

Massively Parallel Text language comparison using methods from vector-space modeling

Michael Cysouw

Philipps-Universität Marburg

cysouw@uni-marburg.de

- There is a deep methodological affinity between the
 - ▶ **Semantic Map** approach from linguistic typology, and
 - ▶ **Vector Space** models from corpus linguistics
- Parallel text offer the possibility to link these two fields

Meaning

- We as humans “know” what meaning is
- But: can we operationalize meaning?
- Logic seemed the way to go for a long time
- Vector algebra seems to be a more practical option

Proposal

- Define the meaning of a linguistic form as the set of all contexts in which it occurs
- Under this strongly extensionalistic definition of meaning, variation in meaning becomes readily measurable

A bit of history ...

Die Bedeutung eines Wortes ist
sein Gebrauch in der Sprache

Ludwig Wittgenstein
Philosophische Untersuchungen, no. 43 (1953)

You shall know a word
by the company it keeps!

John Rupert Firth
Studies in Linguistic Analysis (1957)



... the concept of a translation process in which, in determining meaning for a word, account is taken of the immediate context.

It would hardly be practical to do this by means of a generalized dictionary which contains all possible phrases $2N+1$ words long: for the number of such phrases is horrifying, even to a modern electronic computer

Warren Weaver
Translation (1949)



... if we consider words or morphemes A and B to be more different in meaning than A and C, then we will often find that the distributions of A and B are more different than the distributions of A and C.

**In other words, difference of meaning
correlates with difference of distribution**

Zellig Harris

Methods in Structural Linguistics (1951)



Proposal

- Define the meaning of a linguistic form as the set of all contexts in which it occurs
- Under this strongly extensionalistic definition of meaning, variation in meaning becomes readily measurable

Proposal

- Define the **meaning** of a linguistic form as the set of all contexts in which it occurs
- Under this strongly extensionalistic definition of **meaning**, variation in **meaning** becomes readily measurable

Proposal

- Define the **property M** of a linguistic form as the set of all contexts in which it occurs
- Under this strongly extensionalistic definition of **property M**, variation in **property M** becomes readily measurable

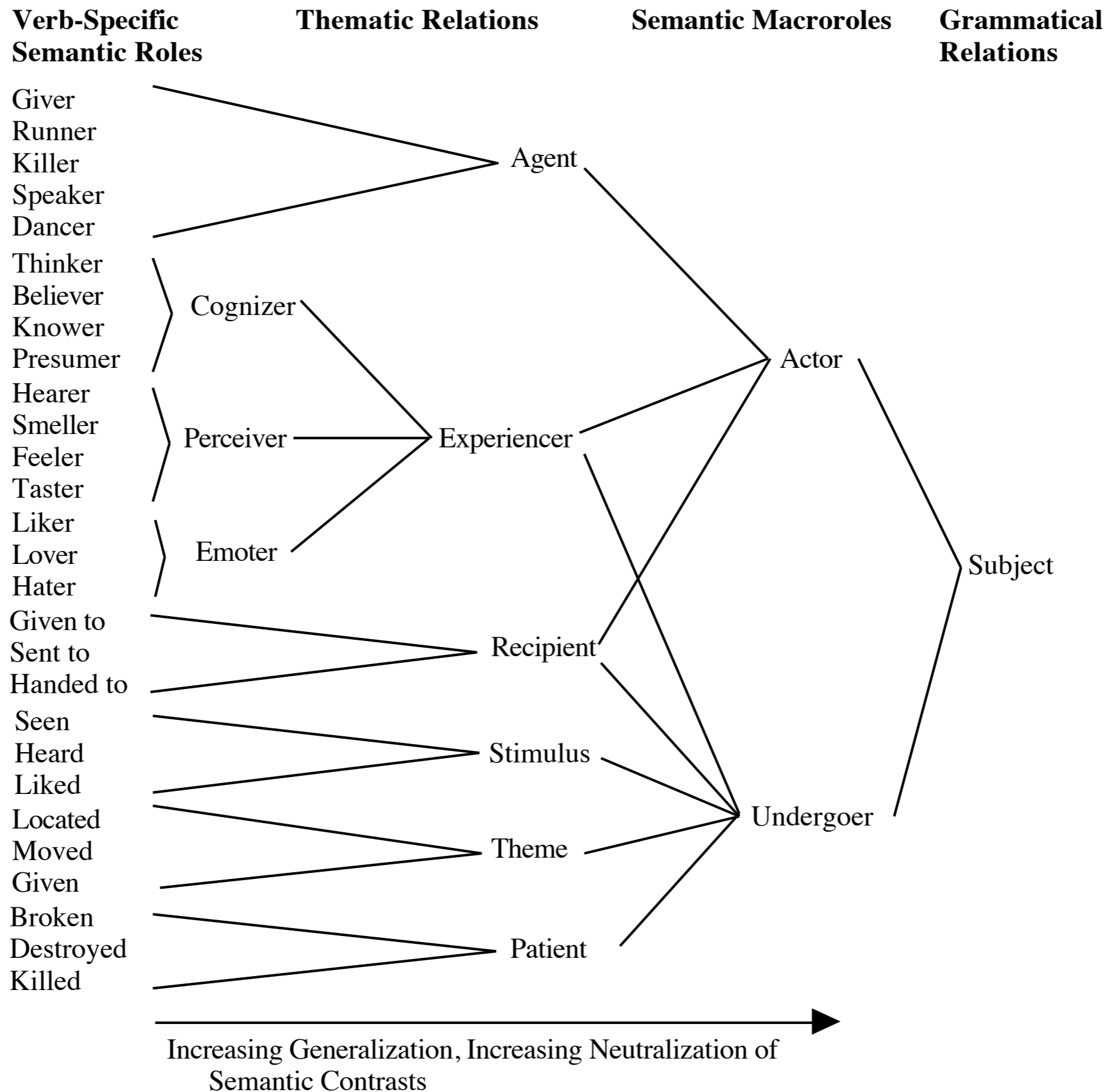
Semantics

- **Exemplar Semantics**
Meaning of linguistic forms is (re-)instantiated and (re-)produced in individual utterances
- **Similarity Semantics**
Meaning of linguistic forms is defined relationally to each other

In practice

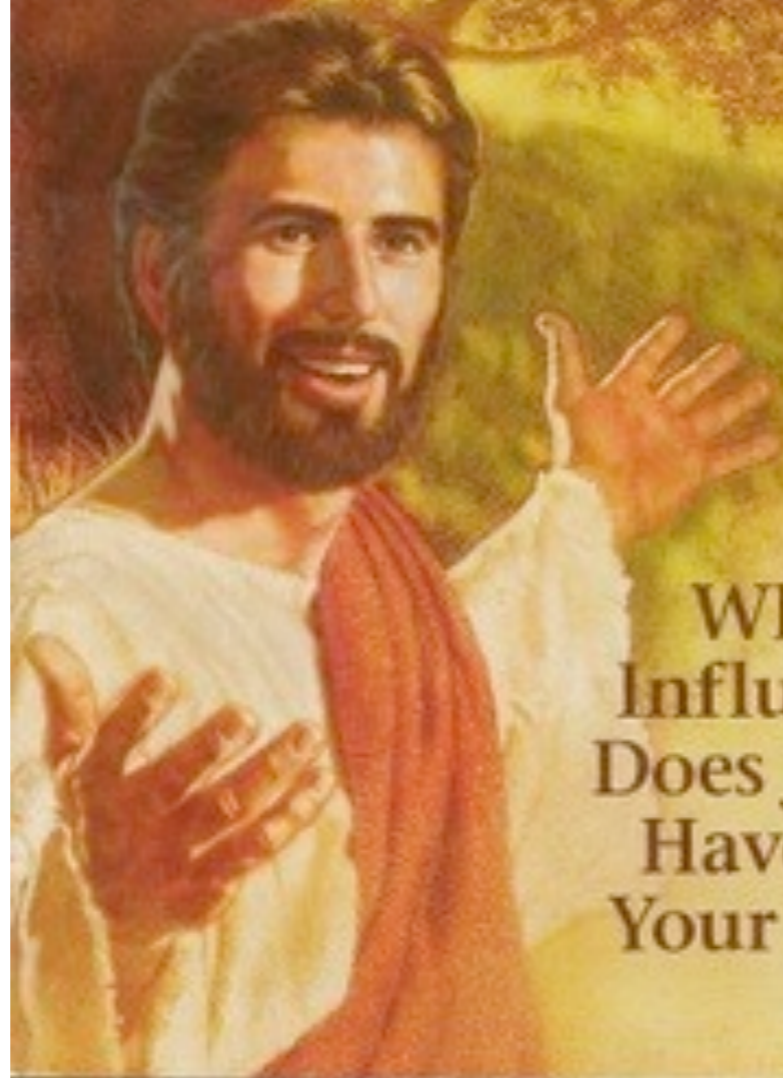
- **Within one language**
investigate co-occurrences of linguistic forms within documents (or sentences, word-windows, etc.)
- **Comparing languages**
investigate the co-occurrence of linguistic forms across languages in parallel texts

Inducing Semantic Roles



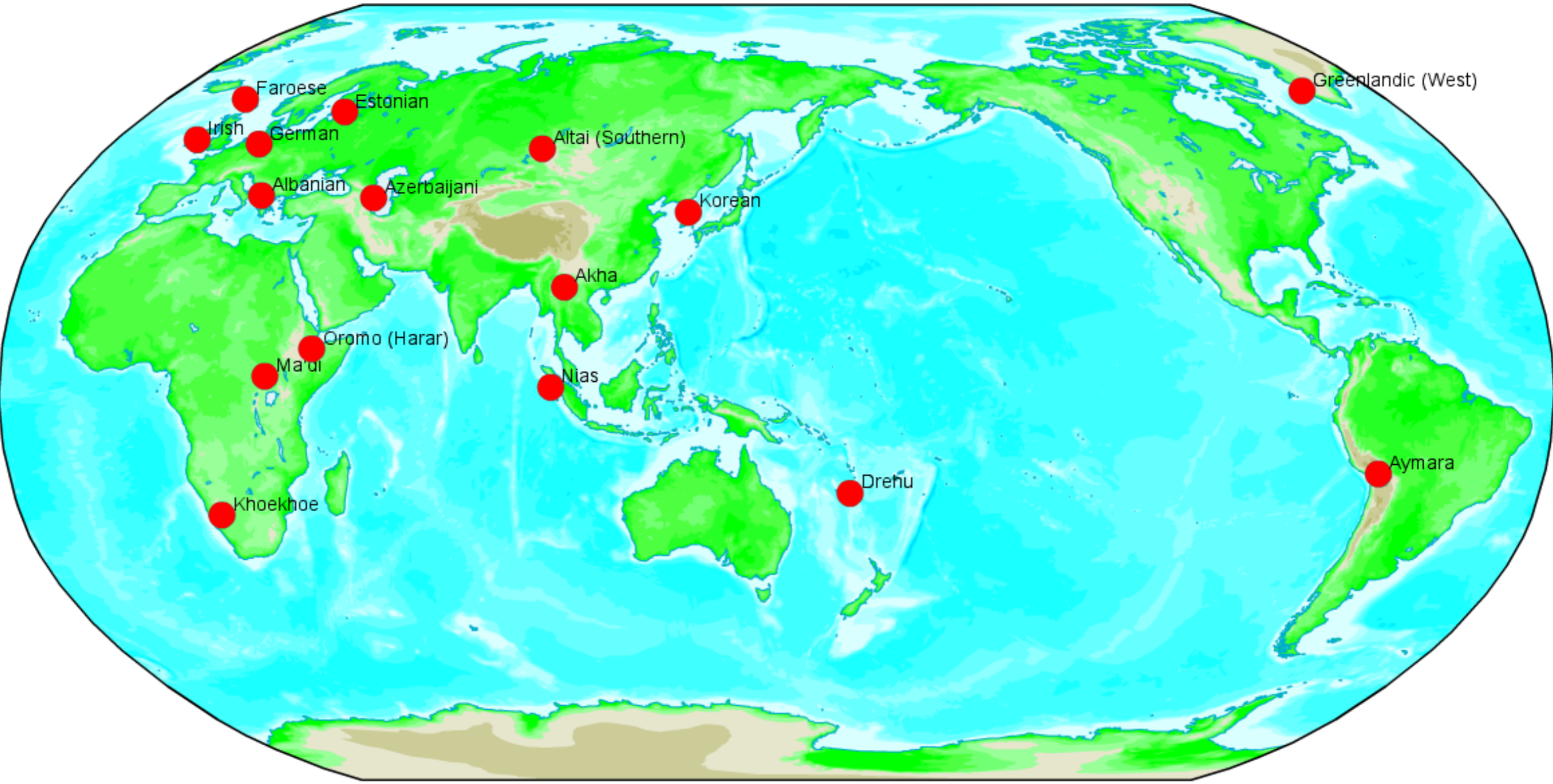
DECEMBER 1, 2011

THE **WATCHTOWER**
ANNOUNCING JEHOVAH'S KINGDOM



What
Influence
Does **Jesus**
Have on
Your Life?

- 1 What important information is contained in the Bible?
- 2 Who is the Bible's author?
- 3 Why should you study the Bible?
- 4 The Bible is a precious gift from God.
- 5 The Bible alone tells us what we must do to please God.
- 6 The Bible was written by some 40 different men over a period of 1,600 years, beginning in 1513 B.C.E.
- 7 So God in heaven, not any human on earth, is the Author of the Bible.
- 8 God made sure that the Bible was accurately copied and preserved.
- 9 More Bibles have been printed than any other book.
- 10 Not everyone will be happy to see you studying the Bible, but do not let that stop you.
- 11 But the Bible tells us that there is only one TRUE God.
- 12 But when the Bible was written, the name Jehovah appeared in it some 7,000 times
- 13 God is a Spirit, says the Bible.
- 14 The Bible reveals Jehovah's personality to us.
- 15 The Bible tells us that he is also merciful, kind, forgiving, generous, and patient.
- 16 We learn about God from creation and from the Bible.
- 17 Another way we can learn about God is by studying the Bible.
- 18 By disobeying God's command, the first man, Adam, committed what the Bible calls sin.
- 19 This is what the Bible refers to as the ransom.
- 20 Some of your loved ones may become very angry because you are studying the Bible.
- 21 What is the Bible's view of separation and of divorce?
- 22 The Bible says that a husband is the head of his family.
- 23 Parents need to spend time with their children and study the Bible with them,
- 24 When marriage mates have problems getting along together, they should try to apply Bible counsel.
- 25 The Bible urges us to show love and to be forgiving.
- 26 But God does not approve of them if they come from false religion or are against Bible teachings.
- 27 The only two birthday celebrations spoken of in the Bible were held by persons who did not worship Jehovah.
- 28 The Bible teaches that only a few people are on the narrow road to life.
- 29 The Bible foretold that after the death of the apostles, ...
- 30 True Christians love one another, respect the Bible, and preach about God's Kingdom.
- 31 Another mark of true religion is that its members have a deep respect for the Bible.
- 32 They try to live by the Bible in their everyday life.
- 33 The Bible is the basis for what is taught.
- 34 By now you have learned many good things from the Bible.



Language	Genus	Family	Macroarea	Alignment
Oromo	Eastern Cushitic	Afro-Asiatic	Africa	Marked nominative
Khoekhoe	Central Khoisan	Khoisan	Africa	Nominative-accusative
Ma'di	Moru-Ma'di	Nilo-Saharan	Africa	No case marking (?)
Albanian	Albanian	Indo-European	Europe	Nominative-accusative
Irish	Celtic	Indo-European	Europe	No case marking (?)
Faroese	Germanic	Indo-European	Europe	Nominative-accusative
Estonian	Finnic	Uralic	Europe	Nominative-accusative
Altai	Turkic	Altaic	Asia	Nominative-accusative
Azerbaijani	Turkic	Altaic	Asia	Nominative-accusative
Korean	Korean	Korean	Asia	Nominative-accusative
Akha	Burmese-Lolo	Sino-Tibetan	Asia	Ergative-absolutive (?)
Drehu	Oceanic	Austronesian	Pacific	Active-inactive
Nias	Sundic	Austronesian	Pacific	Marked absolutive
Greenlandic	Eskimo-Aleut	Eskimo-Aleut	America	Ergative-absolutive
Aymara	Aymaran	Aymaran	America	Marked Nominative (?)

Albanian	Faroese	Estonian	Greenlandic
<i>bibla</i> Nominative	<i>biblian</i> Nominative	<i>piibel</i> Nominative	<i>biibilip</i> Ergative
<i>biblën</i> Accusative	<i>bibliuna</i> Accusative	<i>piiblit</i> Partitive	<i>biibli</i> Absolutive
<i>biblës</i> Genitive/Dative	<i>bibliunnar</i> Genitive	<i>piibli</i> Genitive	<i>biibilmik</i> Instrumental
...	<i>bibliuni</i> Dative	<i>piiblis</i> Inessive	<i>biibilmi</i> Locative
	...	<i>piiblist</i> Elative	...
		...	

