

# **Fast**

- Why is Julia fast?
- · Why is C fast?
- Why are R & Python slow?

# Why is C fast?

- C is compiled
  - · What does it mean?
- A function in C knows
  - the type of variables it uses and
  - the type of outputs
- Based on these information one can optimise code
- In general the more you know about something the better you can optimize algorithms (e.g. categorical sumby - demoed later)

# Why are R & Python slow?

- R & Python are interpreted not compiled
- A function in R or Python does not have the same level of "knowledge" of a C function

### Julia is fast because

- It pretends to be R & Python
- But specialises its compiled code to the type of inputs like C

Julia compiles the same function to specialized code depending on type of input

```
@code_warntype testmul(10, 100)
@code_warntype testmul("b","c")
```

## Implications for writing performant code

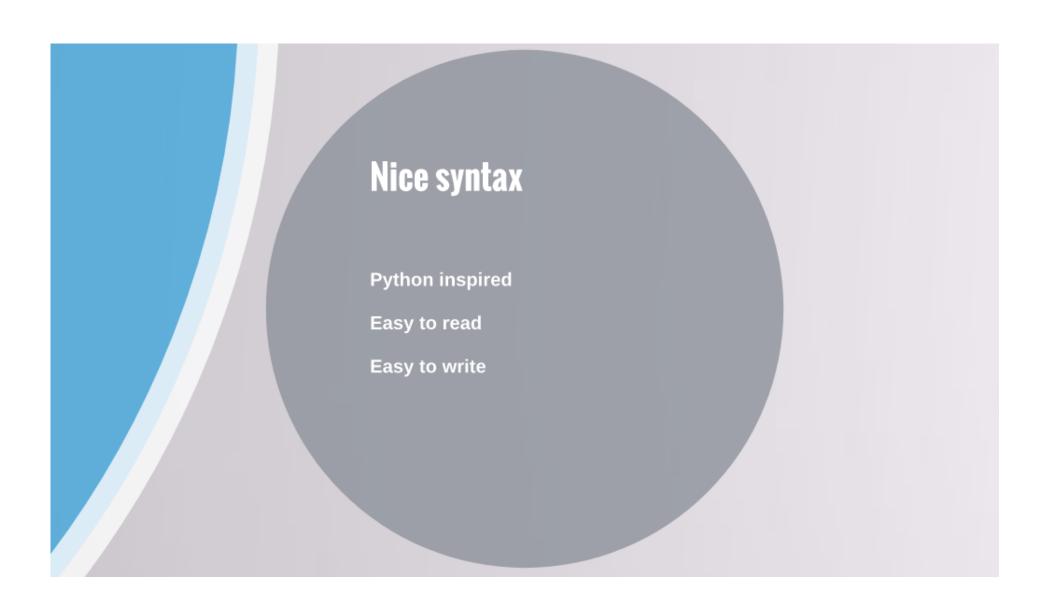
Type-stability is important
Types: Integers, Strings, Floats, etc.

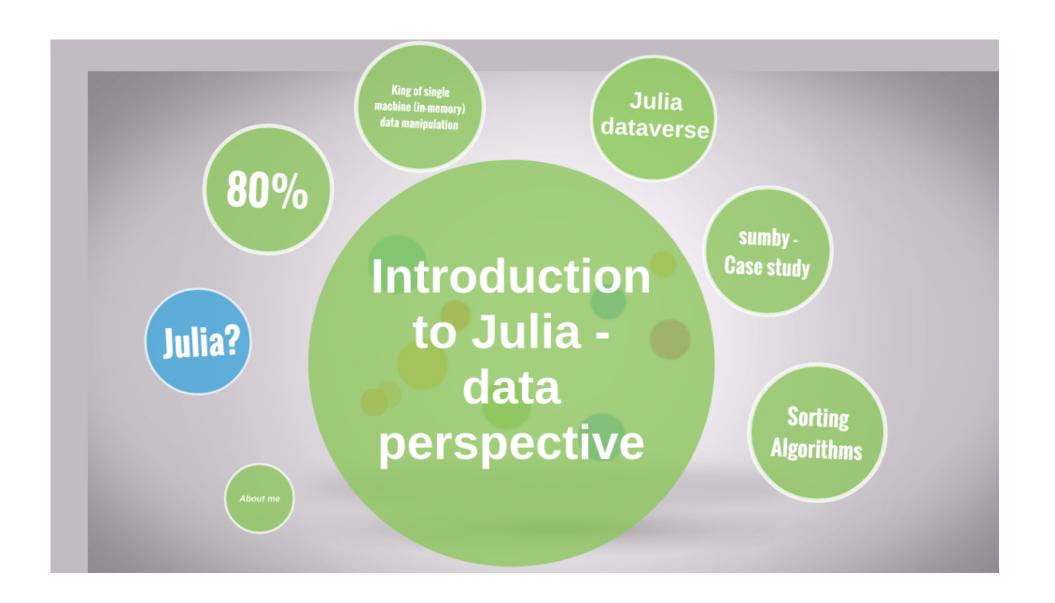
If Julia can't infer type then it will use ANY as type Speed might reduce to R & Python

