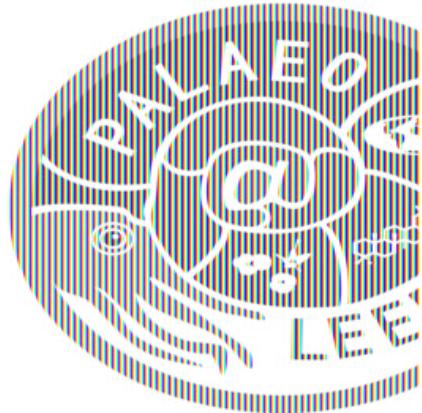


# Sampling bias in the fossil record

Bethany J. Allen | Alex M. Dunhill | Graeme T. Lloyd



# Workshop structure

- Alex talk (video | slides)
- This talk (video | slides)
- Scripts (hands on)
- Ask questions (best to do other bits first!)
  - Drop-in Zoom, Friday 12th June 2pm-3:30pm BST, check your emails for the link
  - Use the '#workshop' channel on Discord
  - Email the organisers

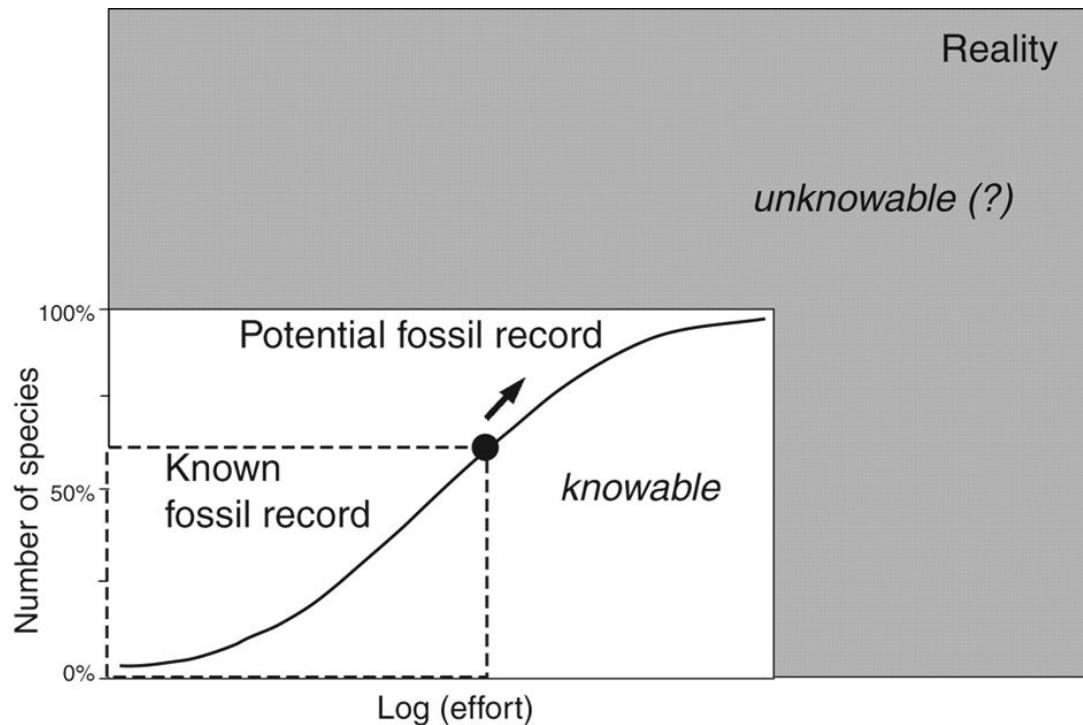
# Outline for this talk

Core concepts

Big data

Methods

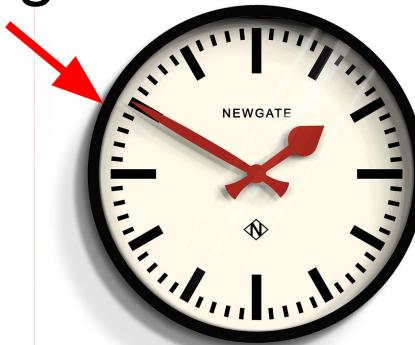
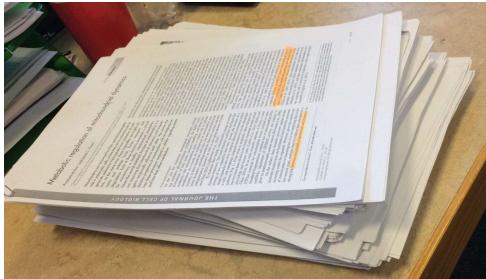
# The fossil record is an incomplete record



# Sampling has many facets



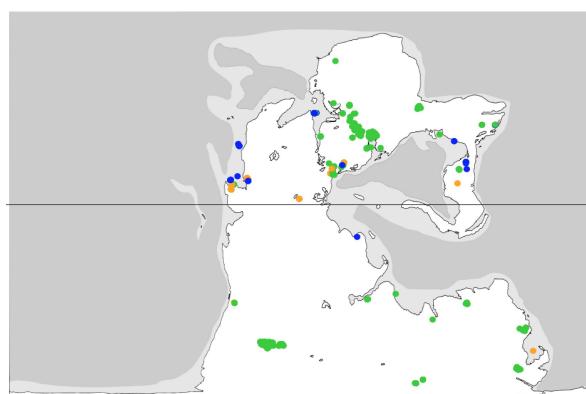
Sampling



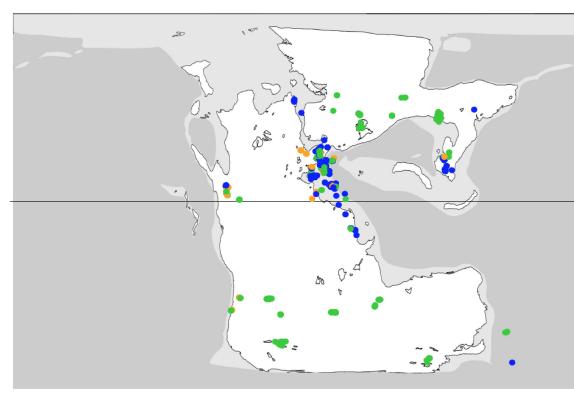
# Sampling varies spatially and temporally



