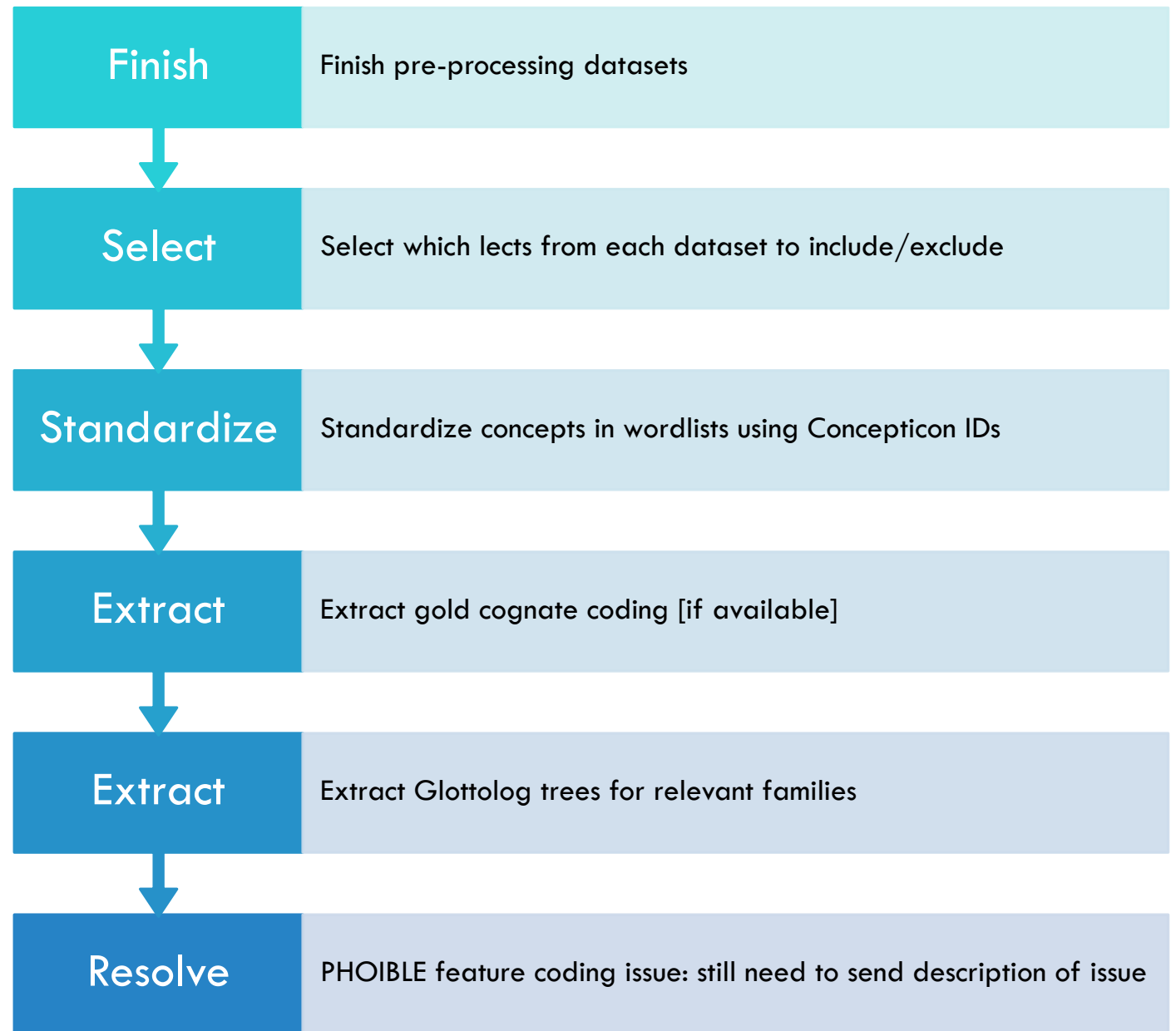




THESIS SEMINAR MEETING

Philip Georgis
June 18, 2021

CURRENT TASKS



NEW STANDARD FORMAT

- Individual language documents → CLDF format dataset document

The image displays four separate windows, each showing a list of words and their phonetic transcriptions in a specific language. The windows are titled 'Fuzhou.txt', 'Belarusian.txt', 'Northern Pashto.txt', and 'Friulian.txt'.

