

THESIS SEMINAR MEETING

Philip Georgis
June 7, 2021

OUTLINE OF PROGRESS

Data entry from Arabic dataset

Challenges with Turkic dataset

New dataset: Polynesian languages

Index of all language varieties

New GitHub repository

• Recall: semi-manual data entry/cleaning required due to source formatting

1.1	Data
	Data

		CA	Mor	Mlt	Cai	Dms	Irq	Skh	AqArb	Сур	Glf	Ymn	Nig	Bux	Nub
1	ALL	kull	koll	Koll u	kull	kəll	kull	ţill	kəll, sa:γi: n	kull	kil	kull	t∫at	kullu	kulu
2	ASH	ra ma: d	r'm ad'	rmi: d	ram a:d° *	rama :d	ruma :d	1	r ^c am a:d	rama t r-m- d	rama ad	1	1	ra'matt	ruman
3	BAR K	qirf at	qe∫ r`a	?o∫r a	?i∫r*	?ə∫ər	gi∫ra	ki∫r	səvi:y e	1	gi∫ra	gi∫r	li:he, girfe	1	*kokobo lataka, (girifa)
4	BEL LY	bat' n	ker ∫	za??	bat [°] n*	bat [°] ə n	bat ^r i n	bat`n , t∫ar∫	dзoof	patn b-t ^r - n	bat [§] n	bat [°] n , kar∫	kir∫	batin	batna
5	BIG	kab i:r	kbi r	kibi r	kibii r	kbi:r	tʃbi:r	tʃabi: r	gəbi: r	k-b- r	Sood	Kabii r	kabi:r	ka'biir	kbir
6	BIR D	t ^c a: 7ir	t ^c ir	Sasf ur	t ^c ee r*	t ^s eer	t'eer	t ^c eer	ku:tʃ ka:ye	2	t ^s eer	t ^s ajr	t ^c e:ra	tayra	ter
7	BITE	۲ad	۲ed	gide	۲ad۲	γadγd	۲að٬ð	Sað'ð	1	۲aðð	۲ad⁵d	lugus	adʻdʻa	yaʕazz	adi