Comparing Roles

Context	Albanian	Faroese	Estonian	Greenlandic
1	bibla	bíbliuni	piibel	biibili
2	biblës	bíbliunnar	piibli	biibilimik
3	biblën	bíbliuna	piiblit	biibili
4	bibla	bíblían	piibel	biibili
5	bibla	bíblían	piibel	biibilip
6	bibla	bíbliuna	piibli	biibili
7	biblës	bíbliunnar	piibli	biibilimut
8	bibla	bíblían	piiblit	biibilip
9	bibla	NA	piiblit	biibili
10	biblën	bíbliuna	piiblit	biibilimik
11	bibla	bíblían	piibel	biibilimili
12	bibla	bíblían	piibel	biibilili
13	bibla	bíblían	piibel	biibilimi
14	bibla	bíblían	piibel	biibilimi
15	bibla	bíblían	piibel	biibilimi
16	bibla	bíbliuni	piibli	biibililu
17	biblën	bíbliuna	piiblit	biibilimik
18	bibla	bíblían	piiblis	biibilip
19	bibla	bíblían	piiblis	biibilimi
20	biblën	bíbliuna	piiblit	biibilimik
21	NA	bíblían	piibel	biibilimi
22	bibla	bíbliuni	piibel	biibili
23	biblën	bíbliuna	piiblit	biibilimillu
24	biblike	bíblían	piibli	biibilimi
25	bibla	bíblían	piibel	biibilimi
26	biblës	bíbliunnar	piibli	biibilimi
27	bibla	bíblían	piiblis	biibilimi
28	bibla	bíblían	piibel	biibilimi
29	bibla	bíblían	piibel	biibilimi
30	biblën	bíbliuna	piiblist	biibilimik
31	biblën	bíbliuni	piibli	biibilimik
32	biblës	bíbliuni	piibli	biibili
33	bibla	bíbliuna	piibel	biibilimik
34	bibla	bíbliuni	piiblist	biibilimeersunik

0.00	0.67	0.50	0.75	0.67	0.75	0.69	0.75	0.73	0.56	0.44	0.63	0.53	0.75	0.53	0.56	0.63	0.50	0.63	0.63	0.60	0.50	0.56	0.69	0.75	0.57	0.44	0.69	0.56	0.63	0.63	0.50	0.75	0.63
0.67	0.00	0.73	0.80	0.73	0.67	0.20	0.73	0.64	0.60	0.80	0.67	0.93	0.73	0.57	0.80	0.73	0.73	0.73	0.67	0.50	0.67	0.73	0.40	0.60	0.15	0.53	0.67	0.67	0.60	0.60	0.38	0.73	0.80
0.50	0.73	0.00	0.88	0.80	0.69	0.81	0.75	0.80	0.25	0.75	0.75	0.73	0.88	0.80	0.69	0.13	0.69	0.81	0.19	0.73	0.81	0.06	0.75	0.88	0.64	0.69	0.88	0.88	0.44	0.63	0.64	0.75	0.69
0.75	0.80	0.88	0.00	0.40	0.31	0.81	0.56	0.53	0.88	0.44	0.44	0.40	0.13	0.47	0.75	0.94	0.63	0.63	0.94	0.53	0.44	0.94	0.88	0.25	0.93	0.81	0.31	0.44	1.00	1.00	0.79	0.50	0.94
0.67	0.73	0.80	0.40	0.00	0.60	0.80	0.60	0.79	0.80	0.47	0.47	0.57	0.47	0.50	0.73	0.87	0.53	0.60	0.87	0.57	0.53	0.80	0.80	0.47	0.77	0.67	0.47	0.47	0.93	1.00	0.85	0.53	0.87
0.75	0.67	0.69	0.31	0.60	0.00	0.69	0.63	0.53	0.75	0.63	0.63	0.67	0.44	0.67	0.81	0.75	0.75	0.75	0.88	0.60	0.56	0.75	0.81	0.56	0.79	0.81	0.63	0.63	0.81	0.81	0.71	0.63	0.94
0.69	0.20	0.81	0.81	0.80	0.69	0.00	0.63	0.60	0.75	0.81	0.75	0.93	0.75	0.67	0.81	0.81	0.81	0.81	0.81	0.60	0.75	0.81	0.50	0.75	0.29	0.63	0.69	0.69	0.69	0.69	0.50	0.81	0.81
0.75	0.73	0.75	0.56	0.60	0.63	0.63	0.00	0.47	0.88	0.56	0.50	0.67	0.63	0.53	0.75	0.75	0.63	0.69	0.81	0.80	0.69	0.75	0.75	0.63	0.79	0.75	0.56	0.56	0.75	0.75	0.86	0.81	0.75
0.73	0.64	0.80	0.53	0.79	0.53	0.60	0.47	0.00	0.73	0.80	0.73	0.86	0.67	0.64	0.80	0.87	0.87	0.87	0.80	0.86	0.67	0.87	0.73	0.67	0.69	0.73	0.67	0.73	0.80	0.80	0.62	0.80	0.80
0.56	0.60	0.25	0.88	0.80	0.75	0.75	0.88	0.73	0.00	0.81	0.75	0.80	0.81	0.80	0.69	0.25	0.69	0.81	0.19	0.67	0.88	0.25	0.69	0.81	0.57	0.63	0.81	0.88	0.44	0.63	0.57	0.63	0.81
0.44	0.80	0.75	0.44	0.47	0.63	0.81	0.56	0.80	0.81	0.00	0.56	0.27	0.38	0.27	0.69	0.81	0.44	0.44	0.81	0.53	0.31	0.75	0.81	0.38	0.86	0.63	0.31	0.19	0.88	0.88	0.86	0.56	0.81
0.63	0.67	0.75	0.44	0.47	0.63	0.75	0.50	0.73	0.75	0.56	0.00	0.53	0.50	0.47	0.69	0.81	0.44	0.56	0.81	0.33	0.63	0.75	0.63	0.50	0.64	0.63	0.56	0.56	0.88	0.88	0.71	0.69	0.81
0.53	0.93	0.73	0.40	0.57	0.67	0.93	0.67	0.86	0.80	0.27	0.53	0.00	0.27	0.36	0.60	0.73	0.40	0.33	0.80	0.43	0.40	0.73	0.73	0.33	0.77	0.60	0.33	0.33	0.80	0.80	0.85	0.60	0.73
0.75	0.73	0.88	0.13	0.47	0.44	0.75	0.63	0.67	0.81	0.38	0.50	0.27	0.00	0.33	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.13	0.79	0.69	0.19	0.31	0.94	0.94	0.79	0.44	0.88
0.53	0.57	0.80	0.47	0.50	0.67	0.67	0.53	0.64	0.80	0.27	0.47	0.36	0.33	0.00	0.73	0.87	0.33	0.27	0.80	0.50	0.20	0.80	0.53	0.20	0.57	0.33	0.20	0.07	0.80	0.80	0.64	0.60	0.73
0.56	0.80	0.69	0.75	0.73	0.81	0.81	0.75	0.80	0.69	0.69	0.69	0.60	0.69	0.73	0.00	0.63	0.63	0.75	0.69	0.80	0.75	0.69	0.75	0.75	0.71	0.75	0.75	0.81	0.69	0.56	0.57	0.63	0.44
0.63	0.73	0.13	0.94	0.87	0.75	0.81	0.75	0.87	0.25	0.81	0.81	0.73	0.88	0.87	0.63	0.00	0.75	0.88	0.19	0.80	0.94	0.13	0.81	0.94	0.71	0.75	0.94	0.94	0.31	0.50	0.79	0.63	0.63
0.50	0.73	0.69	0.63	0.53	0.75	0.81	0.63	0.87	0.69	0.44	0.44	0.40	0.56	0.33	0.63	0.75	0.00	0.25	0.75	0.53	0.44	0.69	0.63	0.56	0.64	0.38	0.56	0.44	0.81	0.81	0.64	0.69	0.69
0.63	0.73	0.81	0.63	0.60	0.75	0.81	0.69	0.87	0.81	0.44	0.56	0.33	0.50	0.27	0.75	0.88	0.25	0.00	0.81	0.47	0.31	0.81	0.56	0.38	0.57	0.31	0.38	0.25	0.81	0.81	0.64	0.69	0.75
0.63	0.67	0.19	0.94	0.87	0.88	0.81	0.81	0.80	0.19	0.81	0.81	0.80	0.88	0.80	0.69	0.19	0.75	0.81	0.00	0.80	0.88	0.19	0.75	0.81	0.64	0.69	0.81	0.88	0.44	0.63	0.64	0.69	0.69
0.60	0.50	0.73	0.53	0.57	0.60	0.60	0.80	0.86	0.67	0.53	0.33	0.43	0.40	0.50	0.80	0.80	0.53	0.47	0.80	0.00	0.60	0.73	0.40	0.40	0.46	0.53	0.47	0.47	0.80	0.80	0.62	0.67	0.87
0.50	0.67	0.81	0.44	0.53	0.56	0.75	0.69	0.67	0.88	0.31	0.63	0.40	0.44	0.20	0.75	0.94	0.44	0.31	0.88	0.60	0.00	0.88	0.75	0.31	0.71	0.56	0.31	0.19	0.88	0.81	0.57	0.56	0.75
0.56	0.73	0.06	0.94	0.80	0.75	0.81	0.75	0.87	0.25	0.75	0.75	0.73	0.88	0.80	0.69	0.13	0.69	0.81	0.19	0.73	0.88	0.00	0.75	0.88	0.64	0.69	0.88	0.88	0.44	0.63	0.71	0.75	0.69
0.69	0.40	0.75	0.88	0.80	0.81	0.50	0.75	0.73	0.69	0.81	0.63	0.73	0.75	0.53	0.75	0.81	0.63	0.56	0.75	0.40	0.75	0.75	0.00	0.63	0.21	0.38	0.63	0.63	0.69	0.63	0.43	0.81	0.75
0.75	0.60	0.88	0.25	0.47	0.56	0.75	0.63	0.67	0.81	0.38	0.50	0.33	0.13	0.20	0.75	0.94	0.56	0.38	0.81	0.40	0.31	0.88	0.63	0.00	0.64	0.56	0.06	0.19	0.88	0.88	0.64	0.50	0.81
0.57	0.15	0.64	0.93	0.77	0.79	0.29	0.79	0.69	0.57	0.86	0.64	0.77	0.79	0.57	0.71	0.71	0.64	0.57	0.64	0.46	0.71	0.64	0.21	0.64	0.00	0.36	0.64	0.64	0.57	0.57	0.31	0.86	0.71
0.44	0.53	0.69	0.81	0.67	0.81	0.63	0.75	0.73	0.63	0.63	0.63	0.60	0.69	0.33	0.75	0.75	0.38	0.31	0.69	0.53	0.56	0.69	0.38	0.56	0.36	0.00	0.56	0.44	0.63	0.63	0.43	0.81	0.69
0.69	0.67	0.88	0.31	0.47	0.63	0.69	0.56	0.67	0.81	0.31	0.56	0.33	0.19	0.20	0.75	0.94	0.56	0.38	0.81	0.47	0.31	0.88	0.63	0.06	0.64	0.56	0.00	0.13	0.88	0.88	0.64	0.50	0.81
0.56	0.67	0.88	0.44	0.47	0.63	0.69	0.56	0.73	0.88	0.19	0.56	0.33	0.31	0.07	0.81	0.94	0.44	0.25	0.88	0.47	0.19	0.88	0.63	0.19	0.64	0.44	0.13	0.00	0.88	0.88	0.71	0.56	0.81
0.63	0.60	0.44	1.00	0.93	0.81	0.69	0.75	0.80	0.44	0.88	0.88	0.80	0.94	0.80	0.69	0.31	0.81	0.81	0.44	0.80	0.88	0.44	0.69	0.88	0.57	0.63	0.88	0.88	0.00	0.31	0.57	0.75	0.56
0.63	0.60	0.63	1.00	1.00	0.81	0.69	0.75	0.80	0.63	0.88	0.88	0.80	0.94	0.80	0.56	0.50	0.81	0.81	0.63	0.80	0.81	0.63	0.63	0.88	0.57	0.63	0.88	0.88	0.31	0.00	0.36	0.81	0.56
0.50	0.38	0.64	0.79	0.85	0.71	0.50	0.86	0.62	0.57	0.86	0.71	0.85	0.79	0.64	0.57	0.79	0.64	0.64	0.64	0.62	0.57	0.71	0.43	0.64	0.31	0.43	0.64	0.71	0.57	0.36	0.00	0.79	0.64
0.75	0.73	0.75	0.50	0.53	0.63	0.81	0.81	0.80	0.63	0.56	0.69	0.60	0.44	0.60	0.63	0.63	0.69	0.69	0.69	0.67	0.56	0.75	0.81	0.50	0.86	0.81	0.50	0.56	0.75	0.81	0.79	0.00	0.81
0.63	0.80	0.69	0.94	0.87	0.94	0.81	0.75	0.80	0.81	0.81	0.81	0.73	0.88	0.73	0.44	0.63	0.69	0.75	0.69	0.87	0.75	0.69	0.75	0.81	0.71	0.69	0.81	0.81	0.56	0.56	0.64	0.81	0.00

0.00	0.67	0.50	0.75	0.67	0.75	0.69	0.75	0.73	0.56	0.44	0.63	0.53	0.75	0.53	0.56	0.63	0.50	0.63	0.63	0.60	0.50	0.56	0.69	0.75	0.57	0.44	0.69	0.56	0.63	0.63	0.50	0.75	0.63	
0.67	0.00	0.73	0.80	0.73	0.67	0.20	0.73	0.64	0.60	0.80	0.67	0.93	0.73	0.57	0.80	0.73	0.73	0.73	0.67	0.50	0.67	0.73	0.40	0.60	0.15	0.53	0.67	0.67	0.60	0.60	0.38	0.73	0.80	
0.50	0.73	0.00	0.88	0.80	0.69	0.81	0.75	0.80	0.25	0.75	0.75	0.73	0.88	0.80	0.69	0.13	0.69	0.81	0.19	0.73	0.81	0.06	0.75	0.88	0.64	0.69	0.88	0.88	0.44	0.63	0.64	0.75	0.69	
0.75	0.80	0.88	0.00	0.40	0.31	0.81	0.56	0.53	0.88	0.44	0.44	0.40	0.13	0.47	0.75	0.94	0.63	0.63	0.94	0.53	0.44	0.94	0.88	0.25	0.93	0.81	0.31	0.44	1.00	1.00	0.79	0.50	0.94	
0.67	0.73	0.80	0.40	0.00	0.60	0.80	0.60	0.79	0.80	0.47	0.47	0.57	0.47	0.50	0.73	0.87	0.53	0.60	0.87	0.57	0.53	0.80	0.80	0.47	0.77	0.67	0.47	0.47	0.93	1.00	0.85	0.53	0.87	
0.75	0.67	0.69	0.31	0.60	0.00			0.53	0.75	0.63	0.63	0.67	0.44	0.67	0.81	0.75	0.75	0.75	0.88	0.60	0.56	0.75	0.81	0.56	0.79	0.81	0.63	0.63	0.81	0.81	0.71	0.63	0.94	
0.69	0.20	0.81	0.81	0.80				0.60	0.75	0.81	0.75	0.93	0.75	0.67	0.81	0.81	0.81	0.81	0.81	0.81	0.60	0.75	0.81	0.50	0.75	0.29	0.63	0.69	0.69	0.69	0.69	0.50	0.81	0.81
0.75	0.73	0.75	0.56	0.60				0.47	0.88	0.56	0.50	0.67	0.63	0.53	0.75	0.75	0.63	0.69	0.81	0.80	0.69	0.75	0.75	0.63	0.79	0.75	0.56	0.56	0.75	0.75	0.86	0.81	0.75	
0.73	0.64	0.80	0.53	0.79	0.53	0.60	0.47	0.00	0.73	0.80	0.73	0.86	0.67	0.64	0.80	0.87	0.87	0.87	0.80	0.86	0.67	0.87	0.73	0.67	0.69	0.73	0.67	0.73	0.80	0.80	0.62	0.80	0.80	
0.56	0.60	0.25	0.88	0.80	0.75	0.75	0.88	0.73	0.00	0.81	0.75	0.80	0.81	0.80	0.69	0.25	0.69	0.81	0.19	0.67	0.88	0.25	0.69	0.81	0.57	0.63	0.81	0.88	0.44	0.63	0.57	0.63	0.81	
0.44	0.80	0.75	0.44	0.47	0.63	0.81	0.56	0.80	0.81	0.00	0.56	0.27	0.38	0.27	0.69	0.81	0.44	0.44	0.81	0.53	0.31	0.75	0.81	0.38	0.86	0.63	0.31	0.19	0.88	0.88	0.86	0.56	0.81	
0.63	0.67	0.75	0.44	0.47	0.63	0.75	0.50	0.73	0.75	0.56	0.00	0.53	0.50	0.47	0.69	0.81	0.44	0.56	0.81	0.33	0.63	0.75	0.63	0.50	0.64	0.63	0.56	0.56	0.88	0.88	0.71	0.69	0.81	
0.53	0.93	0.73	0.40	0.57	0.67	0.93	0.67	0.86	0.80	0.27	0.53	0.00	0.27	0.36	0.60	0.73	0.40	0.33	0.80	0.43	0.40	0.73	0.73	0.33	0.77	0.60	0.33	0.33	0.80	0.80	0.85	0.60	0.73	
0.75	0.73	0.88	0.13	0.47	0.44	0.75	0.63	0.67	0.81	0.38	0.50	0.27	0.00	0.33	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.13	0.79	0.69	0.19	0.31	0.94	0.94	0.79	0.44	0.88	
0.53	0.57	0.80	0.47	0.50	0.67	0.67	0.53	0.64	0.80	0.27	0.47	0.36	0.33	0.00	0.73	0.87	0.33	0.27	0.80	0.50	0.20	0.80	0.53	0.20	0.57	0.33	0.20	0.07	0.80	0.80	0.64	0.60	0.73	
0.56	0.80	0.69	0.75	0.73	0.81	0.81	0.75	0.80	0.69	0.69	0.69	0.60	0.69	0.73	0.00	0.63	0.63	0.75	0.69	0.80	0.75	0.69	0.75	0.75	0.71	0.75	0.75	0.81	0.69	0.56	0.57	0.63	0.44	
0.63	0.73	0.13	0.94	0.87	0.75	0.81	0.75	0.87	0.25	0.81	0.81	0.73	0.88	0.87	0.63	0.00	0.75	0.88	0.19	0.80	0.94	0.13	0.81	0.94	0.71	0.75	0.94	0.94	0.31	0.50	0.79	0.63	0.63	
0.50	0.73	0.69	0.63	0.53	0.75	0.81	0.63	0.87	0.69	0.44	0.44	0.40	0.56	0.33	0.63	0.75	0.00	0.25	0.75	0.53	0.44	0.69	0.63	0.56	0.64	0.38	0.56	0.44	0.81	0.81	0.64	0.69	0.69	
0.63	0.73	0.81	0.63	0.60	0.75	0.81	0.69	0.87	0.81	0.44	0.56	0.33	0.50	There is no difference between context 5 and itself in any language														0.81	0.81	0.64	0.69	0.75		
0.63	0.67	0.19	0.94	0.87	0.88	0.81	0.81	0.80	0.19	0.81	0.81	0.80	0.88															0.44	0.63	0.64	0.69	0.69		
0.60	0.50	0.73	0.53	0.57	0.60	0.60	0.80	0.86	0.67	0.53	0.33	0.43	0.40															0.80	0.80	0.62	0.67	0.87		
0.50	0.67	0.81	0.44	0.53	0.56	0.75	0.69	0.67	0.88	0.31	0.63	0.40	0.44															0.88	0.81	0.57	0.56	0.75		
0.56	0.73	0.06	0.94	0.80	0.75	0.81	0.75	0.87	0.25	0.75	0.75	0.73	0.88															0.44	0.63	0.71	0.75	0.69		
0.69	0.40	0.75	0.88	0.80	0.81	0.50	0.75	0.73	0.69	0.81	0.63	0.73	0.75															0.69	0.63	0.43	0.81	0.75		
0.75	0.60	0.88	0.25	0.47	0.56	0.75	0.63	0.67	0.81	0.38	0.50	0.33	0.13															0.88	0.88	0.64	0.50	0.81		
0.57	0.15	0.64	0.93	0.77	0.79	0.29	0.79	0.69	0.57	0.86	0.64	0.77	0.79	0.57	0.71	0.71	0.64	0.57	0.64	0.46	0.71	0.64	0.21	0.64	0.00	0.36	0.64	0.64	0.57	0.57	0.31	0.86	0.71	
0.44	0.53	0.69	0.81	0.67	0.81	0.63	0.75	0.73	0.63	0.63	0.63	0.60	0.69	0.33	0.75	0.75	0.38	0.31	0.69	0.53	0.56	0.69	0.38	0.56	0.36	0.00	0.56	0.44	0.63	0.63	0.43	0.81	0.69	
0.69	0.67	0.88	0.31	0.47	0.63	0.69	0.56	0.67	0.81	0.31	0.56	0.33	0.19	0.20	0.75	0.94	0.56	0.38	0.81	0.47	0.31	0.88	0.63	0.06	0.64	0.56	0.00	0.13	0.88	0.88	0.64	0.50	0.81	
0.56	0.67	0.88	0.44	0.47	0.63	0.69	0.56	0.73	0.88	0.19	0.56	0.33	0.31	0.07	0.81	0.94	0.44	0.25	0.88	0.47	0.19	0.88	0.63	0.19	0.64	0.44	0.13	0.00	0.88	0.88	0.71	0.56	0.81	
0.63	0.60	0.44	1.00	0.93	0.81	0.69	0.75	0.80	0.44	0.88	0.88	0.80	0.94	0.80	0.69	0.31	0.81	0.81	0.44	0.80	0.88	0.44	0.69	0.88	0.57	0.63	0.88	0.88	0.00	0.31	0.57	0.75	0.56	
0.63	0.60	0.63	1.00	1.00	0.81	0.69	0.75	0.80	0.63	0.88	0.88	0.80	0.94	0.80	0.56	0.50	0.81	0.81	0.63	0.80	0.81	0.63	0.63	0.88	0.57	0.63	0.88	0.88	0.31	0.00	0.36	0.81	0.56	
0.50	0.38	0.64	0.79	0.85	0.71	0.50	0.86	0.62	0.57	0.86	0.71	0.85	0.79	0.64	0.57	0.79	0.64	0.64	0.64	0.62	0.57	0.71	0.43	0.64	0.31	0.43	0.64	0.71	0.57	0.36	0.00	0.79	0.64	
0.75	0.73	0.75	0.50	0.53	0.63	0.81	0.81	0.80	0.63	0.56	0.69	0.60	0.44	0.60	0.63	0.63	0.69	0.69	0.69	0.67	0.56	0.75	0.81	0.50	0.86	0.81	0.50	0.56	0.75	0.81	0.79	0.00	0.81	
0.63	0.80	0.69	0.94	0.87	0.94	0.81	0.75	0.80	0.81	0.81	0.81	0.73	0.88	0.73	0.44	0.63	0.69	0.75	0.69	0.87	0.75	0.69	0.75	0.81	0.71	0.69	0.81	0.81	0.56	0.56	0.64	0.81	0.00	

