

# Typologies in equilibrium

Gerhard Jäger

Tübingen University

SLE 2018, Tallinn University



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

UNIVERSITÄT  
TÜBINGEN



DFG

# Distributional universals

# Distributional universals

- common practice since Greenberg (1963):
  - collect a sample of languages
  - classify them according to some typological feature

⇒ skewed distribution indicates something interesting going on
- Problem: languages are not independent samples
- skewed distribution may reflect
  - skewed diversification rate across families
  - properties of an ancestral bottleneck
- balanced sampling mitigates the first, but not the second problem

# Distributional universals

## 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.”*

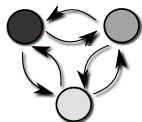
*“In this case, the only way to discover a distributional universal is to **estimate transition probabilities** and as it were to ‘predict’ the stationary distribution on the basis of the equations in (1).”*



# The phylogenetic comparative method

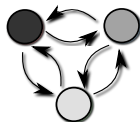
# Modeling language change

**Markov process**



# Modeling language change

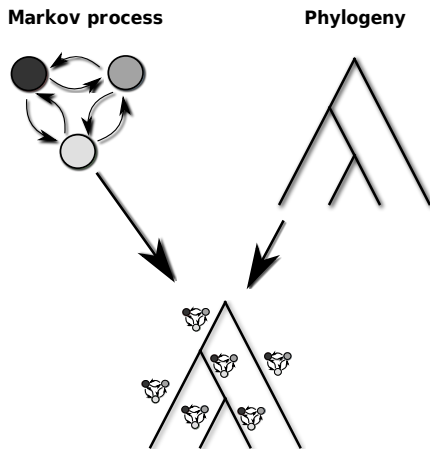
**Markov process**



**Phylogeny**

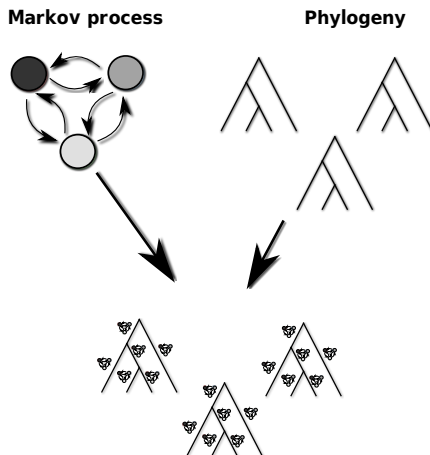


# Modeling language change





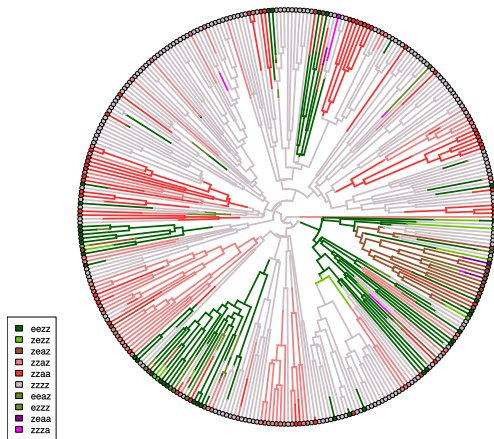
# Modeling language change





# 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



# Major word orders

# Statistics of major word order distribution

- data: WALS intersected with ASJP
- 1,055 languages, 201 lineages, 71 families with at least 3 languages

