

Kahdennekymmenennetensimmäiset kansainväliset kielitieteen olympialaiset

Brasília (Brasilia), 23.–31. heinäkuuta 2024

Joukkuekilpailun tehtävä

Leksikostatistiikka on joukko menetelmiä, joiden tarkoituksena on arvioida, kuinka läheistä sukua kielet ovat toisilleen niiden sanastojen perusteella. Näitä menetelmiä sovelletaan tavallisesti pitkiin sanalistoihin, jotka asiantuntijat merkitsevät manuaalisesti ilmoittamalla, oletetaanko tietyn sanan olevan peräisin samasta lähteestä. Joskus kuitenkin kielitieteilijät käyttävät leksikostatistisia menetelmiä sanalistoihin, jotka on merkitty automaattisesti. Yksi tällainen automaattinen menettely perustuu *konsonanttiluokkien* käsitteeseen, jonka neuvostoliittolais-israelilainen kielitieteilijä Aharon Dolgopolsky esitti vuonna 1964.

P.	p b ɸ β f v	K.	k g x γ q ɣ χ ɥ	Y.	j ç (juuren alussa)	M.	m ɱ
T.	t d ɗ θ ð ʈ ɖ	R.	r ɾ ɽ l ʎ ʝ ʟ ʠ	W.	w ɰ (juuren alussa)	N.	n ɲ ŋ ŋ
S.	s z ʃ ʒ ʂ ʐ ʑ ʒ					Q.	ʈ ɖ
H.	h ʕ ɦ ʕ ʔ h ɦ ʔ, vokaalit ja j ç w ɰ (paitsi juuren alussa)						

Dolgopolskyn konsonanttiluokat

Alla on merkittyjen sanalistojen katkelmia useista maailman kielikunnista. Merkinnät annetaan alaindeksein. Näiden listojen perusteella on luotu kielten sukupuut käyttäen kahta yksinkertaistettua versiota ns. *StarlingNj*-algoritmista, minkä lisäksi jokaiselle sanalle on annettu stabiilisuusindeksi. Yläosan puut ja stabiilisuusindeksit perustuvat manuaalisesti merkittyihin sanalistoihin ja alaosan puut ja stabiilisuusindeksit sanalistoihin, jotka on merkitty automaattisesti. Jokaisen sanalistan perusteella on rakennettu kaksi puuta, jotka seuraavat kahta algoritmiversiota: Algoritmi A ja Algoritmi B. Huomatkaa, että joissakin tapauksissa on useita mahdollisia puuta, jotka vastaavat yhtä sanalista; näissä tapauksissa vain yksi puu valittiin satunnaisesti. Jokaisen puun jokaiselle noodille (solmulle) on määriteltä leksikostatistinen etäisyys. Mitä suurempi etäisyys, sitä läheisempi on kielten välinen suhde. Tarkempi termi olisi siis ”käänteinen leksikostatistinen etäisyys”, ei ”leksikostatistinen etäisyys”. Yksinkertaisuuden vuoksi käytämme tässä tehtävässä termiä ”leksikostatistinen etäisyys”.

Sekä stabiilisuusindeksit että leksikostatistiset etäisyydet pyöristetään kahteen desimaaliin. Jos kolmas numero desimaalipilkun jälkeen on pienempi kuin 5, arvo pyöristetään alaspäin; muuten se pyöristetään ylöspäin. Esimerkiksi, 2,836 pyöristetään 2,84:ään, 0,705 pyöristetään 0,71:een, ja 0,703 pyöristetään 0,70:een. Pyöristäminen pätee vain arvoihin, jotka näytetään ihmislukijoille. Toisin sanoen, tietokone, joka käyttää algoritmeja, ”näkee” pyöristämättömät arvot.