Avoid switching the type of variables within functions

Avoid returning different types of outputs on the same type of inputs

**How does this impact DataFrames?** 



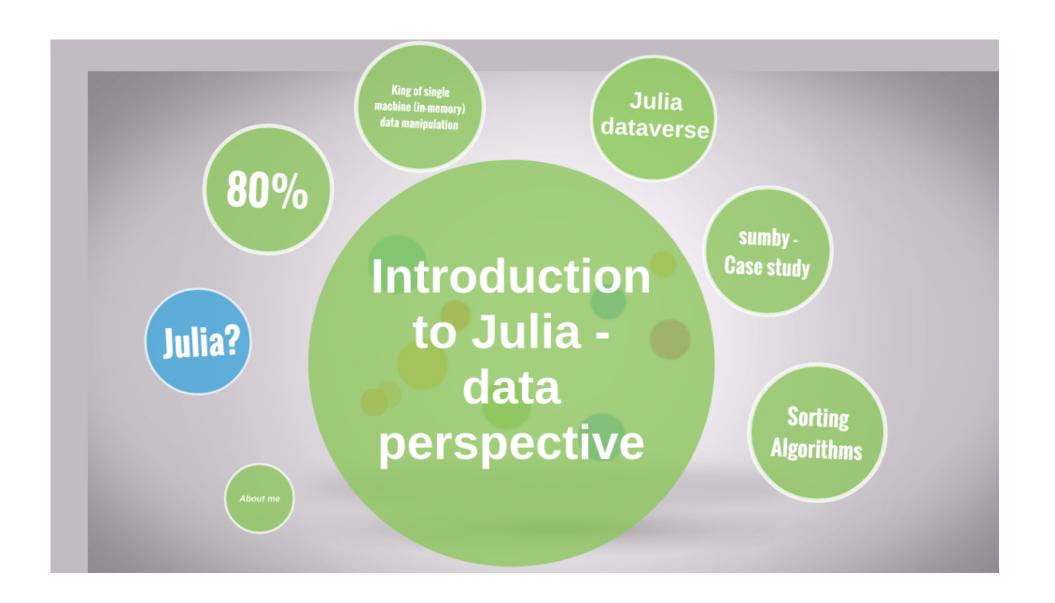




# Life of machine learning expert in banking

"Because finance data is so tricky, Amen suggests that one of the best uses for machine learning in a financial context is simply the cleaning and processing of data."

