

Linguists amass data ...

- Typological databases (WALS)
- Lexical databases (IDS, LWT)
- Textual material (DoBeS, Rausing)

**... and there is even more
available**

- **Full-text searching of grammars/dictionaries**
- **(Online) electronic dictionaries**
- **The web itself as a corpus**

Although there are still some problems ...

- Management of sources
- Orthography
- Terminology



**How can we use all this data
for language comparison ?**

Example 1

**Typology from
language-particular
similarities**

Inchoative - causative verb pairs

- Inchoative
“The stick broke.”
- Causative
“The girl broke the stick.”

Based on data from: Haspelmath, Martin. 1993. "More on the typology of inchoative/causative verb alternations." In: Comrie, B. & Polinsky, M. (eds.) *Causatives and transitivity*. Amsterdam: Benjamins, 87-120.

Verb pairs investigated

begin	dry	melt
boil	fill	open
break	finish	rise/raise
burn	freeze	rock
change	gather	roll
close	get lost/	sink
connect	lose	split
destroy	go out/put	spread
develop	out	stop
die/kill	improve	turn
dissolve	learn/teach	wake up

Arabic**Class A: C / CC**

1. saḥaa / saḥḥaa
8. darasa / darrasa
14. damara / dammara
31. waqafa / waqqafa

Class B: in / ø

2. inkasara / kasara
5. infataḥa / fataḥa
6. inqafala / qafala
13. inṣahara / ṣahara
30. inṣaqqā / ṣaqqā

Class C: t / ?

3. iḥtaraqa / ?aḥraqa
22. intahaa / ?anhaa

Class D: t / ø

9. iltamma / lamma
10. intaṣara / naṣara
17. irtabaṭa / rabaṭa
21. irtafaṭa / rafaṭa
27. imtalaʔa / malaʔa

Class E: ø / ?

11. ġariqa / ?aġraqa
18. ġalaa / ?aġlaa
23. daara / ?adaara
26. ḍaaba / ?aḍaaba

Class F: ta / ø

12. tabaddala / baddala
16. taṭawwara / ṭawwara
19. taʔarjaḥa / ?arjaḥa
24. tadaḥraja / daḥraja
25. tajammada / jammada
28. taḥassana / ḥassana

Singular cases:

4. maata / qatala
7. badaʔa
15. daaṣa / xasira
20. inṭafaʔa / ?aṭfaʔa
29. jaffa / jaffafa

English**Class A: Identical**

1. wake up
2. break
3. burn
5. open
6. close
7. begin
9. gather
10. spread
11. sink
12. change
13. melt
16. develop
17. connect
18. boil
19. rock
22. finish
23. turn
24. roll
25. freeze
26. dissolve
27. fill
28. improve
29. dry
30. split
31. stop

Singular cases:

4. die / kill
8. learn / teach
14. be destroyed / destroy
15. get lost / lose
20. go out / put out
21. rise / raise

Finnish**Class A: ø / ttA**

1. herätä / herättää
3. palaa / polttaa
8. oppia / opettaa
10. levitä / levittää
13. sulaa / sulattaa
18. kiehua / kiehua
19. kiikkua / kiikuttaa
20. sammua / sammuttaa
21. kohota / kohottaa
22. loppua / lopettaa
24. vierä / vierittää
25. jäätyä / jäädyttää
26. liueta / liuottaa
31. pysähtyä / pysähdyttää

Class B: U / A

2. murtua / murtaa
12. muuttua / muuttaa
16. kehittyä / kehittää
23. vääntyä / vääntää
27. täyhtyä / täyttää
28. parantua / parantaa

Class C: UtU / ø

5. avautua / avata
6. sulkeutua / sulkea
14. tuhoutua / tuhota

Singular cases:

4. kuolla / tappaa
7. alkaa / aloittaa
9. kokoontua / koota
11. laskea
15. hukkaantua / hukata
17. yhtyä / yhdistää
29. kuivaa / kuivata
30. haljeta / halkaista

French**Class A: se / ø**

1. se réveiller / réveiller
2. se briser / briser
5. s'ouvrir / ouvrir
6. se fermer / fermer
9. s'assembler / assembler
10. s'étendre / étendre
11. s'enfoncer / enfoncer
15. se perdre / perdre
16. se développer / développer
17. se lier / lier
19. se balancer / balancer
20. s'éteindre / éteindre
21. se lever / lever
23. se tourner / tourner
26. se dissoudre / dissoudre
27. se remplir / remplir
28. s'améliorer / améliorer
30. se fendre / fendre
31. s'arrêter / arrêter