Huomatkaa, että joidenkin sanojen tiedetään tai epäillään olevan lainattu muista kielistä. Esimerkiksi kadiveun kielen sana **jok:i** ’suola’ on lainattu guaranin kielen sanasta **juki** ja ’iipain kielen (Mesa Granden murteen) sana **ʔa:nj** ’vuosi’ on lainattu espanjan kielen sanasta **’ajno**.

Joissakin tapauksissa sanalistoissa annetaan yhdelle merkitykselle useita synonyymejä, jotka erotetaan toisistaan pilkulla. Esimerkkinä mainittakoon vehvosin kielen ’jalkaa’ tarkoittavat sanat.

Alla olevissa tiedoissa kaikki prefiksit on erotettu ”=”-merkillä ja kaikki suffiksit ”-”-merkillä. Joitain sanoja käytetään vain prefiksien kanssa. Nämä alkavat ”=”-merkillä.

Tiedot on kirjoitettu kansainvälisiä foneettisia aakkosia käyttäen. ^ˈ = pääpaino, _ˈ = sivupaino (pääpainoa heikompi), ː = pitkä äänne, ˘ = ylilyhyt äänne, X[˘]Y = X ja Y äännetään yhtenä ääntenä, ˆ = korkea tooni, ˜ = matala tooni, ˆ = laskeva tooni, ˆ˘ = preglottalisoitu äänne (jota edeltää lyhyt ilma-
virran estyminen kurkussa), ˆ˘ = ejiittiivi (konsonantti, jossa kurkunpää suljetaan ja sitten avataan),

◌̥ = soinniton konsonantti, ◌̃ = nasalisoitu äänne (nenän kautta lausuttava äänne), ◌̥ = narinaääni, ⁿ◌ osoittaa ilmapvirtaa nenän kautta konsonantin edellä, ◌^h = aspiroitu (ilmaa puhaltaen lausuttava) konsonantti, ◌^w = labialisoitu (pyöristetyillä huulilla lausuttava) konsonantti, ◌^j = palatalisoitu (kielen keskiosaa kovaan kitalakeen lähestyen lausuttava) konsonantti. **ɑ, æ, ɛ, ɪ, i, ɔ, ʊ, ʉ, ɐ, ʌ, ɒ, ɘ, ɤ, ɞ, ø** ovat vokaaleja. Muut erikoismerkit ovat konsonantteja.

⚠ Minkään tehtävässä mainitun kielen osaamisesta ei ole apua tehtävää ratkaistaessa.

I. osa. Guaikuru-kielet (Argentiina, Brasilia, Paraguay)

	toba (itämurre)	pilagá	mocoví (Chacon murre)	kadiveu
pilvi	l=ʔok ₁	'lo=ʔok ₁	naweyelek ₂	lol:adi ₃
tuli	nodek ₁	'd=oleʔ ₂	norek ₁	n=ol:edi ₂
kala	njaq ₁	'nijaq ₁	naʎin ₂	nij:ogo-ḏzegi ₃
pää	=qajk ₁	'qajk ₁	=qaik ₁	=ak:ilo ₂
tappaa	=alawat ₁	=a'la:t ₁	=alawat ₁	=el:owadi ₁
kuu	ʔawoɔɔjk ₁	ʔa'woʔɔjk ₁	ʃirajɔ ₂	ep:enaj ₃
nenä	=mik ₁	'mik ₁	=mik ₁	=m:iq:o ₁
suola	towe ₁	ol'ɣek ₂	ʔwe ₁	jok:i ₁
kivi	qaʔ ₁	'qaʔ ₁	qaʔ ₁	wet:iɡa ₂
kieli	=atʃ-aʃat ₁	=a'tʃ-aʃat ₁	=oʔley-aʃan-aʃat ₂	=ok:el:i ₃

