

Understanding Language Evolution Using an Event-Based Model

John Huelsenbeck
Department of Integrative Biology
johnh@berkeley.edu

Introduction

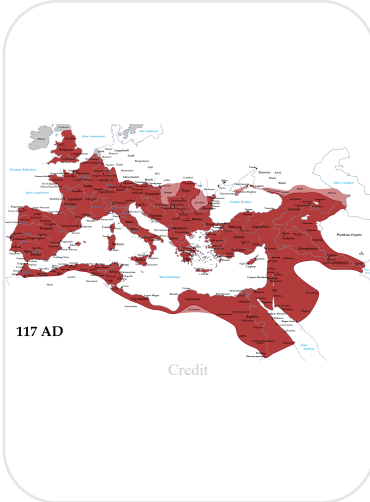
Modern languages are related to one another through a complicated history of divergence and word borrowing. The divergence of languages is caused by the slow change in spoken language as it is passed from parents to offspring. Over time, divergence causes languages to become increasingly different from one another, ultimately to the point where they are mutually unintelligible. Languages that were spoken by the same human group more recently in time are considered to be more closely related to each other than they are to groups that spoke the language more distantly in time; this relatedness information can be depicted by a tree-like diagram called a ‘phylogeny.’ Linguistic borrowing, by contrast, causes languages to become more similar to one another.

Language	IPA	Coding
English	/hænd/	0
German	/hant/	0
French	/mẽ/	1
Spanish	/mano/	1
Italian	/ma:no/	1
Russian	/rɔka/	2
Polish	/rɛŋka/	2

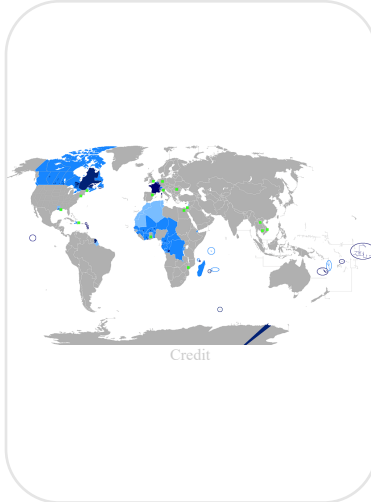
Table 1. Coding of lexical cognates for the word *hand*.