Fuzhou.txt

I/我/ŋwai^{3 2}
all/個郎下/ko^{2 1 2}louŋ^{5 3}ŋa^{2 4 2}
and/共/kɔŋ^{2 4 2}
animal/四肢爬/si^{2 1 2}k^{h a}5⁵βa^{5 3}
ash/灰灰/hwoi^{5 5}hwoi^{5 5}
ash/灰/hwoi^{5 5}
back/背/p^hjaŋ^{5 5}
bad/痞/p^hai^{3 2}
bad/呆/ŋai^{5 3}
bark/树皮/ts^hjeu^{2 1 2}p^hwoi^{5 3}
because/因为/iŋ^{5 5}ŋwoi^{2 4 2}
belly/腹老/pu^{2 4}lo^{3 2}
big/大/twai^{2 4 2}
bird/鳥/tsɛu^{3 2}
bite/咬/ka^{2 4 2}

Belarusian.txt

I/я/ja
all/увесь/uv'ɛsɪ
arm/рука/r'uk'a
ash/попел/p'ɔp'ɛɫ
bark/кара/kar'a
bath/ванна/v'anna
beard/барада/barad'a
belly/жывот/z'ɪv'ɔt
better/лепшы/l'ɛps'ɫɪ
big/вялікі/v'ali'iki
bird/птушка/pt'uška
bite/кусаць/kus'kus
black/чорны/č'orn'ɪ
blood/кроў/kr'ɔw
bone/кось/kɔs'ɫɪ
breast/грудзі/gr'udzi
brother/брат/brat
burn/гарэць/gar'gar
calf/цяля/č'ali'a
chair/стул/stul
chair/крэсла/kresla
child/дзіця/dzi'a
child/рабёнак/rab'ɔnak
cloud/хмара/xm'ara
cloud/воблака/v'ɔblaka
cold/халодны/xalodny
color/колер/k'ɔl'ɛr
come/прыйсці/pr'ijd/pr'ijd
cough/кашаль/k'ašal
cow/карова/kar'ɔva

Northern Pashto.txt

disappear/ورکېدل/wraked'əl
dish/خواره/xwɔɖa
dishes/لوښي/l'ɔʂai
distance/واټن/wat'an
disturb/ماتول/mataw'əl
dive/غوتنه/ɣut'awah'əl
divide/تقسيمول/taqsimaw'əl
divide/وېشل/wej'əl
do/کول/kaw'əl
doctor/طبيب/tab'ib
doctor/داکټر/dakt'ar
dog/سپی/spai
doll/ګودي/gud'əj
door/دروازه/darwaz'a
down/ښکته/škata
drag/وزل/wɖəl
draw/رسمول/rasmaw'əl
dream/خوب/xob
drink/څښل/čsx'əl
drive/بيول/biw'əl
drive/تلل/tləl
drop/قطره/qatr'a
drop/لاسه لويدل/ləl'asawed'əl
drop/شاغکی/š'ašakai
drop, descend/کېدل/ښکته/šk'ataked'əl
dry/وچېدل/wučjed'əl
dry/چا/wəčj
duck/ایلی/il'əj
dust/څا وره/x'awra
dust/ګرد/gard

Friulian.txt

I/jo/jo
all/dut/dut
ash/cinise/çin'ise
bark/scuarce/skw'arçe
belly/panze/p'ançe
big/grant/grant
bird/uciel/uçj'el
bite/muardi/muard/muard
black/neri/n'eri
blood/sanc/sanç
bone/vues/vwes
breast/pet/pet
burn/brusâ/brûs/brus
cloud/nûl/nu:l
cold/frêt/fre:t
come/vigni/vip/vip
day/dî/di
die/murî/mur/mur
dog/cjan/caj
drink/bevi/bev/bev
dry/sec/sek
ear/orele/or'ele
earth/tiere/tj'ere
eat/mangjâ/manç/manç
egg/ûf/u:f
eye/voli/v'oli
father/pari/p'ari
feather/plume/pl'ume
fingernail/ongule/ongul'e
fire/fûc/fu:k

