

ISO	Name	Train char	Dev char	Eval char	Typ / Tok	ISO	Name	Train char	Dev char	Eval char	Typ / Tok
<i>acu</i>	Achuar	1149777	136773	119849	5.902 <sup>-05</sup>	<i>kbh</i>	Camsa	1373946	144068	140039	5.307 <sup>-05</sup>
<i>afz</i>	Afrikaans	3229549	413064	437532	1.985 <sup>-05</sup>	<i>kek</i>	Q'eqchi'	4375494	525831	517455	1.735 <sup>-05</sup>
<i>agr</i>	Aguaruna	991098	118726	103237	6.348 <sup>-05</sup>	<i>lat</i>	Latin	2700731	325553	342365	1.514 <sup>-05</sup>
<i>ake</i>	Akawaio	960849	113905	111793	5.141 <sup>-05</sup>	<i>lav</i>	Latvian	644923	77572	78263	1.161 <sup>-04</sup>
<i>alb</i>	Albanian	3152312	402399	427612	1.808 <sup>-05</sup>	<i>lit</i>	Lithuanian	2531703	313391	309237	2.536 <sup>-05</sup>
<i>amu</i>	Amuzgo	1156128	142241	132194	5.662 <sup>-05</sup>	<i>mam</i>	Mam	1053107	119781	112869	6.222 <sup>-05</sup>
<i>bsn</i>	Barasana	1397953	171482	162042	4.736 <sup>-05</sup>	<i>mri</i>	Maori	3504361	437499	456364	1.501 <sup>-05</sup>
<i>cak</i>	Cakchiquel	1404839	169031	161609	4.033 <sup>-05</sup>	<i>nhg</i>	Nahuatl	1126416	135355	126213	5.548 <sup>-05</sup>
<i>ceb</i>	Cebuano	3985326	509809	536615	1.471 <sup>-05</sup>	<i>nld</i>	Dutch	3224058	392079	416432	2.033 <sup>-05</sup>
<i>ces</i>	Czech	2756308	349505	371027	2.675 <sup>-05</sup>	<i>nor</i>	Norwegian	2941245	374161	392574	2.508 <sup>-05</sup>
<i>cha</i>	Chamorro	641087	66469	67935	1.032 <sup>-04</sup>	<i>pck</i>	Paite	3174091	404462	401042	1.784 <sup>-05</sup>
<i>chq</i>	Chinantec	1548993	174921	164087	4.502 <sup>-05</sup>	<i>plt</i>	Malagasy	3744462	468678	477671	1.705 <sup>-05</sup>
<i>cjp</i>	Cabecar	856441	100035	97246	8.256 <sup>-05</sup>	<i>pol</i>	Polish	2963005	374471	398263	2.088 <sup>-05</sup>
<i>cni</i>	Campa	1149737	133104	120600	3.990 <sup>-05</sup>	<i>por</i>	Portuguese	3010541	380551	404559	2.450 <sup>-05</sup>
<i>dan</i>	Danish	2774922	364278	385352	2.298 <sup>-05</sup>	<i>pot</i>	Potawatomi	212243	25336	24020	1.911 <sup>-04</sup>
<i>deu</i>	German	3195266	391700	417235	2.023 <sup>-05</sup>	<i>ppk</i>	Uma	1050858	115947	110127	5.090 <sup>-05</sup>
<i>dik</i>	Dinka	716411	84429	81572	6.800 <sup>-05</sup>	<i>quc</i>	K'iche'	1153252	131623	127281	5.382 <sup>-05</sup>
<i>dje</i>	Zarma	3126629	372921	405747	1.767 <sup>-05</sup>	<i>quw</i>	Quichua	792834	93791	90930	8.593 <sup>-05</sup>
<i>djk</i>	Aukan	1083303	129770	124682	5.083 <sup>-05</sup>	<i>rom</i>	Romani	818094	91328	89036	7.912 <sup>-05</sup>
