.

## 1.12 Memory allocators

R’s memory allocation is almost all done via routines in `‘src/main/memory.c’`. It is important to keep track of where memory is allocated, as the Windows port (by default) makes use of a memory allocator that differs from `malloc` etc as provided by MinGW. Specifically, there are entry points `Rm_malloc`, `Rm_free`, `Rm_calloc` and `Rm_free` provided by `src/gnuwin32/malloc.c`. This was done for two reasons. The primary motivation was performance: the allocator provided by MSVCRT *via* MinGW was far too slow at handling the many small allocations that the current (since R 1.2.0) allocation system for `SEXPRECs` uses. As a side benefit, we can set a limit on the amount of allocated memory: this is useful as whereas Windows does provide virtual memory it is relatively far slower than many other R platforms and so limiting R’s use of swapping is highly advantageous. The high-performance allocator is only called from `‘src/main/memory.c’`, `‘src/main/regex.c’`, `‘src/extra/pcre’` and `‘src/extra/xdr’`: note that this means that it is not used in packages.

The rest of R should where possible make use of the allocators made available by `‘src/main/memory.c’`, which are also the methods recommended in [section “Memory allocation” in \*Writing R Extensions\*](#) for use in R packages, namely the use of `R_alloc`, `Calloc`, `Realloc` and `Free`. Memory allocated by `R_alloc` is freed by the garbage collector once the ‘watermark’ has been reset by calling `vmaxset`. This is done automatically by the wrapper code calling primitives and `.Internal` functions (and also by the wrapper code to `.Call` and `.External`), but `vmaxget` and `vmaxset` can be used to reset the watermark from within internal code if the memory is only required for a short time.

All of the methods of memory allocation mentioned so far are relatively expensive. All R platforms support `alloca`, and in almost all cases<sup>14</sup> this is managed by the compiler, allocates memory on the C stack and is very efficient.

There are two disadvantages in using `alloca`. First, it is fragile and care is needed to avoid writing (or even reading) outside the bounds of the allocation block returned. Second, it increases the danger of overflowing the C stack. It is suggested that it is only used for smallish allocations (up to tens of thousands of bytes), and that

```
R_CheckStack();
```

is called immediately after the allocation (as R’s stack checking mechanism will warn far enough from the stack limit to allow for modest use of `alloca`). (`do_makeunique` in `‘src/main/unique.c’` provides an example of both points.)

An alternative strategy has been used for various functions which require intermediate blocks of storage of varying but usually small size, and this has been consolidated into the routines in the header file `‘src/main/RBufferUtils.h’`. This uses a structure which contains a buffer, the current size and the default size. A call to

```
R_AllocStringBuffer(size_t blen, R_StringBuffer *buf);
```

---

<sup>14</sup> but apparently not on Windows.

sets `buf->data` to a memory area of at least `blen+1` bytes. At least the default size is used, which means that for small allocations the same buffer can be reused. A call to `R_FreeStringBufferL` releases memory if more than the default has been allocated whereas a call to `R_FreeStringBuffer` frees any memory allocated.

The `R_StringBuffer` structure needs to be initialized, for example by

```
static R_StringBuffer ex_buff = {NULL, 0, MAXELTSIZE};
```

which uses a default size of `MAXELTSIZE = 8192` bytes. Most current uses have a static `R_StringBuffer` structure, which allows the (default-sized) buffer to be shared between calls to e.g. `grep` and even between functions: this will need to be changed if R ever allows concurrent evaluation threads. So the idiom is

```
static R_StringBuffer ex_buff = {NULL, 0, MAXELTSIZE};
...
char *buf;
for(i = 0; i < n; i++) {
    compute len
    buf = R_AllocStringBuffer(len, &ex_buff);
    use buf
}
/* free allocation if larger than the default, but leave
   default allocated for future use */
R_FreeStringBufferL(&ex_buff);
```

### 1.12.1 Internals of `R_alloc`

The memory used by `R_alloc` is allocated as R vectors, of type `RAWSXP` for ‘small’ allocations (less than  $2^{31} - 1$  bytes) and of type `REALSXP` for allocations up to  $2^{34} - 1$  bytes on 64-bit machines. Thus the allocation is in units of 8 bytes, and is rounded up. (Prior to R 2.6.0 `CHARSXP`s were used, and so one byte was added prior to rounding up. This had the effect of over-allocating areas for `doubles` by one and thereby masked several subtle programming errors.)

