

# The Bayesian phylogenetics of grammar

Gerhard Jäger

Tübingen University

University of Amsterdam, April 4, 2018



WORDS BONES GENES TOOLS  
Tracking Linguistic, Cultural, and Biological Trajectories of the Human Past

UNIVERSITÄT  
TÜBINGEN



DFG



European Research Council  
Established by the European Commission

# Major word orders

# Statistics of major word order distribution

- data: WALS intersected with ASJP
- 1,045 languages, 211 lineages, 32 families with at least 5 languages

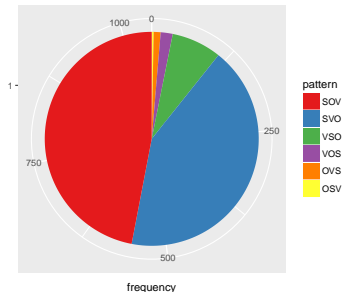
## Raw numbers

SOV	SVO	VSO	VOS	OVS	OSV
491	442	79	19	11	3
47.0%	42.3%	7.6%	1.8%	1.1%	0.3%

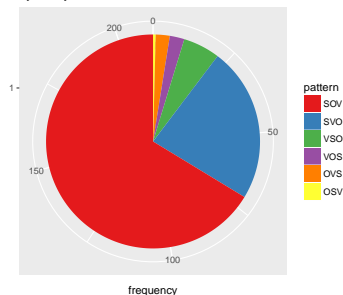
## Weighted by lineages

SOV	SVO	VSO	VOS	OVS	OSV
139.1	49.3	11.8	4.7	4.5	0.8
66.3%	23.4%	5.6%	2.2%	2.1%	0.4%

by language

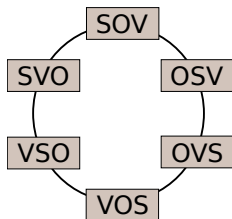


by family



# Previous approaches

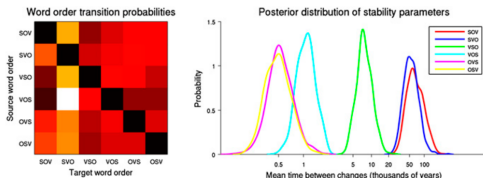
- Gell-Mann and Ruhlen (2011):
  - Proto-world was SOV
  - general pathway:  $SOV \rightarrow SVO \leftrightarrow VSO/VOS$
  - minor pathway:  $SOV \rightarrow OVS/OSV$
  - exceptions due to diffusion
- Ferrer-i-Cancho (2015):



- permutation circle
- transition probability inversely related to path length

# Previous approaches

- Maurits and Griffiths (2014):
  - Bayesian rate estimation, based on five families and NJ-trees

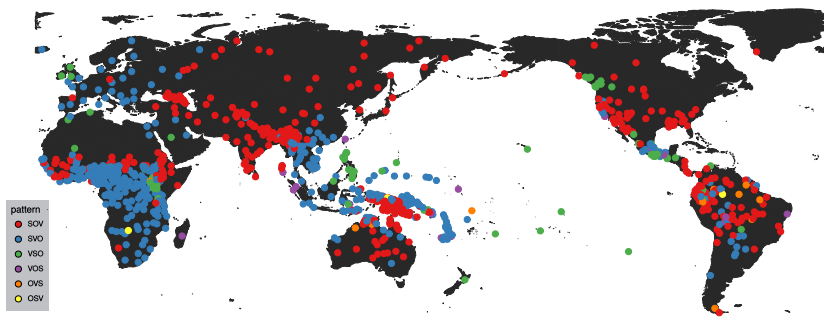


**Fig. 1.** Results of inferring a single mutation matrix  $Q$  for all six language families. (Left) Heat map showing the transition probabilities between word orders. Higher intensity (white, yellow) indicates more-probable transitions compared with lower intensity (red, brown), so SOV is most likely to transition to SVO and SVO to SOV. VSO is much more likely to transition to SVO than to SOV. (Right) Inferred posterior distributions of stability parameters for each word order. The horizontal axis shows the stability parameter, expressed as the mean time between transitions; i.e., higher values indicate a more stable word order.