https://news.efinancialcareers.com/auen/301629/machine-learning-finance/

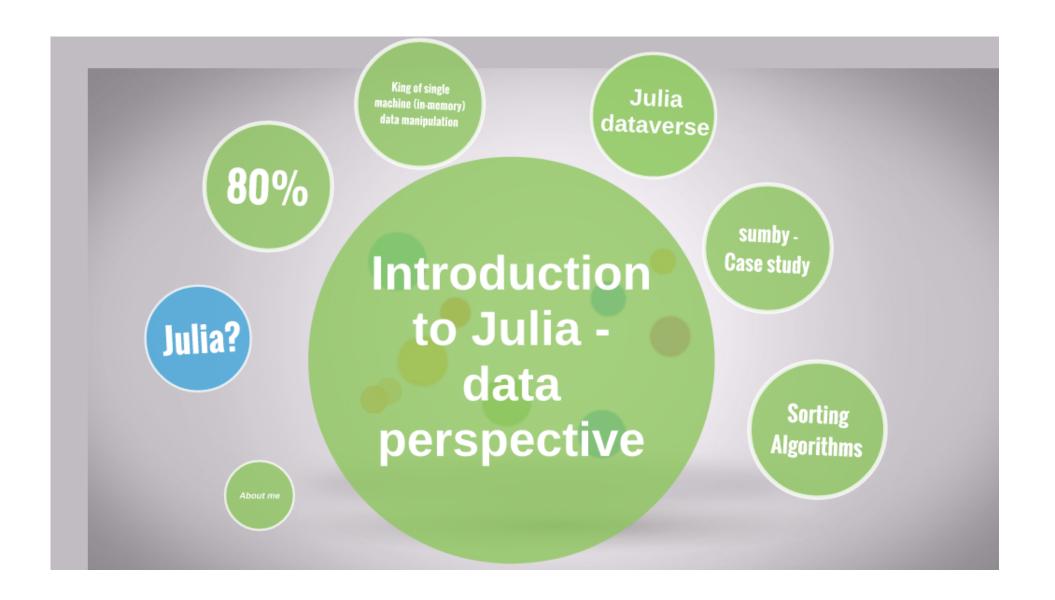


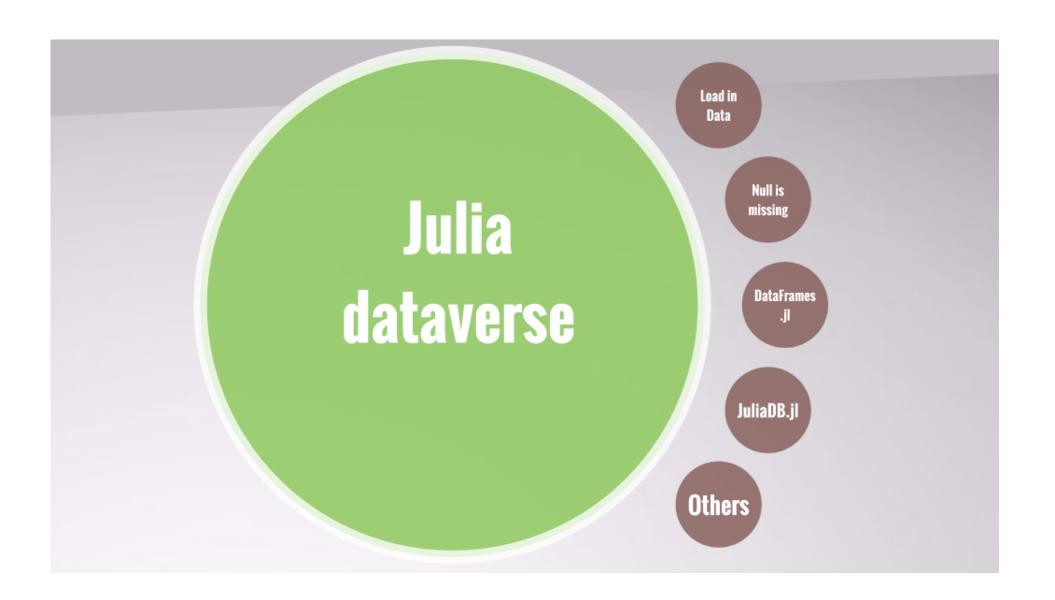
# King of single machine (in-memory) data manipulation