Languages

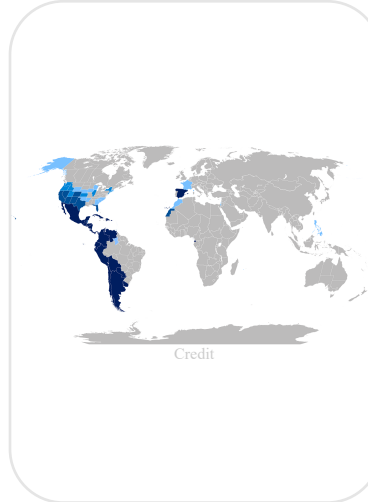
Latin



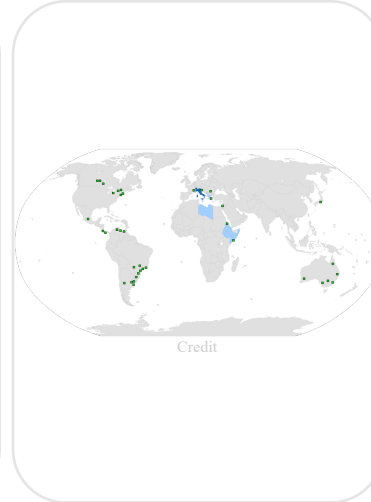
French



Spanish



Italian



Brazilian Portuguese



Portuguese



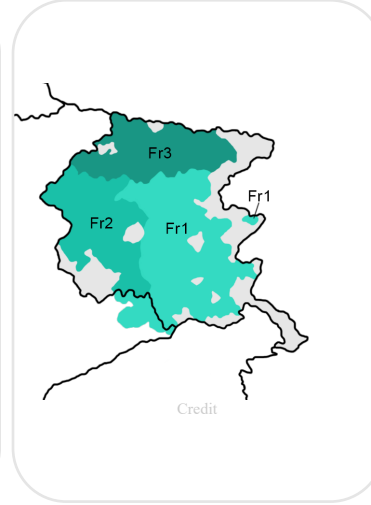
Catalan



Walloon



Friulian

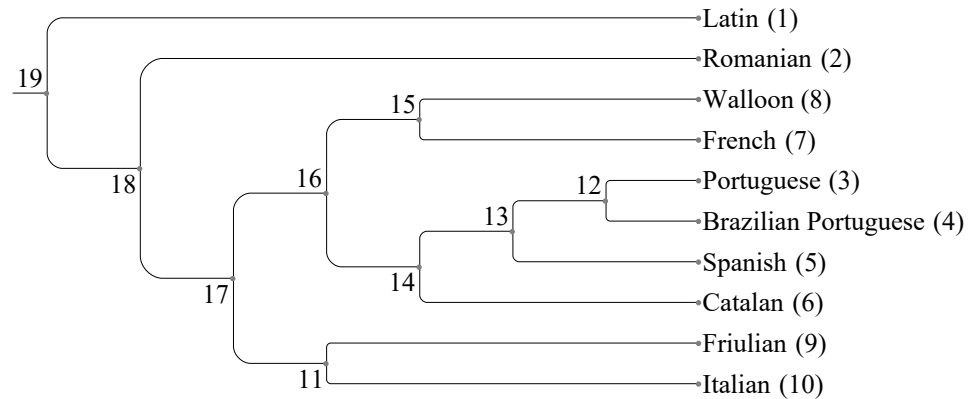


Romanian



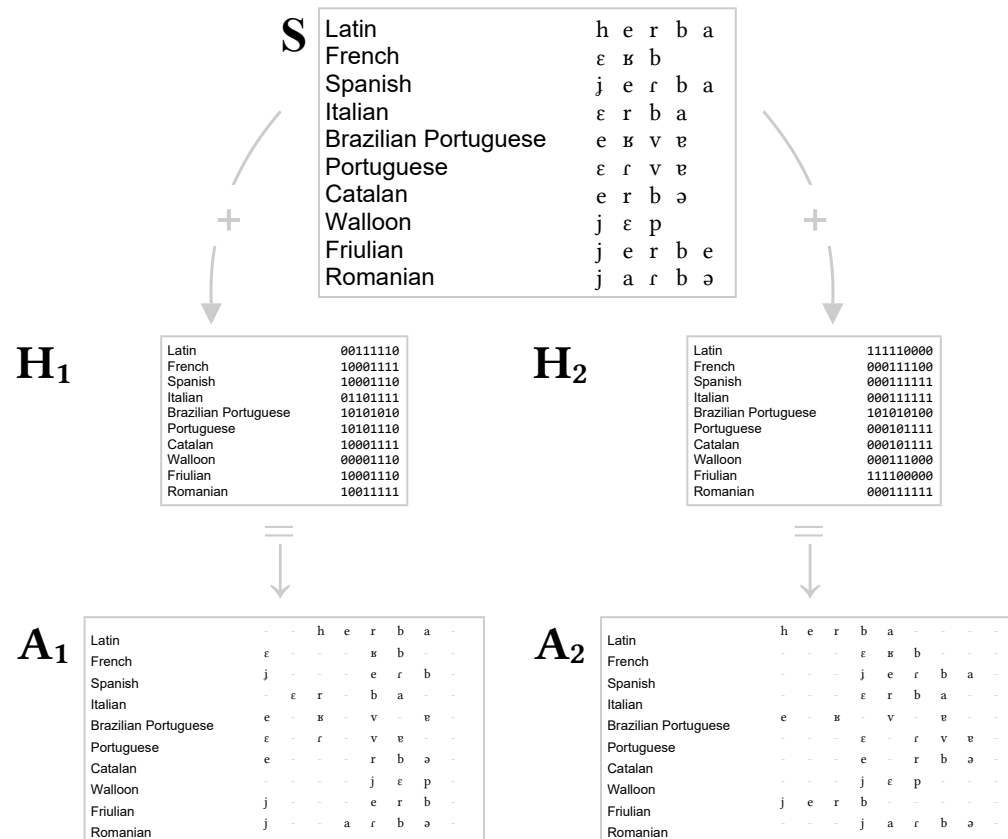
Example Tree

An example tree showing the relationships of $N = 10$ languages.



Alignment

Alignments (**A**) are formed from the observed segments (**S**) and a homology map (**H**).

