

Introduction to molecular dating methods

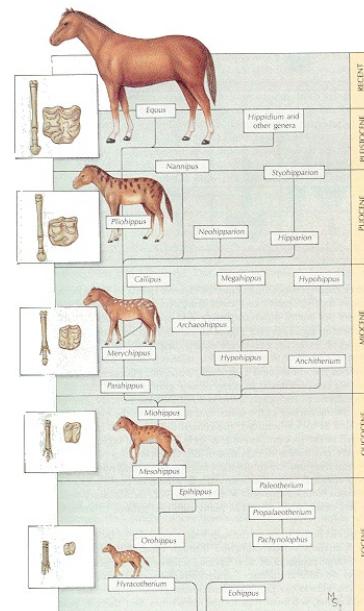
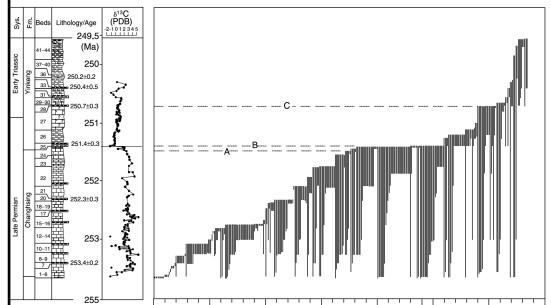
Niklas Wahlberg



LUND
UNIVERSITY

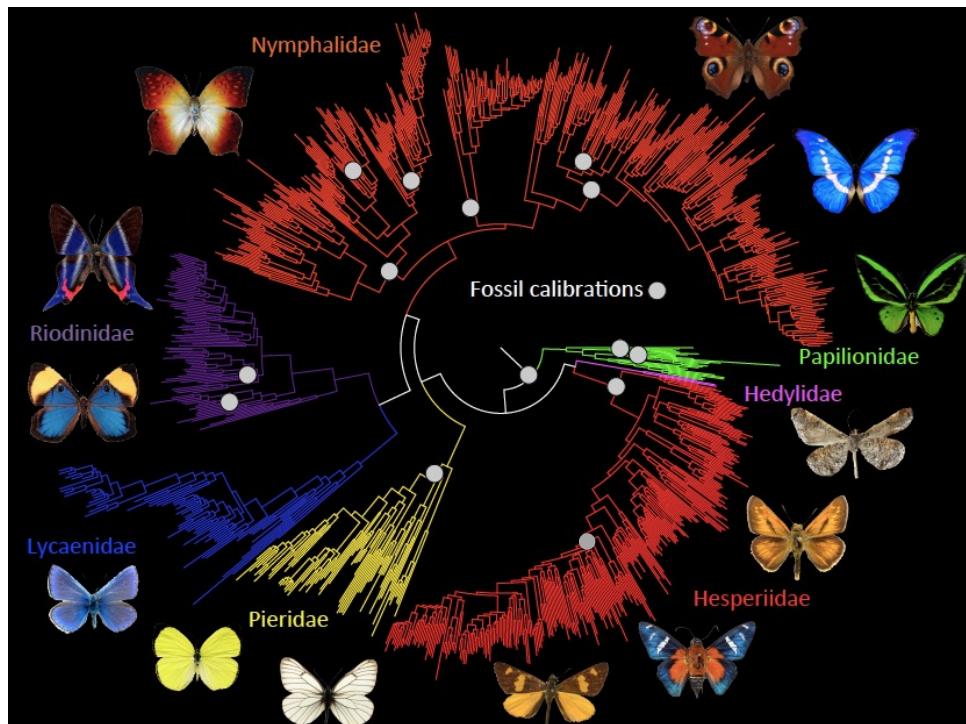
1

The fossil record is the direct evidence of past events and the time at which they occurred

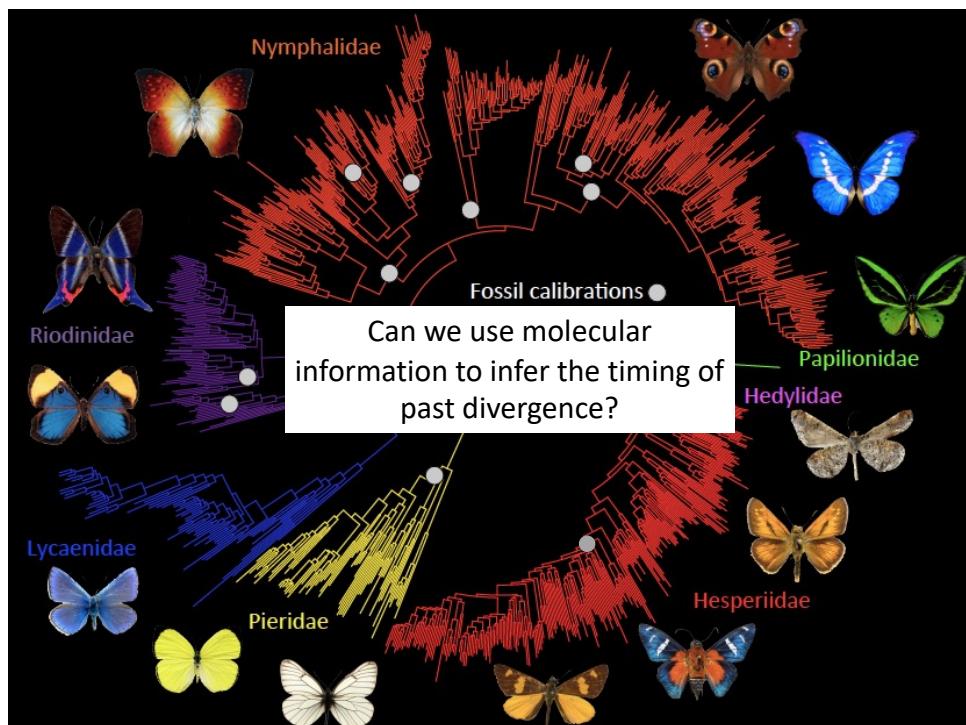


2

1



3



4

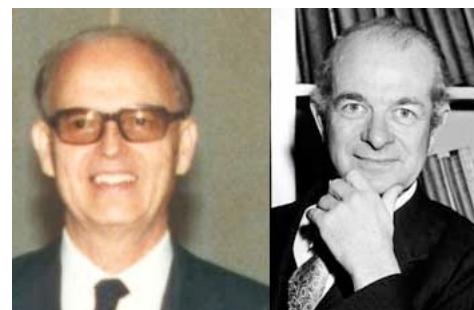
The Molecular Clock

Going back to ancient times

5

Is there a molecular clock?

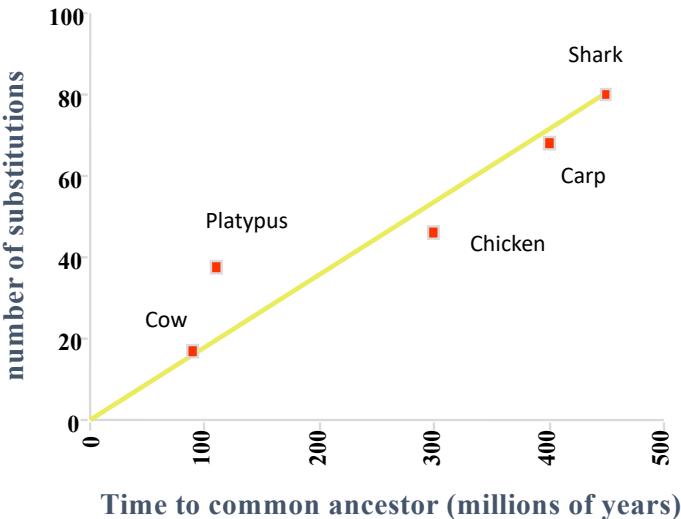
- The idea of a molecular clock was initially suggested by Zuckerkandl and Pauling in 1962 and 1965



6

3

The molecular clock for alpha-globin:
 Each point represents the number of substitutions separating
 each animal from humans



7

Is there a molecular clock?