NEW STANDARD FORMAT

- Individual language documents → CLDF format dataset document

ID	Language_ID	Glottocode	ISO 639-3	Parameter_ID	Value	Form	Segments	Source_Form	Cognate_ID	Loan	Comment	Source
1	Vaeakau-Taumako_EIGHT_1	Vaeakau-Taumako	pile1238	piv	EIGHT	valu	valu	valu	EIGHT_1	FALSE		Hovdhaugen-375-2009
2	Wallisian_EIGHT_1	Wallisian	wall1257	wls	EIGHT	valu	valu	valu	EIGHT_1	FALSE	Pollex 06: Valu.	POLLEX
3	Maori_EIGHT_1	Maori	maor1246	mri	EIGHT	waru	waru	waru	EIGHT_1	FALSE		Biggs-85-2005
4	Kapingamarangi_EIGHT_1	Kapingamarangi	kapi1249	kpg	EIGHT	walu	walu	walu	EIGHT_1	FALSE	Pollex 06: Walu.	POLLEX
5	Tahitian_EIGHT_1	Tahitian	tahi1242	tah	EIGHT	vaʔu	vaʔu	vaʔu	EIGHT_1	FALSE		Clark-173-2005
6	Emae_EIGHT_1	Emae	emae1237	mmw	EIGHT	βaru	βaru	βaru	EIGHT_1	FALSE		52375
7	Rapanui_EIGHT_1	Rapanui	rapa1244	rap	EIGHT	vaʔu	vaʔu	vaʔu	EIGHT_1	FALSE		POLLEX
8	Mangareva_EIGHT_1	Mangareva	mang1401	mrv	EIGHT	varu	varu	varu	EIGHT_1	FALSE		POLLEX
9	Luangiua_EIGHT_1	Luangiua	onto1237	ojv	EIGHT	valu	valu	valu	EIGHT_1	FALSE	Pollex 06: Valu.	POLLEX
10	Tongatongalslands_EIGHT_1	Tongan	tong1325	ton	EIGHT	valu	valu	valu	EIGHT_1	FALSE	Pollex 06: Valu.	117207
11	Tikopia_EIGHT_1	Tikopia	tiko1237	tkp	EIGHT	varu	varu	varu	EIGHT_1	FALSE	Pollex 06: Varu.	POLLEX
12	Sikaiana_EIGHT_1	Sikaiana	sika1261	sky	EIGHT	valu	valu	valu	EIGHT_1	FALSE	Pollex 06: Valu.	POLLEX
13	NorthMarquesan_EIGHT_1	North Marquesan	nort2845	mrq	EIGHT	vaʔu	vaʔu	vaʔu	EIGHT_1	FALSE	used on Nuku Hiv	POLLEX
14	EastFutuna_EIGHT_1	East Futuna	east2447	fud	EIGHT	valu	valu	valu	EIGHT_1	FALSE		POLLEX
15	Pukapuka_EIGHT_1	Pukapuka	puka1242	pkp	EIGHT	valu	valu	valu	EIGHT_1	FALSE		Salisbury-152-2005
16	MeleFila_EIGHT_1	Mele-Fila	mele1250	mxe	EIGHT	eβaru	βaru	βaru	EIGHT_1	FALSE		52375
17	AustralA_EIGHT_1	Austral A	aust1304	aut	EIGHT	vaGu	vagu	vagu	EIGHT_1	FALSE		Tamaititahio-1213-2015
18	Tuamotuan_EIGHT_1	Tuamotuan	tuam1242	pmt	EIGHT	varu	varu	varu	EIGHT_1	FALSE	Pollex 06: Varu.	POLLEX
19	Niuean_EIGHT_1	Niuean	niue1239	niu	EIGHT	valu	valu	valu	EIGHT_1	FALSE	Pollex 06: Valu.	POLLEX
20	AustralB_EIGHT_1	Austral B	aust1304	aut	EIGHT	vaʔu	vaʔu	vaʔu	EIGHT_1	FALSE		Meyer-128-2005
21	FutunaAniwa_EIGHT_1	Futuna-Aniwa	futu1245	fut	EIGHT	varu	varu	varu	EIGHT_1	FALSE	Pollex 06: Varu.	POLLEX
22	Hawaiian_EIGHT_1	Hawaiian	hawa1245	haw	EIGHT	walu	walu	walu	EIGHT_1	FALSE	Pollex 06: Walu.	71458
23	Rarotongan_EIGHT_1	Rarotongan	raro1241	rar	EIGHT	varu	varu	varu	EIGHT_1	FALSE	Pollex 06: Varu.	POLLEX
24	Penrhyn_EIGHT_1	Penrhyn	pen1237	pnh	EIGHT	varu	varu	varu	EIGHT_1	FALSE	Pollex 06: Varu.	POLLEX
25	Nukuria_EIGHT_1	Nukuria	nuku1259	nur	EIGHT	varu	varu	varu	EIGHT_1	FALSE		Davletshin-1212-2015
26	Samoa_EIGHT_1	Samoa	samo1305	smo	EIGHT	valu	valu	valu	EIGHT_1	FALSE	Pollex 06: Valu.	Blust-118-2005
27	RennellBellona_EIGHT_1	RennellBellona	renn1242	mnv	EIGHT	bangu	bangu	bangu	EIGHT_1	FALSE	eight	POLLEX
28	Tuvalu_EIGHT_1	Tuvalu	tuva1244	tlv	EIGHT	valu	valu	valu	EIGHT_1	FALSE		29903
29	Anuta_EIGHT_1	Anuta	anut1237	aud	EIGHT	varu	varu	varu	EIGHT_1	FALSE	eight	POLLEX
30	Vaeakau-Taumako_FIFTY_3	Vaeakau-Taumako	pile1238	piv	FIFTY	gatoaelima	gatoaelima	gatoaelima	FIFTY_3	FALSE		Hovdhaugen-375-2009
31	Maori_FIFTY_8	Maori	maor1246	mri	FIFTY	rimatekau	rimatekau	rimatekau	FIFTY_8	FALSE		Biggs-85-2005
32	Mangareva_FIFTY_1	Mangareva	mang1401	mrv	FIFTY	rima rongoʻuru	rima rongoʻuru	rima rongoʻuru	FIFTY_1	FALSE		POLLEX
33	Tongatongalslands_FIFTY_1	Tongan	tong1325	ton	FIFTY	nimangofulu	nimangofulu	nimaŋoʻfulu	FIFTY_1	FALSE		117207
34	EastFutuna_FIFTY_2	East Futuna	east2447	fud	FIFTY	kaulima	kaulima	kaulima	FIFTY_2	FALSE		POLLEX
35	Pukapuka_FIFTY_4	Pukapuka	puka1242	pkp	FIFTY	tinolima	tinolima	tinolima	FIFTY_4	FALSE	50 people	Salisbury-152-2005
36	Pukapuka_FIFTY_2	Pukapuka	puka1242	pkp	FIFTY	laulima	laulima	lau+lima	FIFTY_2	FALSE	traditional	Salisbury-152-2005
37	Pukapuka_FIFTY_1	Pukapuka	puka1242	pkp	FIFTY	limangaulu	limangaulu	limaŋaʻulu	FIFTY_1	TRUE	(Raro.)	Salisbury-152-2005
38	AustralA_FIFTY_9	Austral A	aust1304	aut	FIFTY	paeʔahuGu	paeʔahugu	paeʔaʻhugu	FIFTY_9	TRUE	Likely a Tahitian	Tamaititahio-1213-2015
39	Niuean_FIFTY_5	Niuean	niue1239	niu	FIFTY	lima fiha	limafiha	lima fiha	FIFTY_5	FALSE		POLLEX
40	AustralB_FIFTY_9	Austral B	aust1304	aut	FIFTY	paeʔaʔuru	paeʔaʔuru	paeʔaʻʔuru	FIFTY_9	TRUE	Tahitian loan wit	Meyer-128-2005
41	Hawaiian_FIFTY_6	Hawaiian	hawa1245	haw	FIFTY	kanalima	kanalima	kanalima	FIFTY_6	FALSE		71458
42	Rarotongan_FIFTY_1	Rarotongan	raro1241	rar	FIFTY	rimaŋauru	rimaŋauru	rimaŋaʻuru	FIFTY_1	FALSE		POLLEX
43	Nukuria_FIFTY_7	Nukuria	nuku1259	nur	FIFTY	tipurima	tipurima	tipu rima	FIFTY_7	FALSE		Davletshin-1212-2015

