

Seamless R Extensions using Rcpp and RInside

Dirk Eddebuettel
Debian & R

Joint work with Romain François

Presentation on March 30, 2010 to
UCLA Department of Statistics (3pm)
Los Angeles R Users Group (6pm)

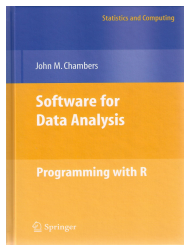


Outline

- 1 Extending R
 - Why ?
 - The standard API
 - Inline
- 2 Rcpp
 - Overview
 - New API
 - Examples
- 3 Rcpp Usage Examples
 - RInside
 - Others
- 4 Summary
 - Key points
 - Resources



Motivation



Chambers. *Software for Data Analysis: Programming with R*. Springer, 2008

Chambers (2008) opens chapter 11 (*Interfaces I: Using C and Fortran*) with these words:

Since the core of R is in fact a program written in the C language, it's not surprising that the most direct interface to non-R software is for code written in C, or directly callable from C. All the same, including additional C code is a serious step, with some added dangers and often a substantial amount of programming and debugging required. You should have a good reason.



Motivation

Chambers (2008) then proceeds with this rough map of the road ahead:

Against:

- It's more work
- Bugs will bite
- Potential platform dependency
- Less readable software

In Favor:

- New and trusted computations
- Speed
- Object references

So is the deck stacked against us?



Outline

- 1 Extending R
 - Why ?
 - The standard API
 - Inline
- 2 Rcpp
 - Overview
 - New API
 - Examples
- 3 Rcpp Usage Examples
 - RInside
 - Others
- 4 Summary
 - Key points
 - Resources



Compiled Code: The Basics

R offers several functions to access compiled code: we focus on `.C` and `.Call` here. (*R Extensions*, sections 5.2 and 5.9; *Software for Data Analysis*).

The canonical example is the convolution function:

```
1 void convolve(double *a, int *na, double *b,
2              int *nb, double *ab)
3 {
4   int i, j, nab = *na + *nb - 1;
5
6   for(i = 0; i < nab; i++)
7     ab[i] = 0.0;
8   for(i = 0; i < *na; i++)
9     for(j = 0; j < *nb; j++)
10      ab[i + j] += a[i] * b[j];
11 }
```



Compiled Code: The Basics cont.

The convolution function is called from R by

```
1 conv <- function(a, b)
2   .C("convolve",
3     as.double(a),
4     as.integer(length(a)),
5     as.double(b),
6     as.integer(length(b)),
7     ab = double(length(a) + length(b) - 1))$ab
```

As stated in the manual, one must take care to coerce all the arguments to the correct R storage mode before calling `.C` as mistakes in matching the types can lead to wrong results or hard-to-catch errors.



Compiled Code: The Basics cont.

Using `.Call`, the example becomes

```
1 #include <R.h>
2 #include <Rdefines.h>
3
4 extern "C" SEXP convolve2(SEXP a, SEXP b)
5 {
6     int i, j, na, nb, nab;
7     double *xa, *xb, *xab;
8     SEXP ab;
9
10    PROTECT(a = AS_NUMERIC(a));
11    PROTECT(b = AS_NUMERIC(b));
12    na = LENGTH(a); nb = LENGTH(b); nab = na + nb - 1;
13    PROTECT(ab = NEW_NUMERIC(nab));
14    xa = NUMERIC_POINTER(a); xb = NUMERIC_POINTER(b);
15    xab = NUMERIC_POINTER(ab);
16    for(i = 0; i < nab; i++) xab[i] = 0.0;
17    for(i = 0; i < na; i++)
18        for(j = 0; j < nb; j++) xab[i + j] += xa[i] * xb[j];
19    UNPROTECT(3);
20    return(ab);
21 }
```



Compiled Code: The Basics cont.

