

Back to the roots

using regular sound correspondences
for linguistic phylogeny
(as one should)

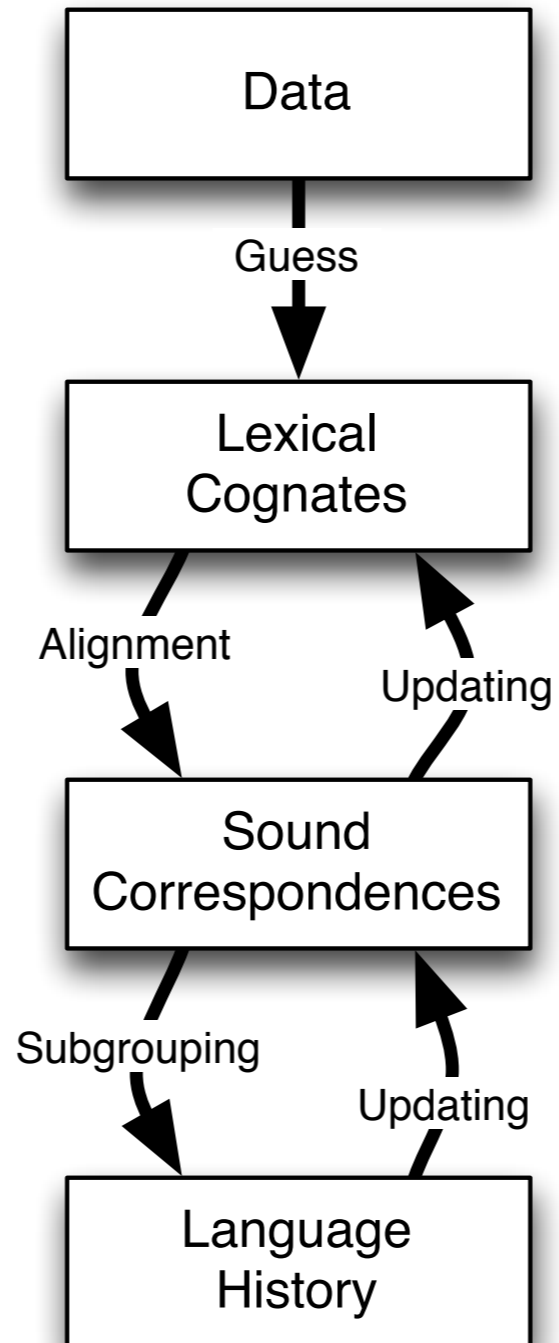
Michael Cysouw
Philipps Universität Marburg

Bern, 24 November 2012

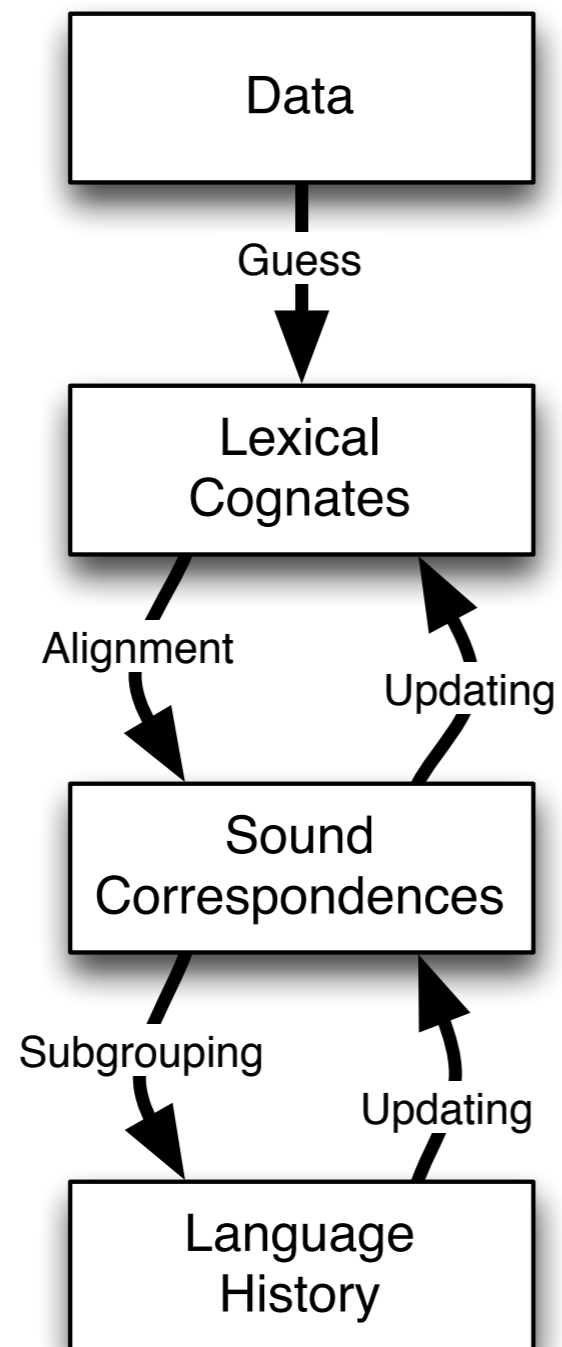
Methods for Lexical Comparison

- Comparative Method
- Swadesh Method
- Black Box Method
- Alignment Method
- Sound Change Method

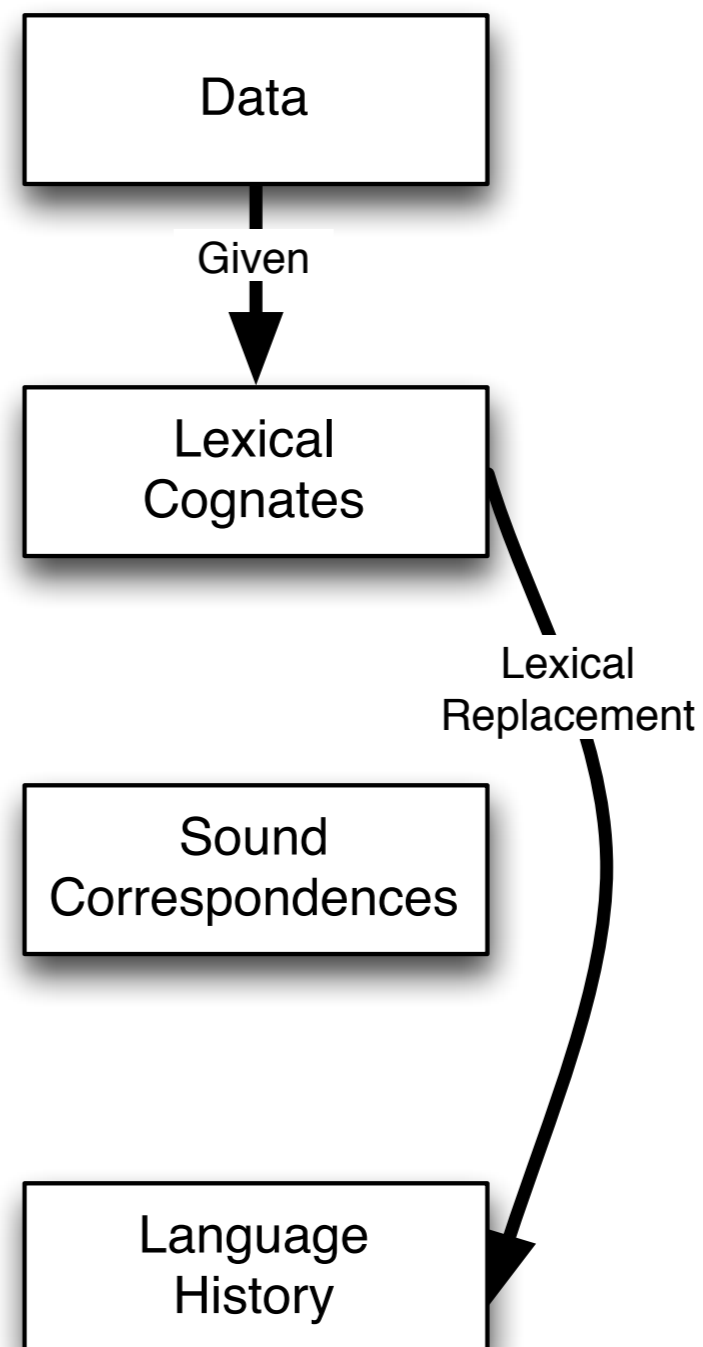
Comparative Method



Comparative Method



Swadesh Method



Swadesh Method

- **Necessary prerequisite**
known cognacy in parallel wordlist
- **Method**
Reconstruct history from
distribution of cognates

Swadesh Method

- **Language distances**
Swadesh (1952)
- **Stochastic model with first tree**
Sankoff (1969), Dobson (1969)
- **Modern replacement model**
Gray & Jordan (2000), (Gleason 1959)

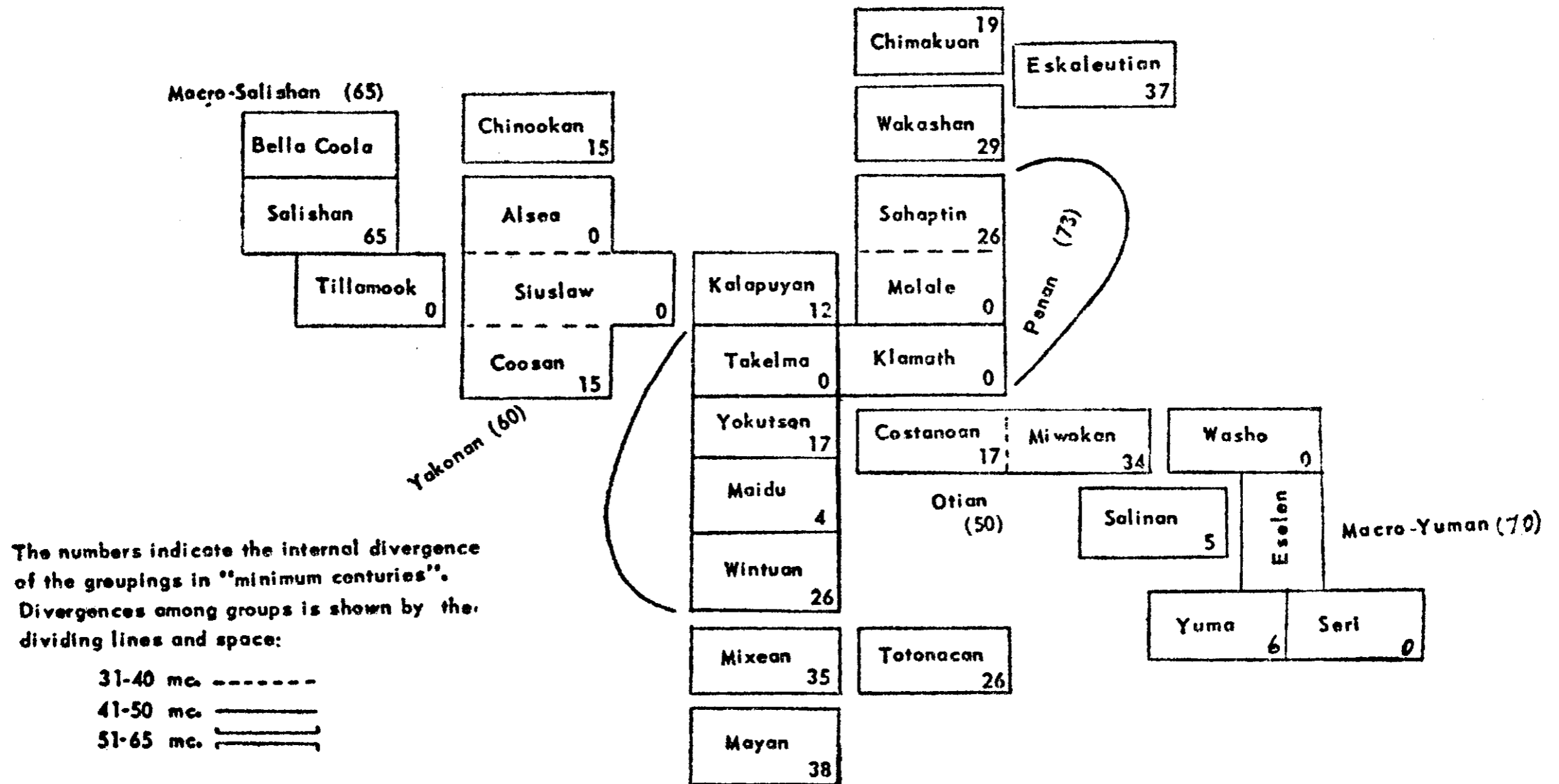


FIG. 2. Internal and external relationships of the Penan micro-phylum.

Swadesh, Morris. 1952. Lexico-Statistic Dating of Prehistoric Ethnic Contacts: With Special Reference to North American Indians and Eskimos. Proceedings of the American Philosophical Society 96(4). 452-463.

Swadesh, Morris. 1959, "Linguistics As An Instrument of Prehistory. Southwestern Journal of Anthropology 15, no. 1: 20-35.



