

**A stochastic Dollo-model for phylogenetic inference:
model elaboration and checking.**

Q Atkinson, GN, D Welch, and R Gray ‘From words to dates: water into wine, mathemagic or phylogenetic inference’, *Trans. Philological Soc.* 103, 193-219 (2005)

GN, RD Gray, ‘Dated ancestral trees from binary trait data and its application to the diversification of languages’, *J. Roy. Statist. Soc. B*, 70:545-566 (2008)

RJ Ryder, GN, ‘Missing data in a stochastic Dollo model for cognate data, and its application to the dating of Proto-Indo-European’, *arXiv:0908.1735* (2009 and JRSSC, to appear).

Data

1. *Ringe et al. 2002* 24 Indo-European languages 4000 distinct traits (*schematically*)

	ALL	AND	ANIMAL			
oldenglish	1000000	100000000000000000	00010000	.	.	.
oldhighgerman	1000000	100000000000000000	00010000			
gothic	1000000	??????????????????	00001000			
latin	0100000	0000000100000000	00000100			
umbrian	0100000	00000000000000010	?????????	.	.	.
...			

New Data

2. *M Rother 2008*, from Sprachatlas von Bayerisch-Schwaben

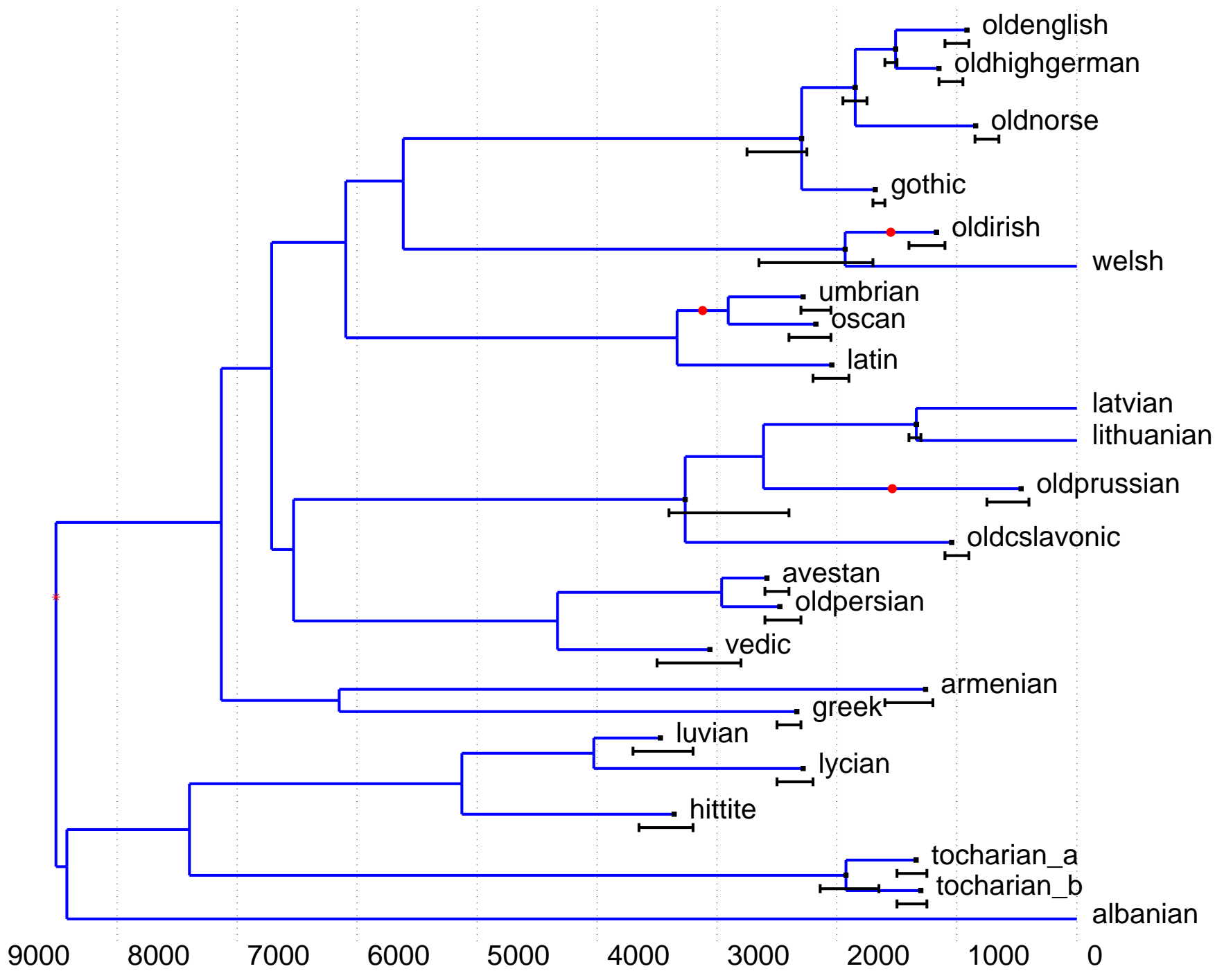
Includes dialects Ries, Ostfranken, Nordbaiern, Ostlech(Sued), Ammersee(Ost), Lechrain,...

