# This is the Way

*Integrating Open Data Science Workflow & Software Carpentry into the Statistical Ecology Classroom*
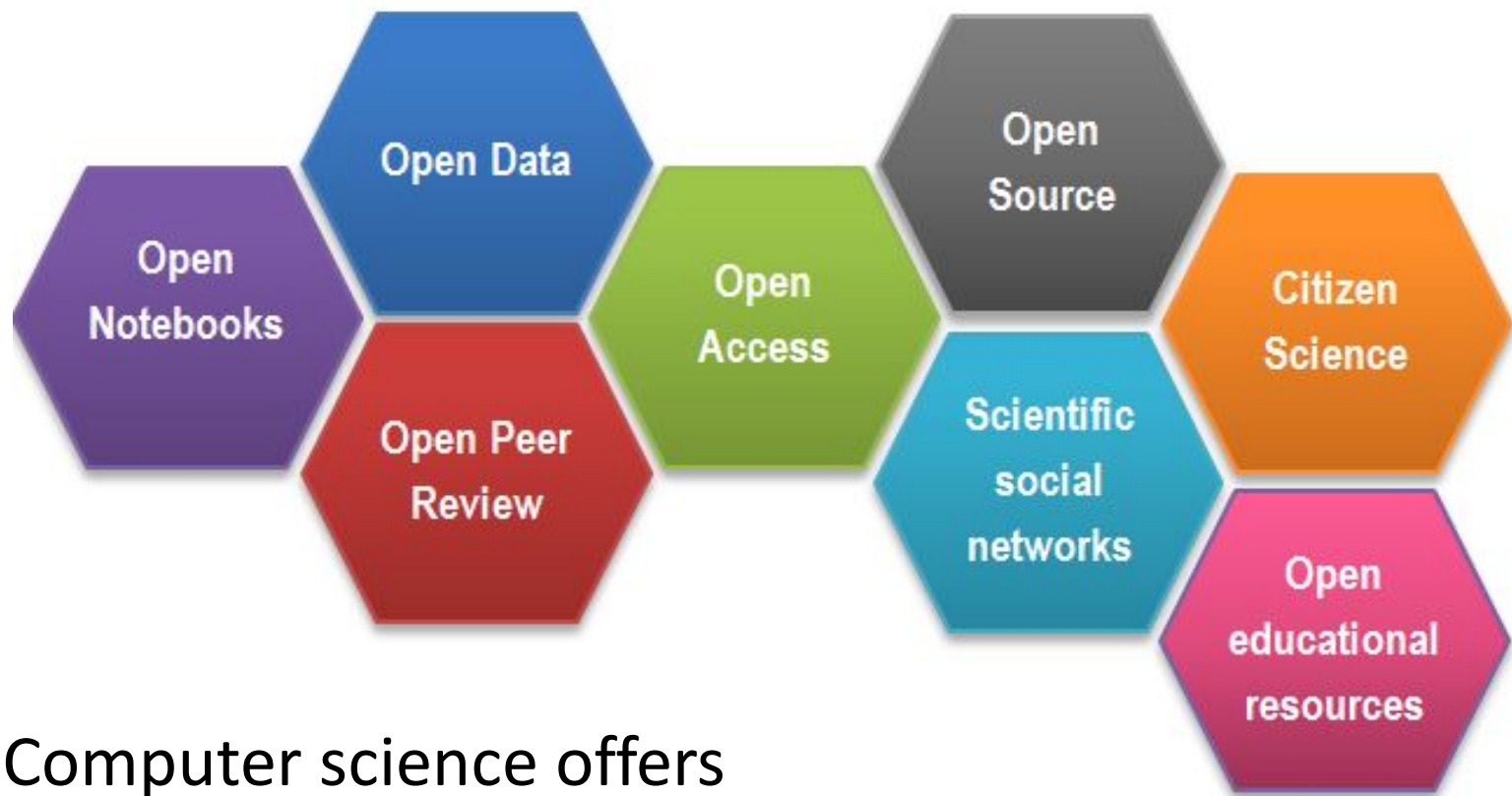
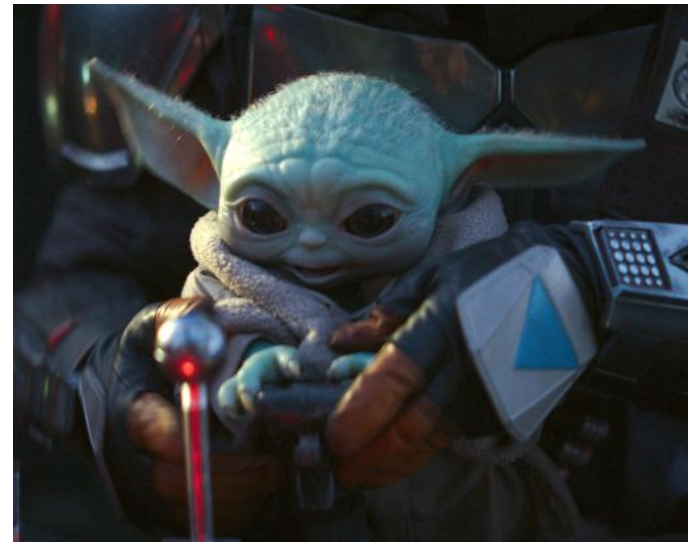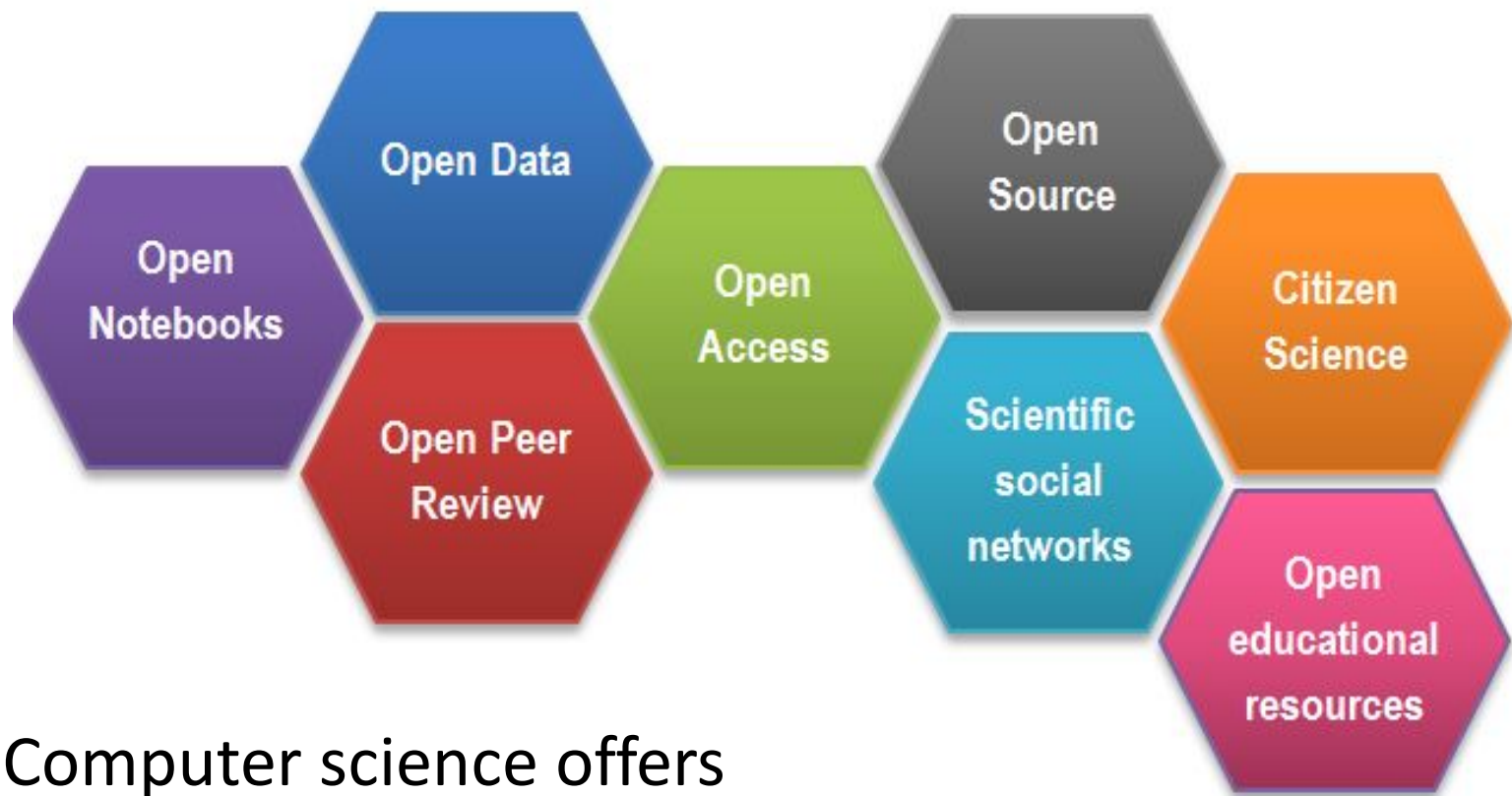