David Sankoff

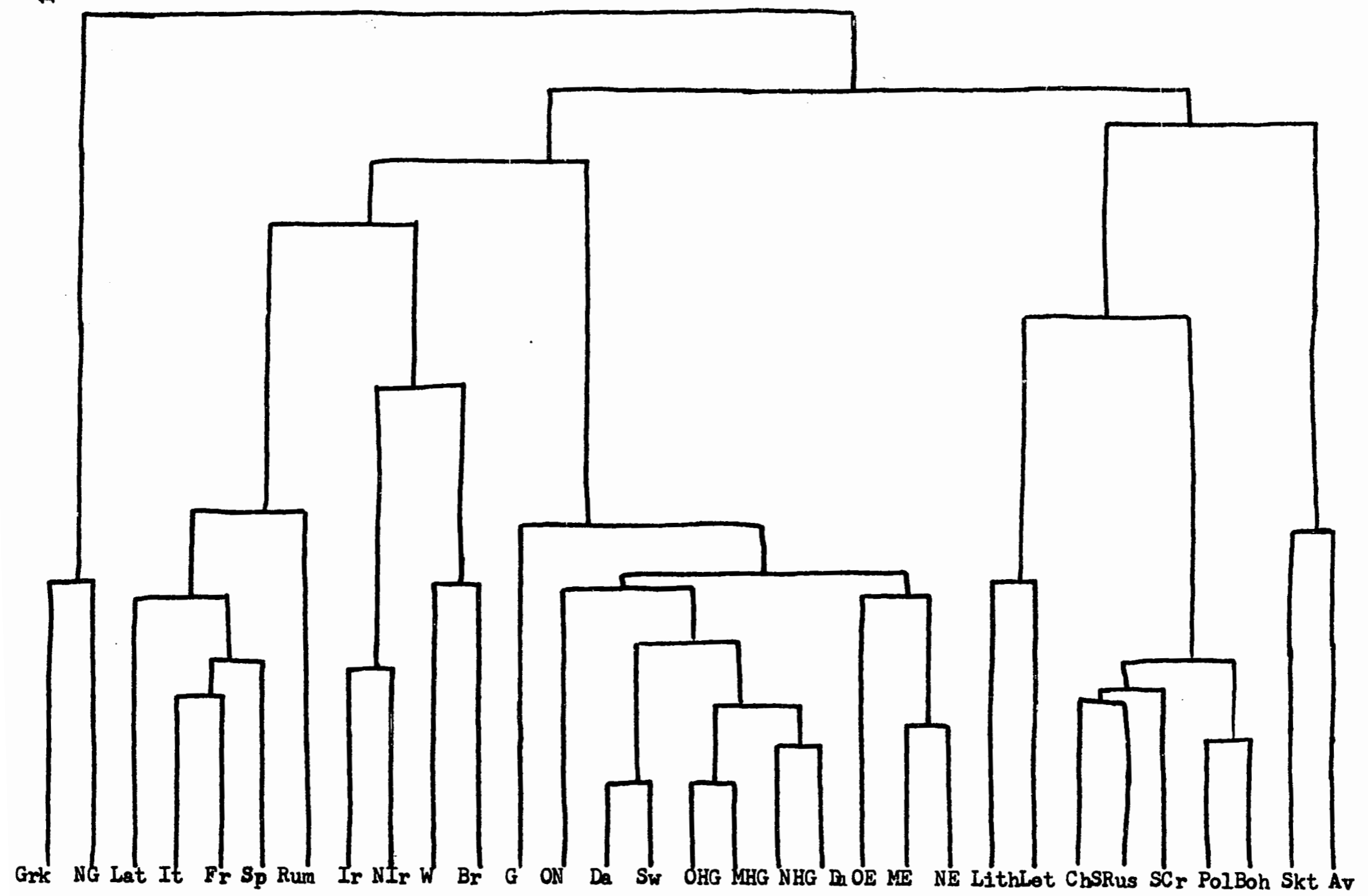


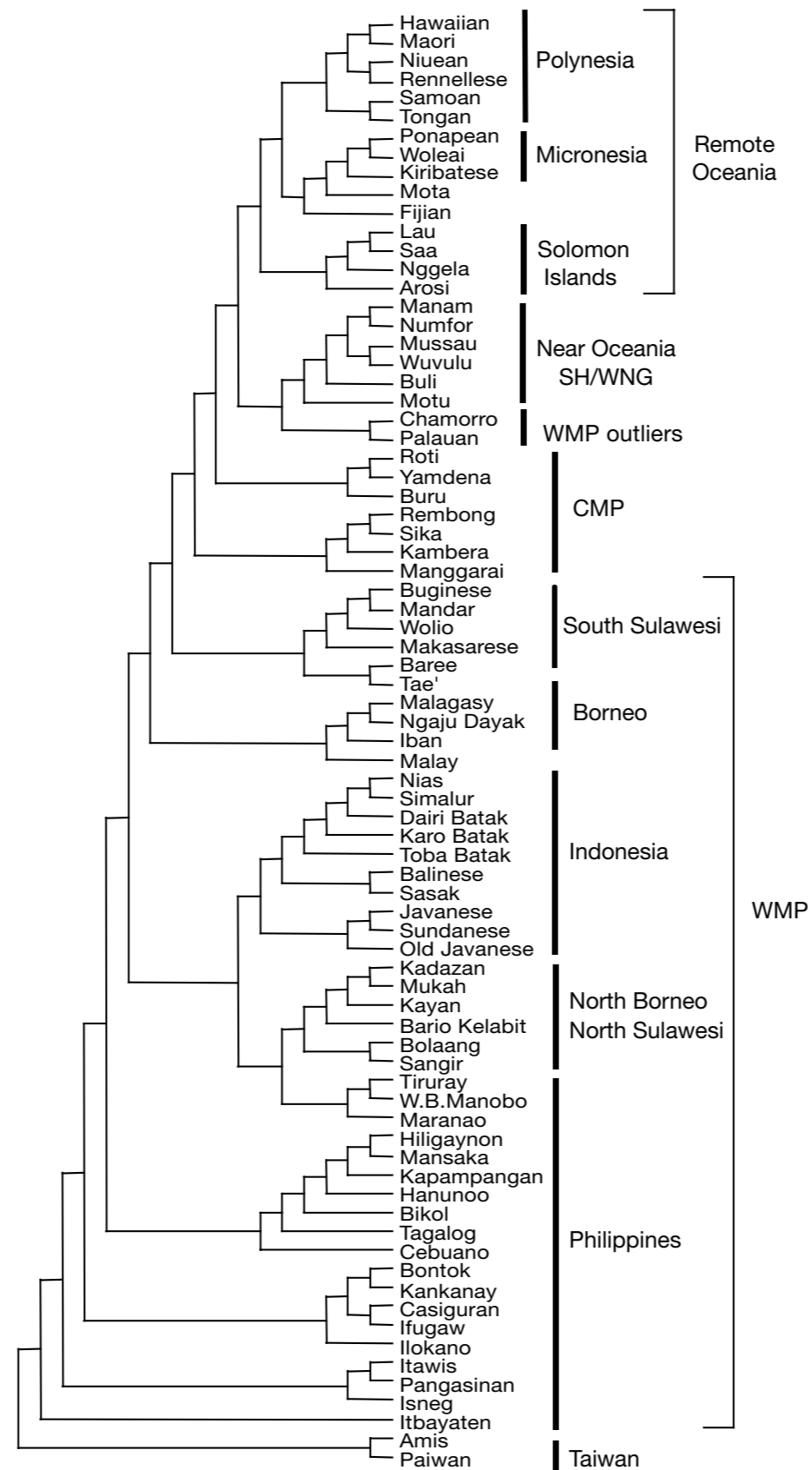
Annette J. Dobson

Dobson, A. J. "Lexicostatistical Grouping." *Anthropological Linguistics* 11, no. 7 (1969): 216-221.
Dobson, Annette J. "Unrooted Trees for Numerical Taxonomy." *Journal of Applied Probability* 11, no. 1 (1974).

Fig. 16. Tree of Indo-European languages as reconstructed from data in Table 4.

114

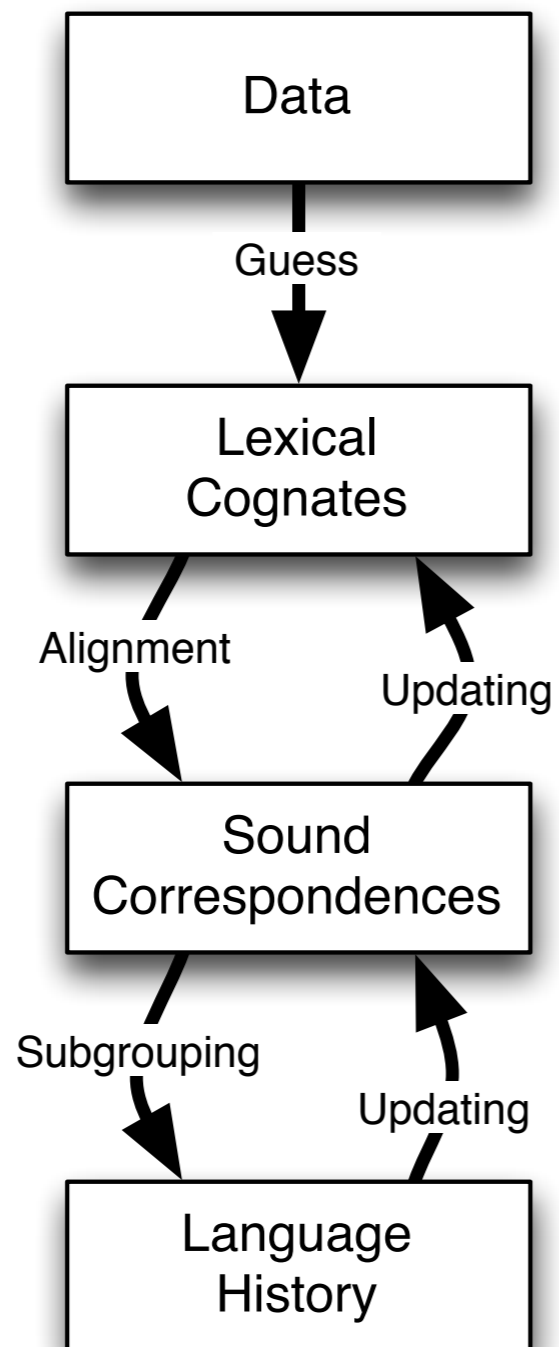




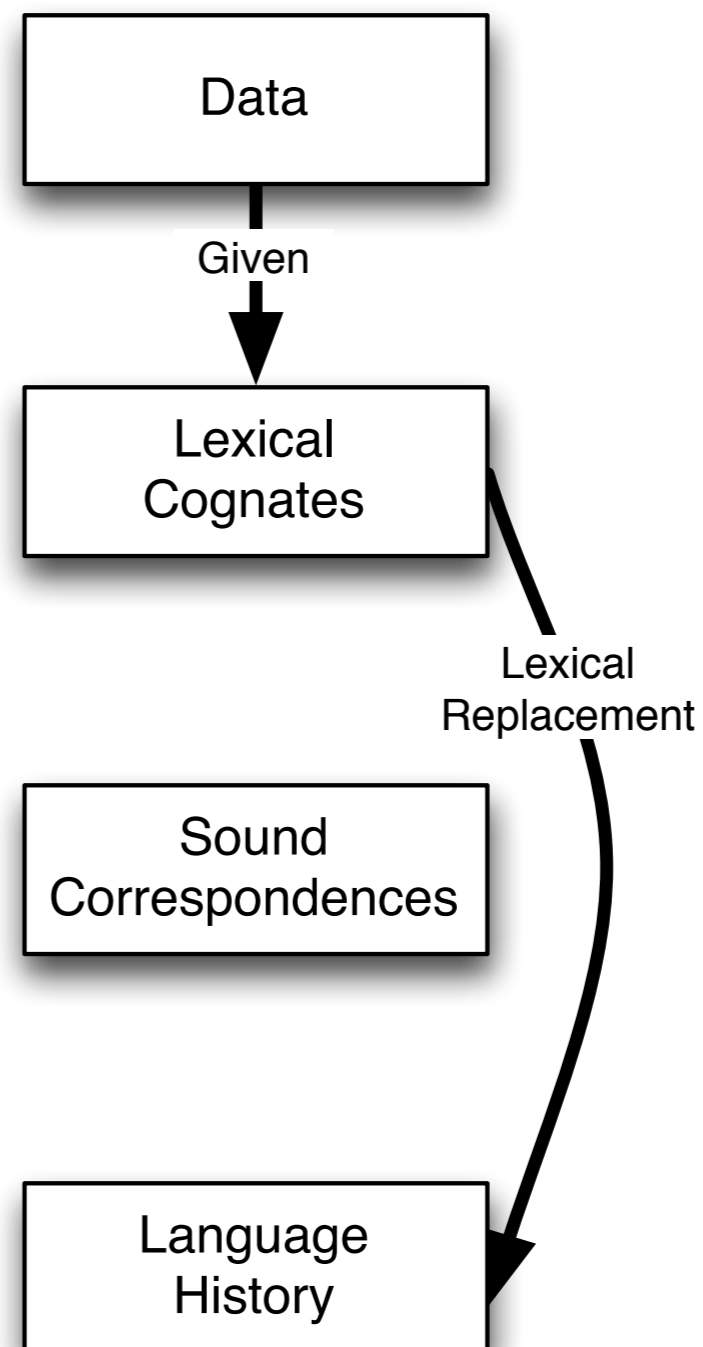
Gray, Russell D, and Fiona M Jordan. "Language Trees Support the Express-train Sequence of Austronesian Expansion." *Nature* 405 (2000): 1052-1055.

Gleason, H A. "Counting and Calculating for Historical Reconstruction." *Anthropological Linguistics* 1, no. 2 (1959): 22-32.