• Recall: semi-manual data entry/cleaning required due to source formatting

Ü	l Dat	CA	Mor	Mlt	Cai	Dms	Irq	Skh	AqArb	Сур	Glf	Ymn	Nig	Bux	Nub									
		Ü, (- Cu.	J5		J.K.I.	,	C, p	J		9	Dux	1144									
1	ALL	kull	koll	19.253	kull	kəll	kull	ffill Index	kəll. Gloss	kull	kil CA	kull Mor	tfat Mlt	kullu Cai	kıılıı Dms	 Irg	Skh	AqArb	Сур	Glf	Ymn	Nig	Bux	Nub
				u				- 4	10 HEART		qalb	qelb	?alb	?alb	?alb	qalb	kal°b	fa:d	kilep	galb	galb	galib	qalb	gelba
10	200			0	0		tr.		11 HORN		qarn	qer ^c n	?arn	?arn	?arən	girin	girin		korne	garn	garn*	garn	ſɔx	*gurun
2	ASH	ra	rˁm	rmi:	ram	rama	ruma	- 4	12 I		?ana:	ana	ji:n	ana	?ana	?a:ni	ana	ana	?-n-a	ana	?ana	ana	anaa	'ana
		ma:	ad°	d	a:d°	:d	:d	- 4	13 KILL		qatala	qtel	?atel	mawwit	?atal	kital	katal	qatal	q-t-l	gital	gatal	katal	qa'tal	katul
		d	uu	"	*			- 4	14 KNEE		rukbat	rokba	rkobba	rukba*	rəkbe	rukba	r°uʧbi			rikba	rukbih*	rukuba		*rukuba
		a			×			- 4	15 KNOW		Sarafa	۲r°ef	jaf	Sirif	Səref	Siraf	Sarsaf	Sərəf	γ-r-f	Saraf	diriSaraf	Sirif	?araf	'arufu
																				dara				
3	BAR	qirf	qe∫	?o∫r	?i∫r*	?ə∫ər	gi∫ra		16 LEAF		waraqat	werqa	wer?a	war?a*	wara?a	warga	waraka	waraqa		wriga	waragih*	waraga	varaq	korofo
	K	at	rˁa	a		2,53	2000		17 LIE		raqada	ttmed	mtedd	itmaddad*	tla??ah	tmaddad	ttsallas			tamaddad	ragad*	ragad	daaγa	
									18 LIVER		kabid	kebda	fwi:d	kibda*	kəbəd	kabda			q-s°-b	tfibd	kabdih*	kibde	kabda	*koobda
_	BEL	21			1 - 15	1	1	- 4	19 LONG		t ^c awi:l	t ^c wil	twil	t ^c awi:l	t°awi:l	t°uwi:l	t ^c awi:l	t°-w-l	t°-w-l	t ^c awiil	t ^c awi:l	t ^c awil	ta'viil	towil
4	LY	bat	ker	za??	bat [°]	bat [°] ə	batʻi		0 LOUSE		qamlah	qemla	?amla	?amla*	?amle	gaml ^c a	kaml°		q-m-l	gamla	gamleh		qamla	
	3730,752	n	l		n*	n	n	5	MAN		radzul	r°azel	radzel	r°a:gil	rədʒdʒa:l	ridʒdʒa:l	zalami	radzəl	Intsan	rayyaal	radzdza:l	radʒil	'zaaker	ragi
								1000															ra'ʒul	
5	BIG	kab	kbi	kibi	kibii	kbi:r	tʃbi:r		MANY		kaθi:r	bezzaf	ħafna	kitiir	kti:r	hwa:ja	ʧθi:r	kəti:r	xtir	waad3id	kaθi:r	kati:r	ka'siir	*milan
1				L	, KIDII	KDI.I	95		3 MEAT		laħm	lħam	laħam	laħma	(laħam)	laħam	laħm	laħəm	l-ħ-m	laħam	laħm	la:ma	læħəm	la:m
	DID	i:r	ı	1	Т		* Around and the		MOON		qamar	gamr°a	?amar	?amar	?amar ^c	gumar	kamar ^c		q-m-r	gamar	gamar	gamar	qamar	*ʃa(h)ari
6	BIR D	t°a:	t°ir	۲asf	t'ee	t ^c eer	t'eer		55 MOUN	TAIN	фabal	3bel	muntanja	gabal*	dʒabal	фibal	tsabal	dзabal	3ipel	фibal	dʒabal*	hadzar	dзabal	фebel
		7ir		ur	r*	_			6 MOUTH	1	fam	fomm	ħal?	bu??	təmm	ħaliq	ħalk	təmm	*xankħ-n-	k ħalʤ	fumm*	gaddu:ma	bæb	*kasma
7	BITE	۲ad	۲ed	gide	۲ad۲	\ad'd	Sað'ð										bu:z			buuz				
				3					7 NAME		ism	smija	isem	ism	?əsəm	?isim	ism	əsəm	?-s-m	isim	?ism	isim	?isim	*asma
									8 NECK		γunq	γenq	Yon?	ra?aba*	ra?be	rugba	rakaba		r-q-b	rguba	rugbeh	ragaba		*ragabtu
									9 NEW		dʒadi:d	3did	dʒdida	gidiid	ժʒdi:d	dʒidi:d	tsadi:d	dʒədi:d	3titd3-d-d	dgadiid	dʒadi:d	dʒadi:d	dgadiid	фedid