- The idea of a molecular clock was initially suggested by Zuckerkandl and Pauling in 1962 and 1965
- They noted that rates of amino acid replacements in animal haemoglobins were roughly proportional to time - as judged against the fossil record

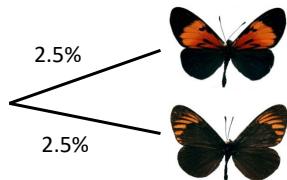
=> implies the existence of a sort of molecular clock ticking faster or slower for different genes but at a more or less constant rate for a genes among different lineages

8

4

The molecular clock hypothesis

- Assumes an equal rate of molecular evolution over time



- A 5% difference between species means they have each diverged 2.5% since their common ancestor
- If a fossil or other evidence will let us calibrate this clock we can convert % difference to years

9

Assumptions of a perfect clock

- Molecular change is a linear function of time with substitutions accumulating following a Poisson distribution - any variation will be stochastic [imagine 1 substitution / million yrs]
- Rate of change is equal across all sites and lineages
- The phylogeny can be estimated without error
- The number of substitutions along each lineage can be estimated without error
- Calibration dates for all times of divergence used to calculate the rate of the molecular clock are known without error
- A regression of time on number of substitutions can be conducted without error

10

Dating with a molecular clock

- “Universal Molecular Clocks”
- Calibrations proposed for various taxa / genes
- eg. mtDNA molecular clock of animals
 - ~ 2% sequence divergence per million years for vertebrates
 - ~ 1% sequence divergence per million years for invertebrates



11

There is no universal molecular clock

- The initial proposal saw the clock as a Poisson process with a constant rate
- Now known to be more complex - differences in rates occur for:
 - different sites in a molecule
 - different genes
 - different regions of genomes
 - different genomes in the same cell
 - different taxonomic groups for the same gene
- **There is no universal molecular clock**

