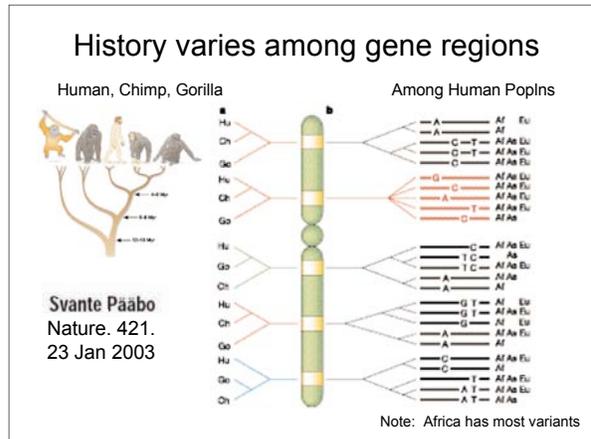
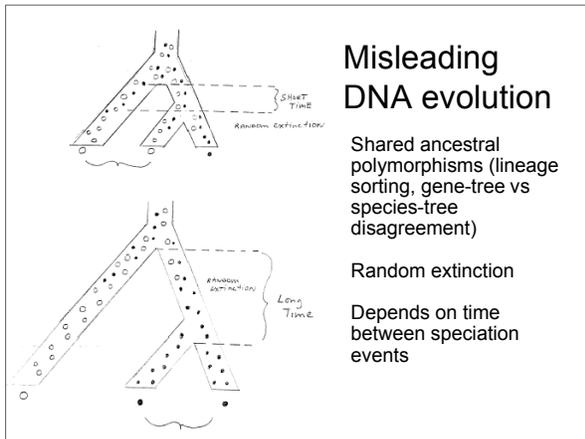


Accommodating Biases in DNA Data

Chris Simon
Ecology & Evolutionary Biology
University of Connecticut

DNA data do not guarantee the correct phylogeny

1. The problem of shared ancestral polymorphisms
2. Multiple substitutions at a single site hide earlier substitutions



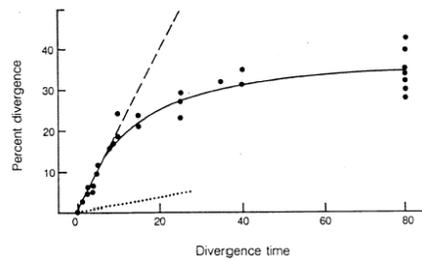
Other reasons for gene tree not matching species tree.

- Hybridization (organelle genomes move across species boundaries and introgress faster than nuclear genes)- recent or ancient
- Horizontal gene transfer

Solution: use methods that create species trees from multiple gene trees (Deep coalescence, gene tree parsimony/consensus tree, Bayesian coalescent gene tree conditioned on species tree, e.g., BEST)

Misleading DNA evolution

2. Multiple substitutions hide previous changes



Corrections for multiple substitutions

Jukes-Cantor (1969) Assumptions

1. $A = T = G = C$ No nucleotide bias
2. Every base changes to every other base with equal probability (no TS/TV bias)
3. All sites change with the same probability (no ASRV)

Also: probability substitution & base composition remains constant over time/across lineages

Jukes-Cantor Assumptions Incorrect

1. Nucleotide bias is common
mtDNA honey bee 84.9% A+T;
D. yakuba 78.6%;
Evanoid wasps 77.8% (low for hymenoptera)

Some bacterial lineages are A+T rich,
others G+C rich.

Accommodating biases

- Can be done by weighting in parsimony analyses
- Using models of evolution in maximum likelihood and Bayesian analyses (models that improve on Jukes Cantor)

Biased substitutions accommodated in parsimony by weighting

Substitution bias: By weighting (step matrices) but very difficult to do well and sophisticated weighting slows MP analysis significantly.

ASRV: By using successive approximations or other iterative weighting but highly dependent on initial tree. Or by down-weighting variable regions.