3. *S Greenhill et al 2008*

400 Austronesian languages (Maori, Indonesian, Malagasy,...), 30000 lexical traits

4. *A Kitchen et al 2009*

25 Semitic languages (Hebrew, Akkadian, Amharic,...), 674 lexical traits



Generic trait-observation model

$$g = (E, t, k) \quad t = (t_1, \dots, t_{2L-1})$$

$$z = (z_1, \dots, z_N) \quad z_a = (\tau_a, i_a) \in [g]$$

λ : set-element birth rate;

μ : set element death rate;

ρ : catastrophe rate, $k = (k_1, \dots, k_{2L-2})$

κ , death prob, $\nu = \lambda\kappa/\mu$ mean births

ξ_i : probability element a visible at leaf i

$$Z \sim \Pi(\lambda; [g]),$$

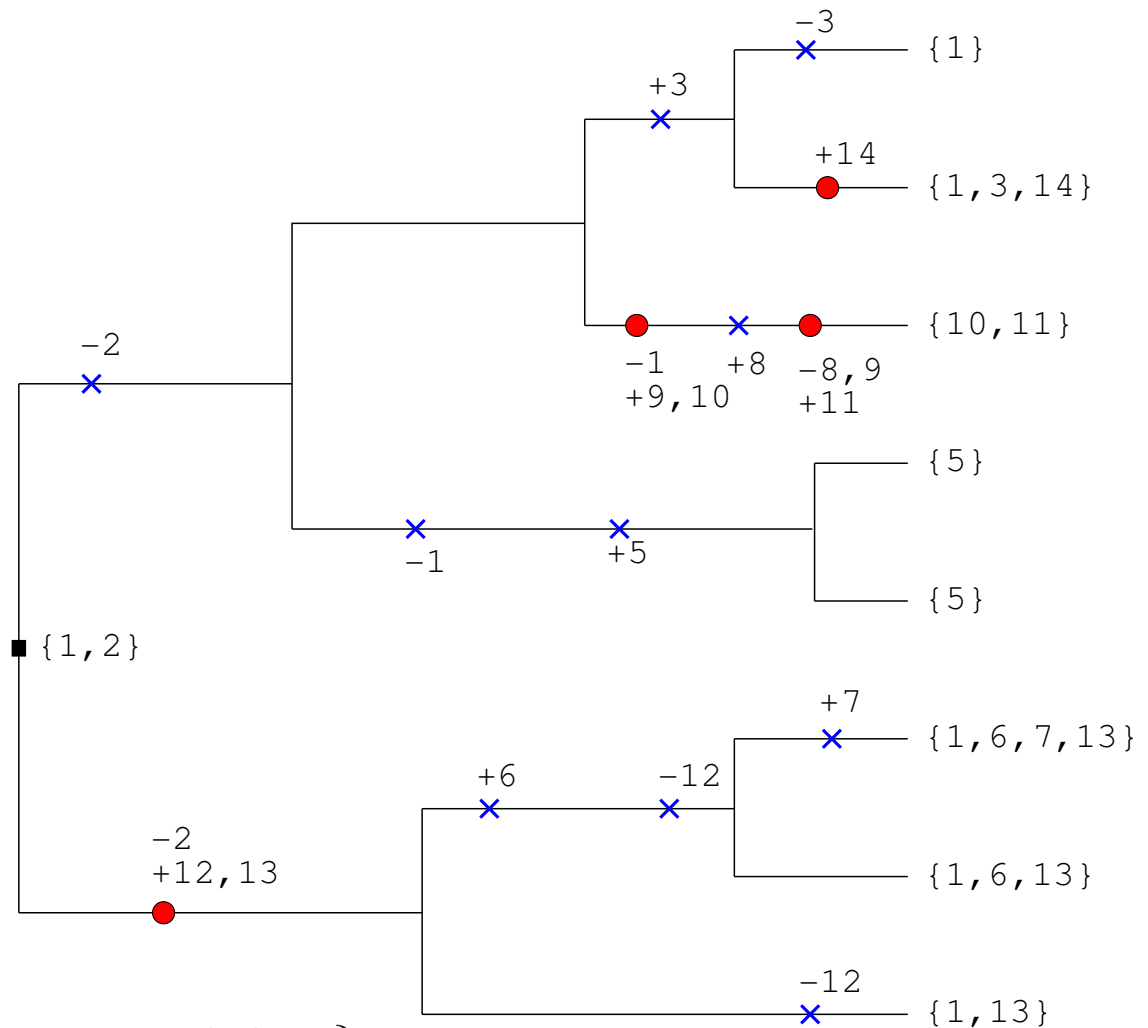
$$\mathcal{E}_{Z_0} = \{\text{birth at } Z_0 \in [g] \text{ yields registered data}\}$$