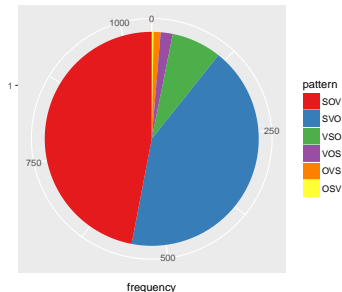
## Raw numbers

SOV	SVO	VSO	VOS	OVS	OSV
497	447	78	20	10	3
47.1%	42.4%	7.4%	1.9%	0.9%	0.3%

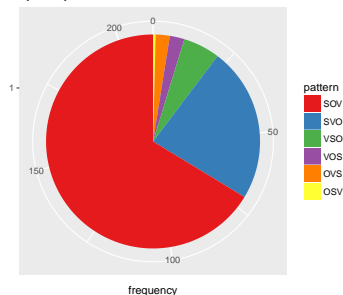
## Weighted by lineages

SOV	SVO	VSO	VOS	OVS	OSV
135.1	46.9	10.5	4.0	3.7	0.8
67.2%	23.3%	5.2%	2.0%	1.8%	0.4%

by language



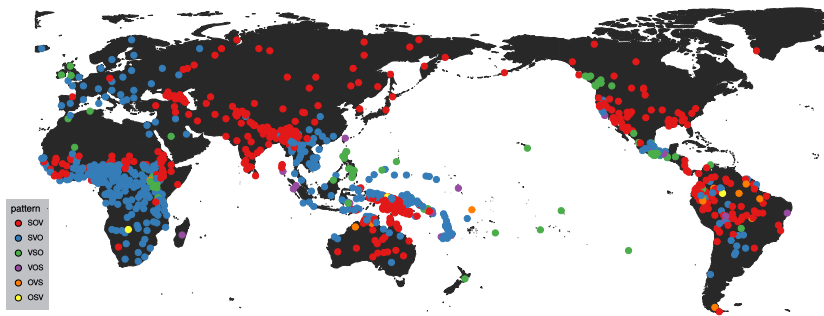
by family



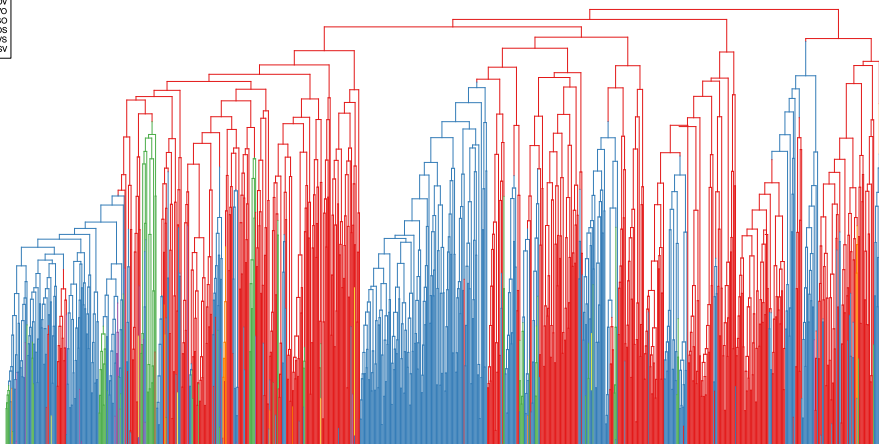
# Phylogenetic non-independence

- 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



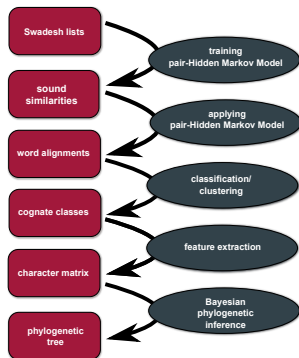
# Phylogenetic non-independence



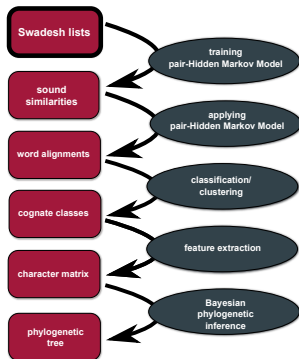
# Inferring trees across many families



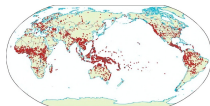
# From words to trees



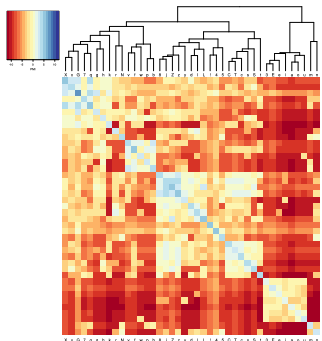
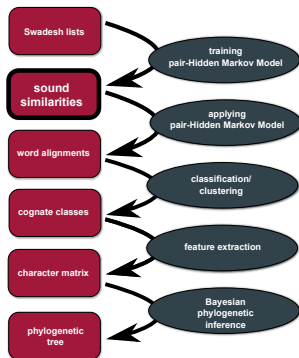
# From words to trees