Simon et al. 1994. Annals ESA

History of ASRV: Weighting to accommodate ASRV

Less variable characters (less homoplasy) given more weight

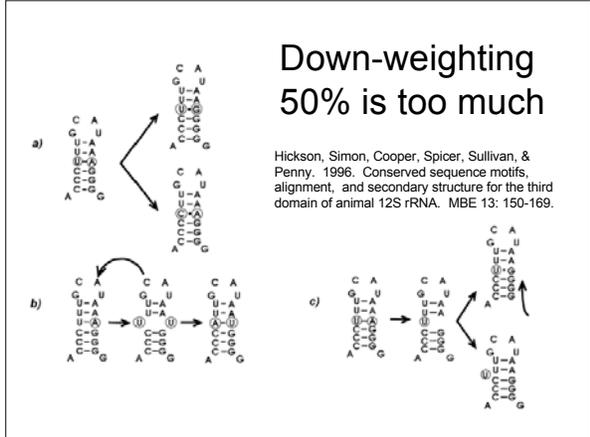
Weight 1st & 2nd position more heavily than third (Irwin et al. 1991) **Clearly important**



Weight stems more heavily than loops (or vice versa) (DeSalle et al. 1987, Vawter & Brown 1993, Springer & Douzery 1996). **DO NOT DO THIS**

History of ASRV: Weighting to accommodate ASRV

Weight paired regions less than unpaired regions due to compensatory changes (Wheeler & Honeycut 1988, Hillis & Dixon 1991) **How much to down-weight? Fifty percent too much. Hillis & Dixon recommended less.**



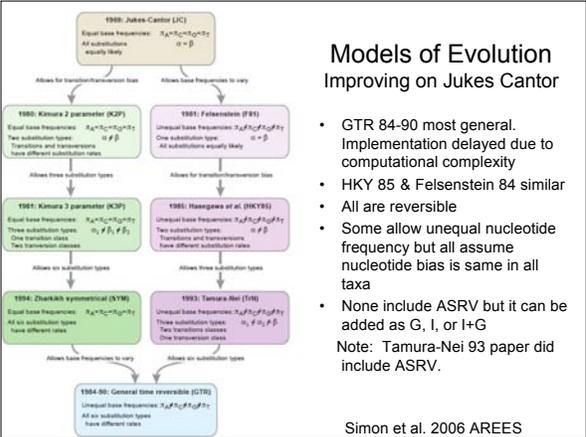
Huelsenbeck & Ronquist include a doublet model in MrBayes

- Accommodates paired, correlated changes
- Kjer 2004 found this model improved results of 18S rRNA phylogeny of insect orders

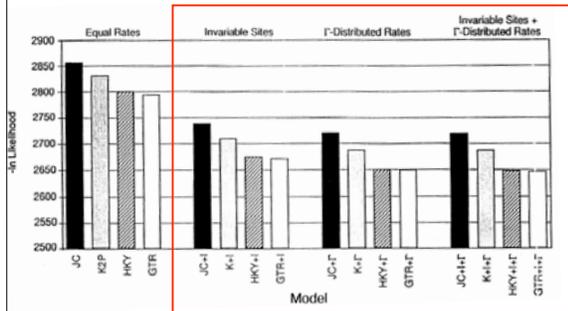
History of ASRV: Weighting to accommodate ASRV

Divide a rRNA molecule into rate variability classes (Simon 1991, Hickson 1993)
 Problem: **Where to divide?**

Calculate probability of substitution at individual sites....

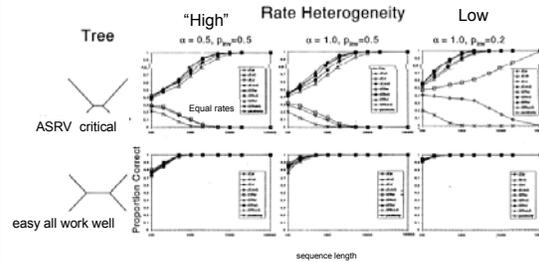


ASRV >> fit improvement than models



Frati, Simon, Sullivan, Swofford. 1997. JME 44:145-158

Model effectiveness in phylogeny construction depends on ASRV & tree shape



Sullivan & Swofford 2001. Syst. Biol. 50: 723-729

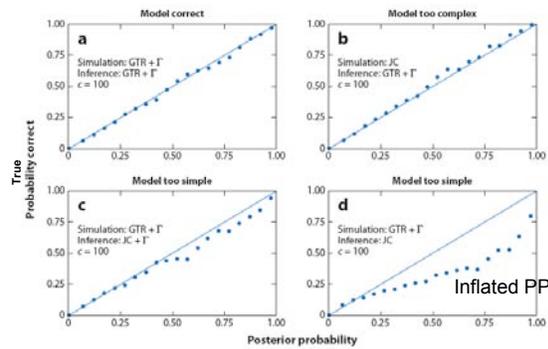
Yang 1996. TREE. Summarized effects of ignoring ASRV

