# A Critique of *R*
## from the Perspective of Programming Language Theory

Sungwoo Park

Pohang University of Science and Technology, Korea

Japanese *R* Use*R*'s meeting

Dec 8, 2006

# First Encounter with R

- A regional workshop on R in May, 2006
  - motto:

    ***"Don't teach SAS. Teach R instead."***

- An invited talk at the workshop
  - supposed to say ***"SAS is bad. R is good."***
  - actually said ***"SAS is really bad, R is also bad."***

- R seemed to have quite a few flaws in its design.

# 'Towards 2020 Science'

- A report on

  *"the role and future of science over the next 14 years"*

  - by the 2020 Science Group
    - over 30 scientists elected for their expertise
    - met over an intense 3 days in July 2005
  - 86 pages
  - sponsored by Microsoft

# Towards 2020 Science: A Draft Roadmap

# R to Be Reckoned With

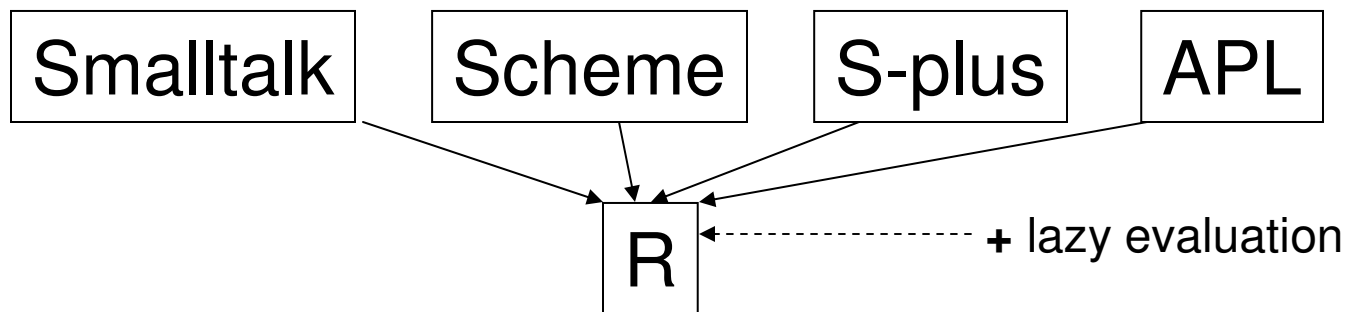| Issues | Emergence of communities that build / share scientific / statistical computing tools (e.g. 'R' statistical computing and graphics project / community) |
|---|---|
| **2005** | |

*"Many niche areas of software development exist where alternatives and/or enhancements of managed platforms are deployed and used by scientists, including ... and __the language R__."*

# R Dissected

- Popularity of R in the statistics community
  - statistical computing
  - high level graphics

    *"Many users will come to R mainly for its graphical facilities." – An introduction to R*

- R as a hybrid language

# Caveat

- A technical debate on
  *"Why is your programming language good/bad?"*
  $\approx$

  A religious debate on
  *"What is the best religion?"*

  $\Rightarrow$ **Take this presentation with a grain of salt.**

- One thing is certain, however:
  *"More features do not always mean
  a better programming language."*

# Outline

- Introduction *V*
- **Programming paradigm for R**
  - **Imperative language?**
  - **Functional language?**
  - **Both?**
  - **Or neither?**
- Lexical scoping
- Further analysis
- A functional language for R users
- Conclusion

# Imperative vs Functional

| Imperative languages | Functional languages |
|---|---|
| • Everything denotes a command. <br> • Variables are mutable. <br> • Functions are not first-class objects. | • Everything denotes a value. <br> • Variables are immutable. <br> • **Functions as first-class objects.** |

Functions are first-class objects in R.
Does this mean that R is a functional language?

# Imperative Languages

- A program consists of **commands**.
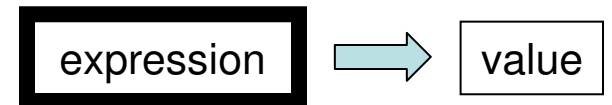  - command = "do something"

- Nothing wrong:

  **<u>if</u> (x == 1)**

      **x ← x + 1**

  **<u>else</u>**

      **x ← x - 1**

- Nothing wrong either:

  **<u>if</u> (x == 1)**

      **x ← x + 1**

# Functional Languages

- A program consists of **expressions**.
  - expression = "obtain a **value**"

| expression | ⟹ | value |
|---|---|---|

- Nothing wrong:

    <u>if</u> **(x == 1)**

    **x + 1**

    <u>else</u>

    **x - 1**

<u>if</u> (1 == -1) 10 <u>else</u> -10

⟹  <u>if</u> (false) 10 <u>else</u> -10

⟹  -10

- But this does not make sense:

    <u>if</u> **(x == 1)**

    **x + 1**

  - What is the value if $x \neq 1$?