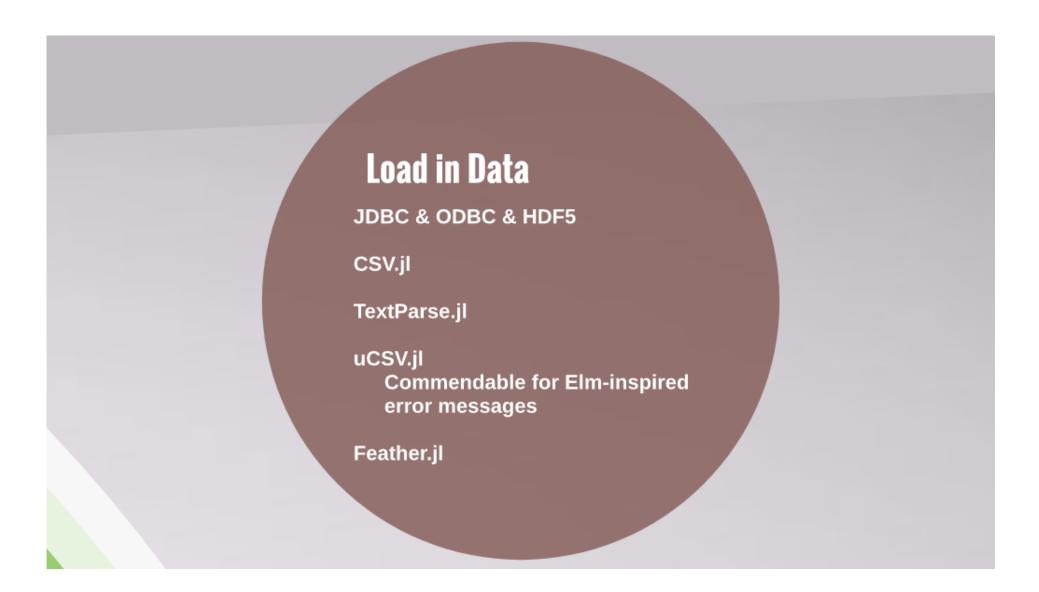
R's data.table

Generally the fastest at almost everything I've benchmarked

Including Group-by, sorting, reading and writing CSVs







#### **Null** is missing

- R has NA which is used for missing values
- · Julia has missing

```
using Missings 

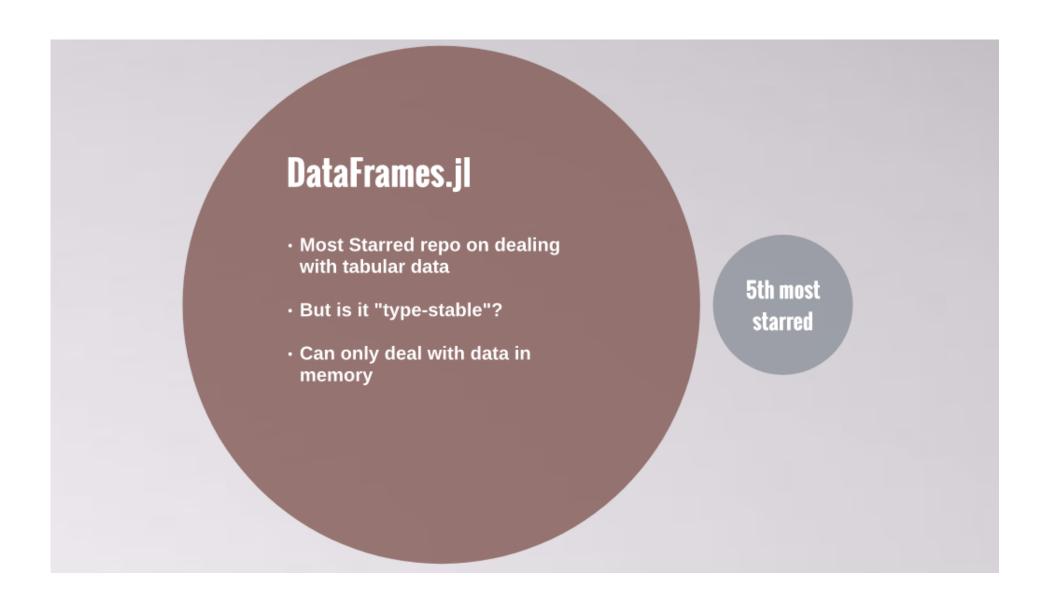
1 + missing | missing | 
true | missing | true | 
true & missing | missing |
```