12

Challenges

- Saturation

Ancest GGC~~G~~CG

Seq 1 AGCG~~A~~G

Seq 2 GC~~G~~GA~~C~~

Number of changes

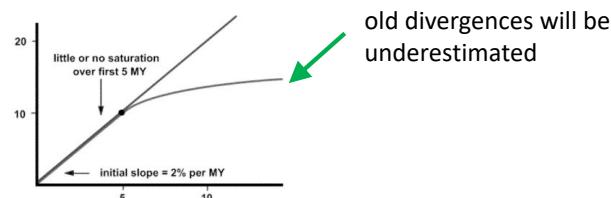
1 2 3

Seq 1 C → G → T → A

Seq 2 C → → → A
 1

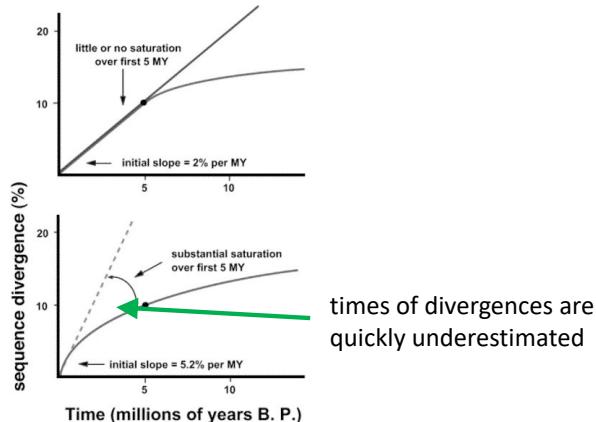
13

Saturation problems



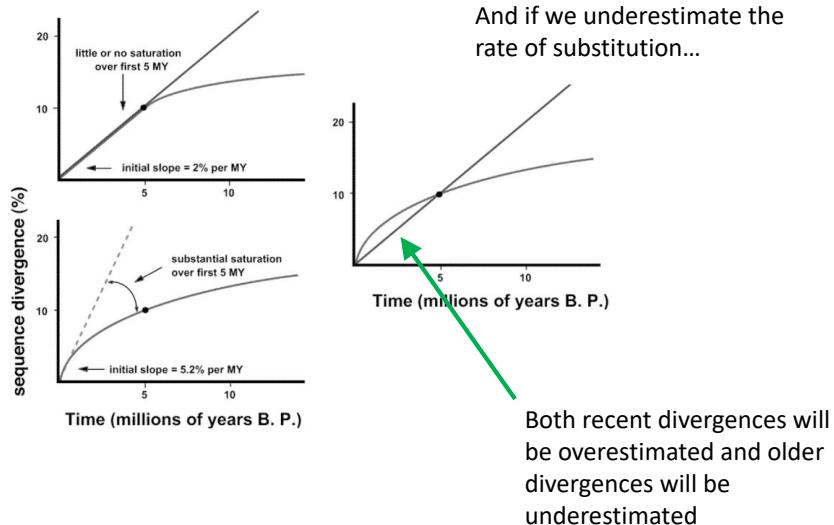
14

Saturation problems



15

Saturation problems



16

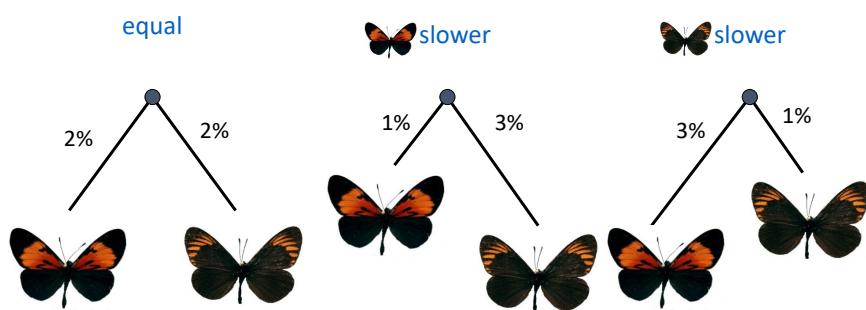
Challenges

- Saturation
- Rate Heterogeneity - violation of homogeneity



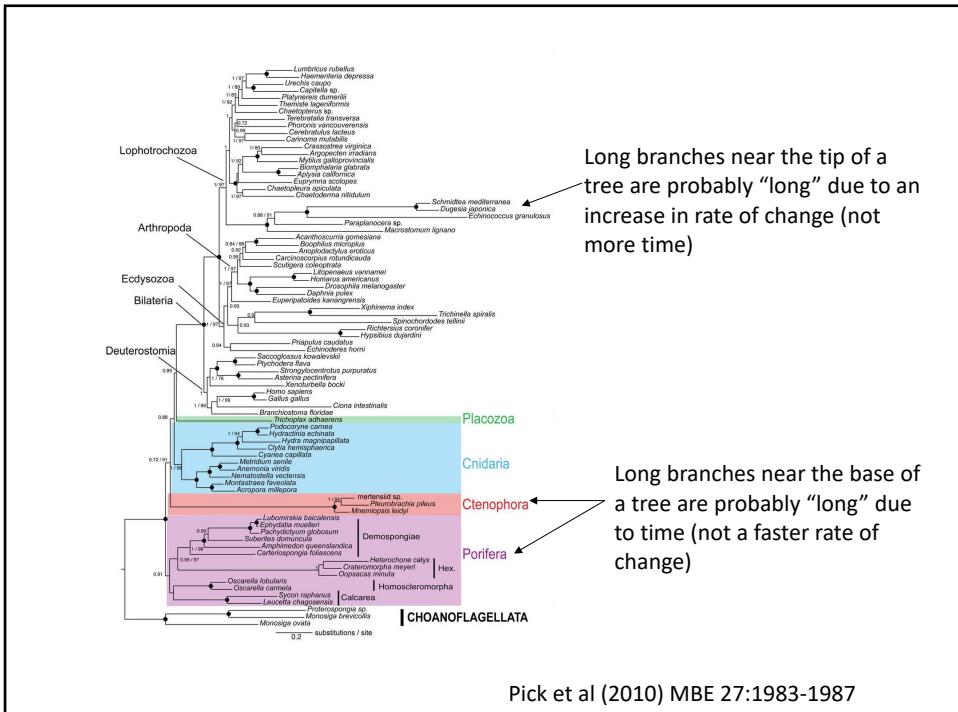
17

No universal molecular clock



Molecular distance from to is the same in all cases

18



19

Teasing apart RATE and TIME

- Branch lengths are proportional to:

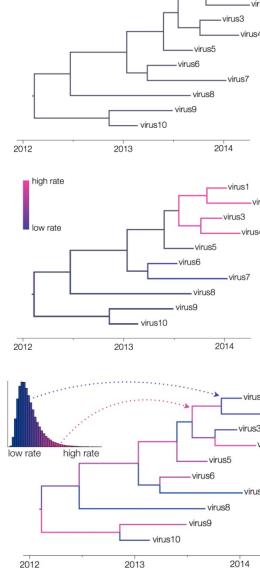
$$\text{RATE} * \text{TIME}$$

- If rates are constant then lengths are proportional to time
- If rates are not constant then we have a hard time relating branch lengths to time

20

Molecular clocks can be relaxed

- Strict or “global” clock
 - Many programs/methods/algorithms
- Local clocks
 - Maximum Likelihood (PAML, QDate)
 - Mean path length (Pathd8)
- Relaxed clocks
 - Non-parametric rate smoothing (r8s)
 - Penalized likelihood (r8s)
 - Bayesian, fixed tree (multidivtime, PhyBayes)
 - Bayesian, tree co-estimated (BEAST, MrBayes)



21

What is a relaxed clock?

- Relaxed clock: rate allowed to vary among branches
 - Autocorrelated relaxed clock: rates in adjacent branches are related
 - Uncorrelated relaxed clock: rates identically and independently distributed among branches



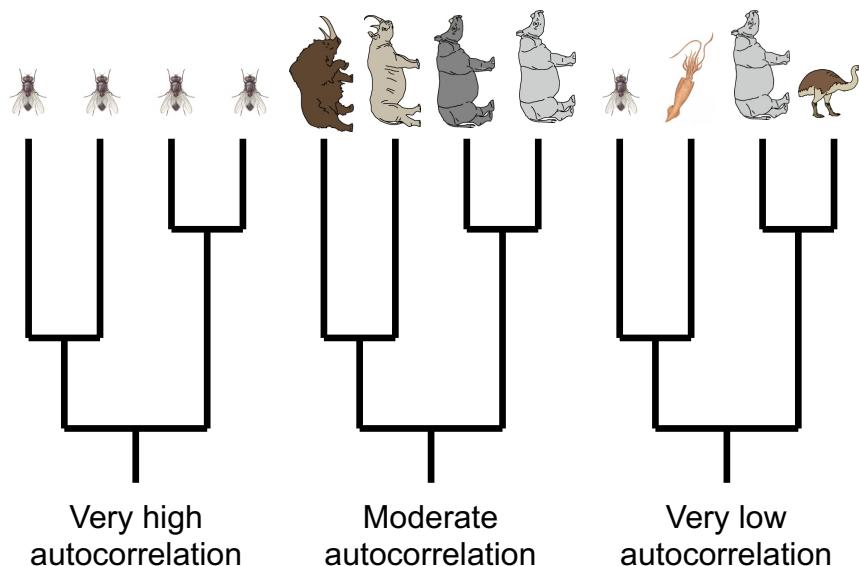
22

Autocorrelated relaxed clocks

- Fixed topologies are input!
- Treat substitution rate as a heritable trait, so that it can ‘evolve’ through the tree
- Rate is assumed to be tied to:
 - Life history traits (e.g., generation time, population size, body size)
 - Cellular/biochemical environment
- Available in r8s, multidivtime, PhyBayes, BEAST, PAML

23

Autocorrelated relaxed clocks



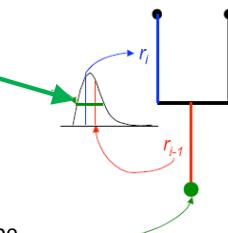
24

24

Modeling autocorrelation

- Model of autocorrelated rate change used to describe prior distribution of rates
- Lognormal
 - $\log(r_i) \sim N(\log(r_{i-1}), v_t)$

v controls the s.d. of the distribution

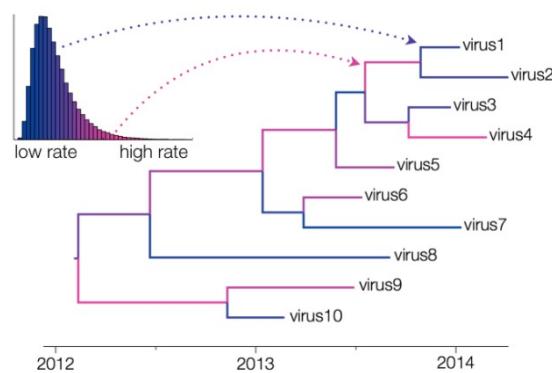


Further assumption needs to be made about rate at the root

25

Uncorrelated relaxed clocks

- Models available in *BEAST*
 - **Lognormal distribution**
Most rates cluster around the mean
 - **Exponential distribution**
Most rates are quite low



26

13

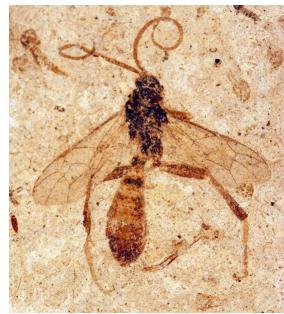
Lognormal uncorrelated relaxed clock

- In the uncorrelated lognormal relaxed clock, two statistics can be obtained:
 - 1. Coefficient of variation of rates**
Measures the rate variation among branches
A value of 0 indicates clocklike evolution
 - 2. Covariance of rates**
Measures autocorrelation of rates between adjacent branches

27

Challenges

- Saturation
- Rate Heterogeneity - violation of homogeneity
- Calibration



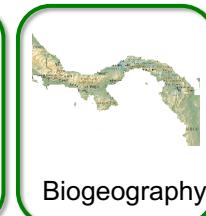
28

Separating rate and time

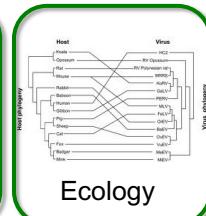
- Information about rate
 - Substitution rate obtained from an independent study
- Information about time – *prior information*:



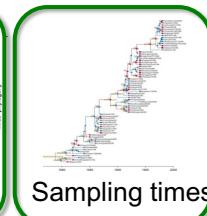
Fossil record



Biogeography



Ecology



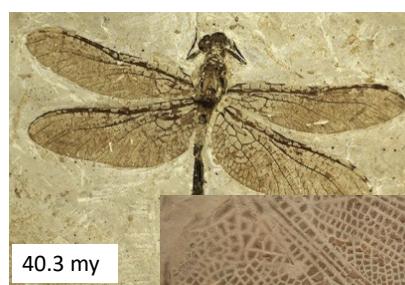
Sampling times

29

Calibration: Fossil record

- Fossil record provides minimum estimates of divergence times

Fossil record



Identified as belonging to
the family Aeshnidae and
genus *Aeshna*

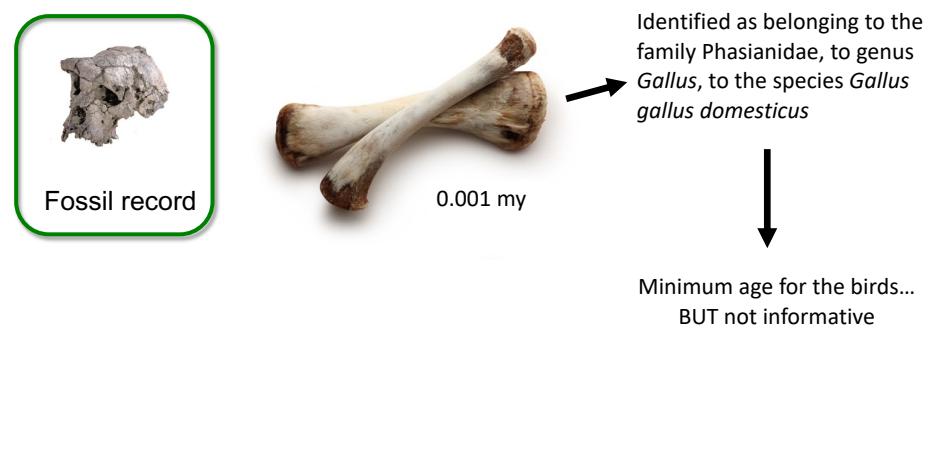


informative

30

Calibration: Fossil record

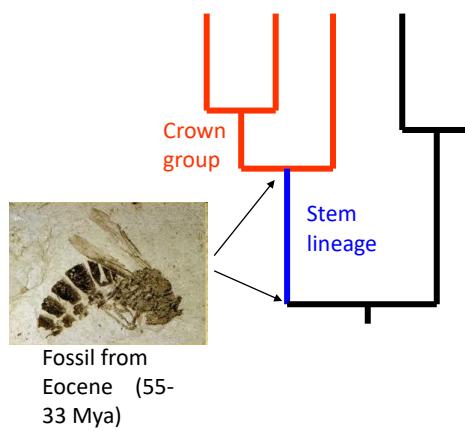
- Fossil record provides minimum estimates of divergence times



31

Problems with fossils

- Incompleteness of fossil record
- Identification
 - Species / Genus / Family?
- Position
 - Stem or crown?
- Which date?
 - Min / Mid / Max of Epoch?



32

Calibration errors

- Preservational bias
 - Hard parts
 - Environment, proximity to water bodies
 - Age
 - Sampling effort
- Taxonomic affinity
 - Fragmentary fossils
 - Extinct, stem lineages
- Stratigraphic and isotopic dating errors

33

Calibration: Biogeography

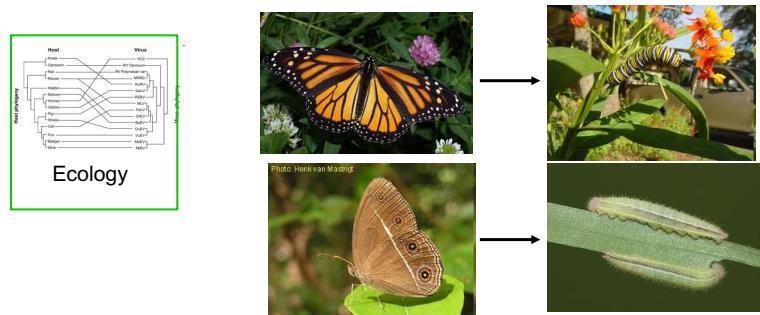
- Biogeographic events can provide maximum estimates of divergence times



34

Calibration: Ecology

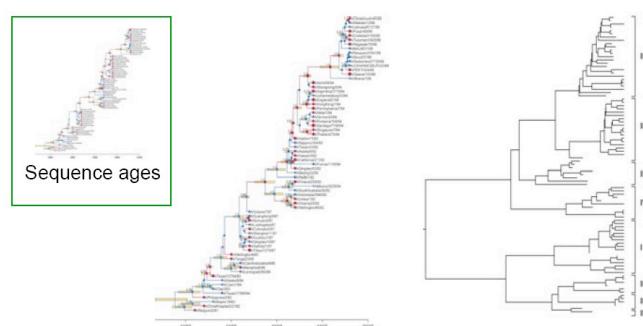
- Knowledge of tight ecological associations can be used to provide maximum estimates of divergence times



35

Calibration: Sequence ages

- Sequence ages provide sufficient age information for e.g. viruses



36

Calibration in Bayesian framework

$$f(\theta|D) = \frac{f(D|\theta)f(\theta)}{\int f(D|\theta)f(\theta)d\theta}$$

θ : model (substitution model(s), tree, etc)

prior: prior expectation we have for parameters of the model

37

Calibration in Bayesian framework

$$f(\theta|D) = \frac{f(D|\theta)f(\theta)}{\int f(D|\theta)f(\theta)d\theta}$$

θ : model (substitution model(s), tree, etc)

prior: prior expectation we have for parameters of the model

For example: age of nodes based on fossil information

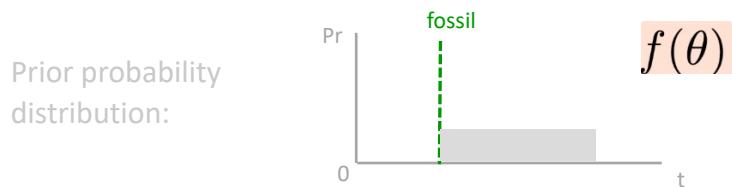
38

19

Calibration in Bayesian framework

$$\text{posterior} \quad \text{data} \quad \text{prior}$$
$$f(\theta|D) = \frac{f(D|\theta)f(\theta)}{\int f(D|\theta)f(\theta)d\theta}$$

For example: age of nodes based on fossil information

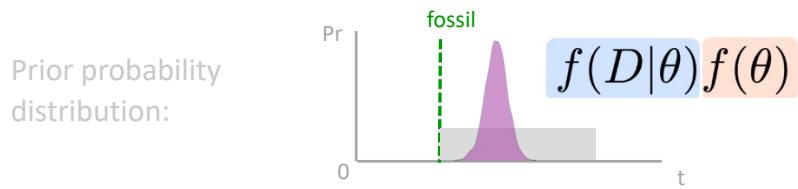


39

Calibration in Bayesian framework

$$\text{posterior} \quad \text{data} \quad \text{prior}$$
$$f(\theta|D) = \frac{f(D|\theta)f(\theta)}{\int f(D|\theta)f(\theta)d\theta}$$