# R: Not Functional

```
> foo
function (x) {
    if (x <= 0) 1
  }
> foo (0)
[1] 1
> foo (1)
> foo (0) + foo (1)
numeric(0)
```

**if (x <= 0) 1** is not an expression:
it does not always evaluate to a value.
**foo** is not a function:
it is not defined on positive integers.

**numeric(0)** is what?

# Variable Binding

```
> x = 1 + 1
> x
[1] 2
```

- A variable *x* is "**bound**" to value **2**.
- From now on, any occurrence of **x** is replaced by **2**.

```
> y = x + x
> y
[1] 4
```

# Variables are NOT Variable?

| Imperative languages | Functional languages |
|---|---|
| • The contents of a variable can change.<br><br>   **> x ← 0**<br><br>   **> x**<br><br>   **[1] 0**<br><br>   **> x ← 1**<br><br>   **> x**<br><br>   **[1] 1** | • The contents of a variable **never** change.<br><br>⇒ You **cannot** assign a new value to a variable.<br><br>• Surprise?<br><br>⇒ nothing special in functional languages |

So, R is an imperative language?

# References in Functional Languages

- There are assignments, but not to variables.

  ⇒ assignments to **references.**

- Reference (≈ pointer in C)

  – points to a heap cell.

```
- val x = ref 0;          // initialization
val x = ref 0 : int ref
- !x;                     // dereferencing
val it = 0 : int
- x := 1;                 // assignment
val it = () : unit
- !x;                     // dereferencing
val it = 1 : int
```



15

# R: Neither Functional Nor Imperative

| Imperative languages | Functional languages |
|---|---|
| • Everything denotes a command. | • Everything denotes a value. |
| • **Variables are mutable.** | • Variables are immutable. |
| • Functions are not first-class objects. | • **Functions as first-class objects** |

- Functions are first-class objects, but
  no clear definition of **commands** or **expressions**
  no distinction between **variables** and **references**
- A fatal design decision

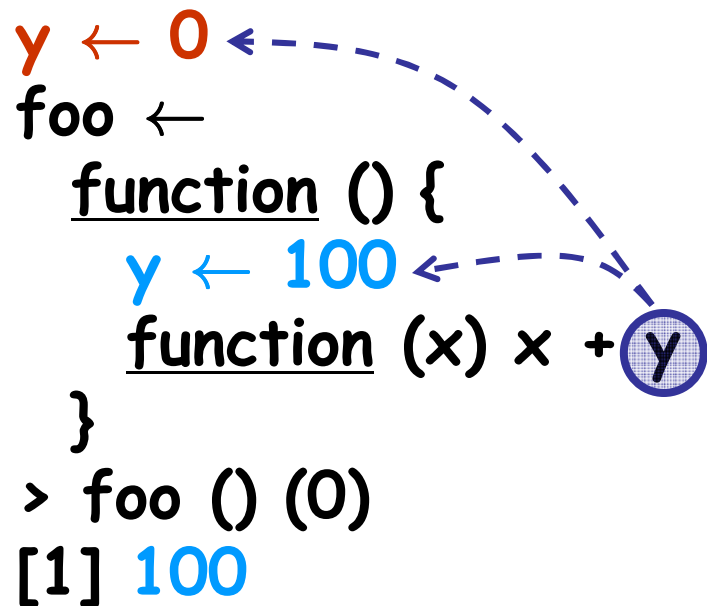  ⇒ engenders many idiosyncrasies in the definition.

# Outline

- Introduction *V*
- Programming paradigm for R *V*
- **Lexical scoping**
- Further analysis
- A functional language for R users
- Conclusion

# Lexical Scoping

- Uses bindings that are active at the time of creating a function.

```
y ← 0
foo ←
  function () {
    y ← 100
    function (x) x + y
  }
> foo () (0)
[1] 100
```
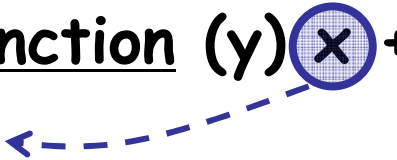
- Useful in R because functions are first-class objects.
- Unfortunately R fails to implement lexical scoping correctly.

