

Lexikostatistik 2.0

Gerhard Jäger

Seminar für Sprachwissenschaft, Tübingen

IDS-Jahrestagung

12. März 2013

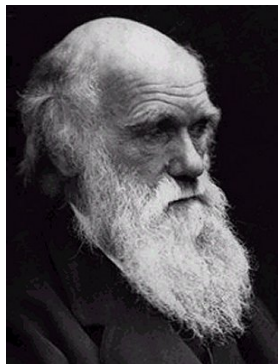


Sprachwandel und Evolution

„The formation of different languages and of distinct species, and the proofs that both have been developed through a gradual process, are curiously parallel. [...] We find in distinct languages striking homologies due to community of descent, and analogies due to a similar process of formation. The manner in which certain letters or sounds change when others change is very like correlated growth. [...] The frequent presence of rudiments, both in languages and in species, is still more remarkable. [...]

Languages, like organic beings, can be classed in groups under groups; and they can be classed either naturally according to descent, or artificially by other characters. Dominant languages and dialects spread widely, and lead to the gradual extinction of other tongues.“

(Darwin, The Descent of Man)



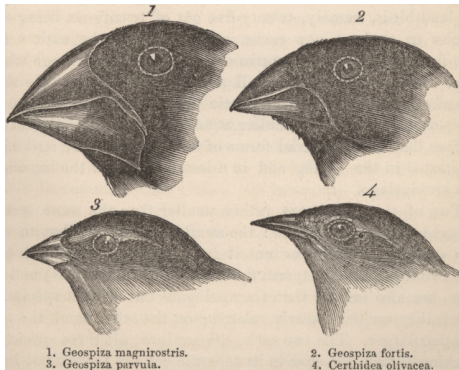
Sprachwandel und Evolution

Vater Unser im Himmel, geheiligt
werde Dein Name

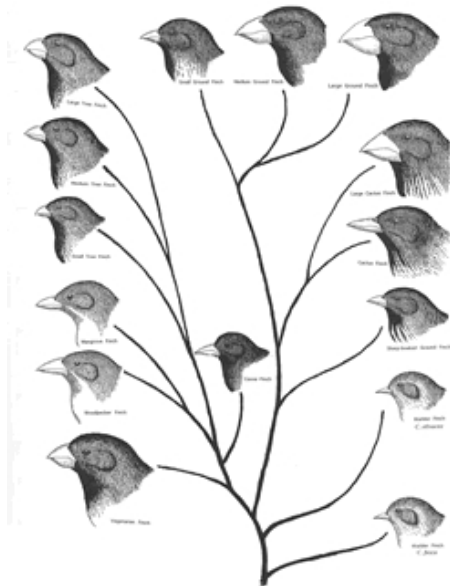
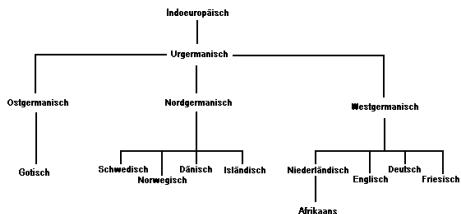
Onze Vader in de Hemel, laat Uw
Naam geheiligt worden

Our Father in heaven, hallowed be
your name

Fader Vor, du som er i himlene!
Helliget vorde dit navn



Sprachwandel und Evolution



Sprachwandel und Evolution

Mittelhochdeutsch:

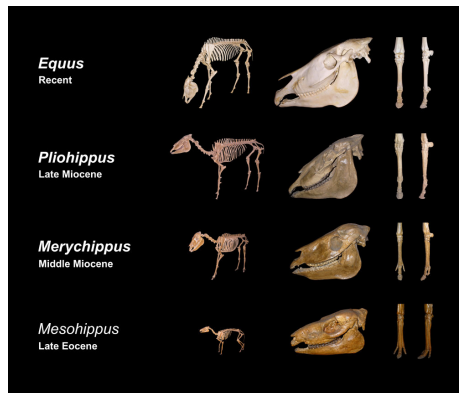
Got vater unser, dâ du bist in dem
himelrîche gewaltic alles des dir ist,
geheiliget sô werde dîn nam

Althochdeutsch:

Fater unser thû thâr bist in himile, si
giheilagôt thîn namo

Gotisch:

Atta unsar þu in himinam, weihnai
namo þein



Höherentwicklung im Sprachwandel

Pidgin- und Kreolsprachen

- eine Indianerin zu einem weißen Verehrer in Pidgin-English:

You silly. You weak. You baby-hand. No catch horse. No kill buffalo. No good but for sit still—read book.

- Satz aus dem Sranan, einer Englisch-basierten Kreolsprache aus Surinam:

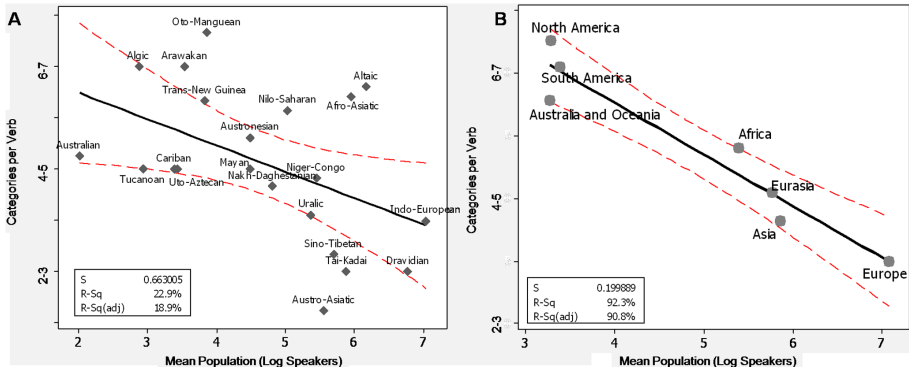
A hondiman datai ben bai wan oso gi en mati.

‘Der Jäger, der ein Haus gekauft hat, gab es seinem Freund.’

