# Python and R Together at Last

*Writing Cross-Language Tools*

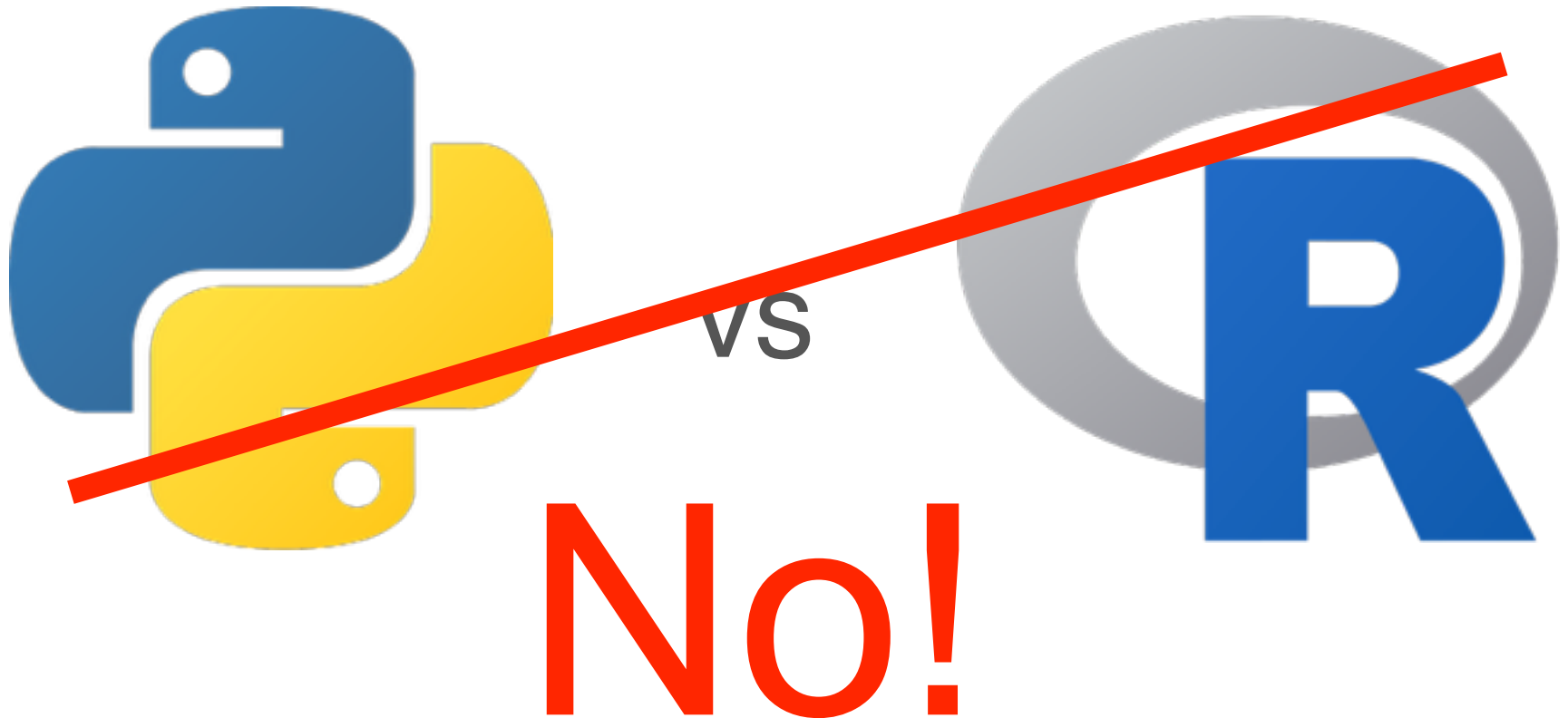**CIVIS**
ANALYTICS

*Building a Data-Driven World[TM]*

VS

vs

No!

*Meet users where they are*

## Prior Knowledge

R is popular in some fields, Python in others. Diverse teams are often polyglot.

## Availability of Key Packages

Important packages are often available in only one language. ***NLTK*** in Python, ***glmnet*** in R. This means a data science workflow often needs to use multiple languages.

## Tradeoffs

Different languages optimize for different things. Python is a general purpose language, R is optimized for statistics/ manipulation of tabular data, Go is a great fit for network services.

Civis Analytics

# Some tools are already cross-language

How?

