

R package development class

Toby Dylan Hocking
Assistant Professor
Northern Arizona University
toby.hocking@ucr-prod.com



talk overview

1. data.table NSF POSE project
2. rOpenSci Statistical Software Project
3. Google Summer of Code

who am i

- BA, MS, PHD in Stats/Math (machine learning)
- Assistant Professor of computer science
- Using R since 2003! 20+ years! Author of `data.table`
- `data.table` user since 2015, contributor
- Principal Investigator, NSF Pathways for Innovation
- Open-Source Ecosystems (POSE)

1/3

data.tableNs
project

data.frame

- 2D *columnar* data structure
- rows and columns
- subset rows — `DF[DF$id != "a",]`
- select columns — `DF[, "val"]`
- subset rows & select columns —
`DF[DF$id != "a", "val"]`
- that's pretty much it...



data.table

- Like data.frame, but with more powerful R code syntax, and C code implementation
- R package on CRAN since 2006
- Created by Matt Dowle, co-author Arun Srinivasan since 2013, 50+ contributors
- 1463 other CRAN packages require data.table (in most popular 0.05% of all CRAN packages)

Comparing data.table

- tidyverse R package 1.0 on CRAN in 2016
- tidyverse packages `tibble + readr + tidyr + dplyr ~ c`
- tidyverse uses `DF |> ... |>`, `data.table` uses `DT[...][...]`
- tidyverse is verbose (lots of code), `data.table` is concise
- example: `tibble |> filter(x=="a") |> group_by(z) |> sum`
vs: `DT[x=="a", .(m=mean(y)), by=z]`
- tidyverse has many dependencies, `data.table` has none
- tidyverse has frequent breaking changes, `data.table` does not
(easier for users to upgrade to new `data.table` version)

Why is data.
popular/power
(efficiency)

two kinds of data tab

- **Efficient R code syntax (saves programming**
 - **Efficient C code implementation (saves time**
- data sets can be analyzed using smaller con

data table R code

- think in terms of — rows, what to do with columns
- Matt's 2014 useR talk <https://youtu.be/ql...>