The vectors allocated are protected via the setting of `R_VStack`, as the garbage collector marks everything that can be reached from that location. When a vector is `R_allocated`, its `ATTRIB` pointer is set to the current `R_VStack`, and `R_VStack` is set to the latest allocation. Thus `R_VStack` is a single-linked chain of vectors currently allocated via `R_alloc`. Function `vmaxset` resets the location `R_VStack`, and should be to a value that has previously be obtained *via* `vmaxget`: allocations after the value was obtained will no longer be protected and hence available for garbage collection.

## 1.13 Internal use of global and base environments

This section notes known use by the system of these environments: the intention is to minimize or eliminate them.

### 1.13.1 Base environment

The graphics devices system maintains two variables `.Device` and `.Devices` in the base environment: both are always set. The variable `.Devices` gives a list of character vectors of the names of open devices, and `.Device` is the element corresponding to the currently active device. The null device will always be open.

There appears to be a variable `.Options`, a pairlist giving the current options settings. But in fact this is just a symbol with a value assigned, and so shows up as a base variable.

Similarly, the evaluator creates a symbol `.Last.value` which appears as a variable in the base environment.

Errors can give rise to objects `.Traceback` and `last.warning` in the base environment.

### 1.13.2 Global environment

The seed for the random number generator is stored in object `.Random.seed` in the global environment.

Some error handlers may give rise to objects in the global environment: for example `dump.frames` by default produces `last.dump`.

The `windows()` device makes use of a variable `.SavedPlots` to store display lists of saved plots for later display. This is regarded as a variable created by the user.

## 1.14 Modules

R makes use of a number of shared objects/DLLs stored in the ‘modules’ directory. These are parts of the code which have been chosen to be loaded ‘on demand’ rather than linked as dynamic libraries or incorporated into the main executable/dynamic library.

For a few of these (e.g. `vfonts`) the issue is size: the database for the Hershey fonts is included in the C code of the module and was at one time an appreciable part of the codebase for a rarely used feature. However, for most of the modules the motivation has been the amount of (often optional) code they will bring in via libraries to which they are linked.

- `internet` The internal HTTP and FTP clients and socket support, which link to system-specific support libraries.
- `lapack` The code which makes use of the LAPACK library, and is linked to ‘libRlapack’ or an external LAPACK library.
- `vfonts` The Hershey font databases and the code to draw from them.
- `X11` (Unix-alikes only.) The `X11()`, `jpeg()` and `png()` devices. These are optional, and link to the `X11`, `jpeg` and `libpng` libraries.
- ‘`Rbitmap.dll`’  
(Windows only.) The code for the BMP, JPEG and PNG devices and for saving on-screen graphs to those formats. This is technically optional, and needs source code not in the tarball.
- ‘`Rhtml.dll`’  
(Windows only.) A link to an ActiveX control that displays Compiled HTML help. This is optional, and only compiled if CHTML is specified.
- ‘`iconv.dll`’  
(Windows only.) A DLL compiled via Visual C++ which contains the routines to convert between character sets.
- ‘`internet2.dll`’  
(Windows only.) An alternative version of the internet access routines, compiled against Internet Explorer internals (and so loads ‘`wininet.dll`’ and ‘`wsock32.dll`’).