$$\lambda(z_0) = \lambda \Pr(\mathcal{E}_{Z_0} | g, \mu, \lambda, \xi, Z_0 = z_0)$$

$$Z | \mathcal{E}_Z \sim \Pi(\lambda(z); [g]),$$

$$\Pr(dz | g, \mu, \lambda, \xi) = f_Z(z) dz$$

$$= \frac{1}{N!} e^{-\Lambda([g])} \prod_{a=1}^N \lambda(z_a) dz_a$$



[D*]	[I*]	[D~]
1000000000000000	1111111111111111	1000000000000000
1010000000000001	0101111111111111	?0?0000000000001
000000000011000	10111001110110	0?000??001?00?
000010000000000	1111111111111111	000010000000000
000010000000000	1111011111111111	0000?0000000000
100001100000010	1011111111111111	1?0001100000010
100001000000010	1111111111111111	100001000000010
1000000000000010	1111111111111100	1000000000000??

[I]	[D]
1 111 1 11	1 000 0 00
0 111 1 11	? 000 0 01
1 100 1 10	0 0?? 1 0?
1 111 1 11	0 100 0 00
1 011 1 11	0 ?00 0 00
1 111 1 11	1 011 0 10
1 111 1 11	1 010 0 10
1 111 1 00	1 000 0 ??

$$Y_a = \sum_{i=1}^L \mathbf{I}_{i,a}^* \mathbf{D}_{i,a}^*$$

$$R_2(\tilde{D}) = (\tilde{D}_a : Y_a > 1)$$

$$\mathcal{E}_{Z_a} = \{\mathbf{D}_a = R_2(\tilde{\mathbf{D}}_a)\}$$

Likelihood

$$\begin{aligned}
P[\mathbf{D} = D \mid g, \mu, \lambda, \xi, \mathbf{D} = R(\tilde{\mathbf{D}})] &= \int P[\mathbf{D} = D \mid g, \mu, \xi, Z = z, \mathbf{D} = R(\tilde{\mathbf{D}})] f_Z(z) dz, \\
&= \frac{1}{N!} \left(\frac{\lambda}{\mu} \right)^N \exp \left(-\frac{\lambda}{\mu} \sum_{\langle i, j \rangle \in E} P[\mathcal{E}_Z \mid Z = (t_i, i), g, \mu, \kappa, \xi] (1 - e^{-\mu(t_j - t_i + k_i T_C)}) \right) \\
&\quad \times \prod_{a=1}^N \left(\sum_{\langle i, j \rangle \in E_a} \sum_{d^* \in \mathcal{D}_a} P[\mathbf{D}_a^* = d^* \mid Z_a = (t_i, i), g, \mu, \kappa] (1 - e^{-\mu(t_j - t_i + k_i T_C)}) \right) \\
&\quad \times \prod_{a=1}^N \prod_{i=1}^L \xi_i^{I_{i,a}} (1 - \xi_i)^{1 - I_{i,a}}
\end{aligned}$$

Posterior

$$\begin{aligned}
p(g, \mu, \lambda, \kappa, \rho, \xi \mid \mathbf{D} = D) &= \frac{1}{\mu \lambda} f_G(g \mid T) p_R(\rho) \frac{e^{-\rho|g|} (\rho|g|)^{k_T}}{k_T!} \\
&\quad \times P[\mathbf{D} = D \mid g, \mu, \lambda, \xi, \mathbf{D} = R(\tilde{\mathbf{D}})]
\end{aligned}$$