General form: **DT[i, j, by]**

SQL: WHERE

SELECT | UPDATE

data.frame(DF) vs data.table(DT)

```
sum(DF[DF$code != "abd", "valA"])
```

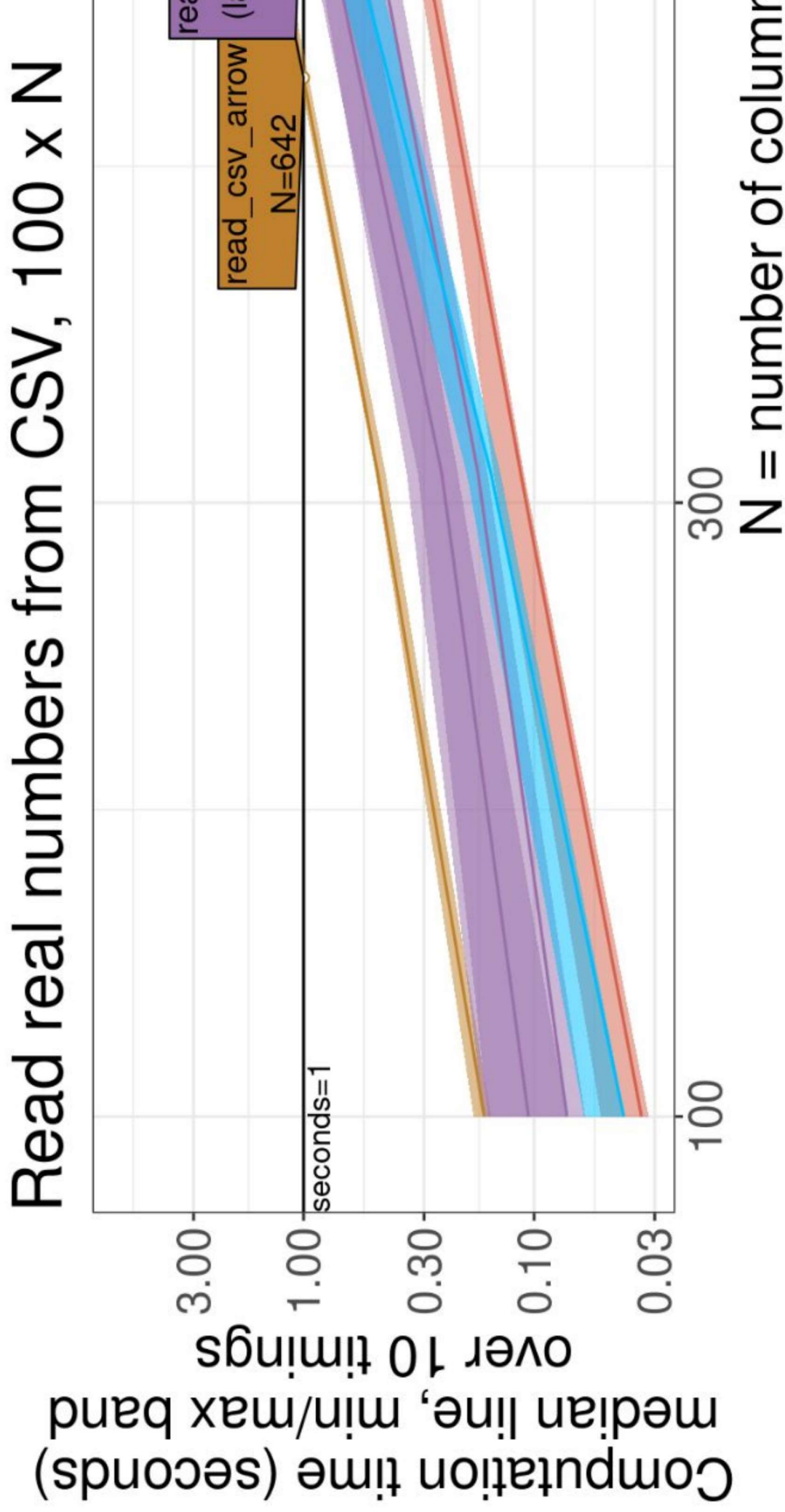
```
DT[code != "abd", sum(valA)]
```

- Consider subset of rows with "abd" in code
- compute sum of values in valA column
- DF needs to be repeated, no repetition
- sum can be placed in the square brackets instead of the sum function than outside with DF.

two kinds of data tab

- Efficient R code syntax (saves programming
- Efficient C code implementation (saves time
larger data sets can be analyzed using smaller
resources)

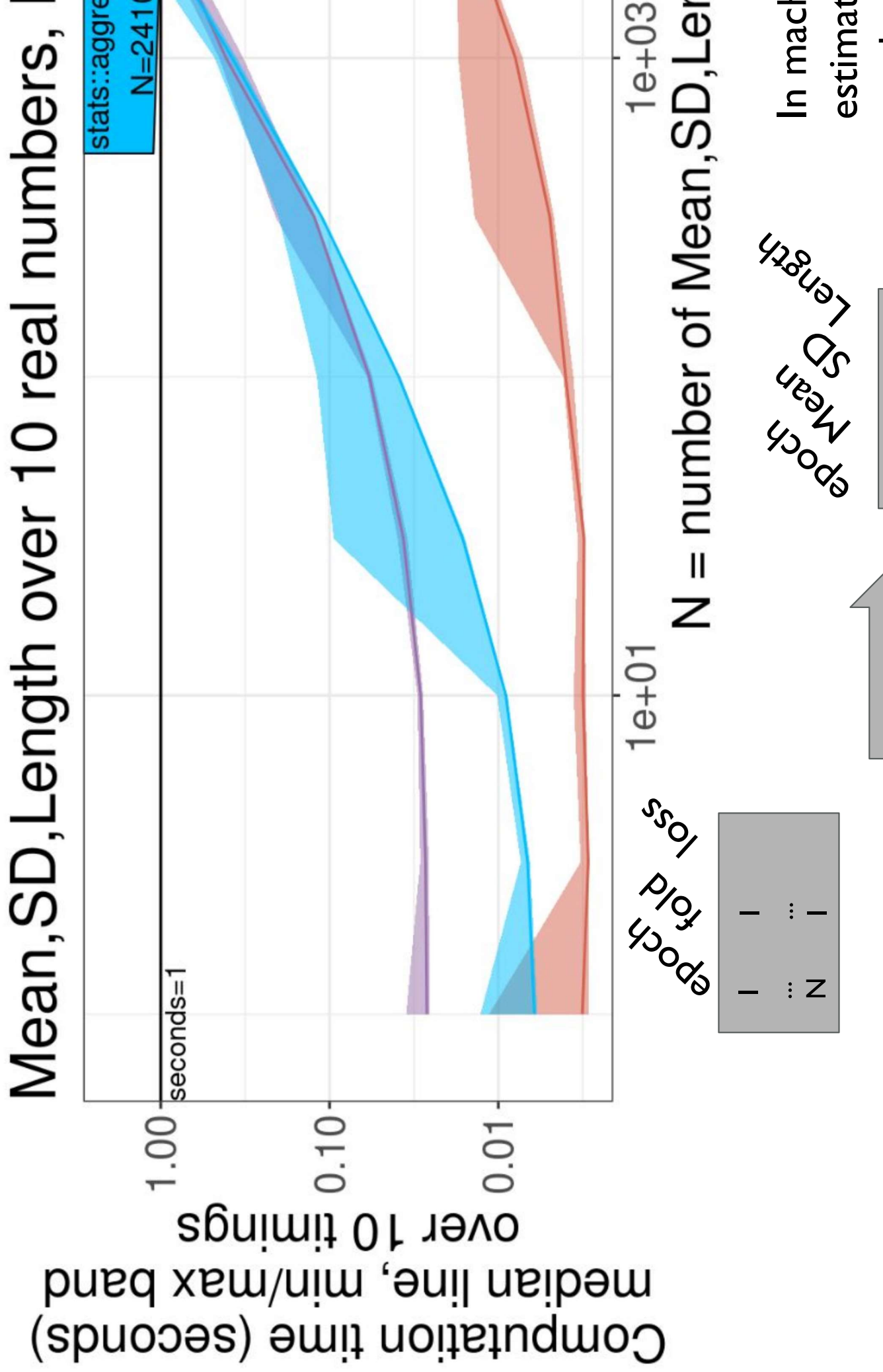
data.table::fread is an extremely efficient