Class B: Identical

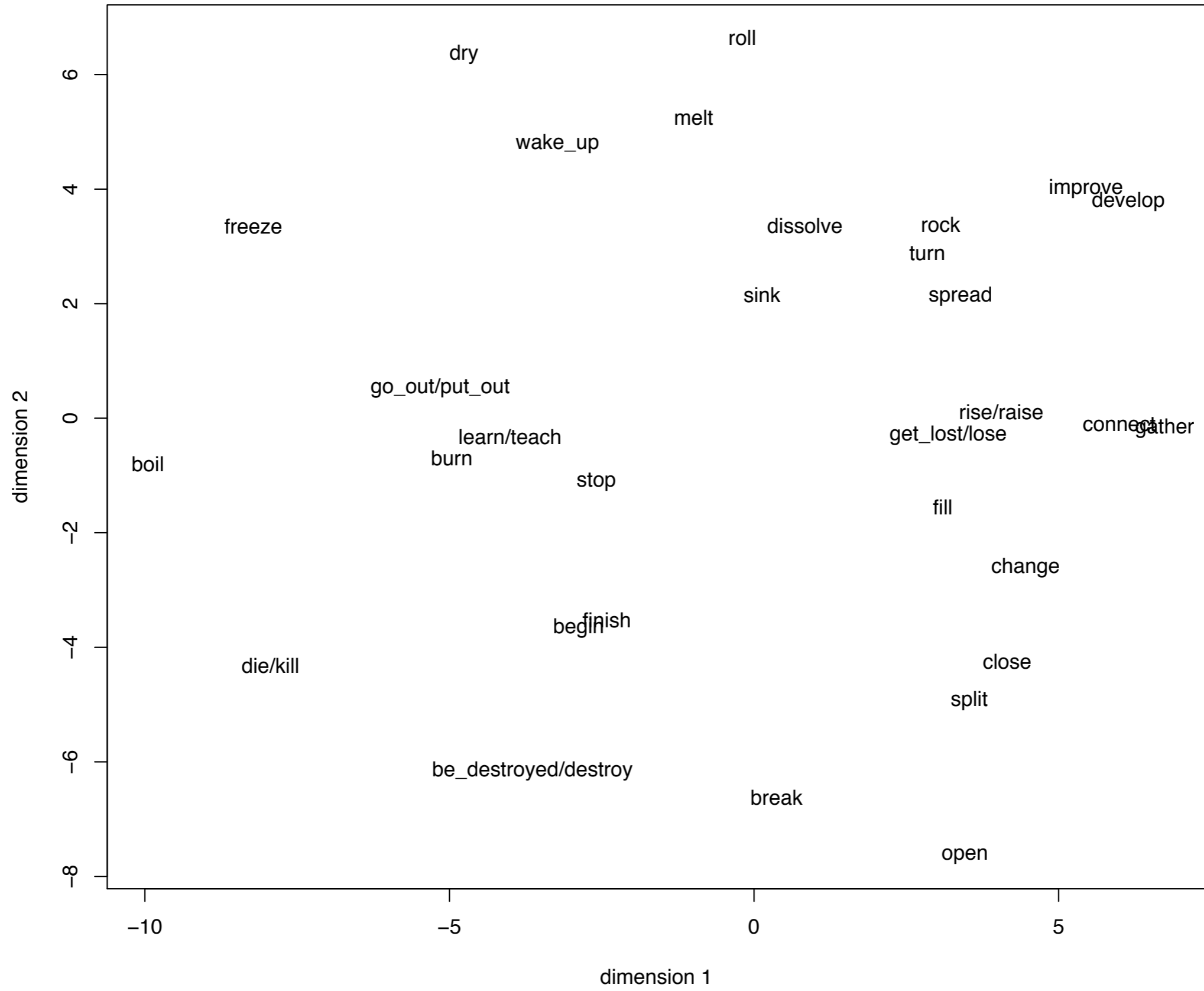
3. brûler
7. commencer
8. apprendre
12. changer
22. finir
24. rouler
25. geler
29. sécher