• Recall: semi-manual data entry/cleaning required due to source formatting

CA Mor Mit Cai Dms Irq Skh AgArb Cyp Gif Ymn Nig Bux Nub ASH/rama:d'rama:d's BARK/rijr*/rijr BARK/rijr*/rijr BARK/rijr*/rijr BARK/rijr*/rijr BARK/rijr*/rijr BARK/rijr*/rijr BARK/rijr*/rijr* BARK/rijr*/rijr*/rijr*/rijr* BARK/rijr*/rijr*/rijr*/rijr* BARK/rijr*/r	1.	1 Dat	a															Calitat
ALL				Mor	Mlt	Cai	Dms	Irq	Skh	AqArb	Сур	Glf	Ymn	Nig	Bux	Nub		
A0 HEART qalb qelb 7alb 7alb 7alb qalb qarn q	1	ALL	kull	koll	10.253	kull	kəll	kull									Irq	BELLY/bat'n*/bat'n
2 ASH ra rm ram rama ruma d2 2 2ans ana an					u					10 HEAR	Т	qalb	qelb	?alb	?alb		-	
2	1					0	10	the state of the s		11 HORN		qarn	qer ^c n	?arn	?arn	?arən		
Max	2	ASH	ra	rˁm	rmi:	ram	rama	ruma		12 I		?ana:	ana	ji:n	ana	?ana		
d			ma:	ad°	d	a:d°	:d			13 KILL		qatala	qtel	?atel	mawwit	?atal		
45 KNOW Sarafa Sree jaf Sirif Saref Sree Siraf BREAST/sidir*/sidir BURN/hara? CLOW/sāhaab/sāha: b COLOW/sāhaab/sāha: b COLOW/sāhaab						*	555	1.50		44 KNEE		rukbat	rokba	rkobba	rukba*	rəkbe	rukb	kb BONE/Sadma/Sadma
3 SAR qirf qef ra a ra a qif qef ra a qif qef at ra a qif qef at ra a qif qef			u			*				15 KNOV	V	Sarafa	۲r°ef	jaf	Sirif	۲əref	۲iraf	raf BREAST/sidir*/sidir
A	_	\sqcup																BURN/ħara?/ħara?
4 8 LIVER kabid kebda fwi:d kibda* kabad kabd COLD/bard/bard COME/ga/ga 4 8 LIVER kabid kebda fwi:d kibda* kabad kabd COLD/bard/bard COME/ga/ga 5 0 LOUSE qamlah qemla ?amla ?amla ?amla ?amla ?amla ?amla a pamla* ?amla pamla* ?amla pamla* ?amla pamla* ?amla pamla* ?amla* ?amla* pamla* ?amla* pamla* ?amla* pamla* ?amla* pamla* ?amla* pamla* pa	3	BAR	qirf	qe∫	?o∫r	?i∫r*	?ə∫ər	gi∫ra	4	16 LEAF		waraqat	werqa	wer?a	war?a*	wara?a		
4 BEL Dat' Ker Za?? Dat' Dat's		\ \	at	rˁa	a		0.000.00	20000000	4	17 LIE		-		mtedd		tla??ah		
Set									_			100041000000	100000000000000000000000000000000000000					
N	1	BEL	hats	kor	7222	hats	hat's	hatsi	_				t ^c wil					
S BIG Kab Kibi Kibii Kibi	4		Dat	Ker	zarr	Secretary Section	Section 1986	000100000000	5	0 LOUS	E	qamlah	qemla	?amla	?amla*	150000000000000000000000000000000000000	gam	DIE/maat/ma:t
S BIG kab kbi kibi kibii kibii kbi:r tfbi:r 52 MANY ka6i:r bezzaf ħafna kitiir kti:r hwa. 53 MEAT laħm lham laħam laħam laħam (laħam) laħam laħ			n	J		n*	n	n	5	MAN		radzul	r°aʒel	radzel	r°a:gil	rədʒdʒa:l	ridge	bd DOG/Kalb/Kalb
S NAME SAN KAD															19.2222222			
i:r r r r r r r r r r r r r r r r r r r	5	BIG	kab	kbi	kibi	kibii	kbi:r	tʃbi:r	_				100.00					
6 BIR Care the the street of t				r	r	r		,	5			laħm						EARTH RIPT (45 (45
7 BITE Sad Sed gide Sad' Sad'd Sað's 56 MOUTH fam fomm hal? bu?? təmm haliq 57 NAME ism smija isem ism ?əsəm ?isin 58 NECK Sunq Senq Son? ra?aba* ra?be rugb 59 NEW dadi:d 3did dadid gidiid dadi:d dadi:d dadi:d fish-samak/samak		BIR			· ·		24	.5	_				-					F. F. F. J. J. J. J.
7 BITE Sad Sed gide Sad' Sad'd Sað'ð 57 NAME ism smija isem ism Sad's rashe rashe rugb 58 NECK Sunq Senq Sons rasaba* rashe rugb 59 NEW dadi:d 3did dadid gidiid dadi:d dadi:d sadid sadi	6			t'ir		t'ee	t'eer	t eer				-	-	-	0	1	_	FCC / h = - 45 - / h = - 45 -
7 BITE Sad Sed gide Sad' Sad'd Sad	Ļ	Щ	/ir		ur	r*			-	e Mon.	ГН	fam	fomm	ħal?	bu??	təmm	ħalio	
58 NECK Sunq Senq Son? ra?aba* ra?be rugb 59 NEW dgadi:d gdiid dgdid gidiid dgdi:d dgdi:d dgdi:d fISH/samak/samak FEATHER/ri:sa*/ri:sa FIRE/na:r/na:r FISH/samak/samak	7	BITE	۲ad	۲ed	gide	۲ad۲	\ad`d	۲að٬ð			_						21 -	EAT N / (-b-//-b-
59 NEW dadi:d 3did dadida gidiid dadi:d dadi:d fIRE/na:r/na:r FIRE/na:r/na:r FISH/samak/samak		1/0							_			1000					10000	FEATHER/ri:[a*/ri:[a
S9 NEW Gadi:d 3did Gadid Gadid Gadid Gadid FISH/samak/samak									_								_	FTRE/parr/parr
									_ 5	9 NEW		dʒadi:d	3did	dʒdida	gidiid	dʒdi:d	dʒidi	