Late Permian

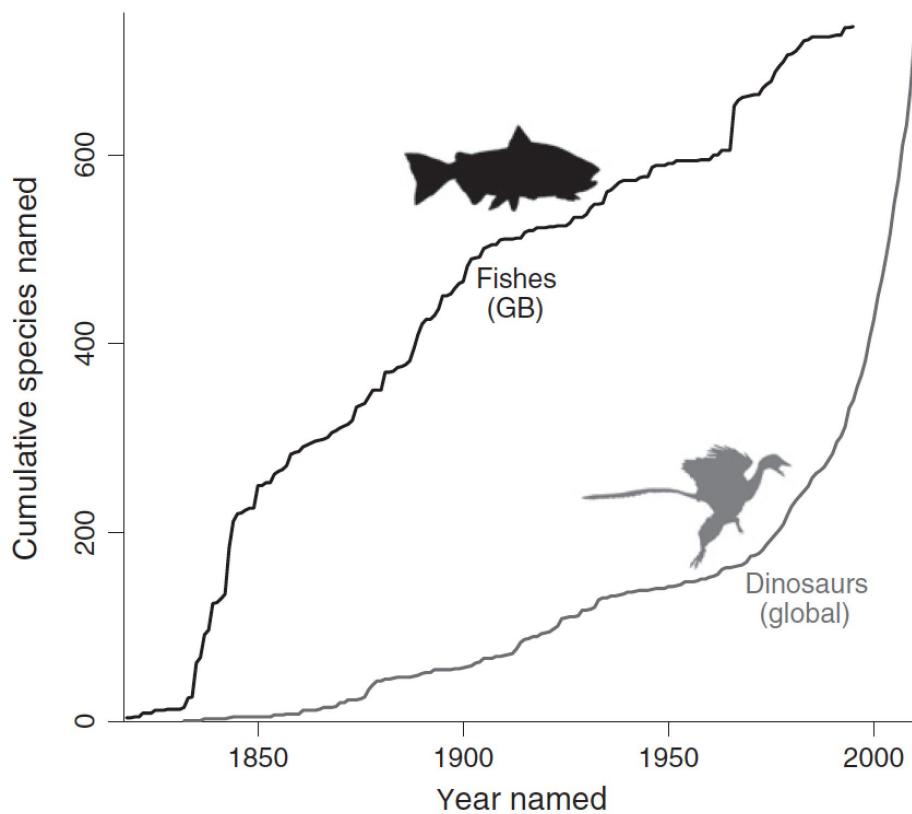


Early Triassic

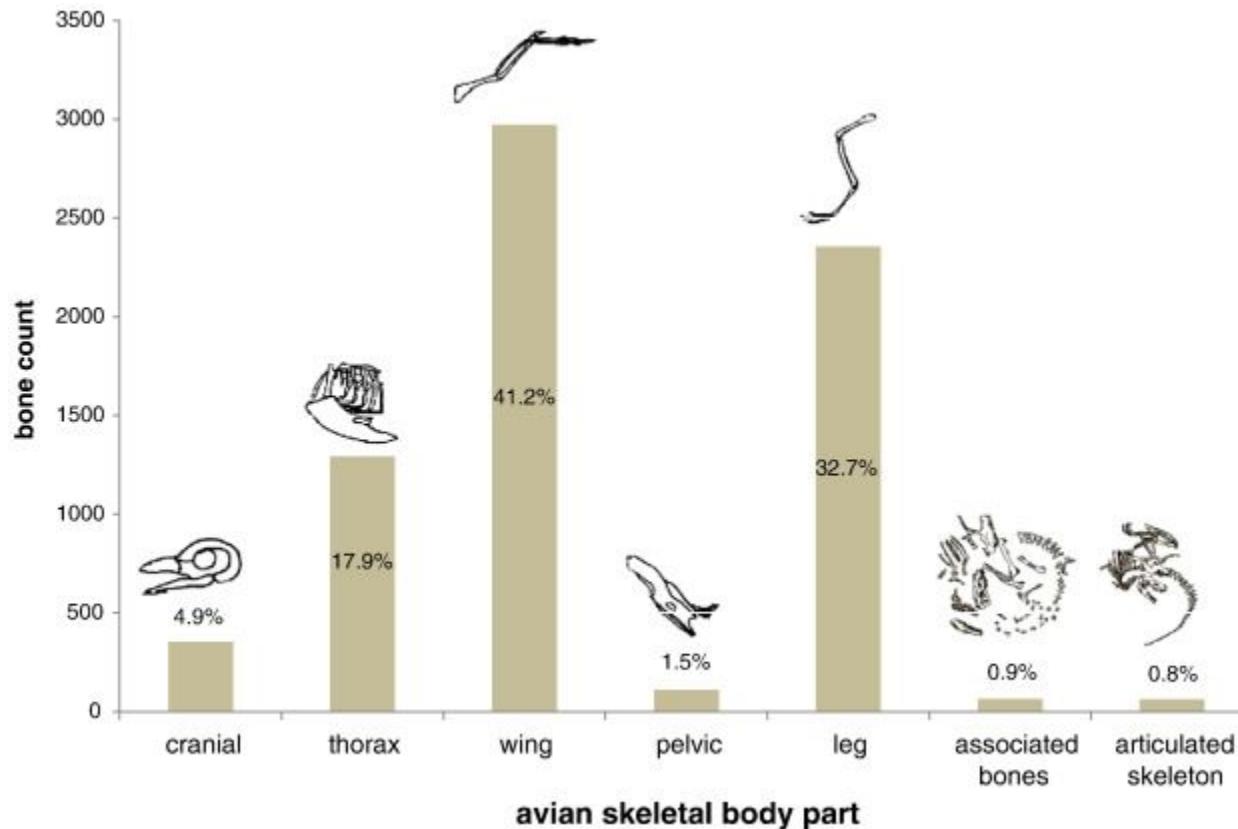


Middle Triassic

# Sampling varies by clade



# Sampling varies anatomically



# Sampling as explanatory variable: patient zero

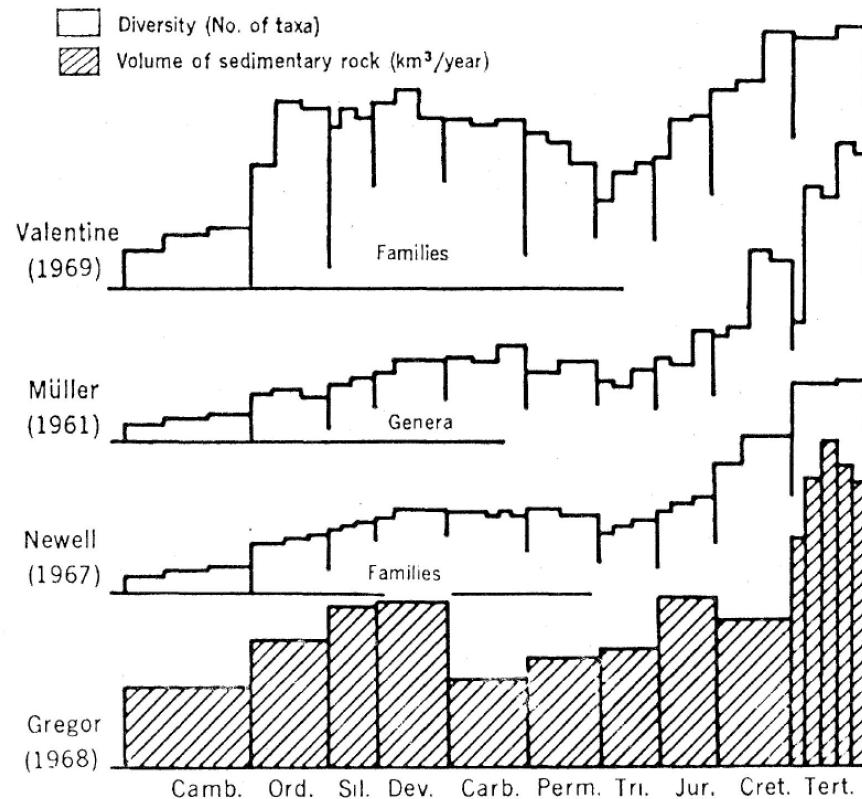


Fig. 1. Comparison of the number of taxa and the volume of sedimentary rock during the Phanerozoic. The diversity data are based mainly on well-skeletonized marine invertebrates (1, 8, 9, 12).

# Sampling can drive pattern



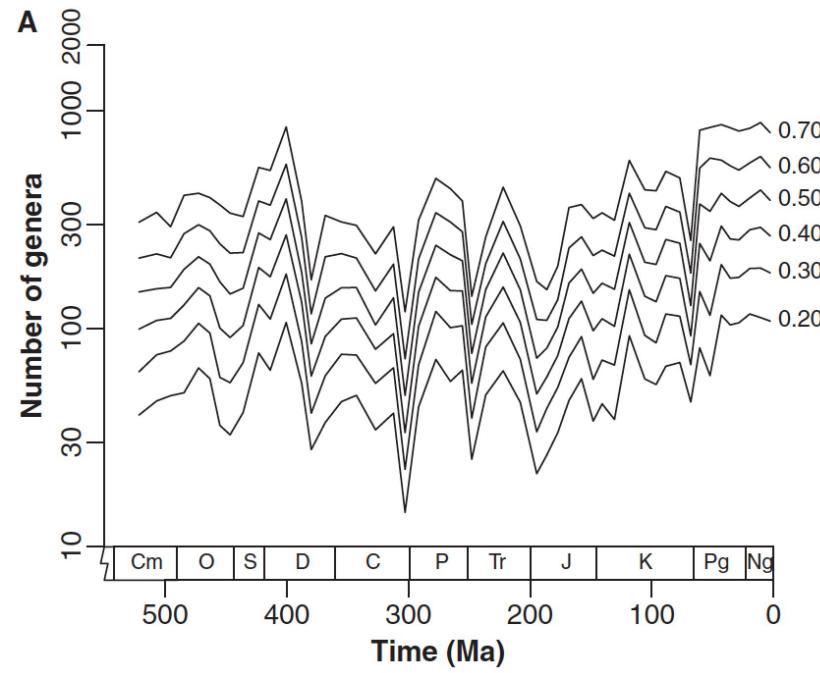
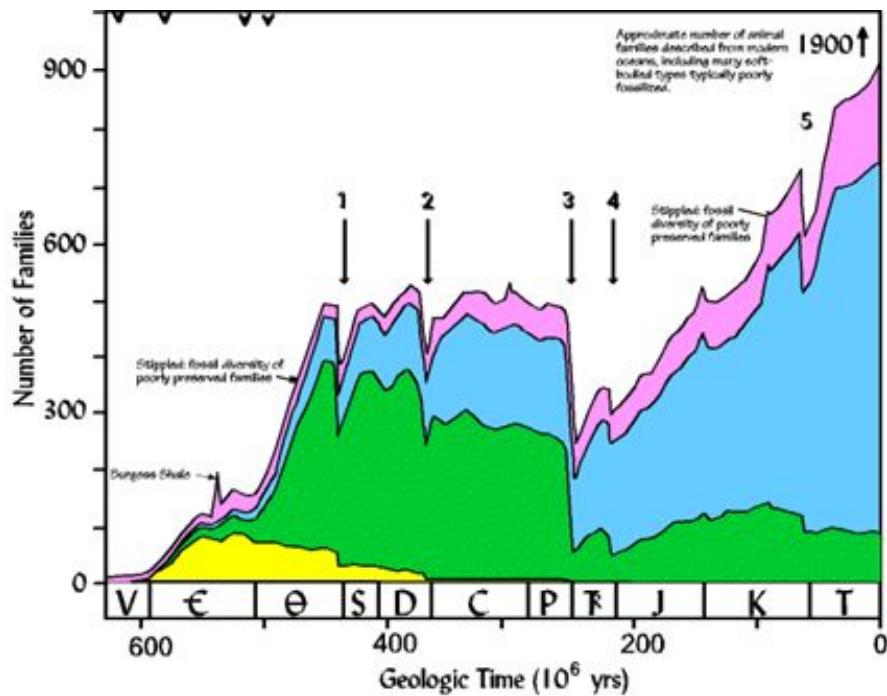
**George Takei**   
@GeorgeTakei

"If we didn't do any testing, we would have very few cases."

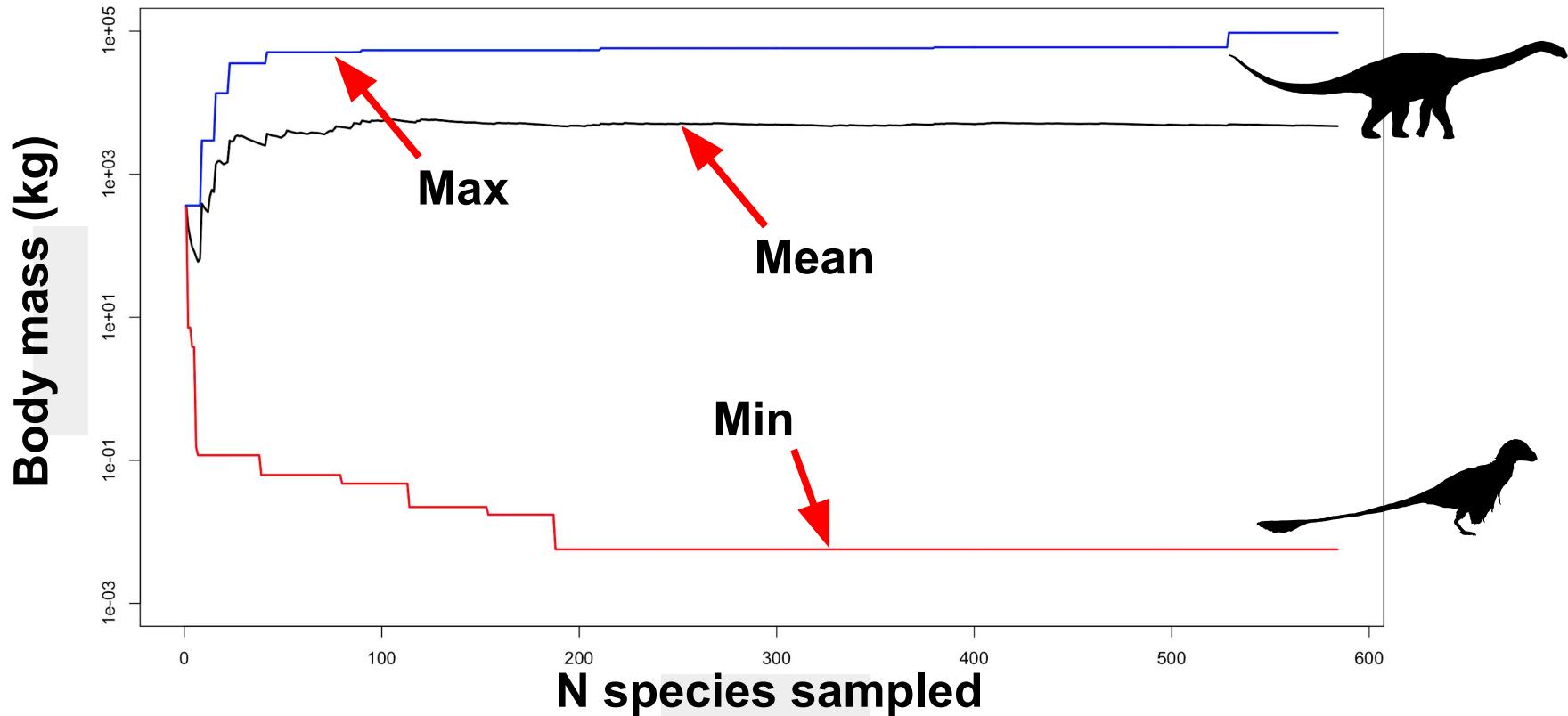
— Donald J. Trump

12:12 AM · May 16, 2020 · Twitter for iPhone

# Removing sampling bias



# The monotony of it all



Data from: Benson et al. 2018; *Palaeontology*

# The monotony of it all [code]

```
# Code to produce plot on previous plot

# Grab the Benson dinosaur body mass data from GitHub:
DinoBodyMasses <-
read.table("https://raw.githubusercontent.com/bethany-j-allen/sampling_bias_workshop/master/Data/DinoBodyMasses.txt?token=ABZPFHSQJZOGHA
KWFKTP73K6ZUEPA", header = TRUE)

# Shuffle the order to randomly resample the data:
x <- sample(DinoBodyMasses[, "Mass_kg"])

# Set random seed (so you will get exactly same plot):
set.seed(9)

# Resample data as minimum, mean and maximum:
Samples <- do.call(rbind, lapply(as.list(1:length(x)), function(y) {z <- x[1:y]; c(min(z), mean(z), max(z))}))

# Add column names to output:
colnames(Samples) <- c("MinMass", "MeanMass", "MaxMass")

# Make plot:
plot(x = 1:length(x), y = Samples[, "MeanMass"], log = "y", xlab = "N species sampled", ylab = "Mass (kg)", type = "l", ylim = c(0.001, 100000), lwd = 2)
points(x = 1:length(x), y = Samples[, "MinMass"], type = "l", col = "red", lwd = 2)
points(x = 1:length(x), y = Samples[, "MaxMass"], type = "l", col = "blue", lwd = 2)
```