(aus John McWhorter, 2003, *The Power of Babel*)

Höherentwicklung im Sprachwandel

Anpassung der Grammatik an soziale Gegebenheiten



(aus G. Lyupan & R. Dale, 2010, PLoS ONE 5(1))

Konvergente Evolution

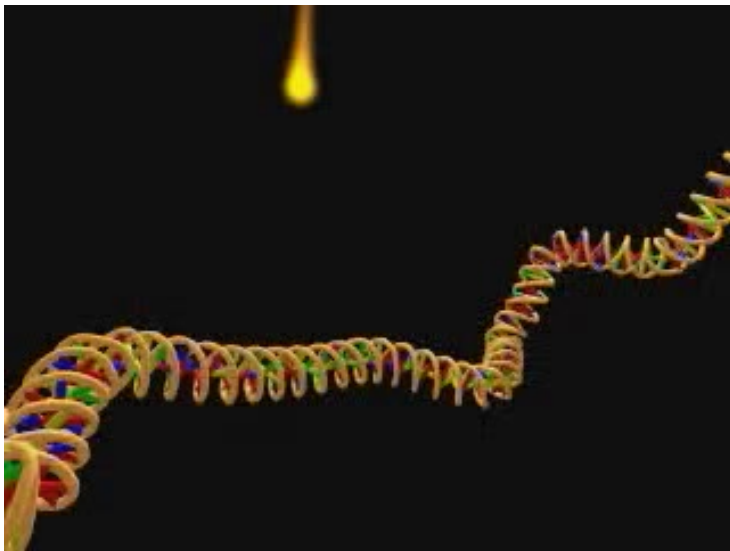
Taking Flight

To take to the air, three very different vertebrates lightened bones and transformed hands into wings.



- Altenglisch *docga* > Englisch *dog* ('Hund')
- Proto-Paman **gudaga* > Mbabaram *dog* ('Hund')

Evolution via Mutation in der Biologie



Lautgesetze

Erste bzw. Germanische Lautverschiebung (Indoeuropäisch → Germanisch)	Phase	Zweite bzw. Hochdeutsche Lautverschiebung (Germanisch → Althochdeutsch)	Beispiele (Neuhochdeutsch)	Jahrhundert	Dialektgebiete
G: /*b/ → /*p/	1	/*p/ → /ff/ → /f/	niederdeutsch: slapen , englisch: sleep → schlafen ; niederdeutsch und englisch: Schipp , ship → Schiff niederdeutsch: scherp , englisch: sharp → scharf	4/5	oberdeutsch und mitteldeutsch
	2	/*p/ → /pf/	niederdeutsch: Peper , englisch: pepper → Pfeffer ; niederdeutsch: Plauch , englisch: plough → Pflug ; niederdeutsch: scherp , englisch: sharp , althochdeutsch: scarph , mittelhochdeutsch: scharpf	6/7	oberdeutsch
G: /*d/ → /*t/	1	/*t/ → /ss/ → /s/	niederdeutsch: dat , wat , eten ; englisch: that , what , eat → das , was , essen	4/5	ober- und mitteldeutsch ¹
	2	/*t/ → /ts/	niederdeutsch: Tiet , englisch: tide (Gezeiten), schwedisch: tid → Zeit ; niederdeutsch: ver-tellen , englisch: tell → er-zählen ; Timmermann → Zimmermann	5/6	ober- und mitteldeutsch
G: /*g/ → /*k/	1	/*k/ → /xx/ → /x/	niederdeutsch: ik , altenglisch: ic → ich ; niederdeutsch und englisch: make n, make → machen ; niederdeutsch: auk → auch	4/5	ober- und mitteldeutsch ²
	2	/*k/ → /kx/	Kind → bairisch: Kchind	7/8	südbairisch, hoch- und höchstalemannisch
G: /*bʰ/ → /*b/ V: /*p/ → /*b/	3	/*b/ → /p/	Berg , bist → bairisch: perg , pist	8/9	teilweise bairisch und alemannisch
G: /*d/ → /*d/ → /*d/ V: /*t/ → /*d/ → /*d/	3	/*d/ → /t/	niederdeutsch: Dag oder Dach , englisch: day → Tag ; niederfränkisch: vader → Vater	8/9	oberdeutsch
G: /*gʰ/ → /*g/ V: /*k/ → /*g/	3	/*g/ → /k/	Gott → bairisch: Kott	8/9	teilweise bairisch und alemannisch
G: /*t/ → /p/ [ð]	4	/p/ → /d/ /ð/ → /d/	englisch: thorn , thistle , through , brother → Dorn , Distel , durch , Bruder	9/10	gesamtes deutsches Dialektkontinuum

Lautgesetze

- Lautgesetze sind spezifisch für eine bestimmte Sprachwandel-Periode
- gelten nahezu universell für alle Instanzen des betroffenen Lautes in der betroffenen Sprache
- im Idealfall gibt es schriftliche Zeugnisse der älteren und der jüngeren Sprachstufe (z.B. Latein/romanischen Sprachen, Althochdeutsch/Mittelhochdeutsch)
- meistens müssen Lautgesetze durch systematischen Vergleich verwandter Sprachen identifiziert werden
- erlaubt partielle Rekonstruktion der gemeinsamen Ursprungssprache

The Indo-European language family

- William Jones 1786:
„The Sanskrit Language, whatever be its antiquity, is of wonderful structure; more perfect than the Greek, more copious than the Latin, and more exquisitely refined than either; yet bearing to both of them a stronger affinity both in the roots of verbs and the forms of grammar, than could possibly have been produced by accident; so strong indeed that no philologer could examine them at all without believing them to have sprung from some common source, which perhaps no longer exists: there is similar reason, so not quite so forcible, for supposing that both the Gothic and the Celtic, though blended with a different idiom, had the same origin with the Sanskrit; and the old Persian might be added to the same family, if this were the place for discussing any question concerning the antiquities of Persia.“

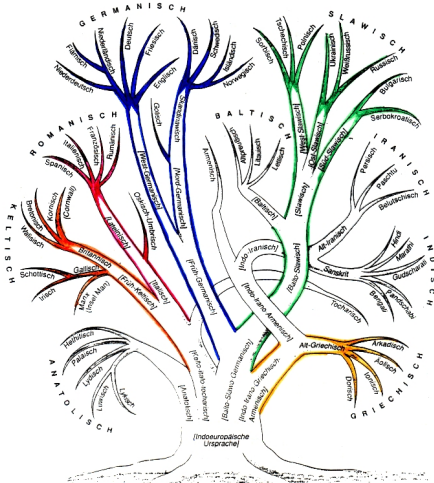
Sprachrekonstruktion durch die komparative Methode

- erste erfolgreiche Anwendung auf *Indo-europäisch* im 19. Jhd.