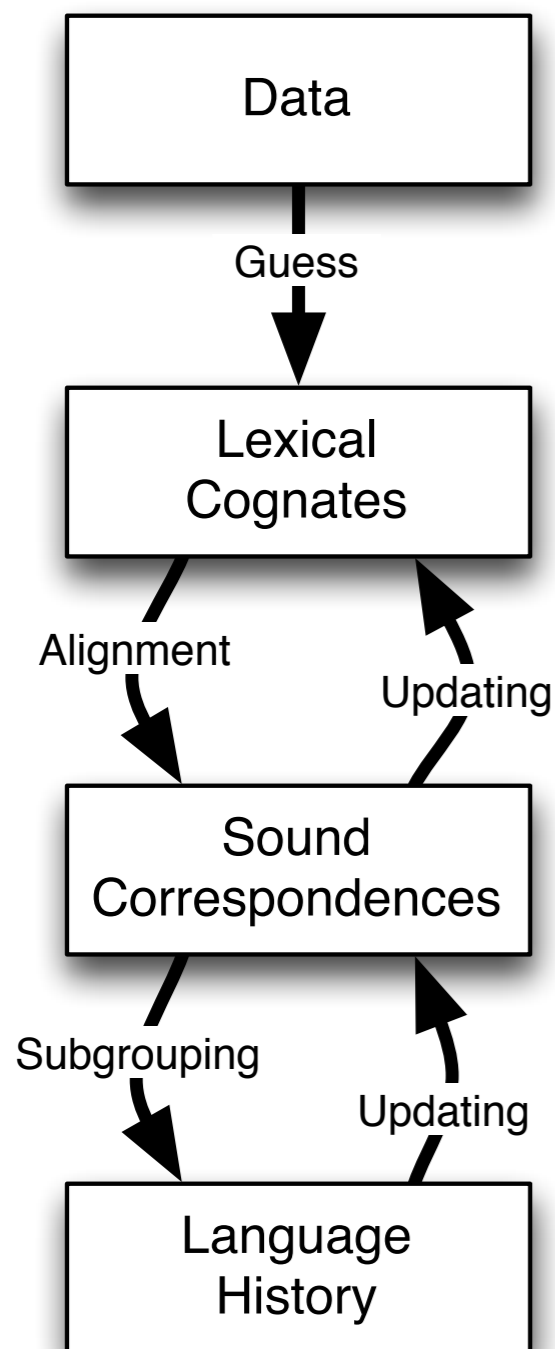
Comparative Method



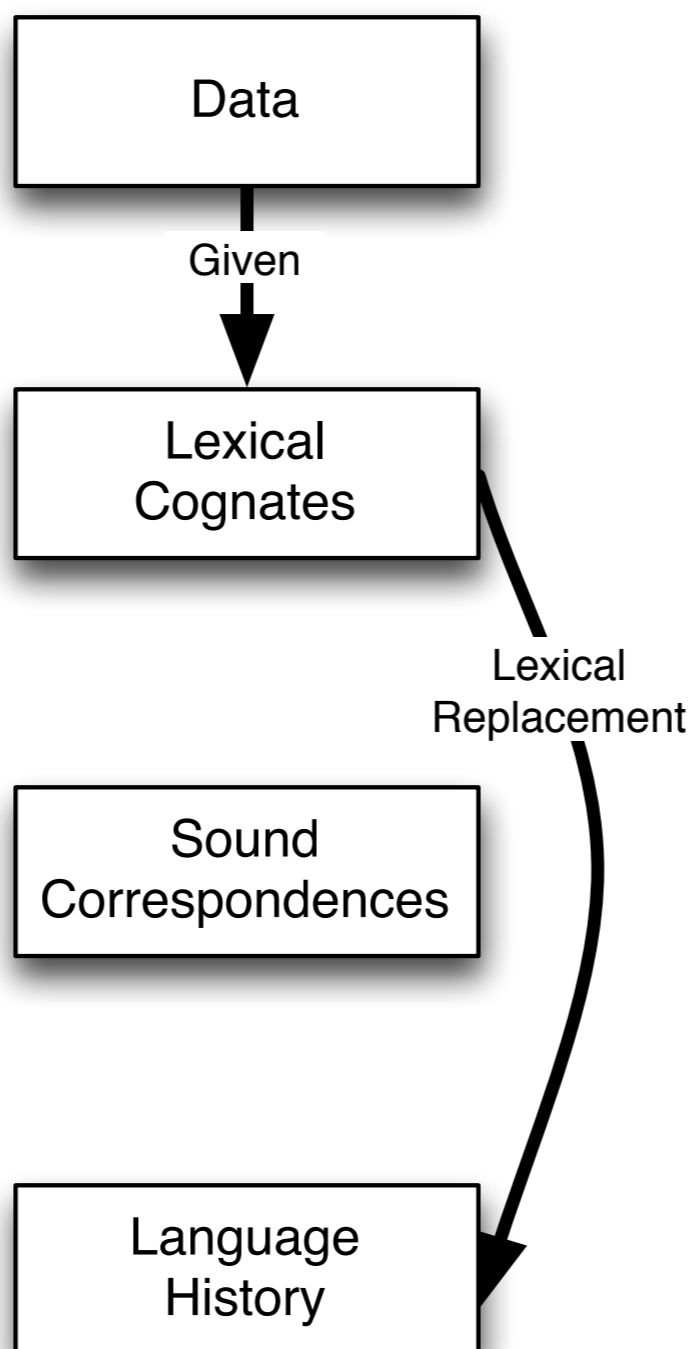
Swadesh Method



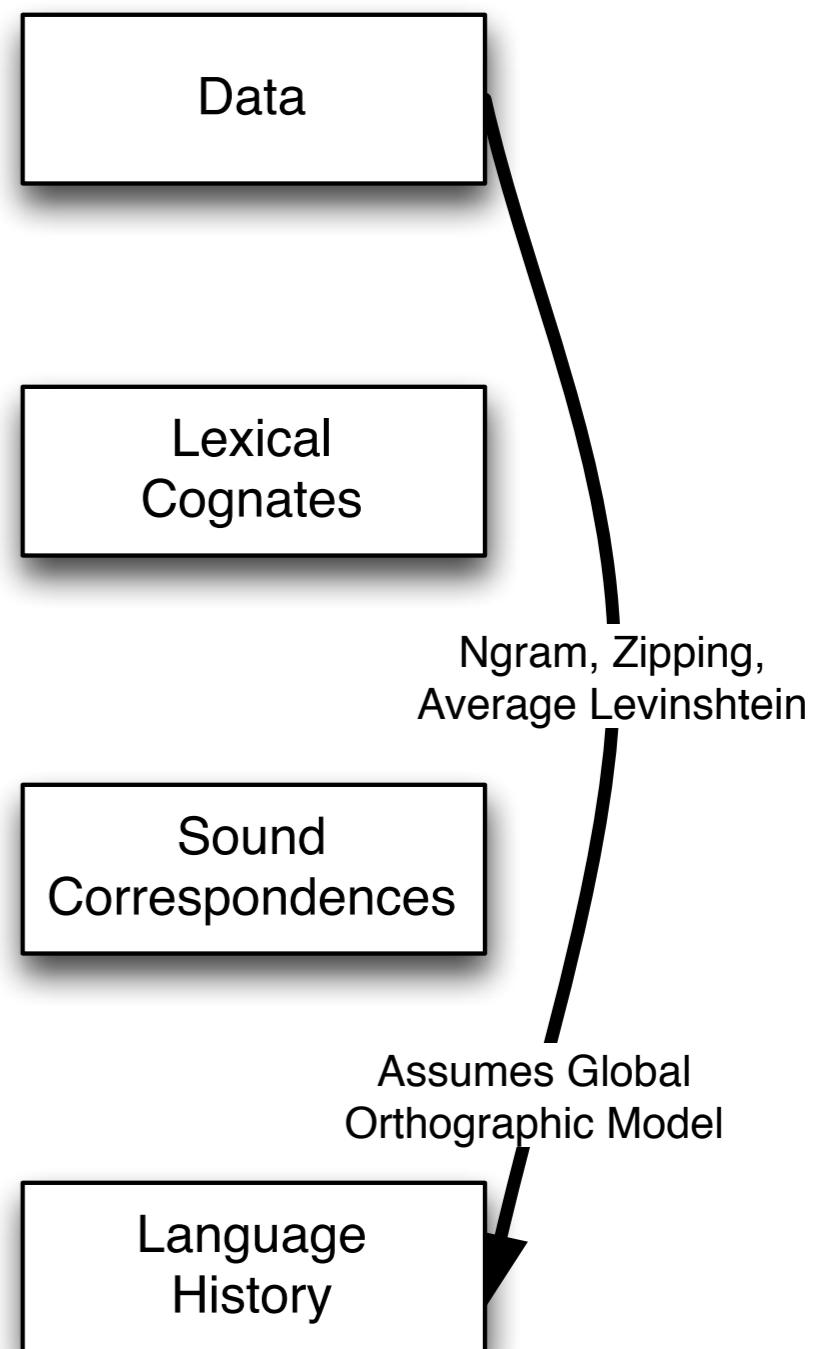
Comparative Method



Swadesh Method



Black Box Method



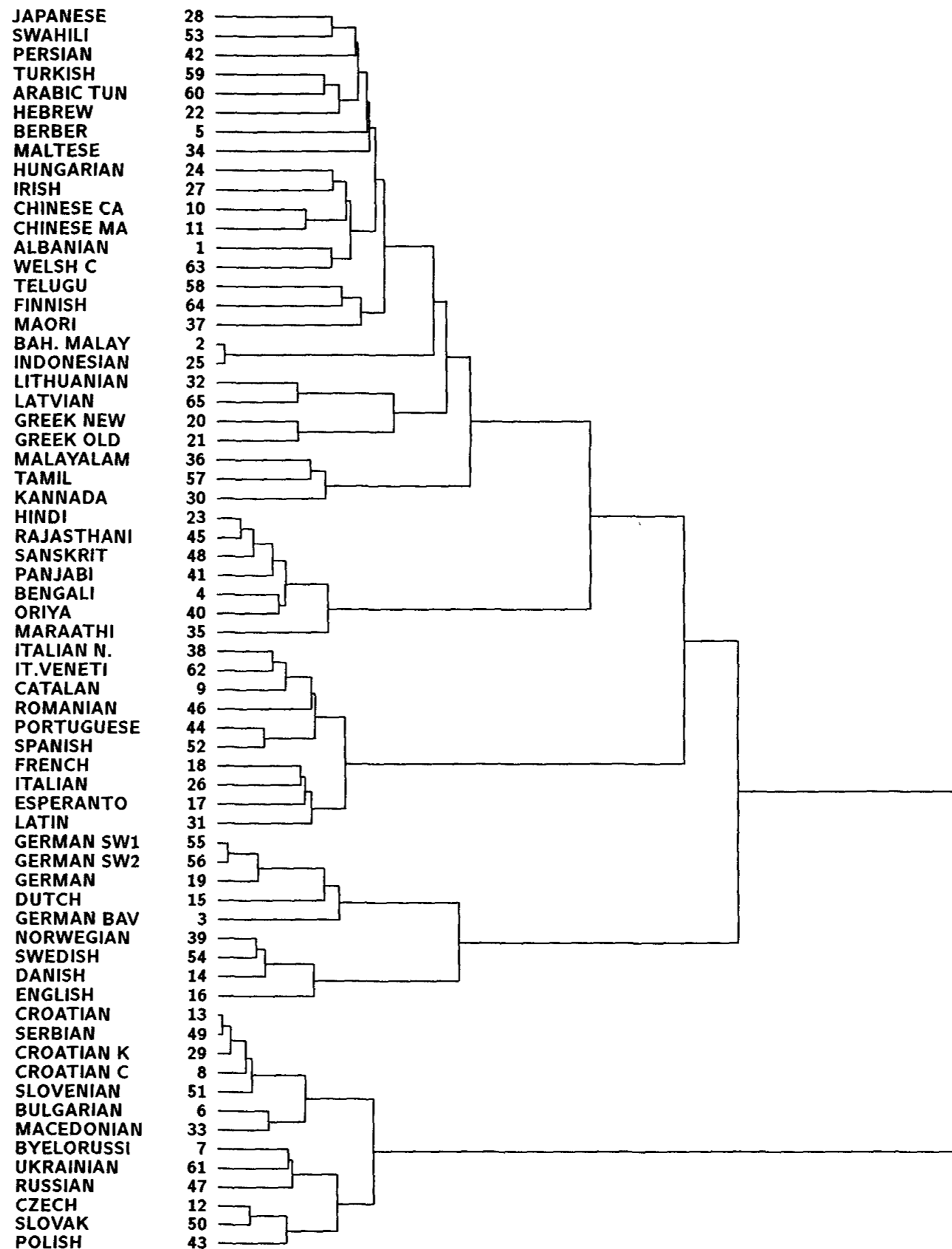
Black Box Methods

- **Necessary prerequisite**
Global orthographic model
("same alphabet for all data")
- **Method**
Compare global orthographic similarity

Black Box Methods

- **Aggregate Levenshtein Distance**
Batagelj et al. (1992), Kessler (1995)
- **N-Gram Similarity**
Huffman (1998)
- **Zippering Distance**
Benedetto et al. (2002)