# Big data



The Paleobiology Database  
revealing the history of life



# Big data



PaleobioGoogling  
@PaleoGoogling



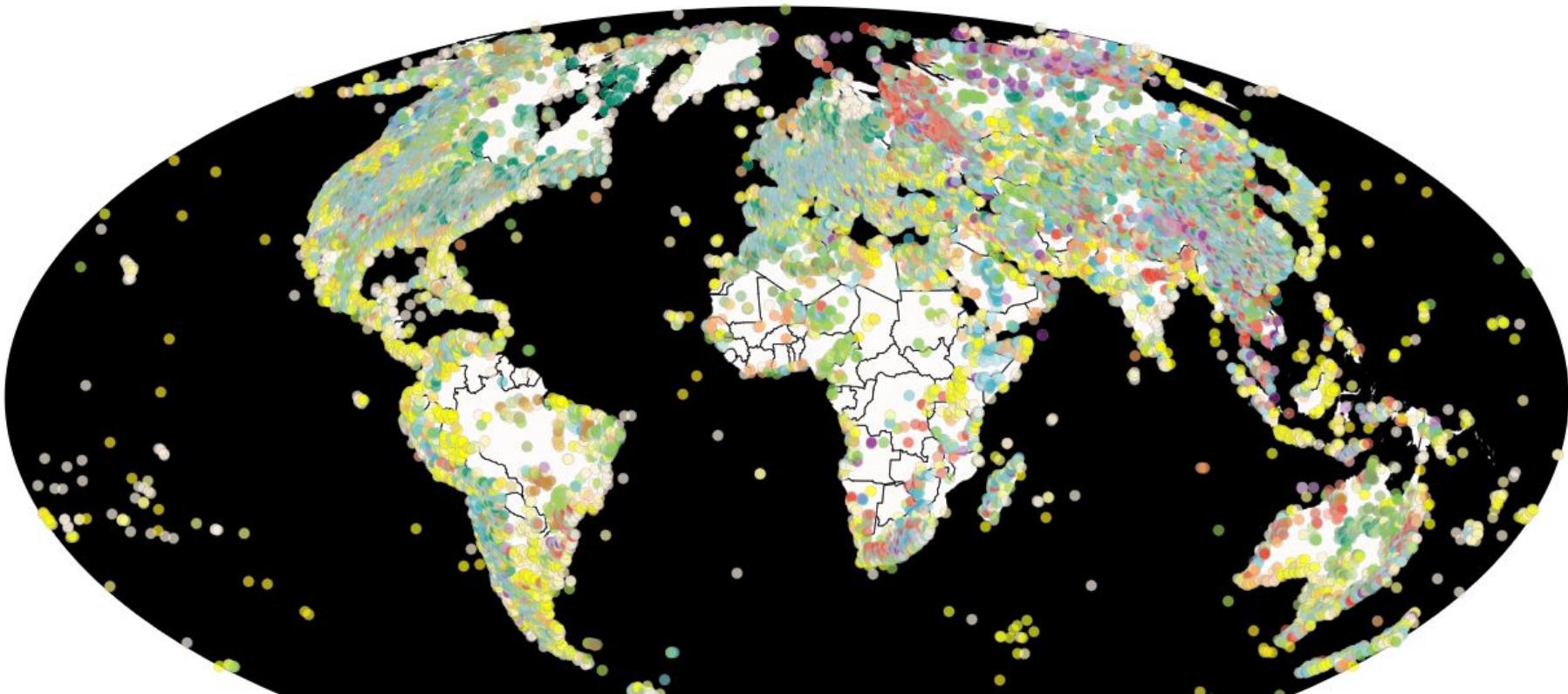
what happens when the entire fossil record is in the  
[@PaleoDB?](#)

3:13 PM · Jan 15, 2020 · [Twitter Web App](#)

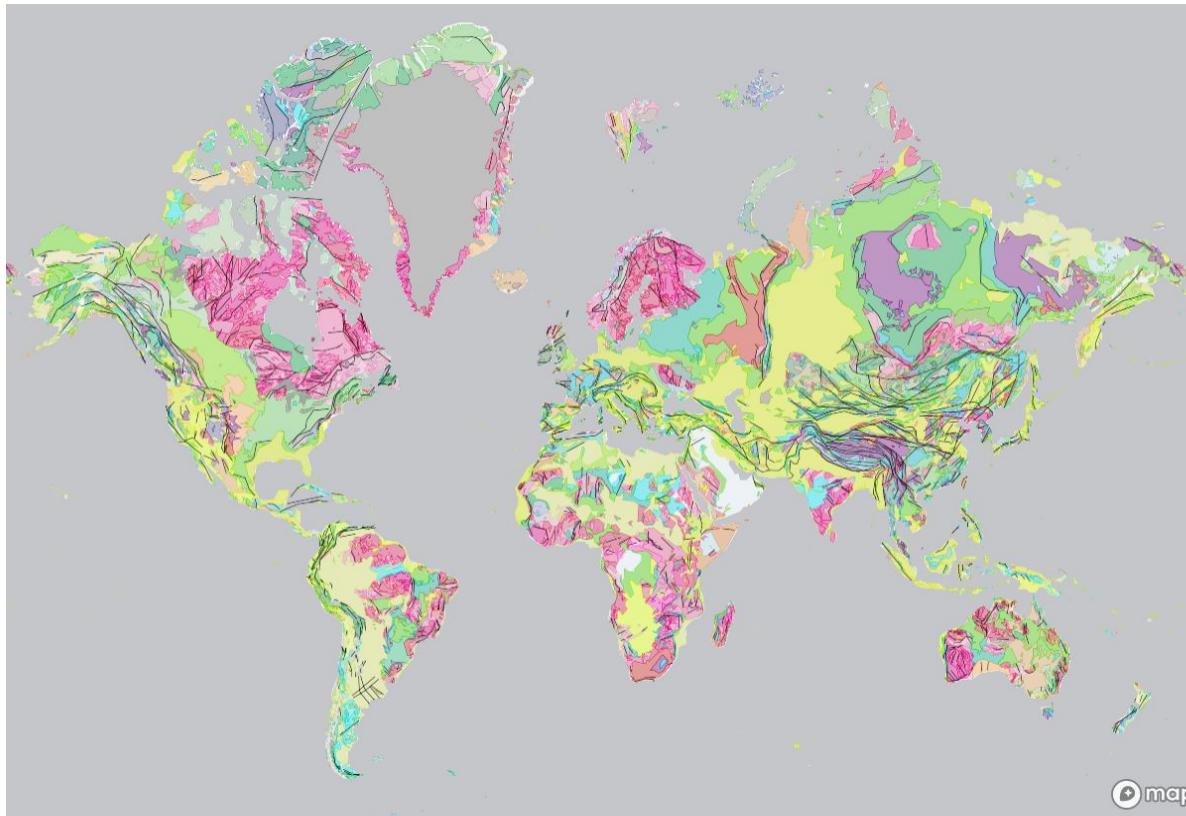
---

3 Retweets 31 Likes

# Paleobiology Database



# Macrostrat



# Help! My data are biased



Pinned Tweet



**PaleobioGoogling**  
@PaleoGoogling



R package to improve poor fossil record

2:46 AM · Apr 1, 2019 · [Twitter Web Client](#)