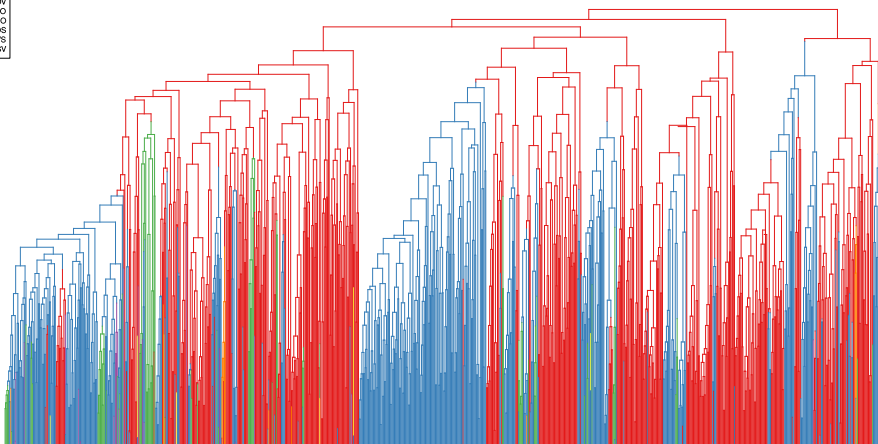
# Phylogenetic non-independence

- languages are phylogenetically structured
- if two closely related languages display the same pattern, these are not two independent data points

⇒ we need to control for phylogenetic dependencies



# Phylogenetic non-independence



# Phylogenetic non-independence

## Maslova (2000):

*“If the A-distribution for a given typology cannot be assumed to be stationary, a distributional universal cannot be discovered on the basis of purely synchronic statistical data.”*

*“In this case, the only way to discover a distributional universal is to **estimate transition probabilities** and as it were to ‘predict’ the stationary distribution on the basis of the equations in (1).”*

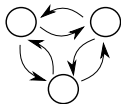




# The phylogenetic comparative method

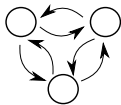
# Modeling language change

**Markov process**



# Modeling language change

Markov process

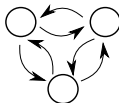


Phylogeny



# Modeling language change

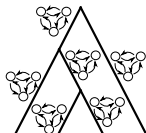
**Markov process**



**Phylogeny**

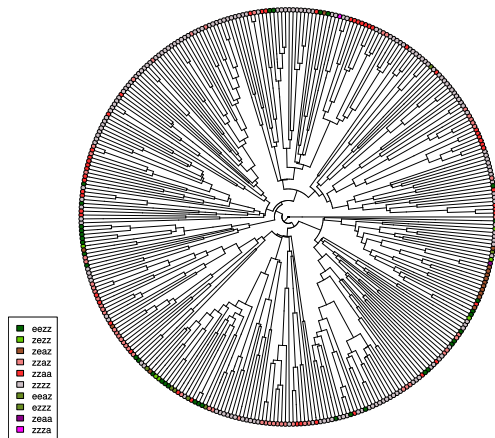


**Branching process**



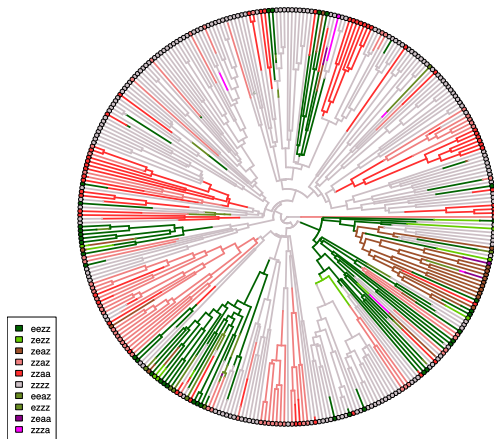
# Estimating rates of change

- if phylogeny and states of extant languages are known...



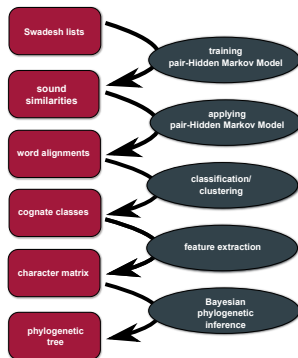
# Estimating rates of change

- if phylogeny and states of extant languages are known...
- ... transition rates and ancestral states can be estimated based on Markov model



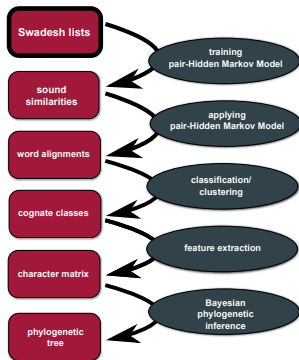
# Inferring a world tree of languages

# From words to trees

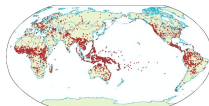




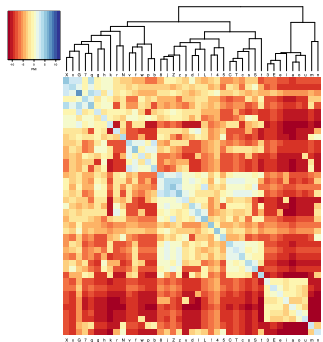
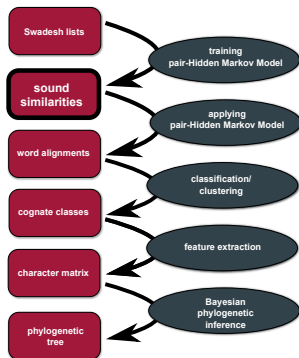
# From words to trees