Class C: ø / faire

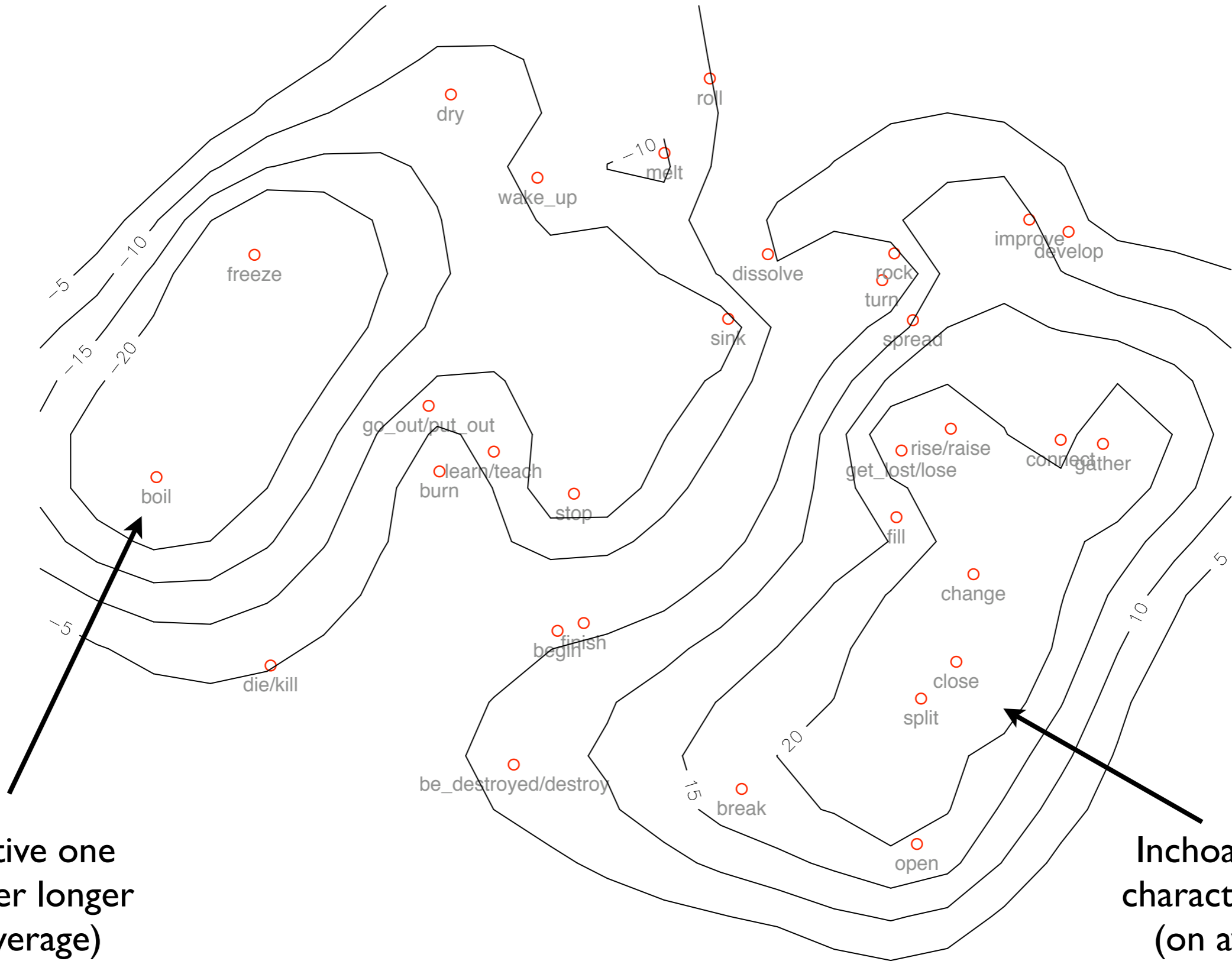
13. fondre / faire fondre
18. bouillir / faire bouillir

Singular cases:

4. mourir / tuer
14. être détruit / détruire



Inchoative-Causative Character Count



Causative one character longer (on average) than inchoative

Inchoative one character longer (on average) than causative

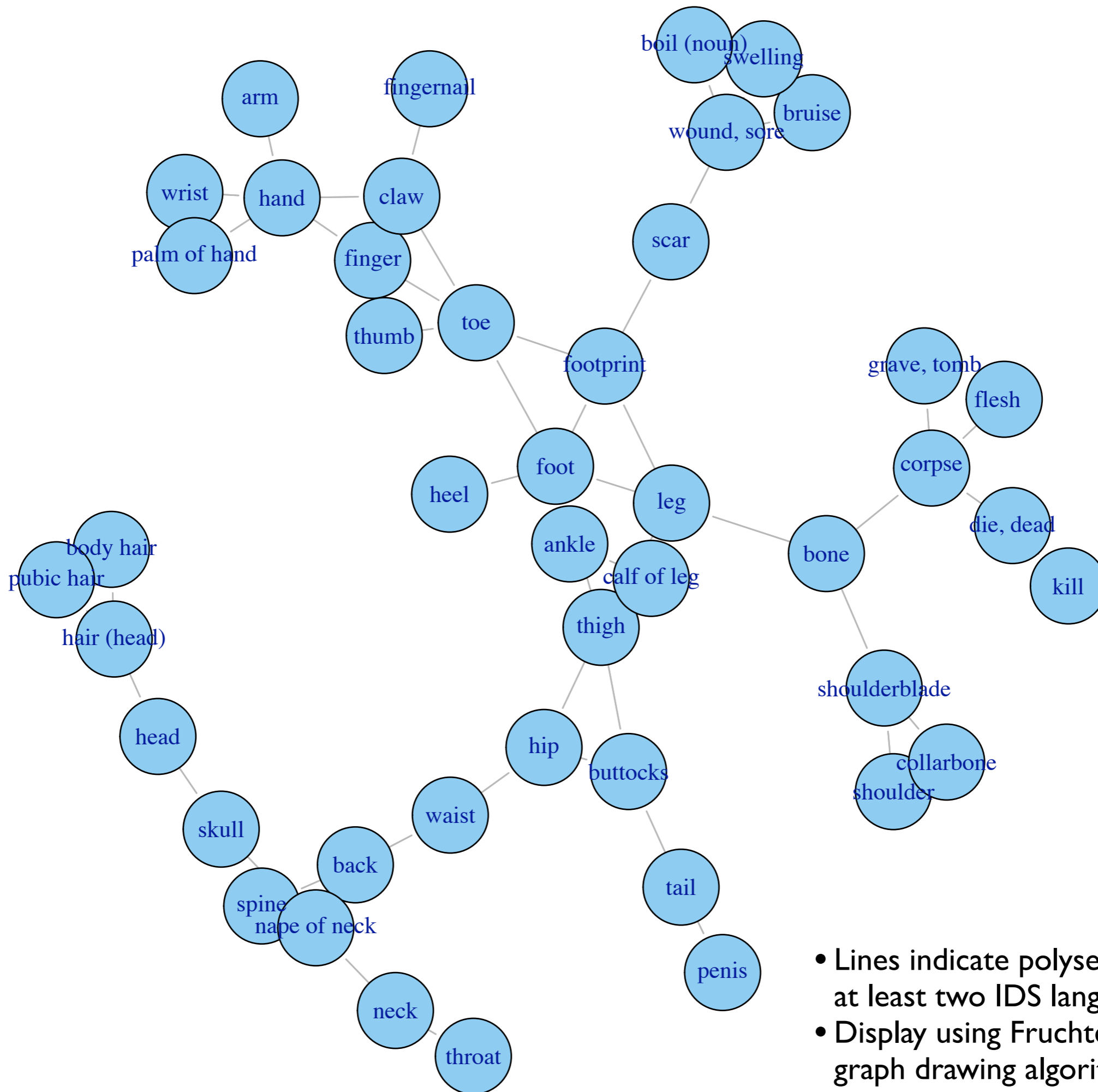
Example 2

Semantics from recurrent polysemy

using the Intercontinental Dictionary Series



chapter_id	entry_id	trans_english	Aché [Phonemic]	Aguaruna [Phonemic]	Ayoreo [Phonemic]
4	110	body	eče	iyaš	gaay; 'uud-ode?
4	120	skin, hide	pire	duwap-ɨ pušihi	a'ooi
4	130	flesh	oʔo	nɨhɨ	a'ooi; a'raganua
4	140	hair (head)	aʔa; tō aʔa	intaš	gaa'terero
4	142	beard	buta	intaš	abu'haa-sii
4	144	body hair	eče aʔa	susu	o'hoi
4	145	pubic hair	ta aʔa		
4	146	dandruff		du'sɨ buuk	
4	150	blood	bu-ɨʔɨ	numpa	i'yoii
4	151	vein, artery	ga-ɨʃɨ	num'pa hinti	'ku-kura-y
4	160	bone	ikā	ukunč	aṅokeei
4	162	rib	rukā	pagaɨ	orotabi'di
4	170	horn	ačī	kaču	u'ei
4	180	tail	buaʔa	uhu'ki	ka'ri
4	190	back	piče	tuntup; maya-tai	giido'booi
4	191	spine	piče ikā	tagkihi	oo'hi
4	200	head	tōʔō	buuk	gaa'toi
4	201	side of head, temple	tōʔō ēbe	kahašik	
4	202	skull	tō gape	buuk sakahu	
4	203	brain	tōʔō	buɛuk	ta'rooi



- Lines indicate polysemy in at least two IDS languages
- Display using Fruchterman-Reingold graph drawing algorithm

Example 3

Language similarities by counting word-edits using the Intercontinental Dictionary Series