# Two Options

## Native/Compiled Extensions (C/C++)

**Pros**
- fast!
- many languages speak C

**Cons**
- takes more code
- difficult

**Examples**
- Stan
- XGBoost

## RPC over TCP/HTTP or IPC

**Pros**
- every language speaks TCP/HTTP
- easy to "wire up" host language

**Cons**
- cost of communication

**Examples**
- Spark
- H2o

Civis Analytics

# Two Options

**Native/Compiled Extensions (C/C++)**

**RPC over TCP/HTTP or IPC**

**Pros**

- fast!
- many languages speak C

**Cons**

- takes more code
- difficult

**Examples**

- Stan
- XGBoost

**Pros**

- every language speaks TCP/HTTP
- easy to "wire up" host language
- cost of communication
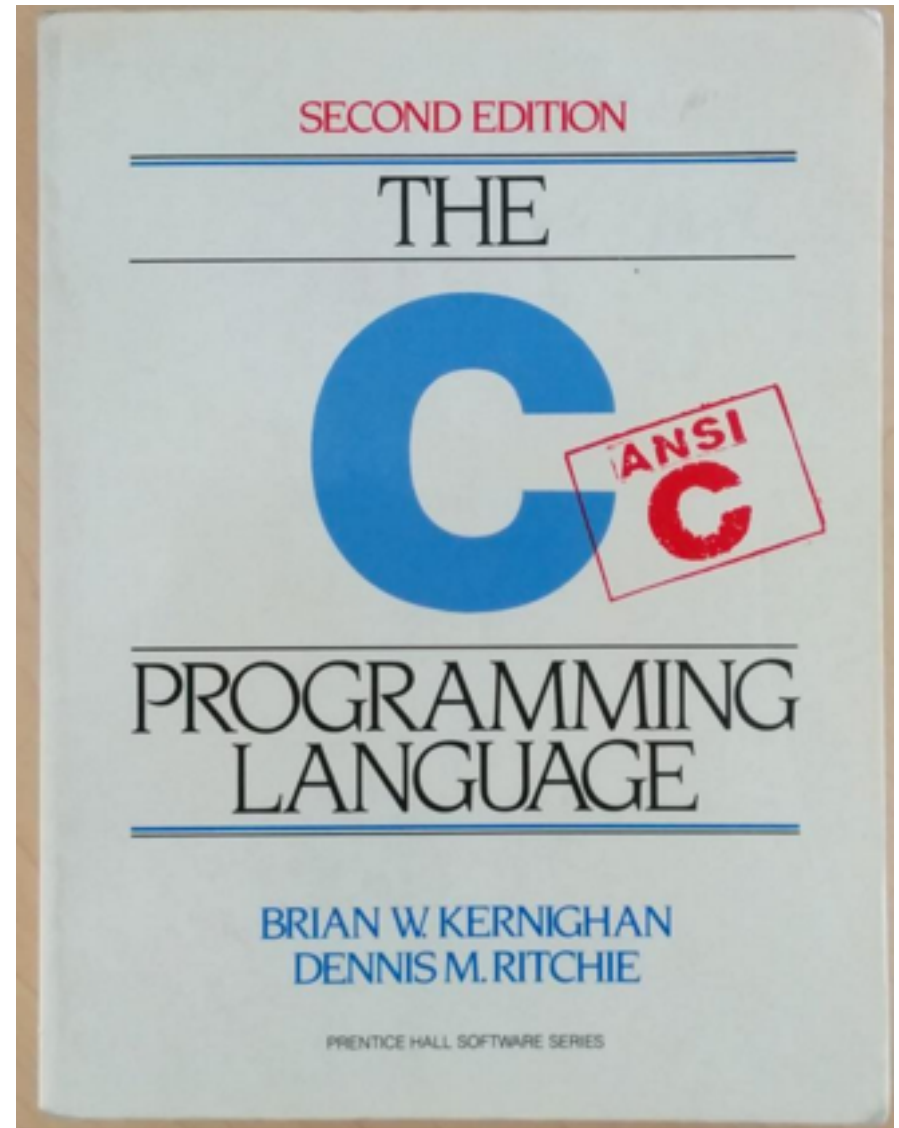
**Examples**

- Spark
- H2o

Our focus for today.

