January 2003

#### **Big Data in Phylogenetics**

#### Dealing with the deluge

thomsonr@hawaii.edu



(1982 - 2005)



Sequences (millions)

# skipping ~30 figures that all show growth of data or decreasing cost



Thomson and Shaffer 2010

#### PERSPECTIVE

#### Big Data: Astronomical or Genomical?

**Growth of DNA Sequencing** 



Stephens et al. 2015

5

#### To what end?

"We think big data is what everyone cares about. It's not. It's stories."

#### - Dr. Jessica Utts President, American Statistical Association

The goal is to gather 'sufficient' data in order to answer a question 'robustly.'

#### To what end?

"We think big data is what everyone cares about. It's not. It's stories."

#### - Dr. Jessica Utts President, American Statistical Association

The goal is to gather 'sufficient' data in order to answer a question 'robustly.' The question is what is interesting.

#### To what end?

"We think big data is what everyone cares about. It's not. It's stories."

#### - Dr. Jessica Utts President, American Statistical Association

The goal is to gather 'sufficient' data in order to answer a question 'robustly.' The question is what is interesting. This is no different than it's always been.

### A case study

- A very specific question:
  - What are the phylogenetic affinities of turtles?

- Brings up more general issues:
  - How do we approach difficult phylogenetic problems?
  - How **should** we approach difficult phylogenetic problems?

# Turtle Phylogenetics

- Overarching problem:
  - Where do turtles sit in the amniote tree?



Early approaches relied on osteology (primarily of the skull)
Anapsid
Diapsid
Synapsid





• Early approaches relied on osteology (primarily of the skull)



#### • Primary issue with this hypothesis



#### More osteology



#### Reippel and deBraga 1996 Nature

#### Molecular Information



#### Molecular Information

• Nuclear data Iwabe et al. 2004 MBE Hugall et al. 2007 Syst Biol

Hedges and Polling 1998 Science

## Summary



#### Turtle Genomics

- 3 genome consortia
- Several more independent studies





National Human Genome Research Institute





## Phylogenomics



#### Shaffer et al. 2013 Genome Biol

# Phylogenomics

• All analyses agree!



#### MicroRNA Result



Lyson et al. 2011 Biol Lett

### Summary

• Ugh...so what do we do?



### Summary

• Ugh...so what do we do?



# Data in Phylogenetics

- Let's take a step back.
- How have we been approaching this (and most other) phylogenetic questions?

4 nuclear genes



11 nuclear genes



Hedges and Polling 1998 Science

# Data in Phylogenetics

- Let's take a step back.
- How have we been approaching this (and most other) phylogenetic questions?



11 nuclear genes



Hedges and Polling 1998 Science

# Phylogenomics

 $\circ~$  Inferences result from  $\underline{both}$  data and the model







Kumar et al. 2012 MBE









 In developing a statistical model for a problem, we inevitably make a tradeoff



1 gene

10 genes



• The point.



• The point.



#### How do we know it's a bigger problem?



#### How do we know it's a bigger problem?



#### Where's the disagreement coming from?

#### How do we know it's a bigger problem?



# 'Big data' turtle studies

- Chiari et al. (2012)
  - 248 transcriptomic loci
  - 12 taxa
- Crawford et al. (2012)
  - 1,145 UCEs
  - 10 taxa
- Fong et al. (2012)
  - 75 Sanger-sequenced loci
  - 129 taxa
- Lu et al. (2013)
  - 1,638 transcriptomic and genomic loci
  - 11 taxa

- Shaffer et al. (2013)
  - 1,955 genomic loci
  - 8 taxa
- Wang et al. (2013)
  - 1,113 genomic loci
  - 12 taxa

#### **Bipartition Bayes Factors**









1/1,000,000,000,000,000,000,000,000

That's 27 zeroes!

If you played a lottery every minute with that chance of winning, you still probably wouldn't win, unless you played for...

the age of the universe\*190,258,751,903









Brown and Thomson 2017



Equivocation about turtle placement across genes



This dataset supports turtles as sister to crocodilians. But what's up with these outliers? How influential are they?

## Both look like paralogs



Brown and Thomson 2017 Bodega2017 BigDataDiscussion Thomson.key - March 14, 2017

# Strong influence



# Strong influence



# Strong influence



#### Take homes

- More data does not necessarily lead to more accuracy, or to consensus
- A lot of phylogenomic **progress** is actually about figuring out how to **model data well**, not collect more data per se

#### Some Possible Ways Forward

• **Embrace** the computational **challenge** 

### Embrace the computation

- Analyses need not finish quickly
- Advances in computation help a lot here
  - parallel architectures and code
  - fast computation libraries
  - availability of compute resources
  - $\circ\,$  new methods on the horizon (HMC, IDR)

#### Embrace the computation



#### Embrace the computation

#### Compute Resources



Name	Status	CPUs	Peak TFlops	Utilization	Running Jobs	Queued Jobs	Other Jobs
Stampede 🖋 User Guide	✓ Healthy	102400	9600.0	61%	562	2310	161
Comet 🛢 🗲 User Guide	✓ Healthy	47616	2000.0	87%	1487	360	379
XStream <b>■</b> ⊁ User Guide	✓ Healthy	1300	1001.7	82%	303	85	415

### XSEDE jobs by field for 2016



2016-01-01 to 2016-12-31 Src: XDCDB. Powered by XDMoD/Highcharts

#### XSEDE users by field for 2016

1. Materials Research 379 2. Biophysics 453 665 3. Advanced Scientific Computing 314 4. Chemistry 263 5. Computer and Computation Research 239 6. Training 215 7. Biochemistry and Molecular Structur... 201 8. Computer and Information Science an... 9. Astronomical Sciences 10. Physical Chemistry 11. Avg of 110 others 0 50 300 350 400 450 500 800 8... 100 150 200 250 550 600 650 700 750 Number of Users

Number of Users: Active: by Field of Science

Systematic and Population Biology

Number of Users: Active

#### New Tools on the Horizon

- $\circ~$  More complex models
- faster estimation of marginal likelihoods. e.g., Inflated Density Ratio
- More efficient sampling. e.g., Hamiltonian Monte Carlo

#### A Comparison of two Markov Chain Monte Carlo samplers

Tamara Broderick (UC Berkeley), David Duvenaud (Cambridge)

#### Some Possible Ways Forward

- **Embrace** the computational **challenge**
- **Get very picky** about our data. Careful and detailed data exploration is your friend.

#### Some Possible Ways Forward

- Embrace the computational challenge
- **Get very picky** about our data. Careful and detailed data exploration is your friend.
- **Carefully consider tradeoffs** between speed and approximation



Leaché and Rannala 2010



Brown and Thomson unpublished figure