	Griechisch	Vedisch	Awestisch	Latein	Wallsisch	Gotisch	Armenisch	Tocharisch A	A. K. Slawisch	Litauisch	Indogermanisch (rekonstruiert)
1	heis (< *hens < *sems)	éka	aēuua	ūnus (Altlatein: oinos)	un	ains	mi	sas	inū	vianas	*oyno-, oyko-, sem-
2	dúō	dvá	duua	duō	dau	twai	erkow	wu	dūva	dù	*duwóh ₁
3	treis	tri	θrāiō	trēs	tri	breis	erek`	tre	trije	trýs	*tréyes
4	téttares	catváras	caθuuārō	quattuor	pedwar	fidwor	čork`	štwar	četyre	keturñ	*k ^w etwóres
5	pénte	pánca	panca	quinque	pump	fimf	hing	pāñ	peṭī	penki	*pénk ^w e
6	héks	ṣát	xšuuas	sex	chwech	saihs	več	šák	šestī	šeši	*swéks
7	heptá	ṣaptá	hapta	septem	saith	sibun	ewt`n	špät	sedmī	septyni	*septrij
8	októ	aṣṭá	ašta	octō	wyth	ahtau	owt`	okät	osmī	aštuoni	*októ
9	ennéa	náva	nauua	novem	naw	niun	inn	ñu	devęti	devyni	*néwn
10	déka	dása	daša	decem	deg	taihun	tasn	šák	desęti	dęsimt	*dékṃ
20	wikati (dorisches)	vimšati	vīšaiti	vigintī	ugeint (Mittelwallsisch)		k`san	wiki			*wikmṭi
100	hekatón	śatám	satəm	centum	cant	hund		kánt	sūto	šimtas	*kṃtóm

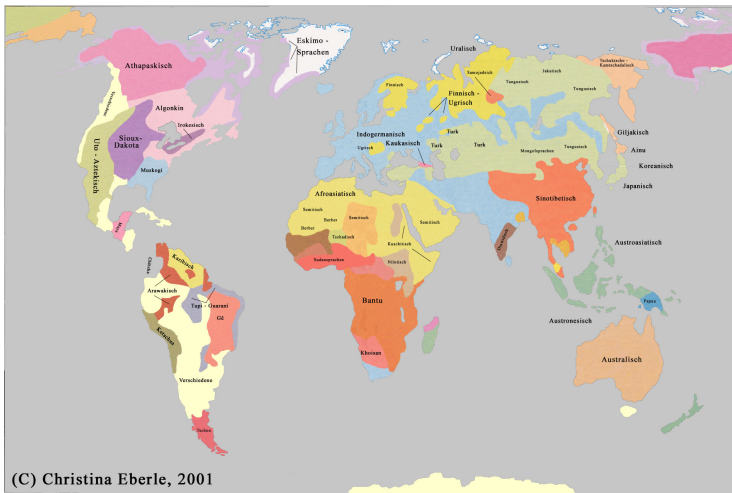
Language trees

- komparative Methode ergibt Abstammungsbaum einer Sprachfamilie



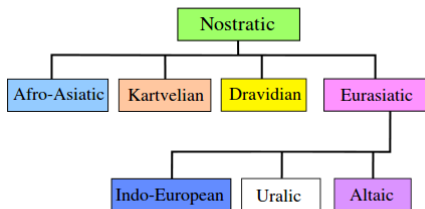
Grenzen der komparativen Methode

- Zeittiefe beschränkt auf 2 000 bis 8 000 Jahre



Tiefe Sprachverwandtschaften

- Vielzahl von Vorschlägen für Meta-Familien
 - Nostratisch:
 - erstmals von Pedersen (1903) vorgeschlagen
 - ursprünglicher Vorschläge: Indo-europäisch, Finno-ugrisch, Samoyedisch, Turk-Sprachen, Mongolisch, Manchu, Yukaghir, Eskimo, Semitisch und Hamitisch
 - weiterentwickelt durch „Moskauer Schule“ in den 1960ern
 - Versuch der Rekonstruktion von Wortschatz



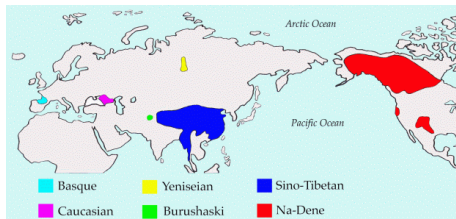
Tiefe Sprachverwandtschaften

- Vielzahl von Vorschlägen für Meta-Familien
 - Eurasiatisch
 - vorgeschlagen von Greenberg (2000)
 - umfasst Indo-europäisch, Uralisch-Yukaghirisch, Altaisch, Tschuktscho-Kamtschadalisch, Eskimo-Aleutisch, Koreanisch-Japanisch-Ainu, Gilyak, Etruskisch
 - diverse Argumente, v.a. Morphologie und Phonologie



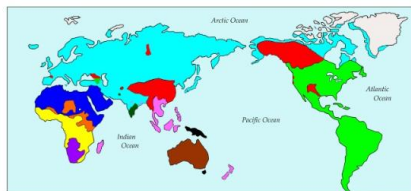
Tiefe Sprachverwandtschaften

- Vielzahl von Vorschlägen für Meta-Familien
 - Dene-Kaukasisch
 - umfasst Ne-Dene, Kaukasisch, Sino-Tibetisch, Jenniseiisch, Burushaski, manchmal auch Baskisch



Tiefe Sprachverwandtschaften

- Vielzahl von Vorschlägen für Meta-Familien
 - Amerindisch
 - vorgeschlagen von Greenberg (1987)
 - umfasst alle Indianersprachen außer Na-Dene