# Why C

- **Python and R "speak" C**

- **Fast!**

- **Portable (mostly)**

- **Simple**

# C++: The Good Parts



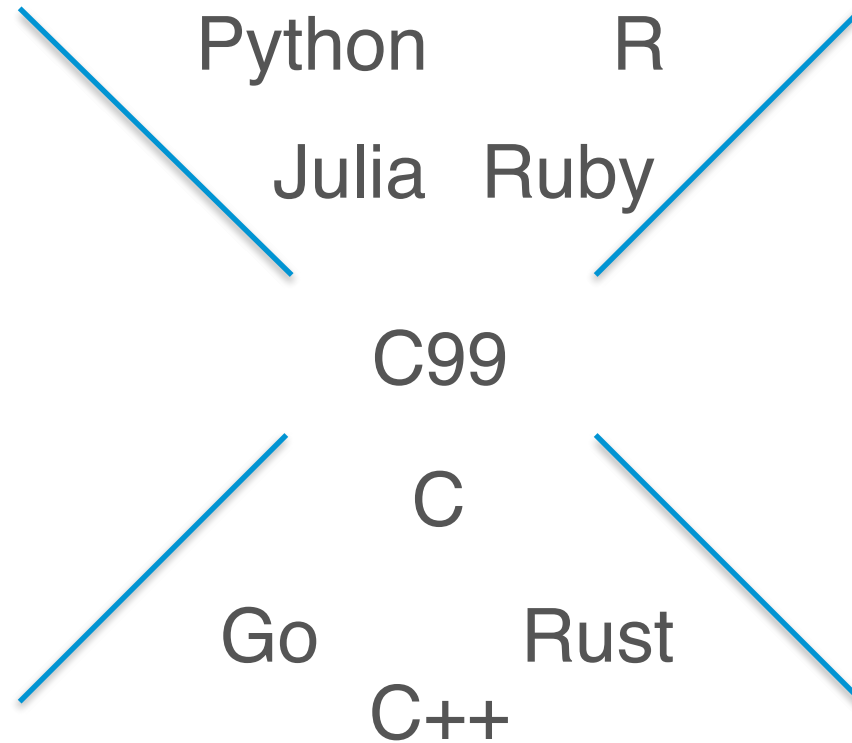Jonathan Adamczewski @twoscomplement · 24 Jul 2015
C++: The Good Parts

13 1.9K    ♥ 1.9K    •••

# Modern C

- **tooling has come a long way**

- **various "sanitizers"**

  - **address/memory sanitizer**

  - **undefined behavior sanitizer**

  - **leak sanitizer**

  - **thread sanitizer**

- **clang gives much better error messages**

# Alternatives: The Hourglass Interface

Python　　　R

Julia　Ruby

C99

C

Go　　　Rust
C++

Credit: Hourglass Interfaces for C++ APIs, Stefanus Du Toit

# Alternatives: The Hourglass Interface

**host language**

Python    R

Julia  Ruby

C99

C

Go    Rust

C++

Credit: Hourglass Interfaces for C++ APIs, Stefanus Du Toit

# Alternatives: The Hourglass Interface

Python     R

Julia  Ruby

C99          **public api**

C

Go     Rust

C++

Credit: Hourglass Interfaces for C++ APIs, Stefanus Du Toit

# Alternatives: The Hourglass Interface

Python       R

Julia  Ruby

C99

**implementation language**          C

Go          Rust

C++

Credit: Hourglass Interfaces for C++ APIs, Stefanus Du Toit

# The Mighty Summation Function

```python
1 def tally(s):
2     total = 0
3     for elm in s:
4         total += elm
5     return total
```

Note: It's best to start development in a language like python.

# Smoke Test

```
In [1]: tally([1, 2, 3])
Out[1]: 6
```

# C Implementation

```c
1  #include <stddef.h>
2
3  double tally(double *s, size_t n) {
4      double total = 0;
5      for (size_t i = 0; i < n; i++) {
6          total += s[i];
7      }
8      return total;
9  }
```

# C Implementation

```
1  #include <stddef.h>
2
3  double tally(double *s, size_t n) {
4      double total = 0;
5      for (size_t i = 0; i < n; i++) {
6          total += s[i];
7          return total;
8      }
9  }
```

need to pass the length

# C/C++ and Python

- **Cython**

- **CFFI**

- **ctypes**

- **C (via the Python C API)**

# The Python C API

```c
1  #include <stdio.h>
2  #include "Python.h"
3  #include "tally.h"
4
5  static PyObject *tally_(PyObject *self, PyObject *args) {
6      // decode/cast the args
7      // call our C function tally
8      // build the result
9  }
10
11 // module method table
12 static PyMethodDef MethodTable[] = {
13     // ...
14 };
15
16 // module def
17 static struct PyModuleDef tally_module = {
18     // ...
19 };
20
21 // module init
22 PyMODINIT_FUNC PyInit_tally_py(void) {
23     return PyModule_Create(&tally_module);
24 }
```