Now the call simplifies to just the function name and the vector arguments—all other handling is done at the C/C++ level:

```
1 conv <- function(a, b) .Call("convolve2", a, b)
```

In summary, we see that

- there are different entry points
- using different calling conventions
- leading to code that may need to do more work at the lower level.



Outline

- 1 Extending R
 - Why ?
 - The standard API
 - **Inline**
- 2 Rcpp
 - Overview
 - New API
 - Examples
- 3 Rcpp Usage Examples
 - RInside
 - Others
- 4 Summary
 - Key points
 - Resources



Compiled Code: inline

`inline` is a package by Oleg Sklyar et al that provides the function `cfunction` which can wrap Fortran, C or C++ code.

```
1 ## A simple Fortran example
2 code <- "
3     integer i
4     do 1 i=1, n(1)
5     1 x(i) = x(i)**3
6 "
7 cubefn <- cfunction(signature(n="integer", x="numeric"),
8                    code, convention=".Fortran")
9 x <- as.numeric(1:10)
10 n <- as.integer(10)
11 cubefn(n, x)$x
```

`cfunction` takes care of compiling, linking, loading, ... by placing the resulting dynamically-loadable object code in the per-session temporary directory used by R.



Outline

- 1 Extending R
 - Why ?
 - The standard API
 - Inline
- 2 Rcpp
 - Overview
 - New API
 - Examples
- 3 Rcpp Usage Examples
 - RInside
 - Others
- 4 Summary
 - Key points
 - Resources



Compiled Code: Rcpp

In a nutshell:

- `Rcpp` makes it easier to interface C++ and R code.
- Using the `.Call` interface, we can use features of the C++ language to automate the tedious bits of the macro-based C-level interface to R.
- One major advantage of using `.Call` is that richer R objects (vectors, matrices, lists, ... in fact most SEXP types incl functions, environments etc) can be passed directly between R and C++ without the need for explicit passing of dimension arguments.
- By using the C++ class layers, we do not need to manipulate the SEXP objects using any of the old-school C macros.
- `inline` eases usage, development and testing.



Rcpp example

The convolution example can be rewritten in the 'Classic API':

```
1 #include <Rcpp.h>
2
3 RcppExport SEXP convolve_cpp(SEXP a, SEXP b)
4 {
5     RcppVector<double> xa(a);
6     RcppVector<double> xb(b);
7
8     int nab = xa.size() + xb.size() - 1;
9
10    RcppVector<double> xab(nab);
11    for (int i = 0; i < nab; i++) xab(i) = 0.0;
12
13    for (int i = 0; i < xa.size(); i++)
14        for (int j = 0; j < xb.size(); j++)
15            xab(i + j) += xa(i) * xb(j);
16
17    RcppResultSet rs;
18    rs.add("ab", xab);
19    return rs.getReturnList();
20 }
```



Outline

- 1 Extending R
 - Why ?
 - The standard API
 - Inline
- 2 Rcpp
 - Overview
 - **New API**
 - Examples
- 3 Rcpp Usage Examples
 - RInside
 - Others
- 4 Summary
 - Key points
 - Resources



Rcpp: The 'New API'

Rcpp was significantly extended over the last few months to permit more natural expressions. Consider this comparison between the R API and the new Rcpp API:

```

1 SEXP ab;
2 PROTECT(ab = allocVector(STRSXP, 2));
3 SET_STRING_ELT( ab, 0, mkChar("foo") );
4 SET_STRING_ELT( ab, 1, mkChar("bar") );
5 UNPROTECT(1);

```

```

1 CharacterVector ab(2) ;
2 ab[0] = "foo" ;
3 ab[1] = "bar" ;

```

Data types, including STL containers and iterators, can be nested. and other niceties. Implicit converters allow us to combine types:

```

1 std::vector<double> vec;
2 [...]
3 List x(3);
4 x[0] = vec;
5 x[1] = "some text";
6 x[2] = 42;

```

```

1 // With Rcpp 0.7.11 or later we can do:
2 std::vector<double> vec;
3 [...]
4 List x = List::create(vec,
5                       "some text",
6                       42);