Gavin Fay

gfay@umassd.edu

 @gavin_fay

UMassD Teaching & Learning Conference

17 January 2020

Computer science offers many tools to facilitate shiny Open Science.

Computer science offers many tools to facilitate shiny Open Science.

But, few scientists are trained in their use.

# Motivation #1



art: @allison_horst

@juliesquid

nature
ecology & evolution

**PERSPECTIVE**
PUBLISHED: 23 MAY 2017 | VOLUME: 1 | ARTICLE NUMBER: 0160

## Our path to better science in less time using open data science tools

Julia S. Stewart Lowndes[1]*, Benjamin D. Best[2], Courtney Scarborough[1], Jamie C. Afflerbach[1], Melanie R. Frazier[1], Casey C. O'Hara[1], Ning Jiang[1] and Benjamin S. Halpern[1,3,4]

Source: Lowndes et al. 2017: ohi-science.org/betterscienceinlesstime

# Motivation #2

# Motivation #2

PERSPECTIVE

## Good enough practices in scientific computing

[1]⊚\*, Jennifer Bryan[2]⊚, Karen Cranston[3]⊚, Justin Kitzes[4]⊚, Lex Nederbragt[5]⊚,
[6]⊚

...pentry Foundation, Austin, Texas, United States of America, 2 RStudio and Department of
...ersity of British Columbia, Vancouver, British Columbia, Canada, 3 Department of Biology,
...y, Durham, North Carolina, United States of America, 4 Energy and Resources Group,
...alifornia, Berkeley, Berkeley, California, United States of America, 5 Centre for Ecological and
...ynthesis, University of Oslo, Oslo, Norway, 6 Data Carpentry, Davis, California, United States