Aggregate Levenshtein Distance

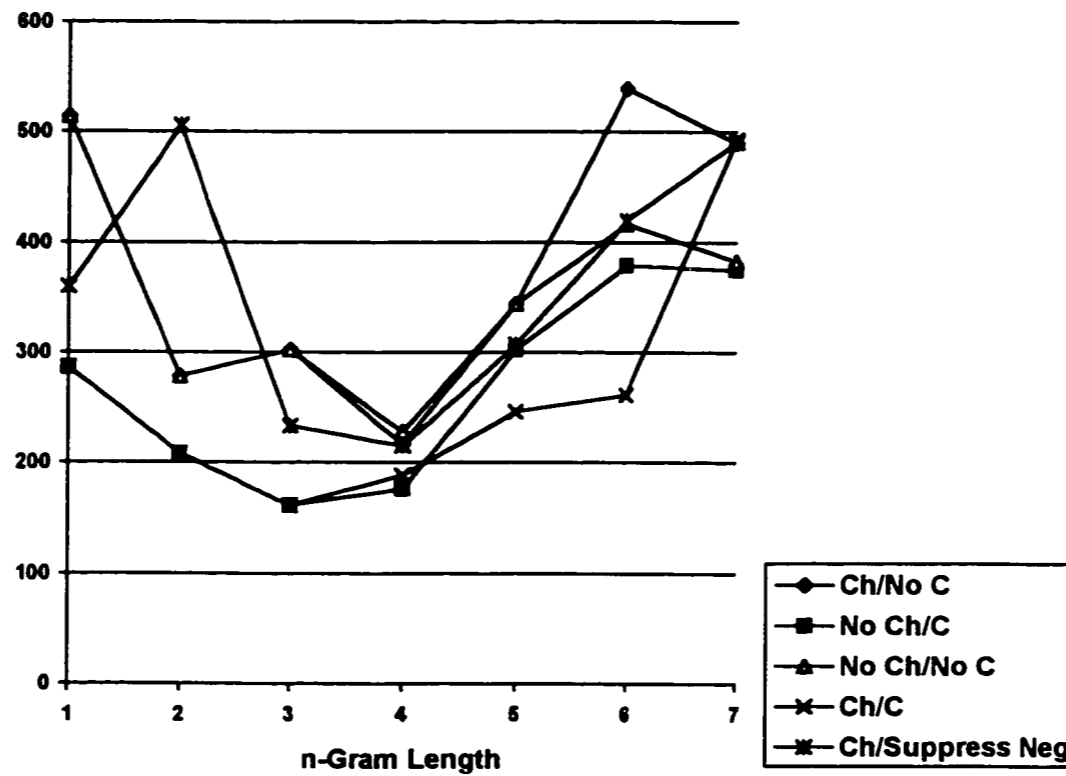


Appendix 9. Divergence from reference classification

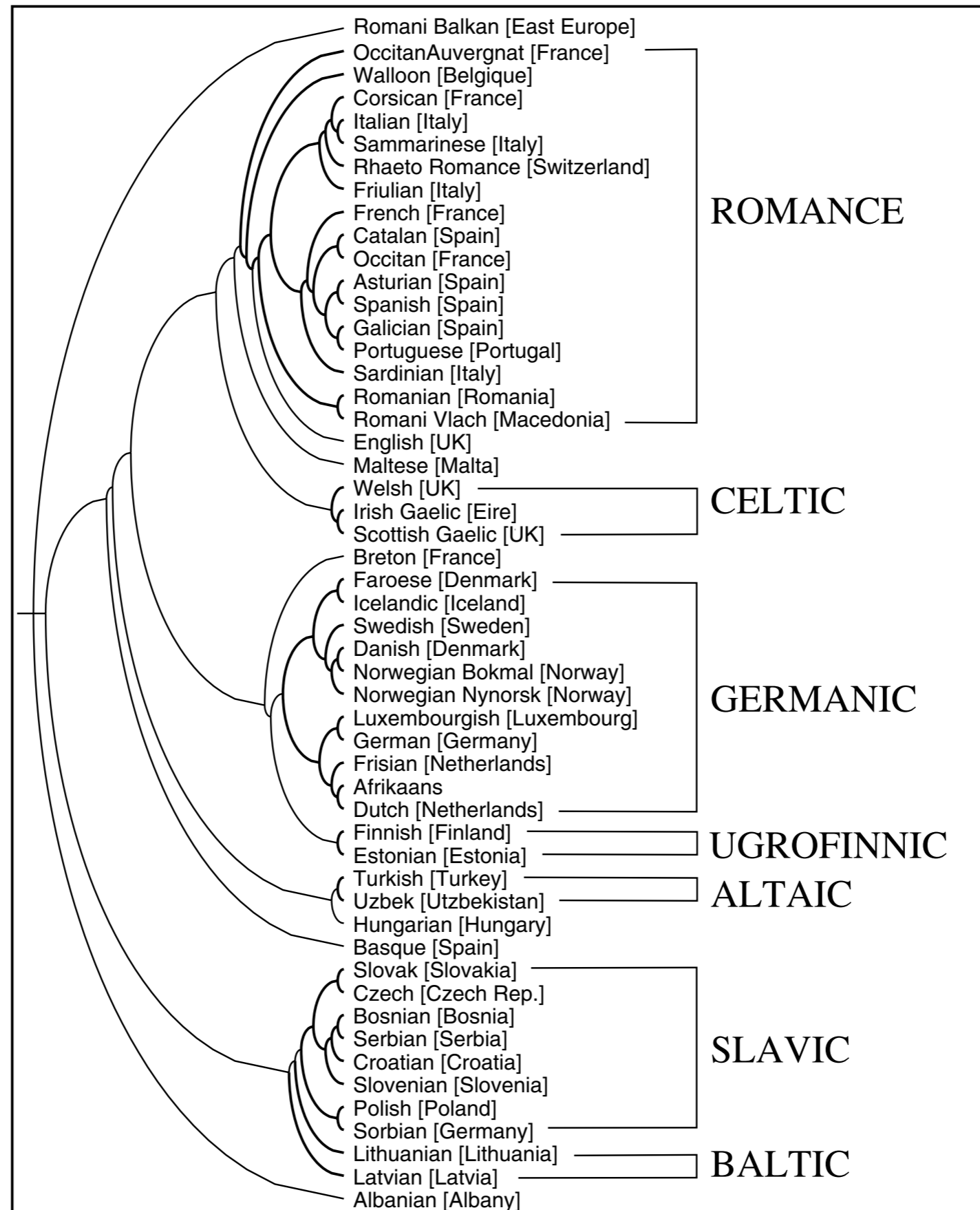
Text display:

N-gram Length	<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>	<u>6</u>	<u>7</u>
Chained, No centroid	514	278	302	217	344	539	490
Not Chained, Centroid	286	208	161	176	302	379	375
Not Chained, No Centroid	514	278	302	228	344	417	383
Chained, Centroid	286	208	161	188	246	261	493
Suppress Negative	360	506	233	215	307	420	491

Divergence from Reference Classification vs. n-Gram Length



Zippering

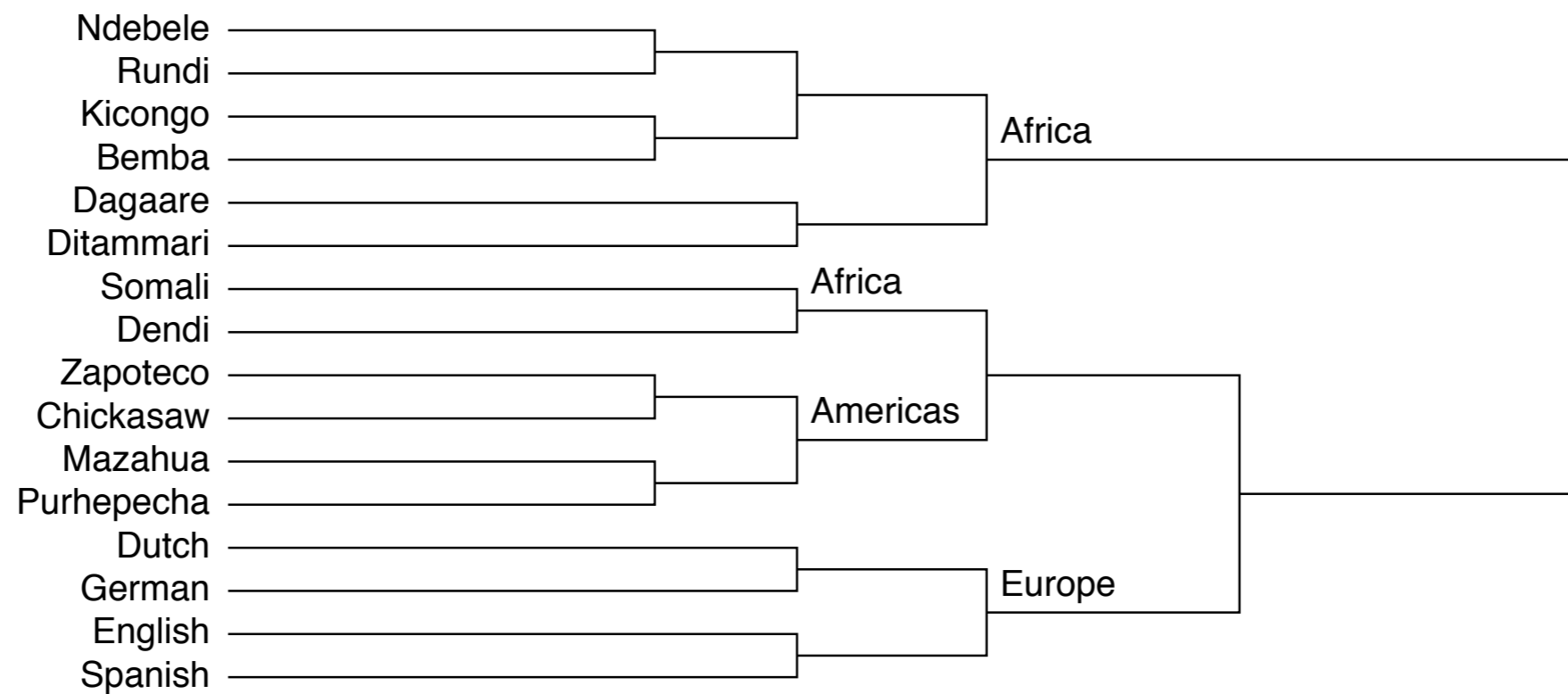


Zippering

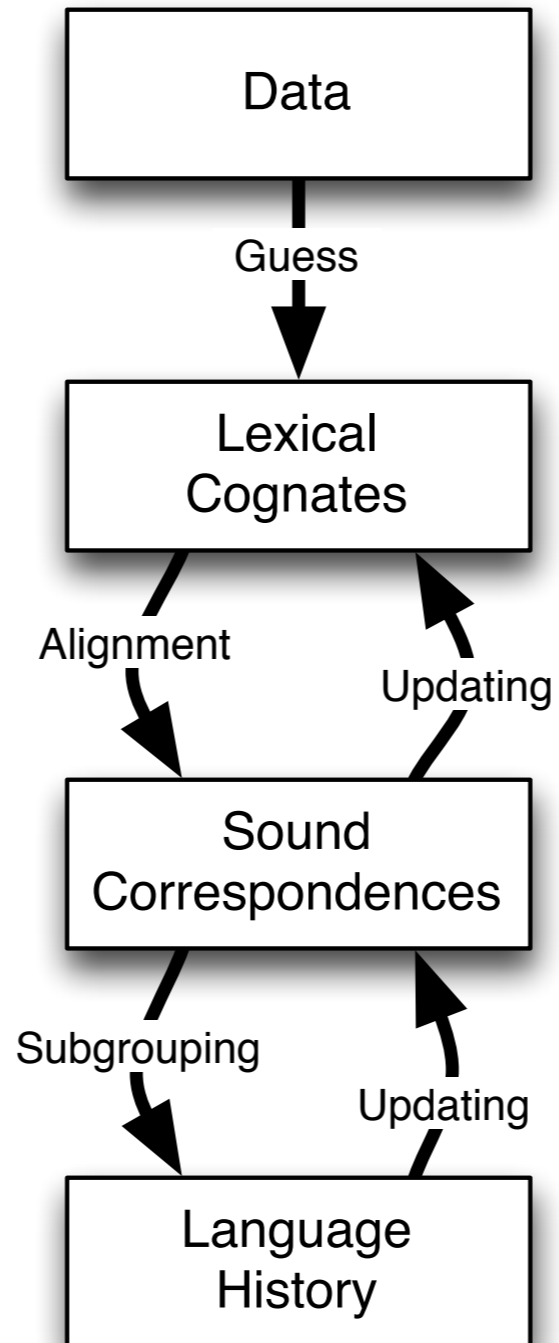
Normalized Compression Distance

(C = compressed size)

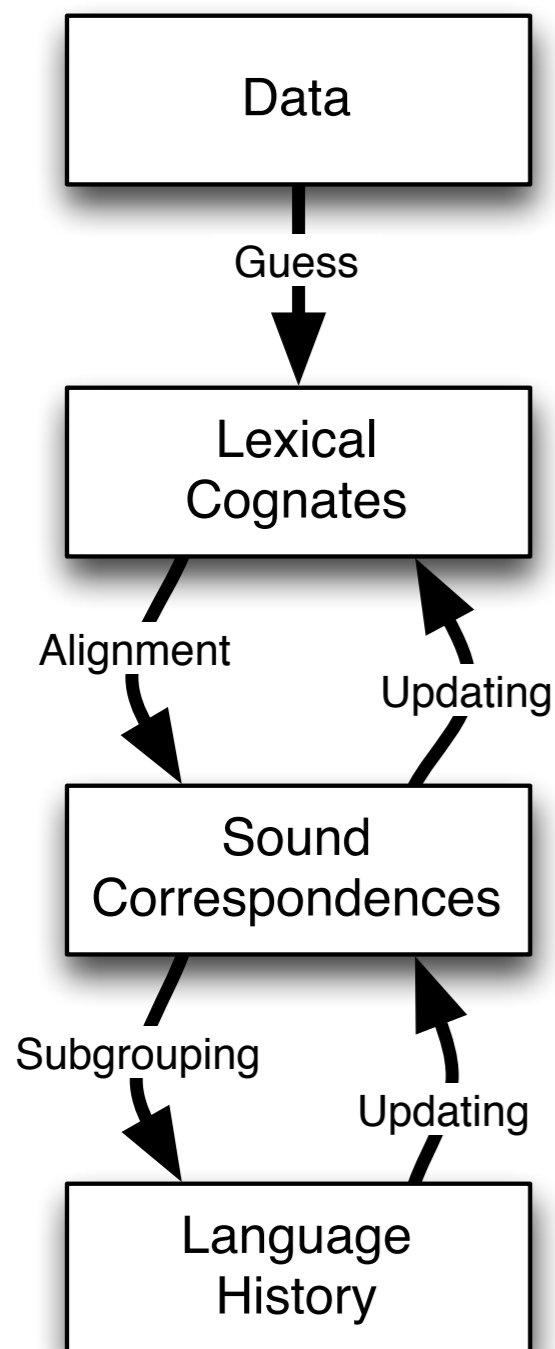
$$NCD(x, y) = \frac{C(xy) - \min\{C(x), C(y)\}}{\max\{C(x), C(y)\}}.$$