NEW STANDARD FORMAT

- Individual language documents → CLDF format dataset document

ID	Language_ID	Glottocode	ISO 639-3	Parameter_ID	Value	Form	Segments	Source_Form	Cognate_ID	Loan	Comment	Source
Vaeakau-Taumako_EIGHT_1	Vaeakau-Taumako	pile1238	piv	EIGHT	valu	valu	va l u	valu	EIGHT_1	FALSE		Hovdhaugen-375-2009
Wallisian_EIGHT_1	Wallisian	wall1257	wls	EIGHT	valu	valu	va l u	valu	EIGHT_1	FALSE	Pollex 06: Valu. :	POLLEX
Maori_EIGHT_1	Maori	maor1246	mri	EIGHT	waru	waru	wa r u	waru	EIGHT_1	FALSE		Biggs-85-2005
Kapingamarangi_EIGHT_1	Kapingamarangi	kapi1249	kpg	EIGHT	walu	walu	wa l u	walu	EIGHT_1	FALSE	Pollex 06: Walu.	POLLEX
Tahitian_EIGHT_1	Tahitian	tahi1242	tah	EIGHT	va'u	vaʔu	va ʔ u	vaʔu	EIGHT_1	FALSE		Clark-173-2005
Emae_EIGHT_1	Emae	emae1237	mmw	EIGHT	βaru	βaru	β a r u	βaru	EIGHT_1	FALSE		52375
Rapanui_EIGHT_1	Rapanui	rapa1244	rap	EIGHT	va'u	vaʔu	va ʔ u	vaʔu	EIGHT_1	FALSE		POLLEX
Mangareva_EIGHT_1	Mangareva	mang1401	mrv	EIGHT	varu	varu	va r u	varu	EIGHT_1	FALSE		POLLEX

- Uses Concepticon glosses as Parameter IDs
- Cross-indexed with Glottocodes and ISO 639-3 codes
- Contains data about cognate sets, borrowings, sources, and comments (when available) all together to facilitate things later on
- Required revisiting all datasets and preprocessing again → revealed that the first cursory round of preprocessing was insufficient in many cases (esp. Italic & Polynesian)

VOCABULARY/COGNATE INDICES

- Using new format dataset files, vocabulary index files can be easily generated in order to visually examine/compare word forms within cognate sets
- Cognate word forms all in the same row, loanwords enclosed in parentheses