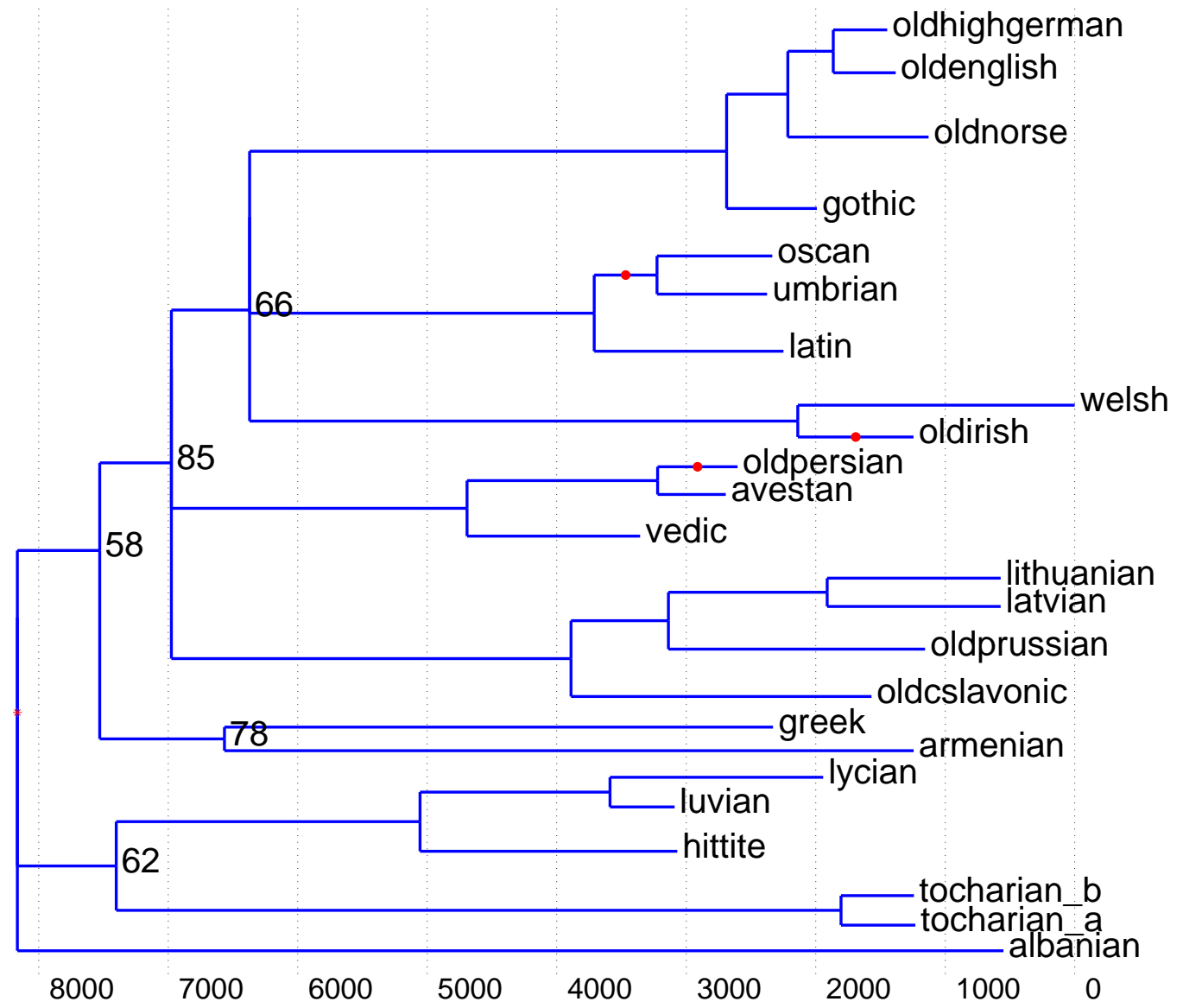
Notation

$\mathcal{D}_a = \{ \text{possible } D_{:,a}^* \text{-vectors given } D_{:,a} \}$

$E_a = \text{edges where birth at } Z_a \text{ could evolve to } D_{:,a}.$

<i>Quantity</i>	<i>95% HPD</i>	<i>units</i>
root age	(7110,9750)	yr BP
death rate μ	(1.7,2)	# per 10^4 yr
cat. rate ρ	(0.31,0.97)	# per 10^4 yr
death prob. κ	(0.3,0.42)	

A catastrophe every 15000yr
each $\simeq 2400$ yr.