# No Lexical Scoping

- R

```
> x ← 1
> foo ← function (y) x + y
> x ← 100
> foo (0)
[1] 100
```
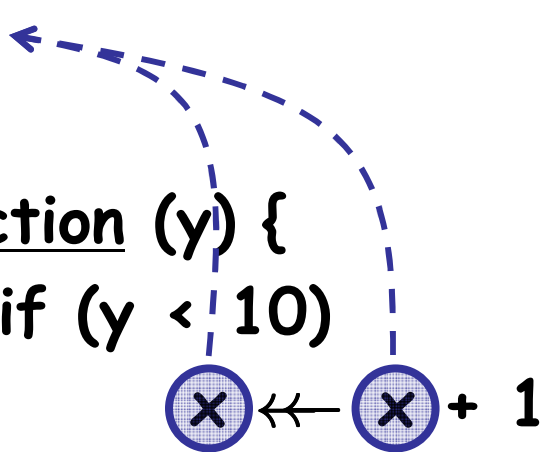
- Standard ML

```
val x = 1
val foo = fn y => x + y
val x = 100
- foo 0;
val it = 1 : int
```
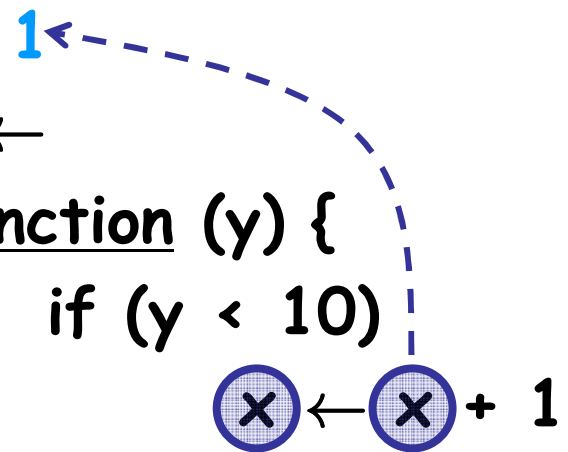
- Dynamic scoping at the top-level
- lexical scoping at inner levels
  - $\Rightarrow$ for the sake of compatibility?

# ↞ VS ←

```
x ← 1
foo ←
    function (y) {
        if (y < 10)
            x ↞ x + 1
    }
> foo (0)
> x
[1] 2
```

```
x ← 1
foo ←
    function (y) {
        if (y < 10)
            x ← x + 1
    }
> foo (0)
> x
[1] 1
```

# Special Top Level?

*"While purely functional languages do not allow assignment,* **_they allow it at top-level_**; *otherwise the user could not define new functions."*

*— Lexical Scope and Statistical Computing*

$\Rightarrow$ Wrong!

- There is nothing special for the top level.

- assignment at the top level?

  - No, it's just a **binding**.

- Due to failure to distinguish between
  **variables** and **references**, or
  **bindings** and **assignments**.

# Lexical Scoping in CS

*"Although the usual definition of static or lexical scope in computer science is that ..., **this definition is not specific enough.** Computer scientists tend not to differentiate as finely because **their concerns are different."***
       *– Lexical Scope and Statistical Computing*

$\Rightarrow$ This is absolutely wrong.

$$\begin{array}{rcl} \text{expression} & e & ::= \ x \mid \lambda x : A.\, e \mid e\ e \\ \text{environment} & \eta & ::= \ \cdot \mid \eta, x \hookrightarrow v \end{array}$$

$$\dfrac{x \hookrightarrow v \in \eta}{\eta \vdash x \hookrightarrow v}\ \mathbf{Var_e} \qquad \dfrac{}{\eta \vdash \lambda x : A.\, e \hookrightarrow [\eta, \lambda x : A.\, e]}\ \mathbf{Lam_e}$$

$$\dfrac{\eta \vdash e_1 \hookrightarrow [\eta', \lambda x : A.\, e] \quad \eta \vdash e_2 \hookrightarrow v_2 \quad \eta', x \hookrightarrow v_2 \vdash e \hookrightarrow v}{\eta \vdash e_1\ e_2 \hookrightarrow v}\ \mathbf{App_e}$$

# Outline

- Introduction $\checkmark$
- Programming paradigm for R $\checkmark$
- Lexical scoping $\checkmark$
- **Further analysis**
- A functional language for R users
- Conclusion

# R TYPE

- Dynamic type binding
  - An R object can change its type during the computation.

$$x \leftarrow c(1.0, 2.0, 3.0)$$
$$x \leftarrow 47$$

  - **typeof** returns the type of an R object.
    - **symbol**, **pairlist**, **closure**, **environment**

- Is it good? $\Rightarrow$ philosophical debate
  - dynamic type binding is good for:
    - quick, small programming tasks
  - static type binding is good for:
    - large programming tasks

# Complex Semantics

- Ex. *Section 3.4 Indexing* in *R Language Definition*

```
x[i]
x[i, j]
x[[i]]
x[[i, j]]
x$a
x$"a"
```

- Why on earth such "a" complex semantics for statistical computing?