100

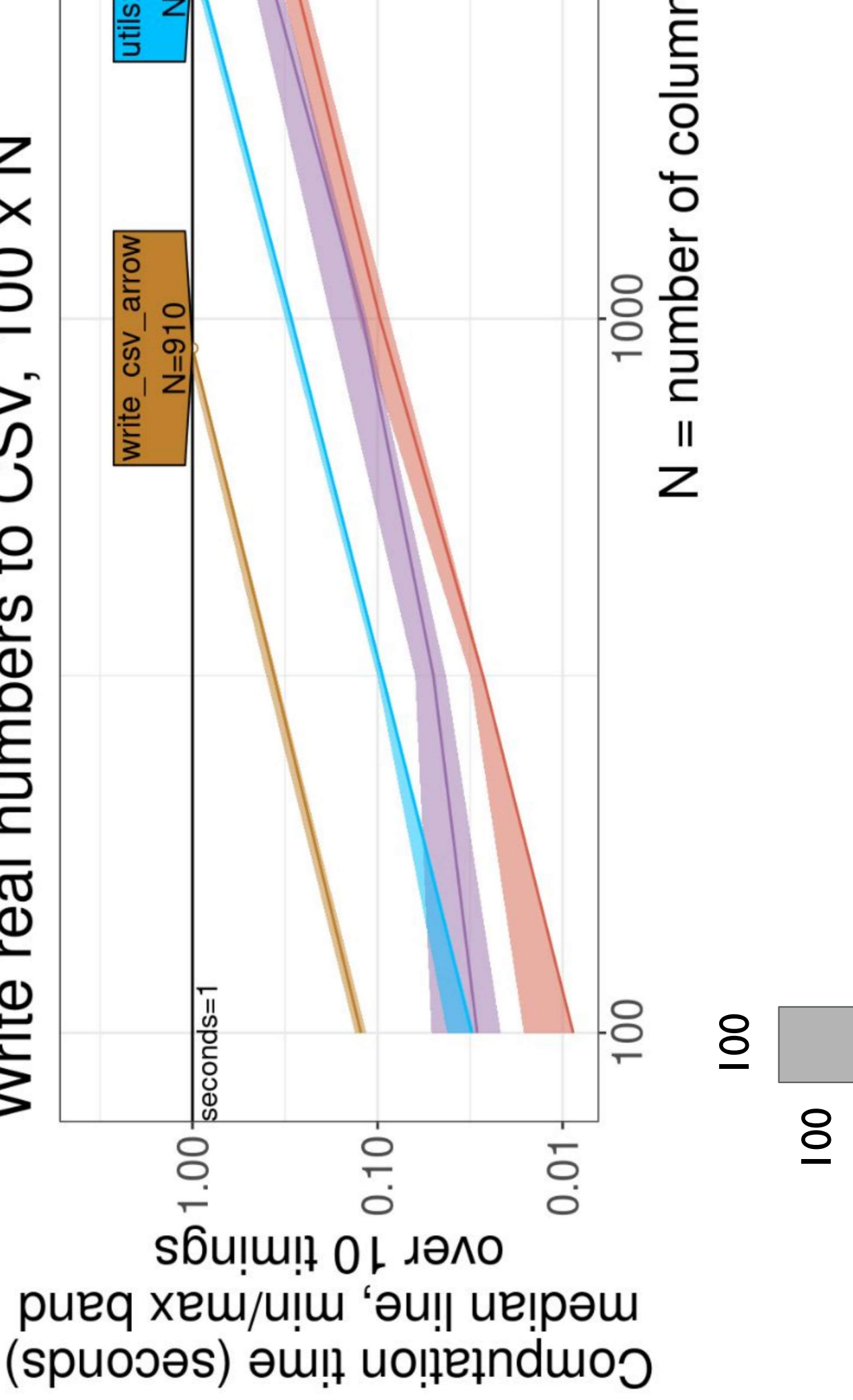
100

data.table computes summaries 100



`data.table::fwrite` is an extremely efficient

Write real numbers to CSV, $100 \times N$



most underrated package



Conor Nash

@conornash

Data.table is the most underrated
has saved me *days* in waiting
complete.

powerful



Alexander Flyax

@aflyax

somebody should just write a version of
`data.table` for `#python`. end of story.
powerful exists at the moment.

data.table data



Joey Reid

@JoeyPReid

data.table

data.table

data.table

data.table

data.table

ggplot2

rstan

knitr

great sadness



Jim Savage

@khakieconomist

With great sadness I was forced
data.table today.