```


Functional programming in both languages

In R, functional programming is easy:

```

1 R> data(faithful); lapply(faithful, summary)
2 $eruptions
3   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
4   1.60   2.16   4.00   3.49   4.45   5.10
5
6 $waiting
7   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
8   43.0   58.0   76.0   70.9   82.0   96.0

```

We can do that in C++ as well and pass the R function down to the data elements we let the STL iterate over:

```

1 src <- 'Rcpp::List input(data);
2       Rcpp::Function f(fun);
3       Rcpp::List output(input.size());
4       std::transform(input.begin(), input.end(), output.begin(), f);
5       output.names() = input.names();
6       return output; '
7 cpp_lapply <- cfunction(signature(data="list", fun = "function"), src, Rcpp = TRUE)

```



Exception handling

Automatic catching and conversion of C++ exceptions:

```
R> library(Rcpp); library(inline)
R> cpp <- '
+   Rcpp::NumericVector x(xs); // automatic conversion from SEXP
+   for (int i=0; i<x.size(); i++) {
+       if (x[i] < 0)
+           throw std::range_error("Non-negative values required");
+       x[i] = log(x[i]);
+   }
+   return x; // automatic conversion to SEXP
+ '
```

```
R> fun <- cfunction(signature(xs="numeric"), cpp, Rcpp=TRUE)
R> fun( seq(2, 5) )

[1] 0.6931 1.0986 1.3863 1.6094

R> fun( seq(5, -2) )

Error in fun(seq(5, -2)) : Non-negative values required

R> fun( LETTERS[1:5] )

Error in fun(LETTERS[1:5]) : not compatible with INTSXP
R>
```



Outline

- 1 Extending R
 - Why ?
 - The standard API
 - Inline
- 2 Rcpp
 - Overview
 - New API
 - **Examples**
- 3 Rcpp Usage Examples
 - RInside
 - Others
- 4 Summary
 - Key points
 - Resources



Rcpp example

The convolution example can be rewritten in the new API:

```
1 #include <Rcpp.h>
2
3 RcppExport SEXP convolve_cpp(SEXP a, SEXP b){
4   Rcpp::NumericVector xa(a); // automatic conversion from SEXP
5   Rcpp::NumericVector xb(b);
6
7   int n_xa = xa.size();
8   int n_xb = xb.size();
9   int nab = n_xa + n_xb - 1;
10
11  Rcpp::NumericVector xab(nab);
12
13  for (int i = 0; i < n_xa; i++)
14    for (int j = 0; j < n_xb; j++)
15      xab[i + j] += xa[i] * xb[j];
16
17  return xab; // automatic conversion to SEXP
18 }
```



Speed comparison

In a recently-submitted paper, the following table summarises the performance of convolution examples:

Implementation	Time in millisec	Relative to R API
R API (as benchmark)	32	
<code>RcppVector<double></code>	354	11.1
<code>NumericVector::operator[]</code>	52	1.6
<code>NumericVector::begin</code>	33	1.0

Table 1: Performance for convolution example

We averaged 1000 replications with two 100-element vectors – see `examples/ConvolveBenchmarks/` in `Rcpp` for details.



Another Speed Comparison Example

- Regression is a key component of many studies. In simulations, we often want to run a very large number of regressions.
- R has `lm()` as the general purposes function. It is very powerful and returns a rich object—but it is not *lightweight*.
- For this purpose, R has `lm.fit()`. But, this does not provide all relevant auxiliary data as *e.g.* the standard error of the estimate.
- For the most recent *Introduction to High-Performance Computing with R* tutorial, I had written a hybrid R/C/C++ solution using the GNU GSL.
- We complement this with a new C++ implementation around the Armadillo linear algebra classes.