# So Many Complex/Special Cases

- From *R Language Definition*
  - *"Another more <u>subtle</u> difference is ..."*
  - *"... evaluated in some <u>unexpected</u> cases."*
  - *"... can lead to <u>surprises</u>."*
  - *"In a <u>very few</u> cases, ..."*
  - *"... in certain (<u>rather rare</u>) circumstances, ..."*
  - *"... are treated <u>specially</u>."*
  - *"... should be done <u>with caution</u>."*
  - *"A couple of <u>special</u> rules apply, though:"*
  - *"... is <u>not guaranteed</u> to hold in all implementations."*
  - *"is <u>not generally</u> handled correctly."*
  - *"The <u>special exception</u> for ... is admittedly <u>peculiar</u>."*

# Evolution or Degeneration?

"R appears to be working fine."

"??? seems often useful, so let's add it to R."

"Now ??? is available, but there is something fishy going on."

- Example of ??? = first-class functions

  *"This ability is rarely used even though it is **potentially very powerful**." – Lexical Scope and Statistical Computing*

  – incorporating first-class functions

  without expressions and bindings

  ⇒ fitting a square peg into a round hole

- The worst example of ??? is yet to come, however.

# Lazy Evaluation

*"A policy of lazy arguments is **<u>very useful</u>** because ... This **<u>can be very useful</u>** for specifying functions or models in symbolic form."*
— *R: A Language for Data Analysis and Graphics*

- Evaluation strategy of R
  - eager evaluation for built-in functions: fully evaluate arguments
  - lazy evaluation for *promise* objects: evaluate only when necessary

- Yes, lazy evaluation is a great idea.
  - Ex. Haskell
  - **<u>But only if all functions are pure mathematical functions</u>**.

- Lazy evaluation + computational effects $\Rightarrow$ total complete mess
  - computational effects (= side effects)
    - **plot**, **print**, vector update, assignments
  - Functions in R are not mathematical functions anyway.
  - Solution from programming language theory = monad

- Besides lazy evaluation in R is **<u>not</u>** really lazy evaluation!

# Meta-programming in R

- **quote** creates unevaluated expressions.
- **eval** treats programs as data.

  ```
  > e <- quote (2 + 2)
  > v <- eval (e)
  ```

- Useful constructs? Yes!
  - implementing compilers, staged computation, and so on

- But do you really need **quote**, **eval**, **deparse**, **substitute** for statistical computing?

*"More frequently, one wants to ... in order to deparse it and use it for **labeling plots**, for instance." – R Language definition*

  $\Rightarrow$ launching a nuclear missile to kill a fly

# Why Not Use First-Class Functions?

- A weird program exploiting lazy evaluation and **eval**

```
curve ← function (expr, from, to) {
    x ← seq (from, to, length=500)
    y ← eval (substitute (expr))
    plot(x, y, type="l")
}
curve (x^2 - 1, -2, 2)
```

This function call does not make sense.
⇒ misunderstanding of lazy evaluation!

- A quick fix = use a first-class function

```
curve ← function (f, from, to) …
curve (function (x) x^2- 1, -2, 2)
```

# Other Minor (Yet Serious) Points

- Maintaining state within functions

  *"The ability to preserve state information between function invocations is a **very useful** feature ..."*
  $\qquad\qquad$ – *R: A Language for Data Analysis and Graphics*

  $\Rightarrow$ a trivial exercise in functional programming

- Confusion between definition and implementation
  *"To understand completely the rules ..., the reader needs to be familiar with the notion of an evaluation frame."*
  $\qquad\qquad\qquad\qquad\qquad$ – *An Introduction to R*

  – Specific implementation strategies are taken as part of the definition.

    - environment, closure, call stack, evaluation frame, ...

# Outline

- Introduction *✔*
- Programming paradigm for R *✔*
- Lexical scoping *✔*
- Further analysis *✔*
- **A functional language for R users**
- Conclusion

# Next Generation R?

- Claim

  1. Admit it or not, R is an ill-designed language.
  2. Nevertheless, R is too juicy to give up:
     - statistical computing
     - high level graphics
  3. R shares a lot in common with functional languages.

- Plan

  – extend an existing functional language with an interface to the R base library.

# Objective CAML with R

- Objective CAML
  - industrial strength functional language
  - rough speed comparison
    - nearly as fast as, or sometimes faster than, C
    - consistently faster than C++
    - about 10 times faster than Matlab
  - strong type system (based on type theory)
    - **<u>significantly</u>** less development time than in C
    - more reliable code than in C
  - huge library contributed by users
  - free!
- Let's develop an Objective CAML interface to R!

# Preliminary Results

# Outline

- Introduction *✔*
- Programming paradigm for R *✔*
- Lexical scoping *✔*
- Further analysis *✔*
- A functional language for R users *✔*
- **Conclusion**

# Summary

- R is great!
  - library for statistical computing
  - library for publication quality graphics
  - the whole statistics community actively contributing new libraries

- R is an ill-designed language, however.

- So, it's time to act.
  - just use programming language theory!

# Thanks a lot!