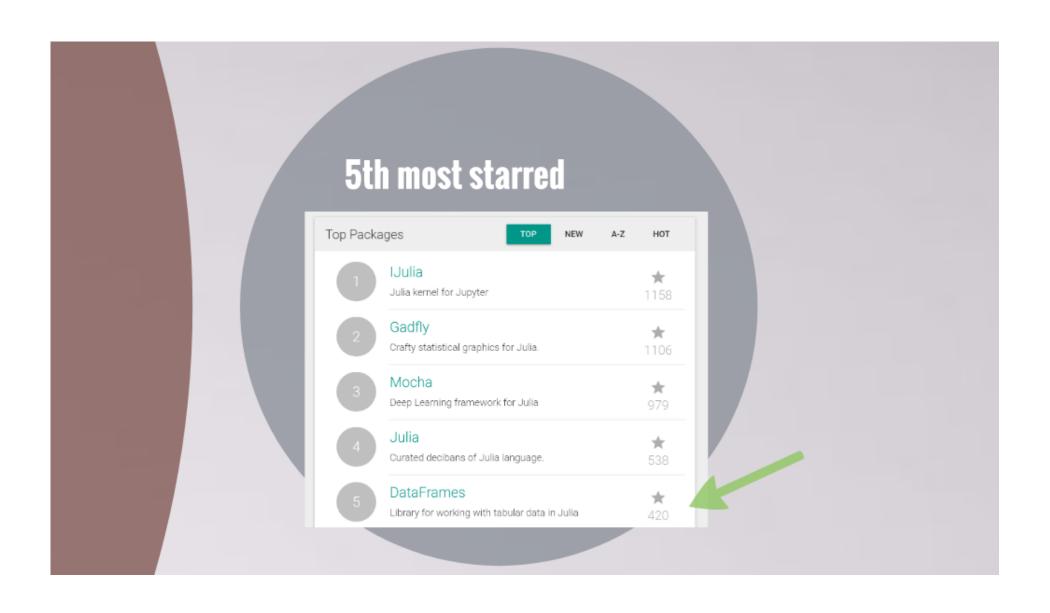
Should use Missings.missing to represent null values

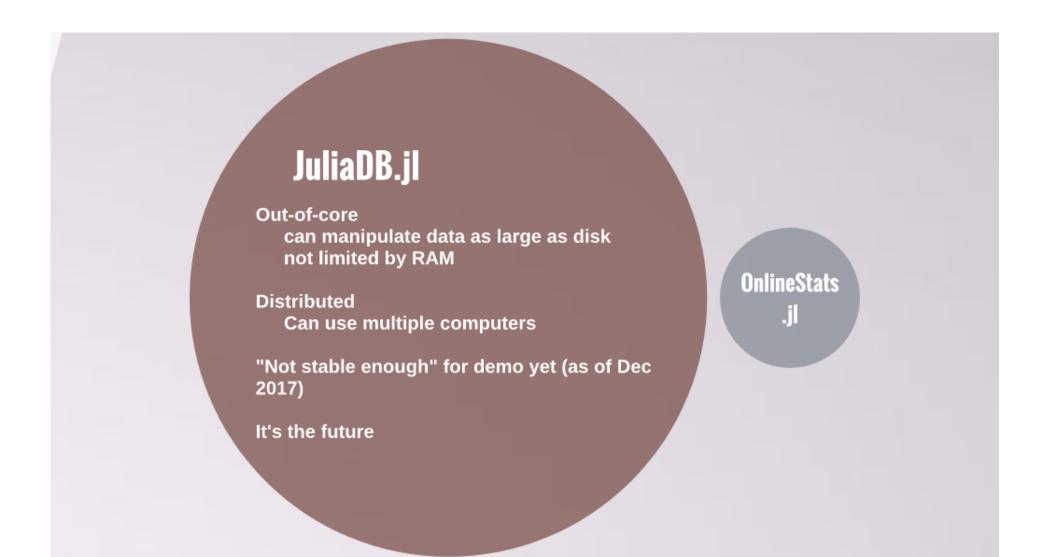
> Nullable and Nulls are outdated



# **Nullable and Nulls are** outdated • Literature before Dec 2017 may refer to Nullable or Nulls · Those are obsolete Use Missings.jl's missing instead







# OnlineStats.jl Algorithms that work on one chunk of data at a time avoid RAM limit issues and • enable parallelism Allows Julia to eventaully take on SAS