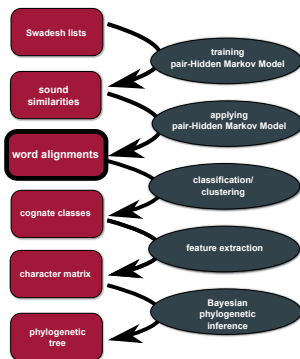
<i>concept</i>	Latin	English
<i>I</i>	ego	Ei
<i>you</i>	tu	yu
<i>we</i>	nos	wi
<i>one</i>	unus	w3n
<i>two</i>	duo	tu
<i>person</i>	persona, homo	pers3n
<i>fish</i>	piskis	fiS
<i>dog</i>	kanis	dag
<i>louse</i>	pedikulus	laus
<i>tree</i>	arbor	tri
<i>leaf</i>	foly~u*	lif
<i>skin</i>	kutis	skin
<i>blood</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

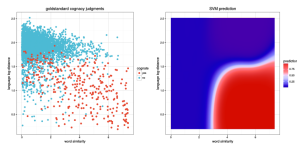
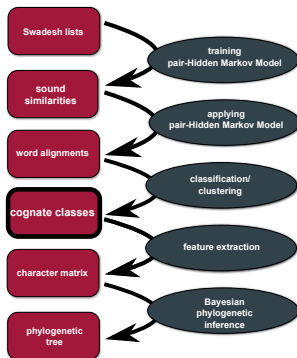


# From words to trees



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-buH-	-----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	aaf-1	gel-e-fu--	--ar-ab-u-n-----
Reta	aab--	nal--buH-	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----

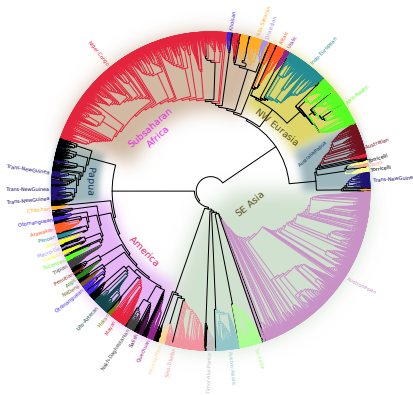
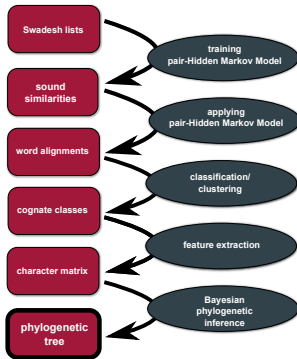
# From words to trees



	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	nosotroes:B	emsi:C	vir:A
<i>one</i>	w3n:A	uno:B	emas:C, ema:C	ains:D
<i>two</i>	tu:A	dos:B	Sy-o:C, Sio:D	cui:E
<i>person</i>	per3n:A	persona:A	an3-ropos: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-o3-C, iLos:C	zon3:A
<i>star</i>	star:A	astroya:A	astari:A, astro:A	S3Er3:A
<i>water</i>	wat3r:A	agn-a3:B	nero:C	van3r:A
<i>stone</i>	ston:A	piedra:B	petra:B	Stain:A
<i>fire</i>	fi3r:A	fuego:B	foty-a3:C	foia:D
<i>path</i>	p3r:A	senda:B	Uromos:C	pf-at3:A, vek3:D
<i>moon/chain</i>	maunt3n:A	sero:B, nonta3a:A	vuno:C, oros:D	b3k3:E
<i>full</i>	ful:A	yeno:B	yematos:C, pliria:D	fol:A
<i>new</i>	nu:A	nuevo:A	neos:A, Ternary-oo3:B	noi:A
<i>name</i>	nea:A	nombre:A	onoma:A	naa3:A