0.00	0.67	0.50	0.75	0.67	0.75	0.69	0.75	0.73	0.56	0.44	0.63	0.53	0.75	0.53	0.56	0.63	0.50	0.63	0.63	0.60	0.50	0.56	0.69	0.75	0.57	0.44	0.69	0.56	0.63	0.63	0.50	0.75	0.63
0.67	0.00	0.73	0.80	0.73	0.67	0.20	0.73	0.64	0.60	0.80	0.67	0.93	0.73	0.57	0.80	0.73	0.73	0.73	0.67	0.50	0.67	0.73	0.40	0.60	0.15	0.53	0.67	0.67	0.60	0.60	0.38	0.73	0.80
0.50	0.73	0.00	0.88	0.80	0.69	0.81	0.75	0.80	0.25	0.75	0.75	0.13			0.69	0.13	0.69	0.81	0.19	0.73	0.81	0.06	0.75	0.88	0.64	0.69	0.88	0.88	0.44	0.63	0.64	0.75	0.69
0.75	0.80	0.88	0.00	0.40	0.31	0.81	0.56	0.53	0.88	0.44	0.44	0.13			0.75	0.94	0.63	0.63	0.94	0.53	0.44	0.94	0.88	0.25	0.93	0.81	0.31	0.44	1.00	1.00	0.79	0.50	0.94
0.67	0.73	0.80	0.40	0.00	0.60	0.80	0.60	0.79	0.80	0.47	0.47	0.13			0.73	0.87	0.53	0.60	0.87	0.57	0.53	0.80	0.80	0.47	0.77	0.67	0.47	0.47	0.93	1.00	0.85	0.53	0.87
0.75	0.67	0.69	0.31	0.60	0.00	0.69	0.63	0.53	0.75	0.63	0.63	0.67	0.44	0.67	0.81	0.75	0.75	0.75	0.88	0.60	0.56	0.75	0.81	0.56	0.79	0.81	0.63	0.63	0.81	0.81	0.71	0.63	0.94
0.69	0.20	0.81	0.81	0.80	0.69	0.00	0.63	0.60	0.75	0.81	0.75	0.93	0.75	0.67	0.81	0.81	0.81	0.81	0.81	0.60	0.75	0.81	0.50	0.75	0.29	0.63	0.69	0.69	0.69	0.69	0.50	0.81	0.81
0.75	0.73	0.75	0.56	0.60	0.63	0.63	0.00	0.47	0.88	0.56	0.50	0.67	0.63	0.53	0.75	0.75	0.63	0.69	0.81	0.80	0.69	0.75	0.75	0.63	0.79	0.75	0.56	0.56	0.75	0.75	0.86	0.81	0.75
0.73	0.64	0.80	0.53	0.79	0.53	0.60	0.47	0.00	0.73	0.80	0.73	0.86	0.67	0.64	0.80	0.87	0.87	0.87	0.80	0.86	0.67	0.87	0.73	0.67	0.69	0.73	0.67	0.73	0.80	0.80	0.62	0.80	0.80
0.56	0.60	0.25	0.88	0.80	0.75	0.75	0.88	0.73	0.00	0.81	0.75	0.80	0.81	0.80	0.69	0.25	0.69	0.81	0.19	0.67	0.88	0.25	0.69	0.81	0.57	0.63	0.81	0.88	0.44	0.63	0.57	0.63	0.81
0.44	0.80	0.75	0.44	0.47	0.63	0.81	0.56	0.80	0.81	0.00	0.56	0.27	0.38	0.27	0.69	0.81	0.44	0.44	0.81	0.53	0.31	0.75	0.81	0.38	0.86	0.63	0.31	0.19	0.88	0.88	0.86	0.56	0.81
0.63	0.67	On average, there is a small difference between the role in context 4 and 14														0.81	0.33	0.63	0.75	0.63	0.50	0.64	0.63	0.56	0.56	0.88	0.88	0.71	0.69	0.81			
0.53	0.93	On average, there is a small difference between the role in context 4 and 14														0.80	0.43	0.40	0.73	0.73	0.33	0.77	0.60	0.33	0.33	0.80	0.80	0.85	0.60	0.73			
0.75	0.73	On average, there is a small difference between the role in context 4 and 14														0.88	0.40	0.44	0.88	0.75	0.13	0.79	0.69	0.19	0.31	0.94	0.94	0.79	0.44	0.88			
0.53	0.57	On average, there is a small difference between the role in context 4 and 14														0.80	0.50	0.20	0.80	0.53	0.20	0.57	0.33	0.20	0.07	0.80	0.80	0.64	0.60	0.73			
0.56	0.80	On average, there is a small difference between the role in context 4 and 14														0.69	0.80	0.75	0.69	0.75	0.75	0.71	0.75	0.75	0.81	0.69	0.56	0.57	0.63	0.44			
0.63	0.73	On average, there is a small difference between the role in context 4 and 14														0.19	0.80	0.94	0.13	0.81	0.94	0.71	0.75	0.94	0.94	0.31	0.50	0.79	0.63	0.63			
0.50	0.73	On average, there is a small difference between the role in context 4 and 14														0.75	0.53	0.44	0.69	0.63	0.56	0.64	0.38	0.56	0.44	0.81	0.81	0.64	0.69	0.69			
0.63	0.73	On average, there is a small difference between the role in context 4 and 14														0.81	0.47	0.31	0.81	0.56	0.38	0.57	0.31	0.38	0.25	0.81	0.81	0.64	0.69	0.75			
0.63	0.67	On average, there is a small difference between the role in context 4 and 14														0.00	0.80	0.88	0.19	0.75	0.81	0.64	0.69	0.81	0.88	0.44	0.63	0.64	0.69	0.69			
0.60	0.50	On average, there is a small difference between the role in context 4 and 14														0.80	0.00	0.60	0.73	0.40	0.40	0.46	0.53	0.47	0.47	0.80	0.80	0.62	0.67	0.87			
0.50	0.67	On average, there is a small difference between the role in context 4 and 14														0.88	0.60	0.00	0.88	0.75	0.31	0.71	0.56	0.31	0.19	0.88	0.81	0.57	0.56	0.75			
0.56	0.73	On average, there is a small difference between the role in context 4 and 14														0.19	0.73	0.88	0.00	0.75	0.88	0.64	0.69	0.88	0.88	0.44	0.63	0.71	0.75	0.69			
0.69	0.40	0.75	0.88	0.80	0.81	0.50	0.75	0.73	0.69	0.81	0.63	0.73	0.75	0.53	0.75	0.81	0.63	0.56	0.75	0.40	0.75	0.75	0.00	0.63	0.21	0.38	0.63	0.63	0.69	0.63	0.43	0.81	0.75
0.75	0.60	0.88	0.25	0.47	0.56	0.75	0.63	0.67	0.81	0.38	0.50	0.33	0.13	0.20	0.75	0.94	0.56	0.38	0.81	0.40	0.31	0.88	0.63	0.00	0.64	0.56	0.06	0.19	0.88	0.88	0.64	0.50	0.81
0.57	0.15	0.64	0.93	0.77	0.79	0.29	0.79	0.69	0.57	0.86	0.64	0.77	0.79	0.57	0.71	0.71	0.64	0.57	0.64	0.46	0.71	0.64	0.21	0.64	0.00	0.36	0.64	0.64	0.57	0.57	0.31	0.86	0.71
0.44	0.53	0.69	0.81	0.67	0.81	0.63	0.75	0.73	0.63	0.63	0.63	0.60	0.69	0.33	0.75	0.75	0.38	0.31	0.69	0.53	0.56	0.69	0.38	0.56	0.36	0.00	0.56	0.44	0.63	0.63	0.43	0.81	0.69
0.69	0.67	0.88	0.31	0.47	0.63	0.69	0.56	0.67	0.81	0.31	0.56	0.33	0.19	0.20	0.75	0.94	0.56	0.38	0.81	0.47	0.31	0.88	0.63	0.06	0.64	0.56	0.00	0.13	0.88	0.88	0.64	0.50	0.81
0.56	0.67	0.88	0.44	0.47	0.63	0.69	0.56	0.73	0.88	0.19	0.56	0.33	0.31	0.07	0.81	0.94	0.44	0.25	0.88	0.47	0.19	0.88	0.63	0.19	0.64	0.44	0.13	0.00	0.88	0.88	0.71	0.56	0.81
0.63	0.60	0.44	1.00	0.93	0.81	0.69	0.75	0.80	0.44	0.88	0.88	0.80	0.94	0.80	0.69	0.31	0.81	0.81	0.44	0.80	0.88	0.44	0.69	0.88	0.57	0.63	0.88	0.88	0.00	0.31	0.57	0.75	0.56
0.63	0.60	0.63	1.00	1.00	0.81	0.69	0.75	0.80	0.63	0.88	0.88	0.80	0.94	0.80	0.56	0.50	0.81	0.81	0.63	0.80	0.81	0.63	0.63	0.88	0.57	0.63	0.88	0.88	0.31	0.00	0.36	0.81	0.56
0.50	0.38	0.64	0.79	0.85	0.71	0.50	0.86	0.62	0.57	0.86	0.71	0.85	0.79	0.64	0.57	0.79	0.64	0.64	0.64	0.62	0.57	0.71	0.43	0.64	0.31	0.43	0.64	0.71	0.57	0.36	0.00	0.79	0.64
0.75	0.73	0.75	0.50	0.53	0.63	0.81	0.81	0.80	0.63	0.56	0.69	0.60	0.44	0.60	0.63	0.63	0.69	0.69	0.69	0.67	0.56	0.75	0.81	0.50	0.86	0.81	0.50	0.56	0.75	0.81	0.79	0.00	0.81
0.63	0.80	0.69	0.94	0.87	0.94	0.81	0.75	0.80	0.81	0.81	0.81	0.73	0.88	0.73	0.44	0.63	0.69	0.75	0.69	0.87	0.75	0.69	0.75	0.81	0.71	0.69	0.81	0.81	0.56	0.56	0.64	0.81	0.00

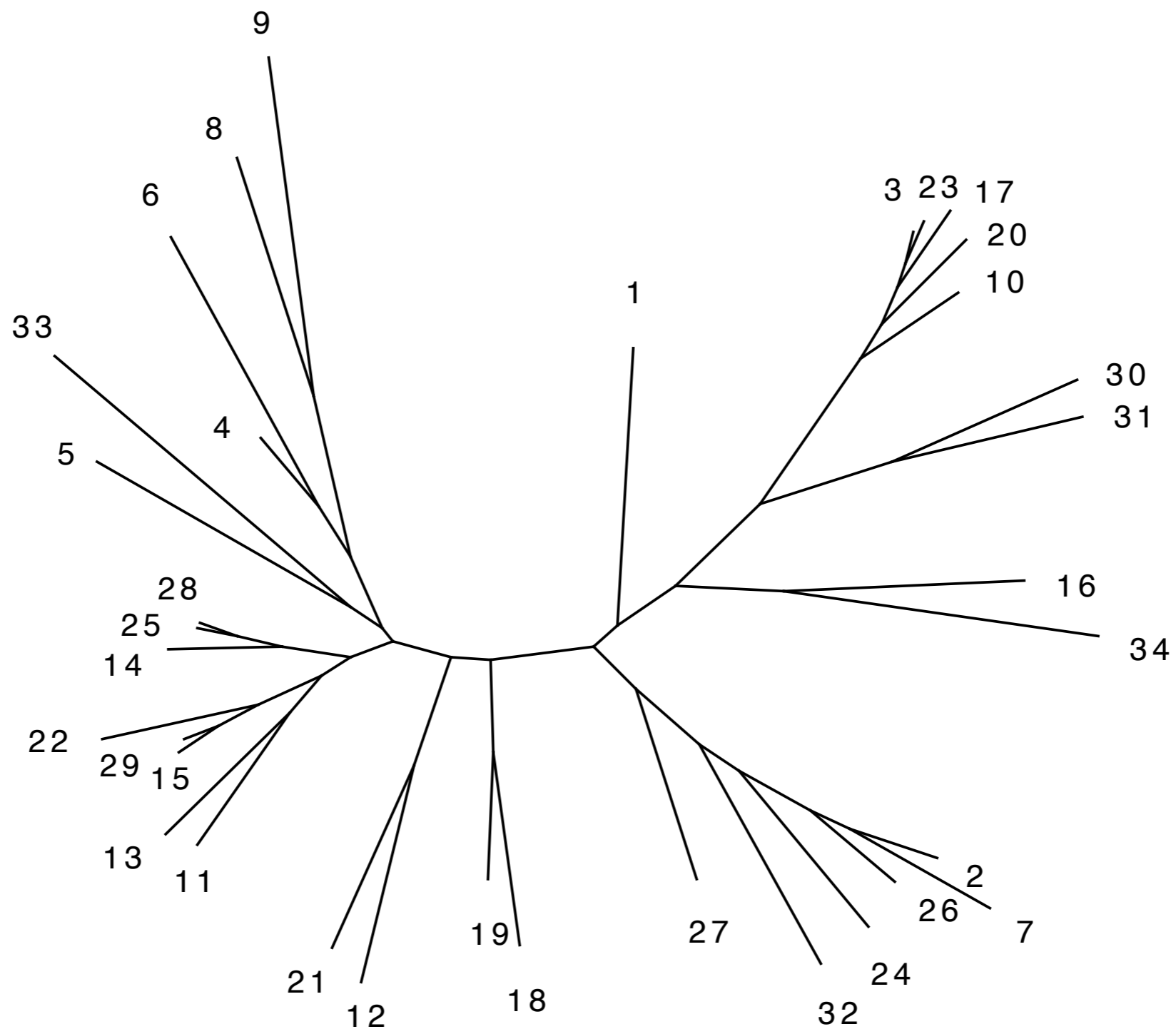
On average, there is a small difference between the role in context 4 and 14