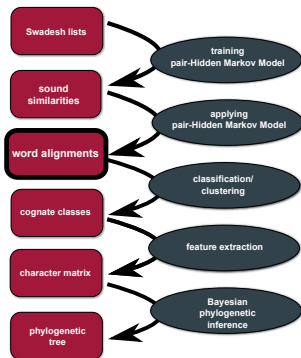
<i>concept</i>	Latin	English
<i>I</i>	ego	Ei
<i>you</i>	tu	yu
<i>we</i>	nos	wi
<i>one</i>	unus	w3n
<i>two</i>	duo	tu
<i>person</i>	persona, homo	pers3n
<i>fish</i>	piskis	fiS
<i>dog</i>	kanis	dag
<i>louse</i>	pedikulus	laus
<i>tree</i>	arbor	tri
<i>leaf</i>	foly~u*	lif
<i>skin</i>	kutis	skin
<i>blood</i>	saNgw~is	bl3d
<i>bone</i>	os	bon
<i>horn</i>	kornu	horn
<i>ear</i>	auris	ir
<i>eye</i>	okulus	Ei



# From words to trees

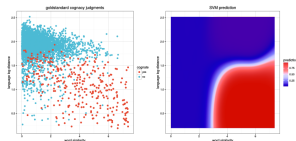
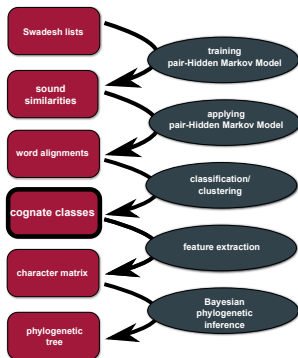


# From words to trees



Language	<i>ftsh:z</i>	<i>tongue:l</i>	<i>smoke:l</i>
Abui-Atangmelang	-af-u		
Abui-Fuimelang	-af-u	tal-l-fi--	
Adang	aab--	tal-E-b---	awai--b-a-n-o-7o-
Blagar-Bakalang	-ab--	--j-e-bur-	--ad--b-a-n-aKka-
Blagar-Bama	aab--	teg-e-bur-	-----b-e-n-a-xa-
Blagar-Kulijahi	-ab--	tej-e-bur-	-----b-e-n-aKka-
Blagar-Nule	aab--	tej-e-bur-	--ad--b-e-n-aKka-
Blagar-Tuntuli	aab--	tej-e-bur-	a-adge-b-a-n-a-q--
Blagar-Warsalelang	-ab--	tel-e-bur-	a-ad--b-a-n-a-x--
Bunaq			-----b-o-t-o-h--
Deing	haf--		-----buu-n-----
Hamap	7ab--	nar-g-bull-	-----b-a-n-o-7--
Kabola	hab--	tal-e-b---	awal--b-e-n-e-7o-
Kaera-Padangsul	-ab--	talee-b---	a-ad--b-e-naa-x--
Kafoa	-afU1	tal-l-p---	-----f-o-n-a-----
Kamang	-ap-1	nal--pu--	-----p-u-n-----a-
Kiraman	-Eb-	nal-l-bar-	--ar--b-a-n-o-kan
Klon	-eb-1	gel-E-b---	--ed-ab-o-n-----
Kui	-eb-	tal-l-ber-	--ar--b-o-n-o-k--
Kula	-ap-1	-il-l-p---	-----p--n--ekka-
Nedebang	asf-1	gel-e-fu--	--ar-ab-u-n-----
Reta	aab--	nal--bul-	a-ad--b-o-n-a----
Sar-Adiabang	haf--	--p-e-fal-	--ar--buu-n-----
Sar-Nule	haf--	nal-e-faj-	
Sawila	-ap-1	gal-impuru	-----p-u-n-a-ka-
Teiwa-Madar	xaf--	gel-i-vi--	-----buu-n-----
Wersing	-ap-1	nej-e-bur-	--ad-ap-u-n-a-k--
Wpantar	hap--	nal-e-bu--	-----b-unn-a----