# From words to trees



# Estimating word-order transition patterns



# Workflow

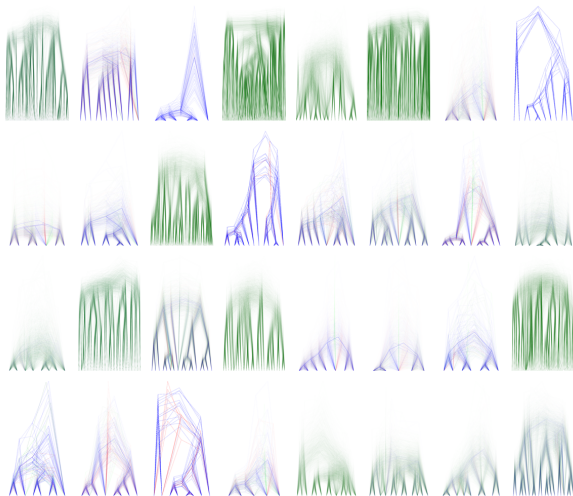
(data from all 77 families with  $\geq 3$  languages in data base; 924 languages in total)

- estimate posterior tree distributions with MrBayes for each family, using Glottolog as constraint tree
- estimate transition rates
- 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

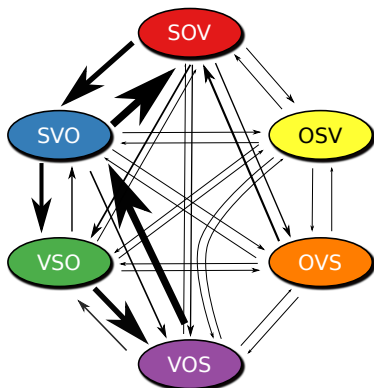
# Phylogenetic tree sample



# Estimating transition rates

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

expected strength of flow



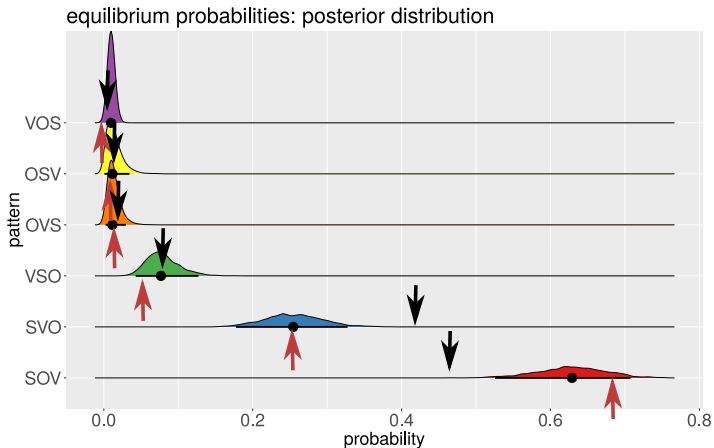
# Reconstruction history with SIMMAP

- estimated frequency of mutations within the 77 families under consideration (posterior mean and 95% HPD, 100 simulations)

	SOV		SVO		VSO		VOS		OVS		OSV	
<b>SOV</b>	–		51.5	[19; 82]	10.2	[1; 19]	7.5	[0; 29]	5.8	[0; 14]	4.2	[0; 13]
<b>SVO</b>	83.8	[31; 131]	–		22.3	[2; 42]	10.4	[0; 30]	2.8	[0; 8]	3.9	[0; 12]
<b>VSO</b>	1.4	[0; 5]	8.3	[0; 24]	–		29.0	[5; 45]	3.0	[0; 9]	1.1	[0; 5]
<b>VOS</b>	4.3	[0; 15]	141.9	[115; 188]	30.9	[17; 47]	–		2.1	[0; 9]	1.0	[0; 3]
<b>OVS</b>	11.1	[0; 28]	0.8	[0; 4]	1.8	[0; 8]	0.4	[0; 3]	–		0.8	[0; 5]
<b>OSV</b>	4.2	[0; 15]	0.4	[0; 3]	1.9	[0; 11]	1.1	[0; 7]	1.1	[0; 9]	–	