abcdefgh \longleftrightarrow bacdxefgh

one addition

abcdefgh

abcd~~x~~efgh

one deletion

abcdx~~e~~fgh

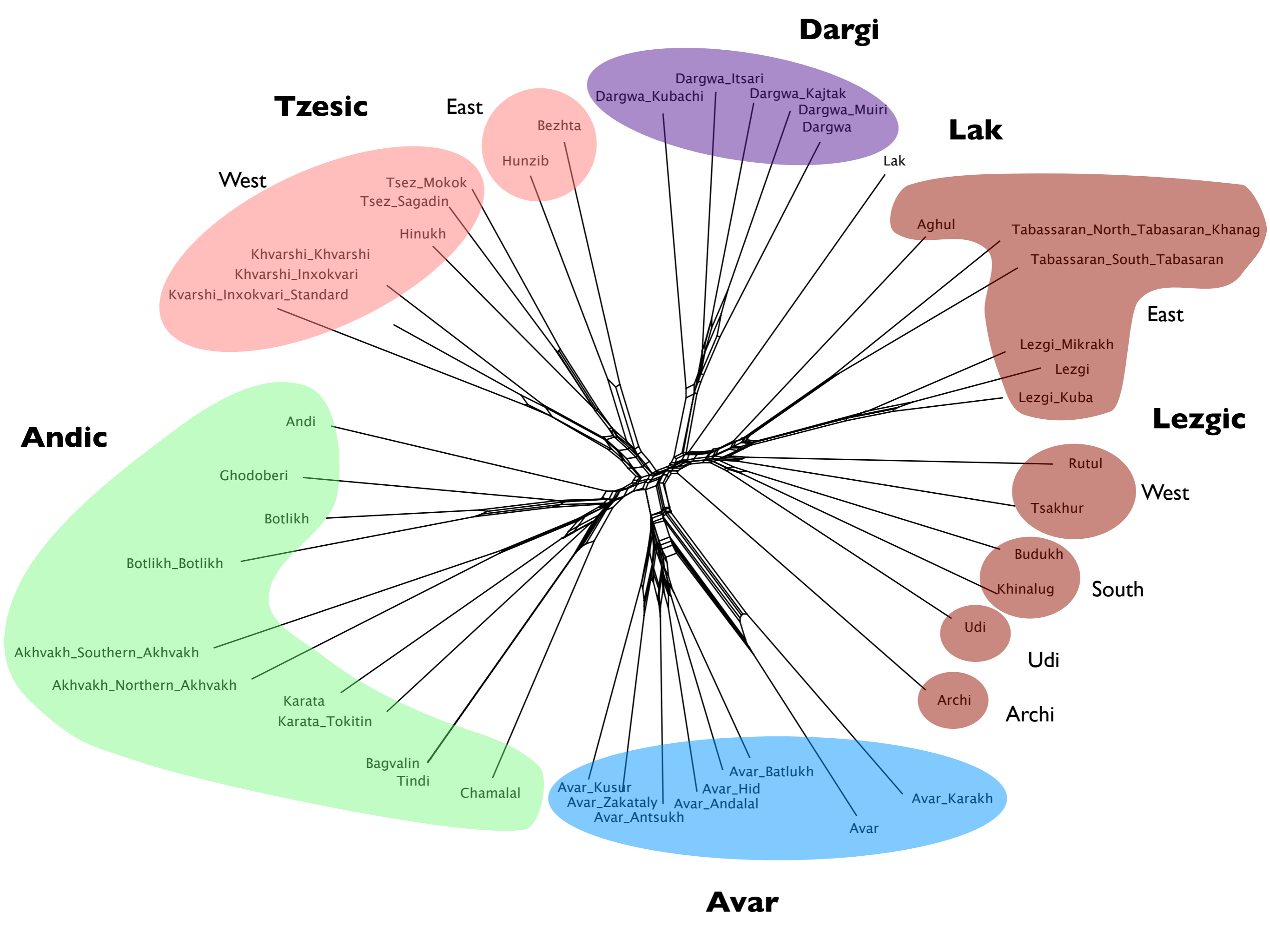
abcdxfgh

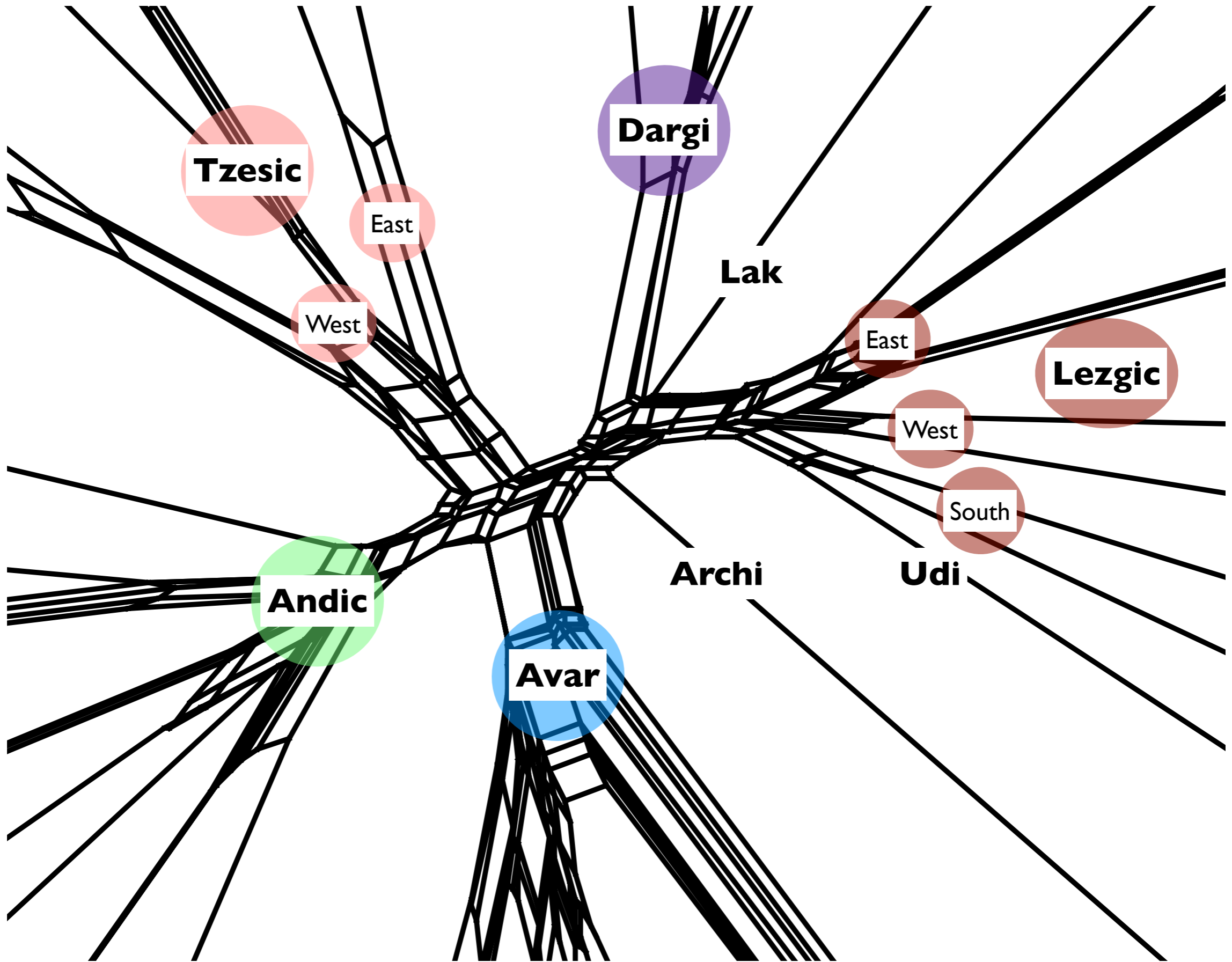
one inversion

abcdxfgh