	A-algoritmi	B-algoritmi	
manuaalinen	<p>leksikostatistinen etäisyys</p>		Stabiilisuusindeksit: pilvi 0,50 tuli 0,50 kala 0,50 pää 0,75 tappaa 1,00 kuu 0,50 nenä 1,00 suola 0,67 kivi 0,75 kieli 0,50
automaattisoitu			Stabiilisuusindeksit: pilvi 0,50 tuli 0,50 kala 0,75 pää 0,75 tappaa 1,00 kuu 0,50 nenä 1,00 suola 0,25 kivi 0,75 kieli 0,50

II. osa. Nubialainen kielikunta (Egypti, Sudan)

	dongolau	kenuzi	dilling	kadaru	debri	birgid
tappaa	^h bɛ:₁	be:₁	hur₂	wur-i₂	wur-i₂	fila:l-e₁
kuu	u¹n-at-t₁	an-at-ti₁	nɔn-ti₁	nɔn-tu₁	nɔn-to₁	ma:l₂
vesi	¹ess₁	essi₁	ɔti₁	ɔto₁	ɔtu₁	eji₁
antaa	¹tir₁	tir₁	ti₁	ti₁	ti₁	te:-n₁
hyvä	¹sɛrɛ:₁	sere:₁	ken₂	kɛn₂	kɛŋ₂	azze-n₃
tuuli	¹turug₁	turug₁	irf-i₂	irf-o₂	irf-o₂	kurr-i₃
tukka	¹dil-ti₁	si:r₂	tɛl-ti₁	til-tu₁	til-tu₁	ur=dill-e₁
vatsa	¹tu:₁	tu:₁	te-te₂	to₁	to₁	tu:₁
nukkua	¹nɛ:r₁	ne:r₁	jer₁	dwallɛli₂	jer-i₁	ne:r-i₁
aurinko	¹masil₁	masil₁	ɛj₂	aju₂	ɛŋgal-to₃	ʔi:zi₂

	A-algoritmi	B-algoritmi	
manuaalinen			Stabiilisuusindeksit: tappaa 0,50 kuu 0,83 vesi 1,00 antaa 1,00 hyvä 0,50 tuuli 0,50 tukka 0,83 vatsa 0,83 nukkua 0,83 aurinko 0,50
automaattisoitu			Stabiilisuusindeksit: tappaa 0,33 kuu 0,50 vesi 0,50 antaa 0,67 hyvä 0,50 tuuli 0,50 tukka 0,83 vatsa 1,00 nukkua 0,50 aurinko 0,50

- (A) (2 pistettä) \mathfrak{R} lausutaan kuin ranskan r , kielen takaosassa. Mihin Dolgopolsky-luokkaan se kuuluu ja miten saitte sen selville?
- (B) (2 pistettä) Vasemmassa yläkulmassa oleva nubialaisten kielten puu on vain yksi kahdesta mahdollisesta puusta tässä algoritmin ja merkintätyyppin yhdistelmässä. Piirtäkää toinen mahdollinen puu.
- (C) (2 pistettä) Vasemmassa alakulmassa oleva nubialaisten kielten puu on vain yksi kahdesta mahdollisesta puusta tässä algoritmin ja merkintätyyppin yhdistelmässä. Piirtäkää toinen mahdollinen puu.
- (D) (2 pistettä) Leksikostatistinen etäisyys 0,49 (joka on määritetty nubialaisten kielten oikeassa yläkulmassa olevan puun juurelle) on pyöristetty kahteen desimaaliin, kuten jotkut muutkin etäisyydet tässä tehtävässä. Mikä on tarkka etäisyys?