# From words to trees



	English	Spanish	Modern Greek	Standard German
<i>I</i>	E1:A	yo:B	exo:C	ix:D
<i>you</i>	yu:A	ustet:B, tu:C	esi:D	du:E
<i>we</i>	wi:A	nosotro:B	emi:C	vir:A
<i>one</i>	w3n:A	uno:B	ena:C, ena:C	ain:D
<i>two</i>	tu:A	dos:B	Sy-o:C, Sio:D	cui:E
<i>person</i>	per3n:A	persona:A	an3-rope:B	m3n:S:C
<i>fish</i>	fi3:A	pezkado:A, pes:A	peari:B	fi3:A
<i>dog</i>	dag:A	pero:B	aTili:C, #Tiloo:C	hunt:D
<i>come</i>	k3e:A	veni:B	er3-o:C	kh-om3n:A
<i>sun</i>	s3n:A	sol:B	ily-3n:C, iloo:C	zon3:A
<i>star</i>	star:A	astroya:A	astari:A, astro:A	S3Er:A
<i>water</i>	wat3r:A	agn-a:B	nero:C	van3r:A
<i>stone</i>	ston:A	piedra:B	petra:B	Stain:A
<i>fire</i>	fi3r:A	fuego:B	foty-a:C	foia:D
<i>path</i>	p3r:A	senda:B	Uromos:C	pf-at:A, vek:D
<i>moon/chain</i>	maunt3n:A	sero:B, nonta3a:A	vuno:C, oros:D	b3rk:E
<i>full</i>	ful:A	yeno:B	yematos:C, pliria:D	fol:A
<i>new</i>	nu:A	nuevo:A	neos:A, Ternary-3n:B	noi:A
<i>name</i>	nea:A	nombre:A	onoma:A	naa3:A





# Estimating word-order transition patterns

# Workflow

(data from all 32 families with  $\geq 5$  languages in data base; 778 languages in total)

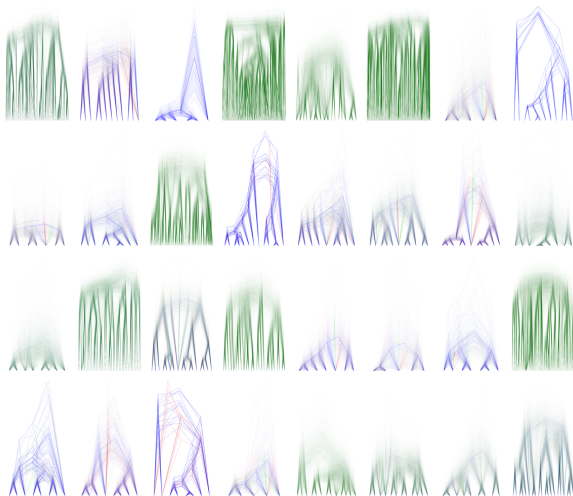
- estimate posterior tree distributions with MrBayes for each family, using Glottolog as constraint tree
- test whether universal or lineage-specific model gives a better fit
- estimate transition rates with best model
- estimate stationary distribution of major word order categories
- apply *stochastic character mapping* (SIMMAP; Bollback 2006)
- estimate expected number of mutations for each transition type



# Estimating posterior tree distributions

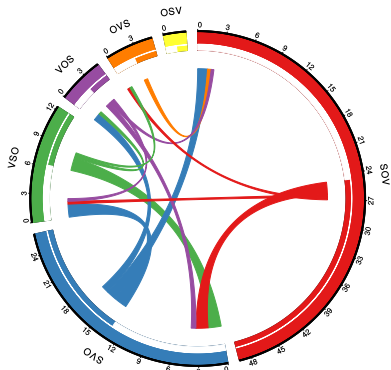
- using characters extracted from ASJP data (Jäger 2018)
- Glottolog as constraint tree
- $\Gamma$ -distributed rates
- ascertainment bias correction
- relaxed molecular clock (IGR)
- uniform tree prior
- stop rule: 0.01, samplefreq=1000
- if convergence later than after 1,000,000 steps, sample 1,000 trees from posterior

# Phylogenetic tree sample



# Estimating transition rates

- totally unrestricted model, all 30 transition rates are estimated independently
- implementation using RevBayes (Höhna et al., 2016)



# Reconstruction history with SIMMAP

- estimated frequency of mutations within the 32 families under consideration (posterior mean, 100 iterations)

	<b>SOV</b>	<b>SVO</b>	<b>VSO</b>	<b>VOS</b>	<b>OVS</b>	<b>OSV</b>
<b>SOV</b>	–	20.2	3.2	0.5	3.3	0.4
<b>SVO</b>	17.6	–	23.9	14.5	1.5	1.1
<b>VSO</b>	1.5	19.9	–	2.5	1.8	0.4
<b>VOS</b>	1.0	5.4	2.3	–	0.9	0.3
<b>OVS</b>	2.8	0.9	0.6	0.4	–	0.2
<b>OSV</b>	0.5	0.5	0.4	0.3	0.5	–