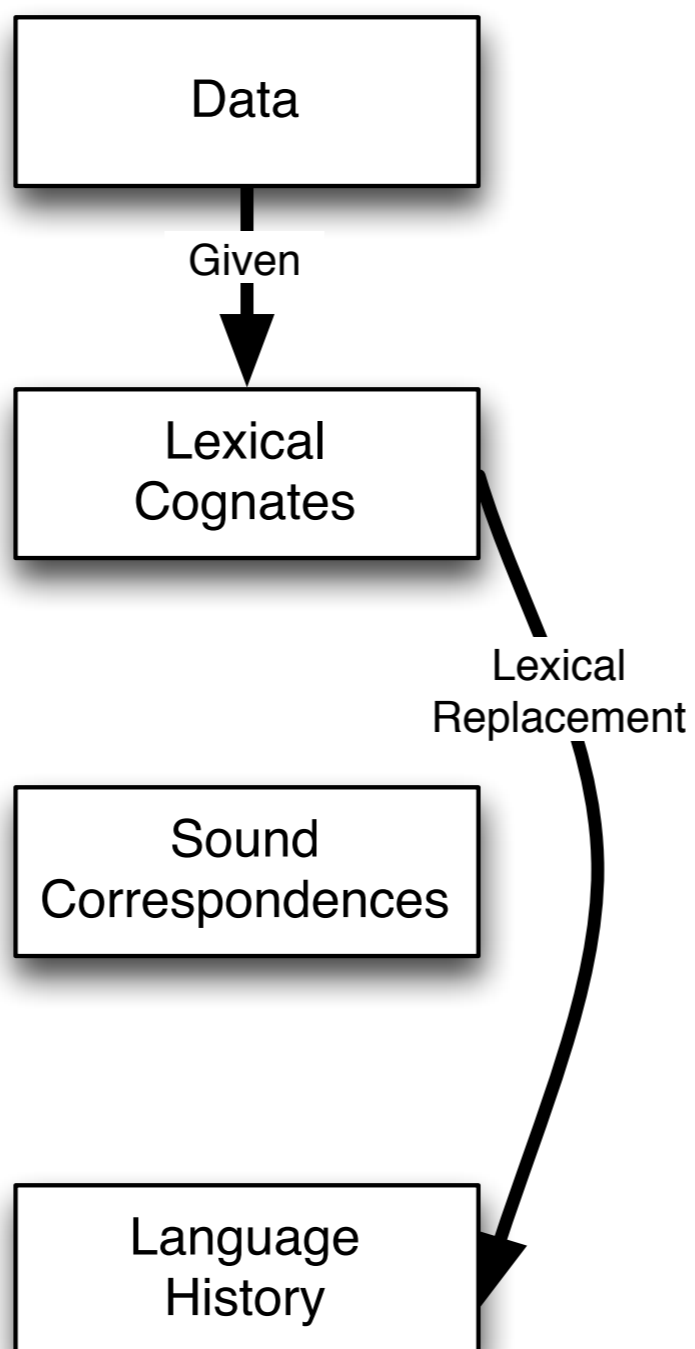
Comparative Method



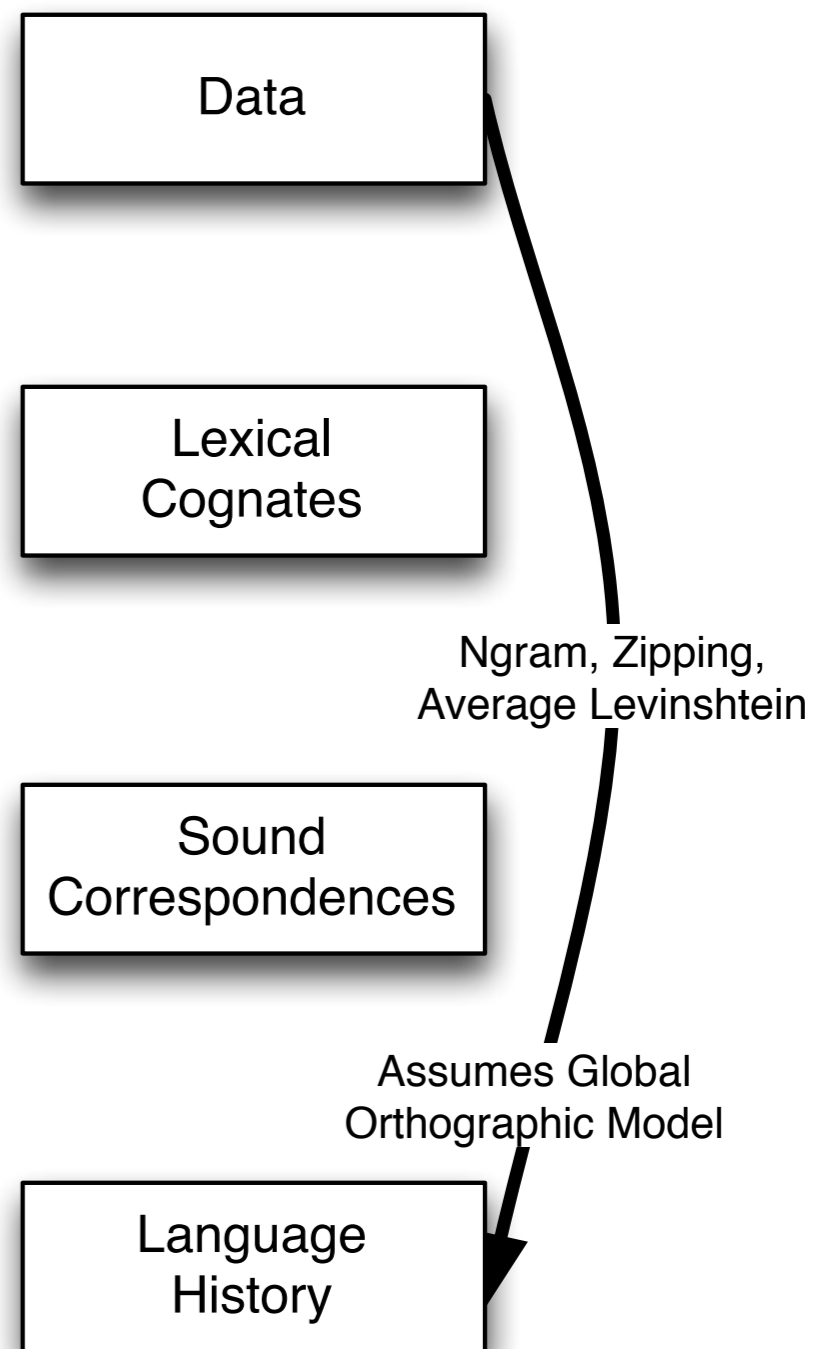
Comparative Method



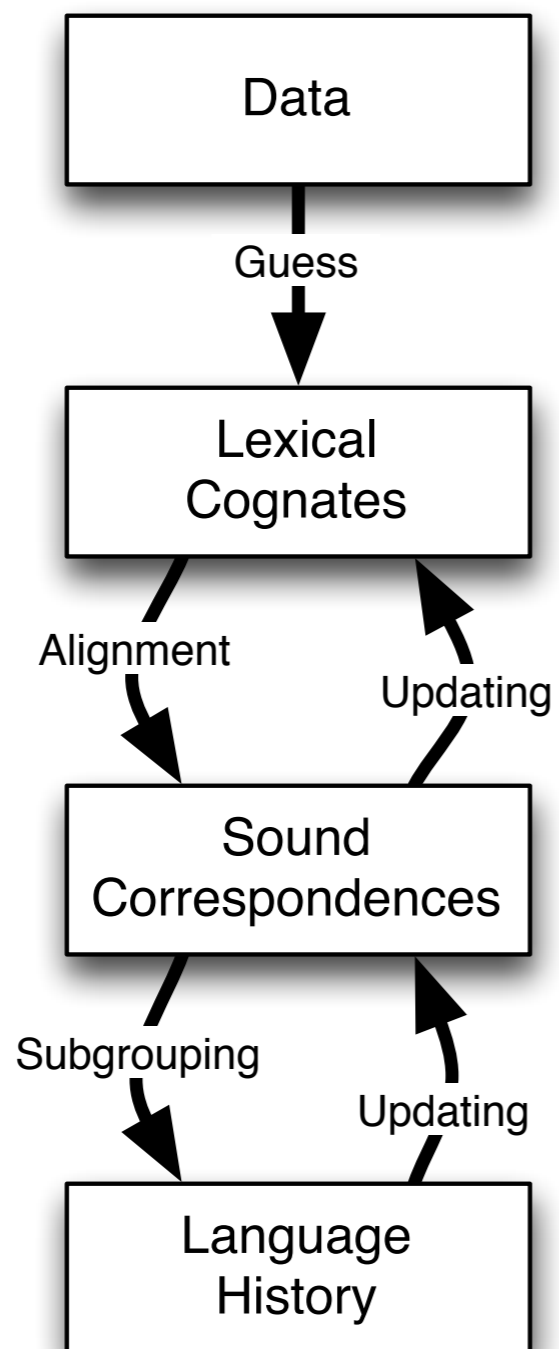
Swadesh Method



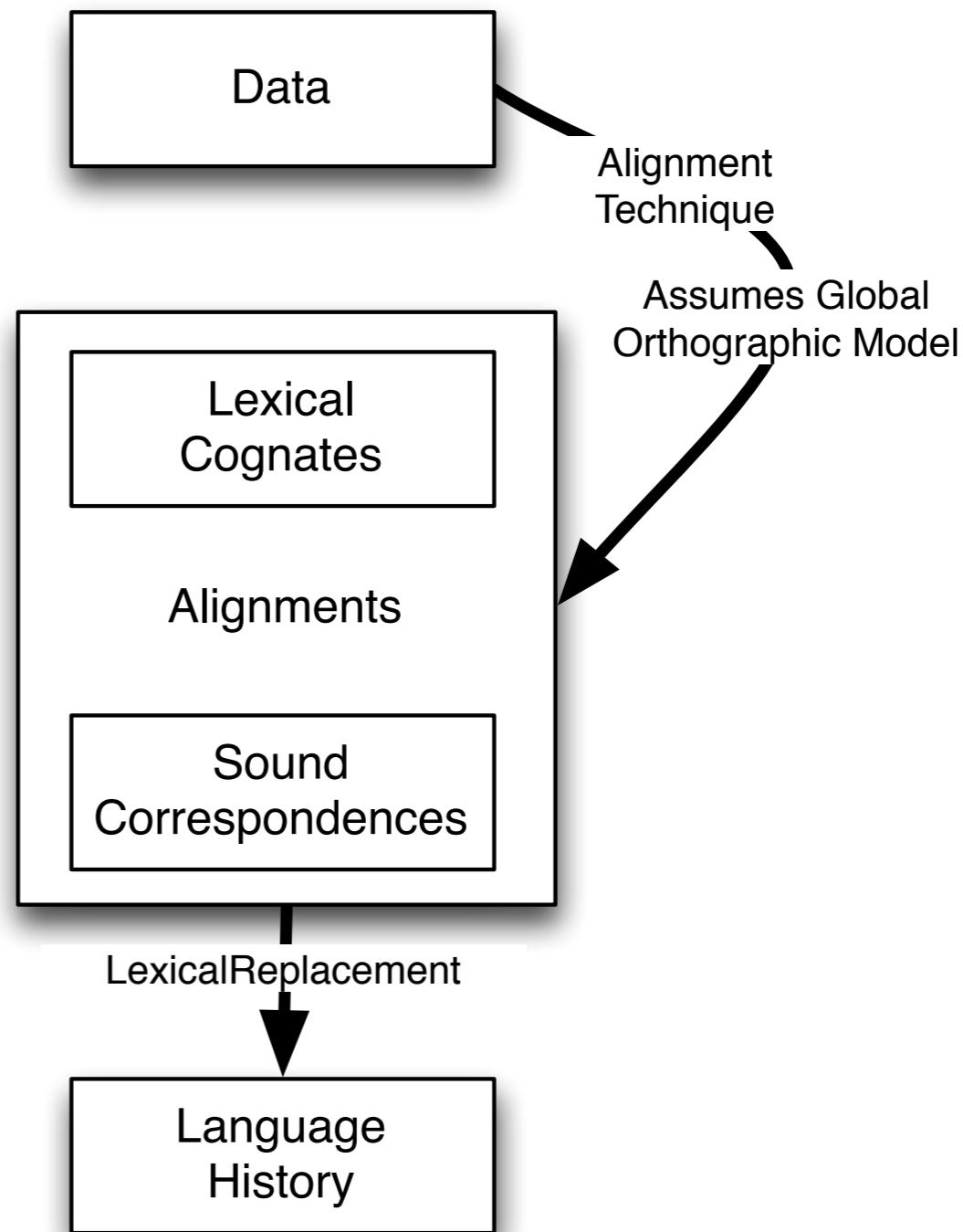
Black Box Method



Comparative Method



Alignment Method



Alignment Method

- **Necessary prerequisite**
Global orthography model and parallel wordlist
- **Method**
Discover sound correspondences, possibly together with decision on cognacy