III. osa. Mataguaiolainen kielikunta (Argentiina, Bolivia, Paraguay)

	vitši (Bermejon alajuoksun murre)	vitši (Rivadavian murre)	vehvos	veehnajekki	ijohvaaha	manhui	nivakle (alajuoksun murre)	nivakle (yläjuoksun murre)	maká
tuli	ʔitoχ ₁	ʔitəχ ₁	ʔitah ₁	ʔi:taχ ₁	ʔhwat ₂	ʔeite ₁	ʔitax ₁	ʔitax ₁	feʔt ₂
kala	ʔwahat ₁	wahat ₁	wahat ₁	ʔwa:hat ₁	siʔjus ₋₁	ʃiʔjus ₋₁	saxetʃ ₋₁	saxetʃ ₋₁	sehets ₋₁
jalka	=patʃ _{u1}	=qəɓ ₂	=patʃ _{o1} , =kala ₂	=pa:kʔ _{o1}	=ʔsat ₃	=kaʔla ₂	=φo ₄	=φo ₄	=fʔi ₅
vesi	ʔinot ₁	ʔinət ₁	wah ₂	ʔina:t ₁	ʔiʔnat ₁	ʔaʔnat ₁	jinaʔt ₁	jinaʔt ₁	iweli ₃
antaa	=ʔwen _{u1}	=wen _{u1}	=ʔwen _{o1}	=ʔwen _{o1}	=ʔwehn-a ₂	=ʔhaj ₃ , =ʔwen ₂	=xut ₄	=xut-ej ₄	tis-ix ₅
hyvä	ʔis ₁	ʔis ₁	ʔis ₁	ʔis ₁	ʔes ₁	ʔeis ₁	ʔis ₁	ʔis ₁	t=ejkʔun-ej ₂
tuuli	ʔinwok ^w ₁	ʔinwək ₁	ʔihwok ^w ₁	=ja:ʔ ₂ , =x ^w ox ^w ₃	ʔhlahwu ₄	ʔhlahwu ₄	ʔaβiʔm ₅	ʔaβiʔm ₅	tʔunikʔi ₆
puu	haʔlo ₁	hal ₁	haʔla ₁	haʔla ₁	ʔaʔla ₁	ʔaʔla-k ₁	ʔaʔkxi-juk ₂	jiʔkla ₁	naxka-k ₃
tukka	=ʔwule-j ₁	=wule-j ₁	=ʔwole-j ₁	=ʔwo:le-ç ₁ , hi:lenaχ ₂	=ʔwole ₁	=ʔwole-j ₁	=ʔateʔtʃ ₃	=jeʔs ₄	=ʔewkux-its ₅
tappaa	=lon ₁	=lən ₁	=lan ₁	=la:ŋ ₁	=ʔlaʔan ₁	=ʔlan ₁	=kla ₁	=kla ₁	=lan ₁

	A-algoritmi	B-algoritmi	
manuaalinen			Stabiilisuusindeksit: tuli 0,78 kala 1,00 jalka 0,33 vesi 0,78 antaa 0,44 hyvä 0,89 tuuli 0,33 puu 0,78 tukka 0,67 tappaa 1,00
automaatioitu			Stabiilisuusindeksit: tuli 0,78 kala 0,44 jalka 0,33 vesi 0,56 antaa 0,67 hyvä 0,89 tuuli 0,22 puu 0,67 tukka 0,67 tappaa 1,00

IV. osa. Mongolilainen kielikunta (Kiinan kansantasavalta, Mongolia, Venäjä)

(E) (10 pistettä) Tarkastelkaa seuraavaa sanalista. Laskekaa sekä manuaalisia että automatisoituja merkintöjä vastaavat stabiilisuusindeksit.

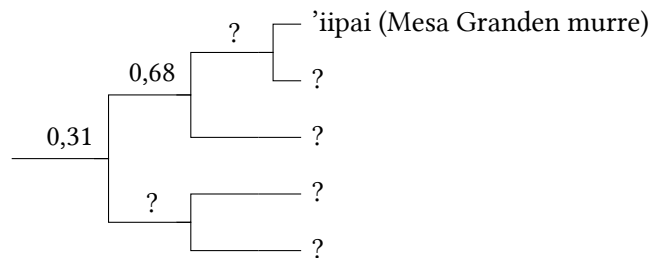
Auttaaksemme teitä, olemme jo laskeneet molemmat stabiilisuusindeksit 'kaikki'-sanalle. Satunnaisessa järjestyksessä nämä ovat 0,36 ja 0,40.