Gloss	Aromanian	Asturian	Catalan (Ca)	Catalan (Ce)	Catalan (M)	Catalan (M)	Catalan (Nc)	Catalan (Va)	Dalmatian	Emiliano (C)	Emiliano (F)	Emiliano (R)	Franco-Pro	French	Friulian	Galician	Istro-Roma	Italian (Fol)	Italian (Gro)	Italian (Sta)	Ladin (Fass)	Ladin (Garc)	Latin (Arch)	Latin (Late)	Ligurian (G)	Ligurian (I)
FULL_1	pl'inũmpl	en'eno	pʔe	pʔe	pʔe	pʔe	pʔe	pʔe	plajn	pi:n	pin	pi:n	pʌ:	plɛ̃	plɛŋ	tʃ'eo	pʌir	p'jenu~p'ir	p'jeno	p'jeno	pjen	pləŋ	pʔ'e:nus	pʔ'e:nus	piŋ	piŋ
GIVE_1	daw	dar							dwor	dɛ:r	dar	dɛ:r			da:	dar	do	da	d'are	d'are	der	dɛ	d'arɛ	d'arɛ	da:	da
GIVE_2			don'ar	dun'a	don'ar	dun'a	don'a	don'ar						dɔne												
GIVE_3													baʎ'i:													
GOOD_1	b'unũ	boŋ	bɔ~bɔn	bɔ~bɔn	bɔ~bɔn	bɔ~bɔn	bɔ~bɔn	bɔ~bɔn	buŋ	bɔwn	bon	bɔwn	bɔ	bɔ	bɔŋ	bo	bur	b'onu	b'ɔno	b'wɔno	boŋ	boŋ	b'ɔnɔs	b'ɔnɔs	buŋ	buŋ
GREEN_1	v'garde	b'erde	vert	bert	bert	vert	bert	bert	(v)jard	verd	verd	vejrd	vɛ	vɛʁ	vert	b'erde		v'erde	v'erde	v'erde	vert	vart	w'irɔɪs	v'irɔɪs	v'erde	v'erde
GREEN_2																	(zel'en)									
HAIR_1			kaβ'eʎ	kaβ'ɛʎ		kaβ'ɛj	kaβ'eʎ	(kaβ'eʎ)	kap'ej	kav'i:	kav'i	kav'i:	(ʃev'ø:)	ʃəvø	cav'ɛj	kaβ'elo		k'apil:u	kap'el:i	kap'el:o	çav'ej	çav'əj	k'apɪt:ɔs	k'apɪt:ɔs	kav'el:i	kav'el:i
HAIR_2	p'erũ	p'elo														p'elo	per									
HAIR_3					m'ɔŋo																					
HAND_1	m'ina	m'ano	ma	ma	ma	ma	ma	ma	mwoŋ	maŋ	man	ma:n	mã	mɛ̃	maŋ	maŋ	mər	(m'ano)	m'ano	m'ano	maŋ	maŋ	m'anɔs	m'anɔs	maŋ	maŋ
HEAD_1	k'apũ	kab'eθa	kap	kap	kap	kap	kap	kap	kup		ko				ca:f	kaβ'eθa	kop	kap'ɔtʃ:a			tʃɛf	tʃɛ	k'apɔt	k'apɔt		
HEAD_2		t'jesta							(t'jasta)	t'esta		t'esta	t'eta	tɛt					t'esta	t'esta					t'esta	t'esta
HEAD_3																										
HEART_1		koraθ'ɔŋ	kɔr	kɔr	kɔr	kɔ	kɔr	kɔr	kwor	ko:r	kɔr	ko:r	kø:	kœʁ	ku:r	koraθ'ɔŋ		k'ɔre	k'ɔre	k'wɔre	ker	kwer	kɔr	kɔr	kø	kø
HEART_2																										
HEART_3	'inima																j'irimæ									

ARABIC DATASET: RATCLIFFE (2020)

- Known issues
 - Poor formatting of dataset (scraped from PDF)
 - Unclear transcription conventions
 - Obscure dialects of unclear classification
 - Cypriot Arabic forms often only as triconsonantal root

1.1 Data

		CA	Mor	Mlt	Cai	Dms	Irq	Skh	AqArb	Cyp	Glf	Ymn	Nig	Bux	Nub
1	ALL	kull	koll	Koll u	kull	kəll	kull	tʃill	kəll, sa:yi: n	kull	kil	kull	tʃat	kullu	kulu
2	ASH	ra ma: d	r ^ʕ m ad ^ʕ	rmi: d	ram a:d ^ʕ *	rama :d	ruma :d	l	r ^ʕ am a:d	rama t r-m- d	rama ad	l	l	ra'matt	ruman
3	BAR K	qirf at	qeʃ r ^ʕ a	ʔofr a	ʔifr*	ʔəʃər	gɪfra	kɪfr	səvi:y e	l	gɪfra	gɪfr	li:he, girfe	l	*kokobo lataka, (girifa)
4	BEL LY	bat ^ʕ n	ker ʃ	zaʔʔ	bat ^ʕ n*	bat ^ʕ ə n	bat ^ʕ i n	bat ^ʕ n , tʃarʃ	dʒoof	patn b-t ^ʕ - n	bat ^ʕ n	bat ^ʕ n , karʃ	kirʃ	batin	batna
5	BIG	kab i:r	kbi r	kibi r	kibii r	kbi:r	tʃbi:r	tʃabi: r	gəbi: r	k-b- r	ʃood	Kabii r	kabi:r	ka'biir	kbir
6	BIR D	t ^ʕ a: ʔir	t ^ʕ ir	ʃasf ur	t ^ʕ ee r*	t ^ʕ eer	t ^ʕ eer	t ^ʕ eer	ku:tʃ ka:ye	2	t ^ʕ eer	t ^ʕ ajr	t ^ʕ e:ra	tayra	ter
7	BITE	ʃad	ʃed	gide	ʃad ^ʕ	ʃad ^ʕ d	ʃað ^ʕ ð	ʃað ^ʕ ð	l	ʃaðð	ʃad ^ʕ d	lugus	ad ^ʕ d ^ʕ a	yaʃazz	adi

ARABIC DATASET: RATCLIFFE (2020)

- New issue
- Cognate coding only given wrt Classical Arabic
- Only certain word forms cognate coded
- Formatting means it would still largely need to be entered manually

1.2 Cognacy evaluations dialects and CA

Note: Ordinary black 0 represents non-cognation. Red 0 indicates non-attestation. In calculating word stability, dialects for which a particular item is unattested are factored out. When this type of pro-rating is used in calculating word stability, the result is also indicated in red.