# Refining the model with Reversibly Jump MCMC

- Estimating 30 transition rates is a tall order, given that the data possibly only reflect about 130 transition events
- hand-crafted sub-model construction: time consuming, subjective and error prone
- solution: posterior sampling over sub-models using *Reversible Jump Markov Chain Monte Carlo* (RJMCMC, Green 1995)

## RJMCMC

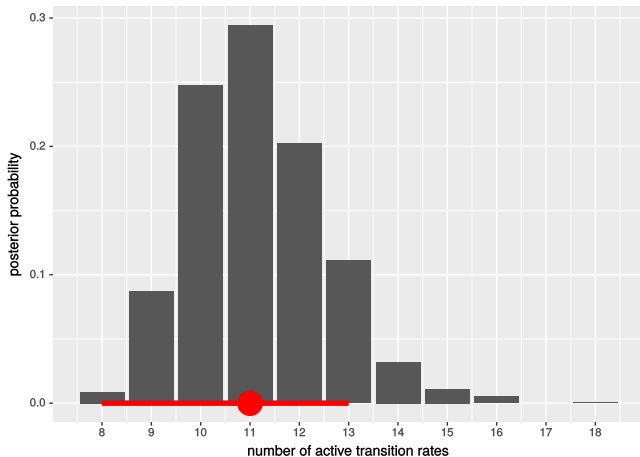
RJMCMC assumes a prior distribution over sub-models (where some transition rates are set to 0) and simultaneously samples from the set of sub-models and the parameter spaces of the sub-models.

# Model comparison

model	marginal likelihood	AICM
<i>lineage-specific</i>	$-423.0 \pm 0.08$	$926.4 \pm 0.5$
<i>circular GTR</i>	$-420.0 \pm 1.72$	$851.7 \pm 1.6$
<i>circular</i>	$-414.2 \pm 0.72$	$851.6 \pm 2.1$
<i>RJ/GTR</i>	$-413.4 \pm 2.96$	$855.9 \pm 4.7$
<i>unrestricted</i>	$-406.7 \pm 0.78$	$846.4 \pm 2.5$
<i>unrestricted GTR</i>	$-404.4 \pm 0.89$	$843.5 \pm 3.6$
<i>RJ</i>	$-398.0 \pm 0.57$	$827.2 \pm 2.1$

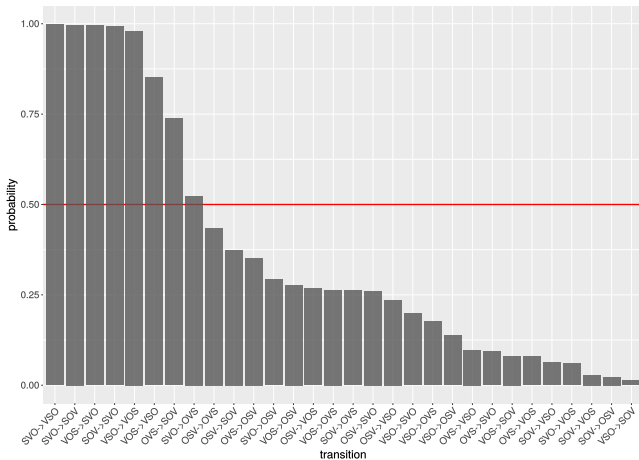
# Refining the model with Reversibly Jump MCMC

Number of active transition rates: posterior distribution



# Refining the model with Reversibly Jump MCMC

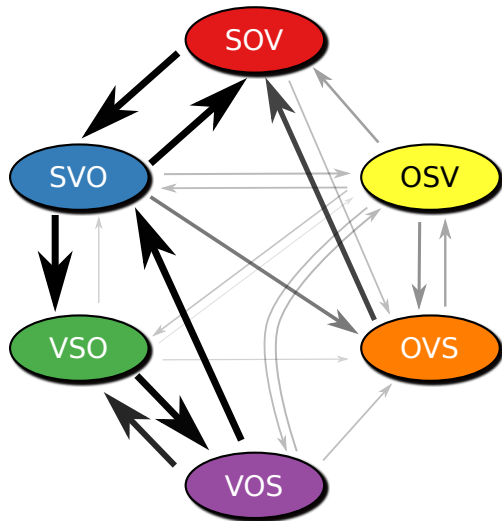
## Probabilities of active transition rates: posterior distribution