Model Error: Predicting Calibrations

Predict ages of nodes for which calibration data is available.

Quantify goodness-of-fit using Bayes factors - would we reject the truth?

$\Gamma^{(c)}$ space of all trees respecting c 'th constraint, and $\Gamma = \bigcap_{c=0}^C \Gamma^{(c)}$.

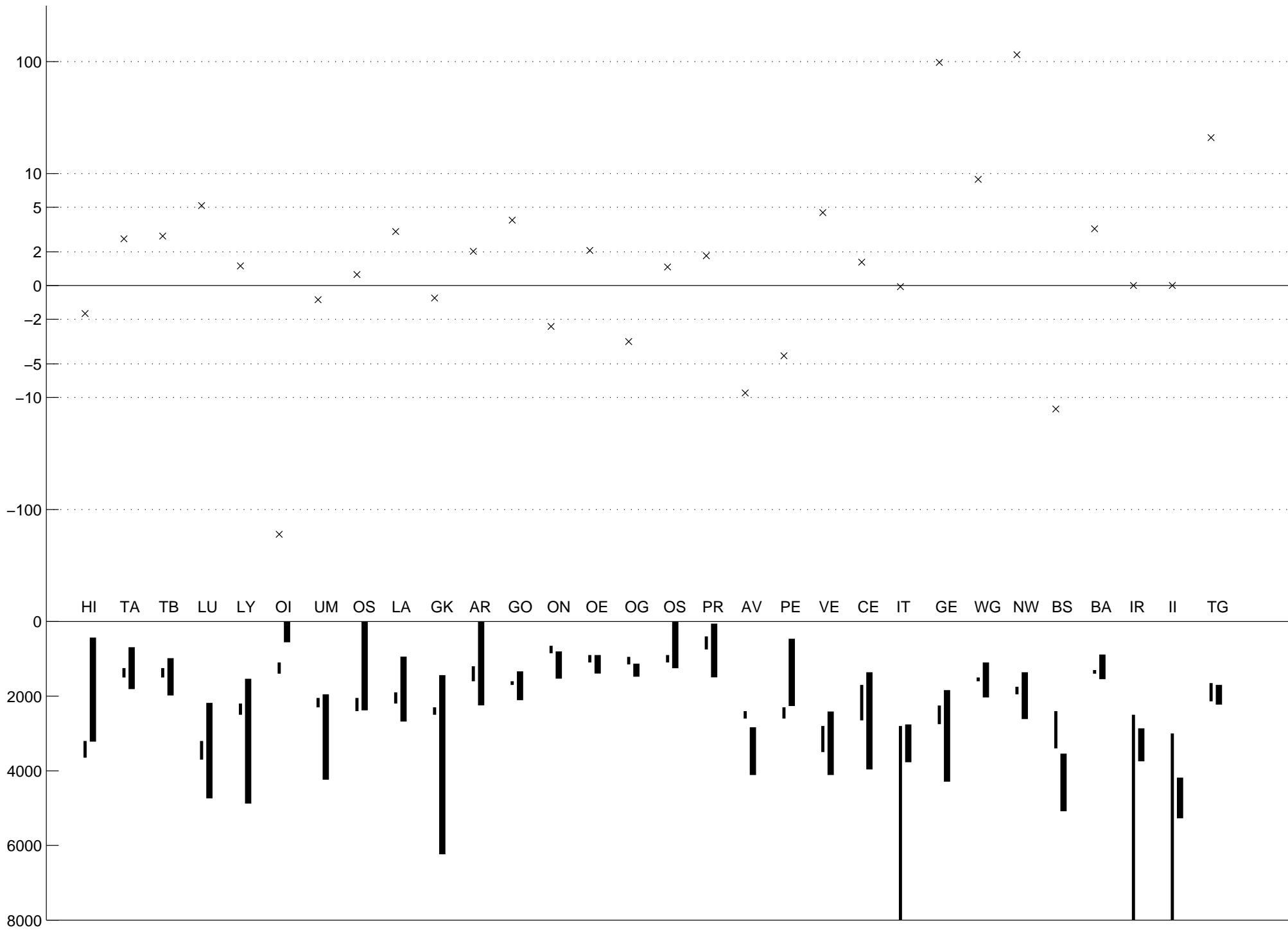
For $c = 1, 2, \dots, C$ let

$$\Gamma^{-c} = \bigcap_{\substack{c'=0 \\ c' \neq c}}^C \Gamma^{(c')}.$$

$M_0 : g \in \Gamma$, $M_1 : g \in \Gamma^{-c}$. Evidence for c 'th constraint

$$\begin{aligned} B_c &= \frac{P(D|g \in \Gamma)}{P(D|g \in \Gamma^{-c})} \\ &= \frac{P(D|g \in \Gamma \cap \Gamma^{-c})}{P(D|g \in \Gamma^{-c})} \\ &= \frac{P(g \in \Gamma|D, g \in \Gamma^{-c})}{P(g \in \Gamma|g \in \Gamma^{-c})}. \end{aligned}$$