<i>dop</i>	Lukpa	864347	96068	95094	6.442 <sup>-05</sup>	<i>ron</i>	Romanian	3107966	394280	419241	1.760 <sup>-05</sup>
<i>eng</i>	English	3238389	418697	450324	1.509 <sup>-05</sup>	<i>shh</i>	Tachelhit	738833	82470	79965	5.659 <sup>-05</sup>
<i>epo</i>	Esperanto	3029361	391874	409980	1.514 <sup>-05</sup>	<i>slk</i>	Slovak	2821379	361684	385812	2.662 <sup>-05</sup>
<i>est</i>	Estonian	778681	177680	5329	8.007 <sup>-05</sup>	<i>slv</i>	Slovene	2883854	369397	381124	2.366 <sup>-05</sup>
<i>eus</i>	Basque	801072	94904	93712	7.073 <sup>-05</sup>	<i>sna</i>	Shona	3013970	384548	407708	2.154 <sup>-05</sup>
<i>ewe</i>	Ewe	807990	97081	94738	9.315 <sup>-05</sup>	<i>som</i>	Somali	3750398	468849	498051	1.314 <sup>-05</sup>
<i>fin</i>	Finnish	3195802	402386	426808	1.789 <sup>-05</sup>	<i>spa</i>	Spanish	3082800	388079	410641	2.190 <sup>-05</sup>
<i>fra</i>	French	3246315	404464	435820	2.129 <sup>-05</sup>	<i>srp</i>	Serbian	2503088	319878	341496	2.402 <sup>-05</sup>
<i>gbi</i>	Galela	1347199	144215	129606	4.442 <sup>-05</sup>	<i>ssw</i>	Swahili	763781	84855	82704	6.979 <sup>-05</sup>
<i>gla</i>	Gaelic	68110	6802	7167	5.848 <sup>-04</sup>	<i>swe</i>	Swedish	3206128	403712	427604	1.932 <sup>-05</sup>
<i>glv</i>	Manx	392897	58690	48239	1.360 <sup>-04</sup>	<i>tgl</i>	Tagalog	3848347	477839	505730	1.283 <sup>-05</sup>
<i>hat</i>	Creole	3332162	400606	440487	1.989 <sup>-05</sup>	<i>tmh</i>	Tuareg	270666	30146	31636	2.436 <sup>-04</sup>
<i>hrv</i>	Croatian	2594494	340781	359694	2.549 <sup>-05</sup>	<i>tur</i>	Turkish	2719803	318179	345666	2.601 <sup>-05</sup>
<i>hun</i>	Hungarian	3020721	376738	408697	2.391 <sup>-05</sup>	<i>usp</i>	Uspanteco	1134539	131631	125891	5.747 <sup>-05</sup>
<i>ind</i>	Indonesian	3528757	405822	454277	1.823 <sup>-05</sup>	<i>vie</i>	Vietnamese	3194226	379697	417496	4.635 <sup>-05</sup>
<i>isl</i>	Icelandic	2968652	376562	389091	2.517 <sup>-05</sup>	<i>wal</i>	Wolaytta	837506	98141	96951	6.004 <sup>-05</sup>
<i>ita</i>	Italian	2979890	388587	409629	2.091 <sup>-05</sup>	<i>wol</i>	Wolof	683480	80261	77575	8.201 <sup>-05</sup>
<i>jak</i>	Jakalteko	1116611	131793	122853	5.615 <sup>-05</sup>	<i>xho</i>	Xhosa	3005476	377342	399338	1.692 <sup>-05</sup>
<i>jiv</i>	Shuar	888886	98309	97483	5.624 <sup>-05</sup>	<i>zul</i>	Zulu	690644	80944	77975	8.946 <sup>-05</sup>
<i>kab</i>	Kabyle	798503	91677	87964	7.770 <sup>-05</sup>						