Khoisan	Dravidian	Austric
Niger-Kordofanian	Kartvelian	Indo-Pacific
Nilo-Saharan	Eurasiatic	Australian
Afro-Asiatic	Dene-Caucasian	Amerind

Language Families of the World (after Greenberg)

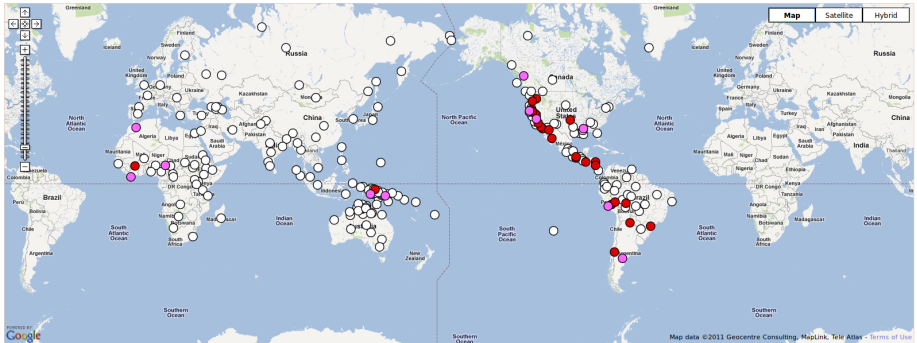


Tiefe Sprachverwandtschaften

- Merritt Ruhlen, ein Schüler von Greenberg, behauptet sogar, „Proto-World“ z.T. rekonstruieren zu können, z.B. das Wort *akwa* für Wasser (das sich faszinierenderweise von Adam und Eva über Cicero bis zu Umberto Eco im Indoeuropäisch/Italisch/Lateinisch/Italienischen Zweig nicht verändert hat)
- derartige Vorschläge basieren häufig auf geographischen Häufungen einzelner Merkmale, wie z.B. Pronominalformen

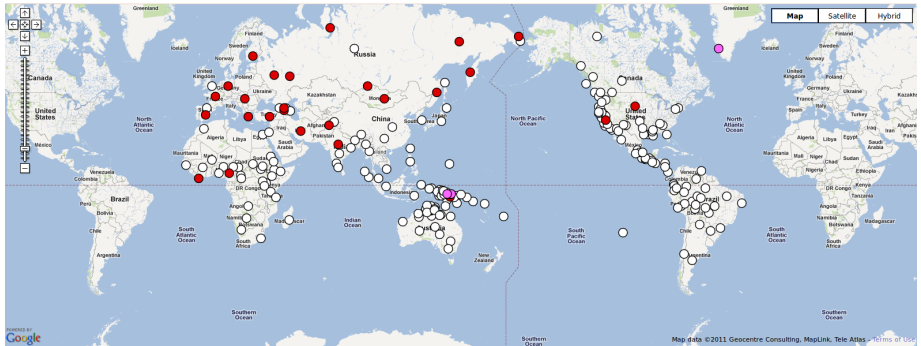
Tiefe Sprachverwandtschaften

- N/M-Pronomina



Tiefe Sprachverwandtschaften

● M/T-Pronomina



Phylogenetische Rekonstruktion in der Bioinformatik

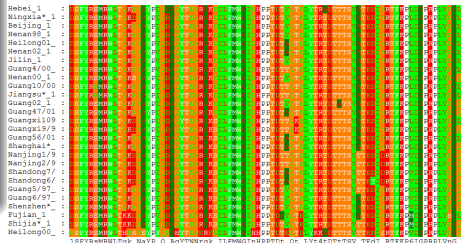
Sequenzalinierung

- Algorithmus findet optimale Alinierung zwischen Sequenzen
- Anzahl der Mutationen wird somit abgeschätzt
- ergibt Abschätzung des evolutionären Abstands zwischen den entsprechenden Organismen

```

Hebei_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Ningxia*_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Beijing_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Henan98_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Heilong01 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Henan02_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Jilin_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Guang4/00 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Henan00_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Guang10/00 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Jiangsu*_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Guang9_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Guang47/01 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Guangxi109 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Guangxi9/9 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Guang56/01 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Shanghai* : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Nanjing1/9 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Nanjing2/9 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Shandong7 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Shandong6 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Guang5/97 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Guang6/97 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Shenzhen* : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Fujian_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Shijia*_1 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
Heilong00 : DSFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG
ISFYRMRMLCKRNAYVFCQAAYTNNRGGILFPMGIIHFFPTDQQLMLVRRDITISVTEIIRRFYFELIGREPLVWG

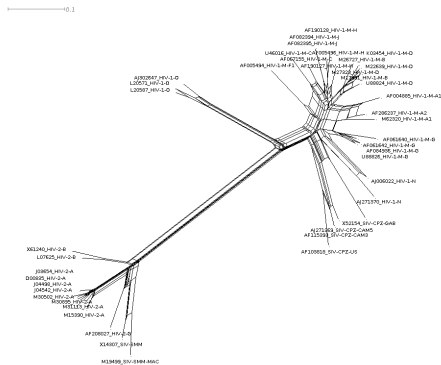
```



Phylogenetische Rekonstruktion in der Bioinformatik

Phylogenetische Bäume

- statistische Verfahren zur Rekonstruktion des wahrscheinlichsten Stammbaums
- häufig konfligierende Information wegen:
 - konvergenter Evolution
 - Rück-Mutation
 - lateraler Gen-Transfer
- Darstellung alternativer Rekonstruktionen in Netzwerk-Strukturen



SplitsTree Software, Huson & Bryant, MatNat-Fakultät

Phylogenetische Rekonstruktion in der Bioinformatik

Alternative: Cluster-Karten

- Organisation aller Datenpunkte (=Molekularsequenzen) in 2- oder 3-dimensionalen Raum
 - größere Ähnlichkeit entspricht (simulierter) physikalischer Anziehungskraft und umgekehrt
 - Algorithmus findet Energie-Minimum
- ⇒ verwandte Sequenzen bilden Cluster