For example: age of nodes based on fossil information



40

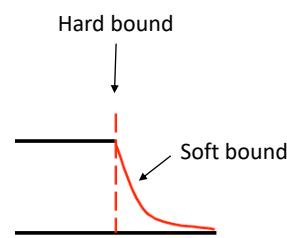
Calibration types

- Point calibrations
- Hard minimum/maximum bounds
- Soft minimum/maximum bounds
- Parametric prior distributions
 - Normal distribution
 - Lognormal distribution
 - Exponential distribution

41

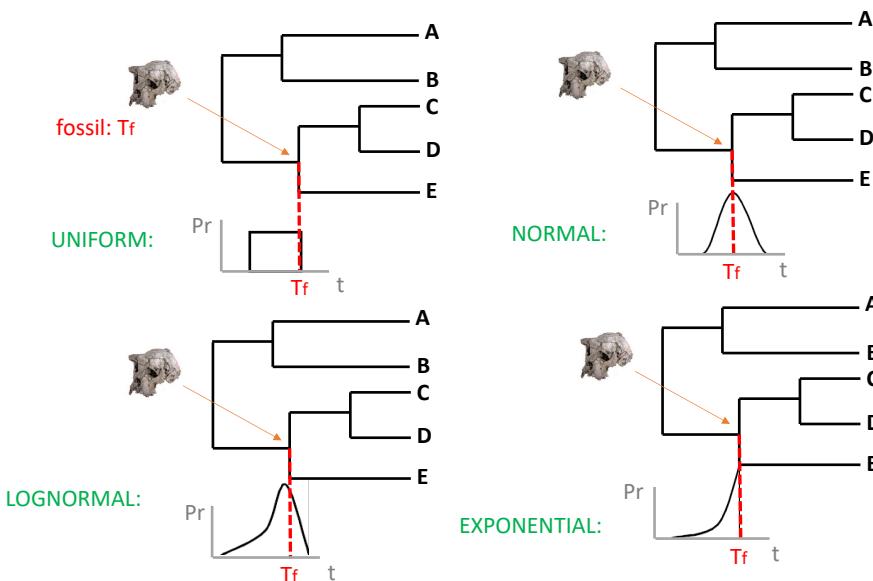
Hard/Soft Bounds

- Extension of hard bounds
- Soft:
 - Assign non-zero probability to values outside bound
 - Able to forgive calibration errors

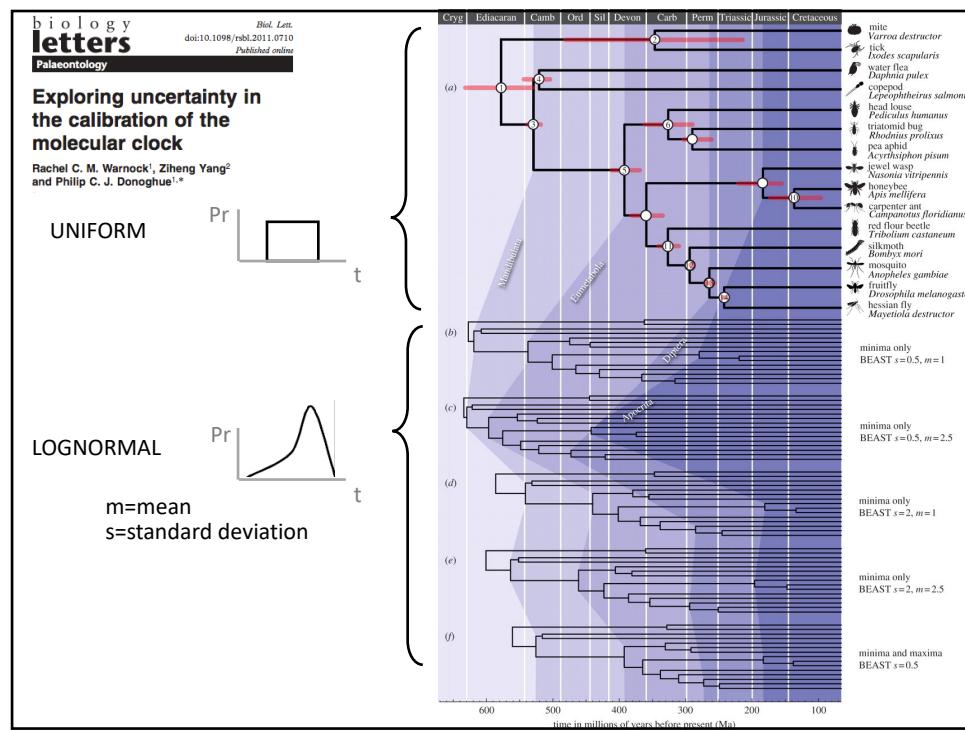


42

Prior distributions

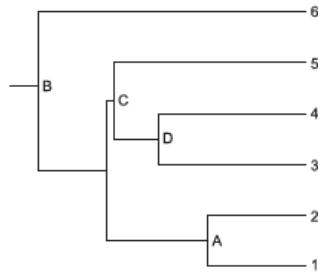


43



44

Multiple calibrations

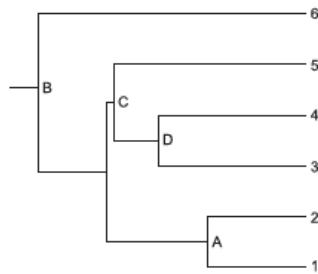


- ▶ Molecular-clock estimates can be sensitive to the positions of the calibrations in the phylogenetic tree, especially when only a single or very few calibrations are available
- ▶ a small number of calibrations can lead to a biased estimate of the substitution rate if there is substantial among-lineage rate variation

5

45

Multiple calibrations



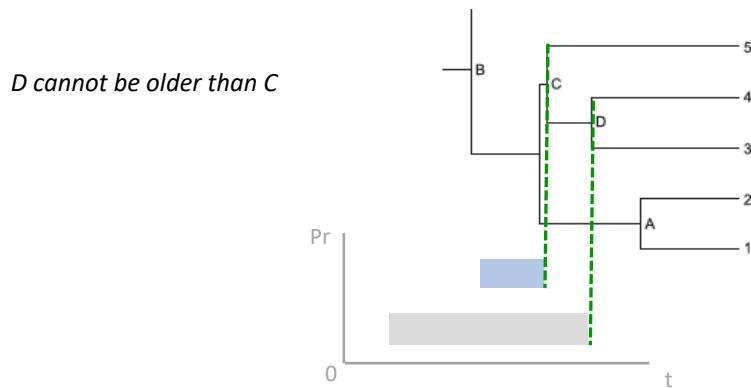
- ▶ can improve the accuracy of date estimates in the presence of taxon undersampling
- ▶ substitution rate is primarily estimated from the branches between the calibrating nodes and the tips
=> deeper calibrations capture a larger proportion of the overall genetic variation.