Alignment Method

- **Alignment**
Covington (1996, 1998)
- **Relation to Levenshtein**
Kondrak (2000)
- **Multiple Alignment**
Bhargava & Kondrak (2009), Prokić et al. (2009)

Spanish/Italian/French 'five':

θ i ŋ k - o

c i ŋ k w e

s ã - k - -

Covington's alignments

ALINE's alignments

three : trēs

θ r i y
t r ē s

|| θ r iy ||
|| t r ē || s

blow : flāre

b l - - o w
f l ā r e -

|| b l o || w
|| f l ā || re

full : plēnus

f - - - u l
p l ē n u s

|| f u l ||
|| p - l || ēnus

fish : piscis

f - - - i š
p i s k i s

|| f i š ||
|| p i s || kis

I : ego

- - a y
e g o -

|| ay ||
|| e || go

tooth : dentis

- - - t u w θ
d e n t i - s

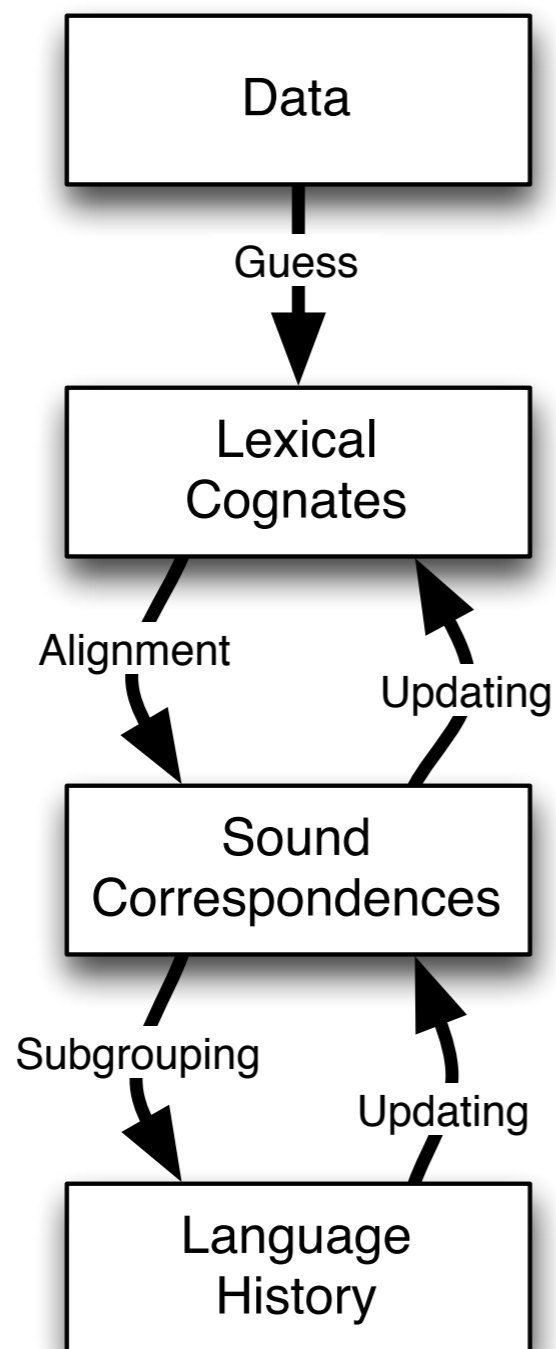
|| t uw θ ||
den || t i s ||

D--E--N-
Z--E--N-
DZIE--N-
DI--E--NA
D--I--A-
D--I--E-
Z--U--E-
J--O--UR
DJ--O--U-

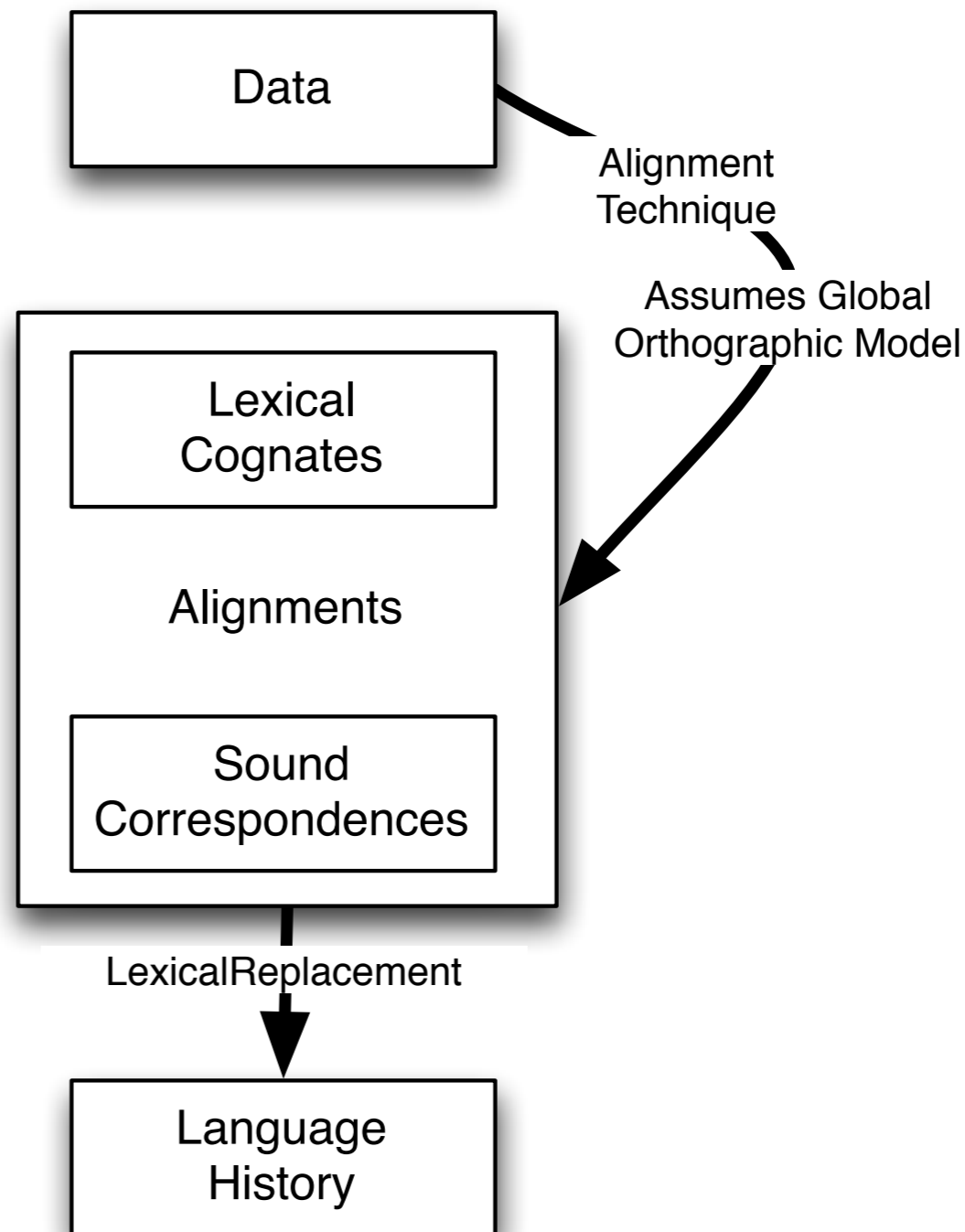
j	'a	-	-	-	-
j	'a	z	e	-	-
-	'a	s	-	-	-
j	'a	s	-	-	-
j	'a	z	e	k	a
j	'ε	-	-	-	-
-	'b	s	-	-	-

Prokić, Jelena, Martijn Wieling, and John Nerbonne. "Multiple Sequence Alignments in Linguistics." In Proceedings of the EACL 2009 Workshop on Language Technology and Resources for Cultural Heritage, Social Sciences, Humanities, and Education (LaTeCH-SHELT&R 2009). Association for Computational Linguistics, 2009.