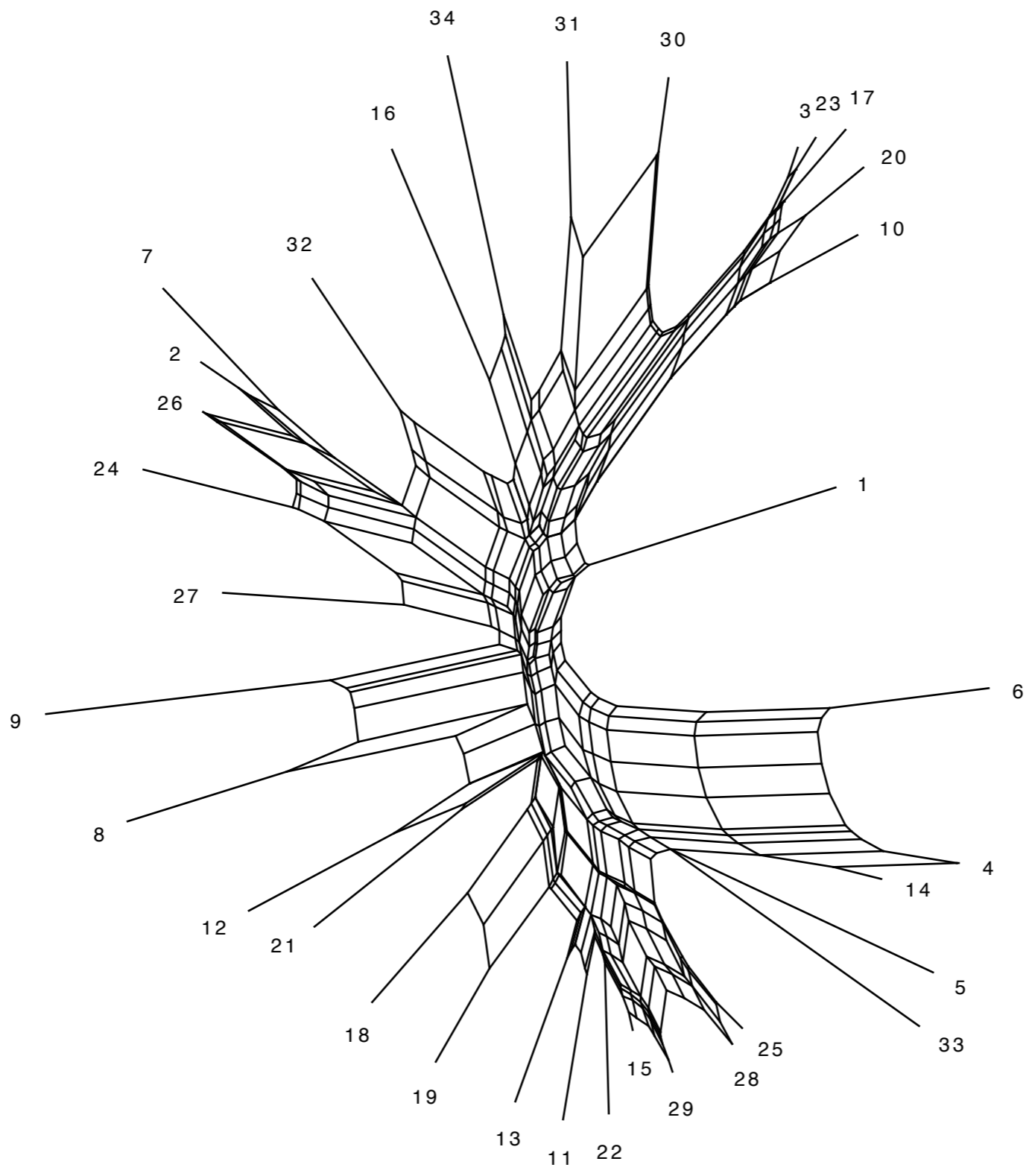
4: The Bible is a precious gift from God.
 14: The Bible reveals Jehovah's personality to us.

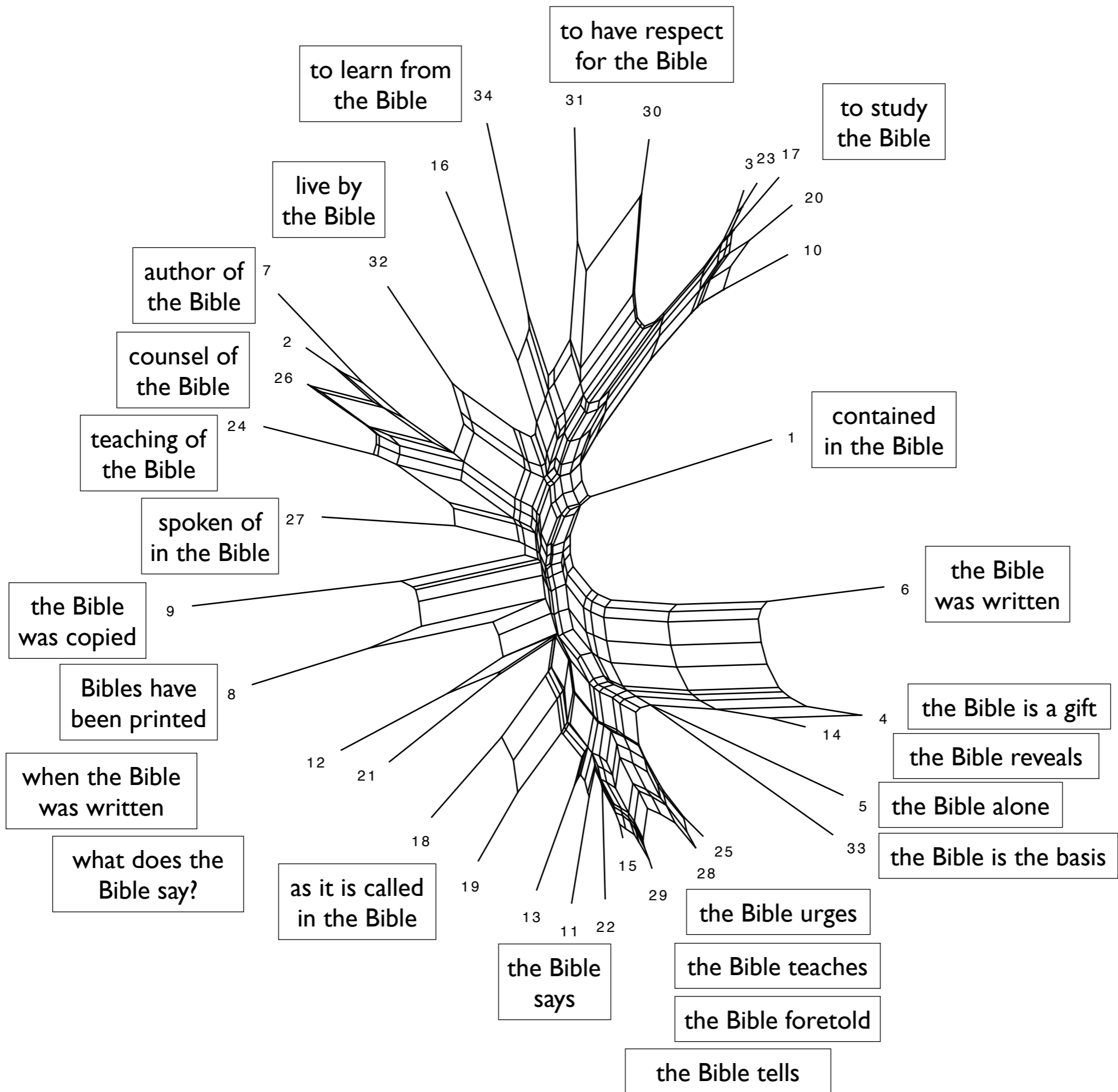
0.00	0.67	0.50	0.75	0.67	0.75	0.69	0.75	0.73	0.56	0.44	0.63	0.53	0.75	0.53	0.56	0.63	0.50	0.63	0.63	0.60	0.50	0.56	0.69	0.75	0.57	0.44	0.69	0.56	0.63	0.63	0.50	0.75	0.63	
0.67	0.00	0.73	0.80	0.73	0.67	0.20	0.73	0.64	0.60	0.80	0.67	0.93	0.73	0.57	0.80	0.73	0.73	0.73	0.67	0.50	0.67	0.73	0.40	0.60	0.15	0.53	0.67	0.67	0.60	0.60	0.38	0.73	0.80	
0.50	0.73	0.00	0.88	0.80	0.69	0.81	0.75	0.80	0.25	0.75	0.75	0.73	0.88	0.80	0.69	0.13	0.69	0.81	0.19	0.73	0.81	0.06	0.75	0.88	0.64	0.69	0.88	0.88	0.44	0.63	0.64	0.75	0.69	
0.75	0.80	0.88	0.00	0.40	0.31	0.81	0.56	0.53	0.88	0.44	0.44	0.40	0.13	0.47	0.75	0.94	0.63	0.63	0.94	0.53	0.44	0.94	0.88	0.25	0.93	0.81	0.31	0.44	1.00	1.00	0.79	0.50	0.94	
0.67	0.73	0.80	0.40	0.00	0.60	0.80	0.60	0.79	0.80	0.47	0.47	0.57	0.47	0.50	0.73	0.87	0.53	0.60	0.87	0.57	0.53	0.80	0.80	0.47	0.77	0.67	0.47	0.47	0.93	1.00	0.85	0.53	0.87	
0.75	0.67	0.69	0.31	0.60	0.00	0.69	0.63	0.53	0.75	0.63	0.63	0.67	0.44	0.67	0.81	0.75	0.75	0.75	0.88	0.60	0.56	0.75	0.81	0.56	0.79	0.81	0.63	0.63	0.81	0.81	0.71	0.63	0.94	
0.69	0.20	0.81	0.81	0.80	0.69	0.00	0.63	0.88			0.75	0.93	0.75	0.67	0.81	0.81	0.81	0.81	0.81	0.81	0.60	0.75	0.81	0.50	0.75	0.29	0.63	0.69	0.69	0.69	0.69	0.50	0.81	0.81
0.75	0.73	0.75	0.56	0.60	0.63	0.63	0.00				0.50	0.67	0.63	0.53	0.75	0.75	0.63	0.69	0.81	0.80	0.69	0.75	0.75	0.63	0.79	0.75	0.56	0.56	0.75	0.75	0.86	0.81	0.75	
0.73	0.64	0.80	0.53	0.79	0.53	0.60	0.47				0.73	0.86	0.67	0.64	0.80	0.87	0.87	0.87	0.80	0.86	0.67	0.87	0.73	0.67	0.69	0.73	0.67	0.73	0.80	0.80	0.62	0.80	0.80	
0.56	0.60	0.25	0.88	0.80	0.75	0.75	0.88	0.73	0.00	0.81	0.75	0.80	0.81	0.80	0.69	0.25	0.69	0.81	0.19	0.67	0.88	0.25	0.69	0.81	0.57	0.63	0.81	0.88	0.44	0.63	0.57	0.63	0.81	
0.44	0.80	0.75	0.44	0.47	0.63	0.81	0.56	0.80	0.81	0.00	0.56	0.27	0.38	0.27	0.69	0.81	0.44	0.44	0.81	0.53	0.31	0.75	0.81	0.38	0.86	0.63	0.31	0.19	0.88	0.88	0.86	0.56	0.81	
0.63	0.67	0.75	0.44	0.47	0.63	0.75	0.50	0.73	0.75	0.56	0.00	0.53	0.50	0.47	0.69	0.81	0.44	0.56	0.81	0.33	0.63	0.75	0.63	0.50	0.64	0.63	0.56	0.56	0.88	0.88	0.71	0.69	0.81	
0.53	0.93	0.73	0.40	0.57	0.67	0.93	0.67	0.86	0.80	0.27	0.53	0.00	0.27	0.36	0.60	0.73	0.40	0.33	0.80	0.43	0.40	0.73	0.73	0.33	0.77	0.60	0.33	0.33	0.80	0.80	0.85	0.60	0.73	
0.75	0.73	0.88	0.13	0.47	0.44	0.75	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.53	0.57	0.80	0.47	0.50	0.67	0.67	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.56	0.80	0.69	0.75	0.73	0.81	0.81	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.63	0.73	0.13	0.94	0.87	0.75	0.81	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.50	0.73	0.69	0.63	0.53	0.75	0.81	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.63	0.73	0.81	0.63	0.60	0.75	0.81	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.63	0.67	0.19	0.94	0.87	0.88	0.81	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.60	0.50	0.73	0.53	0.57	0.60	0.60	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.50	0.67	0.81	0.44	0.53	0.56	0.75	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.56	0.73	0.06	0.94	0.80	0.75	0.81	0.63	0.67	0.81	0.28	0.50	0.27	0.00	0.22	0.69	0.88	0.56	0.50	0.88	0.40	0.44	0.88	0.75	0.12	0.78	0.69	0.18	0.21	0.84	0.84	0.79	0.44	0.88	
0.69	0.40	0.75	0.88	0.80	0.81	0.50	0.75	0.73	0.69	0.81	0.63	0.73	0.75	0.53	0.75	0.81	0.63	0.56	0.75	0.40	0.75	0.75	0.00	0.63	0.21	0.38	0.63	0.63	0.69	0.63	0.43	0.81	0.75	
0.75	0.60	0.88	0.25	0.47	0.56	0.75	0.63	0.67	0.81	0.38	0.50	0.33	0.13	0.20	0.75	0.94	0.56	0.38	0.81	0.40	0.31	0.88	0.63	0.00	0.64	0.56	0.06	0.19	0.88	0.88	0.64	0.50	0.81	
0.57	0.15	0.64	0.93	0.77	0.79	0.29	0.79	0.69	0.57	0.86	0.64	0.77	0.79	0.57	0.71	0.71	0.64	0.57	0.64	0.46	0.71	0.64	0.21	0.64	0.00	0.36	0.64	0.64	0.57	0.57	0.31	0.86	0.71	
0.44	0.53	0.69	0.81	0.67	0.81	0.63	0.75	0.73	0.63	0.63	0.63	0.60	0.69	0.33	0.75	0.75	0.38	0.31	0.69	0.53	0.56	0.69	0.38	0.56	0.36	0.00	0.56	0.44	0.63	0.63	0.43	0.81	0.69	
0.69	0.67	0.88	0.31	0.47	0.63	0.69	0.56	0.67	0.81	0.31	0.56	0.33	0.19	0.20	0.75	0.94	0.56	0.38	0.81	0.47	0.31	0.88	0.63	0.06	0.64	0.56	0.00	0.13	0.88	0.88	0.64	0.50	0.81	
0.56	0.67	0.88	0.44	0.47	0.63	0.69	0.56	0.73	0.88	0.19	0.56	0.33	0.31	0.07	0.81	0.94	0.44	0.25	0.88	0.47	0.19	0.88	0.63	0.19	0.64	0.44	0.13	0.00	0.88	0.88	0.71	0.56	0.81	
0.63	0.60	0.44	1.00	0.93	0.81	0.69	0.75	0.80	0.44	0.88	0.88	0.80	0.94	0.80	0.69	0.31	0.81	0.81	0.44	0.80	0.88	0.44	0.69	0.88	0.57	0.63	0.88	0.88	0.00	0.31	0.57	0.75	0.56	
0.63	0.60	0.63	1.00	1.00	0.81	0.69	0.75	0.80	0.63	0.88	0.88	0.80	0.94	0.80	0.56	0.50	0.81	0.81	0.63	0.80	0.81	0.63	0.63	0.88	0.57	0.63	0.88	0.88	0.31	0.00	0.36	0.81	0.56	
0.50	0.38	0.64	0.79	0.85	0.71	0.50	0.86	0.62	0.57	0.86	0.71	0.85	0.79	0.64	0.57	0.79	0.64	0.64	0.64	0.62	0.57	0.71	0.43	0.64	0.31	0.43	0.64	0.71	0.57	0.36	0.00	0.79	0.64	
0.75	0.73	0.75	0.50	0.53	0.63	0.81	0.81	0.80	0.63	0.56	0.69	0.60	0.44	0.60	0.63	0.63	0.69	0.69	0.69	0.67	0.56	0.75	0.81	0.50	0.86	0.81	0.50	0.56	0.75	0.81	0.79	0.00	0.81	
0.63	0.80	0.69	0.94	0.87	0.94	0.81	0.75	0.80	0.81	0.81	0.81	0.73	0.88	0.73	0.44	0.63	0.69	0.75	0.69	0.87	0.75	0.69	0.75	0.81	0.71	0.69	0.81	0.81	0.56	0.56	0.64	0.81	0.00	