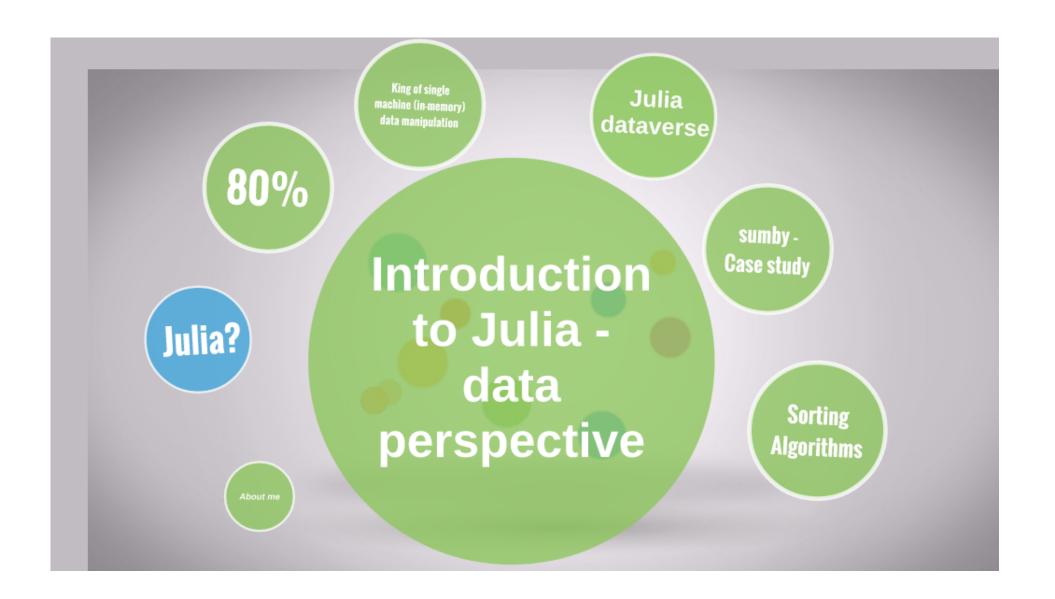
#### **Others**

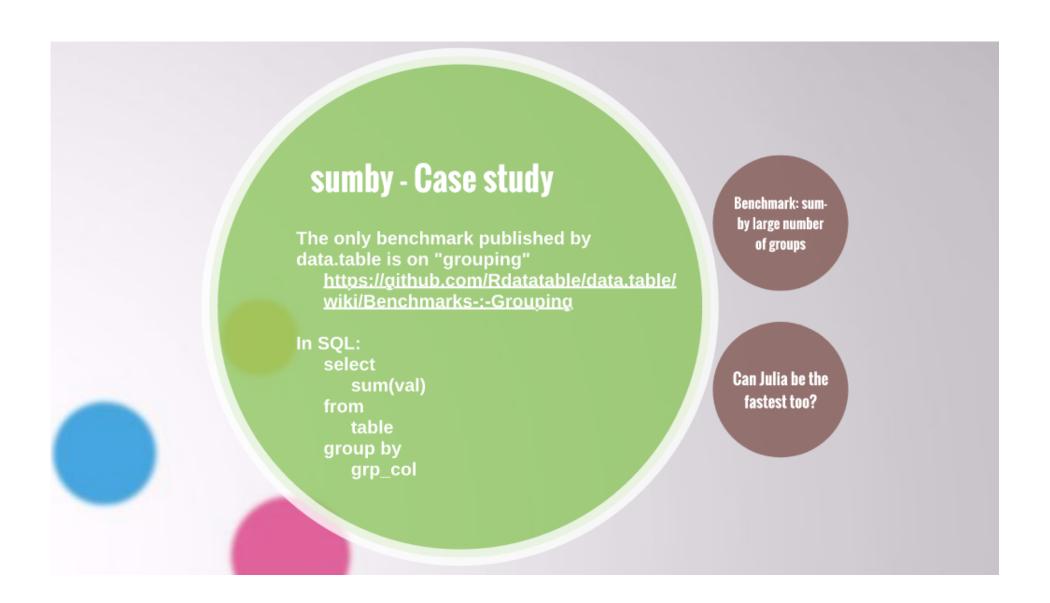
ODBC & JDBC supported by JuliaComputing

IterableTables.jl convert between table-types easily

IndexedTables.jl type-stable? tables

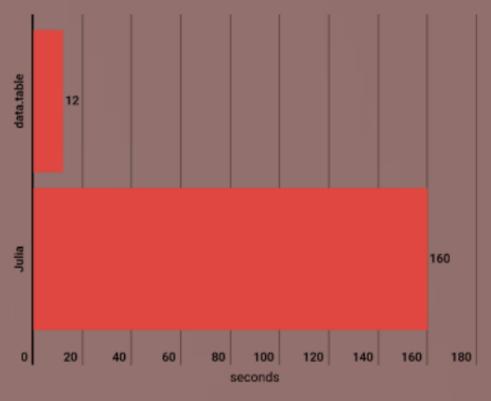
Query.jl a LINQ or dplyr-esque framework seems "slow"