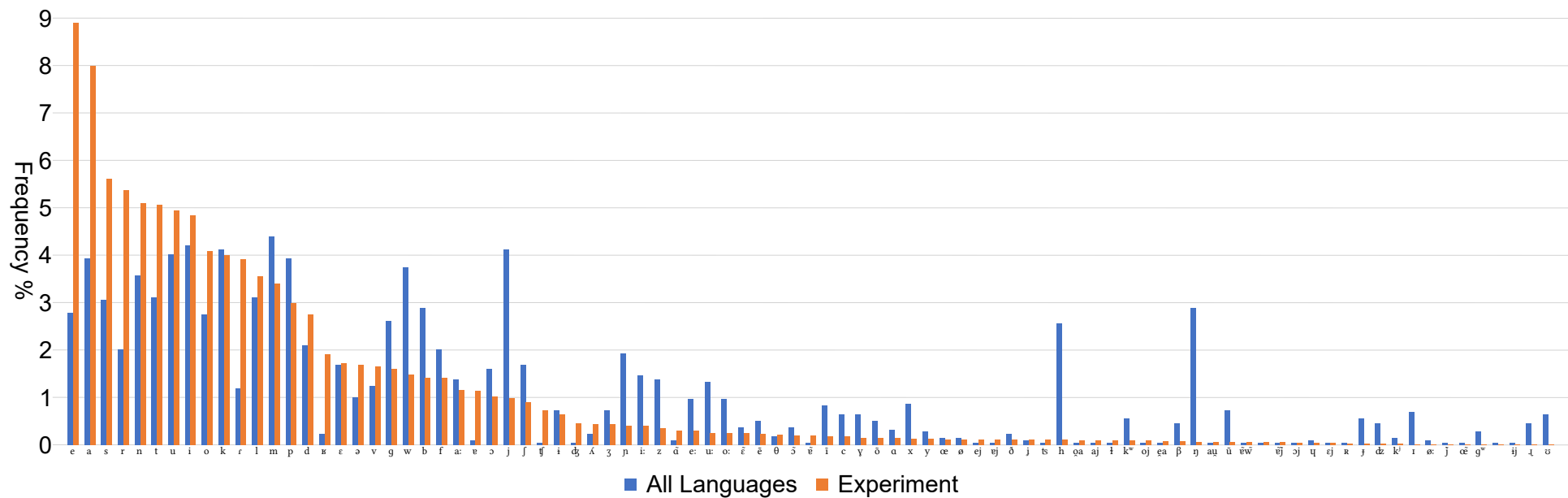


Character Assignments

Each segment gets a different number

1	f	2	e:	3	m	4	i	5	n	6	a	7	e	8	b	9	r	10	ɐ	11	j	12	œ	13	ɛ	14	u	15	l	16	r	17	x	18	o	19	ʎ	20	ʁ
21	i:	22	d	23	ɲ	24	ɔ	25	ə	26	ɔa	27	t	28	ã	29	ĩ	30	ẽ	31	ʈʂ	32	ɨ	33	ẽj	34	s	35	p	36	z	37	ʃ	38	ð	39	o:		
40	k	41	ts	42	a:	43	ẽj	44	aj	45	θ	46	ej	47	ɐj	48	w	49	õ	50	h	51	õ	52	v	53	au	54	ẽ	55	ẽw	56	c						
57	ij	58	ẽ	59	ɬ	60	ɑ	61	u:	62	ɥ	63	y	64	g	65	iw	66	ɔj	67	dz	68	j	69	dʒ	70	oj	71	g ^w	72	ŋ	73	ø	74	β				
75	ɕa	76	ʀ	77	ɛj	78	k ^j	79	ɪ	80	ʒ	81	ʃ	82	ɣ	83	k ^w	84	ũ	85	ɭ	86	ʝ	87	œ	88	ø:	89	ʊ										