Software: Frickey & Lupas, MPI für Entwicklungsbiologie

Die Daten des Automated Similarity Judgment Project

- Projekt am MPI EVA in Leipzig um Sören Wichmann
- erfasst inzwischen über 5 000 Sprachen
- für jede Sprache Grundwortschatz von 40 Wörtern in (vereinfachter) phonetischer Umschrift
- frei elektronisch verfügbar

verwendete Konzepte: *I, you, we, one, two, person, fish, dog, louse, tree, leaf, skin, blood, bone, horn, ear, eye, nose, tooth, tongue, knee, hand, breast, liver, drink, see, hear, die, come, sun, star, water, stone, fire, path, mountain, night, full, new, name*

Automated Similarity Judgment Project

<i>Konzept</i>	Deutsch	Englisch
<i>I</i>	iX	Ei
<i>you</i>	du	yu
<i>we</i>	vir	wi
<i>one</i>	ains	8is
<i>two</i>	cvai	8Et
<i>person</i>	mEnS	pers3n
<i>fish</i>	fiS	fiS
<i>dog</i>	hunt	dag
<i>louse</i>	laus	laus
<i>tree</i>	baum	tri
<i>leaf</i>	blat	lif
<i>skin</i>	haut	skin
<i>blood</i>	blut	bl3d
<i>bone</i>	knoX3n	bon
<i>horn</i>	horn	horn
<i>ear</i>	XXX	ir
<i>eye</i>	aug3	Ei

<i>Konzept</i>	Deutsch	Englisch
<i>nose</i>	naz3	nos
<i>tooth</i>	ch~an	tu8
<i>tongue</i>	ch~uN3	t3N
<i>knee</i>	kni	ni
<i>hand</i>	hant	hEnd
<i>breast</i>	brust	brEst
<i>liver</i>	leb3r	liv3r
<i>drink</i>	triNk3n	drink
<i>see</i>	ze3n	si
<i>hear</i>	her3n	hir
<i>die</i>	Sterb3n	dEi
<i>come</i>	kh~om3n	k3m
<i>sun</i>	zon3	s3n
<i>star</i>	StErn	star
<i>water</i>	vas3r	wat3r
<i>stone</i>	Stain	ston
<i>fire</i>	foia	fEir

Einfache Sequenz-Alinierung

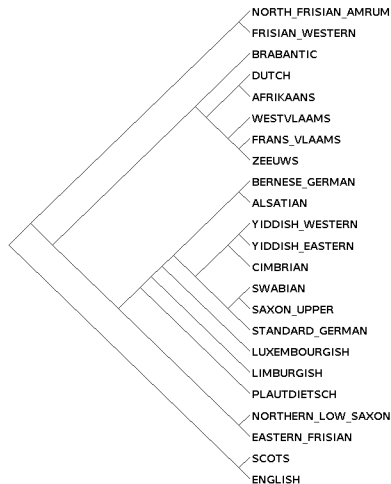
- z.B. Deutsch ↔ Latein

h	o	r	n	
k	o	r	n	u

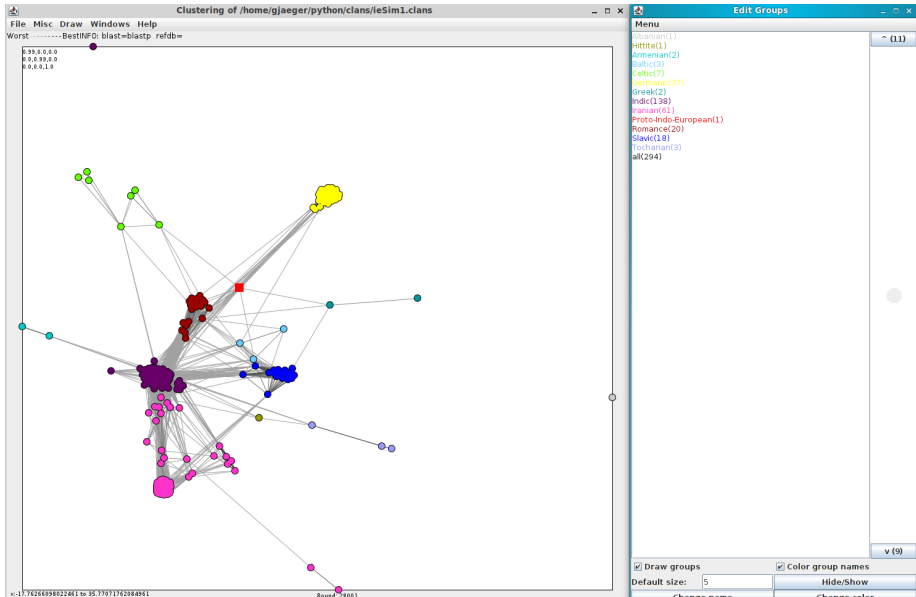
- Anzahl der Unterschiede: 2
- durchschnittliche Abweichung = $2/5$
- Ähnlichkeit zwischen zwei Sprachen wird durch durchschnittliche Abweichung der jeweiligen Übersetzungspaare geschätzt

Einfache Sequenz-Alinierung

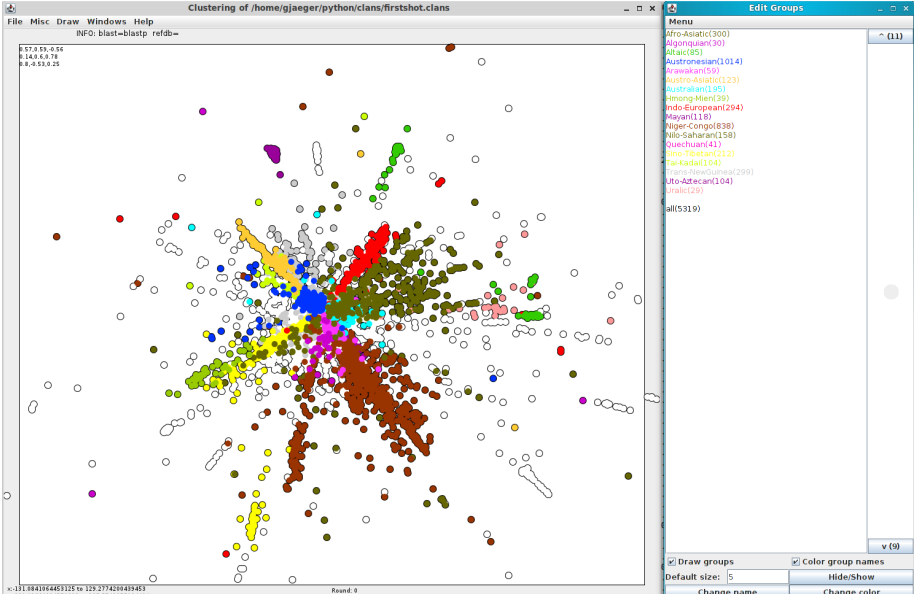
- für die westgermanischen Sprachen/Dialekte ergibt sich die Phylogenie