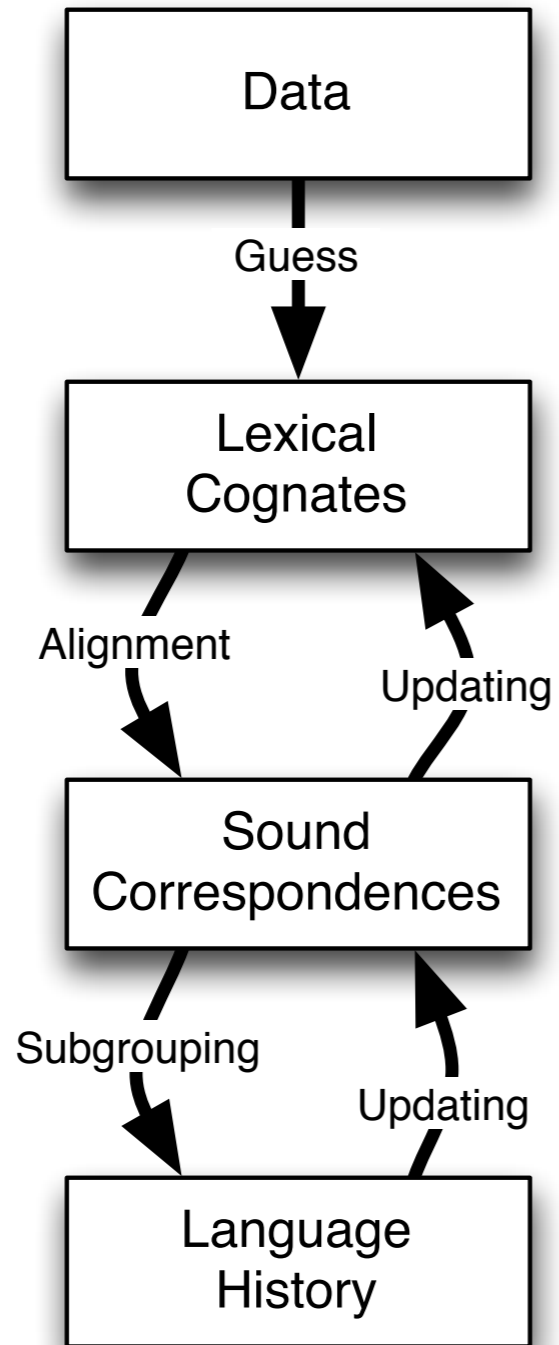
Comparative Method



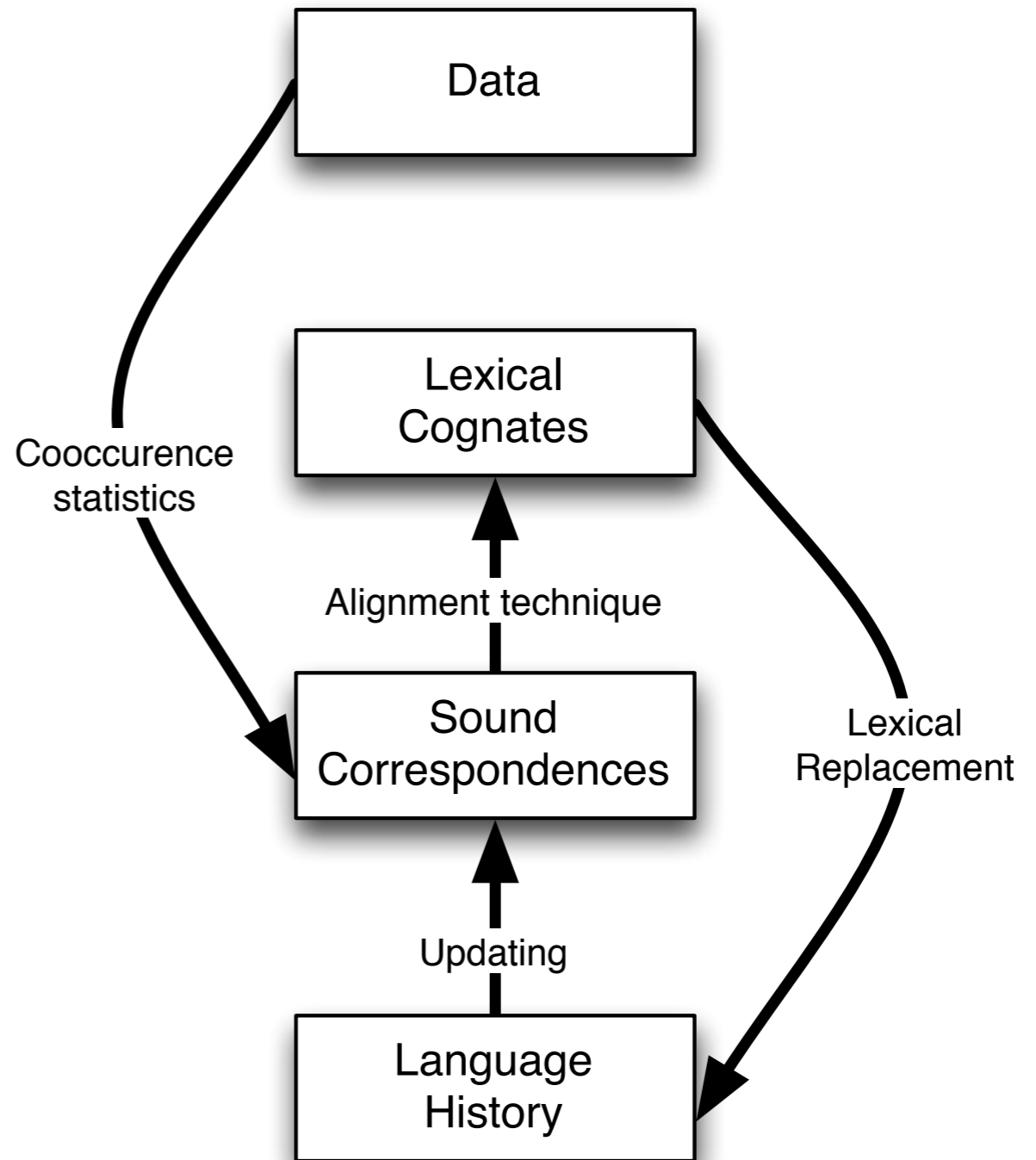
Alignment Method



Comparative Method



Sound Change Method

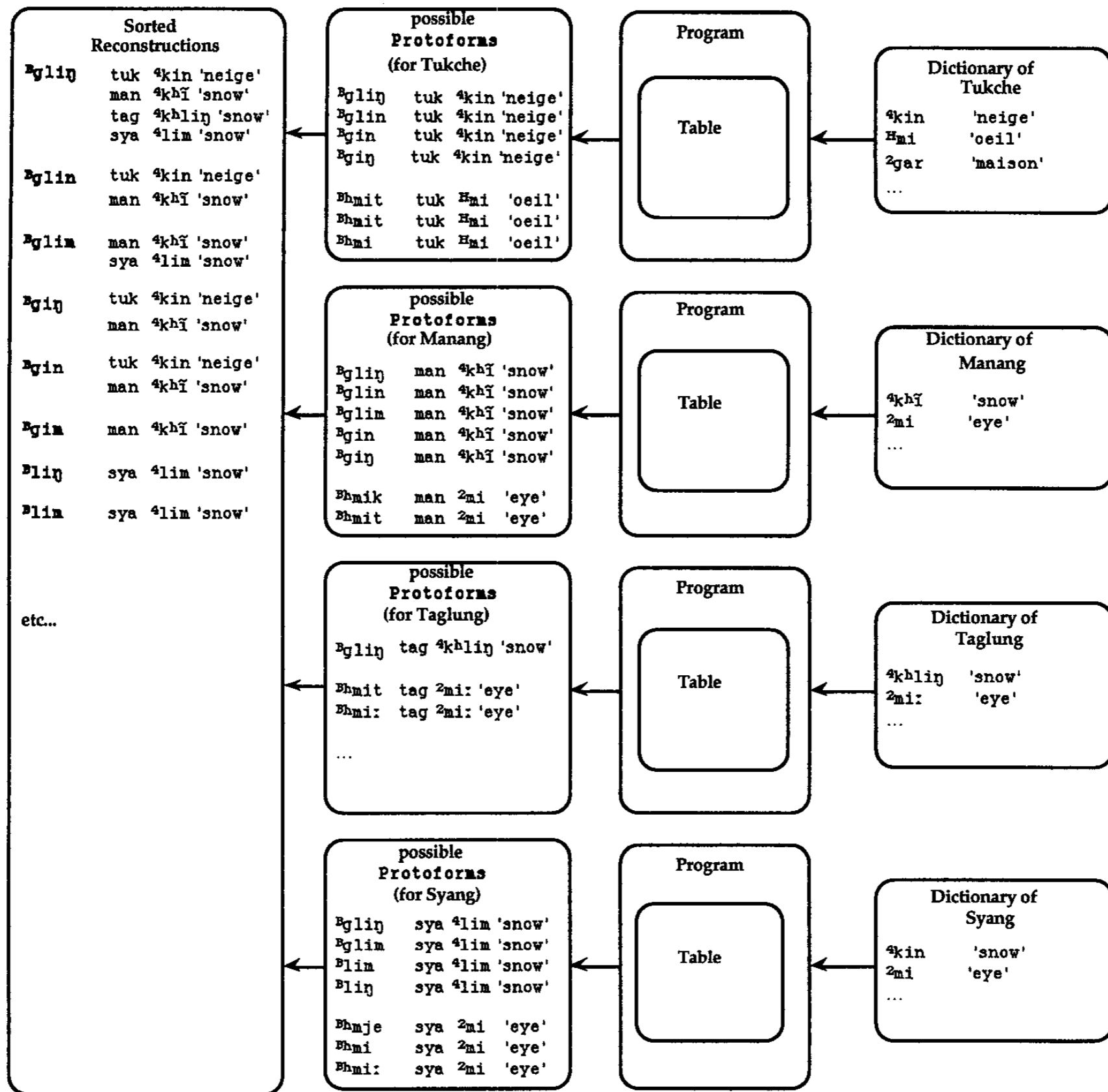


Sound Change Method

- **Necessary prerequisite**
(Large) parallel wordlist,
but **no** global harmonized orthography
- **Method**
Discover sound correspondences through
regularity, and use that to discover cognacy

Sound Change Method

- **Rule-based automatic reconstruction**
(Hewson 1973, Lowe & Mazaudon 1994)
- **Sound co-occurrence**
Ross (1947)
- **Cross-script mapping**
Cysouw & Jung (2007)



Lowe, John Brandon, and Martine Mazaudon. "The Reconstruction Engine: A Computer Implementation of the Comparative Method." *Computational Linguistics* 20, no. 3 (1994): 381-417.

Now let me put the matter more rigidly. Suppose I take 1000 common ideas, numerals, parts of the body, names of relatives, etc. and express them in English and German. I now prepare a Table of the following kind:

English ↓

← German

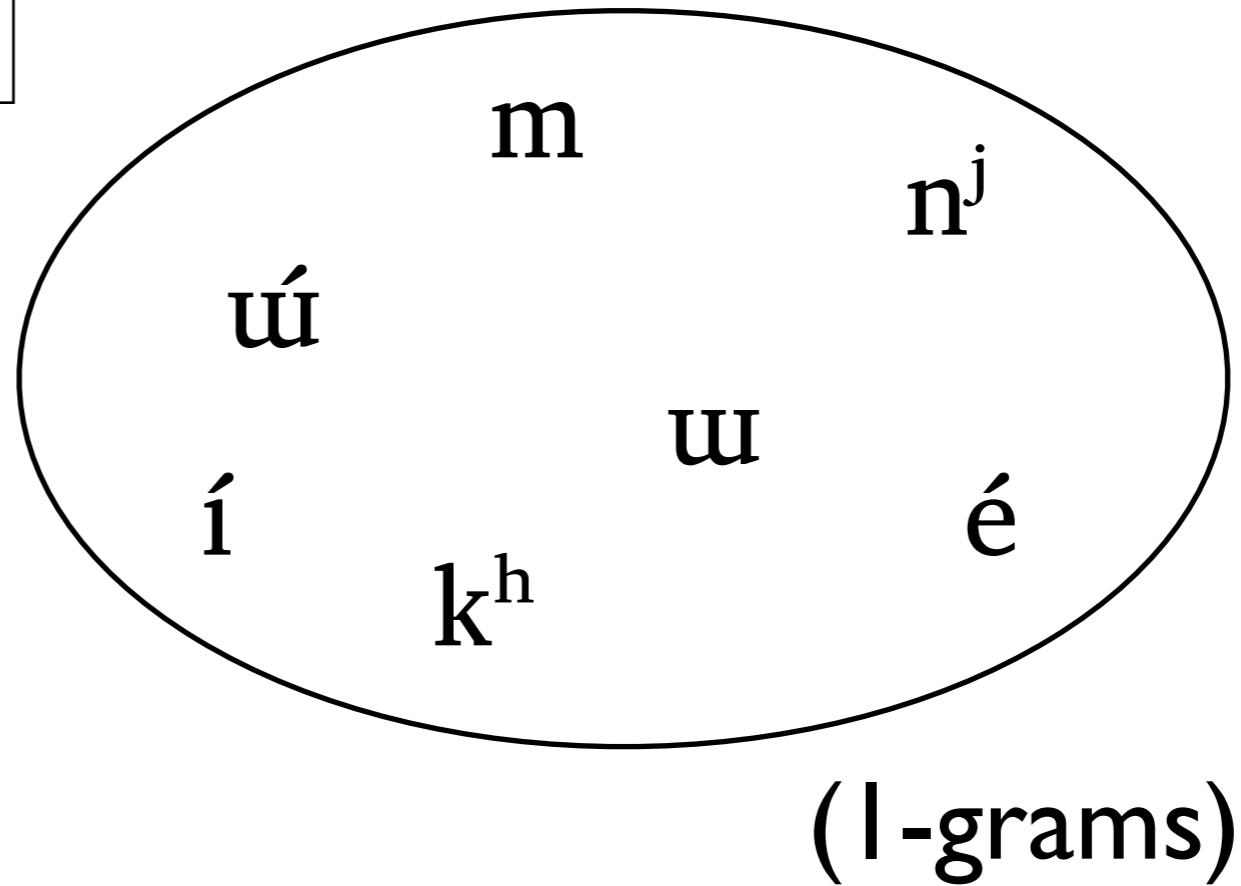
	[pf]	[ts]	[t]	[f]	[h]...
[p]	x x x x x		x			
[f]		x		x x x ¹ x x x		
[t]	x	x x x ² x x x x			x	
[d]			x x x ³ x		x ⁴	
.						
.						
.						
.						
.						
.						
.						

Cross-script mapping

E	R	freq	dice
r	р	184	0.88874745
n	н	115	0.8461936
l	л	104	0.79646295
s	с	114	0.7927922
t	т	165	0.7701921
m	м	47	0.7699933
o	о	184	0.7510106
k	ть	21	0.74458015
p	п	50	0.7388723
i	и	102	0.7034591
a	а	221	0.6866478
u	у	40	0.6449104
c	к	77	0.6251676
e	е	219	0.59066784
b	б	32	0.525643
w	в	46	0.46787763
d	д	42	0.381996
⋮	⋮	⋮	⋮

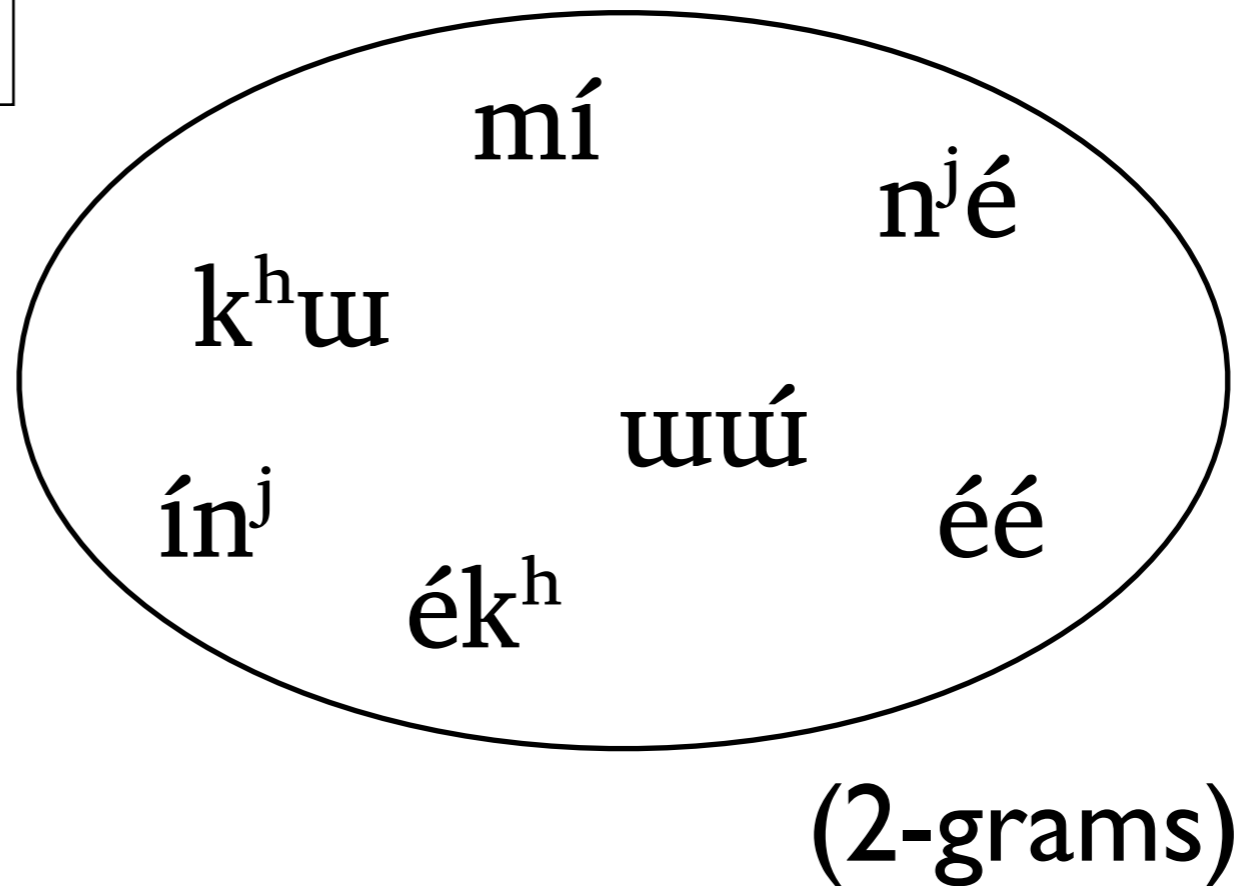
‘bag of symbol’ approach

mín^jéék^huú



‘bag of symbol’ approach

mín^jéék^hwú



Cross-script mapping

Cross-script mapping

- Ignore linear structure of words
“bag of symbols” approach

Cross-script mapping

- Ignore linear structure of words
“bag of symbols” approach
- Use parallel wordlist to estimate
co-occurrences of n-grams

Cross-script mapping

- Ignore linear structure of words
“bag of symbols” approach
- Use parallel wordlist to estimate co-occurrences of n-grams
- N-grams that have a high probability of co-occurrence in parallel meaning are interesting for historical linguistics

	Bora	Muinane
down	tʃin ^j e, paári	báari, gíino
bee	íimúʔóexp ^h i, téʔts ^h ipa	nîbiri, mîbiriʔi
sharp	ts ^h úʔxiβáne	sîxéβano
...