...s contributed equally to this work.
...ftware-carpentry.org

https://doi.org/10.1371/journal.pcbi.1005510

Approximations to 'Best' can still be OK

# TheCarpentries.org

HOME  ABOUT ▾  TEACH ▾  LEARN ▾  JOIN US ▾  OUR COMMUNITY ▾  CONNECT ▾  DONATE  SEARCH  CONTACT

## THE CARPENTRIES

**We teach foundational coding and data science skills to researchers worldwide.**

DATA CARPENTRY

library Carpentry

software carpentry

### What we do
The Carpentries teaches foundational cod-

### Who we are
Our diverse, global community includes

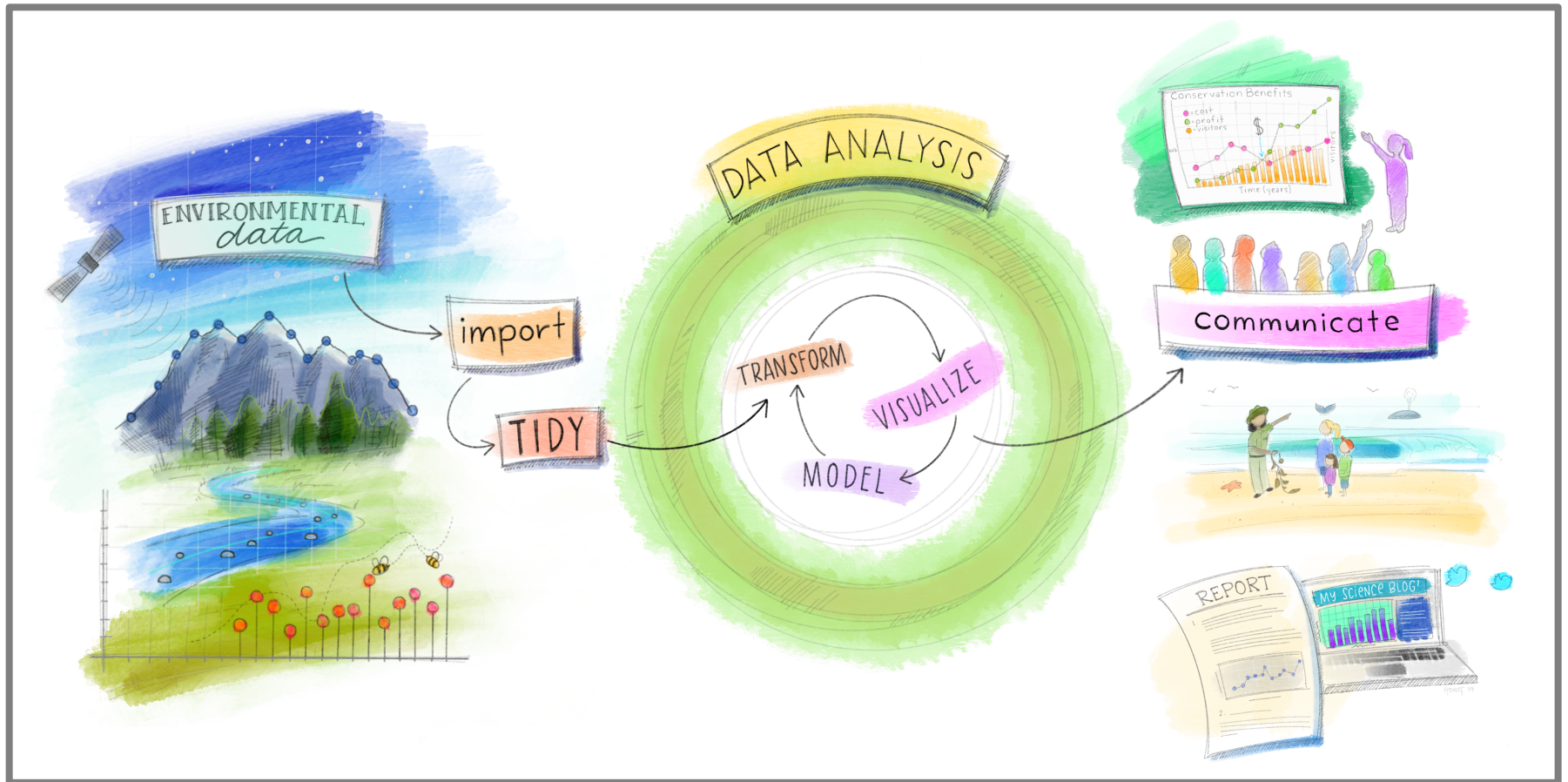### Get involved
See all the **ways you can engage** with The

# Why R ?

- Reproducible Science!
- Full data science workflow
- Language & environment for statistical computing & graphic
- Open source & Free
- Lots of scientists use it!
- AMAZING online community
- Works well with other tools



Artwork by @allison_horst

# Why **R** ?

# Why R ?

- Reproducible Science!
- Full data science workflow
- Language & environment for statistical computing & graphics
- Open source & Free
- Lots of scientists use it!
- AMAZING online community
- Works well with other tools



Artwork by @allison_horst

# MAR 536: Biological Statistics II

*Statistical analysis for biological science graduate students*

Computer labs: Intro to statistical analysis in R

- Partial introduction of tidyverse functions in 2018

Extra credit for using R Markdown in assignments

- 1 student did this in 2017
- All but 1 student did it in 2018
- *Spring 2020 ???*

GF2019

Readable code that remains consistent across tasks


Half in English...half in squiggly


tidyverse
www.rstudio.com

```
 5
 6  fishdata %>%
 7    group_by(year, season, species) %>%
 8    summarise(avg_numbers = mean(abundance),
 9              max_depth = max(depth),
10              avg_temp = mean(surftemp))
11
```

*HT @dataandme*

# Readable code that remains consistent across tasks


Half in English...half in squiggly


tidyverse
www.rstudio.com

```
 5
 6   fishdata %>%          ⭕ "and then"
 7      group_by(year, season, species) %>%
 8      summarise(avg_numbers = mean(abundance),
 9                max_depth = max(depth),
10                avg_temp = mean(surftemp))
11
```

*HT @dataandme*

# MAR 580: Advanced Population Modeling

Fitting ecological models in R & Template Model Builder

*2015*

- separate lectures & computer labs

- many lab assignments

- course materials shared through github repository

# MAR 580: Advanced Population Modeling

Fitting ecological models in R & Template Model Builder

**_2015_**

- separate lectures & computer labs

- many lab assignments

- course materials shared through github repository


I SHOULD NOT HAVE DONE THIS.