	daguuri (hailarin murre)	hamnigani (mantšun murre)	burjaatti (horin murre)	uus- barguutti	ööldi	hošuutti	kalmukki	halha	ordos	šira- juguuri	baoan
kaikki	hɔ:₁	bɔlt₂	bɔxi:₃	bygd₄	tsug₅	lug₅	tsuk₅, xamak₋₁	pux₃, pugt₄, xamăg₋₁	pyyyte₄, xamukʰ₋₁	tʃʰuq₅	hanə₂
kaarna	hails₁	qalihon₁	χoltɔhɔn₂	xalʰhu:₁	xolts₂	xalis₁	dursn₃	xɔɣtʰɔs₂	turusu₃	χalsən₁	arasun₄
vatsa	ke:li₁	gɔtɔhɔn₂	gedehen₂	gedy:₂	ge:s₂	gets₂	gesn₂	gitis₂, xiwɣij₋₁	ketysy₂	ketesən₂	kele₁
lintu	dəgi₋₁	eiwan₁	ʃubu:n₁	ʃuwu:₁	ʃuvu:₁	ʃuwu:₁	ʃowun₁	ʃuwu₁	ʃuβu:₁	ʃu:n₁, peltʃɔr₂	bendzər₂
tuli	gali₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	qal₁	qal₁	χal₁
tie	terg-u:l₁	qargɔi₂	χargi₂, zam₋₁	zam₋₁	dzam₋₁	dzam₋₁	xa:-lɔə₃	tsam₋₁	tʃam₋₁	mør₄	mor₄
suola	hata:₁	dawhɔn₂	dabhan₂	dawuhu:₂	daws₂	daws₂	dawsn₂	tawsă₂	taβusu₂	ta:psən₂	dabsuŋ₂
uida	unpa-du₁	ɔmba₋₁	tʰamar₋₂	umb₋₁	sele₋₃	umba₋₁	us-tɛi₋₄, ɔ:m₋₅	siɣi₋₃	usu-tʃʰi-la₋₄	umpa₋₁	mba₋₁
vesi	ɔsɔ₁	ɔxɔn₁	uhan₁	u:ha₁	usn₁	us₁	usn₁	ɔsɔ₁	usun₁	qʰusun₁	sə₁
tuuli	kein₁	halkin₂	halxin₂	halxi₂	salʰxin₂	salkʲi₂	salʲkn₂	salɣxi₂	kʰi:₁	kʰi:₁	ki₁

V. osa. Jumalainen kielikunta (Meksiko, Yhdysvallat)

(F) (8 pistettä) Tarkastelkaa seuraavaa sanalista. Alla on puu, joka on rakennettu saman sanalistan perusteella. Jotkut tiedot (kielten nimet ja leksikostatistiset etäisyydet) puuttuvat. Täydentäkää aukot. Ilmoittakaa, onko puu manuaalinen vai automaattinen, sekä onko se luotu algoritmilla A vai B.