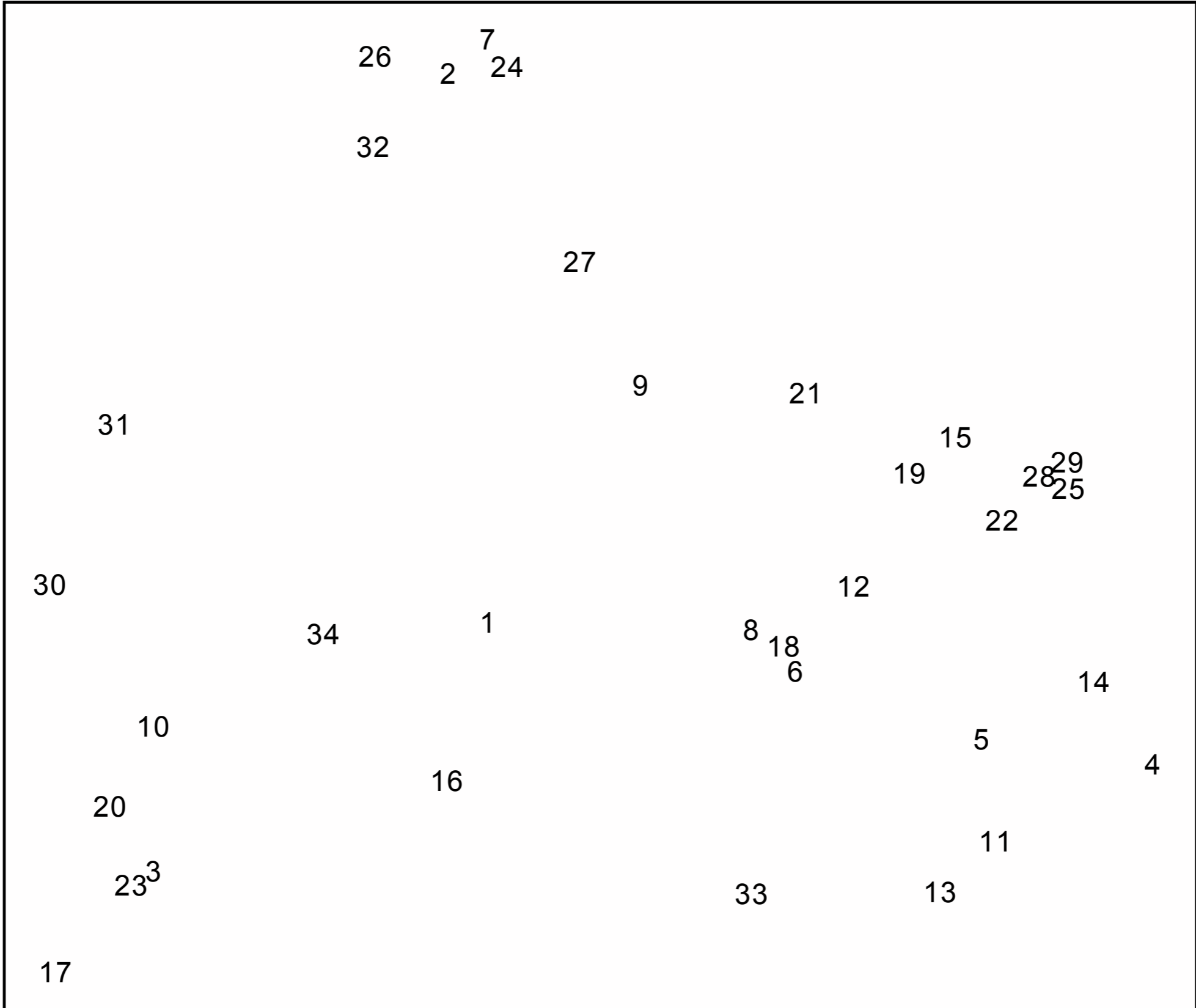
On average, there is a large difference between the role in context 8 and 10

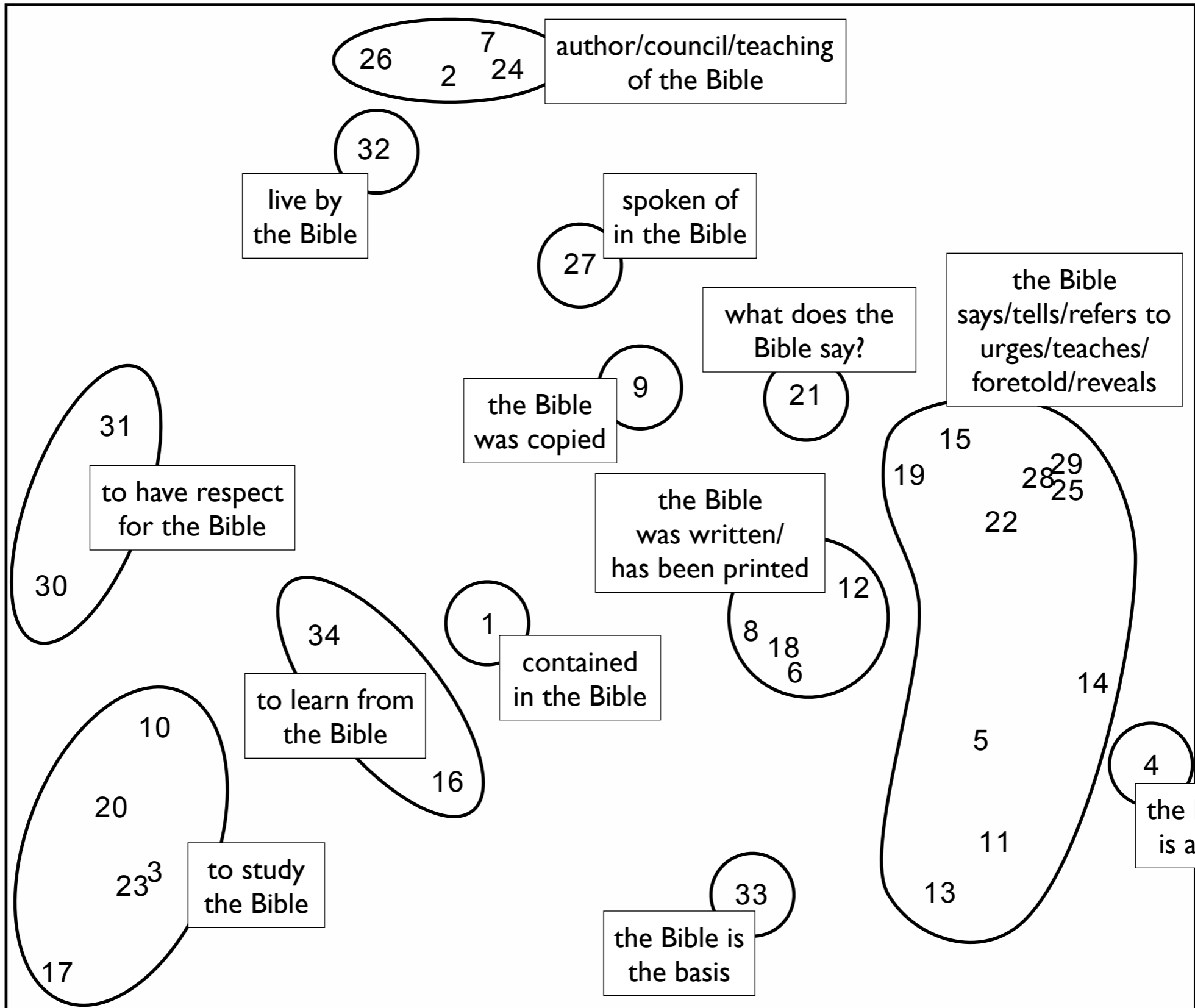
8: God made sure that the Bible was accurately copied and preserved.
 10: Not everyone will be happy to see you studying the Bible.