---

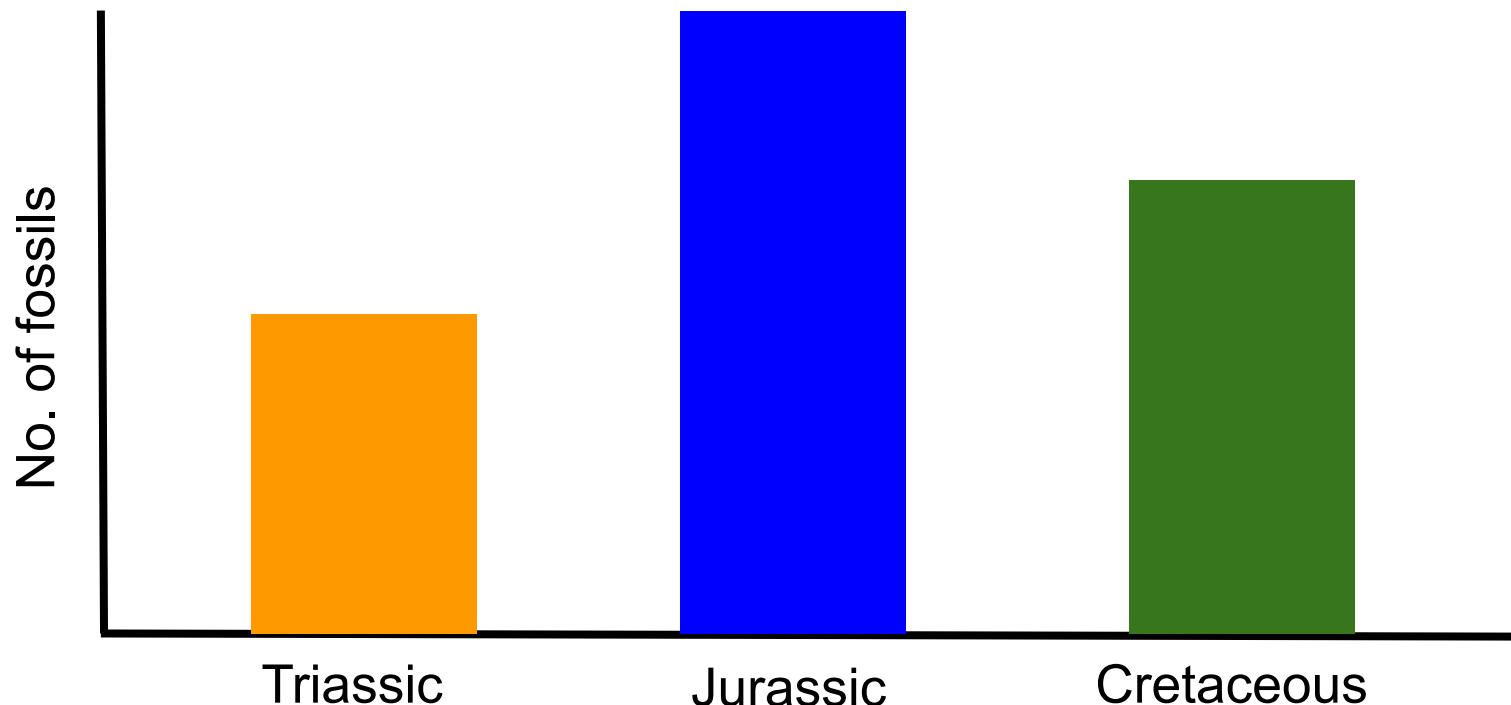
**36** Retweets    **155** Likes

---

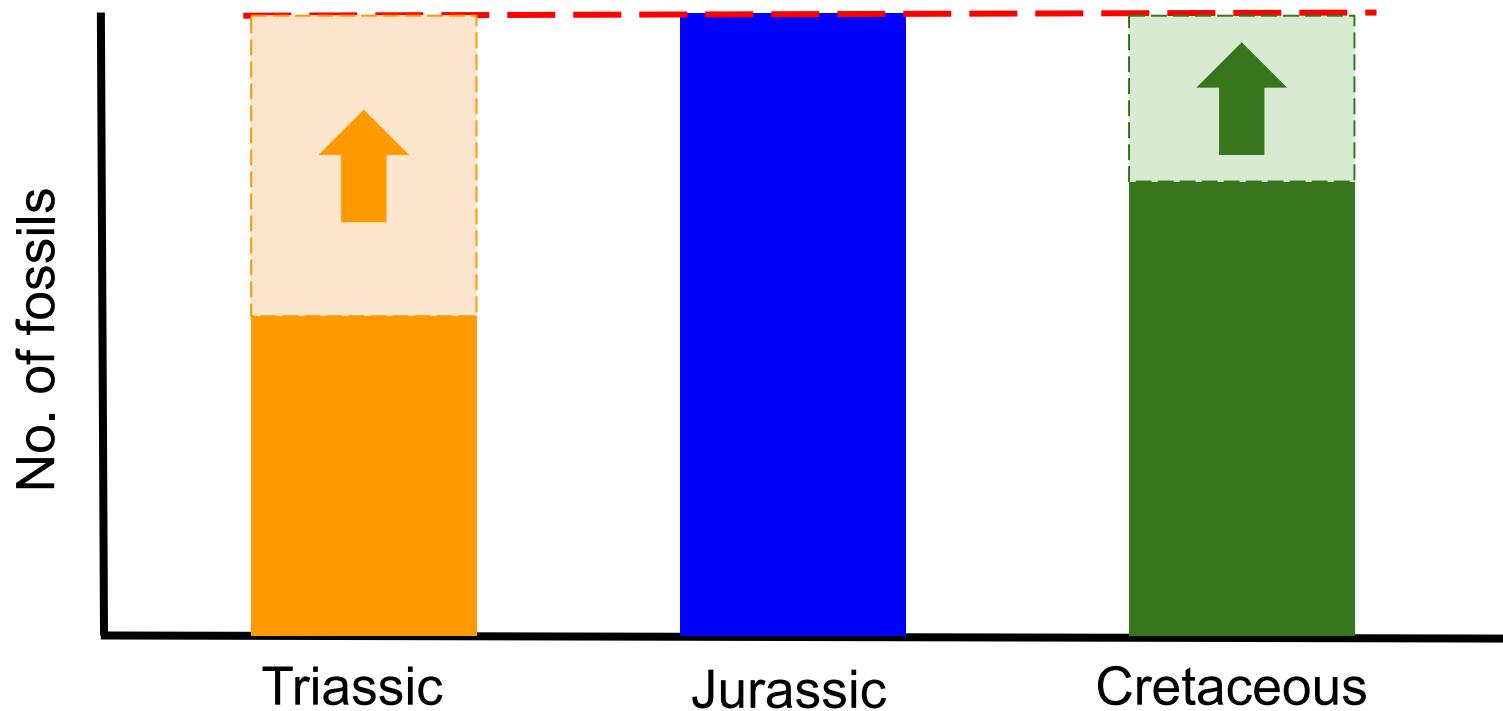
# Methods



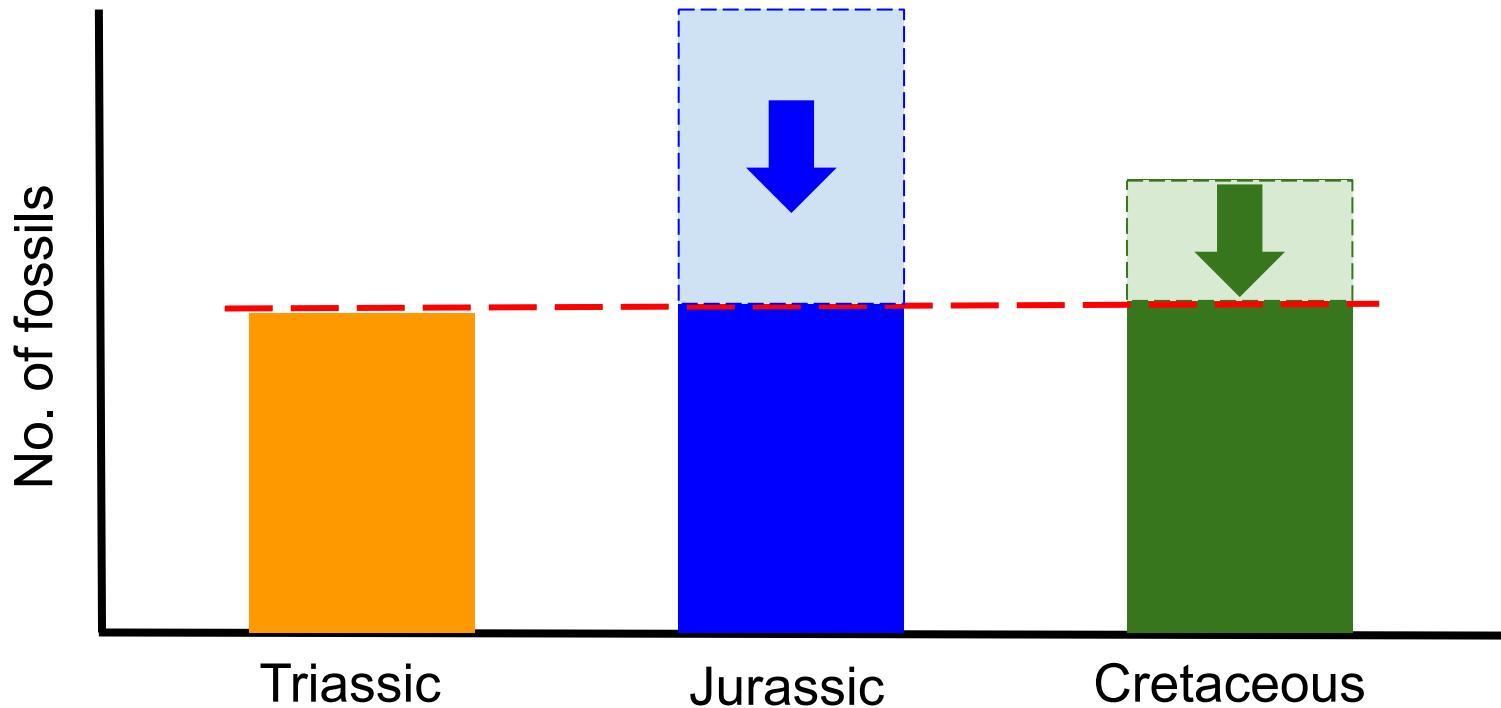
# Keeping it simple?



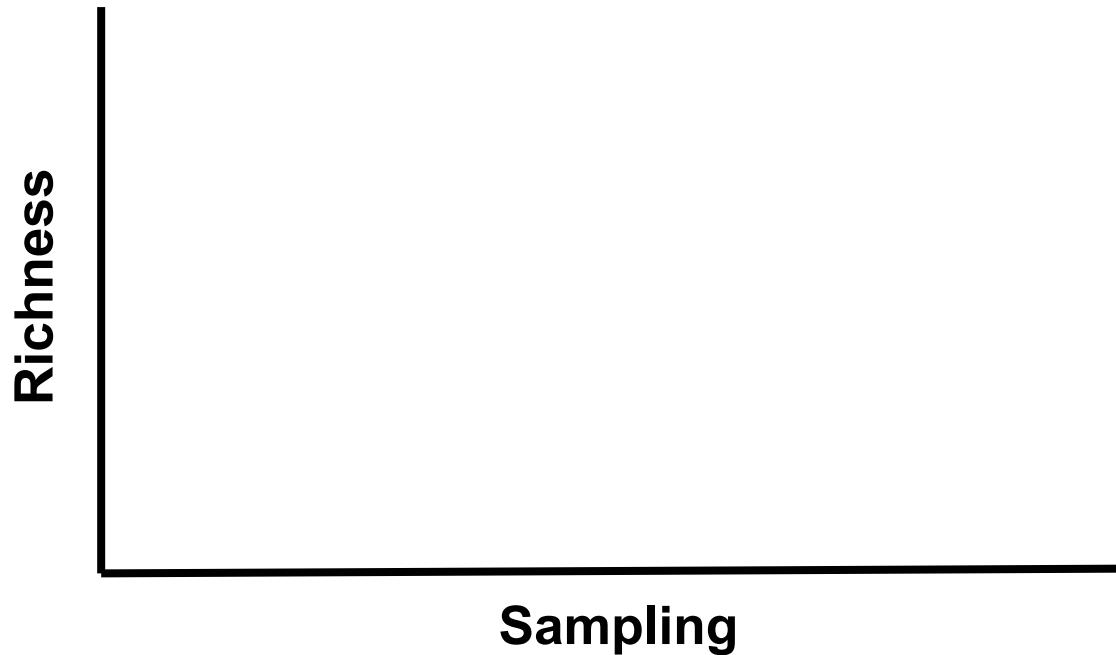
# Extrapolation



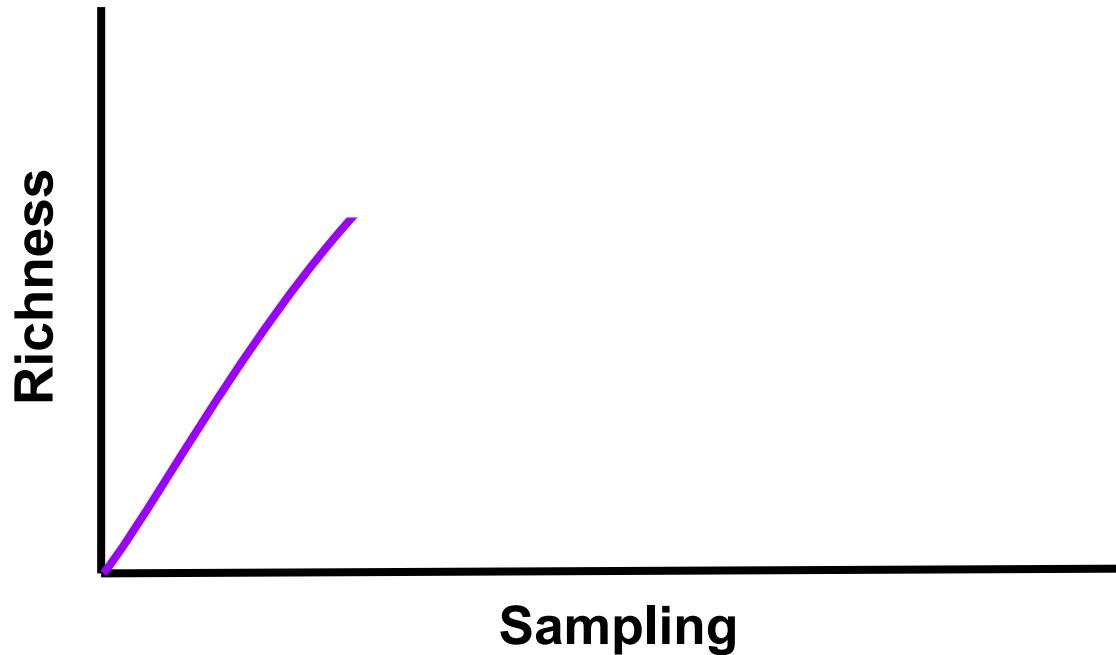
# Subsampling: rarefaction and bootstrapping