	mohave	kokopa	javapai	tiipai (Jamulin murre)	'iipai (Mesa Granden murre)
lyhyt	wena=wen-a ₁	'xʌ=ʔut ₂	'tʃkr=ot-i ₂	lə='ʔuŋ ₁	mə=put-k ₃
lintu	ʔitʃ=i=jer ₁	'ʃa ₂	'ʔ=tʃ=sa ₂	aʔ='ʃa ₂	ʔa:=ʃa:₂
luu	ŋ=a=s=ak ₁	'ŋ=j=a:k ₁	'tʃ=j=a:k-a ₁	'ak ₁	aq ₁
kuiva	i=ro:-v-k ₁	'ʃ=ʔar ₂	'ru-β-i ₁	's=ʔa:j ₃	sa:j ₃
liha	k ^w i:k ^w ay ₁	ʔi='ma:tʃ ₂	'k ^w e:='θo-β-a ₃	'k ^w ak ₄	kuk ^w a:j-p ₁
kaula	maʎaqe ₁	'm=puk ₂	'mlq ₁	i:='puk ₂	i:=puk ₂
nähdä	i=ju:-k ₁	'wi:₂	'ʔu:₁	'wi:w ₂	ə=wu:w ₂
häntä	i:=ʔar ₁	'ʃ=juʎ ₂	'β=hé ₃	ʃə='juʎ ₂	xə=juʎ ₂
kaksi	havik-k ₁	'x=wak ₁	'h ^w ák-i ₁	xə='wak ₁	xə=wak ₁
vuosi	hu:ðe ₁	'mat-'ka:m ₂	'ʔ=tʃ ^h ur-a ₃	mat-'wam ₂	ʔa:n ⁱ -₁



(G) (20 pistettä) On luotu jotkut muut jumalaisten kielten puut, joita vastaavat seuraavat puun juuressa (jokaisesta puusta kauimpana vasemmalla) olevat leksikostatistiset etäisyydet:

1. 0,20
2. 0,23
3. 0,24

Piirtäkää jokainen näistä puista. Ilmoittakaa kunkin puun osalta, onko se manuaalinen vai automaattinen, sekä onko se luotu algoritmilla A vai B.

(H) (3 pistettä) Kaksi (G)-alitehtävässä luetelluista etäisyyksistä on pyöristetty kahteen desimaaliin: 0,23 on pyöristetty 0,225:stä. Mikä muu etäisyys on pyöristetty ja mikä on sen tarkka arvo?

(I) (4 pistettä) Selittäkää, miten stabiilisuusindeksit lasketaan.

(J) (5 pistettä) Selittäkää, miten leksikostatistiset etäisyydet lasketaan.

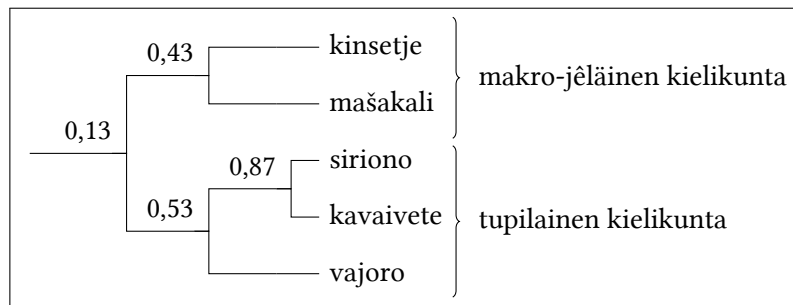
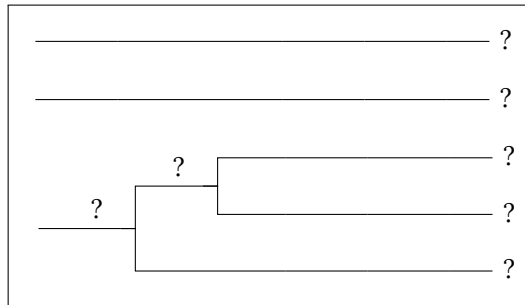
(K) (4 pistettä) Selittäkää A- ja B-algoritmien välinen ero.

VI. osa. Makro-jêläinen kielikunta ja tupilainen kielikunta (Brasilia, Bolivia)

(L) (28 pistettä) Makro-jêläiset ja tupilaiset kielet muodostavat kaksi suurta kielikuntaa, joihin kuuluvia kieliä puhutaan Etelä-Amerikassa. Jotkut kielitieteilijät uskovat niiden olevan kaukaisia sukulaisia. Tarkastelkaa seuraavia sanalistoja.