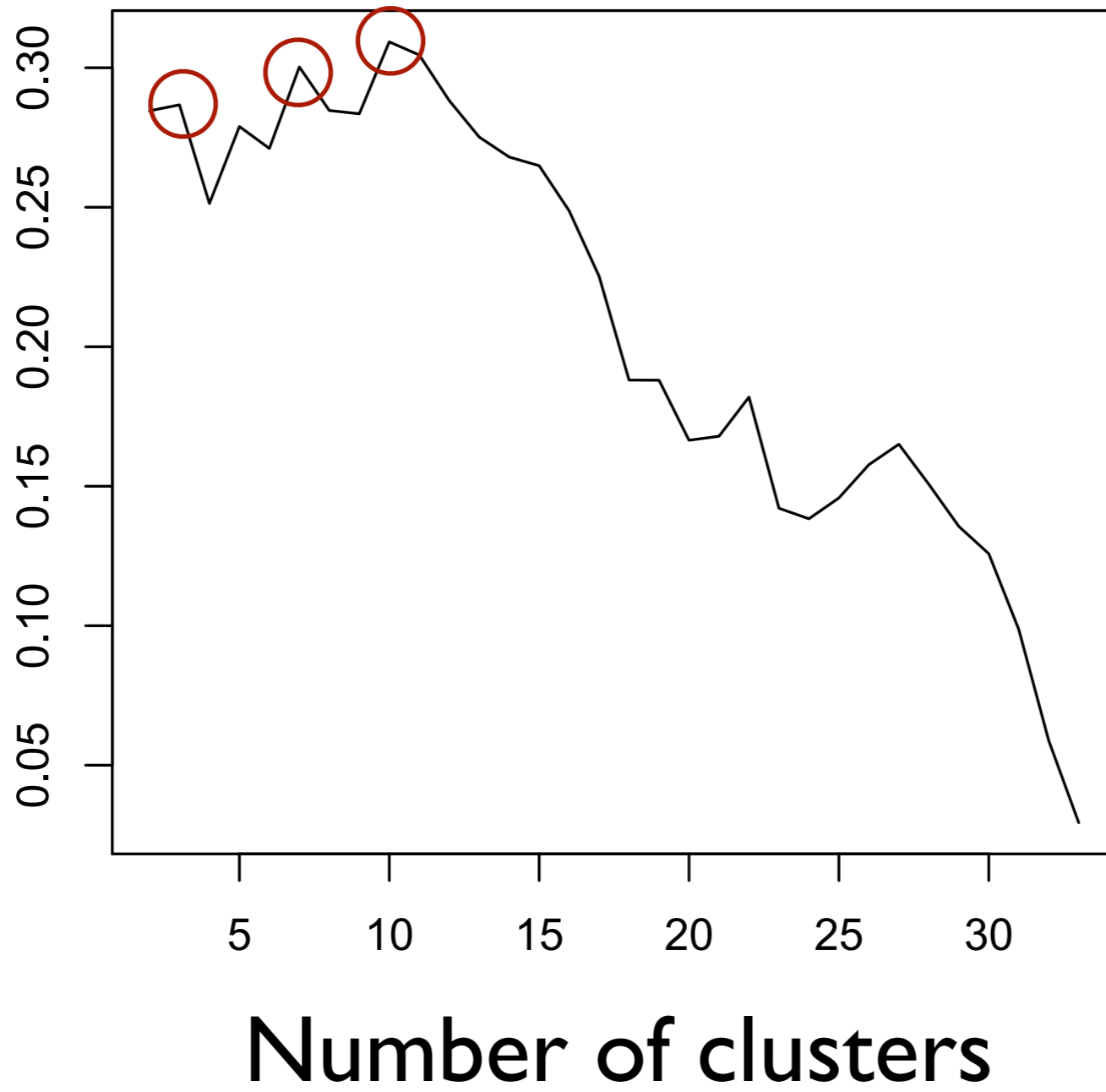


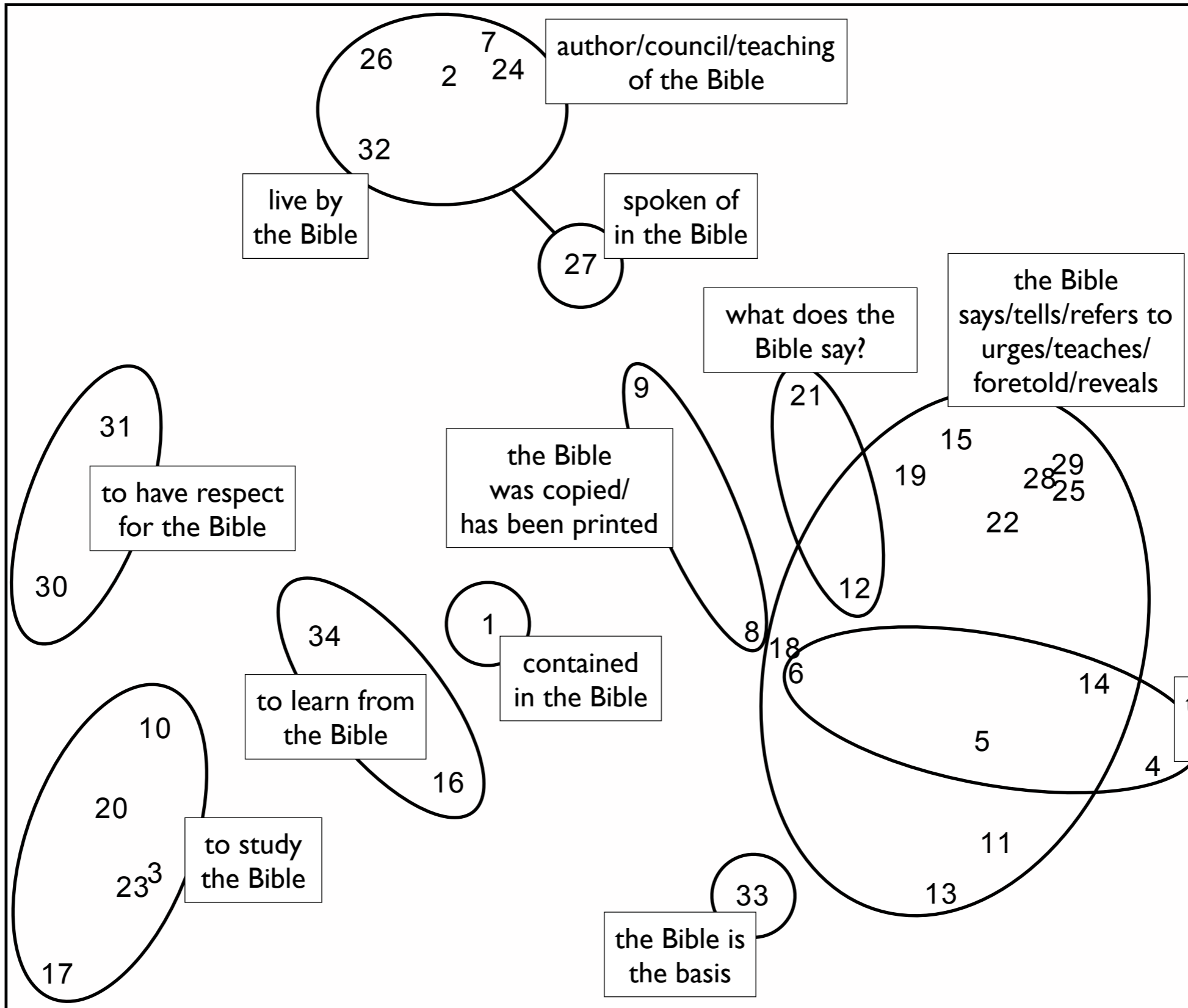


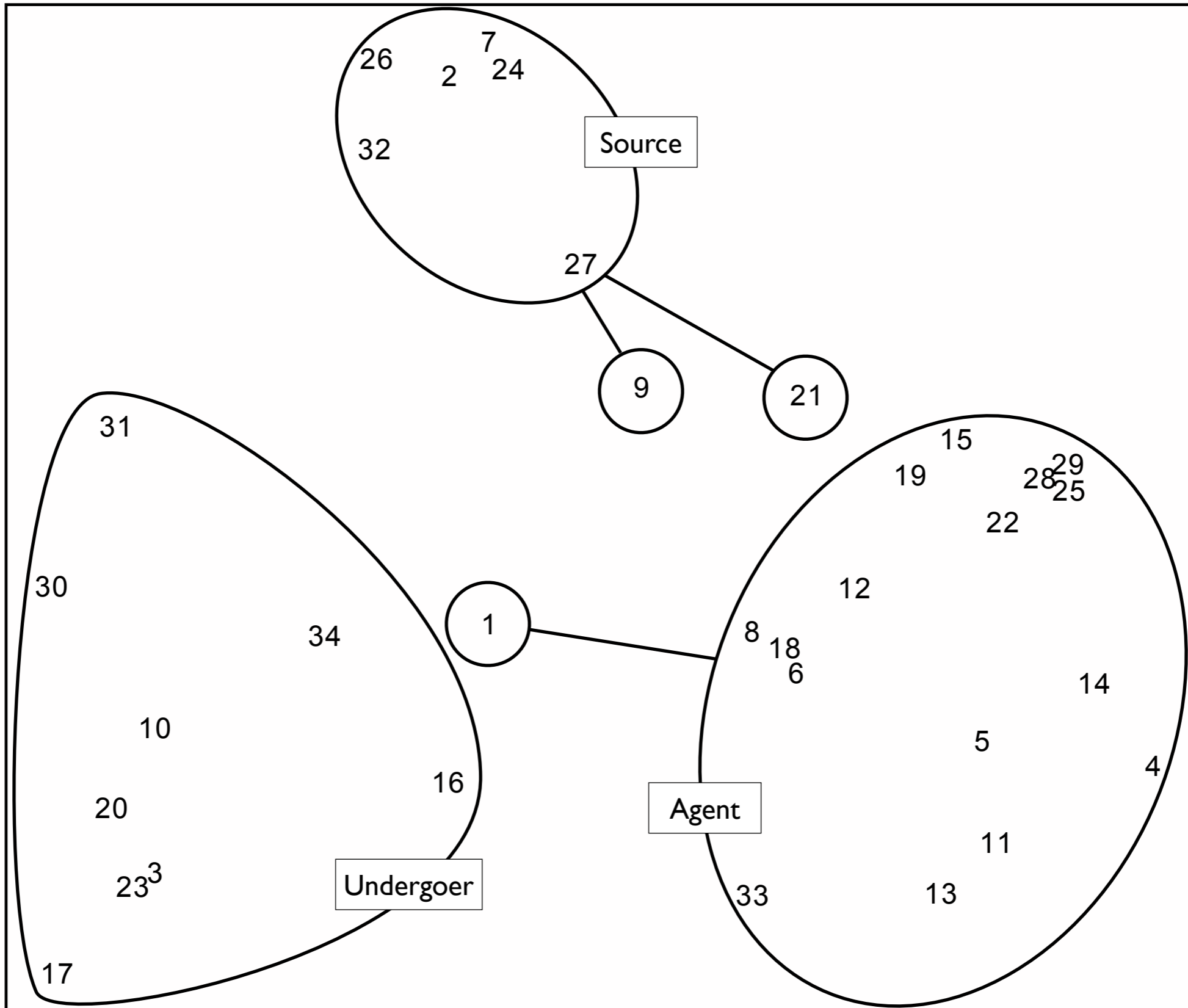




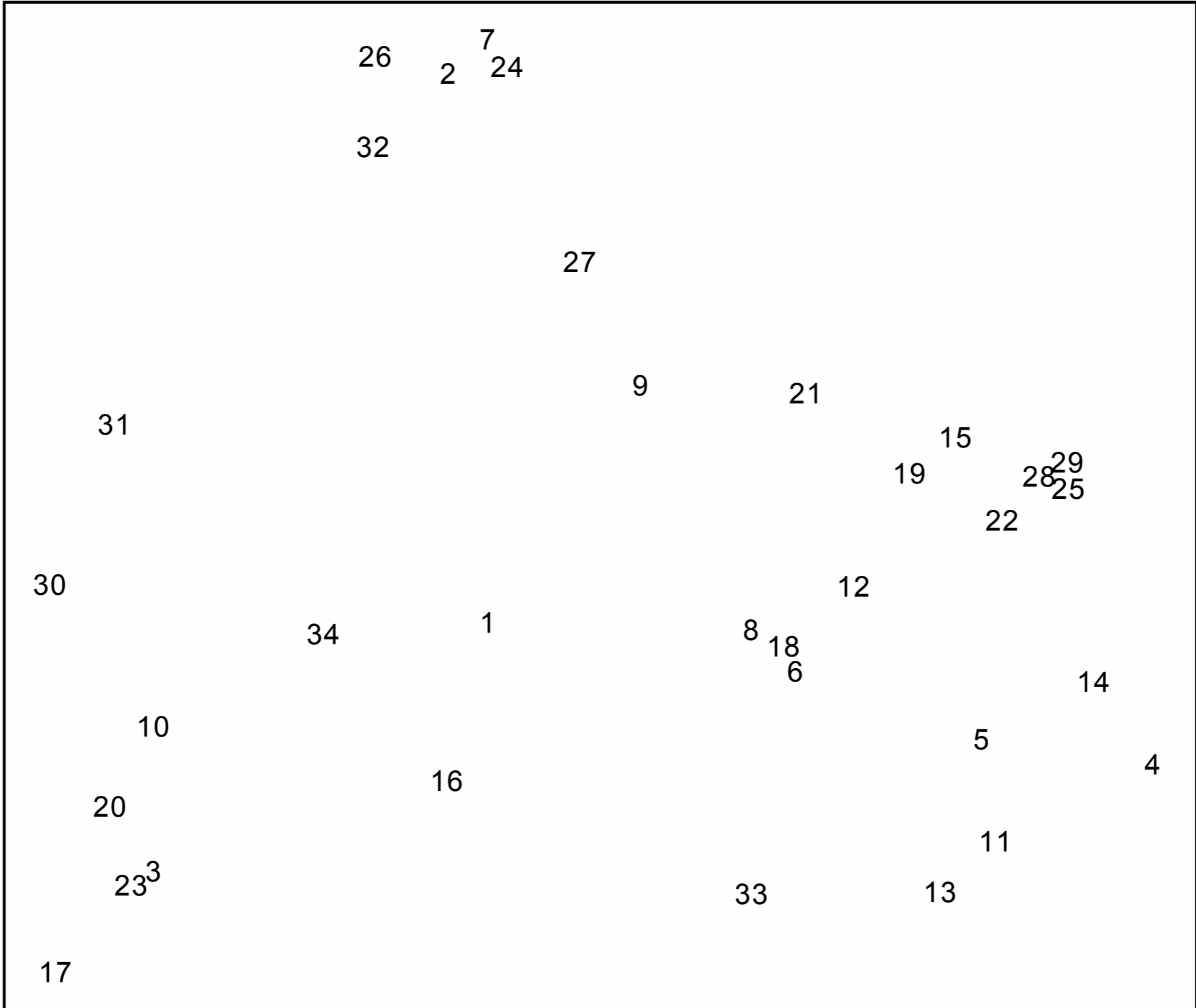
“Fit” of clustering



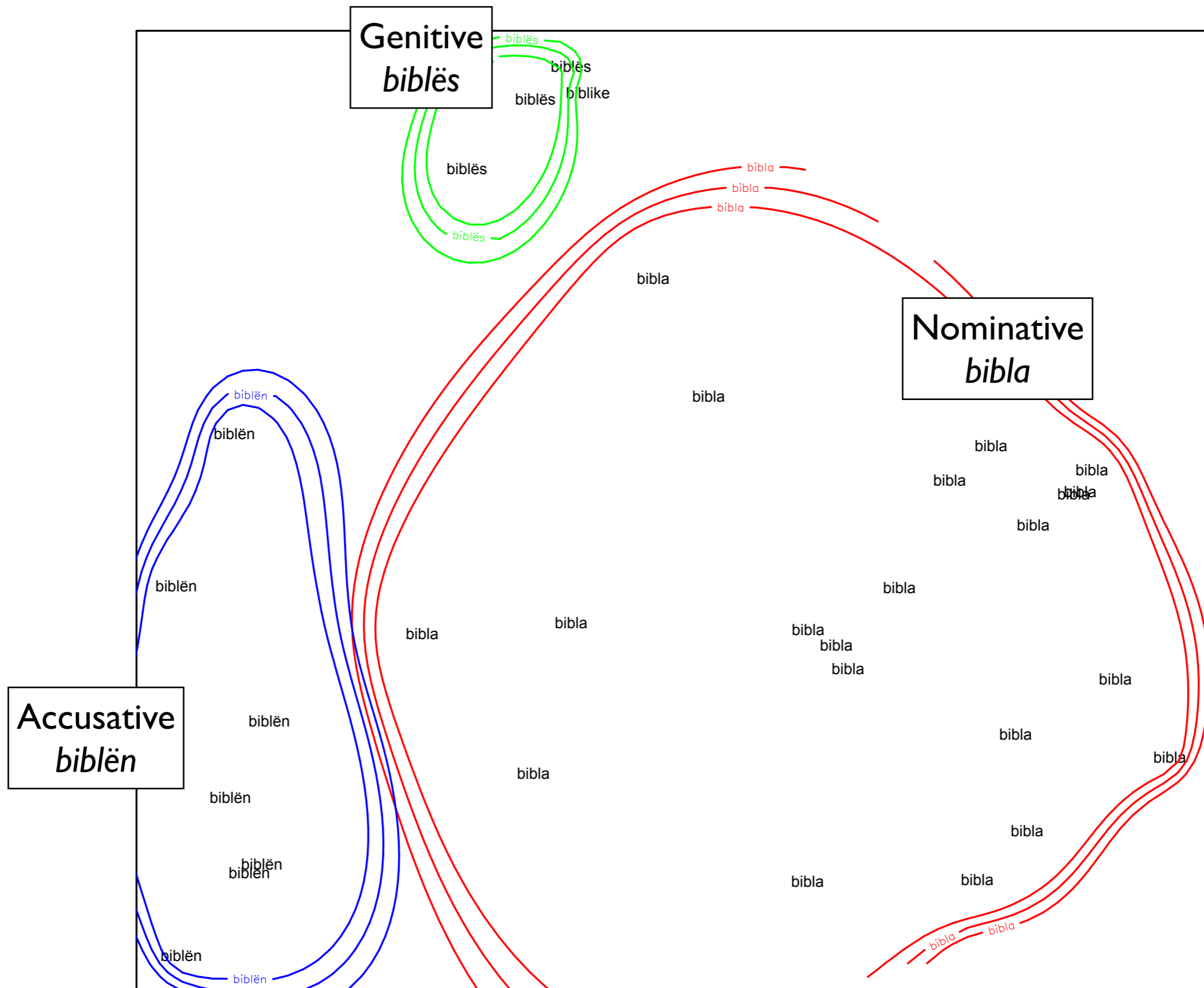




Comparing Languages



albanian



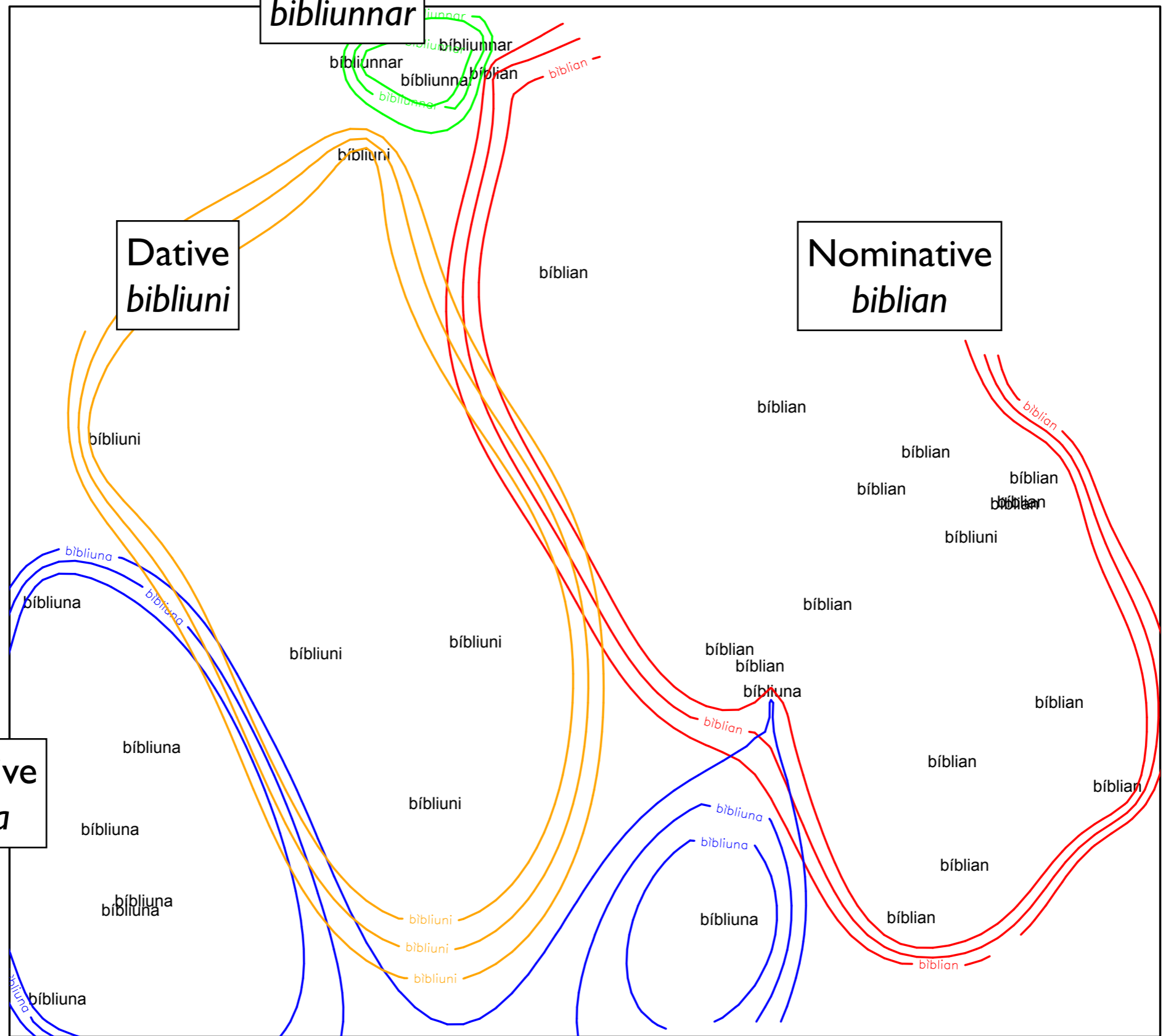
faroese

Genitive
bibliunnar

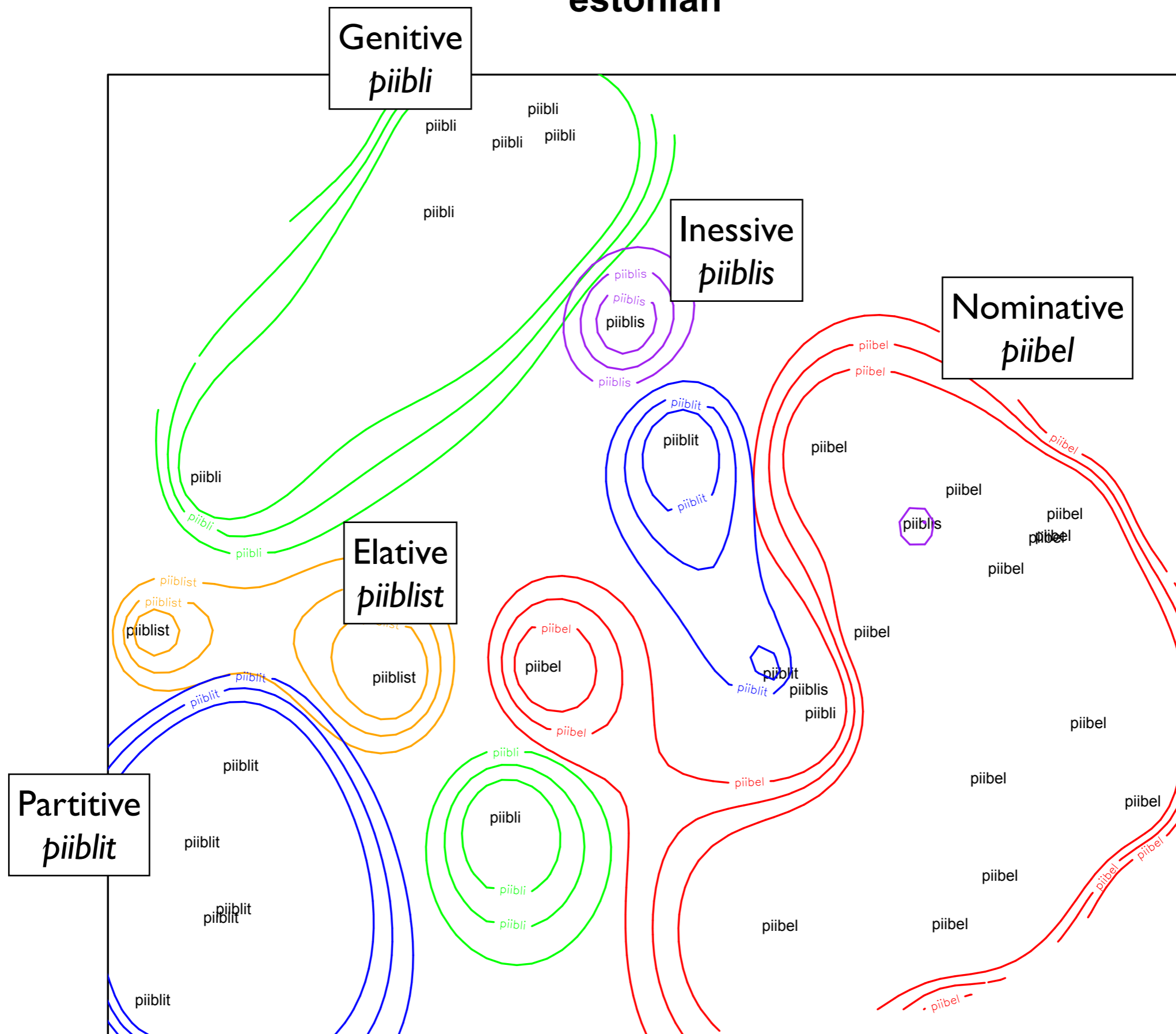
Dative
bibliuni

Nominative
biblian

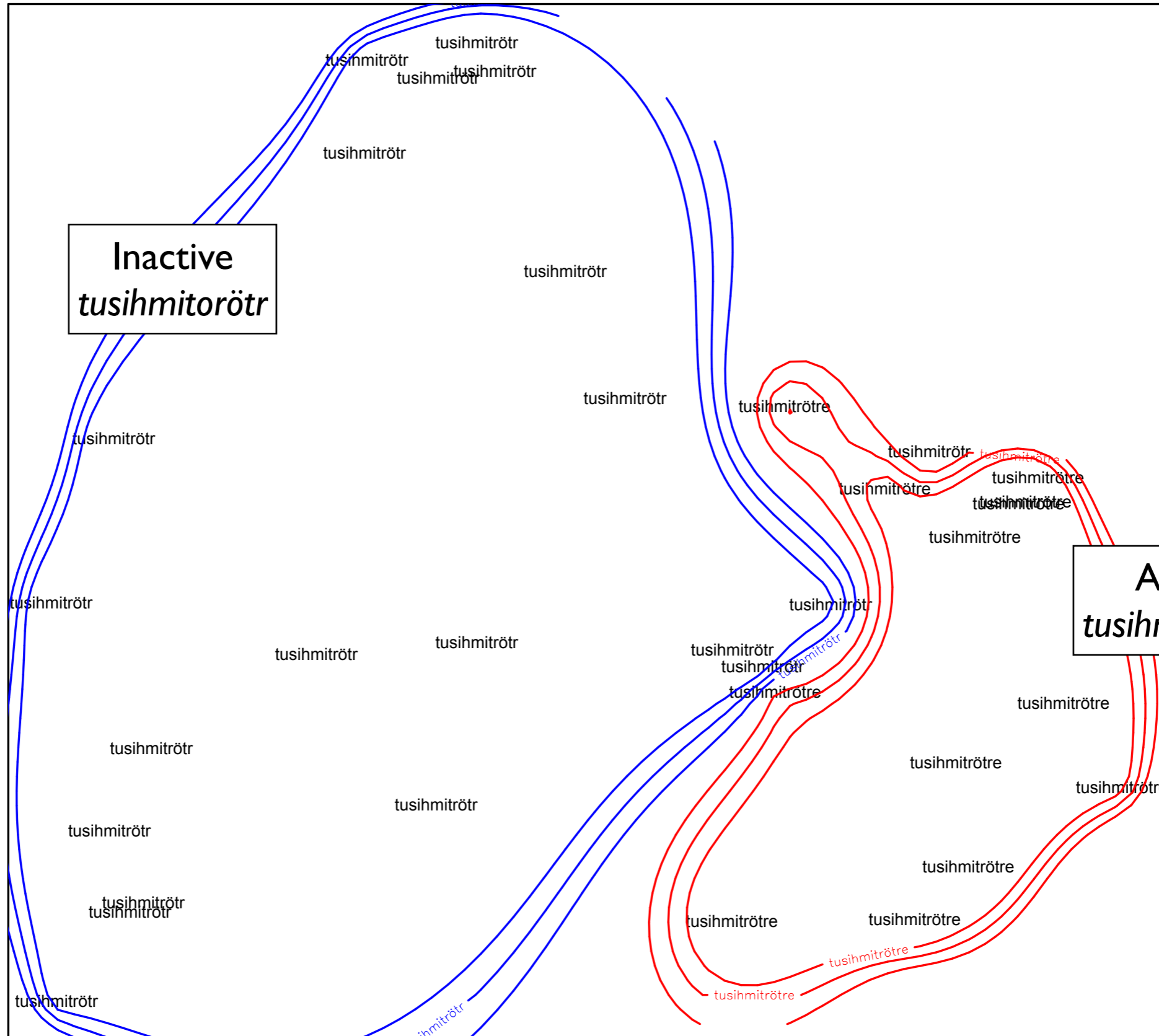
Accusative
bibliuna



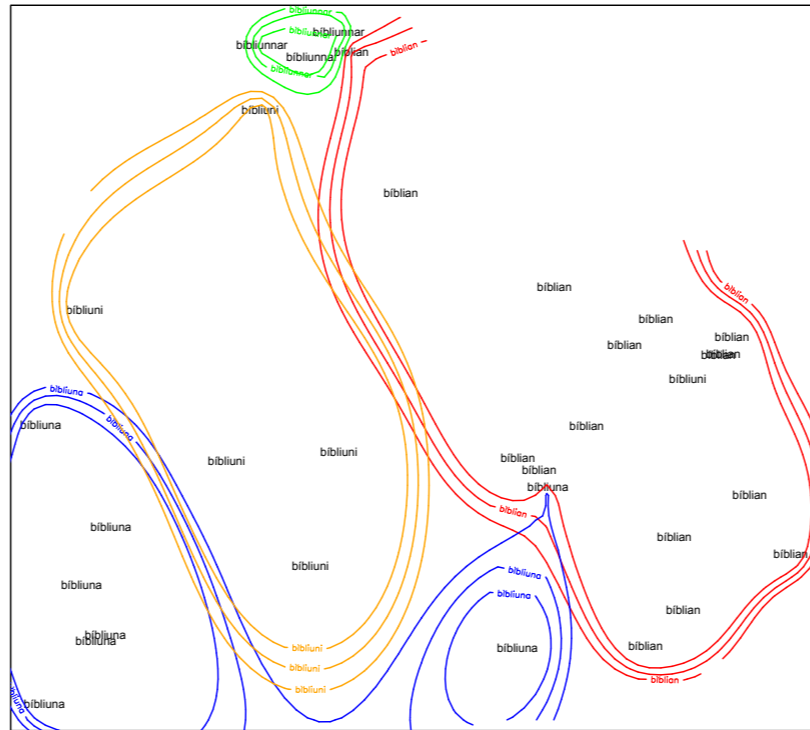
estonian



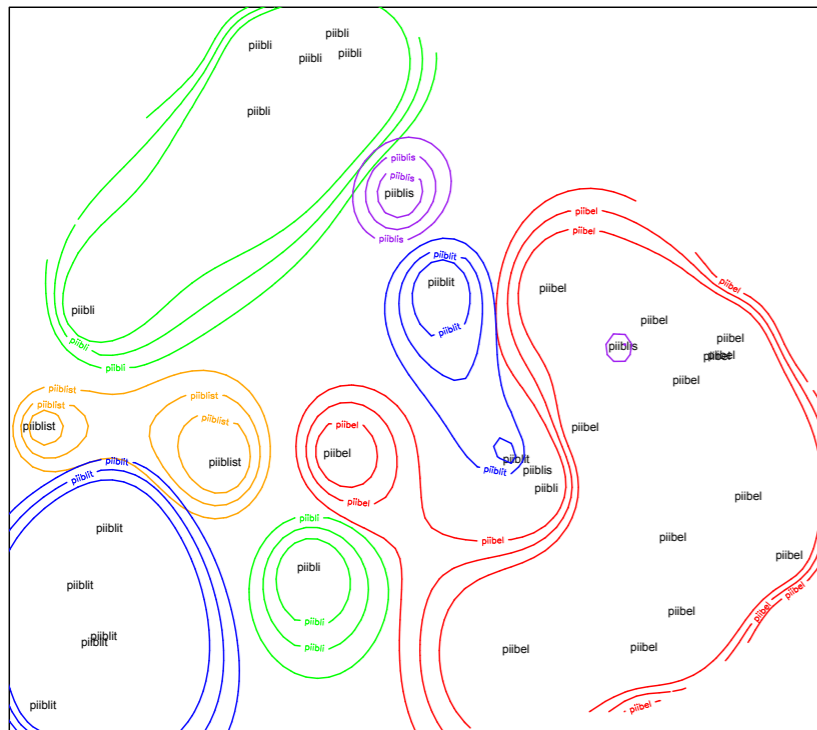
drehu



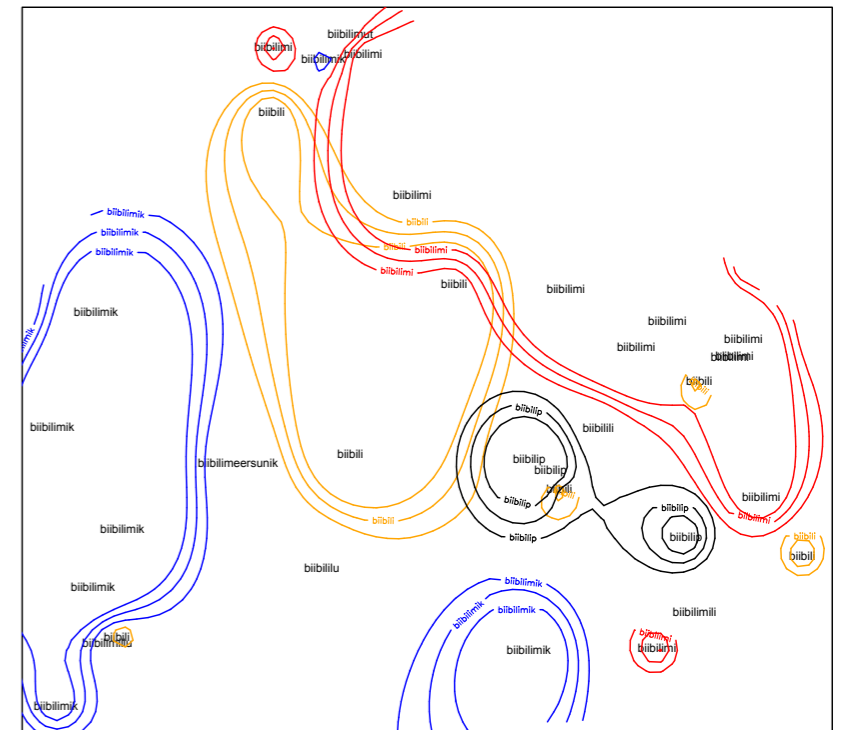
faroese



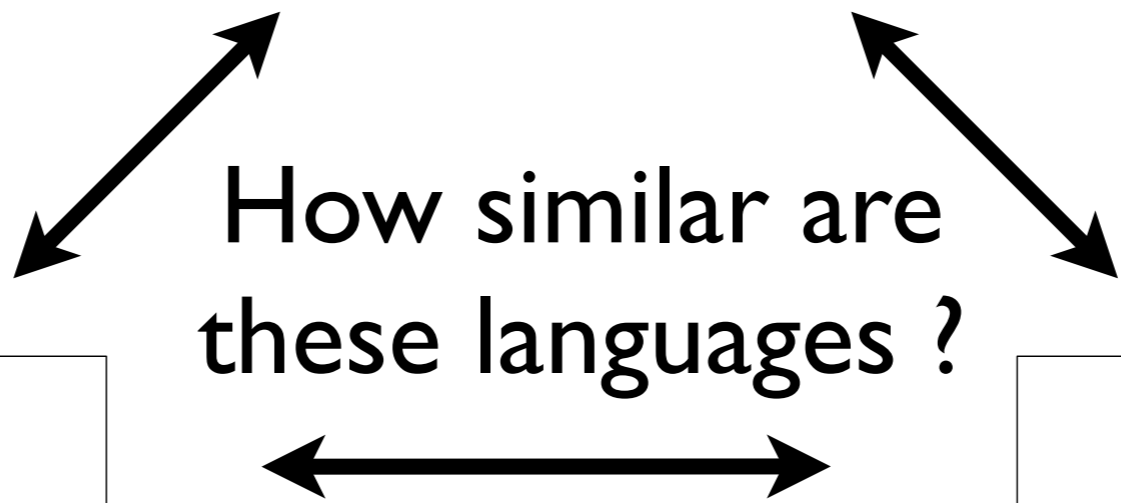
estonian



greenlandic



How similar are these languages ?



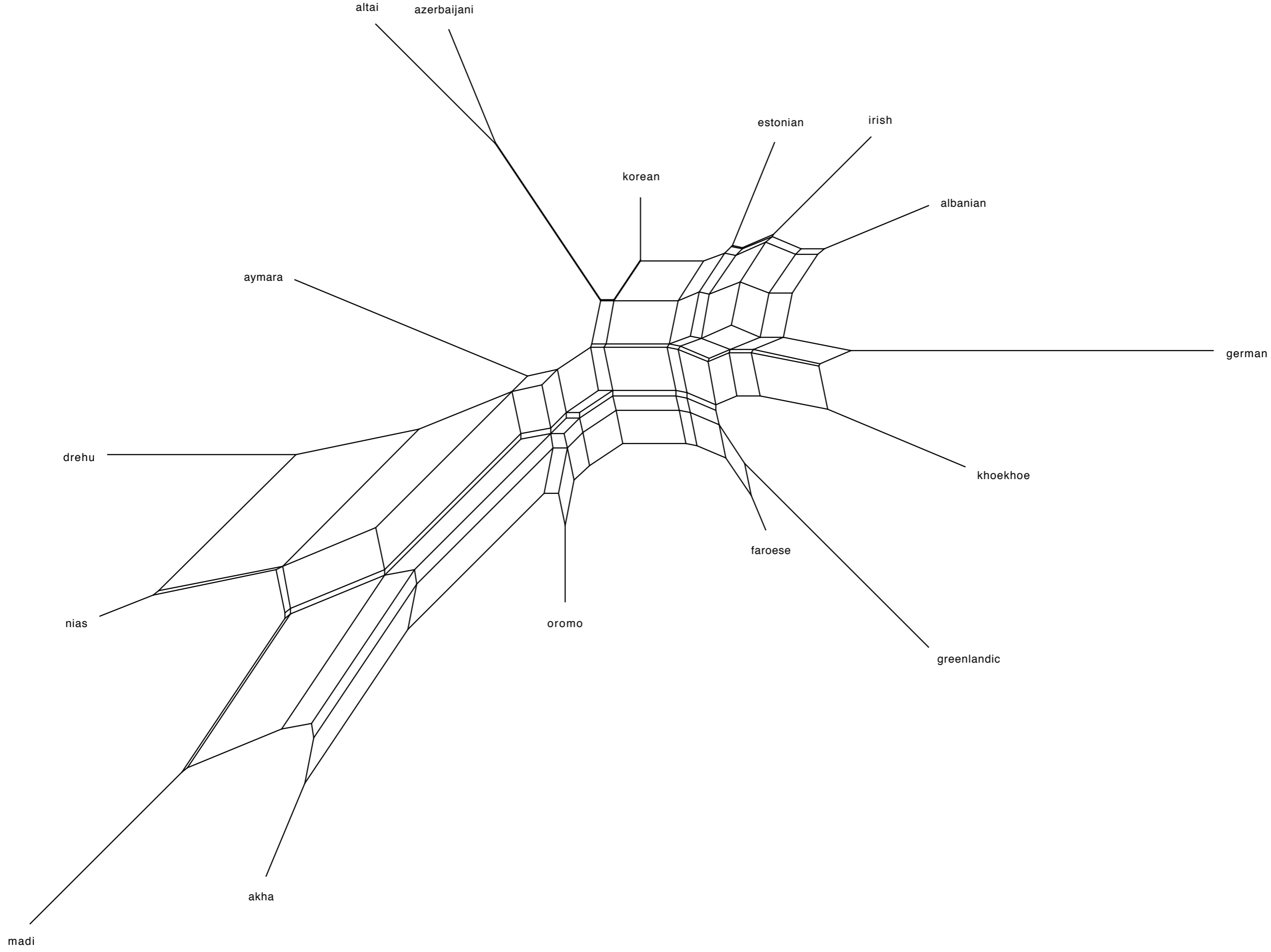
Context	Albanian	Faroese	Estonian	Greenlandic
1	bibla	bíbliuni	piibel	biibili
2	biblës	bíbliunnar	piibli	biibilimik
3	biblën	bíbliuna	piiblit	biibili
4	bibla	bíblían	piibel	biibili