## 2 `.Internal` vs `.Primitive`

C code compiled into R at build time can be called “directly” or via the `.Internal` interface, which is very similar to the `.External` interface except in syntax. More precisely, R maintains a table of R function names and corresponding C functions to call, which by convention all start with `do_` and return a SEXP. Via this table (`R_FunTab` in file `src/main/names.c`) one can also specify how many arguments to a function are required or allowed, whether the arguments are to be evaluated before calling or not, and whether the function is “internal” in the sense that it must be accessed via the `.Internal` interface, or directly accessible in which case it is printed in R as `.Primitive`.

R’s functionality can also be extended by providing corresponding C code and adding to this function table.

In general, all such functions use `.Internal()` as this is safer and in particular allows for transparent handling of named and default arguments. For example, `axis` is defined as

```
axis <- function(side, at = NULL, labels = NULL, ...)
  .Internal(axis(side, at, labels, ...))
```

However, for reasons of convenience and also efficiency (as there is some overhead in using the `.Internal` interface wrapped in a function closure), there are exceptions which can be accessed directly. Note that these functions make no use of R code, and hence are very different from the usual interpreted functions. In particular, `args`, `formals` and `body` return `NULL` for such objects, and argument matching is purely positional (with two exceptions described below).

The list of these “primitive” functions is subject to change: currently, it includes the following.

1. “Special functions” which really are *language* elements, however exist as “primitive” functions in R:

```
{      (      if      for      while repeat break next
return function quote on.exit
```

2. Language elements and basic *operators* (i.e., functions usually *not* called as `foo(a, b, ...)`) for subsetting, assignment, arithmetic and logic. These are the following 1-, 2-, and  $N$ -argument functions:

```
<-  <<-  =      [      [[      $      @
      [<-  [[<-  $<-

+    -    *    /    ^    %%    %*%  %/%
<    <=   ==   !=   >=   >
|    ||   &    &&   !
```

3. “Low level” 0- and 1-argument functions which belong to one of the following groups of functions:

- a. Basic mathematical functions with a single argument, i.e.,

```
abs      sign      sqrt
floor    ceiling

exp      expm1
log2     log10    log1p
cos      sin      tan
acos     asin     atan
cosh     sinh     tanh
acosh    asinh    atanh

gamma    lgamma   digamma  trigamma
```

```
cumsum  cumprod cummax  cummin
```

```
Im Re Arg Conj Mod
```

`log` is a function of one or two arguments, but was made primitive as from R 2.6.0 and so has named rather than positional matching for back compatibility.

`trunc` is a difficult case: it is a primitive that can have zero or more arguments: the default method handled in the primitive has only one.

- b. Functions rarely used outside of “programming” (i.e., mostly used inside other functions), such as

```
nargs          missing
interactive    is.xxx
.Primitive    .Internal
globalenv     baseenv      emptyenv      pos.to.env
unclass
seq_along     seq_len
```

(where *xxx* stands for 27 different notions, such as `function`, `vector`, `numeric`, and so forth, but not `is.loaded`).

- c. The programming and session management utilities

```
debug  undebug  browser  proc.time  gc.time
tracemem  retracemem  untracemem
```

4. The following basic replacement and extractor functions

```
length      length<-
class       class<-
oldClass    oldClass<-
attr        attr<-
attributes  attributes<-
names       names<-
dim         dim<-
dimnames    dimnames<-
environment<-
levels<-
storage.mode<-
```

Note that optimizing `NAMED = 1` is only effective within a primitive (as the closure wrapper of a `.Internal` will set `NAMED = 2` when the promise to the argument is evaluated) and hence replacement functions should where possible be primitive to avoid copying (at least in their default methods).

5. The following few *N*-argument functions are “primitive” for efficiency reasons:

```
:          ~          c          list
call       as.call    as.character  as.complex  as.double
as.integer as.logical  as.raw
expression substitute  as.environment
UseMethod  invisible    standardGeneric
.C         .Fortran    .Call        .External
.Call.graphics  .External.graphics
.subset    .subset2    .primTrace   .primUntrace
rep        seq.int
lazyLoadDBfetch
```

`rep` and `seq.int` manage their own argument matching and so do work in the standard way.