Einfache Sequenz-Alinierung

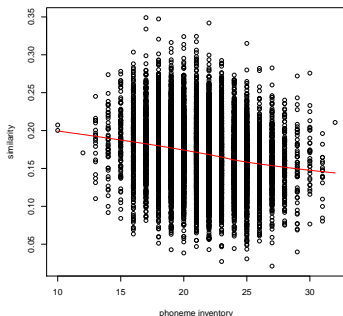


Einfache Sequenz-Alinierung



Einfache Sequenz-Alinierung

- Störeffekt: bei Sprachen mit kleineren Lautinventaren ergeben sich mehr Zufallsähnlichkeiten also bei Sprachen mit vielen verschiedenen Lauten
- Daher erscheinen Sprachen mit wenigen Lauten einander ähnlicher, als sie es tatsächlich sind.



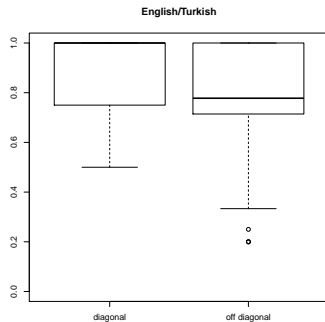
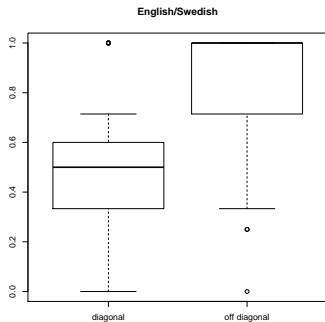
Kalibrierte Alinierung

Englisch / Swedisch

	Ei	yu	wi	w3n	tu	fiS	...
yog	1	2/3	1	1	1	1	
du	1	1/2	1	1	1/2	1	
vi	1/2	1	1/2	1	1	2/3	
et	1	1	1	1	1	1	
tvo	1	1	1	1	2/3	1	
fisk	3/4	1	3/4	1	1	1/2	
	:						

- je näher zwei Sprachen verwandt sind, umso mehr unterscheiden sich die Werte auf der Diagonale von den restlichen Werten

Kalibrierte Alinierung



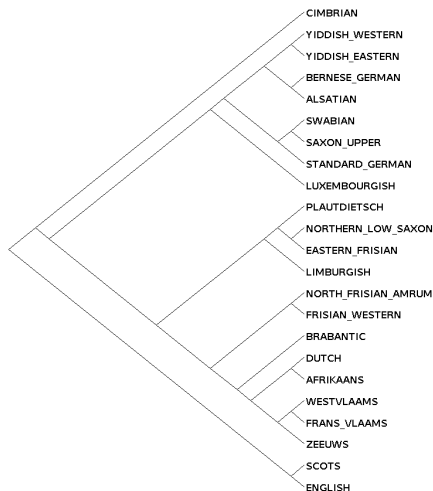
Kalibrierte Alinierung

- Ähnlichkeit des Standard-Deutschen zu:

Schwäbisch	26,13
Zimbrisch	20,28
Niederländisch	23,75
Englisch	17,45
Ur-Indoeuropäisch	10,26
Latein	9,23
Spanisch	8,95
Hindi	8,70
Russisch	8,36
Türkisch	6,33
Ungarisch	6,84

Kalibrierte Alinierung

- für die westgermanischen Sprachen/Dialekte ergibt sich die Phylogenie

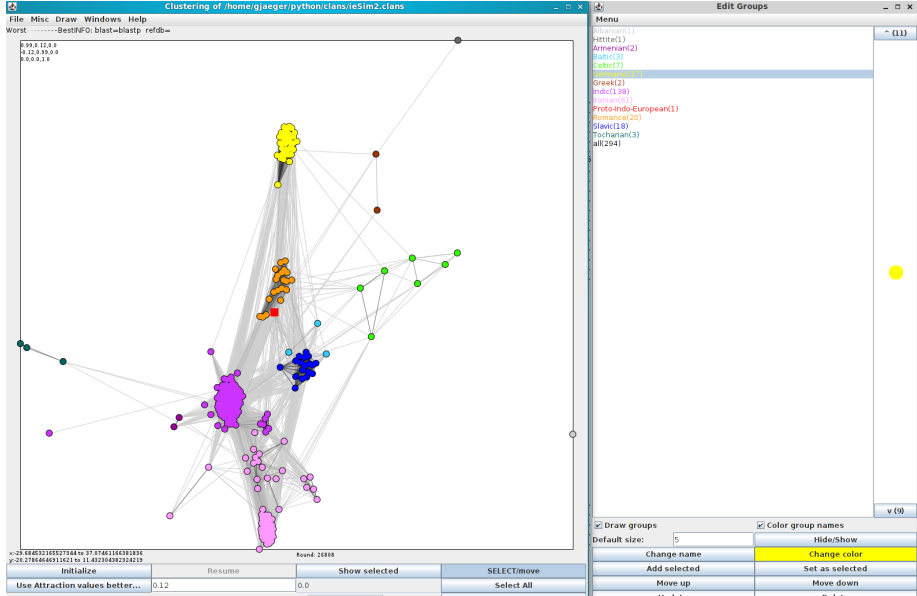


Problemfall: Zimbrisch

<i>Konzept</i>	Deutsch	Zimbrisch
<i>I</i>	iX	ix
<i>you</i>	du	du
<i>we</i>	vir	bar
<i>one</i>	ains	XXX
<i>two</i>	cvai	sben
<i>person</i>	mEnS	menEs
<i>fish</i>	fiS	XXX
<i>dog</i>	hunt	hunt
<i>louse</i>	laus	laus
<i>tree</i>	baum	pom
<i>leaf</i>	blat	placa
<i>skin</i>	haut	XXX
<i>blood</i>	blut	plut
<i>bone</i>	knoX3n	poan
<i>horn</i>	horn	horn
<i>ear</i>	XXX	oar
<i>eye</i>	aug3	ogh~E