Using data.ts efficient big data

See <https://github.com/tdhock/2023-10-Latin>
hour tutorial presentation slides, with code,

Contributing to
we need you

Community blog a

- data.table mascot is a sea lion, which ba
- data.table community has a new blog, T

<https://rdatatable-community.github.io/>

sea lions often float together on the ocean
called "rafts." - Marine Mammal Center

GitHub repos

- data.table has an active issue/Pull Request(PR) thread
- <https://github.com/Rdatatable/data.table/>
- 1000+ open issues, 100+ open PRs
- if you have any time/interest, we could use your help
- easy first contribution: try reproducing an issue (very helpful to know if an issue is reproducible)
- very inclusive community -- after you submit your PR, you can be a maintainer
- join the github group!
- now is a very exciting time to get involved, as we are

Translation Awards

- In 2023-2025, National Science Foundation has been supporting the expansion of the ecosystem of users and contributors
- 20 translation awards, US\$500 each, in order to make messages more accessible, ideas:
 - Translate errors/warnings/messages (potools)
 - Translate most important vignettes (intro, important messages)
 - Translate other documentation (cheat sheets, tutorials, ...)
- Priority: Portuguese, Spanish, Chinese, French

Travel awards

- In 2023-2025, National Science Foundation has expanding the ecosystem of users and contributors
- Eight travel awards, US\$2700 each
- Candidates should give a talk about data.table relevant audience (potential data.table users or
- Call for proposals on <https://rdatatable-commu>

Summary of data-

- concise, consistent syntax
- time and memory efficient
- No dependencies (easy to install)
- No breaking changes (easy to upgrade)
- Inclusive user/developer community to contribute:
- translation awards 115\$500 each

2/3

OpenSci Sts
Software Peer

Sharing R package

- R package is the formal unit of code sharing
- Package includes documentation, tests, vignettes
- CRAN is the Comprehensive R Archive Network, the most widely used package repository
- CRAN regularly checks every package to ensure it works correctly (and compatible with all other packages) to benefit for the R community.

Example of review

Please rather use the Authors@R field and declare Maintainers and Contributors with their appropriate roles with packages, e.g. something like:

```
Authors@R: c(person("Alice", "Developer", role = c("MAINT", "AUTHOR"),  
email = "alice.developer@some.domain.net"),  
person("Bob", "Dev", role = "aut") )
```

Please always write package names, software names and programming interface) names in single quotes in titles, e.g.: --> 'C'; 'PCRE'; ...

Please note that package names are case sensitive.

Peer-reviewed research

- Peer-reviewed paper is the formal unit of
- Peer reviewers typically read the manuscript about
- But typically peer reviewers do not read the whole paper. They read the abstract and the conclusions. They also read the
- the computations in those papers. There are many other things that they read. For example Journal of Statistical Software.

Example of peer

(**) Section 1 and abstract: The proposed model uses along the data to constrain the inference during the programming algorithm. I don't understand the reason discussing train and test data. Training is typically coefficients of variables. In this case, the learning intervals with the presence or absence of a change. by experts using already known data (training set). ring to apriori information would be less confusing.

...

(**) Page 10 Algorithm and line 48: "... by the Deco We would need more details about the recovering of v tracking step the reason why the authors used notation description of the costs? Do we have vector equality

rOpenSci Guide and

- Goal is to provide thorough peer review of R packages (code/documentation/etc), and do for peer-reviewed research papers.
- <https://stats-devguide.ropensci.org/>
- case study of canaper R package peer review <https://github.com/ropensci/software-review>

3/3

Google Summe

What is GSC

- 3 month free/open-source coding project
- You don't have to be a coding expert to participate
- Goal is to teach you how to contribute to open source software projects

- I have been co-administrator of R project for 10+ years, and I have mentored many students (typically college students)

Guide and examples

- List of all organizations

<https://summerofcode.withgoogle.com/projects>

- R project idea list:

<https://github.com/rstats-qsoc/qsoc2024-proposed%20coding%20projects>

Thank you Questions

Toby Dylan Hocking
Assistant Professor
Northern Arizona University
toby.hocking@ucr-prod