5	bibla	bíblían	piibel	biibilip

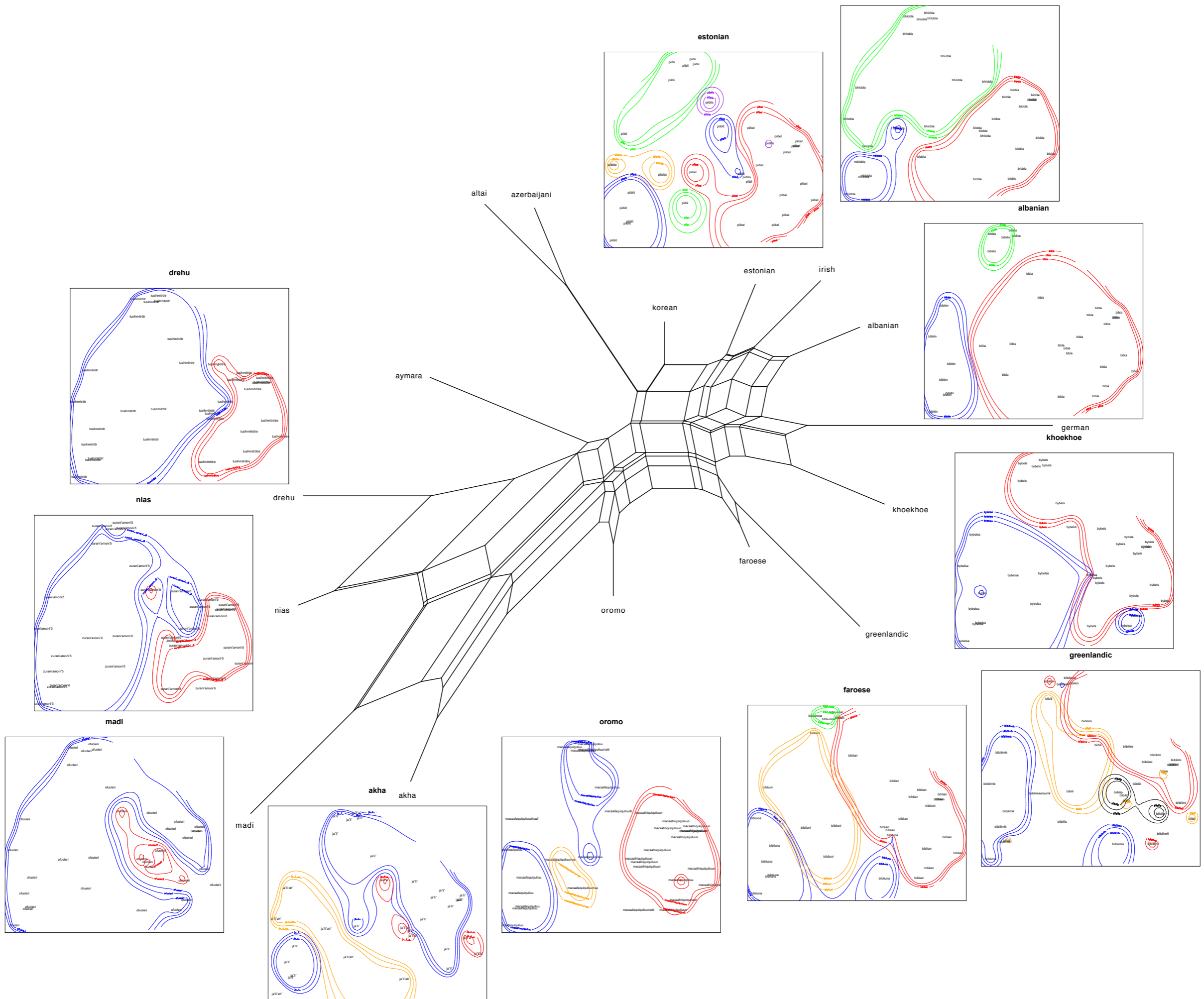
	1	2	3	4	5	...
1	0	1	1	0	0	
2	1	0	1	1	1	
3	1	1	0	1	1	
4	0	1	1	0	0	
5	0	1	1	0	0	
...						

bíbliuna
bíbliunnar
bíblían
NA
bíbliuna
bíblían
bíblían
bíblían
bíblían
bíbliuni
bíbliuna
bíblían
bíblían

	1	2	3	4	5	...
1	0	1	0	0	1	
2	1	0	1	1	1	
3	0	1	0	0	1	
4	0	1	0	0	1	
5	1	1	1	1	0	
...						

20	biblën	bíbliuna	piiblit	biibilimik
21	NA	bíblían	piibel	biibilimi
22	bibla	bíbliuni	piibel	biibili
23	biblën	bíbliuna	piiblit	biibilimillu
24	biblike	bíblían	piibli	biibilimi
25	bibla	bíblían	piibel	biibilimi
26	biblës	bíbliunnar	piibli	biibilimi
27	bibla	bíblían	piiblis	biibilimi
28	bibla	bíblían	piibel	biibilimi
29	bibla	bíblían	piibel	biibilimi
30	biblën	bíbliuna	piiblist	biibilimik
31	biblën	bíbliuni	piibli	biibilimik
32	biblës	bíbliuni	piibli	biibili
33	bibla	bíbliuna	piibel	biibilimik
34	bibla	bíbliuni	piiblist	biibilimeersunik





How to compare like with like?