bacdxefgh

Total of three changes needed



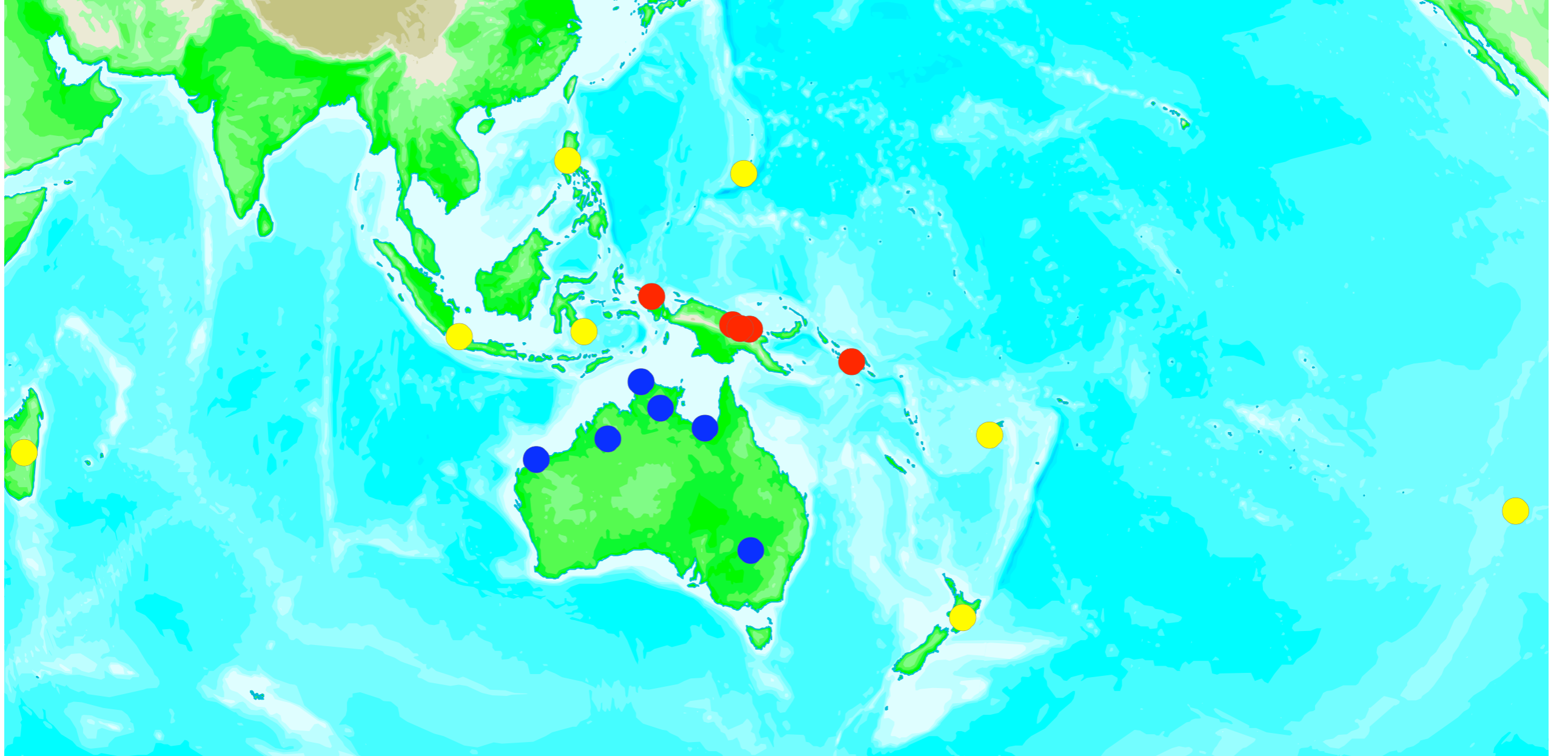


Example 4

Language similarities

using the World Atlas of Language Structures







● Maybrat

● Tukang Besi

● Tiwi

● Kobon