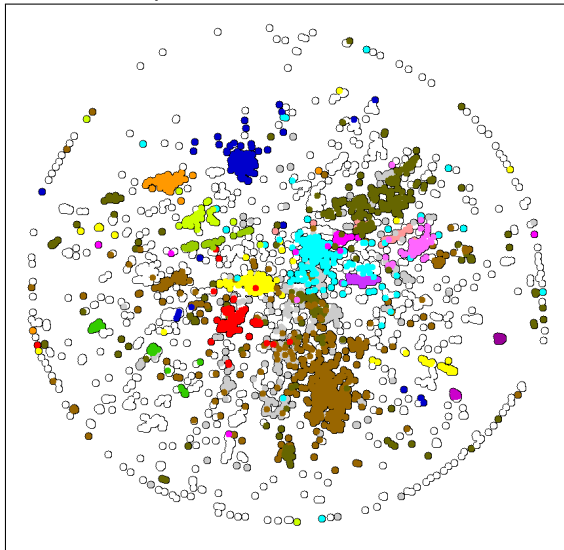
<i>Konzept</i>	Deutsch	Zimbrisch
<i>nose</i>	naz3	naza
<i>tooth</i>	ch~an	XXX
<i>tongue</i>	ch~uN3	suNa
<i>knee</i>	kni	XXX
<i>hand</i>	hant	hant
<i>breast</i>	brust	prust
<i>liver</i>	leb3r	IEbara
<i>drink</i>	triNk3n	trinkh~
<i>see</i>	ze3n	zeg
<i>hear</i>	her3n	hor
<i>die</i>	Sterb3n	sterb
<i>come</i>	kh~om3n	kh~Em
<i>sun</i>	zon3	zuna
<i>star</i>	StErn	stErna
<i>water</i>	vas3r	basar
<i>stone</i>	Stain	stoan
<i>fire</i>	foia	boar

Kalibrierte Alinierung



Kalibrierte Alinierung

Alle 5 000 Sprachen:



Edit Groups

Menu

- Afro-Asiatic(300)
- Algonquian(30)
- Altaic(85)
- Austronesian(1023)
- Arawakan(59)
- Austra-Asiatic(123)
- Australian(198)
- Hmong-Mien(39)
- Indo-European(294)
- Mayan(118)
- Niger-Congo(838)
- Nilo-Saharan(159)
- Quechuan(41)
- Sino-Tibetan(212)
- Uralic(29)**
- Trans-New Guinea(299)
- Uto-Aztecan(104)
- Uralic(29)

all(5381)

v (9)

Draw groups Color group names

Default si... 5 Hide/Show

Change name	Change color
Add selected	Set as selected
Move up	Move down
Update	Delete

Kalibrierte Alinierung

- etablierte Sprachfamilien bilden stabile Cluster
- keine darüber hinausgehenden sichtbaren Muster

Kalibrierte Alinierung

- Methode ist relativ grobkörnig

h	a	n	t	h	a	n	t
h	E	n	d	m	a	n	o

- Ähnlichkeit ist in beiden Fällen 50%
- Korrespondenz $a \sim E$, $t \sim d$ sind nach linguistischen Kriterien viel natürlicher als $h \sim m$ or $t \sim o$
- Deutsch/Englisch und Deutsch/Spanisch erscheinen hier äquidistant, obwohl die Ähnlichkeit zwischen Deutsch und Englisch intuitiv viel größer ist

Needleman-Wunsch-Algorithmus

- Analogie zur Bioinformatik: Mutationen zwischen verschiedenen Aminosäuren-Paaren sind unterschiedlich wahrscheinlich
- Algorithmus sucht *wahrscheinlichste* Übereinstimmungen zwischen Sequenzen
- $a \sim E$, $d \sim t$ sind im Sprachwandel wahrscheinlicher als $t \sim o$

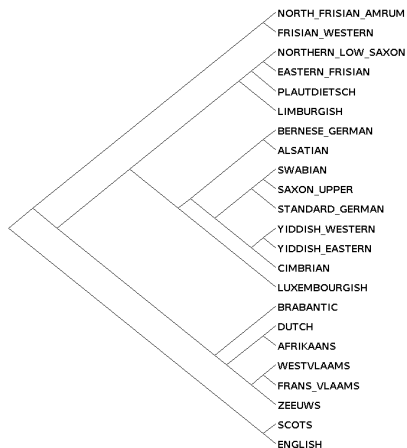
Gewichtete Alinierung

- Ähnlichkeit des Standard-Deutschen zu:

	ungewichtet	gewichtet
Schwäbisch	26,13	35,44
Zimbrisch	20,28	31,86
Niederländisch	23,75	29,76
Englisch	17,45	22,14
Ur-Indoeuropäisch	10,26	15,86
Latein	9,23	12,54
Spanisch	8,95	9,48
Hindi	8,70	12,35
Russisch	8,36	11,89
Türkisch	6,33	5,76
Ungarisch	6,84	7,57