[illegible]

ARABIC DATASET: RATCLIFFE (2020)

- New issue
- Cognate coding only given wrt Classical Arabic
- Only certain word forms cognate coded
- Formatting means it would still largely need to be entered manually
- Looked into contacting Robert Ratcliffe about questions/issues...
 - He died in 2017, paper and dataset published posthumously in 2020

1.2 Cognacy evaluations dialects and CA

Note: Ordinary black 0 represents non-cognation. Red 0 indicates non-attestation. In calculating word stability, dialects for which a particular item is unattested are factored out. When this type of pro-rating is used in calculating word stability, the result is also indicated in red.

[illegible]

NEW ARABIC DATASET

“Varieties of Arabic Swadesh lists” [Wiktionary]

- <https://en.wiktionary.org/w/index.php?title=Appendix:Varieties of Arabic Swadesh lists&oldid=62250595>

English	Arabic (MSA)	Najdi (Riyadh)	Gulf (Emirati)	North Levantine (Central/Beirut)	Basra (Southern Iraq)	Baghdadi (Central Iraq)	Moslawi (Northern Iraq)	Sanaani (Sana'a)	Dhofari (Dhofar)	Cypriot (Kormakitis)	Egyptian (Delta)	Sudanese (Khartoum)	Chadian (Western Sudanic)	Tunisian (Tunis)	Moroccan (Casablanca)	Moroccan (Northern pre- Hilalian)	Hassaniya (Mauritanian)	Maltese
I	ʔanā	ana	āna	ʔana, ʔana	ʔanah	ʔānī	anā	ʔane(h)	ānā	ana	ana, dana (statements)	ana	ana	ēne, āna	ana, anaya	āna	āna	jien
you (singular)	ʔanta (m), ʔanti(i)	ənta (m), ənti (f)	enta (m), enti (f)	enta (m), ente (f)	atah (m), atī (f)	ente (m), enti (f)	ānta (m), āntē (f)	ant (m), anti (f)	inta (m), inti (f)	int (m), inti (f)	enta (m), enti (f)	inta (m), inti (f)	inta, inte (m), inti (f)	inti	nta, ntaya (m), nti, ntiya (f)	ntina	(ā)nta (m), (ā)ntiyya (f)	int, inti
he	huwa	huww(a)	(ə)hu	huwwe	hūh	huwwa	ōwe, -we	hū	hō	uo	howwa	hū	hū	huwwā	həwá	huwwa	huwwa, hūwa	huwa
we	nahnu, ʔinnā	əhna, (h)ənnā	nehen	naħnā	aħnh	ehne	niħnā	əħnə(h)	naħana	naxni	ehna	(n)ehna	ʔanīna	(a)ħənā	ħna, hnaya	ħna	(n)āħna (m), (n)āħnāti (f)	ahna
you (plural)	ʔantum (m), ʔantunna (f)	əntum (m), ənten (f)	entu (m), enten (f)	ʔento	antm (m), antn (f)	entu (m), enten (f)	āntim	anto	intū (m), intēn (f)	intu	entu	intu	intu(m)	ntūma	ntuma	ntūma	(ā)ntūma (m), ntūmāti (f)	intom
they	hum (m), hunna (f)	hum (m), hənn (f)	(ə)hum (m), (ə)hen (f)	hennē	hamh (m), hanh (f)	humṣa (m), heṇ (f)	hāyyəm, imme	hom (m), hen (f)	hum (m), hən (f)	innen	homma	hum (m), hun (f, rare)	human (m/f), hinna (f, rare)	hūma	huma	hūma	hūma (m), hūmāti (f)	huma
this	hāḍa (m), hāḍihī (f)	(hā)ḍa (m), (hā)ḍi (f)	hāḍa (m), hāḍi (f)	hayda (m), hayde (f)	h'āḡh (m), h'āih (f)	hāḍa (m), hīrātā (f)	hāḍa (m), hāḍī (f)	ḡayyāh (m), tayyih (f)	hāḍī	aḍa (m), aḍī (f)	da(h) (m), di(h) (f)	da (m), de(h) (f)	dā (m), dī (f)	hāḍa (m), hāḍi (f)	hada (m), hadi (f)	hāda (m), hādi (f)	(hā)ḍa (m), (hā)ḍi (f)	dan (m), (f)
that	ḍālika (m), tilka (f), tika (f)	(ha)ḍāk (m), (ha)ḍīk (f)	haḍāk (m), hāḍīk (f)	haydīk (m), hīdek (f)	hḡ'ākḡh (m), hḡī'ch (f)	ḍāka (m), ḍīe'ce (f)	hāḍēk (m), hāḍīk (f)	ḡayyik (m), tayyik (f)	hāḍāk (m), hāḍīk (f)	āḍāk (m), āḍīk (f)	dawwat (m), dayyat (f), (variants) dawwan, dawwak, dawwa, etc.	dāk (m), dek (f)	dāk (m), dīk (f)	hāḍāka (m), hāḍīka (f)	hadak, dak (m), hadik, dik (f)	hadāk (m), hadīk (f)	ḍāk, ḍowk- huwwa (m), ḍīk, ḍek-hiyya (f)	dak (m), (f)
here	hunā	fiḍā, həna	əhni	hown ~ hōn	hn'āh	əhnāna, əhnā	hūnī, hāwn	hāna	hinnī	awnā, annaxula	hena; aho, henaho(wwat) (m), ahe, hənnə(wwat)	hina	hine	hūnī, hnē	hna, hnaya	hna	hūn, hūnāti(yya)	hawn