● Amele

● Lavukaleve

● Alamblak

● Mangarrayi

● Gooniyandi

● Kayardild

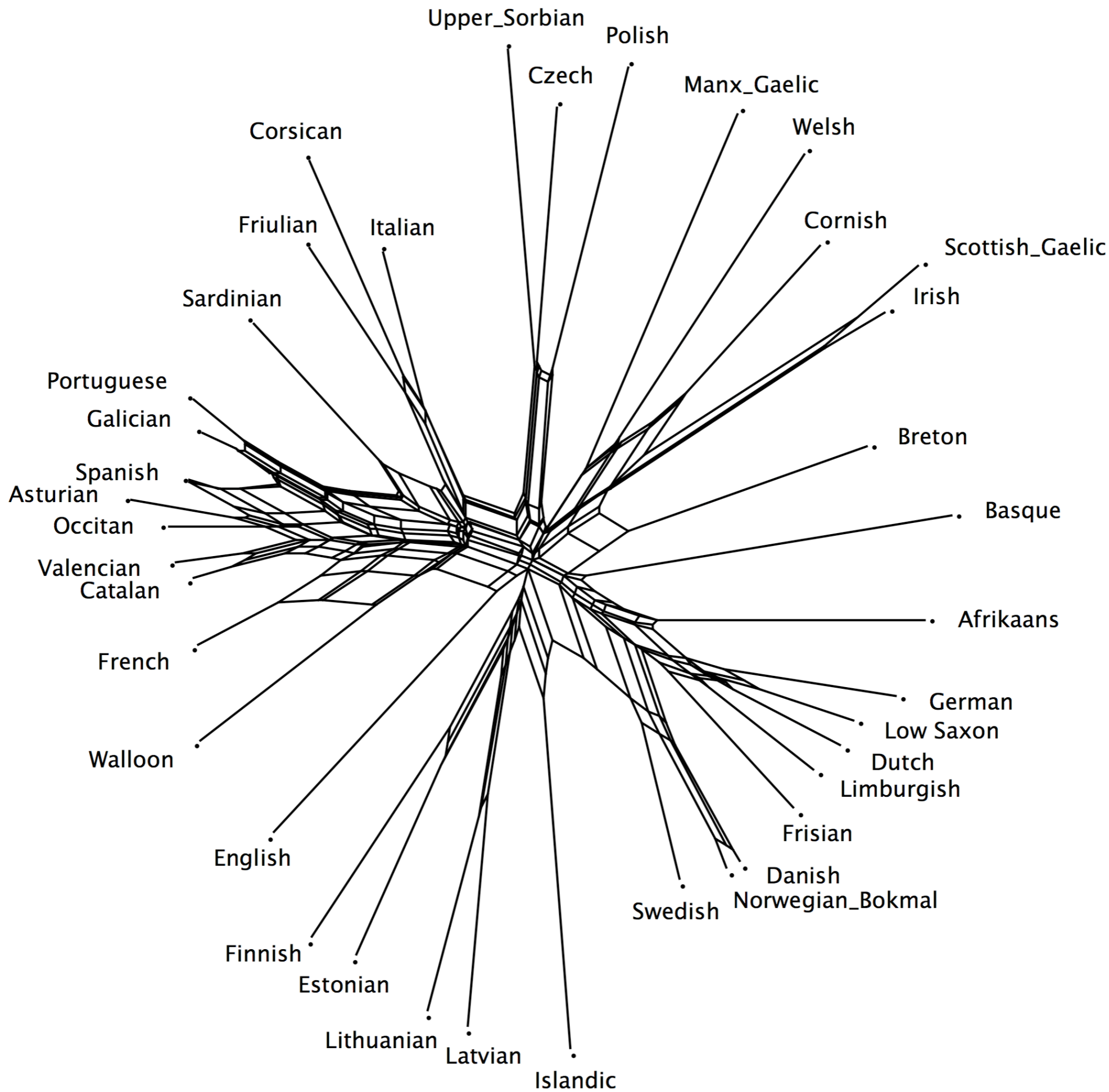
● Ngiyambaa

● Martuthunira



Example 5

Language similarities
by graphotactics
using unannotated web-corpora



The future ?!

- Preparation of data is still very laborious
- Applicability of quantitative methods and visualisations has to be explored for linguistics
- Suitable approaches have to be adapted for linguistic data