6

46

Multiple calibrations

- ▶ Be careful: priors interact with each others
- ▶ For example, node orders



47

47

Multiple calibrations

- ▶ Be careful: priors interact with each others
- ▶ For example, node orders
- ▶ Marginal priors resulting from prior interactions can differ from the initial user prior
 - This can be visualized by removing the data and running the same analysis

$$f(\theta|D) = \frac{f(D|\theta)f(\theta)}{\int f(D|\theta)f(\theta)d\theta}$$

48

48

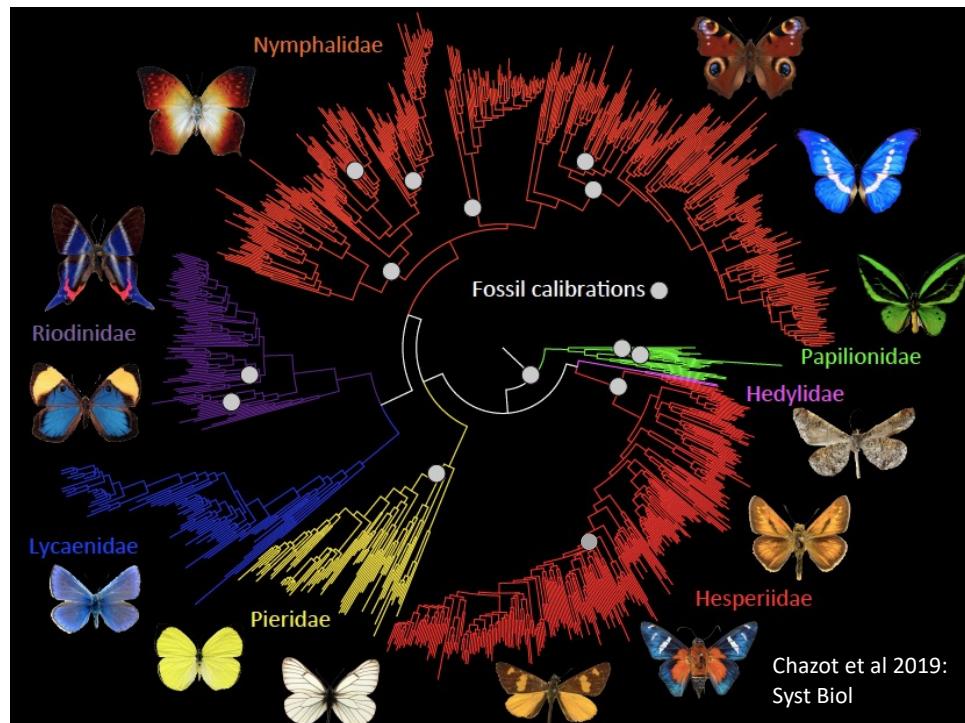
Multiple calibrations

- ▶ Be careful: priors interact with each others
- ▶ For example, node orders
- ▶ Marginal priors resulting from prior interactions can differ from the initial user prior
 - This can be visualized by removing the data and running the same analysis

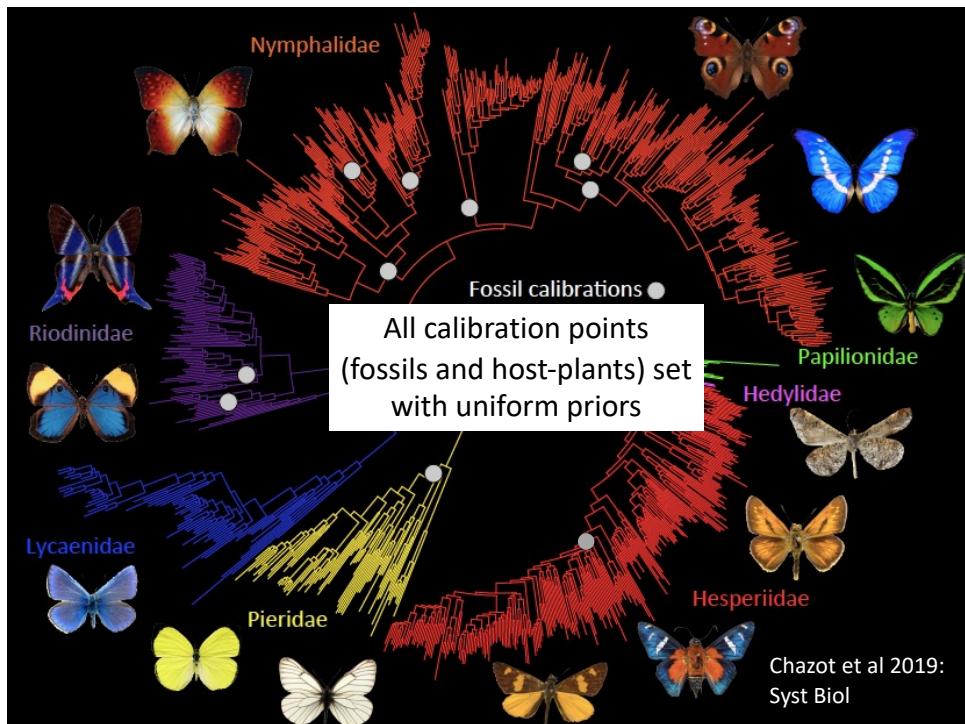
$$f(\theta|D) = \frac{f(\theta)}{\int f(D|\theta)f(\theta)d\theta}$$