## 2.1 Special primitives

A small number of primitives are *specials* rather than *builtins*, that is they are entered with unevaluated arguments. This is clearly necessary for the language constructs and the assignment operators. `&&` and `||` conditionally evaluate their second argument, and `~`, `.Internal`, `call`, `expression` and `missing` do not evaluate their arguments.

`rep` and `seq.int` are special as they evaluate some of their arguments conditional on which are non-missing. `c` is special to allow it to be used with language objects.

The subsetting, subassignment and `@` operators are all special. (For both extraction and replacement forms, `$` and `@` take a symbol argument, and `[` and `[[` allow missing arguments.)

`UseMethod` is special to avoid the additional contexts added to calls to builtins when profiling (via `Rprof`).

## 2.2 Special internals

There are also special `.Internal` functions: `switch`, `Recall`, `cbind`, `rbind` (to allow for the `deparse.level` argument), `lapply`, `eapply` and `NextMethod`.

## 2.3 Prototypes for primitives

As from R 2.5.0, prototypes are available for the primitive functions and operators, and there are used for printing, `args` and package checking (e.g. by `tools::checkS3methods` and by package `codetools`). There are two environments in the `base` package (and name space), `.GenericArgsEnv` for those primitives which are internal S3 generics, and `.ArgsEnv` for the rest. Those environments contain closures with the same names as the primitives, formal arguments derived (manually) from the help pages, a body which is a suitable call to `UseMethod` or `NULL` and environment the base name space.

The C code for `print.default` and `args` uses the closures in these environments in preference to the definitions in base (as primitives).

The QC function `undoc` checks that all the functions prototyped in these environments are currently primitive, and that the primitives not included are better thought of as language elements (at the time of writing

```
$ $<- && ( : @ [ [[ [[<- [<- { || ~ <- <<- =
break for function if next repeat return while
```

. One could argue about `~`, but it is known to the parser and has semantics quite unlike a normal function. And `:` is documented with different argument names in its two meanings.)

The QC functions `codoc` and `checkS3methods` also make use of these environments (effectively placing them in front of base in the search path), and hence the formals of the functions they contain are checked against the help pages by `codoc`. However, there are two problems with the generic primitives. The first is that many of the operators are part of the S3 group generic `Ops` and that defines their arguments to be `e1` and `e2`: although it would be very unusual, an operator could be called as e.g. `"+"(e1=a, e2=b)` and if method dispatch occurred to a closure, there would be an argument name mismatch. So the definitions in environment `.GenericArgsEnv` have to use argument names `e1` and `e2` even though the traditional documentation is in terms of `x` and `y`: `codoc` makes the appropriate adjustment via `tools:::make_S3_primitive_generic_env`. The second discrepancy is with the `Math` group generics, where the group generic is defined with argument list `(x, ...)`, but most of the members only allow one argument when used as the default method (and `round` and `signif` allow two as default methods): again fix-ups are used.

Those primitives which are in `.GenericArgsEnv` are checked (via `'tests/primitives.R'` to be generic *via* defining methods for them, and a check is made that the remaining primitives are probably not generic, by setting a method and checking it is not dispatched to (but this can fail for other reasons). However, there is no certain way to know that if other `.Internal` or primitive functions are not internally generic except by reading the source code.

## 3 Internationalization in the R sources

The process of marking messages (errors, warnings etc) for translation in an R package is described in [section “Localization” in \*Writing R Extensions\*](#), and the standard packages included with R have (with an exception in **grDevices**) been internationalized in the same way as other packages.

### 3.1 R code

Internationalization for R code is done in exactly the same way as for extension packages. As all standard packages which have R code also have a namespace, it is never necessary to specify domain, but for efficiency calls to `message`, `warning` and `stop` should include `domain = NA` when the message is constructed *via* `gettextf`, `gettext` or `ngettext`.

For each package, the extracted messages and translation sources are stored under package directory ‘po’ in the source package, and compiled translations under ‘inst/po’ for installation to package directory ‘po’ in the installed package. This also applies to C code in packages.