Cai.txt

- Will likely need to exclude some of the (smaller) varieties from comparison, e.g.
 - Cypriot Arabic
 - Much data missing, most only in form of triconsonantal roots vs. full word forms in other lects
 - Soukhne Arabic (rural Syrian dialect from between Damascus and Baghdad)
 - Cannot find a (unique) Glottolog entry
 - Unclear whether it is a Levantine or Eastern Arabic variety
 - Found ASJP wordlist for "Soukhne Syrian Arabic", but cross-referenced same as Iraqi Arabic (https://asjp.clld.org/languages/SOUKHNE_SYRIAN_ARABIC)

TURKIC DATASET: CHALLENGES

Mostly pre-processed and loaded the Turkic dataset (Savelyev & Robeets, 2020)

Α	В	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	
M	Dood	Gagauz		KaraKalpak		KarachayBalka	ı	Karaim		Kazakh		Khakas	
Meaning	Root	Standard rom	ar IPA	Standard roma	Standard romar IPA		r IPA	Standard romar IPA		Standard rom	ar IPA	Standard rom	ar IPA
fire (n.)	*o:t			ot	ot	ot	ot	ot	ot	ot	ot	ot	ot
1-10/01 117	*kaŋ-												
nose (n.)	*burUn	burnu	burnu	murïn	murum	burun	burun	burun	burun	murïn	murun	purun	pur
	*tumčuk												
	*bar-			? bar	? bar	bar	bar	bar	bar	bar	bar	par	par
go (v.)	*jorï- / *jör(i)-												
0.00 (0	*kejt-	git	git	? ket	? ket			ket	ket	ket	ket		
water (n.)	*sib	su	su	suw	suw	suw	suw	suv	suv	su	su	suy	suy
0.00 0.000 0.000	*agiŕ	a:z	a:z	awïz	awwz	awuz	awuz	avuz	avuz	awïz	awuz	aas	aas
mouth (n.)	*ańak												
	*jü:ŕ												
tongue (n.)	*tII	dil	dil	til	til	til	til	til	til	tĭl	let	let	let
blood (n.)	*kia:n	kan	kan	ķan	qan	ķan	qan	kan, (k) kan	kan, (k) qan	ķan	qan	χan	χan
L ()	*siŋök			süyek	syjek	süyek	syjek	sivek, (k) süyek	sivek, (k) syjek	süyek	syjek	söök	søø
bone (n.)	*kemük	kemik	kemik										
266	*sę												
2SG pronoun	*sen	sän	sen	sen	sen	sen	sen	sen	sen	sen	sen	sin	sin
	*tamor			tamïr	tamur	tamïr	tamur	tamur	tamur	tamïr	tamur		
	*jïldïŕ												
	*taŕil = *tasil												
root (n.)	*kök	kök	køk										
	*tüp												
	*ö:ŕ												
come (v.)	*gel	gel	gel	kel	kel	kel	kel	kel	kel	kel	kel	kil	kil
9 19 19 19 19 19 19 19 19 19 19 19 19 19	*köküR(ek)	gü:s	gy:s	kökürek	køkyrek	kökürek	køkyrek		k køkrεk, kekrek,		køkrek	kögəs, köksə	køg
breast (n.) //	*tiö:Í			? tös	? tøs	töš	tø∫	töš, tes	tøſ, tes	tös	tøs		
(chest) (n.)	*kebde			gewde	gewde					kewde	kewde		
	*jag-	ya:mur	ja:mur	žawin, (žamyir)	зажшп, (затуш	r ĭanur. (Balk.) ĭa	w dzanur. (Balk.)	d vamvur	jamyur	žaŋbïr, žawïn	zaŋbuir, zawuin	nanmïr	naŋı

TURKIC DATASET: CHALLENGES

- Possible issues with some of the IPA transcriptions
 - Orthographic symbols sometimes used in place of IPA characters
 - → Possible to replace wherever clear what the intended IPA symbols were

e.g.
$$\langle \tilde{s} \rangle = /J/$$
, $\langle \tilde{w} \rangle = /w/$, or Chuvash $\langle \tilde{o} \rangle = /\partial^w/$

...but some are less clear, and resources may not exist to disambiguate.

- Turkish transcriptions call into question the validity/quality of transcriptions
 - Transcribes phonemes which do not exist in standard Turkish, e.g. /q/, /y/
 - Some words are mistaken: e.g. *üštünde /yʃtynde/ 'above'
 - > Maybe some non-standard dialect of Turkish?
 - Could correct the issues in Turkish, but I don't know enough about other varieties to check/correct them
- Annotations within transcriptions: need to spend more time to sort these out

TURKIC DATASET: CHALLENGES

- Could alternatively use a different dataset for Turkic:
 - NorthEuraLex
 - Only includes 8 Turkic languages
 - Turkish, Azeri, Uzbek, Kazakh, Bashkir, Tatar, Sakha, Chuvash
 - "Culture Words for Turkic" from Diachronic Atlas of Comparative Linguistics (DiACL)
 - Only includes 9 Turkic languages
 - Turkish, Gagauz, Turkmen, Azeri, Crimean Tatar, Kazakh, Kumyk, Tuvan, Southern Altai
 - Not a Swadesh list
 - ... vs. 31 Turkic languages in Savelyev & Robeets (2020) dataset, which also includes cognate coding.
 - → would strongly prefer to find solution with this dataset

POLYNESIAN DATASET (WALWORTH, 2018)

- One final dataset including 30 modern Polynesian languages + Proto-Polynesian reconstructed forms
- Source is well formatted/machine readable, no issues in IPA transcriptions, includes Glottolog and ISO references, cognate coding
- No major issues found → fully preprocessed