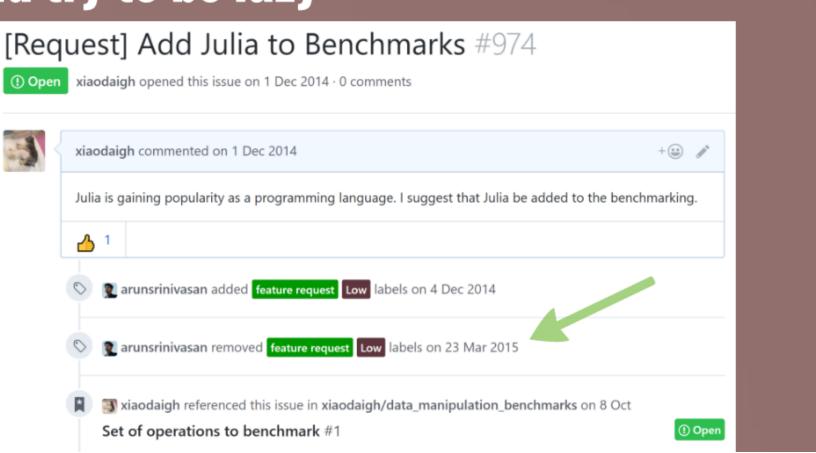


# Benchmark: sum-by large number of groups

- · 200 million rows
- 2 million distinct groups



# I did try to be lazy



### data.table is C backed

Basically it's Julia vs R via C

data.table is generally hard to beat

Shout out to @MattDowle and @arun\_sriniv + 20 others for such an amazing package





# **Shoulders of giants**

Read somewhere that data.table uses radixsort

Found radixsort in SortingAlgorithms.jl

Radixsort was written entirely in Julia; so easy to pick-up

Re-purposed the algorithm for sum-by

Now it's implemented in FastGroupBy.jl

## Radixsort & Multithreading

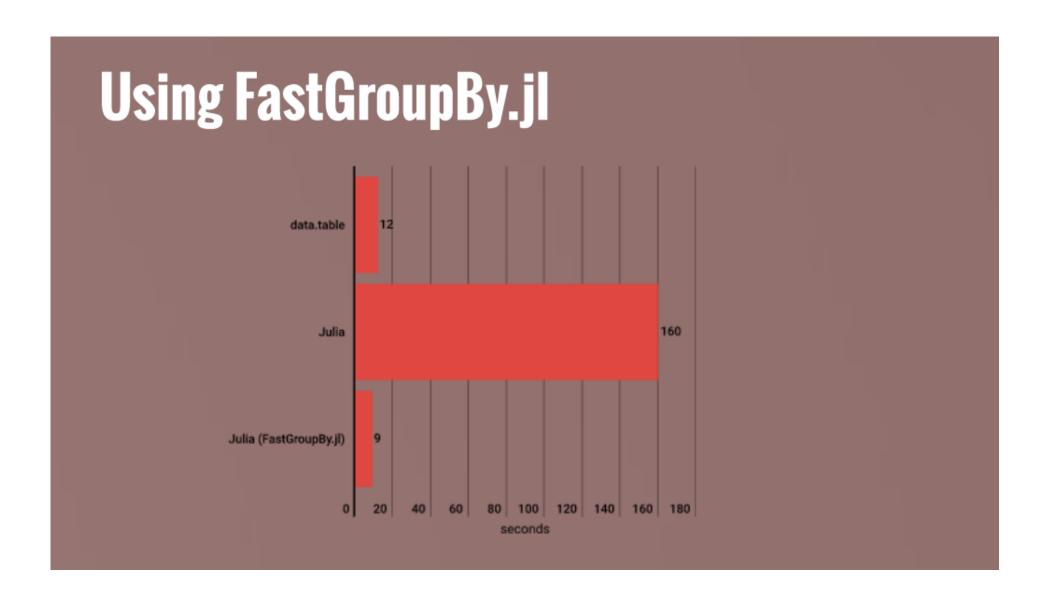
Can sort "bits" type in O(N) time

Uses more memory as it makes a copy of the original vector

It's not a comparison-based sort

Other comparison-based algorithms are O(N\*logN)

Used multithreading to work on mutually-exclusive chunks of the array to further speed-up