Table 1: Stats for all languages in our corpus: for each ISO 639-3 code, we list (from left to right) the language name, the character count (in the train, development, and evaluation sets), and the type-to-token ratio.

## Appendix A: list of ISO 639-3 codes and language names

In Tab. 1 the ISO 639-3 codes for each language are associated with the corresponding language name. In addition to this information, Tab. 1 provides the total count of characters for the three data splits, and the type-to-token ratio.

## Appendix B: Typological Features

The 245 binarized typological features from Littell et al. (2017) that characterize the general properties of each language are plotted as a heat map in Fig. 1. Features are related to syntax if their name starts with *S*, to phonology if it starts with *P*, and to phonemic inventories if it starts with *INV*. Note how some values are so rare that they belong exclusively to a single language in the sample, e.g. the vowel /ə/ for Thai.

## References

Patrick Littell, David R. Mortensen, Ke Lin, Katherine Kairis, Carlisle Turner, and Lori Levin. 2017. *URIEL and lang2vec: Representing languages as typological, geographical, and phylogenetic vectors*. In *Proceedings of EACL*, pages 8–14.

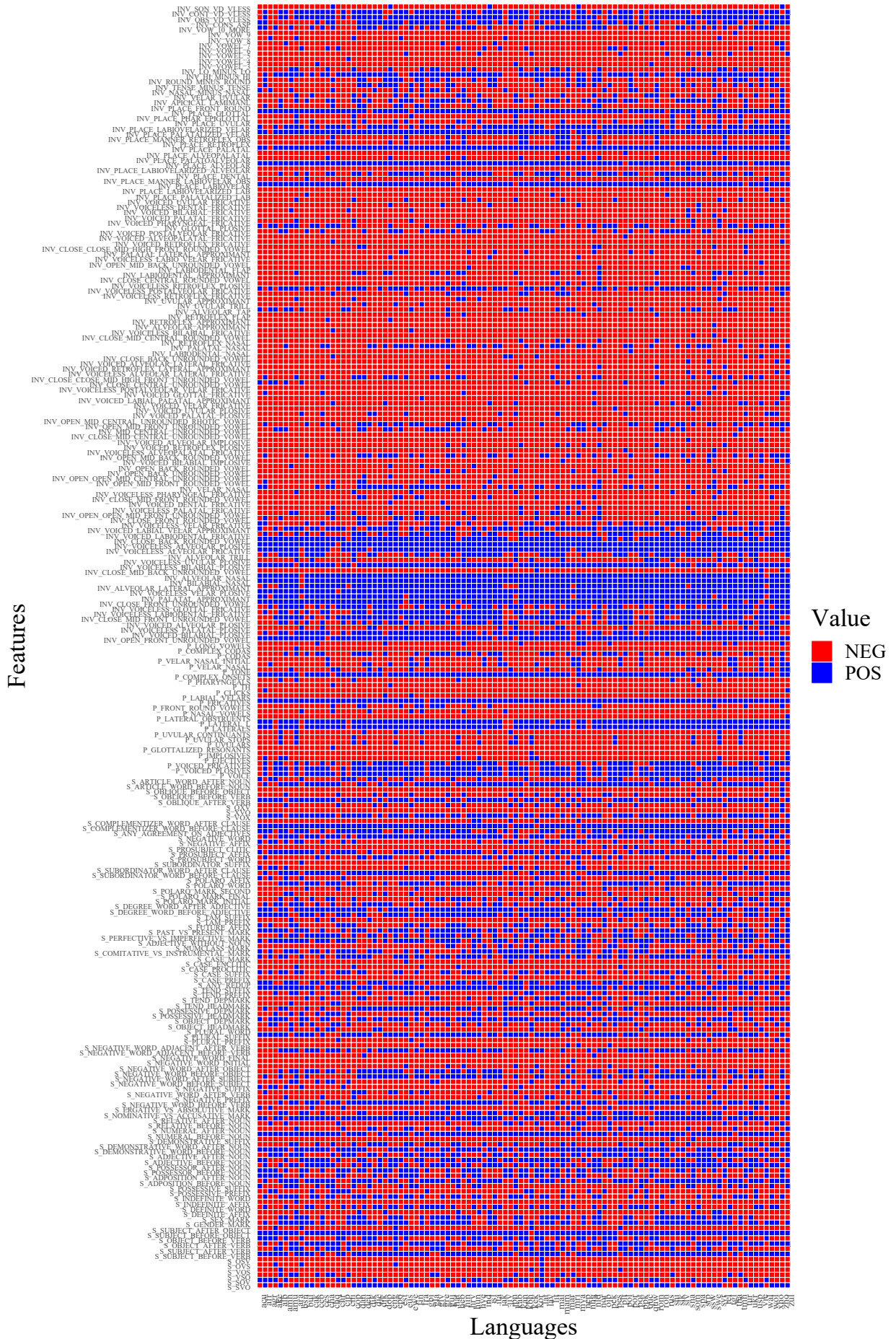


Figure 1: Binary values of the typological features from Littell et al. (2017) (y-axis) for each language (x-axis).