# Refining the model with Reversibly Jump MCMC

Probabilities of active transition rates: posterior distribution



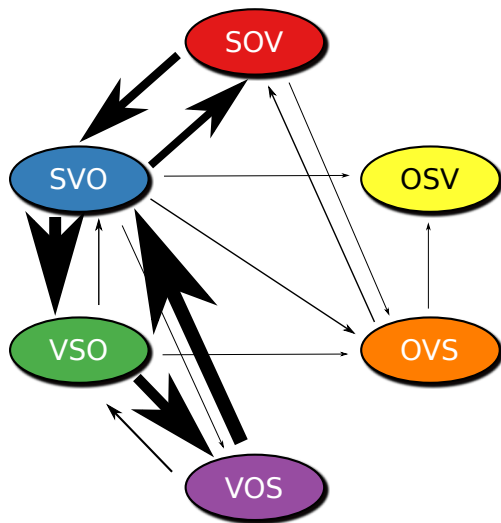
# Reconstruction history with SIMMAP

- estimated frequency of mutations within the 32 families under consideration (posterior mean, 99 iterations)

	SOV		SVO		VSO		VOS		OVS		OSV	
<b>SOV</b>	–		23.1	[14; 30]	0.5	[0; 6]	0.1	[0; 0]	1.9	[0; 9]	0.1	[0; 0]
<b>SVO</b>	20.3	[16; 28]	–		33.0	[20; 45]	2.2	[0; 29]	3.4	[0; 11]	1.2	[0; 7]
<b>VSO</b>	0.0	[0; 0]	3.8	[0; 25]	–		29.7	[0; 46]	1.5	[0; 9]	0.5	[0; 4]
<b>VOS</b>	0.1	[0; 0]	38.3	[19; 54]	6.2	[0; 13]	–		0.9	[0; 5]	0.4	[0; 2]
<b>OVS</b>	4.0	[0; 10]	0.5	[0; 3]	0.9	[0; 6]	0.2	[0; 1]	–		1.1	[0; 6]
<b>OSV</b>	0.7	[0; 6]	0.3	[0; 3]	0.4	[0; 3]	0.6	[0; 5]	0.9	[0; 7]	–	