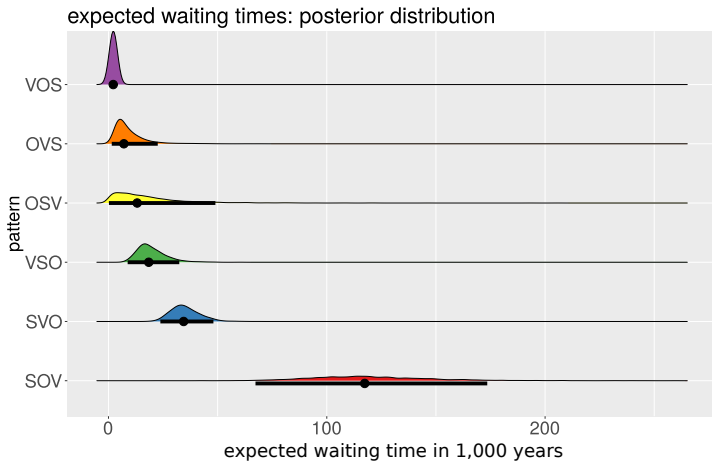
# Posterior distributions

## Empirical vs. estimated distribution



# Posterior distributions

## Waiting times



# Differential case marking



# Statistics of differential case marking distribution

- data: Autotyp intersected with ASJP
- 343 languages, 115 lineages, 29 families with at least 3 languages

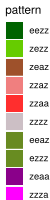
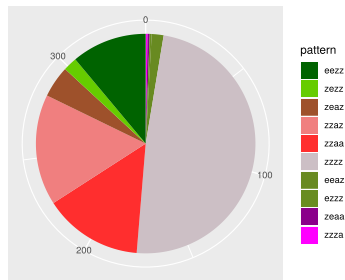
## Raw numbers

ZZZZ	ZZaa	ZZaz	eeZZ	zeaz	eeaz	eZZZ	zeZZ	zeaa	ZZZa
167	50	56	38	16	6	1	7	1	1
48.7%	14.6%	16.3%	11.1%	4.7%	1.7%	0.3%	2.0%	0.3%	0.3%

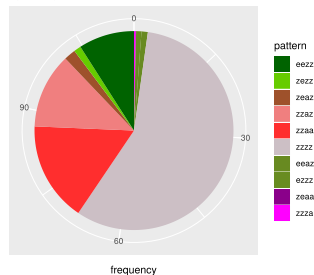
## Weighted by lineages

ZZZZ	ZZaa	ZZaz	eeZZ	zeaz	eeaz	eZZZ	zeZZ	zeaa	ZZZa
65.8	18.6	14.0	10.4	2.2	1.2	1.0	1.3	0.1	0.3
57.2%	16.2%	12.2%	9.1%	1.9%	1.0%	0.9%	1.2%	0.0%	0.3