- Under estimates ts/tv rate ratio
- Lower phylogenetic informativeness
- Causes Tajima's D for neutrality to mimic patterns of population expansion.
- Results in under- or over-estimates of divergence times

More effects of ignoring ASRV

- Strongly affects bootstrap support (Buckley et al. 2001 Syst. Biol., Cunningham & Buckley 2002. MBE)
- Strongly affects Bayesian posterior probabilities in simulations (Lemmon & Moriarty 2004. Syst. Biol.)
- Reduces effectiveness of models of AA evolution (Susko et al. 2003. Syst. Biol.)
- Results in incorrect inferences of early, rapid cladogenesis (Ravell et al. 2005. Syst. Biol.)

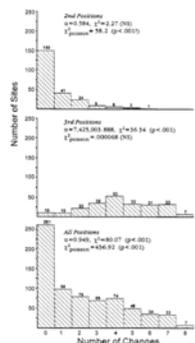
Ignoring ASRV in too-simple model --> inflated PP (Huelsenbeck & Rannala 2004)



Problem: Modeling all genes combined creates a single average model that doesn't fit any one partition very well

Partitioned models avoid "average ASRV" problem

- Partitioned models w/ ASRV improved branch supports & tree likelihood
 - Nylander et al. 2004
 - Castoe et al. 2004
 - Brandley et al. 2004
- In mtDNA gene boundaries can be ignored



Partitioned Models

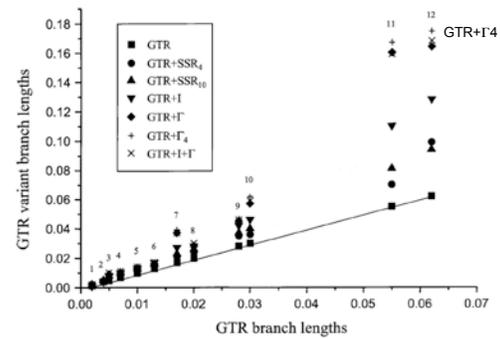
- An attempt to accommodate different rate variation in different partitions.
- Partitioned likelihood available early in PAML but slow
- Partitioned likelihood could only be done in PAUP* by summing across likelihood runs*
- RaxML allows partitions
- Partitioning in MrBayes since 2003

*suggested by Swofford and described in Schwartz et al. 2003.

Partitioned Models: early days

- SSR models in PAUP* accommodate different rate variation in different partitions.
- But have restrictive assumption- the rate of evolution of all sites within each partition is equal
- Newer SSR-gamma or partitioned gamma calculates gamma parameters separately for each partition

How ASRV is accommodated matters



Buckley, Simon & Chambers. 2001. Syst. Biol.

Buckley et al. 2001. Conclusions

- Important to correct for ASRV
- Large difference among corrections in branch lengths & bp
- SSR model branch lengths not much different from uncorrected
- The more distantly related the taxa, the more the ASRV correction matters

Difficult to know how to partition rRNA

Solution: Mixture Models. Pagel & Meade 2004

- Characters viewed as generated by one of a number of distributions
- Each of these distributions contributes to the likelihood during analysis
- Parameters of each model distribution & weights assigned to them estimated from the data
- At end of analysis, each character assigned a probability of belonging to each model

Pagel & Meade 2004

TABLE 6. Numbers of stem and loop sites fitted best by respective rate matrices: mammalian 12S data.

Secondary structure	Q1	Q2	Q3	Q4
Stem	76	21	71	296
Loop	133	112	83	249

Each Q matrix was GTR + G (partitioned & not); parameters varied among the four

- No one model fitted stems vs loops
- Instead, various stem & loop sites were distributed across 4 models with model 4 fitting best the largest number of sites

Pagel & Meade 2004

TABLE 6. Numbers of stem and loop sites fitted best by respective rate matrices: mammalian 12S data.

Secondary structure	Q1	Q2	Q3	Q4
Stem	76	21	71	296
Loop	133	112	83	249

Advantages of Mixture Models

- Patterns are emergent- not pre-specified
- Not constrained
- Not dependent on knowledge of secondary structure

Mixture Models vs Model Averaging*

- Model averaging- several models. Weights and models fixed (e.g., from AIC scores). Models often differ in number of parameters.
- Mixture models- Weights are estimated during analysis; multiple models often the same but differ in the parameter values

* Thanks to Paul Lewis