Linear regression via GSL: lmGSL()

```

1  lmGSL <- function() {
2    src <- '
3
4    RcppVectorView<double> Yr(Ysexp);
5    RcppMatrixView<double> Xr(Xsexp);
6
7    int i, j, n = Xr.dim1(), k = Xr.dim2();
8    double chi2;
9
10   gsl_matrix *X = gsl_matrix_alloc(n,k);
11   gsl_vector *y = gsl_vector_alloc(n);
12   gsl_vector *c = gsl_vector_alloc(k);
13   gsl_matrix *cov = gsl_matrix_alloc(k,k);
14
15   for (i = 0; i < n; i++) {
16     for (j = 0; j < k; j++) {
17       gsl_matrix_set (X, i, j, Xr(i,j));
18     }
19     gsl_vector_set (y, i, Yr(i));
20   }
21
22   gsl_multifit_linear_workspace *wk =
23     gsl_multifit_linear_alloc (n,k);
24   gsl_multifit_linear(X,y,c,cov,&chi2,wk);
25   gsl_multifit_linear_free (wk);
26   RcppVector<double> StdErr(k);
27   RcppVector<double> Coef(k);

```

```

28   for (i = 0; i < k; i++) {
29     Coef(i) = gsl_vector_get(c, i);
30     StdErr(i) =
31       sqrt(gsl_matrix_get(cov, i, i));
32   }
33
34   gsl_matrix_free (X);
35   gsl_vector_free (y);
36   gsl_vector_free (c);
37   gsl_matrix_free (cov);
38
39   RcppResultSet rs;
40   rs.add("coef", Coef);
41   rs.add("stderr", StdErr);
42
43   return = rs.getReturnList();
44   '
45   ## turn into a function that R can call
46   ## args redundant on Debian/Ubuntu
47   fun <-
48     cfunction (signature (Ysexp="numeric",
49       Xsexp="numeric"), src,
50     includes=
51       "#include <gsl/gsl_multifit.h>",
52     Rcpp=TRUE,
53     cppargs="-I /usr/include",
54     libargs="-lgsl -lgslcblas")
55 }

```

Linear regression via Armadillo: ImArmadillo example

```

1  ImArmadillo <- function() {
2    src <- '
3    Rcpp::NumericVector yr(Ysexp);
4    Rcpp::NumericVector Xr(Xsexp);           // actually an n x k matrix
5    std::vector<int> dims = Xr.attr("dim");
6    int n = dims[0], k = dims[1];
7    arma::mat X(Xr.begin(), n, k, false);    // use advanced armadillo constructors
8    arma::colvec y(yr.begin(), yr.size());
9    arma::colvec coef = solve(X, y);         // model fit
10   arma::colvec resid = y - X*coef;         // to comp. std.errr of the coefficients
11   arma::mat covmat = trans(resid)*resid/(n-k) * arma::inv(arma::trans(X)*X);
12
13   Rcpp::NumericVector coefr(k), stderrestr(k);
14   for (int i=0; i<k; i++) {                // with RcppArmadillo template converters
15     coefr[i] = coef[i];                    // this would not be needed but we only
16     stderrestr[i] = sqrt(covmat(i,i));     // have Rcpp.h here
17   }
18
19   return Rcpp::List::create( Rcpp::Named( "coefficients", coefr), // Rcpp 0.7.11
20                             Rcpp::Named( "stderr", stderrestr));
21
22
23   ## turn into a function that R can call
24   fun <- cfunction(signature(Ysexp="numeric", Xsexp="numeric"),
25                    src, includes="#include <armadillo>", Rcpp=TRUE,
26                    cppargs="-I/usr/include", libargs="-larmadillo")
27 }

```


Linear regression via Armadillo: RcppArmadillo

`fastLm` in the new `RcppArmadillo` does even better:

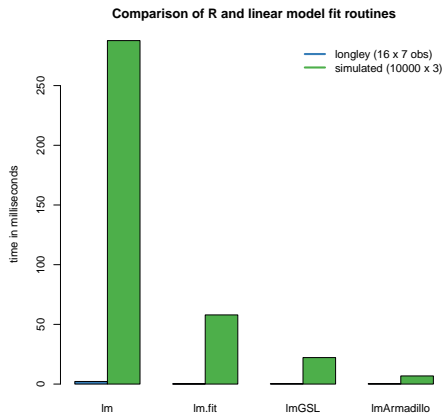
```

1 #include <RcppArmadillo.h>
2
3 extern "C" SEXP fastLm(SEXP ys, SEXP Xs) {
4   Rcpp::NumericVector yr(ys); // creates Rcpp vector from SEXP
5   Rcpp::NumericMatrix Xr(Xs); // creates Rcpp matrix from SEXP
6   int n = Xr.nrow(), k = Xr.ncol();
7
8   arma::mat X(Xr.begin(), n, k, false); // reuses memory and avoids extra copy
9   arma::colvec y(yr.begin(), yr.size(), false);
10
11   arma::colvec coef = arma::solve(X, y); // fit model y ~ X
12   arma::colvec resid = y - X*coef; // residuals
13
14   double sig2 = arma::as_scalar( arma::trans(resid)*resid/(n-k) ); // std. err est
15   arma::colvec sdest = arma::sqrt(sig2*arma::diagvec(arma::inv(arma::trans(X)*X)));
16
17   return Rcpp::List::create( // requires Rcpp 0.7.11
18     Rcpp::Named("coefficients") = coef,
19     Rcpp::Named("stderr") = sdest
20   );
21 }

```



Rcpp Example: Regression timings



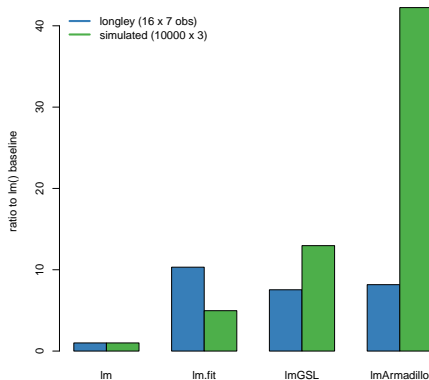
The small `longley` example exhibits less variability between methods, but the larger data set shows the gains more clearly.

For the small data set, all three appear to improve similarly on `lm`.

Source: Our calculations, see `examples/FastLM/` in `Rcpp`.

Another Rcpp example (cont.)

Comparison of R and linear model fit routines



By dividing the `lm` time by the respective times, we obtain the 'possible gains' from switching.

One caveat, measurements depends critically on the size of the data as well as the cpu and libraries that are used.

Source: Our calculations, see `examples/FastLM/` in `Rcpp`.



Outline

- 1 Extending R
 - Why ?
 - The standard API
 - Inline
- 2 Rcpp
 - Overview
 - New API
 - Examples
- 3 Rcpp Usage Examples
 - RInside
 - Others
- 4 Summary
 - Key points
 - Resources



From RApache to littler to RInside

Jeff Horner's work on [RApache](#) lead to joint work in [littler](#), a scripting / cmdline front-end. As it embeds R and simply 'feeds' the REPL loop, the next step was to embed R in proper C++ classes: [RInside](#).

```
1 #include <RInside.h> // for the embedded R via RInside
2
3 int main(int argc, char *argv[]) {
4     RInside R(argc, argv); // create an embedded R instance
5
6     R["txt"] = "Hello, world!\n"; // assign a char* (string) to 'txt'
7
8     R.parseEvalQ("cat(txt)"); // eval the init string, ignoring any returns
9
10    exit(0);
11 }
12 }
```



Another simple example

This example shows some of the new assignment and converter code:

```
1
2 #include <RInside.h>                // for the embedded R via RInside
3
4 int main(int argc, char *argv[]) {
5     RInside R(argc, argv);         // create an embedded R instance
6
7     R["x"] = 10 ;
8     R["y"] = 20 ;
9
10    R.parseEvalQ("z <- x + y") ;
11
12    int sum = R["z"];
13
14    std::cout << "10 + 20 = " << sum << std::endl ;
15    exit(0);
16 }
17
```