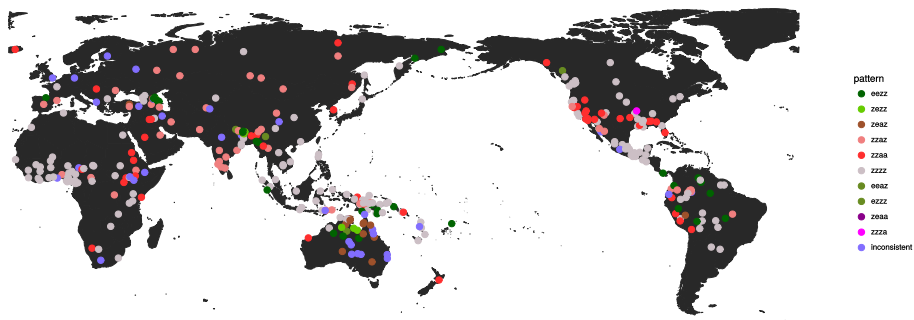
by language



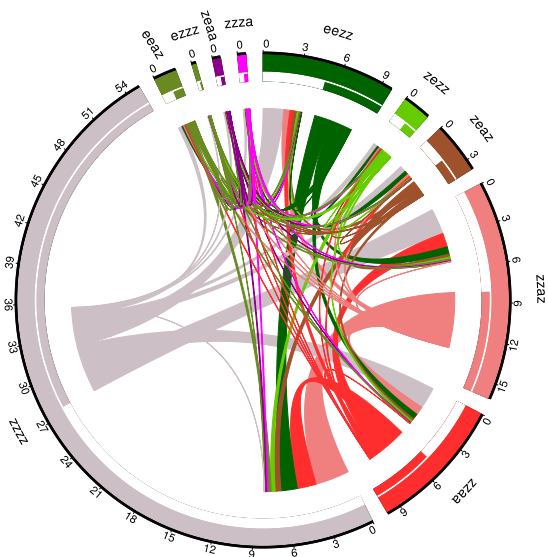
by family



# Geographic distribution



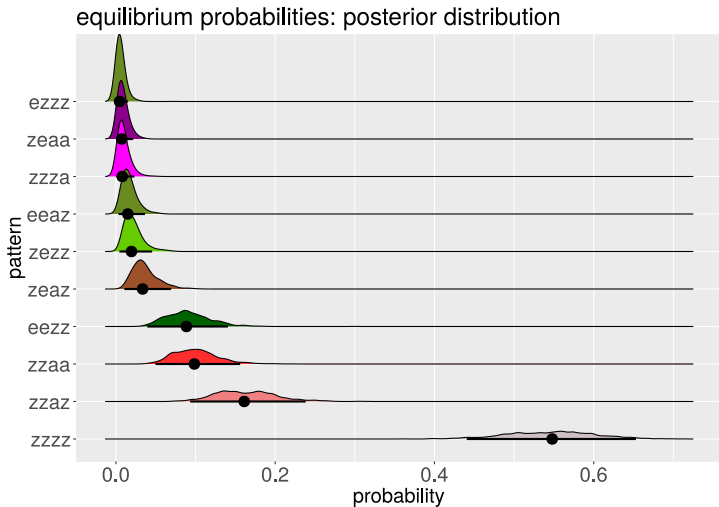
# Estimated rates



- green: nominative/accusative
- red: ergative/absolutive
- grey: no case marking
- brown/purple: mixture of accusative and ergative patterns