- Translations offer a suitable approach to select **functionally** parallel data
- Then, form similarity is established **within each language separately**
- Average **form similarity across many languages** is a proxy to general functional similarity
- Distribution of **form similarity within a language** gives a language classification

Technicalities

Table III Term-document matrix

<i>Term</i>	<i>Doc1</i>	<i>Doc2</i>	<i>Doc3</i>	<i>Doc4</i>	<i>Doc5</i>	<i>Doc6</i>
Passenger traffic volume	1	1	0	5	2	0
Decrease	1	2	1	0	0	0
Increase	0	2	0	0	0	0
Passengers carried	5	1	0	0	0	0
Personal traffic tools	1	0	0	0	0	0
Grow up	4	1	6	0	0	0
Million	4	1	0	0	0	0
Hundred	0	0	0	0	1	0
FAST rapid transit system	0	2	0	0	0	0
Finished	0	1	0	0	0	0
A1 station	0	0	0	5	4	4
B1 station	0	0	0	1	5	0
C1 station	0	0	0	1	0	0
D1 station	0	0	0	1	0	1
E1 station	0	0	0	1	0	2
Passenger-Kilometers	0	1	7	0	0	0
Columniation	0	0	0	0	2	0
Check the number	0	0	0	0	2	0
Ticket Revenues	0	0	0	0	0	7

Trans-cooccurrences

$$\mathbf{O} = \mathbf{WS} \cdot \mathbf{WS}^T$$

$$\mathbf{E} = \mathbf{WS} \cdot \frac{\mathbf{1}_{SS}}{n} \cdot \mathbf{WS}^T$$

$$\mathbf{WW} = -\log \left[\frac{\mathbf{E}^{\mathbf{O}} \exp(-\mathbf{E})}{\mathbf{O}!} \right]$$

$$= \mathbf{E} + \log \mathbf{O}! - \mathbf{O} \log \mathbf{E}$$

who are god s true worshipers on earth today

who will be resurrected

who will rule with jesus

who created all living things

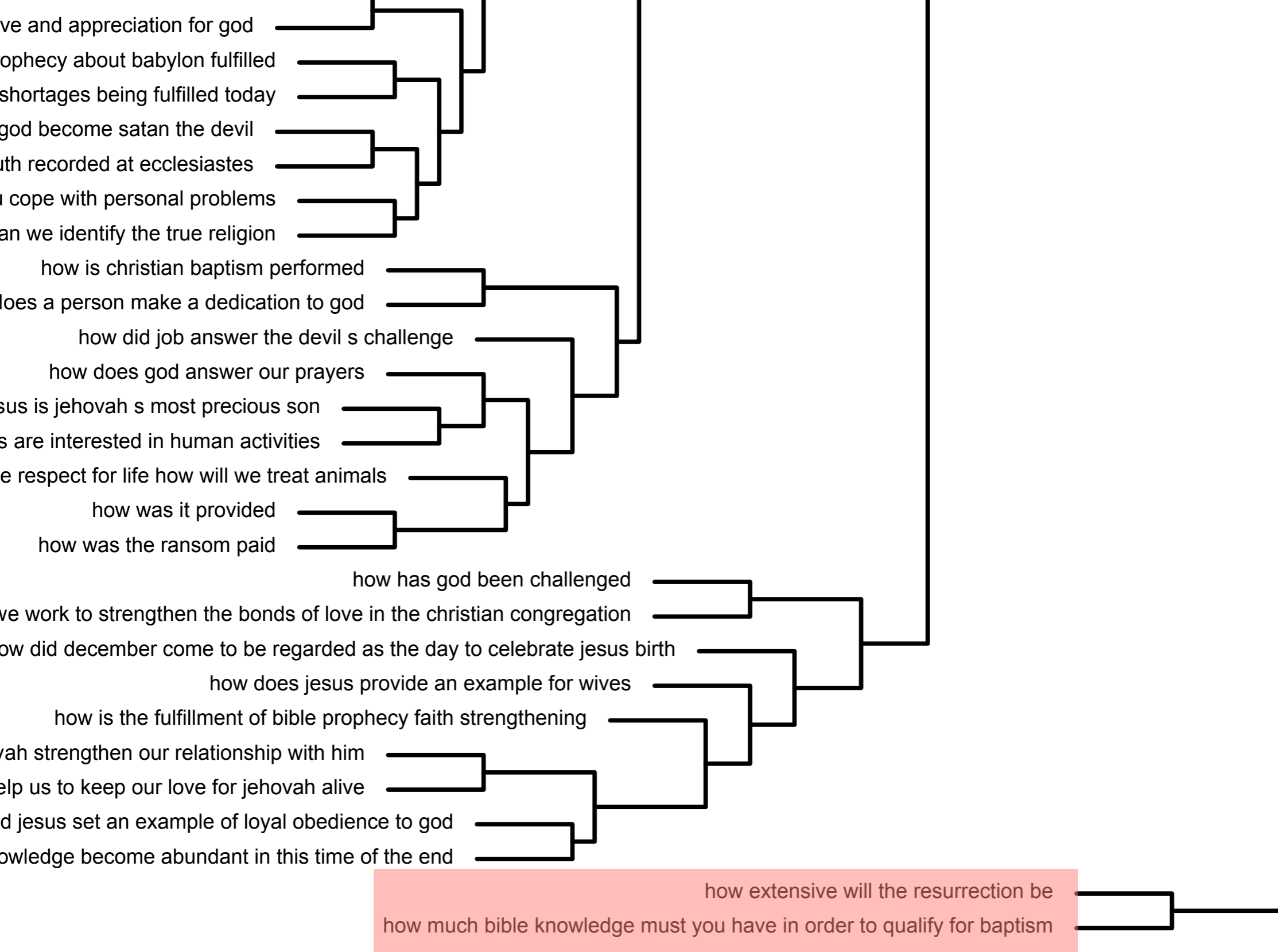
who is jesus christ

who is michael the archangel

**single 'who'
word**

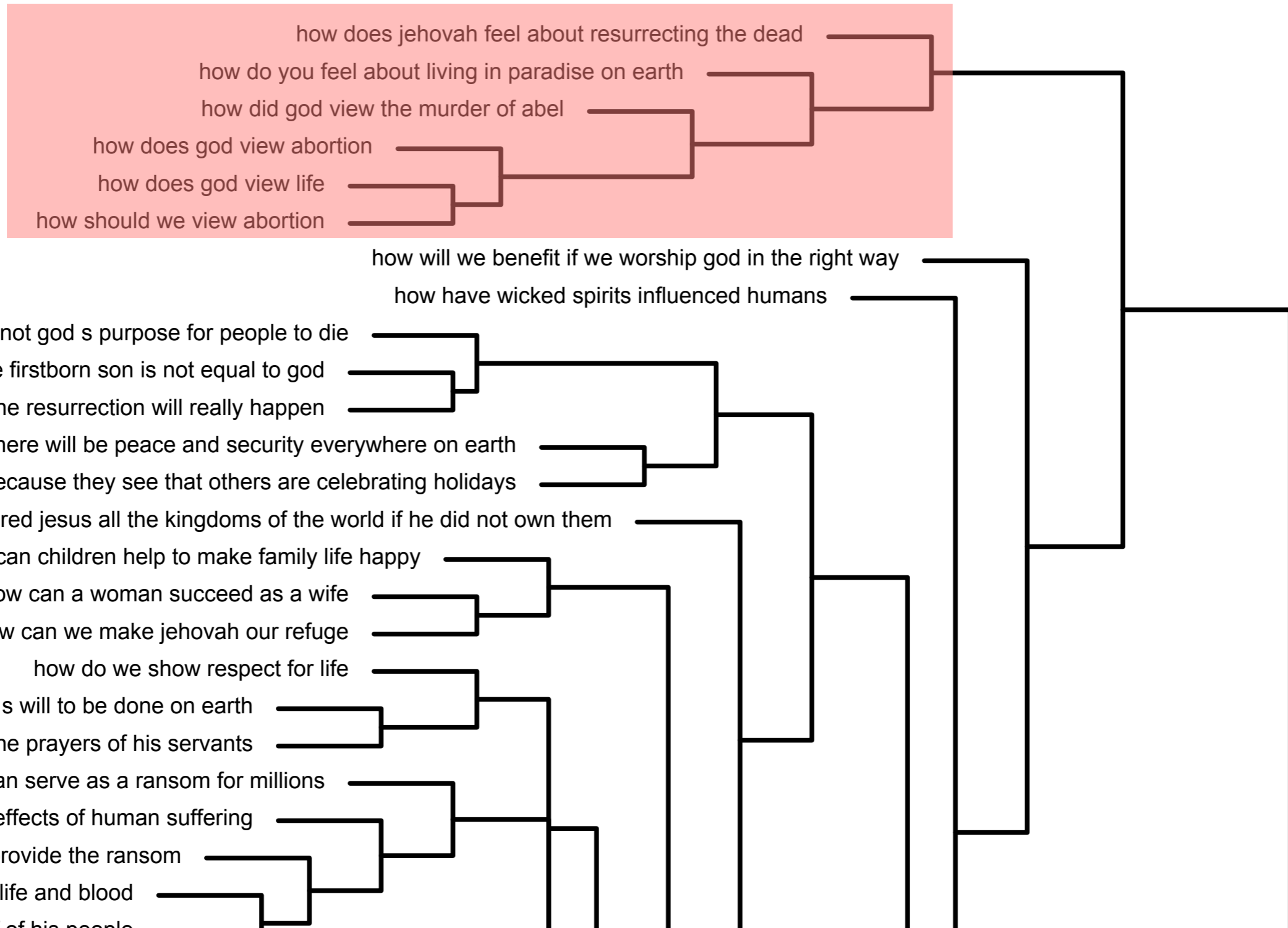
**singular vs.
plural 'who'
words**

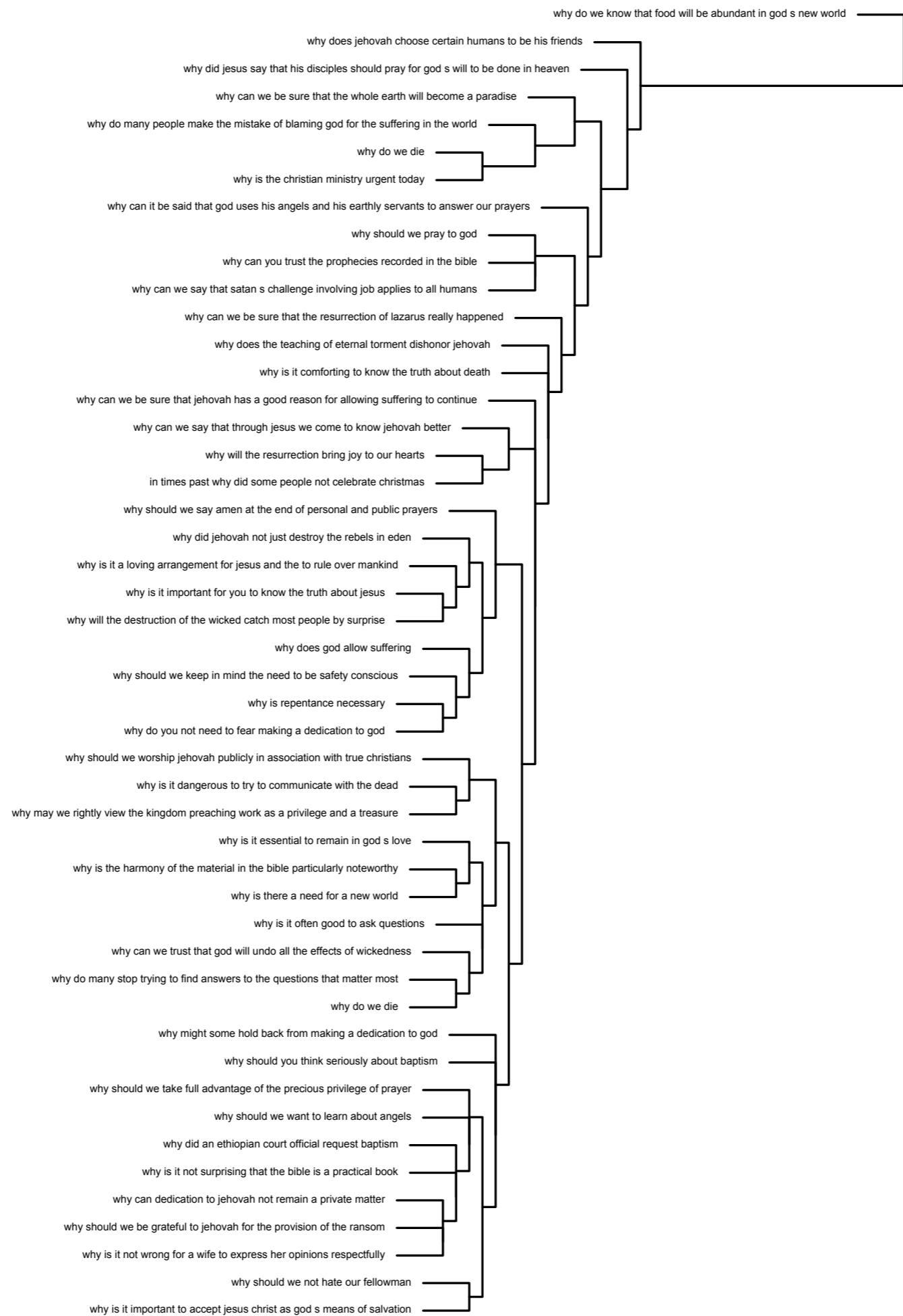




‘to which extent’

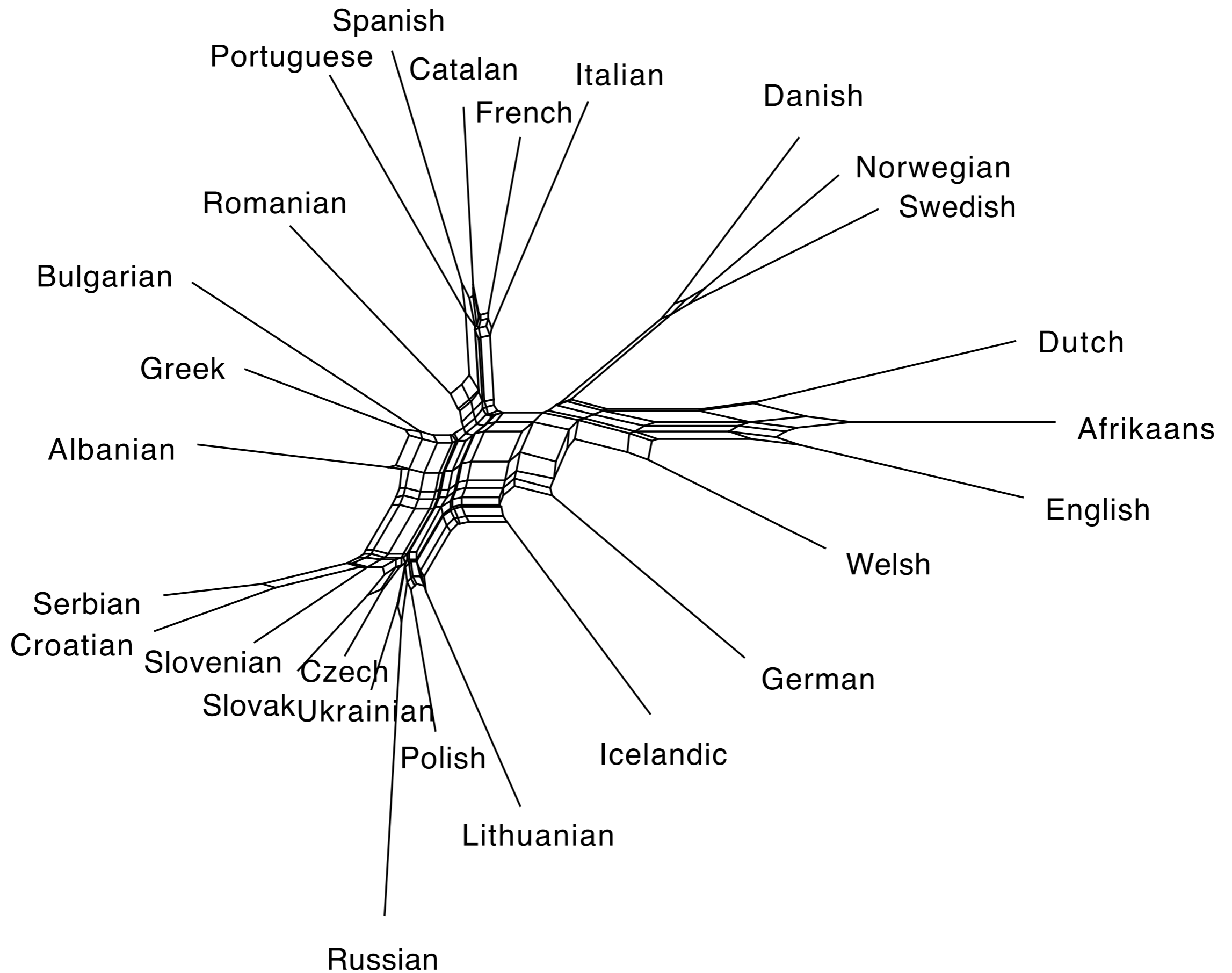
'how do you feel about', 'what do you think'





Sparse alignment

- Structure of languages are really different
- So: do not force an alignment across all languages
- Only align those words that really form a coherent group, 'sparse alignment'
 - ▶ we used affinity propagation clustering in each sentence
- The more words are aligned between two languages, the more similar are these languages



Prospects

- We are organizing large massively parallel texts for automatic processing
- Main problem: copyright
- Linking to monolingual corpora in cooperation with <http://wortschatz.uni-leipzig.de>
- Adding information will be open
 - ▶ morpheme separation, multi-word recognition, etc.
- **Typology will never be the same again!**