	A	B	Γ	Δ	E
kaarna	e='e-ke	h ^w i='k ^h Λ	kup='pe	mĩβm='tɛaj	= 'pe
vatsa	'e=rje	= 't ^h igi	=ã'ũn	= 'tɛj	=rɛ'wɛk
veri	e='ruki	=ka' ⁿ bɾo	=d̥z=a'ɛ	= 'hɛβp	=ru'i
polttaa	= 'rai	=rɔ='k ^h ɹɔ	=po'k ^w a	mũ=...='haβp	=ra'pi
rasva	e='kira	= 't ^h wɔmi	= 'd̥z=ap	= 'tuβp	= 'kap
jalka	'e=i	= 'h ^w aji	= 'βi	=po'ta	= 'pi
käsi	'e=o	=nĩ'k ^h ɹa	= 'βo	= 'nĩβm	= 'pɔ
raskas	e='usi	=wi't ^h i	=po'ti	=βp'tɔj	=pɔ'ij
maksa	'e=ja	= 'nba	=pi'a	=tɛiβpkĩ'nɔj	=pi'ʔa
uusi	e='jasu	= 'ndiwi	=pa'gop	= 'tiβp	=pia'u
juuri	e='rao	=ja'ɾe	kup=kujɔ'pe	mĩβm=nĩβm=tɛa'tiə	=ra'pɔ
iho	'e=i	= 'k ^h Λ	= 'pe	= 'tɛaj	= 'pit
häntä	e='rokoi	= 'nbi	=d̥z=o'k ^w aj	=nã:='kiβp	= 'raj
valkoinen	'e=fĩ	=ja'k ^h a	=d̥zi'ra	=βp'douɥ	= 'sĩɲ
siipi	e='heo	=ja'ɾa	=pe'o	=nĩ'mãuɥ	=pe'pɔ, =ji'wa

Alla on kaksi puuta, jotka on rakennettu samojen sanalistojen perusteella. Jotkut tiedot (kielten nimet ja leksikostatistiset etäisyydet) puuttuvat. Täydentäkää aukot. Ilmoittakaa kunkin puun osalta, onko se manuaalinen vai automaattinen, sekä onko se luotu algoritmilla A vai B.



A	B	Γ	Δ	E
?	?	?	?	?

⚠ Manuaaliset merkinnät ja stabiilisuusindeksit jätettiin tarkoituksella pois tästä alitehtävästä.

(M) (10 pistettä) Dolgopolsky-luokkiin perustuvat automaattiset menettelyt voivat tuottaa väärää tuloksia. Tässä esimerkissä automaattinen menettely havaitsee enemmän yhtäläisyyksiä sirionon kielen ja tietyn makro-jäläisen kielen (kinsetjen kielen) välillä kuin sirionon ja muiden tupilaisten kielten välillä. Ehdottakaa muokattua automaattista menettelyä, joka tuottaisi oikean makro-jäläisten ja tupilaisten kielten luokittelun, jos sitä sovellettaisiin yllä oleviin sanalistoihin, ja kuvatkaa tämä *lyhyesti*.

△ Tämä alitehtävä arvostellaan ainoastaan siinä tapauksessa, että parhaat pisteet saavuttaneet joukkueet päätyvät tasapisteisiin.

Kiitokset Alejandra Vidalille, Maria Konoshenkolle, Ilja Gruntoville ja Jamthô Suyälle muutamia kieliä koskeviin kysymyksiin vastaamisesta. —*Andrei Nikulin, Milena Veneva*

Toimittajat: Ivan Deržanski (tekn. toim.), Hugh Dobbs, Stanislav Gurevitš, Boris Iomdin, Liam McKnight, Andrei Nikulin (päätoimittaja), Aleksejs Peguševs, Jan Petr, Aleksandr Piperski, Marija Rubinštein, Milena Veneva, Elysia Warner.

Suomenkielinen teksti: Andrei Nikulin, Seppo Kittilä.

Onnea kilpailuun!