OVERVIEW OF DATASETS

Family	Source Name	Reference	Number of Varieties
Arabic	Glottometrics of Arabic	Ratcliffe (2020)	14
Italic	Global Lexicostatistic Database	Saenko (2016)	58
(Balto-)Slavic	NorthEuraLex	Dellert et al. (2019)	9 (+2 Baltic)
Uralic	NorthEuraLex	Dellert et al. (2019)	26
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

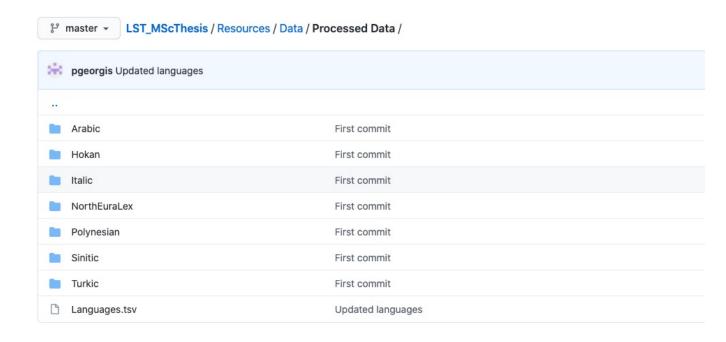
INDEX OF LANGUAGE VARIETIES

- Cross-referenced all varieties from datasets with Glottolog and ISO codes in order to be able to compare results with gold standard trees
- Some varieties could not be found \rightarrow exclude from study (e.g. specific Italic dialects)