49

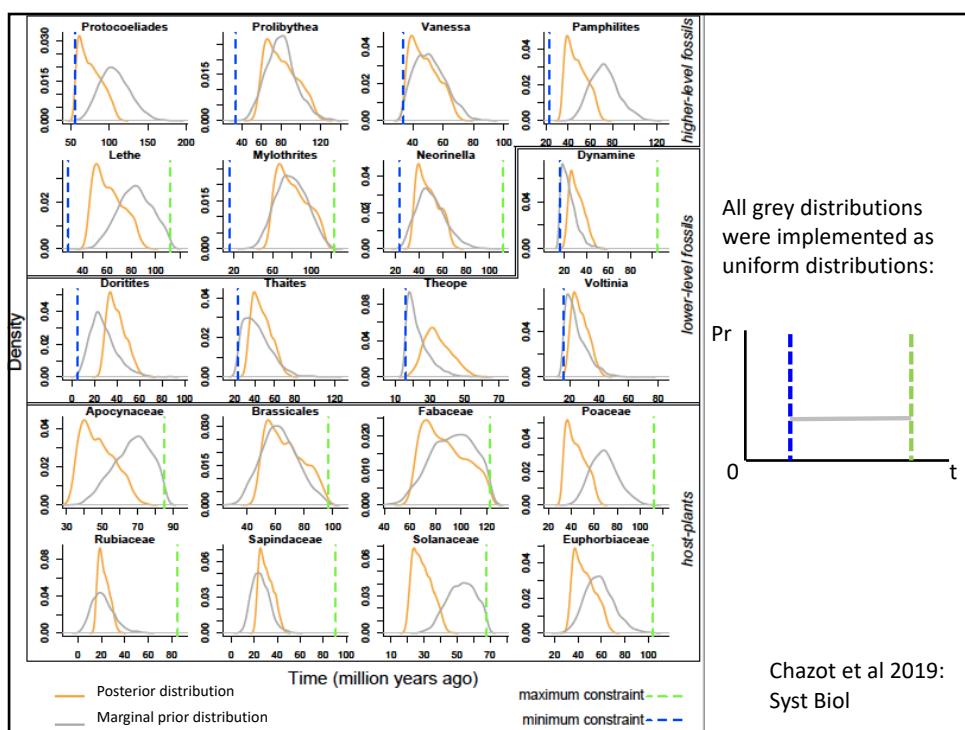
49



50

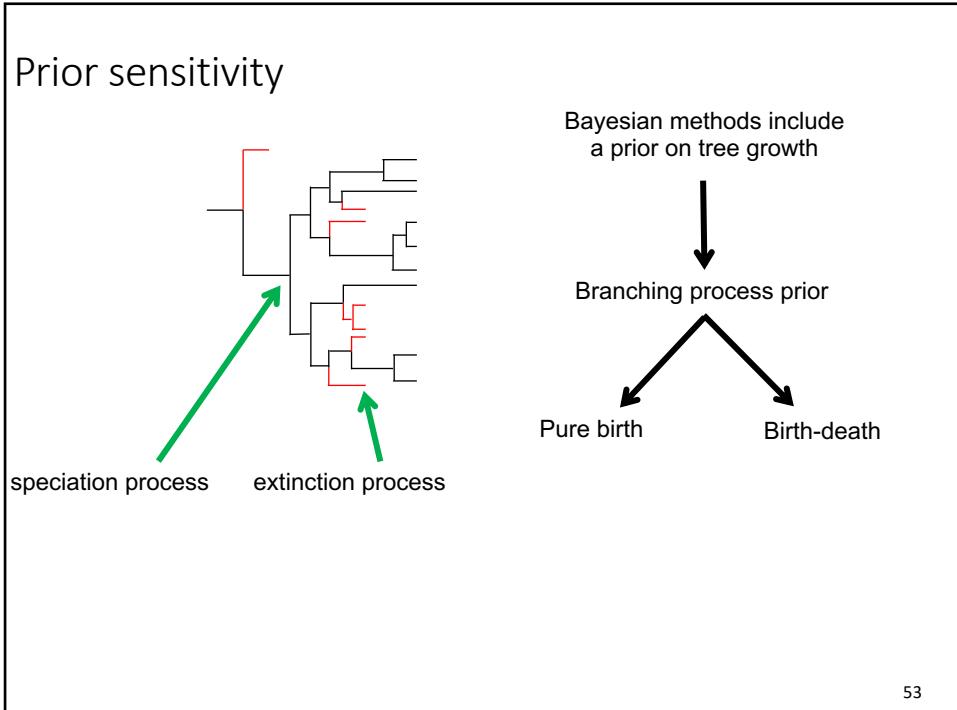


51



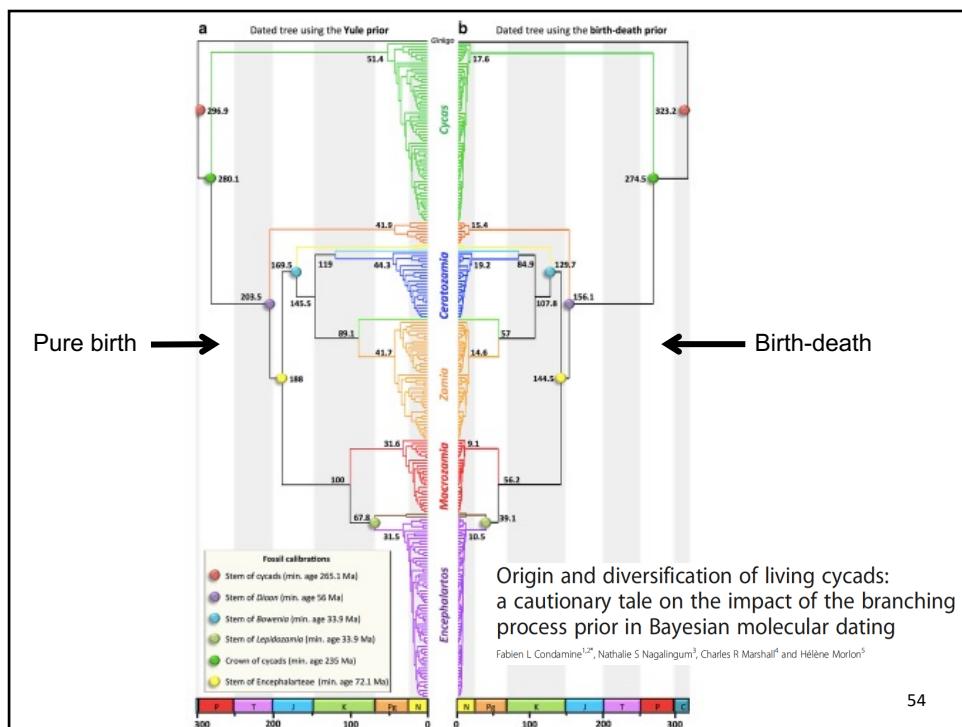
52

Prior sensitivity

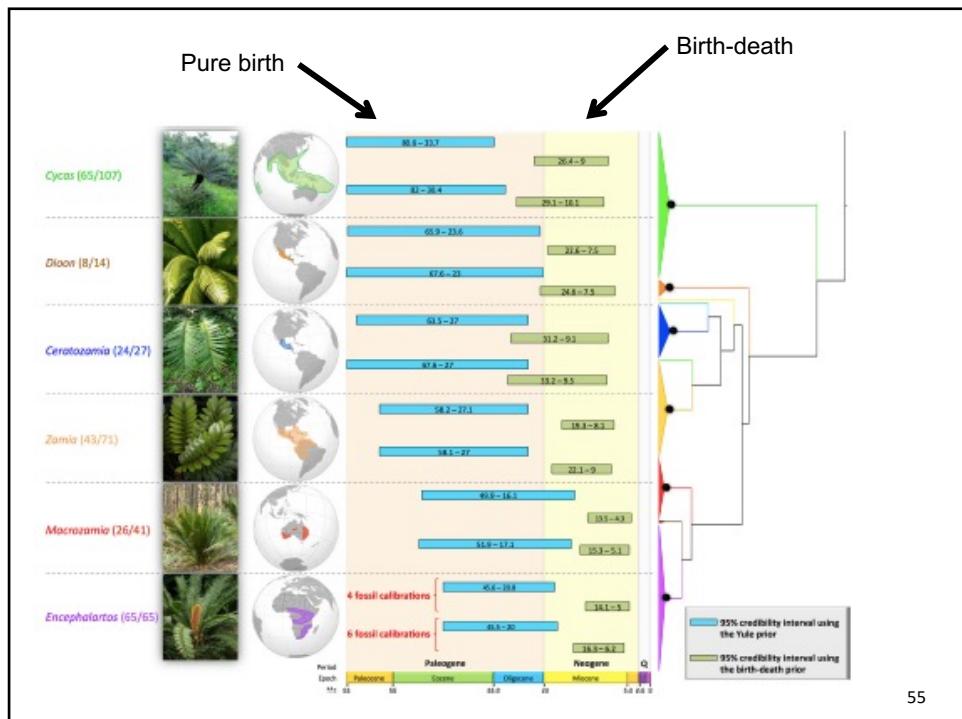


53

53



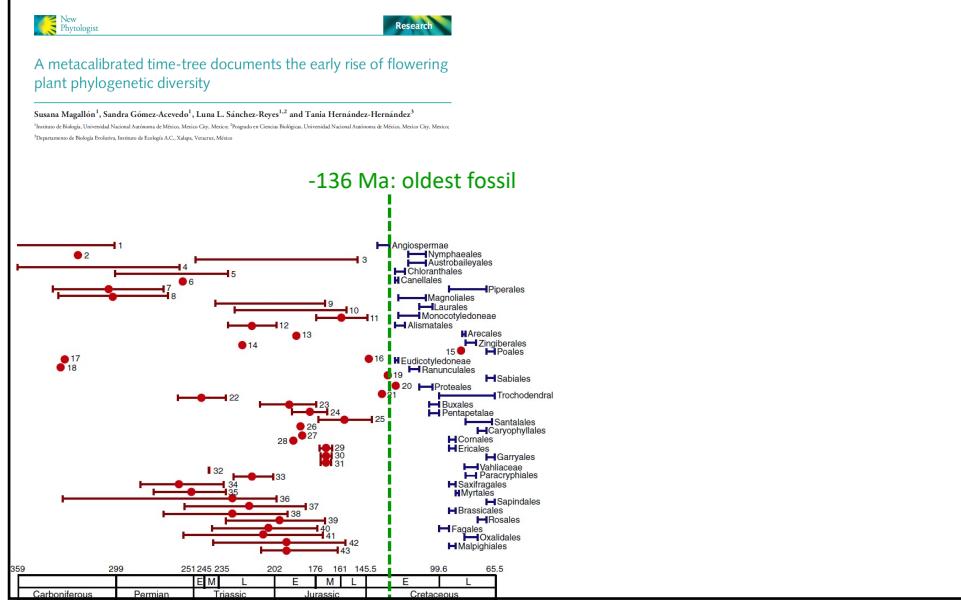
54



55

55

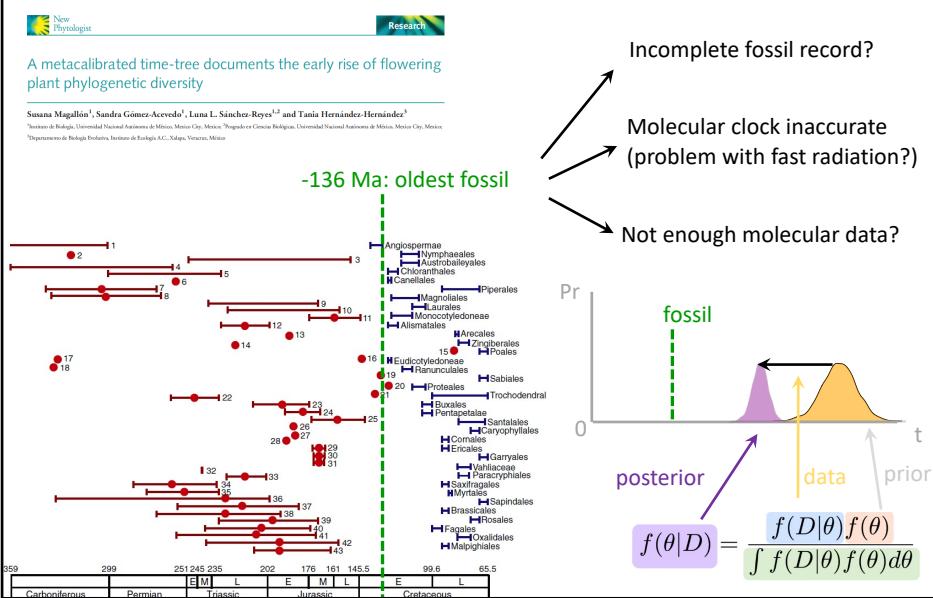
Timing the origin of Angiosperms



56

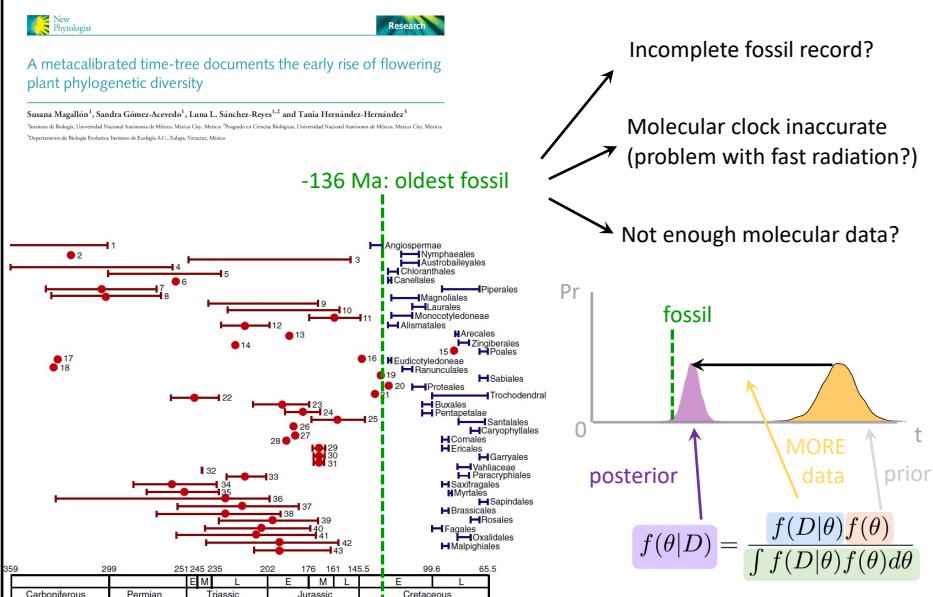
28

Timing the origin of Angiosperms



57

Timing the origin of Angiosperms



58

Timing the origin of Angiosperms

Eocene lantern fruits from Gondwanan
Patagonia and the early origins of Solanaceae

Peter Wilf^{1,*}, Mónica R. Carvalho², María A. Gandolfo², N. Rubén Cúneo³
♦ See all authors and affiliations

Science 06 Jan 2017:
Vol. 355, Issue 6320, pp. 71-75
DOI: 10.1126/science.aag2737



Physalis infinemundi
Physalis
tomatillo group - 9 to 11 My
Nightshades - 35 to 51 My



59

Timing the origin of Angiosperms