# The Python C API: Buffer API

```c
static PyObject *tally_(PyObject *self, PyObject *args) {
    PyObject *buf;
    if (!PyArg_ParseTuple(args, "O", &buf)) {
        return NULL;
    }

    Py_buffer view;
    int buf_flags = PyBUF_ANY_CONTIGUOUS | PyBUF_FORMAT;
    if (PyObject_GetBuffer(buf, &view, buf_flags) == -1) {
        return NULL;
    }

    if (strcmp(view.format,"d") != 0) {
        PyErr_SetString(PyExc_TypeError, "we only take floats :(");
        PyBuffer_Release(&view);
        return NULL;
    }

    double result = tally(view.buf, view.shape[0]);
    PyBuffer_Release(&view);
    return Py_BuildValue("d", result);
}
```

# The Python C API: Buffer API

```
1  static PyObject *tally_(PyObject *self, PyObject *args) {
2      PyObject *buf;
3      if (!PyArg_ParseTuple(args, "O", &buf)) {
4          return NULL;
5      }
6
7      Py_buffer view;
8      int buf_flags = PyBUF_ANY_CONTIGUOUS | PyBUF_FORMAT;
9      if (PyObject_GetBuffer(buf, &view, buf_flags) == -1) {
10         return NULL;
11     }
12
13     if (strcmp(view.format,"d") != 0) {
14         PyErr_SetString(PyExc_TypeError, "we only take floats :(");
15         PyBuffer_Release(&view);
16         return NULL;
17     }
18
19     double result = tally(view.buf, view.shape[0]);
20     PyBuffer_Release(&view);
21     return Py_BuildValue("d", result);
22 }
```

# The Python C API: Buffer API

```c
1  static PyObject *tally_(PyObject *self, PyObject *args) {
2      PyObject *buf;
3      if (!PyArg_ParseTuple(args, "O", &buf)) {
4          return NULL;
5      }
6
7      Py_buffer view;
8      int buf_flags = PyBUF_ANY_CONTIGUOUS | PyBUF_FORMAT;
9      if (PyObject_GetBuffer(buf, &view, buf_flags) == -1) {
10         return NULL;
11     }
12
13     if (strcmp(view.format,"d") != 0) {
14         PyErr_SetString(PyExc_TypeError, "we only take floats :(");
15         PyBuffer_Release(&view);
16         return NULL;
17     }
18
19     double result = tally(view.buf, view.shape[0]);
20     PyBuffer_Release(&view);
21     return Py_BuildValue("d", result);
22 }
```

# The Python C API: Buffer API

```c
static PyObject *tally_(PyObject *self, PyObject *args) {
    PyObject *buf;
    if (!PyArg_ParseTuple(args, "O", &buf)) {
        return NULL;
    }

    Py_buffer view;
    int buf_flags = PyBUF_ANY_CONTIGUOUS | PyBUF_FORMAT;
    if (PyObject_GetBuffer(buf, &view, buf_flags) == -1) {
        return NULL;
    }

    if (strcmp(view.format,"d") != 0) {
        PyErr_SetString(PyExc_TypeError, "we only take floats :(");
        PyBuffer_Release(&view);
        return NULL;
    }

    double result = tally(view.buf, view.shape[0]);
    PyBuffer_Release(&view);
    return Py_BuildValue("d", result);
}
```

# The Python C API

```
1  static PyObject *tally_(PyObject *self, PyObject *args) {
2      PyObject *buf;
3      if (!PyArg_ParseTuple(args, "O", &buf)) {
4          return NULL;
5      }
6
7      Py_buffer view;
8      int buf_flags = PyBUF_ANY_CONTIGUOUS | PyBUF_FORMAT;
9      if (PyObject_GetBuffer(buf, &view, buf_flags) == -1) {
10         return NULL;
11     }
12
13     if (strcmp(view.format,"d") != 0) {
14         PyErr_SetString(PyExc_TypeError, "we only take floats :(");
15         PyBuffer_Release(&view);
16         return NULL;
17     }
18
19     double result = tally(view.buf, view.shape[0]);
20     PyBuffer_Release(&view);
21     return Py_BuildValue("d", result);
22 }
```