### 3.2 Main C code

The main C code (e.g. that in ‘src/\*/\*.c’ and in the modules) is where R is closest to the sort of application for which ‘gettext’ was written. Messages in the main C code are in domain R and stored in the top-level directory ‘po’ with compiled translations under ‘share/locale’.

The list of files covered by the R domain is specified in file ‘po/POTFILES.in’.

The normal way to mark messages for translation is via `_("msg")` just as for packages. However, sometimes one needs to mark passages for translation without wanting them translated at the time, for example when declaring string constants. This is the purpose of the `N_` macro, for example

```
{ ERROR_ARGTYPE,          N_("invalid argument type")},
from ‘src/main/errors.c’.
```

A macro

```
#ifdef ENABLE_NLS
#define P_(StringS, StringP, N) ngettext (StringS, StringP, N)
#else
#define P_(String, StringP, N) (N > 1 ? StringP: String)
#endif
```

as a wrapper for `ngettext`: however in some cases the preferred approach has been to conditionalize (on `ENABLE_NLS`) code using `ngettext`.

The macro `_("msg")` can safely be used in ‘src/appl’; the header for standalone ‘`nmath`’ skips possible translation. (This does not apply to `N_` or `P_`).

### 3.3 Windows-GUI-specific code

Messages for the Windows GUI are in a separate domain ‘RGui’. This was done for two reasons:

The translators for the Windows version of R might be separate from those for the rest of R (familiarity with the GUI helps), and

Messages for Windows are most naturally handled in the native charset for the language, and in the case of CJK languages the charset is Windows-specific. (It transpires that as the `iconv` we ported works well under Windows, this is less important than anticipated.)

Messages for the ‘RGui’ domain are marked by `G_("msg")`, a macro that is defined in ‘`src/gnuwin32/win-nls.h`’. The list of files that are considered is hardcoded in the `RGui.pot-update` target of ‘`po/Makefile.in.in`’: note that this includes ‘`devWindows.c`’ as the menus on the `windows` device are considered to be part of the GUI. (There is also `GN_("msg")`, the analogue of `N_("msg")`.)

The template and message catalogs for the ‘RGui’ domain are in the top-level ‘`po`’ directory.

## 3.4 MacOS X GUI

This is handled separately: see <http://developer.r-project.org/Translations.html>.

## 3.5 Updating

See ‘`po/README`’ for how to update the message templates and catalogs.

## 4 R coding standards

R is meant to run on a wide variety of platforms, including Linux and most variants of Unix as well as 32-bit Windows versions and on MacOS X. Therefore, when extending R by either adding to the R base distribution or by providing an add-on package, one should not rely on features specific to only a few supported platforms, if this can be avoided. In particular, although most R developers use GNU tools, they should not employ the GNU extensions to standard tools. Whereas some other software packages explicitly rely on e.g. GNU make or the GNU C++ compiler, R does not. Nevertheless, R is a GNU project, and the spirit of the GNU *Coding Standards* should be followed if possible.

The following tools can “safely be assumed” for R extensions.

- An ISO C99 C compiler. Note that extensions such as POSIX 1003.1 must be tested for, typically using Autoconf unless you are sure they are supported on all mainstream R platforms (including Windows and MacOS X). Packages will be more portable if written assuming only C89, but this should not be done where using C99 features will make for cleaner or more robust code.
- A FORTRAN 77 compiler (but not Fortran 9x).
- A simple `make`, considering the features of `make` in 4.2 BSD systems as a baseline.

GNU or other extensions, including pattern rules using ‘%’, the automatic variable ‘\$^’, the ‘+=’ syntax to append to the value of a variable, the (“safe”) inclusion of makefiles with no error, conditional execution, and many more, must not be used (see Chapter “Features” in the GNU *Make Manual* for more information). On the other hand, building R in a separate directory (not containing the sources) should work provided that `make` supports the `VPATH` mechanism.