NEW ARABIC DATASET

“Varieties of Arabic Swadesh lists” [Wiktionary]

- https://en.wiktionary.org/w/index.php?title=Appendix:Varieties_of_Arabic_Swadesh_lists&oldid=62250595
- Contains 16 dialects with sufficient coverage, mostly overlap with the 14 in Ratcliffe’s dataset
- Full word forms for Cypriot Arabic vs. only triconsonantal roots
- Explicit transcription conventions in most cases
- Data (more) easily extractable automatically
 - Scraped table from Wiktionary
 - Preprocessed into proper IPA and removed all annotations
 - Maltese: automatically converted orthography → IPA using transcriptions given at Maltese Wiktionary entries (with Wiktionary parser)

TURKIC DATASET (SAVELYEV & ROBEETS, 2020)

- Extremely messy dataset format, inconsistent annotations and transcriptions
- Some mistakes in transcriptions

	A	B	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	
1	Meaning	Root	Gagauz		KaraKalpak		KarachayBalka		Karaim		Kazakh		Khakas	
2			Standard roman	IPA	Standard roman	IPA	Standard roman	IPA	Standard roman	IPA	Standard roman	IPA	Standard roman	IPA
3	fire (n.)	*o:t			ot	ot	ot	ot	ot	ot	ot	ot	ot	ot
4		*kaŋ-												
5	nose (n.)	*burUn	burnu	burnu	murin	murun	burun	burun	burun	burun	murin	murun	purun	purun
6		*tumčuk												
7		*bar-			? bar	? bar	bar	bar	bar	bar	bar	bar	par	par
8	go (v.)	*jorī- / *jör(i)-												
9		*kejt-	git	git	? ket	? ket			ket	ket	ket	ket		
10	water (n.)	*sib	su	su	suw	suw	suw	suw	suv	suv	su	su	suγ	suγ
11		*agif	a:z	a:z	awiz	awuz	awuz	awuz	avuz	avuz	awiz	awuz	aas	aas
12	mouth (n.)	*aŋak												
13		*jü:f												
14	tongue (n.)	*til	dil	dil	til	til	til	til	til	til	til	təl	təl	təl
15	blood (n.)	*kia:n	kan	kan	қан	qan	қан	qan	kan, (k) қan	kan, (k) qan	қан	qan	қан	қан
16	bone (n.)	*sijök			süyek	syjek	süyek	syjek	sivek, (k) süyek	sivek, (k) syjek	süyek	syjek	söök	söök
17		*kemük	kemik	kemik										
18		*se												
19	2SG pronoun	*sen	sän	sen	sen	sen	sen	sen	sen	sen	sen	sen	sin	sin
20		*tamor			tamir	tamur	tamir	tamur	tamur	tamur	tamir	tamur		
21		*jildif												
22		*tafil = *tasil												
23	root (n.)	*kök	kök	kök										
24		*tup												
25		*ö:f												
26	come (v.)	*gel	gel	gel	kel	kel	kel	kel	kel	kel	kel	kel	kil	kil
27	breast (n.) //	*köküR(ek)	gü:s	gy:s	kökürek	kökyrek	kökürek	kökyrek	kökräk, kekrek, k	kökrek, kekrek, k	kökrek	kökrek	kögös, köksə	kögös,
28	(chest) (n.)	*tiö:ł			? tös	? tös	töš	təf	töš, tes	təf, tes	tös	tös		
29		*kebde			gewde	gewde					kewde	kewde		
30		*jag-	ya:mur	ja:mur	žawin, (žamyr)	zawun, (žamyr)	jaŋur, (Balk.)	jaŋur, (Balk.)	daŋamur	jamur	žanbir, žawin	žanbur, zawun	naŋmür	naŋmu

TURKIC DATASET (SAVELYEV & ROBEETS, 2020)

