

Solving the n -language problem: A ecologist's guide to learning Julia

Michael D. Catchen^{1,2} [Timothée Poisot](#)^{3,2}

¹ McGill University ² Québec Centre for Biodiversity Sciences ³ Université de Montréal

Correspondance to:

Michael D. Catchen — michael.catchen@mcgill.ca

This work is released by its authors under a CC-BY 4.0 license



Last revision: *September 10, 2022*

Julia is a good language, ecologists should learn it.

1 Outline

2 An ecologists guide to learning Julia

3 • Why should ecologists learn Julia?

- 4 – No point ignoring the elephant in the room: R is by far the most used language in ecology
- 5 – But many of the problems in contemporary ecology require functionality that R just can't
- 6 handle (mass individual-based models, big data management, quick numerical methods)
- 7 – `julia` has particular features that solve a broad chunk of the difficulties of computational
- 8 analysis of ecological data
 - 9 * and improves user experience by enabling code that is *idiomatic*.
- 10 – this makes `julia` well suited to become "[a platform for ecological prediction and forecasting];
- 11 McIntyre et al, Urban et al

12 • In what order should people approach new topics

- 13 – understand the core concepts that make `julia` different from other languages first
 - 14 * Types
 - 15 * multiple dispatch
- 16 – Learn data management in Julia (emphasize analogies to tidyverse)
 - 17 * DataFrames
- 18 – Learn stats/ml in Julia (emphasize analogies to GLM in R and scikit-learn/tensorflow in
- 19 python)
- 20 – Learn visualization (Makie)

21 • Discussion

- 22 – Talk about adjacent scientific computing course
- 23 – What does the future of computing in ecology look like?

-

Why
should
ecologists
learn
julia?

-

Well,
there
are
the
criteria
that
are
directly
measurable
that
make
it
better
than

R/Python:

- fast

-

native
support
on

GPUs

- But

there

are

also

criteria

Abstract

Introduction

In order to measure, understand, and mitigate the consequences of anthropogenic change on ecosystems the services they provide, ecologists need a modern set of computational tools (Urban *et al.* 2022).

These tools must be performant, but crucially modular and interfaceable (McIntire *et al.* 2022)

The so-called “two-language problem” in computational science, where it is easier for a researcher to develop a prototype of a model/simulation in a high-level language, like Python or R, and later have to port the model to a lower-level compiled language because the performance of these compiled languages (e.g. C++/Fortran) is orders of magnitude faster than high level interpreted languages. In fact, many of the most popular tools in higher-level languages are actually thin wrappers around a compiled (often C++) base (e.g. tidyverse, keras, numpy, TensorFlow, scikit-learn, pandas, etc.). However, the skills required to use or debug—let alone write—scientific software in these lower level languages is not often taught.

Ecological data is often difficult to access and reuse (Gonzalez & Peres-Neto 2015; Poisot *et al.* 2019). Synthesizing data into a single product suitable for analysis often remains tedious as data are not in formats that can be easily combined or interfaced.

Here we propose that we can solve this problem through standardization (Zimmerman 2008)—developing a common definition such that data collected in a variety of contexts can be assimilated while minimizing the overhead of data cleaning and wrangling. A common representation of ecological data will have three primary benefits: it will **1**) enable new forms of analysis by making it easier to combine data from different sources (Heberling *et al.* 2021), **2**) enable continuous integration of new data for next-generation biodiversity monitoring (Kühl *et al.* 2020), and **3**) aid in open sharing and reproducibility of published results (Zimmerman 2008; Borregaard & Hart 2016). Here, we briefly review approaches to data standardization developed in other fields, in order to determine what makes an open standard succeed in

48 promoting data sharing, and what doesn't. Based on the properties of good standards we identify, we
49 propose building a living standard for ecological data in the Julia programming language, and argue this
50 is necessary to obtain the three primary benefits of standardization mentioned earlier.

51 We propose that that Julia has certain properties absent in other popular languages for scientific
52 computing that make it particularly suited for the development of a cohesive, modular, and extendible set
53 of tools ideal for the development of a platform for ecological analysis [4].

54 **The nature of computation in Julia**

55 [Figure 1 about here.]

56 **Types**

57 Why is this useful for ecologists? Often times in ecology, the same information is represented in different
58 formats. Two packages in R might not agree on what the "correct" format to represent information is.

59 At the core of the Julia language is its *type system*. Type systems can often be alienating to those who
60 learned programming in so-called *dynamically* typed languages (like R, python, and JavaScript). In
61 dynamically-typed languages, `x = 5`, and `x = "hello world"` and the language won't care that you
62 changed the type of information that was stored in `x` from a number to a string. Practically, this form of
63 dynamic-typing was adopted because it is far more convenient to write code like that above than defining
64 variables with explicit types, e.g. how you would in C: `char c = "a";` and `int x = 5;`.

65 Julia doesn't require explicit type declarations, meaning `x = 5` is perfectly valid code, but internally Julia is
66 doing the bookkeeping of what type of information is stored in `x`, from an `Int64` to a `String` in the above
67 example.