Windows-specific makefiles can assume GNU `make` 3.75 or later, as no other `make` is viable on that platform.

- A Bourne shell and the “traditional” Unix programming tools, including `grep`, `sed`, and `awk`.

There are POSIX standards for these tools, but these may not fully be supported. Baseline features could be determined from a book such as *The UNIX Programming Environment* by Brian W. Kernighan & Rob Pike. Note in particular that ‘|’ in a regexp is an extended regexp, and is not supported by all versions of `grep` or `sed`. The Open Group Base Specifications, Issue 6, which is technically identical to ISO/IEC 9945 and IEEE Std 1003.1 (POSIX), 2004, are available at <http://www.opengroup.org/onlinepubs/009695399/mindex.html>.

Under Windows, most users will not have these tools installed, and you should not require their presence for the operation of your package. However, users who install your package from source will have them, as they can be assumed to have followed the instructions in “the Windows toolset” appendix of the “R Installation and Administration” manual to obtain them. Redirection cannot be assumed to be available via `system` as this does not use a standard shell (let alone a Bourne shell).

In addition, the following tools are needed for certain tasks.

- Perl version 5 is needed for converting documentation written in Rd format to plain text, HTML, , and to extract the examples. In addition, several other tools, in particular `check` and `build` require Perl.

The R Core Team has decided that Perl (version 5) can safely be assumed for building R from source, building and checking add-on packages, and for installing add-on packages from source. On the other hand, Perl cannot be assumed at all for installing *binary* (pre-built) versions of add-on packages, or at run time.

- Makeinfo version 4.7 or later is needed to build the Info files for the R manuals written in the GNU Texinfo system. (Future distributions of R may contain the Info files.)

It is also important that code is written in a way that allows others to understand it. This is particularly helpful for fixing problems, and includes using self-descriptive variable names, commenting the code, and also formatting it properly. The R Core Team recommends to use a basic indentation of 4 for R and C (and most likely also Perl) code, and 2 for documentation in Rd format. Emacs users can implement this indentation style by putting the following in one of their startup files. (For GNU Emacs 20: for GNU Emacs 21 or later use customization to set the `c-default-style` to "bsd" and `c-basic-offset` to 4.)

```
;;; C
(add-hook 'c-mode-hook
  (lambda () (c-set-style "bsd")))

;;; ESS
(add-hook 'ess-mode-hook
  (lambda ()
    (ess-set-style 'C++)
    ;; Because
    ;;
    ;; DEF GNU BSD K&R C++
    ;; ess-indent-level          2  2  8  5  4
    ;; ess-continued-statement-offset  2  2  8  5  4
    ;; ess-brace-offset            0  0 -8 -5 -4
    ;; ess-arg-function-offset      2  4  0  0  0
    ;; ess-expression-offset        4  2  8  5  4
    ;; ess-else-offset             0  0  0  0  0
    ;; ess-close-brace-offset       0  0  0  0  0
    (add-hook 'local-write-file-hooks
      (lambda ()
        (ess-nuke-trailing-whitespace))))))
(setq ess-nuke-trailing-whitespace-p 'ask)
;; or even
;; (setq ess-nuke-trailing-whitespace-p t)
;;; Perl
(add-hook 'perl-mode-hook
  (lambda () (setq perl-indent-level 4)))
```

(The 'GNU' styles for Emacs' C and R modes use a basic indentation of 2, which has been determined not to display the structure clearly enough when using narrow fonts.)

## 5 Testing R code

When you (as R developer) add new functions to the R base (all the packages distributed with R), be careful to check if *make test-Specific* or particularly, *cd tests; make no-segfault.Rout* still works (without interactive user intervention, and on a standalone computer). If the new function, for example, accesses the Internet, or requires GUI interaction, please add its name to the “stop list” in ‘*tests/no-segfault.Rin*’.

[To be revised: use *make check-devel*, check the write barrier if you change internal structures.]

## Function and variable index

(Index is nonexistent)

## Concept index

(Index is nonexistent)