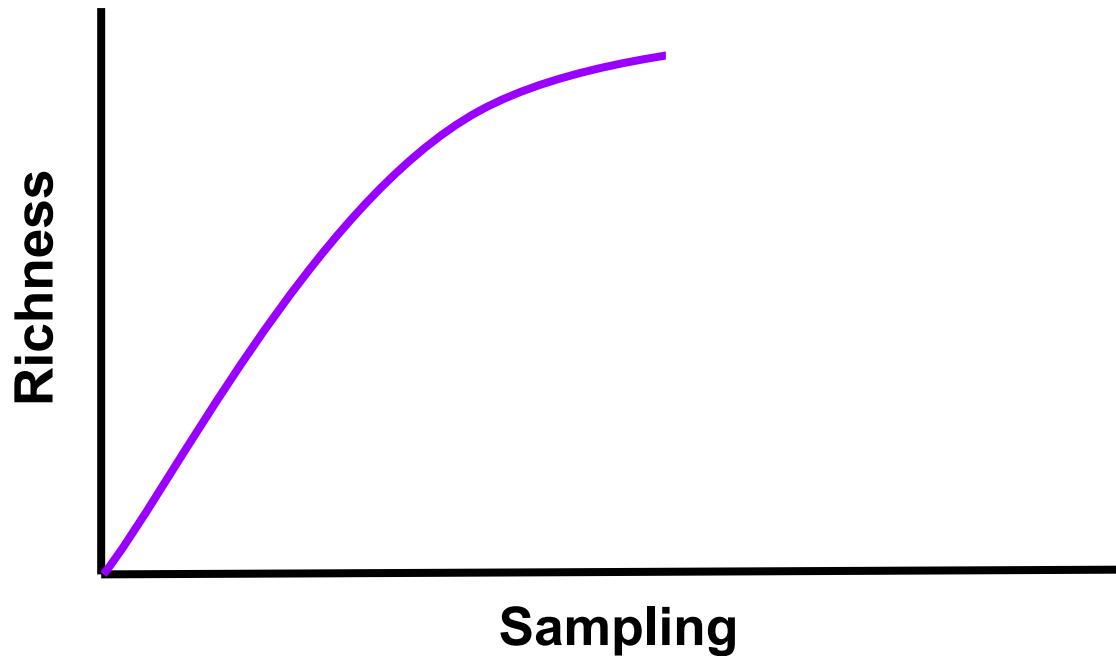
# Collector's / Species Accumulation Curves



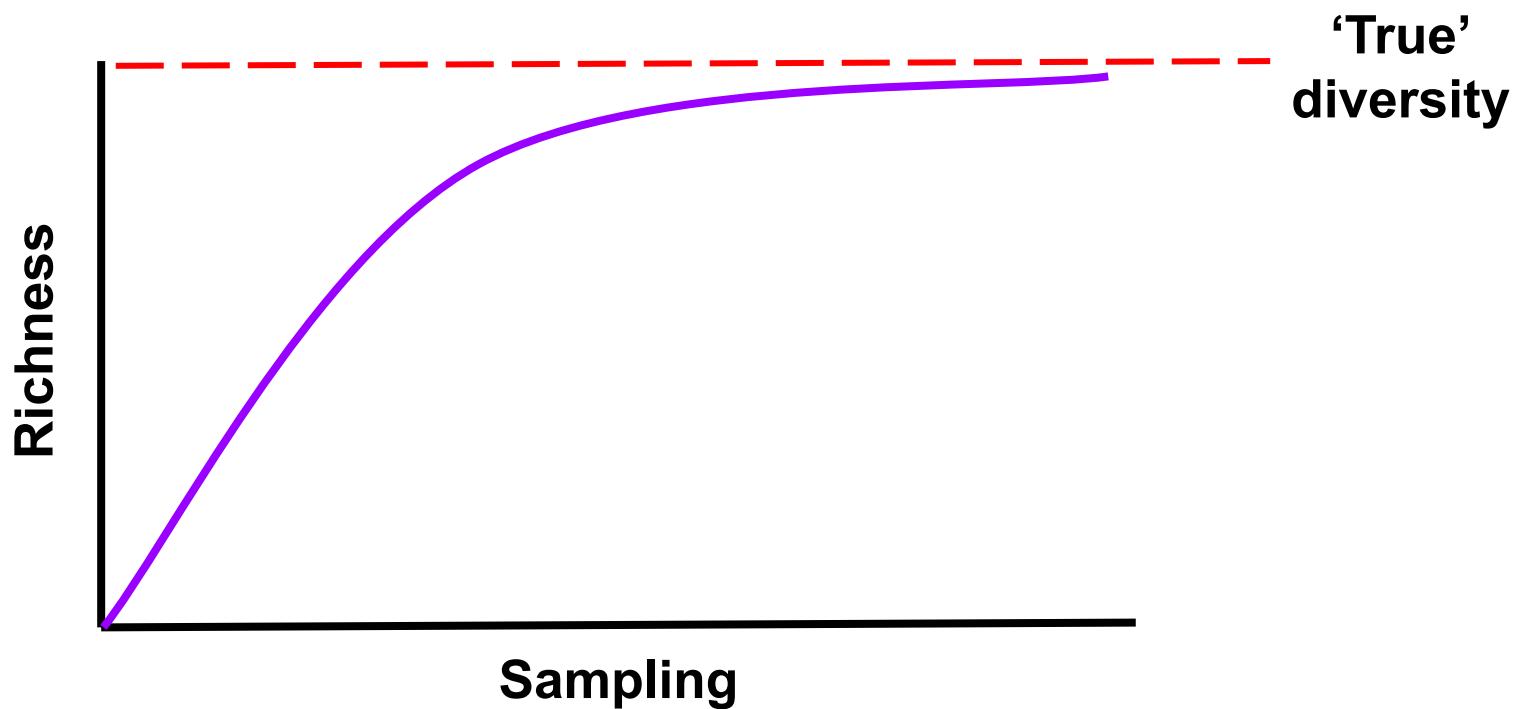
# Collector's / Species Accumulation Curves



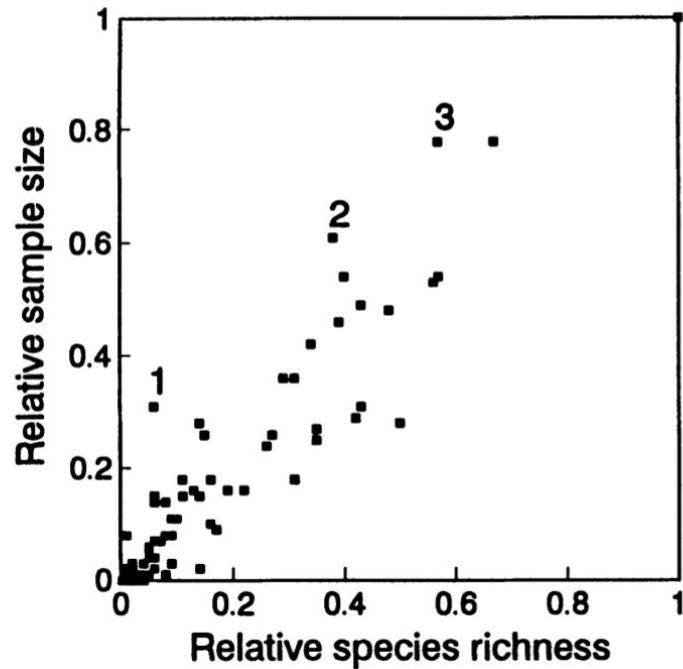
# Collector's / Species Accumulation Curves



# Collector's / Species Accumulation Curves

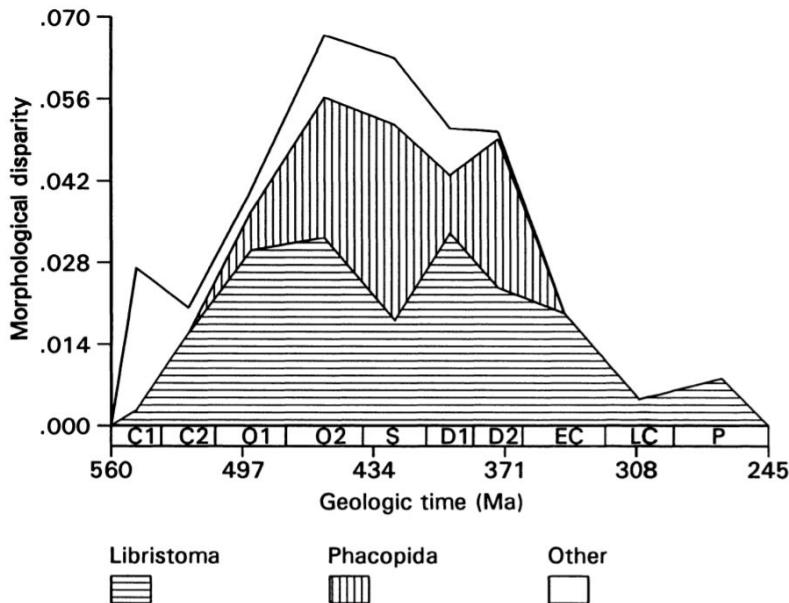


# Bootstrapping example

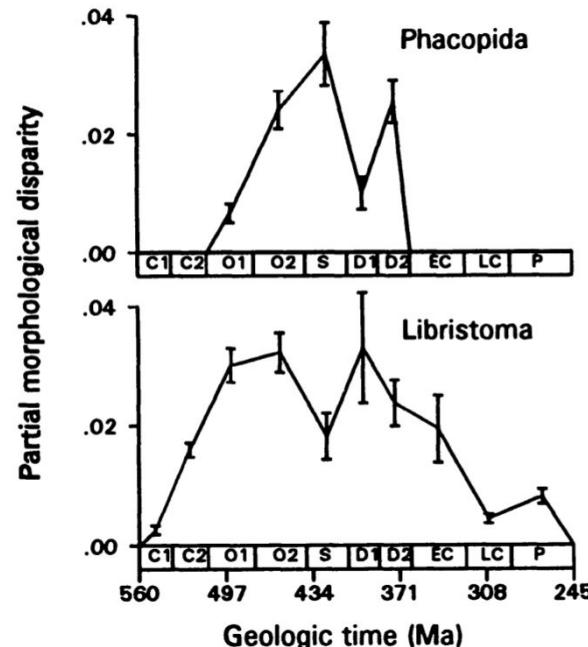


The relationship between sampling and richness in combinations of trilobite subgroups and stratigraphic intervals (Fig. 1)