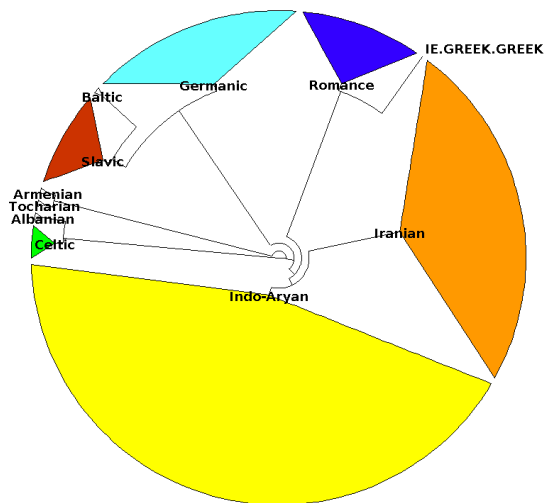
Gewichtete Alinierung

- für die westgermanischen Sprachen/Dialekte ergibt sich die Phylogenie



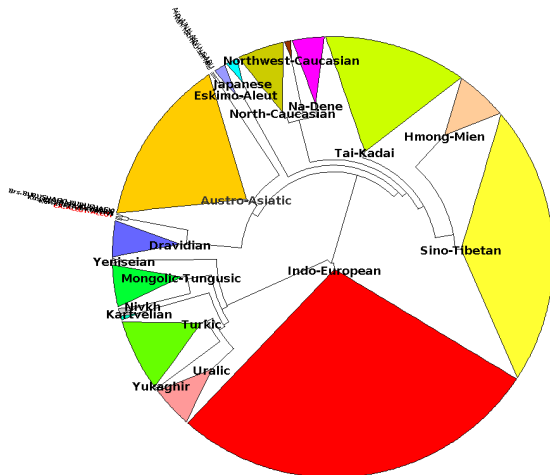
Gewichtete Alinierung

automatische Klassifikation der indo-europäischen Sprachen



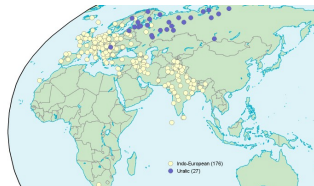
Gewichtete Alinierung

Sprachen Eurasiens (ohne Afro-Asiatisch) + Na-Dene + Eskimo-Aleutisch



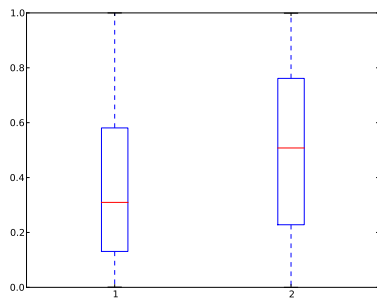
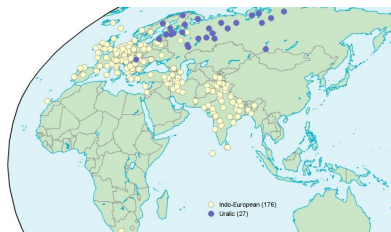
Gewichtete Alinierung

- einige interessante Meta-Verwandtschaften werden sichtbar, v.a.
 - Indo-europäisch/Uralisch
 - Austronesisch/Tai-Kadai



Indo-europäisch/Uralisch

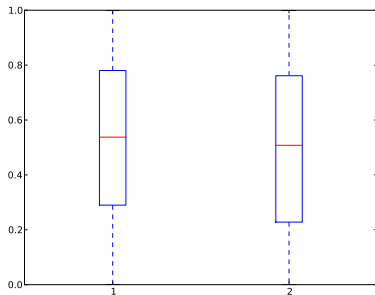
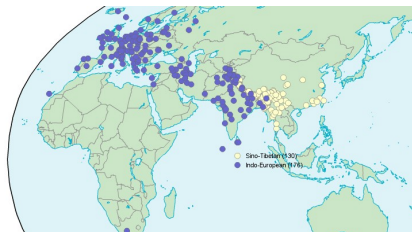
- p -Werte für Vergleich Ähnlichkeiten IE/Ura vs. Zufallspaarungen



p -Wert: $1,5 \times 10^{-20}$

Indo-europäisch/Uralisch

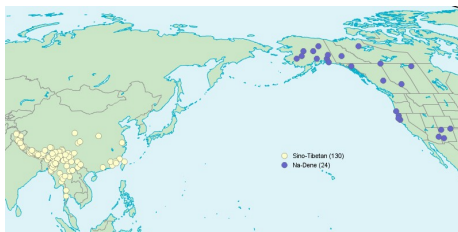
- p -Werte für Vergleich Ähnlichkeiten IE/Sino-Tibetisch vs. Zufallspaarungen



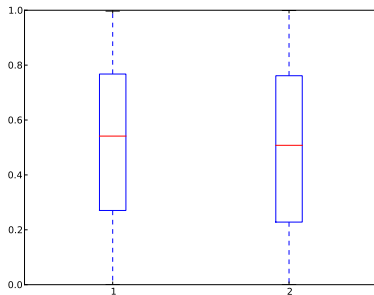
p -Wert: 1

Indo-europäisch / Uralisch

- p -Werte für Vergleich Ähnlichkeiten Sino-Tibetisch/Na-Dene vs. Zufallspaarungen



p -Wert: 1

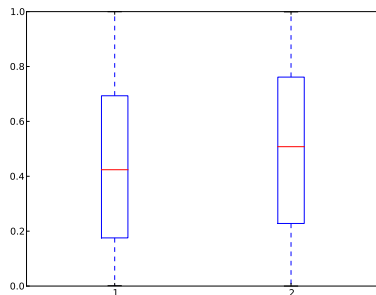


Indo-europäisch/Uralisch

- p -Werte für Vergleich Ähnlichkeiten Tai-Kadai/Austronesisch vs. Zufallspaarungen



p -Wert: 5×10^{-5}



Zusammenfassung

- Erfolgsbilanz:
 - von ca. 1.000 eurasiatischen Sprachen + Na-Dene + Eskimo-Aleutisch wird genau eine Sprache falsch klassifiziert (Aleutisch)
 - Altaisch wird nicht als Einheit erkannt, ist aber auch unter Experten umstritten
 - für alle 5,000 Sprachen ist die Fehlklassifikations-Quote bei 4%

Lexikostatistik funktioniert.

- suggestive Hinweise auf tiefe Verwandtschaften (Nostratisch, Dene-Kaukasisch etc.)
- statistisch wesentlich weniger robust als etablierte Sprachfamilien
- Einfluss von Sprachkontakt?

- starke lexikostatistische Evidenz für die Klassifikationseinheiten, die unter Experten unkontrovers sind
 - schwache Evidenz für umstrittene Klassifikationen
- ⇒ Lexikostatistik detektiert die selbe Evidenz, die auch von der komparativen historischen Linguistik genutzt wird