# The Python C API

```
1  static PyObject *tally_(PyObject *self, PyObject *args) {
2      PyObject *buf;
3      if (!PyArg_ParseTuple(args, "O", &buf)) {
4          return NULL;
5      }
6
7      Py_buffer view;
8      int buf_flags = PyBUF_ANY_CONTIGUOUS | PyBUF_FORMAT;
9      if (PyObject_GetBuffer(buf, &view, buf_flags) == -1) {
10         return NULL;
11     }
12
13     if (strcmp(view.format,"d") != 0) {
14         PyErr_SetString(PyExc_TypeError, "we only take floats :(");
15         PyBuffer_Release(&view);
16         return NULL;
17     }
18
19     double result = tally(view.buf, view.shape[0]);
20     PyBuffer_Release(&view);
21     return Py_BuildValue("d", result);
22 }
```

# The Python C API: Method Table

```c
1 static PyMethodDef MethodTable[] = {
2     {"tally", &tally_, METH_VARARGS, "Compute the sum of an array."},
3     { NULL, NULL, 0, NULL}
4 };
5
6 static struct PyModuleDef tally_module = {
7     .m_base = PyModuleDef_HEAD_INIT,
8     .m_name = "tally_py",
9     .m_size = -1,
10     .m_methods = MethodTable
11 };
12
13 PyMODINIT_FUNC PyInit_tally_py(void) {
14     return PyModule_Create(&tally_module);
15 }
```

# C/C++ and R

- **Rcpp**
- **C (via the R C API)**

# The R C API

```c
 1 #include <R.h>
 2 #include <Rinternals.h>
 3 #include <R_ext/Rdynload.h>
 4 #include "tally.h"
 5
 6 SEXP tally_(SEXP x_) {
 7     // cast/decode the input
 8     // call our tally function
 9     // build the output
10 }
11
12 // method table
13 static R_CallMethodDef callMethods[] = {
14     // ...
15 };
16
17 // module/package init
18 void R_init_tally_r(DllInfo *info) {
19   R_registerRoutines(info, NULL, callMethods, NULL, NULL);
20 }
```

Civis Analytics

# The R C API

```
 1  SEXP tally_(SEXP x_) {
 2    double *x = REAL(x_);
 3    int n = length(x_);
 4
 5    SEXP out = PROTECT(allocVector(REALSXP, 1));
 6    REAL(out)[0] = tally(x, n);
 7    UNPROTECT(1);
 8
 9    return out;
10  }
```

# The R C API

```
 1  SEXP tally_(SEXP x_) {
 2      double *x = REAL(x_);
 3      int n = length(x_);
 4
 5      SEXP out = PROTECT(allocVector(REALSXP, 1));
 6      REAL(out)[0] = tally(x, n);
 7      UNPROTECT(1);
 8
 9      return out;
10  }
```

# The R C API

```
 1 SEXP tally_(SEXP x_) {
 2   double *x = REAL(x_);
 3   int n = length(x_);
 4
 5   SEXP out = PROTECT(allocVector(REALSXP, 1));
 6   REAL(out)[0] = tally(x, n);
 7   UNPROTECT(1);
 8
 9   return out;
10 }
```

# The R C API

```
 1 SEXP tally_(SEXP x_) {
 2   double *x = REAL(x_);
 3   int n = length(x_);
 4
 5   SEXP out = PROTECT(allocVector(REALSXP, 1));
 6   REAL(out)[0] = tally(x, n);
 7   UNPROTECT(1);
 8
 9   return out;
10 }
```

# The R C API: Function Registration

```
1 static R_CallMethodDef callMethods[] = {
2   {"tally_", (DL_FUNC)&tally_, 1},
3   {NULL, NULL, 0}
4 };
5
6 void R_init_tally_r(DllInfo *info) {
7   R_registerRoutines(info, NULL, callMethods, NULL, NULL);
8 }
```

# Miscellaneous

**Dependencies**

Don't depend on APIs from host languages, i.e., numpy, rmath

**Errors**

Use error codes to signal problems. Don't call abort or exit as these will quit the process running the host language.

**Memory**

Typically best to make the host language responsible for allocation and deallocation. It's challenging to transfer ownership over the boarder.

**Logging/Verbosity**

At the very least, make this optional.

**Compiler**

Trust the compiler it's smarter than all of us. Ensure your code compiles without warnings.

# Parting Thoughts

1. Meet users where they are
2. Reach a larger audience
3. Make a bigger impact

Thank You

Bill Lattner
**twitter**: @wlattner
**github**: github.com/wlattner