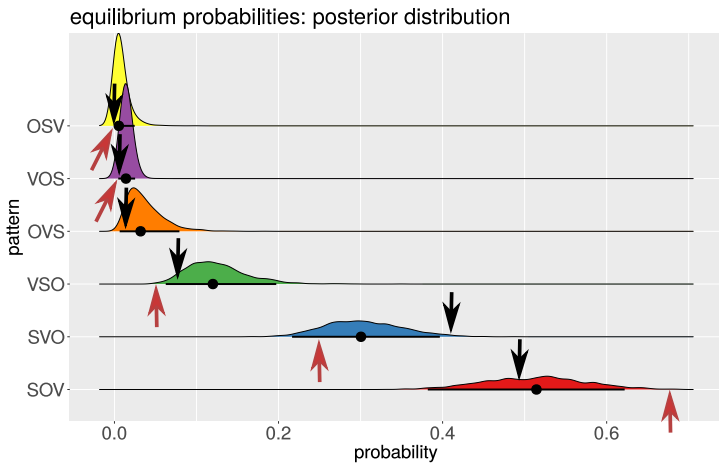
# Reconstruction history with SIMMAP

Expected frequencies of transitions: posterior mean



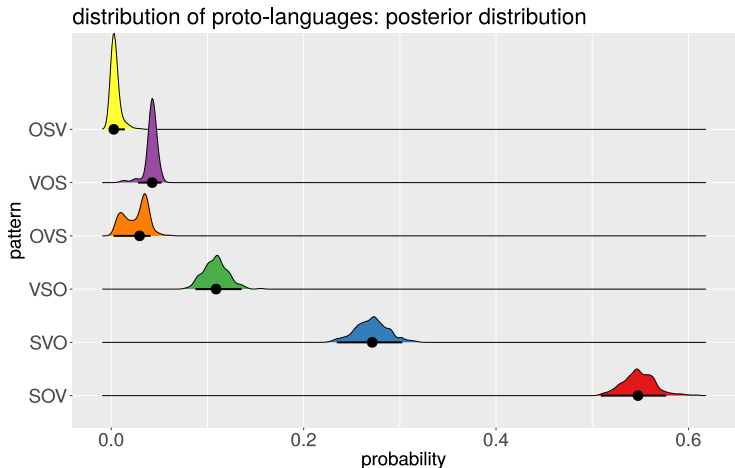
# Posterior distributions

## Empirical vs. estimated distribution



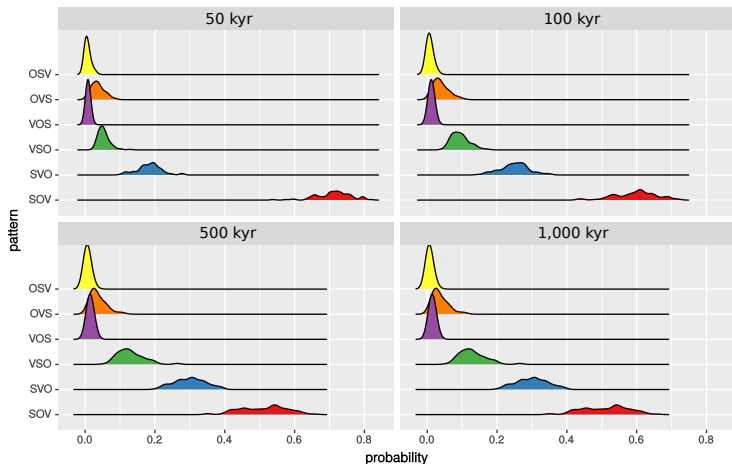
# Posterior distributions

## Expected distribution of Proto-languages



# Posterior distributions

Expected probabilities of Proto-World, given that we can demonstrate SOV for all proto-languages