Bora	Muinane	Bora	Muinane	Bora	Muinane
#k	#k ^h	#i	#i	#n	#n
ki	k ^h u	#a	#a	#m	#m
se	ts ^h i	di	ti	mi	mu
xe	xi	du	to	ni	nu
ga	k ^w a	#d	#t	us	ts ^h i
ba	pa	#s	#ts ^h	#t	#t ^h
#b	#p	gi	tʃi	ig	uk ^w
e#	i#	ni	ni	#ϕ	#p ^h

Using bigram matching for cognate detection

Using bigram matching for cognate detection

- Bora 'two': mɪn^jéék^hwú
Muinane 'two': míínokɪ

Using bigram matching for cognate detection

- Bora 'two': mɪn^jéék^hwú
 Muinane 'two': mɪínokɪ
- Extension of Levenshtein (1966):
 Needleman-Wunsch (1970)

	#m	mi	ii	in	no	ok	kɨ	ɨ#
#m								
mi								
in ^j								
n ^j e								
ee								
ek ^h								
k ^h ʊ								
ʊʊ								
ʊ#								

Levenshtein

	#m	mi	ii	in	no	ok	k†	†#
#m	1							
mi		1						
in ^j			?					
n ^j e				?				
ee					?			
ek ^h								
k ^h ω								
ωω								
ω#								

Levenshtein, V I. "Binary Codes Capable of Correcting Deletions, Insertions, and Reversals." Soviet Physics Doklady 10 (1966): 707–710.

Needleman-Wunsch

	#m	mi	ii	in	no	ok	k†	†#
#m	22	3	2	2	2	2	2	2
mi	4	12	2	2	5	1	1	1
in ^j	2	1	5	9	3	1	1	2
n ^j e	1	1	5	5	4	1	1	2
ee	3	3	3	3	6	2	2	2
ek ^h	1	2	1	1	4	2	3	2
k ^h ω	2	2	2	2	2	1	23	2
ωω	2	2	3	3	2	2	4	4
ω#	2	2	3	2	3	1	3	4

Needleman, S B, and C D Wunsch. "A General Method Applicable to the Search for Similarities in the Amino Acid Sequence of Two Proteins." *Journal of Molecular Biology* 48, no. 3 (1970): 443-453.

	#m	mi	ii	in	no	ok	k†	†#
#m	22	3	2	2	2	2	2	2
mi	4	12	2	2	5	1	1	1
in ^j	2	1	5	9	3	1	1	2
n ^j e	1	1	5	5	4	1	1	2
ee	3	3	3	3	6	2	2	2
ek ^h	1	2	1	1	4	2	3	2
k ^h ω	2	2	2	2	2	1	23	2
ωω	2	2	3	3	2	2	4	4
ω#	2	2	3	2	3	1	3	4

	#m	mi	ii	in	no	ok	k†	†#
#m	22	3	2	2	2	2	2	2
mi	4	12	2	2	5	1	1	1
in ^j	2	1	5	9	3	1	1	2
n ^j e	1	1	5	5	4	1	1	2
ee	3	3	3	3	6	2	2	2
ek ^h	1	2	1	1	4	2	3	2
k ^h ω	2	2	2	2	2	1	23	2
ωω	2	2	3	3	2	2	4	4
ω#	2	2	3	2	3	1	3	4

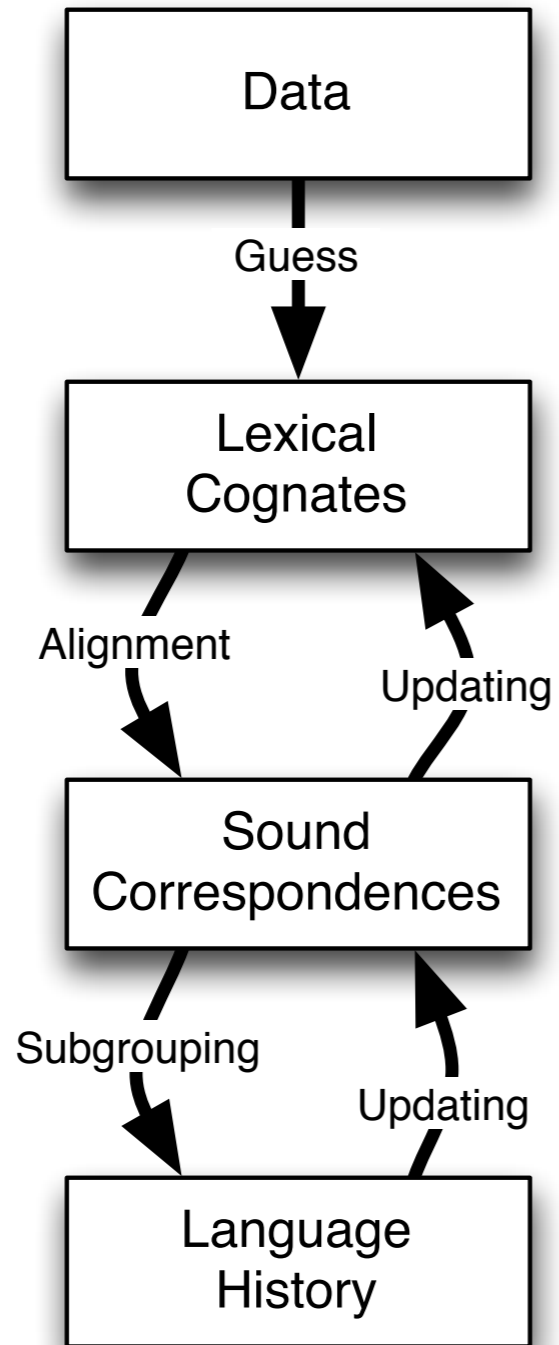
#m	mi	ii	in	no		ok	kt	i#	
#m	mi	in ^j		n ^j e	ee	ek ^h	k ^h ω	ωω	ω#

	#m	mi	ii	in	no	ok	kt	i#
#m	22	3	2	2	2	2	2	2
mi	4	12	2	2	5	1	1	1
in ^j	2	1	5	9	3	1	1	2
n ^j e	1	1	5	5	4	1	1	2
ee	3	3	3	3	6	2	2	2
ek ^h	1	2	1	1	4	2	3	2
k ^h ω	2	2	2	2	2	1	23	2
ωω	2	2	3	3	2	2	4	4
ω#	2	2	3	2	3	1	3	4

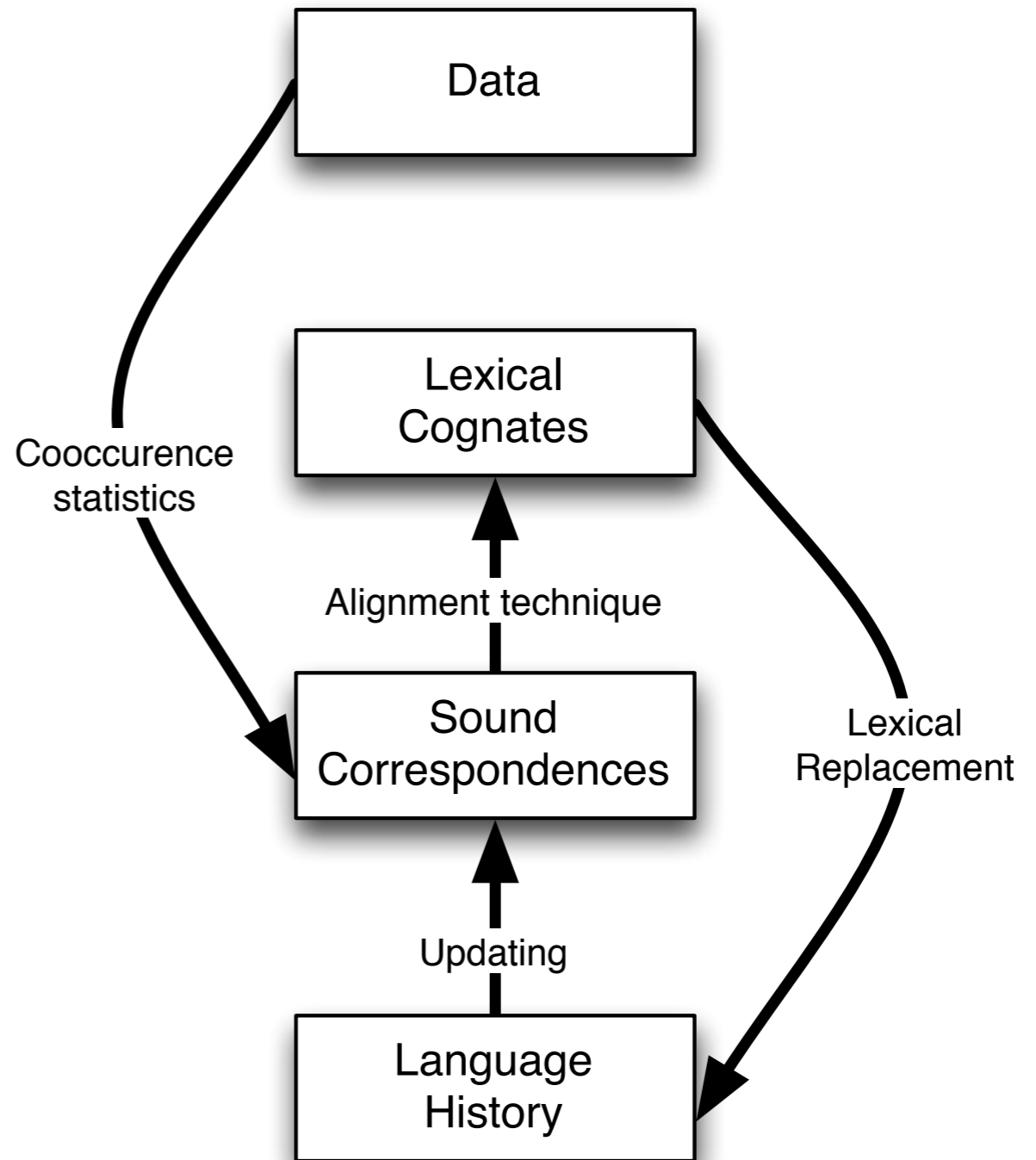
	Ocaina	Witoto Murui
HAND	oŋõõ	onoɕɯ
WE	xoxo	koko
HOUSE	φoo	φo
DOG	hõ?xo	hwko
JAGUAR	hõ?xo	hwko
FATHER	mõõ	moo
HUMMINGBIRD	φα?τίí?t'io	φiθido
TREE	amũũɲa	amena
STICK	amũũɲa	amena
WHO	bó	bu
SLEEP	uíúnõ	wnw
AGOUTI	φúúút'io	φwdo
THIS	bĩ	bie
THIS	ba?i	bie
NAME	maamɯ	mamekw
DAY	moona	aremona
BOW	tsipóxat'a	θwkuira
HEAR	xaaxa	kakade
DAY	moɲamó	aremona

	Ocaina	Witoto Murui
GREASE	φahĩ	φare
YOU (PLURAL)	mõ?	omow
THIS	bu	bie
ARROW	owd'áát'a	dukuraθw
SPEAR	owd'áát'a	dukuraθw
LIP	φα?óó?ko	φue igow
GREEN	moxóoso	mokorede
I	xõ	kue
ONE	t'a	dahe
WE	xo	kaw
TOOTH	a?tii?t'io	iθido
MOUTH	φoow	φue
BELLY	gááho	hebe
FATHER	mõõhõ	moo
YOU (PLURAL)	mõ?xo	omwko
SWAMP	xonuíúβaga	kunere
RAT	mɯɲóóko	mijwe
PATH, TRAIL	naahõ	nawθo
OWL	móóŋõhõ	monuiθw

Comparative Method



Sound Change Method



Conclusion

Conclusion

- *Regular sound change* is a very powerful notion to investigate language history

Conclusion

- *Regular sound change* is a very powerful notion to investigate language history
- **We knew that for a long time !**
we just seem to have forgotten about it in computational approaches

Conclusion

- *Regular sound change* is a very powerful notion to investigate language history
- **We knew that for a long time !**
we just seem to have forgotten about it in computational approaches
- *Regular symbol correspondence* can relatively easily be discovered statistically **before** cognate identification

Conclusion

- *Regular sound change* is a very powerful notion to investigate language history
- **We knew that for a long time !**
we just seem to have forgotten about it in computational approaches
- *Regular symbol correspondence* can relatively easily be discovered statistically **before** cognate identification
- **Reversing the comparative method**
first sound change, then cognacy judgement