Prior Segment Frequencies



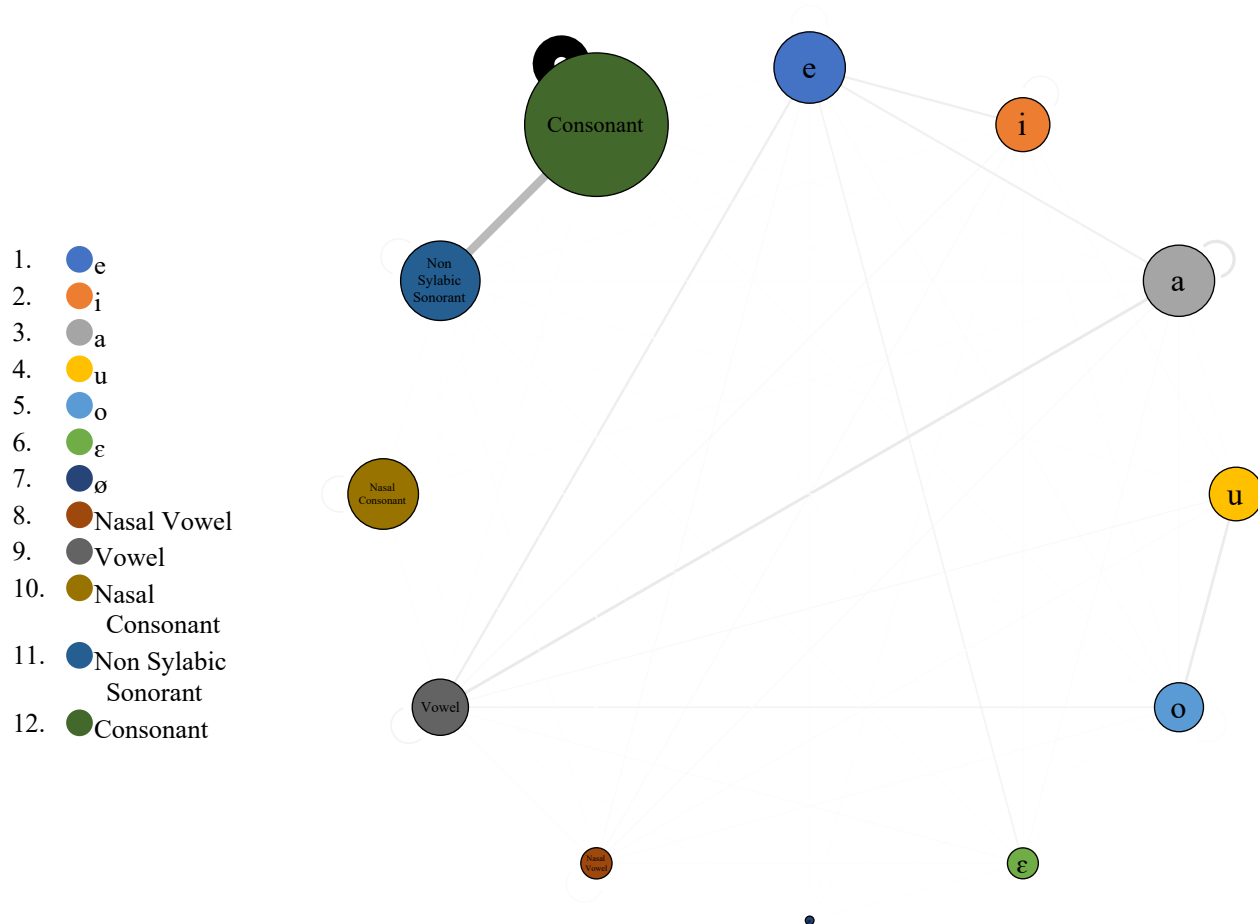
Partition Assignments

Model: “Linguistically Informed”

1) e	e: e
2) i	i i:
3) a	a a:
4) u	u u:
5) o	o o:
6) ε	ε
7) ø	ø ø:
8) Nasal Vowel	ã ĩ ẽ ĩ̃ ẽ̃ ã̃ õ̃ ẽ̃ ẽ̃ ẽ̃ ẽ̃ ẽ̃
9) Vowel	ɐ œ ɔ ə ɔ̃ a ɨ aj ej ɐj au ij ɔ y iw ɔj oj ɛa ej ɪ ʊ
10) Nasal Consonant	m n ɲ ŋ
11) Non Sylabic Sonorant	j l r w
12) Consonant	f b ɾ x ʎ ɸ d t ʈ s p z ʃ ʈ k ts θ h v c ɭ ɥ g dz ɟ ɡʷ β ɾ kʲ ʒ ʃ ʎ kʷ ɭ

Prior Transition Rates

For the 'Linguistically Informed' model, states were grouped into 12 sets. Here, the area of the circles is proportional to the occurrence frequencies for each group. The width of the lines is proportional to the rates of transition between each partition.



The majority rule consensus tree