Constraints with $2 \log(B_c) \lesssim -5$ conflict other model and data. Reject 3 of 30.



Model error: targeted checks

- rate heterogeneity (space-time)
 - explicit (catastrophe) model shows some STRH
 - fit simulated data Gamma rate heterogeneity
- rate heterogeneity (across traits)
 - split data by columns into groups SW100/200, and refit (core 1-100 slower than 101-200)
 - exploratory check based on predictive replicates (tail of very slow traits)
 - fit simulated data
- borrowing/lateral transfer
 - fit simulated data (successful)
 - split data into two groups of language, refit (evidence for some LT)
- several others (*eg* meaning-category non-empty, data missing in blocks)

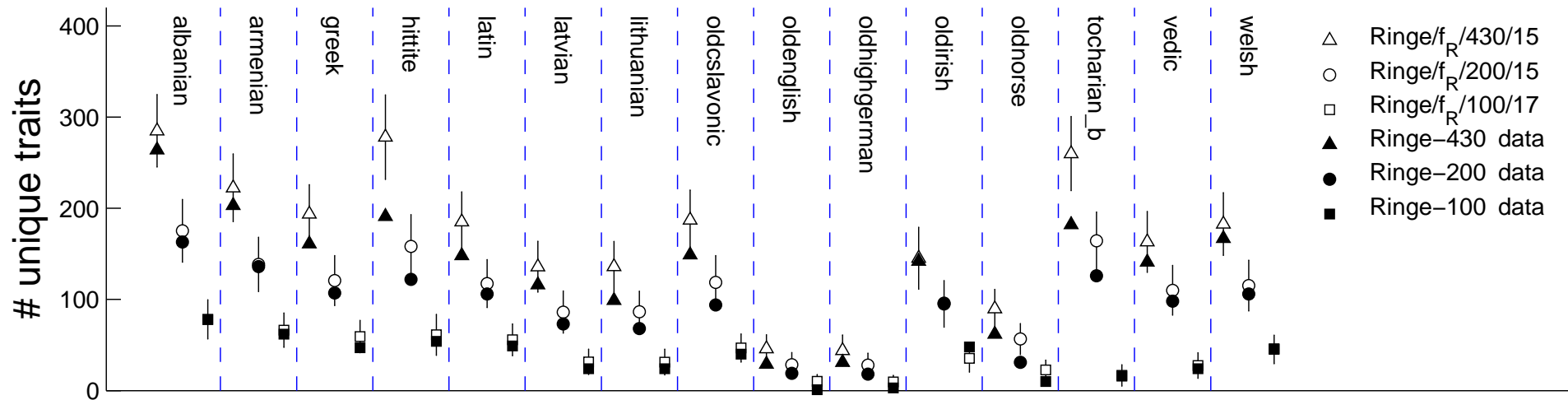
Model error: generic checks for misfit

External Data: remove singletons and predict # for each leaf.

D^* reserved data, $f(D)$ function of data, ψ parameters, predictive DBN

$$P(D'|D) = \int P(D'|\psi)p(\psi|D)d\psi$$

and estimate $P(f(D^*) > f(D')|D)$ (relevant here).



Posterior predictive distributions for the number of unique states, from the Ringe02 data: \triangle 328 meaning categories; \circ Swadesh-200 list; \square Swadesh 100 list.

Predictive replicates: estimate $P(f(D) > f(D')|D)$ for $f(D)$ function of data.

$Y^{(n)} = \#$ traits displayed at exactly n leaves.

y -axis $E(\tilde{Y}^{(n)}|D) - Y^{(n)}(D)$ and $\pm 2\text{std}(\tilde{Y}^{(n)}|D)$. x -axis $n = 2, 3, \dots, L$.

(TL) $S/MH50/U200$ (TR) $S/MH25/U200$ (LL) Ringe/ $f_R/328/15$ (LR) Ringe/ $f_R/100/17$.