# Demo

https://github.com/Julia-Sydney/Julia-Sydney-Talks/tree/master/20171207

#### Misc

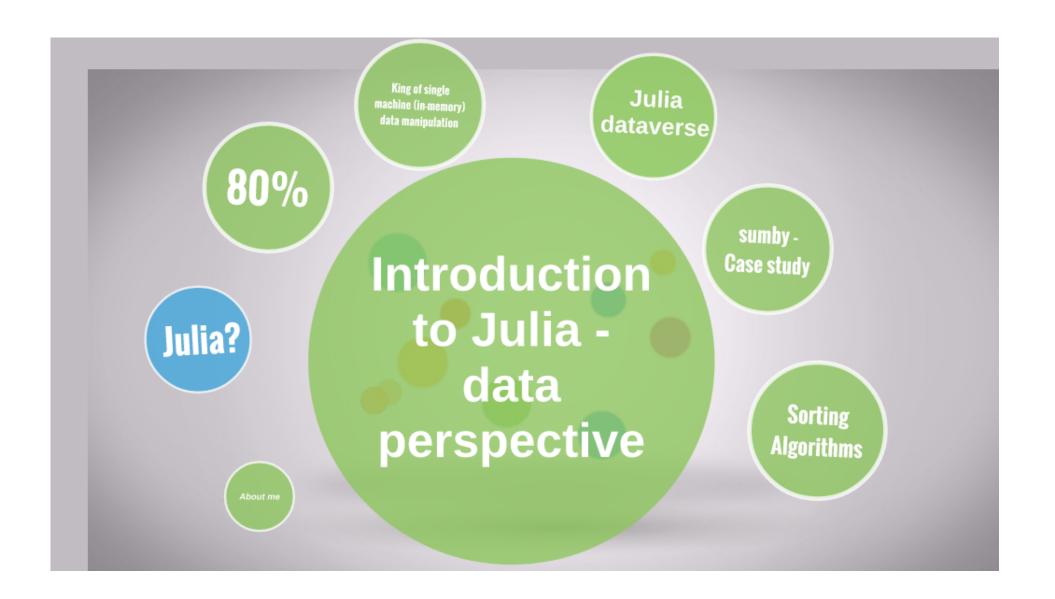
- FastGroupBy.jl's sumby is faster when
  - the by-vector is categorical and we know the categories before hand
  - the by-vector is already sorted (or more precisely just grouped)
- FastGroupBy.jl's sumby is slower when
  - the vector is not 100 million+ in length (but "no par")
  - by-vector is not bits type e.g. String
- Blog: <a href="https://www.codementor.io/zhuojiadai/an-empirical-study-of-group-by-strategies-in-julia-dagnosell">https://www.codementor.io/zhuojiadai/an-empirical-study-of-group-by-strategies-in-julia-dagnosell</a>

### Can Julia be the fastest too?

I think it can but...

It will take time for people to find inefficiencies and implement better algorithms

We have data.table as the beacon



## **Sorting Algorithms**

Julia's sorting algorithms are not yet well optimized

Comparison-based sort is O(nlogn)

- let n be size fo input
- time\_to\_sort(n) = f(nlogn) for some linear function f
- time\_to\_sort(n) = a + b\*nlogn; size of a and b matters!

Julia comes with 3 sorting algorithms

- Quicksort One of best "on average"
- Mergesort Fast but requires more memory
- InsertionSort very quick for small data size

# **Joking Algorithms**

**Random Sort - "Wanna play pick up 52 cards?** 

**Sleep Sort - Interesting!** 

# String sorting/grouping is a key pillar of data science

Group By operations is one of the most common

Being able to groups items together is key

#### **Group By in SQL**

```
    select
        grps,
        sum(val) as tot_val
        from
        some_table
        group by
        grps;
```

# String-sorting in R very fast

#### R

sort(x, method="radix")

#### Julia

- No string radix. I implemented radixsort(x)
- · about 20% slower than R

## Why is R so fast?

R's strings unlike Julia's are stored "centrally" in a cache (look up table)

Using a clever algorithm, it only needs to sort the unique strings

If the proportion of uniques to total is large; then R's advantage vanishes

#### The future

- InternedStrings.jl Similar to how srings are stored in R
- Use CategoricalArrays.jl specialised fast algorithm for groupby
  - equivalent of R's factor type
  - · Grouping known categories before hand
  - Can program fast paths for easy functions like sum
    - · this is done in data.table

# Sorting Algorithms. jl

SortingAlgorithms.jl provides 3 (soon to be 4) more

- TimShot the default in Python, exploits common sequences
- HeapSort uses a heapstructure
- RadixSort R's default; very fast if implemented right
- StringRadixSort My contributions sorts strings fast