- Extremely messy dataset format, inconsistent annotations and transcriptions
- Some mistakes in transcriptions
- Fully cleaned and standardized now
- Solutions to some issues:
 - **Turkish:**
 - /ɣ, q/ → /ɯ, k/
 - Corrected any individual vocabulary entries which I noticed were wrong (yɟtynde → ystynde)
 - **Turkmen:**
 - /s, z/ → /θ, ð/
- Checked with NorthEurLex's Turkic and with Turkic dataset at <http://turkic.elegantlexicon.com/> when in doubt for other languages

URALIC DATASETS

- Originally planned on using NorthEuraLex's Uralic dataset (26 languages)
- Problem: doesn't include gold cognate set coding, I'm not familiar enough with Uralic languages to attempt it myself
- Found new Uralic dataset: De Heer et al., Syrjänen [submitted manuscript]
 - <https://zenodo.org/record/4777568#.YMXcH5MzbiA>
 - Includes cognate coding
 - CLDF format
 - Language coverage is more or less the same as NorthEuraLex's Uralic
 - Still need to preprocess transcriptions

COGNATE CODING: ITALIC

- Words already coded into cognate sets, but all words marked as borrowings lumped together into one set, regardless of source (code = -1)
- Manually re-coded the words marked as borrowings (dataset notes provided the source of borrowing in most cases)
- New cognate ID is the negative cognate set ID which it was borrowed from
- Important in case loanwords are not to be excluded from analysis

Language	Word_ID	Gloss	New CognateID	IPA	Notes	Source Notes
Istro-Romanian	IstroRomanian_BURN (SOMETHING)_-1	BURN (SOMETHING)	-5	pal'i		Borrowed from Croatian {paliti} 'to burn'. Distinct from 'orde {ârde} 'to burn'.
Occitan	ProvençalOccitan_BURN (SOMETHING)_-1	BURN (SOMETHING)	-1	bryl'a		There are two expressions for 'to burn': brûl-'a {brula} and krem-'a {crema}.
Istro-Romanian	IstroRomanian_CLOUD_-1	CLOUD	-2	obl'øk		There is no word for 'cloud' either in Kovačec's dictionary or in the text.
Istro-Romanian	IstroRomanian_DOG_-1	DOG	-4	brek		Probably borrowed from Italian {bracco} 'hound' of Germanic origin.
Occitan	ProvençalOccitan_DOG_-1	DOG	-1	tʃiŋ		Borrowed from French {chien} 'dog'.
Catalan (Castelló de la Plana)	Castellóde laPlanaCatalan_EAT_-1	EAT	-2	mendʒ'ar		Borrowed from Old French {mangier} 'to eat'.
Catalan (Central)	CentralCatalan_EAT_-1	EAT	-2	mənʒ'a		Borrowed from Old French {mangier} 'to eat'.
Catalan (Manises)	ManisesCatalan_EAT_-1	EAT	-2	mənʃ'ar		Borrowed from Old French {mangier} 'to eat'.
Catalan (Minorcan)	MinorcanCatalan_EAT_-1	EAT	-2	mənʒ'a		Borrowed from Old French {mangier} 'to eat'.

COGNATE CODING: SLAVIC AND ARABIC DATASETS

- No gold cognate set coding included for NorthEuraLex Slavic or Wiktionary Arabic datasets
- Idea: semi-automatic with manual correction
 - Try using automatic cognate detection prototype, manually correct results
 - Familiar enough with Slavic that I could do this (have done it previously)
 - Could attempt same for Arabic, generally clear when words are related across Arabic dialects; also can use Ratcliffe's cognate coding as partial reference
 - BUT would need to limit the number of concepts if manual, max 100-200 (NorthEuraLex has >900 concepts)

OVERVIEW OF DATASETS

Family	Source Name	Reference	Number of Varieties
Arabic	Varieties of Arabic Swadesh lists	Wiktionary	16
Italic	Global Lexicostatistic Database	Saenko (2016)	58
(Balto-)Slavic	NorthEuraLex	Dellert et al. (2019)	9 (+2 Baltic)
Uralic	Uralic basic vocabulary with cognate and loanword information	De Heer et al.; Syrjänen	27
Polynesian	Polynesian Segmented Data	Walworth (2018)	31
Sinitic	Collection of basic vocabulary words and characteristic dialect words in modern Chinese dialects	Líu et al. (2007)	19
Turkic	Basic vocabulary datasets for the Turkic languages	Savelyev & Robbeets (2020)	31

STATUS OF DATASETS

Dataset	Fully Preprocessed Transcriptions	Concepticon Cross-Reference	Standardized Format	Gold Cognate Sets Included	Extracted Glottolog Tree
Arabic (Ratcliffe)	✓	✓	✗	partially	✗
Arabic (Wiktionary)	✓	✓	✓	✗	✗
Italic	✓	✓	✓	✓	✗
(Balto-)Slavic	✓	✓	✓	✗	✗
Polynesian	✓	✓	✓	✓	✗
Sinitic	✓	✓	✓	✓	✗
Turkic	✓	✓	✓	✓	✗
Uralic (NorthEuraLex)	✓	✓	✓	✗	✗
Uralic (De Heer, et al.; Syrjänen)	✗	✓	✓	✓	✗

STILL TO DO...