# Bootstrapping example

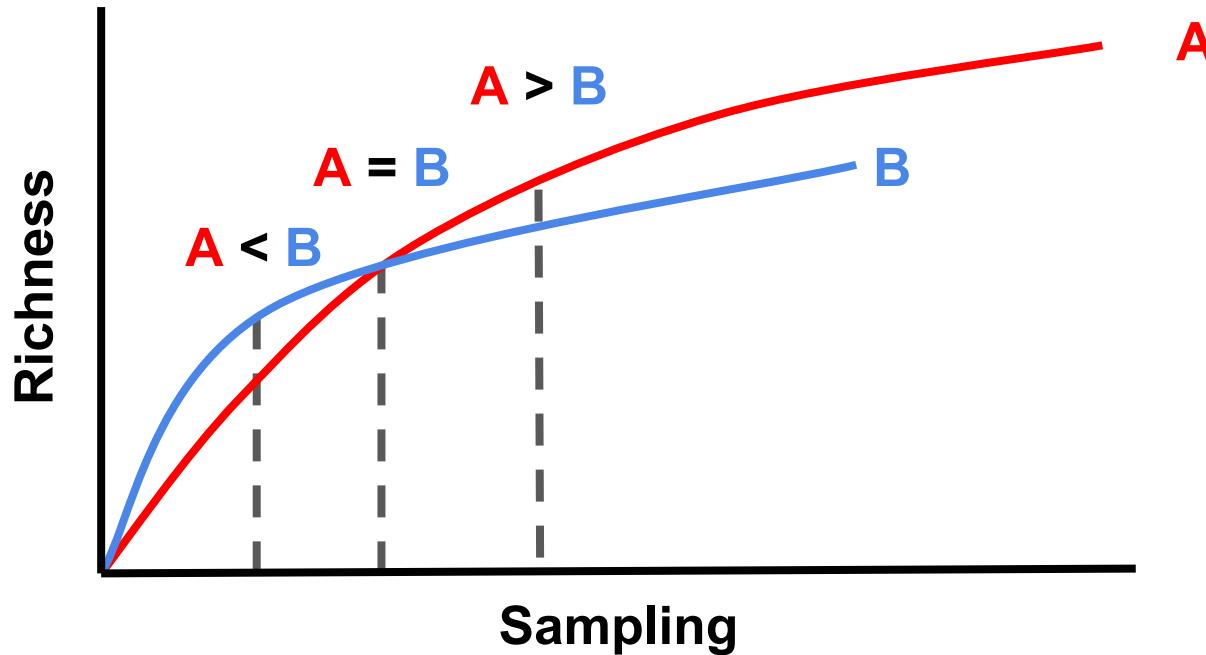


The contribution of different trilobite groups to overall morphological diversity through time (Fig. 2)

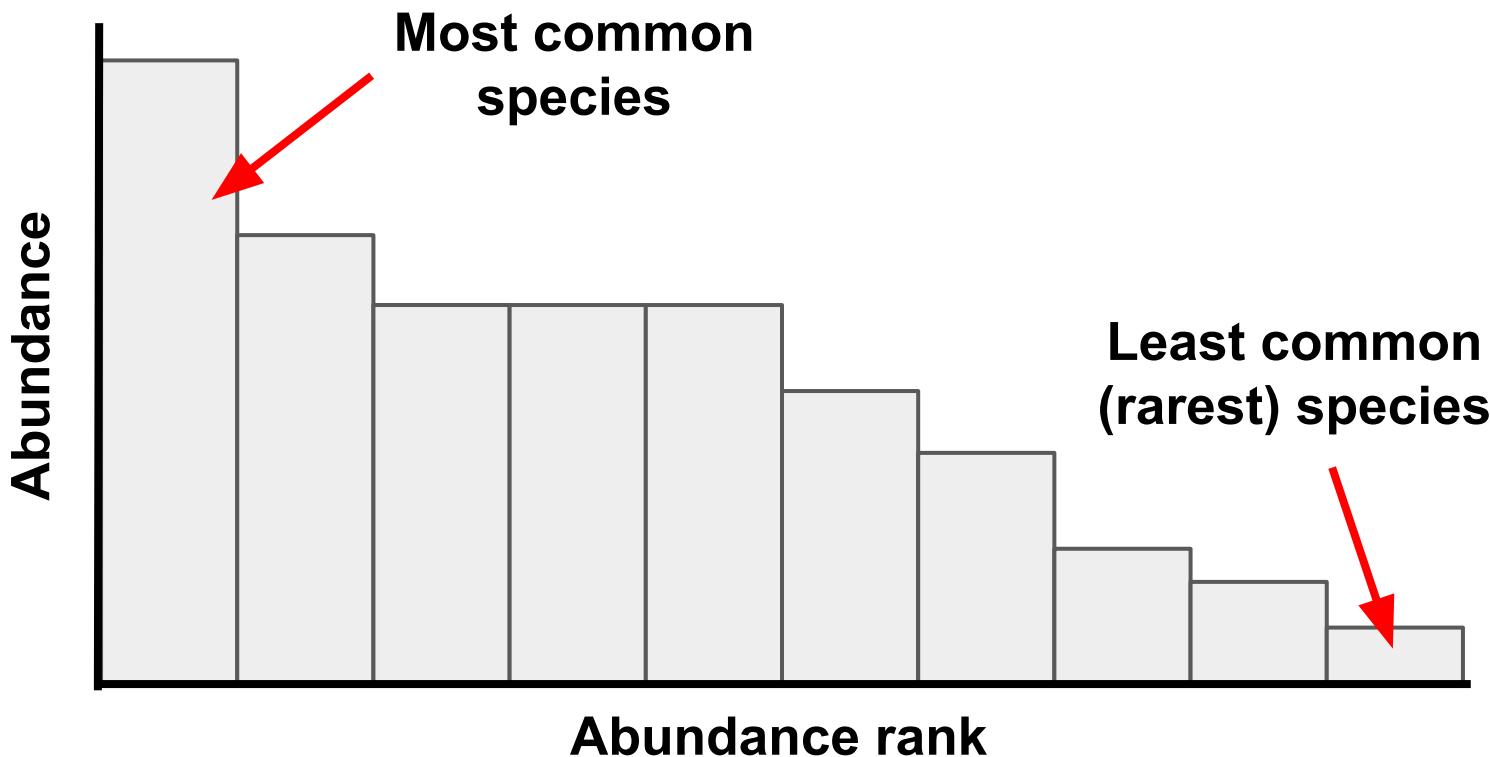


The contributions of Phacopida and Libristoma to overall trilobite disparity, taking into account sampling through a bootstrap regime (Fig. 3)

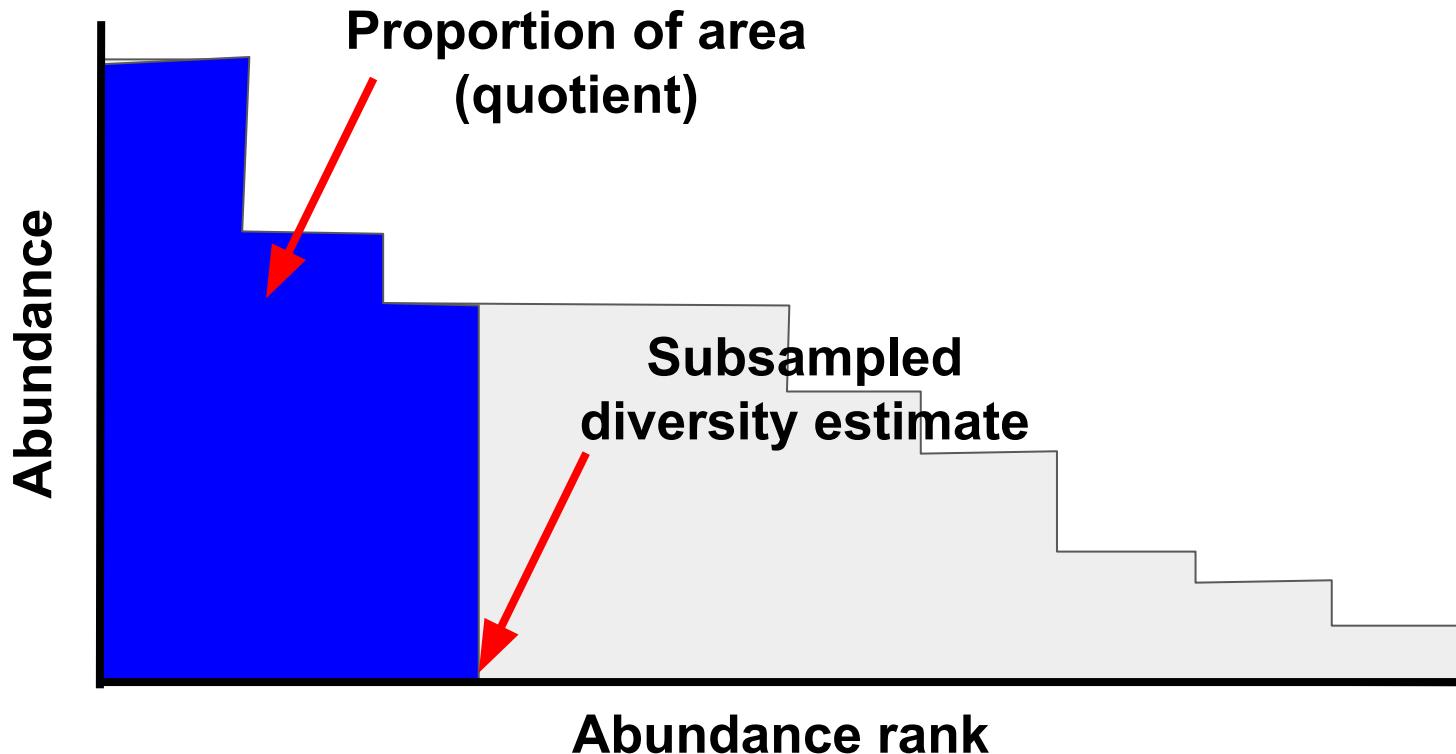
# Subsampling and “The Ghostbuster’s Problem”