Eocene lantern fruits from Gondwanan
Patagonia and the early origins of Solanaceae

Peter Wilf^{1,*}, Mónica R. Carvalho², María A. Gandolfo², N. Rubén Cúneo³
♦ See all authors and affiliations

Science 06 Jan 2017:
Vol. 355, Issue 6320, pp. 71-75
DOI: 10.1126/science.aag2737



Physalis infinemundi
Physalis
tomatillo group - 9 to 11 My
Nightshades - 35 to 51 My



60

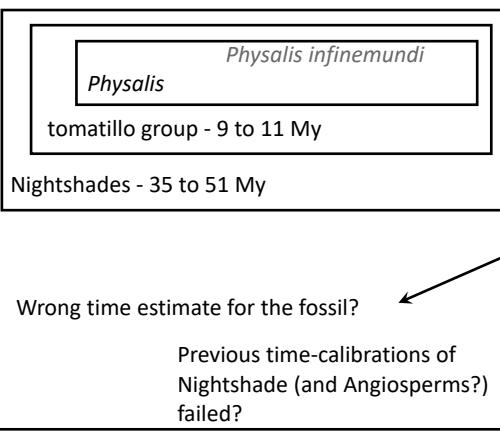
30

Timing the origin of Angiosperms

Eocene lantern fruits from Gondwanan
Patagonia and the early origins of Solanaceae

Peter Wilf^{1,*}, Mónica R. Carvalho², María A. Gandolfo³, N. Rubén Cúneo³
• See all authors and affiliations

Science 06 Jan 2017:
Vol. 355, Issue 6320, pp. 71-75
DOI: 10.1126/science.aag2737



-52 Ma !!!

Wrong time estimate for the fossil?

Previous time-calibrations of
Nightshade (and Angiosperms?)
failed?

Over-estimation of the
phylogenetic position?

61