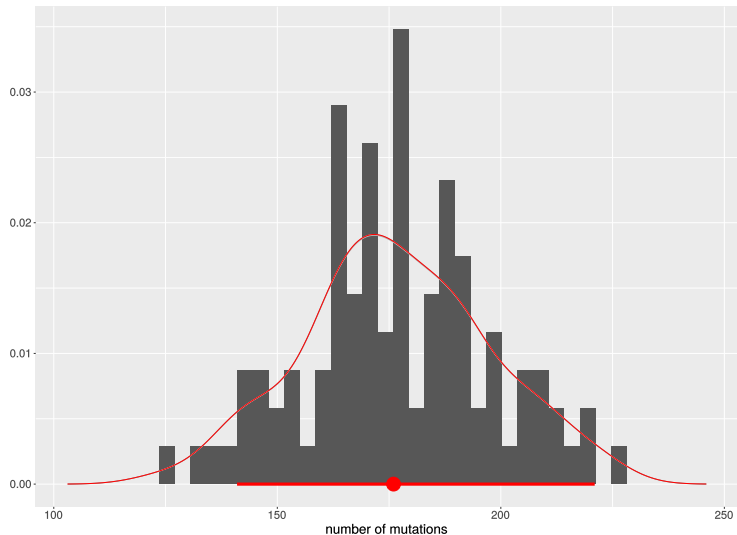
# Posterior distributions

## Waiting times



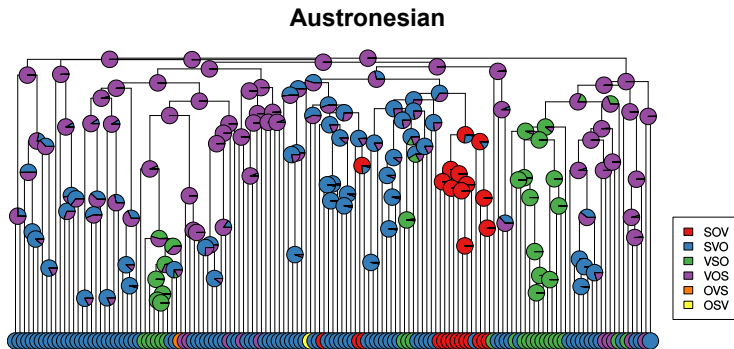
# Posterior distributions

## Number of state changes



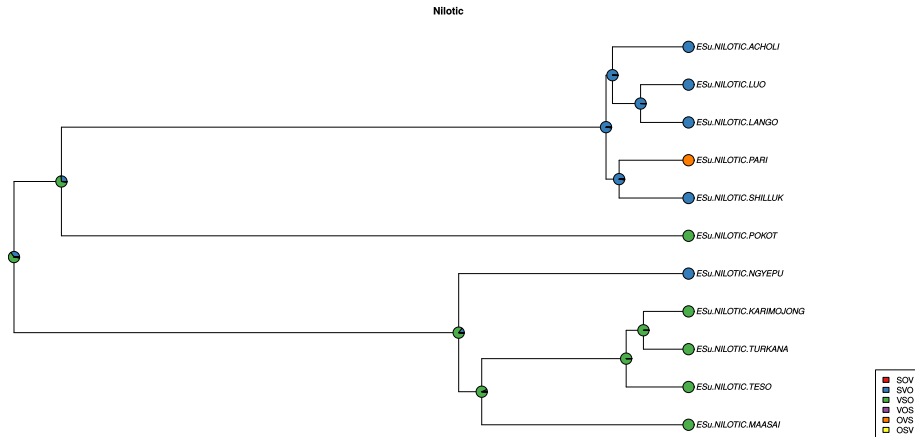


# Ancestral state reconstruction



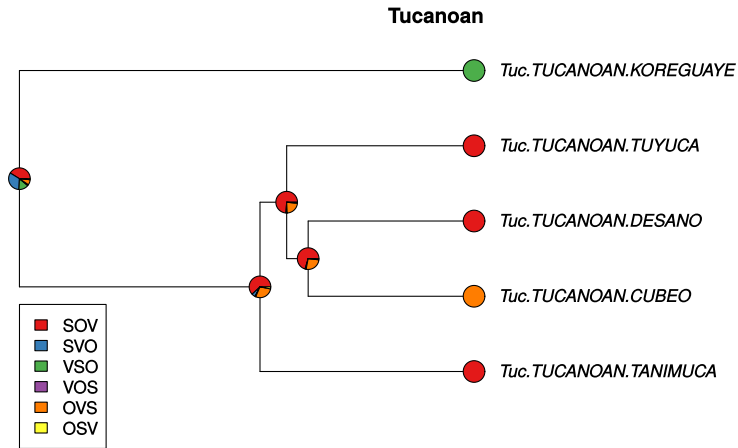
# Examples for unexpected transitions

SVO  $\rightarrow$  OVS



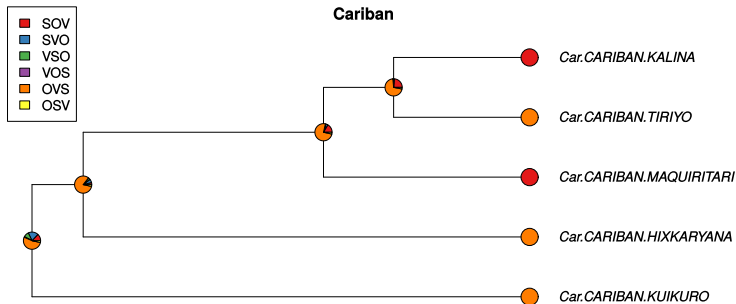
# Examples for unexpected transitions

OVS → SOV



# Examples for unexpected transitions

OVS → SOV



# Summary

- no evidence for general preference of SOV  $\rightarrow$  SVO over the reverse
- SVO is currently over-represented due to recent spread of Austronesian and Atlantic-Congo, but not excessively so
- multiple counter-evidence to Ramon-i-Ferrer's and Gell-Mann & Ruhlen's models

- Jonathan P. Bollback. SIMMAP: stochastic character mapping of discrete traits on phylogenies. *BMC Bioinformatics*, 7(1):88, 2006.
- Ramon Ferrer-i-Cancho. Kauffman's adjacent possible in word order evolution. arXiv preprint arXiv:1512.05582, 2015.
- Murray Gell-Mann and Merritt Ruhlen. The origin and evolution of word order. *Proceedings of the National Academy of Sciences*, 108(42):17290–17295, 2011.
- Peter J. Green. Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika*, 82(4): 711–732, 1995.
- Sebastian Höhna, Michael J. Landis, Tracy A. Heath, Bastien Boussau, Nicolas Lartillot, Brian R. Moore, John P. Huelsenbeck, and Frederik Ronquist. Revbayes: Bayesian phylogenetic inference using graphical models and an interactive model-specification language. *Systematic biology*, 65(4):726–736, 2016.
- Gerhard Jäger. Phylogenetic inference from word lists using weighted alignment with empirically determined weights. *Language Dynamics and Change*, 3(2):245–291, 2013.
- Gerhard Jäger. Support for linguistic macrofamilies from weighted sequence alignment. *Proceedings of the National Academy of Sciences*, 112(41):12752–12757, 2015. doi: 10.1073/pnas.1500331112.
- Gerhard Jäger. Global-scale phylogenetic linguistic inference from lexical resources. arXiv:1802.06079, 2018.
- Gerhard Jäger and Søren Wichmann. Inferring the world tree of languages from word lists. In S. G. Roberts, C. Cuskley, L. McCrohon, L. Barceló-Coblijn, O. Feher, and T. Verhoef, editors, *The Evolution of Language: Proceedings of the 11th International Conference (EVLANG11)*, 2016. Available online: <http://evolang.org/neworleans/papers/147.html>.
- Elena Maslova. A dynamic approach to the verification of distributional universals. *Linguistic Typology*, 4(3):307–333, 2000.
- Luke Maurits and Thomas L. Griffiths. Tracing the roots of syntax with Bayesian phylogenetics. *Proceedings of the National Academy of Sciences*, 111(37):13576–13581, 2014.
- Søren Wichmann, Eric W. Holman, and Cecil H. Brown. The ASJP database (version 17). <http://asjp.clld.org/>, 2016.