And another *parallel* example

```
1 // MPI C++ API version of file contributed by Jianping Hua
2
3 #include <mpi.h> // mpi header
4 #include <RInside.h> // for the embedded R via RInside
5
6 int main(int argc, char *argv[]) {
7
8     MPI::Init(argc, argv); // mpi initialization
9     int myrank = MPI::COMM_WORLD.Get_rank(); // obtain current node rank
10    int nodesize = MPI::COMM_WORLD.Get_size(); // obtain total nodes running.
11
12    RInside R(argc, argv); // create an embedded R instance
13
14    std::stringstream txt;
15    txt << "Hello from node " << myrank // node information
16        << " of " << nodesize << " nodes!" << std::endl;
17    R.assign( txt.str(), "txt"); // assign string to R variable 'txt'
18
19    std::string evalstr = "cat(txt)"; // show node information
20    R.parseEvalQ(evalstr); // eval the string, ign. any returns
21
22    MPI::Finalize(); // mpi finalization
23
24    exit(0);
25 }
```



RInside workflow

- C++ programs compute, gather or aggregate raw data.
- Data is saved and analysed before a new 'run' is launched.
- With `RInside` we now skip a step:
 - collect data in a vector or matrix
 - pass data to `R` — easy thanks to `Rcpp` wrappers
 - pass one or more short 'scripts' as strings to `R` to evaluate
 - pass data back to C++ programm — easy thanks to `Rcpp` converters
 - resume main execution based on new results
- A number of simple examples ship with `RInside`



Outline

- 1 Extending R
 - Why ?
 - The standard API
 - Inline
- 2 Rcpp
 - Overview
 - New API
 - Examples
- 3 Rcpp Usage Examples
 - RInside
 - **Others**
- 4 Summary
 - Key points
 - Resources



Users of Rcpp

- RInside uses Rcpp for object transfer and more
- RcppArmadillo (which contains fastLM())
- RcppExamples is a 'this is how you can do it' stanza
- RProtoBuf is what got Romain and me here, it may get rewritten to take more advantage of Rcpp
- RQuantLib is where Rcpp originally started
- highlight is Romain's first re-use of Rcpp
- mvabund, sdcTable, bifactorial, minqa are truly external users which are all on CRAN
- upcoming: pcaMethods (BioC), phylobase, possibly lme4
- *Your package here next?*



Outline

- 1 Extending R
 - Why ?
 - The standard API
 - Inline
- 2 Rcpp
 - Overview
 - New API
 - Examples
- 3 Rcpp Usage Examples
 - RInside
 - Others
- 4 **Summary**
 - **Key points**
 - Resources



Wrapping up

This presentation has tried to convince you that

- While the deck way be stacked against you (when adding C/C++ to R), you can still pick where to play
- R can be extended in many ways; we focus on something that allows us write extensions
 - that are efficient: we want speed and features
 - that correspond to the R object model
 - that also allow us to embed R inside C++
- And all this while retaining 'high-level' STL-alike semantics, templates and other goodies in C++
- Using C++ abstractions wisely can keep the code both clean and readable – yet very efficient



Outline

- 1 Extending R
 - Why ?
 - The standard API
 - Inline
- 2 Rcpp
 - Overview
 - New API
 - Examples
- 3 Rcpp Usage Examples
 - RInside
 - Others
- 4 **Summary**
 - Key points
 - **Resources**



Some pointers

- <http://dirk.eddelbuettel.com/code/rcpp.html>
- <http://romainfrancois.blog.free.fr/index.php?category/R-package/Rcpp>
- <http://cran.r-project.org/package=Rcpp>
- <http://r-forge.r-project.org/projects/rcpp/>
- and likewise for RInside, RProtoBuf and more.



The end

Thank you!