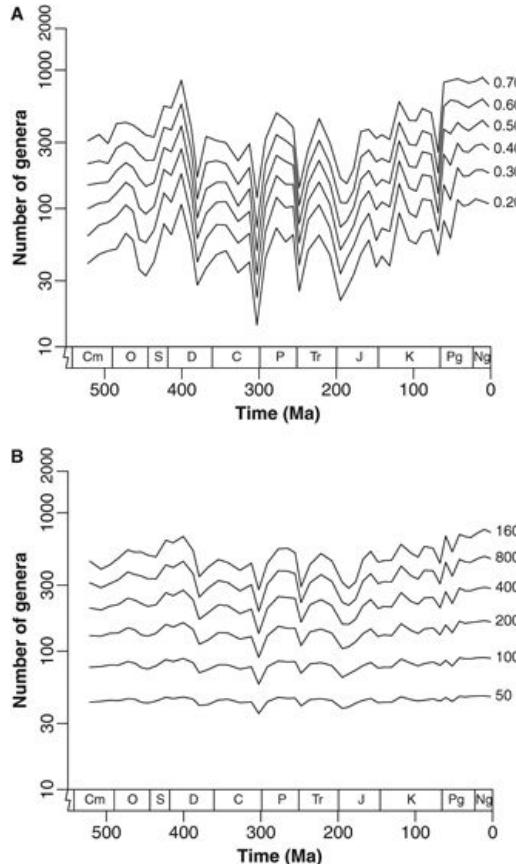
# Rank order abundance



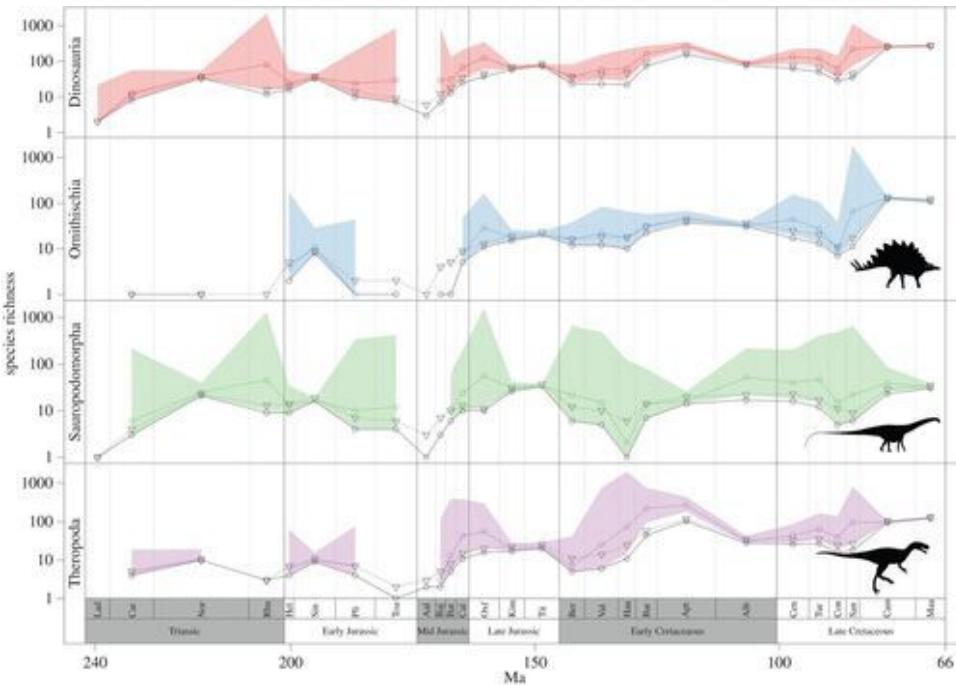
# Coverage-based methods



# Coverage-based methods

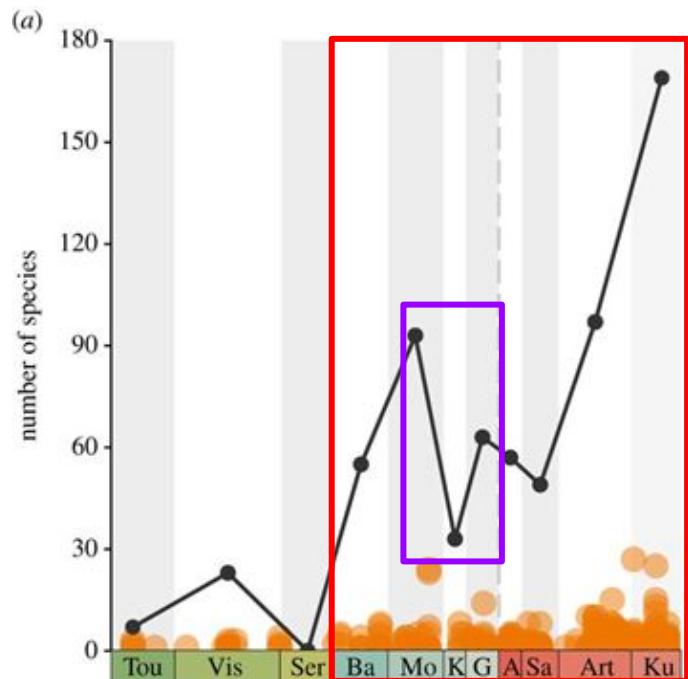


Alroy 2010; *Palaeontology*  
Shareholder Quorum  
Subsampling (SQS)

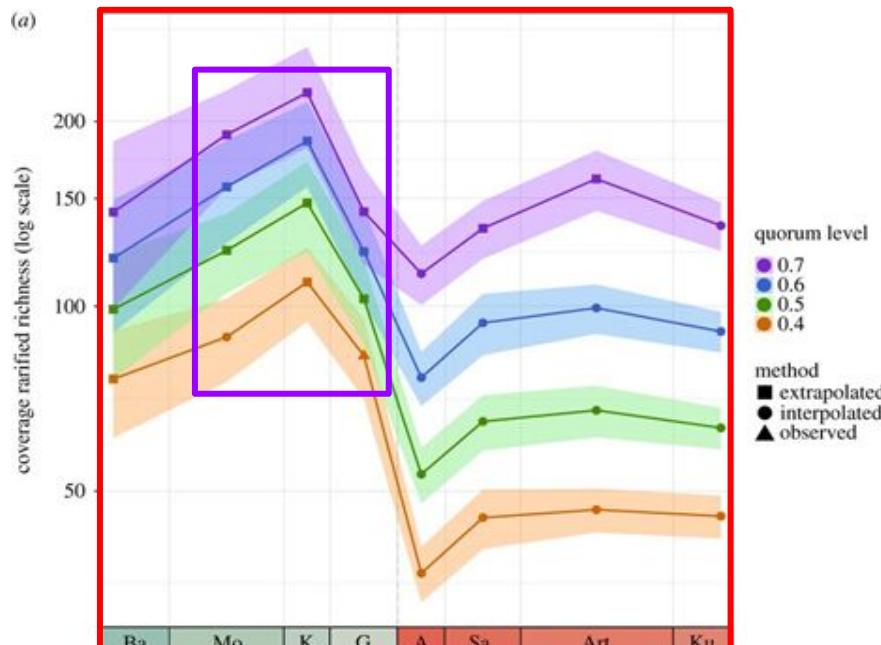


Starrfelt & Liow 2016; *Proc. Roy. Soc. B*  
True Richness estimated using a Poisson  
Sampling model (TriPS)

# Coverage-based examples

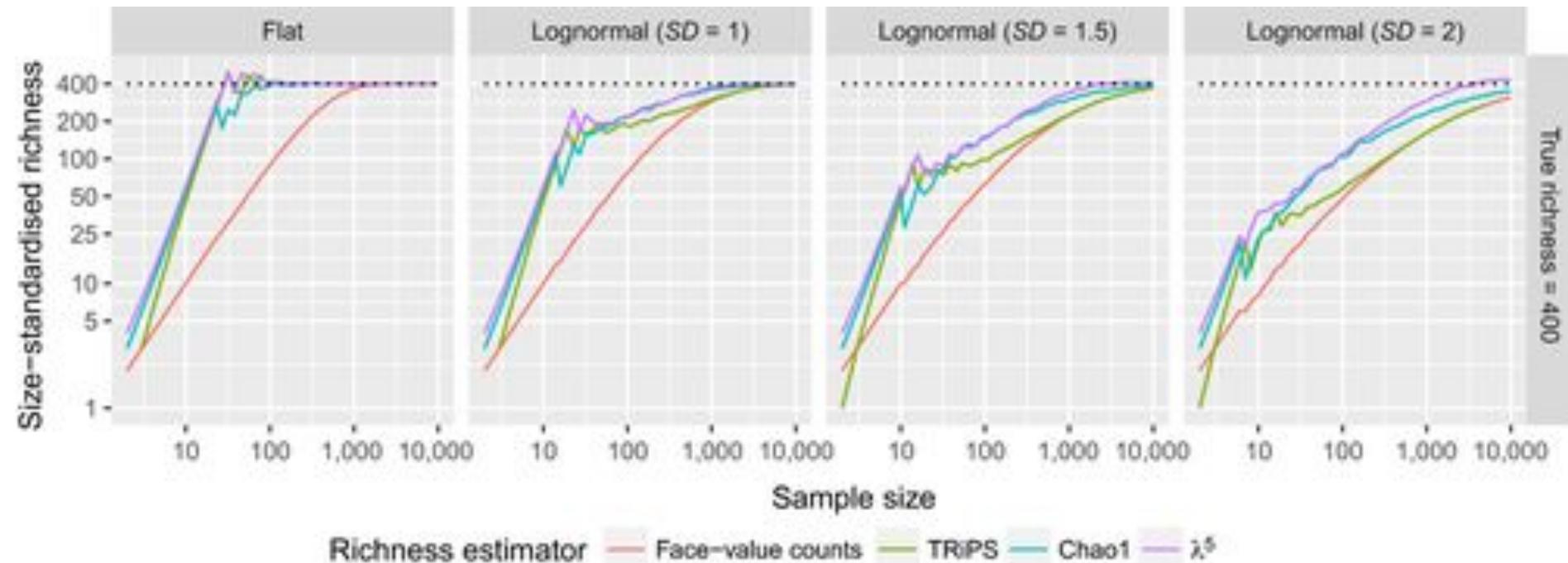


Raw tetrapod diversity curve for the Carboniferous-early Permian (Fig. 1a)

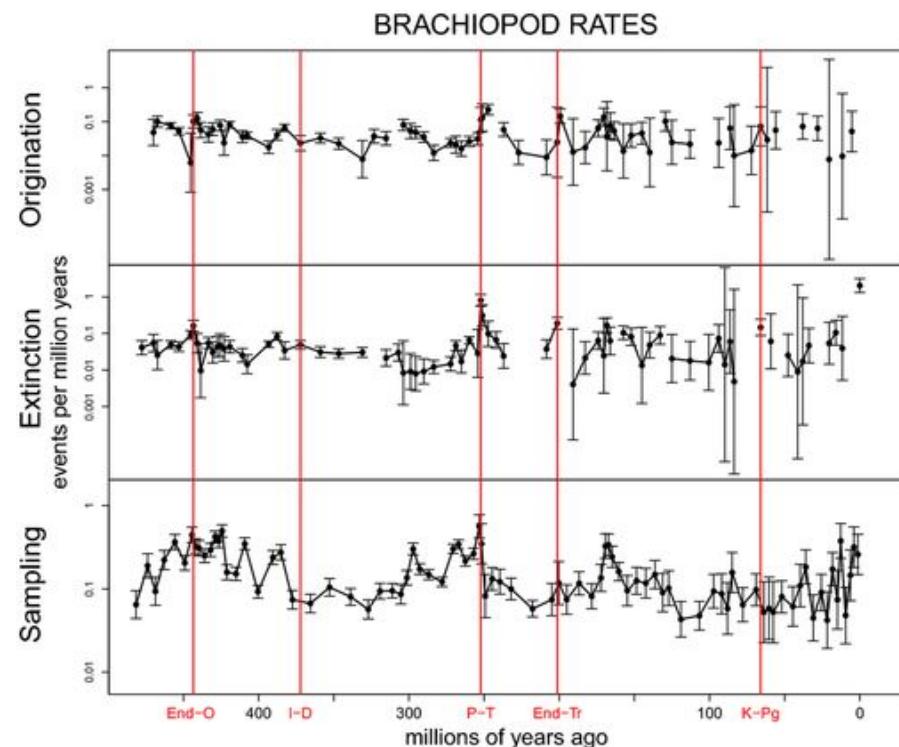
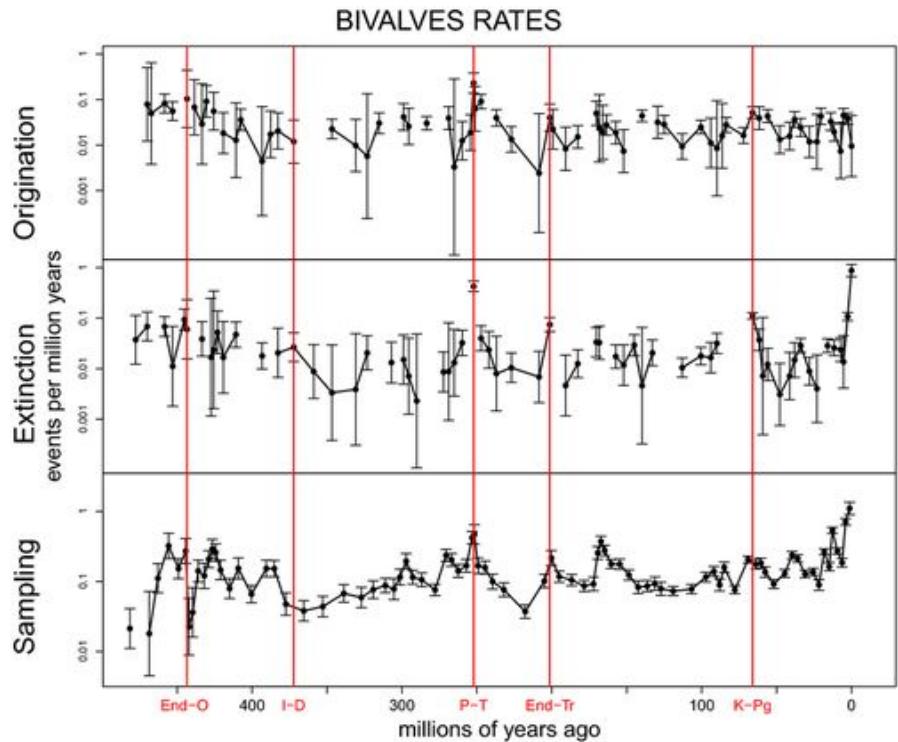