# MAR 580: Advanced Population Modeling

Fitting ecological models in R & Template Model Builder

## *2015*

- separate lectures & computer labs

- many lab assignments

- course materials shared through github repository


## *2019*

- students using R Markdown for assignments

- live coding during mixed lab/lectures

- AirMedia to share student screens to class: debugging aid

- course materials shared via Google Drive

# #quantfish woRkshops

# #quantfish woRkshops

*tutorials for beginner and intermediate R users*

- Students & postdocs lead 1.5 hr sessions
- Live coding
- Learning R by doing useful things straight away
- Less is more
- Sharing of materials via GoogleDocs
- Materials version-controlled using git and github
- Feedback asked for (& acted on) often

GF2019

# A GoogleDoc for each workshop....

**quantfishwoRkshop**
**Introduction to R**

**November 13 2018**

**Rstudio interface**
**https://rstudio.cloud/project/134344**

**Materials (continuously updated)**
**https://github.com/thefaylab/quantfishR_01_introR**

**R for data science (online book): https://r4ds.had.co.nz/**

**Sign-in sheet (name, institution, email)**
- Gavin Fay (UMassD), gfay@umassd.edu
- Margot Wilsterman (UMassD), mwilsterman@umassd.edu
- Danielle Lavoie (UMassD), dlavoie1@umassd.edu
- Megan Winton (UMassD), mwinton@umassd.edu
- Renee Halloran (UMassD), rhalloran@umassd.edu
- Greg DeCelles (MDMF), gregory.decelles@mass.gov
- Nicole Danaher-Garcia (UMassD), ndanaher1@umassd.edu
- Flynn Casey (UMassD), acasey1@umassn.edu
- Beth Larson (UMassD), elarson1@umassd.edu
- Jonathan Cummings (UMassD), jcummings@umassd.edu
- Debra Duarte (UmassD), dduarte@umassd.edu
- Susan Inglis(UmassD), singlis@umassd.edu
- Chang Liu (UMassD), cliu@umassd.edu
- Alex Hansell (UMassD), ahansell@umassd.edu
- Brooke Wright (UMassD), brooke.wright@umassd.edu
- Tammy Silva (UMassD), silva@umassd.edu
- Robert Wildermuth (UMassD), rwildermuth@umassd.edu
- Dave Martins (MA DMF), dave.martins@mass.gov
- Marjorie Lyssikatos (UMassD), mlyssikatos@umassd.edu
- Harriet Booth (MA DMF), harriet.booth@mass.gov

**Topics you would like to see**
- **More on tidyverse**
- **Ggplot**
- **Spatial analysis**
- **Time series analysis**

# Standard environment: RStudio & RStudio Cloud

# Code online in github repository
# Auto-updated during workshop



*HT @rachelss*

# What's next?

More conversion of MAR 536 R labs to the tidyverse.
Course Management using R Studio Cloud / github
rstudio::conf


Tuesday Feb 25
Special Seminar at UMassD-SMAST
**"*R* and teamwork for better science in less time"**
*Dr. Julia Stewart Lowndes, NCEAS*

# Thank you!

[gfay@umassd.edu](mailto:gfay@umassd.edu)

[thefaylab.com](http://thefaylab.com)

 [@gavin_fay](https://twitter.com/gavin_fay)

To be added to
#quantfish email list:
[anovak@umassd.edu](mailto:anovak@umassd.edu)

These slides:

[bit.ly/fay_tlearnconf2020](http://bit.ly/fay_tlearnconf2020)

"This is the best! The kind of intro content I've been looking for. I'm really happy this exists!"

"... let us know beforehand which packages we need (my computer is very slow)"

# Intro to R: Analyze US fish data

- Take data from spreadsheet to visualization

- Data wrangling

- Summarizing data by species over time

- Mapping of fish distributions

# My courses that use R

**MAR 536: Biological Statistics II**

- *statistical analysis for biological science graduate students*

**MAR 580: Advanced Population Modeling**

- *fitting ecological models to data*

**Quantfish WoRkshops**

- *tutorials for beginner and intermediate R users*

**MAR 338: Ecological and Environmental Data Analysis in R**

- *coming 2021 ?*

Artwork by @allison_horst

https://vimeo.com/178485416