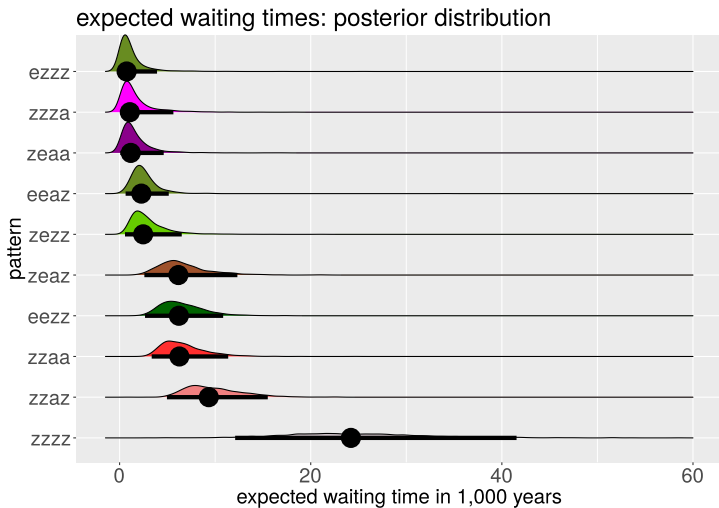
# Posterior distributions

## Estimated distribution



# Posterior distributions

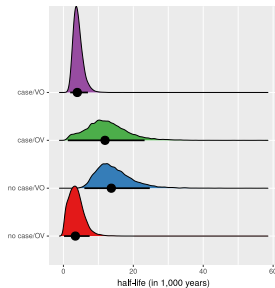
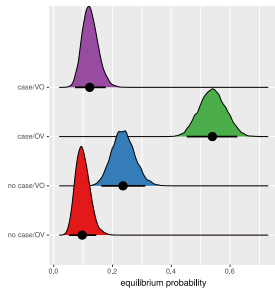
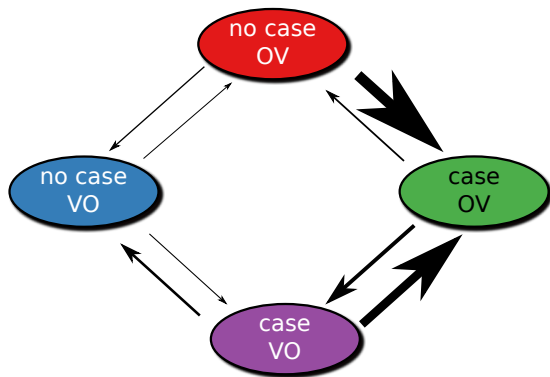
## Waiting times



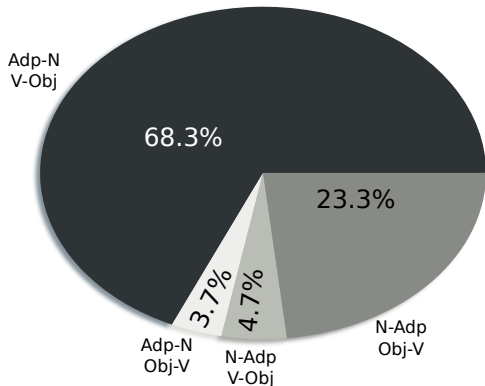
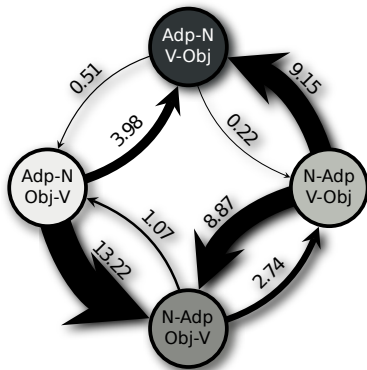
# Posterior distribution

- with 99.7% confidence, accusative languages are more likely than ergative languages
- with 59.6% confidence, consistent accusative is more likely than consistent ergative
- with 100% confidence, differential object marking is more likely than anti-DOM
- with 95.2% confidence, differential subject marking is more likely than anti-DSM

# Further variables



## Further variables





# Conclusion

- Maslova's program can be carried out with phylogenetic comparative method
- future research:
  - equilibrium distributions generally resemble family-wise weighted distributions — bug or feature?
  - hierarchical models instead of one Markov process for all lineages?
  - more data!!! (but there are never enough of them)
  - better methods for feature selection? (Bayes factor test rejected RJ-MCMC)

- Jonathan P. Bollback. SIMMAP: stochastic character mapping of discrete traits on phylogenies. *BMC Bioinformatics*, 7(1):88, 2006.
- Joseph Greenberg. Some universals of grammar with special reference to the order of meaningful elements. In *Universals of Language*, pages 73–113. MIT Press, Cambridge, MA, 1963.
- 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. Global-scale phylogenetic linguistic inference from lexical resources. arXiv:1802.06079, 2018.
- Elena Maslova. A dynamic approach to the verification of distributional universals. *Linguistic Typology*, 4(3):307–333, 2000.
- Søren Wichmann, Eric W. Holman, and Cecil H. Brown. The ASJP database (version 18). <http://asjp.clld.org/>, 2018.