Dataset	Family	Name	Source Name	Glottolog Name	Glottocode	ISO 639-3	ISO 639-3 Name	Reference	Notes
1 NorthEuraLex	Uralic, Finnic	Estonian	Estonian	Estonian	esto1258	ekk	Standard Estonian	Dellert et al., 2019	
2 NorthEuraLex	Uralic, Finnic	Finnish	Finnish	Finnish	finn1318	fin	Finnish	Dellert et al., 2019	
3 NorthEuraLex	Uralic, Finnic	Livonian	Livonian	Liv	livv1244	liv	Liv	Dellert et al., 2019	
4 NorthEuraLex	Uralic, Finnic	North Karelian	North Karelian	Karelian	kare1335	krl	Karelian	Dellert et al., 2019	
NorthEuraLex	Uralic, Finnic	Olonets Karelian	Olonets Karelian	Livvi	livv1243	olo	Livvi	Dellert et al., 2019	
NorthEuraLex	Uralic, Finnic	Veps	Veps	Veps	veps1250	vep	Veps	Dellert et al., 2019	
NorthEuraLex	Uralic, Hungarian	Hungarian	Hungarian	Hungarian	hung1274	hun	Hungarian	Dellert et al., 2019	
NorthEuraLex	Uralic, Khantyic	Northern Khanty	Northern Khanty	Kazym-Berezover-Suryskarer Khanty	khan1273	kca	Khanty	Dellert et al., 2019	
NorthEuraLex	Uralic, Mansic	Northern Mansi	Northern Mansi	Northern Mansi	mans1258	mns	Mansi	Dellert et al., 2019	
NorthEuraLex	Uralic, Mari	Hill Mari	Hill Mari	Western Mari	west2392	mrj	Western Mari	Dellert et al., 2019	
NorthEuraLex	Uralic, Mari	Meadow Mari	Meadow Mari	Eastern Mari	east2328	mhr	Eastern Mari	Dellert et al., 2019	
NorthEuraLex	Uralic, Mordvin	Erzya	Erzya	Erzya	erzy1239	myv	Erzya	Dellert et al., 2019	
NorthEuraLex	Uralic, Mordvin	Moksha	Moksha	Moksha	moks1248	mdf	Moksha	Dellert et al., 2019	
NorthEuraLex	Uralic, Permian	Komi-Permyak	Komi-Permyak	Komi-Permyak	komi1269	koi	Komi-Permyak	Dellert et al., 2019	
NorthEuraLex	Uralic, Permian	Komi-Zyrian	Komi-Zyrian	Komi-Zyrian	komi1268	kpv	Komi-Zyrian	Dellert et al., 2019	
NorthEuraLex	Uralic, Permian	Udmurt	Udmurt	Udmurt	udmu1245	udm	Udmurt	Dellert et al., 2019	
7 NorthEuraLex	Uralic, Saami	Inari Sami	Inari Sami	Inari Saami	inar1241	smn	Inari Sami	Dellert et al., 2019	
NorthEuraLex	Uralic, Saami	Kildin Sami	Kildin Sami	Kildin Saami	kild1236	sjd	Kildin Sami	Dellert et al., 2019	
NorthEuraLex	Uralic, Saami	Lule Sami	Lule Sami	Lule Saami	lule1254	smj	Lule Sami	Dellert et al., 2019	