SQS tetrapod diversity curve for the Carboniferous - early Permian (Fig. 3a)

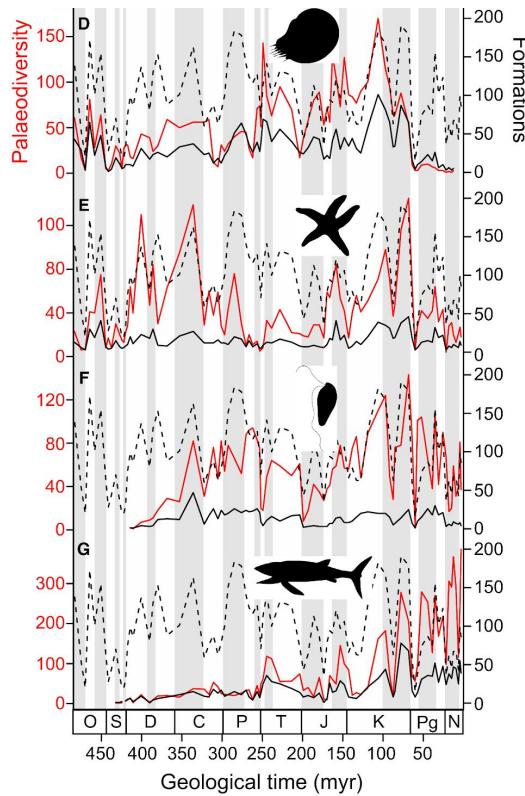
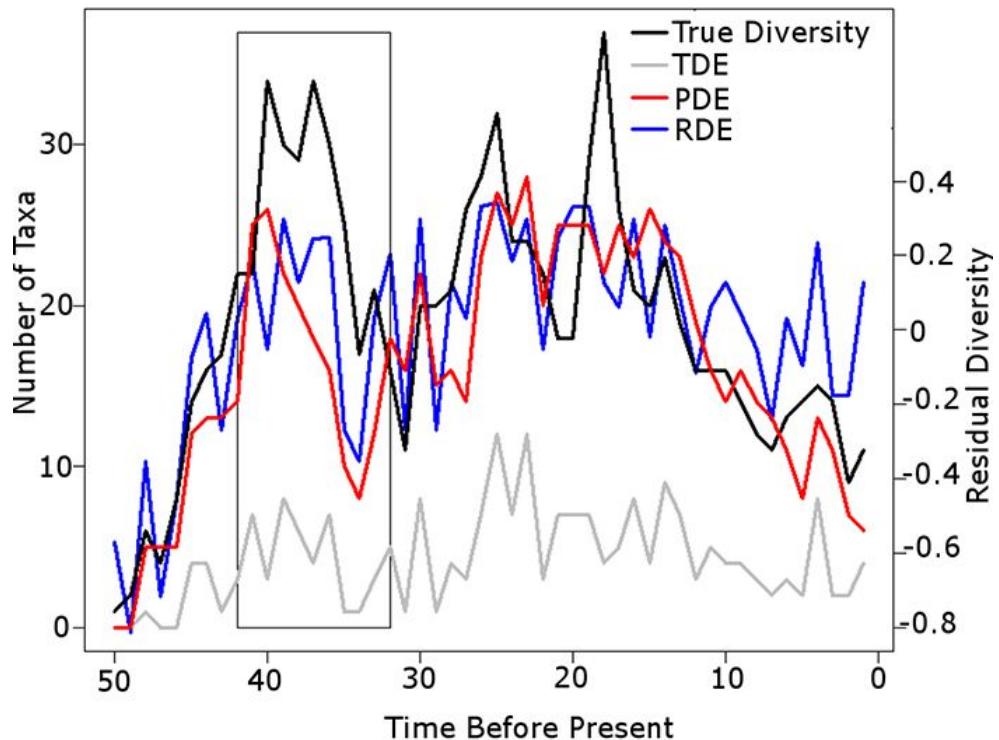
# Which is best?



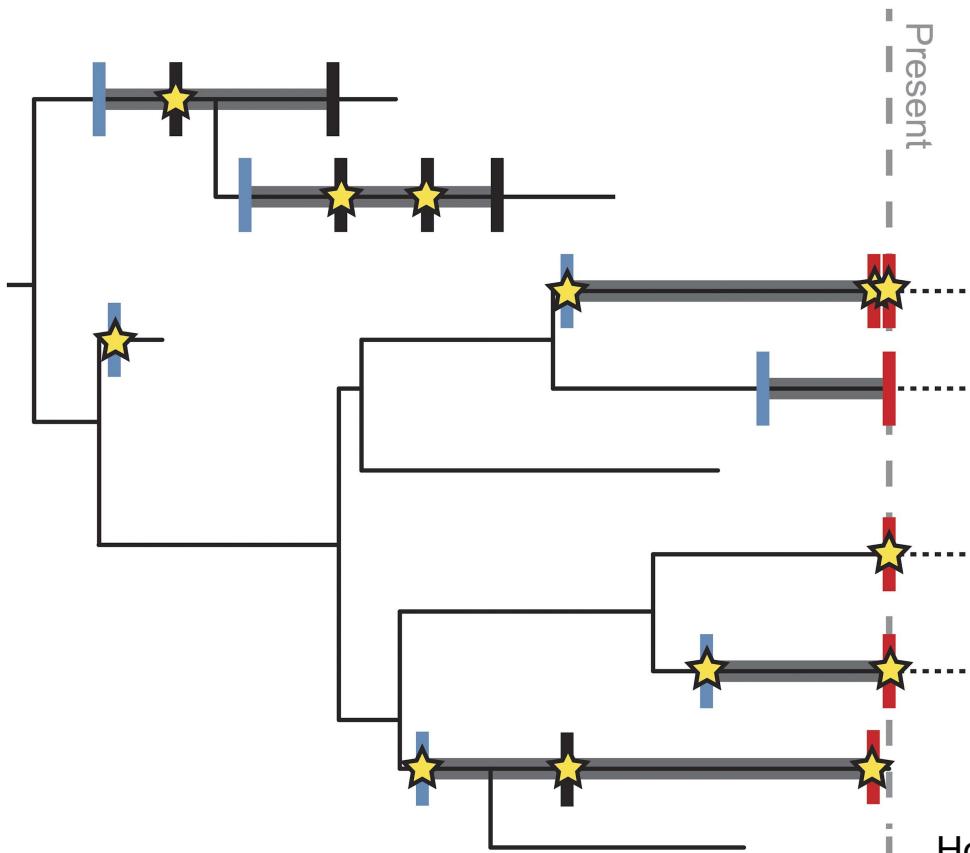
# Capture-mark-recapture



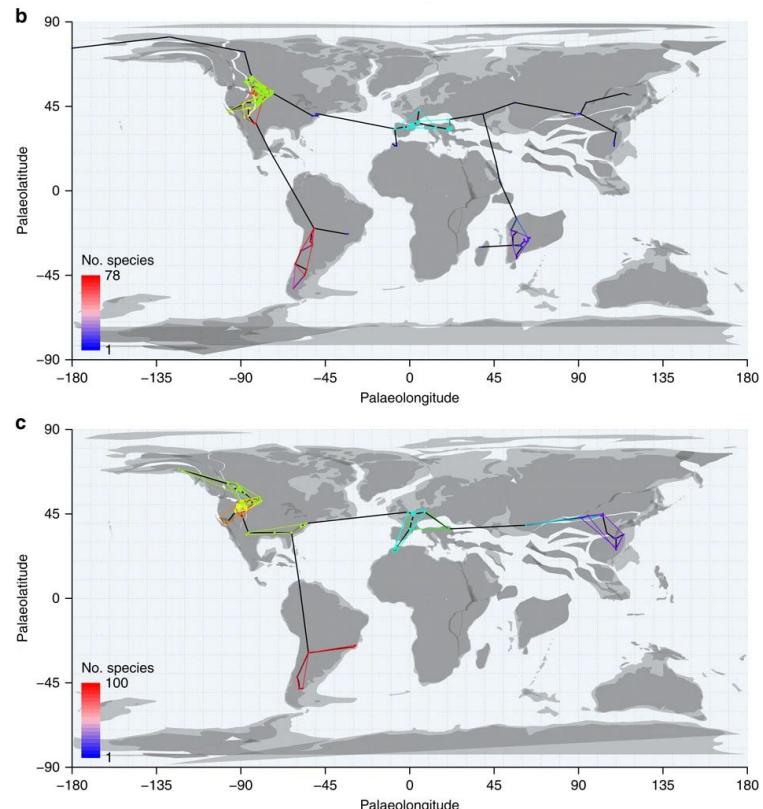
# Sampling proxy correction



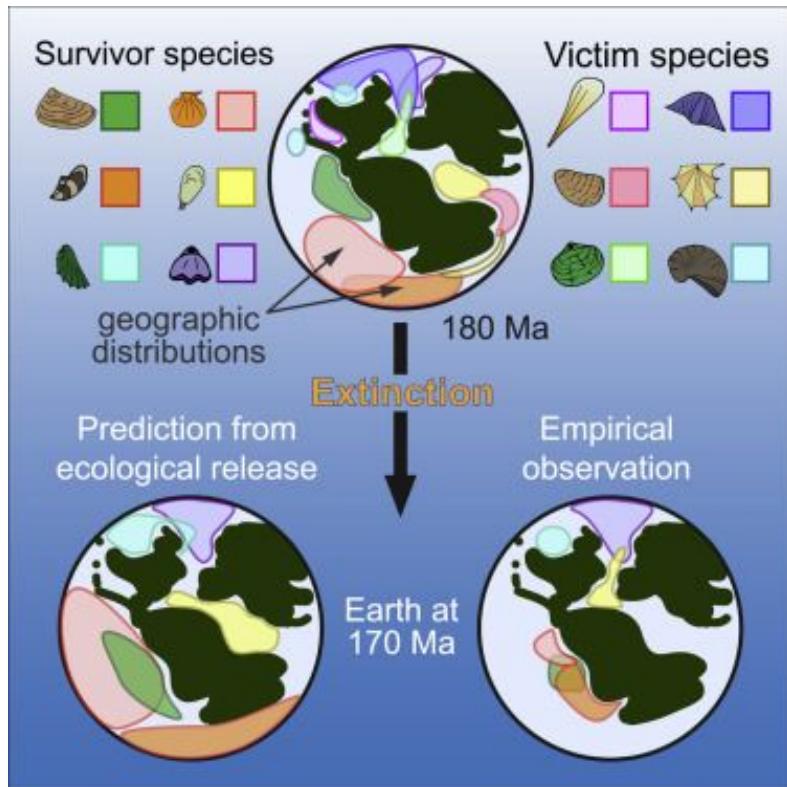
# Who needs fossils when you have trees?



# Space vs time



Close et al. 2017; *Nature Comms.*



Antell et al. 2020; *Curr. Biol.*