68 Using explicit types is central to Julia's speed, but also enables much of its most unique and user-friendly
69 functionality, primarily the use of a *multiple-dispatch* system.

70 **Dispatch**

71 *Dispatch* refers to the way a computer program decides what function to call.

72 In many statically-typed languages, you are allow to use the same function name more than once.

73 **Doing computational science in Julia**

74 **Managing Data**

75 DataFrames.jl and DFMeta.

76 **Doing statistics and machine learning**

77 7. Learn about the statistics ecosystem: StatsBase, Statistics, GLM, MLJ, Flux, Turing

78 **Doing simulation**

79 8. Learn about the simulation libraries (DiffEq, DynamicGrids)

80 9. Learn how various statistics/simulation libraries work together

81 **Discussion**

82 Defining a living standard for ecological data in Julia will make it easier to combine data from different
83 sources by splitting the process of data aggregation from the process of analysis. Integrating data from a
84 particular study, or a new database, would be as simple as implementing the interface from the data
85 source to the standardized types. Data from individual studies could be incorporated into public
86 repositories containing both the raw data and the interface to Julia data structures, and this combined
87 data/interface package is all that is needed to either reproduce the results or incorporate that particular
88 study's data into analysis. This will make combining data from multiple sources easier, and yield benefits
89 for the development and implementation of novel methods, as the software for analysis becomes separate
90 from the software for data cleaning and aggregation.

91 We envision a modern set of tools for ecology in Julia based around the standardized types. Far outside of
92 ecology, the term “ecosystem” is used metaphorically to describe a set of software tools that work together.
93 We imagine multiple “trophic-levels” of packages for ecological science in Julia based around the “basal”

94 set of standardized types — a modular set of tools that can be chained together create arbitrarily complex
 95 analysis pipelines. that can be scaled to meet the needs of next-generation biodiversity monitoring.

96 Borregaard, M.K. & Hart, E.M. (2016). Towards a more reproducible ecology. *Ecography*, 39, 349–353.

97 Gonzalez, A. & Peres-Neto, P.R. (2015). Act to staunch loss of research data. *Nature*, 520, 436–436.

98 Heberling, J.M., Miller, J.T., Noesgaard, D., Weingart, S.B. & Schigel, D. (2021). Data integration enables
 99 global biodiversity synthesis. *Proceedings of the National Academy of Sciences*, 118.

100 Köhl, H.S., Bowler, D.E., Bösch, L., Bruehlheide, H., Dauber, J., Eichenberg, David., *et al.* (2020). Effective
 101 Biodiversity Monitoring Needs a Culture of Integration. *One Earth*, 3, 462–474.

102 McIntire, E.J.B., Chubaty, A.M., Cumming, S.G., Andison, D., Barros, C., Boisvenue, C., *et al.* (2022).
 103 PERFICT: A Re-imagined foundation for predictive ecology. *Ecology Letters*, 25, 1345–1351.

104 Poisot, T., Bruneau, A., Gonzalez, A., Gravel, D. & Peres-Neto, P. (2019). Ecological Data Should Not Be So
 105 Hard to Find and Reuse. *Trends in Ecology & Evolution*, 34, 494–496.

106 Urban, M.C., Travis, J.M.J., Zurell, D., Thompson, P.L., Synes, N.W., Scarpa, A., *et al.* (2022). Coding for
 107 Life: Designing a Platform for Projecting and Protecting Global Biodiversity. *BioScience*, 72, 91–104.

108 Zimmerman, A.S. (2008). New Knowledge from Old Data: The Role of Standards in the Sharing and
 109 Reuse of Ecological Data. *Science, Technology, & Human Values*, 33, 631–652.

1 Defining the types

```
abstract type Pet end
struct Dog <: Pet
  name
end
struct Cat <: Pet
  name
end
```

2 Defining the methods

```
meet(🐶::Dog, 🐱::Cat) = "$(🐶.name) meets $(🐱.name) and barks"
meet(🐶::Dog, 🐾::Dog) = "$(🐶.name) meets $(🐾.name) and sniffs"
meet(🐱::Cat, 🐶::Dog) = "$(🐱.name) meets $(🐶.name) and hisses"
meet(🐱::Cat, 🐱::Cat) = "$(🐱.name) meets $(🐱.name) and slinks"
```

3 Creating instances of types

```
fido = Dog("Fido")
sparky = Dog("Sparky")

tabby = Cat("Tabby")
panko = Cat("Panko")
```

4 Calling the methods

```
meet(sparky, tabby)
> Sparky meets Tabby and barks
meet(fido, sparky)
> Fido meets Sparky and sniffs
meet(panko, fido)
> Panko meets Sparky and hisses
meet(tabby, panko)
> Tabby meets Panko and slinks
```

Figure 1: TODO: caption. Adapted from Karpinski 2019 “The unreasonable effectiveness of multiple dispatch”