NorthEuraLex	Uralic, Saami	Northern Sami	Northern Sami	North Saami	nort2671	sme	Northern Sami	Dellert et al., 2019	
NorthEuraLex	Uralic, Saami	Skolt Sami	Skolt Sami	Skolt Saami	skol1241	sms	Skolt Sami	Dellert et al., 2019	
NorthEuraLex	Uralic, Saami	Southern Sami	Southern Sami	South Saami	sout2674	sma	Southern Sami	Dellert et al., 2019	
NorthEuraLex	Uralic, Samoyedic	Forest Enets	Forest Enets	Forest Enets	fore1265	enf	Forest Enets	Dellert et al., 2019	
NorthEuraLex	Uralic, Samoyedic	Nganasan	Nganasan	Nganasan	ngan1291	nio	Nganasan	Dellert et al., 2019	
NorthEuraLex	Uralic, Samoyedic	Northern Selkup	Northern Selkup	Selkup	selk1253	sel	Selkup	Dellert et al., 2019	
6 NorthEuraLex	Uralic, Samoyedic	Tundra Nenets	Tundra Nenets	Tundra Nenets	nene1249	yrk	Nenets	Dellert et al., 2019	
7 NorthEuraLex	Yeniseian, Northern Yeniseian	Ket	Ket	Ket	kett1243	ket	Ket	Dellert et al., 2019	
8 NorthEuraLex	Yukaghir, Kolymic	Southern Yukaghir	Southern Yukaghir	Southern Yukaghir	sout2750	yux	Southern Yukaghir	Dellert et al., 2019	
NorthEuraLex	Yukaghir, Northern Yukaghir	Northern Yukaghir	Northern Yukaghir	Northern Yukaghir	nort2745	ykg	Northern Yukaghir	Dellert et al., 2019	
Polynesian	Austronesian, Polynesian	Proto-Polynesian	Polynesian	Polynesian	poly1242			Walworth, 2018	
Polynesian	Austronesian, Polynesian, Anuta	Anuta	Anuta	Anuta	anut1237	aud	Anuta	Walworth, 2018	
Polynesian	Austronesian, Polynesian, Carolinean Outlier Polynesian	Kapingamarangi	Kapingamarangi	Kapingamarangi	kapi1249	kpg	Kapingamarangi	Walworth, 2018	
3 Polynesian	Austronesian, Polynesian, East Futuna	East Futuna	EastFutuna	East Futuna	east2447	fud	East Futuna	Walworth, 2018	
4 Polynesian	Austronesian, Polynesian, East Polynesian	Austral A	AustralA	Austral	aust1304	aut	Austral	Walworth, 2018	
5 Polynesian	Austronesian, Polynesian, Fast Polynesian	Austral B	AustralB	Austral	aust1304	aut	Austral	Walworth, 2018	

NEW GITHUB REPOSITORY

- Created new GitHub repository under my personal account: pgeorgis/LST_MScThesis
- Currently private, but I can grant access
- Will use to store update presentations, processed data, future scripts, etc.



NEXT STEPS

