

# Relation between Word Order of Languages and the Entropy of Y-chromosome Haplogroup Distribution of the Speakers' Population

## 言語の語順とその話者集団のY染色体ハプログルーブ 分布のエントロピーとの関係

Terumasa EHARA

江原暉将

Ehara NLP Research Laboratory

江原自然言語処理研究室

<http://www.ne.jp/asahi/eharate/eharate/>

### 1 Introduction

We are investigating the relation between the word order of languages and the speakers' thought pattern. We have approached it through suicide rate and homicide rate (Ehara, 2015, 2017). Previous works concluded that head final language speakers tend to have higher suicide rate and lower homicide rate and head initial language speakers tend to have lower suicide rate and higher homicide rate. Higher suicide rate relates to intropunitive thought pattern and higher homicide rate relates to extrapunitive thought pattern.

In this paper we use another metric concerning to a thought pattern. It is the entropy of Y-chromosome haplogroup (YHg) distribution of the speakers' population (shortly "entropy"). Higher entropy value means the population is rich in diversity and it relates to peaceful thought pattern. Lower entropy value mean the population is poor diversity and it relates to warlike thought pattern (Sakitani, 2008).

Our conjecture is that head final language speakers' population has higher entropy value and head initial language speakers' population has lower entropy value.

### 2 Data

Data for the word order features are obtained from the WALS online database (Dryer, 2013). We use two dominant word order features:

- Order of Object (O) and Verb (V),
- Order of Adjective (A) and Noun (N).

WALS online provides O and V feature values for 1519 languages and A and N feature values for 1366 languages.

YHg data of populations are obtained from Wikipedia's page of "Y-chromosome haplogroups by populations" (Wikimedia Foundation, 2015). After preprocessing (Ehara, 2016), we get the YHg data from 452 populations. The number of languages are 196 in this data. In the Wikipedia pages, the granularity of YHg is different page by page. We adjust the granularity to the coarsest 20 groups from A to T (Karafet et al., 2008). The entropy is calculated by

$$H = \sum_{g \in G} p(g) \log_{10} p(g)$$

where  $G$  is the set of YHg = {A, B, ..., T} and  $p(g)$  is the probability (relative frequency) of YHg  $g$  in a population. Appendix 1 shows the base data used in the research sorted by the entropy.

### 3 Analysis and results

We conduct t-test for O and V word order groups and A and N word order groups. We discard "no dominant order" data. Results are shown in Table 1. Both O and V case and A and N case, mean values are different significantly. The entropy of OV type language speakers' population is tend to be higher than that of VO type language speakers' population. The entropy of AN type language speakers' population is tend

to be higher than that of NA type language speakers' population. This results shows our conjecture is significantly true.

**Table 1: Results of t-test for the entropy**

	OV	VO
n	83	82
mean	0.453	0.363
standard dev.	0.209	0.192
t	2.909	
p	0.0041	

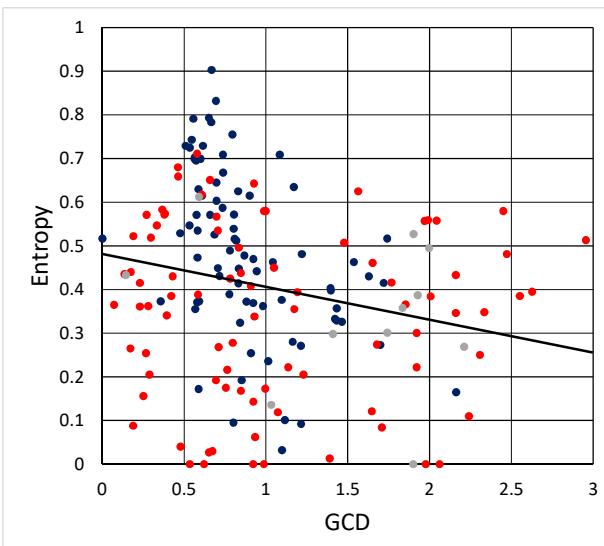
	AN	NA
n	79	86
mean	0.486	0.348
standard dev.	0.191	0.199
t	4.533	
p	0.00001	

#### 4 Distance from Africa

Homo sapiens was born in Africa and spread to all over the world. Because the diversity of haplogroup is due to the change in base sequence, the diversity is rich in near location from Africa and poor in far location from Africa. We consider it in our analysis using the linear regression. We measure the geometrical distance of languages from Africa by the great circular distance (GCD) calculated by the position data (latitude and longitude) in the WALS. The origin is set to the Amharic position. GCD values are listed in Appendix 1. From the result of linear regression analysis, we obtain the regression formula calculating the entropy (H) from GCD :

$$H = -0.0711 * GCD + 0.4819$$

Scattering graph between GCD and H of languages is shown in Figure 1 with the liner regression line.



**Figure 1: Linear regression result**  
(red:NA, blue:AN, gray:No dominant order or  
No data)

We conduct again t-test for the residuals of the regression. The results are shown in Table 2. We recognize again OV and VO groups and AN and NA groups both have the significantly different mean values for the residuals.

**Table 2: Results of t-test for the residuals**

	OV	VO
n	83	82
mean	0.040	-0.044
standard dev.	0.196	0.198
t	2.750	
p	0.0066	

	AN	NA
n	79	86
mean	0.067	-0.058
standard dev.	0.181	0.205
t	4.159	
p	0.00005	

#### 5 Conclusion

Relation between word order (Object and Verb and Adjective and Noun) of languages and the entropy of Y-chromosome haplogroup distribution of the speakers' population is examined. T-test results show OV and AN word order language speakers' population tend to have higher entropy value than VO and NA word order language speakers' population.

One of the remaining issues is to clarify the relation between word order and the entropy of the speakers' mitochondrial DNA haplogroup distribution which reflects diversity of populations with female line.

#### References

- Matthew S. Dryer. 2013. Order of Object and Verb and Order of Adjective and Noun, In: Dryer, Matthew S. & Haspelmath, Martin (eds.) *The World Atlas of Language Structures Online*. Leipzig: Max Planck Institute for Evolutionary Anthropology.  
(<http://wals.info/chapter/83> and 87, Accessed on 2015-3-23).
- Terumasa Ehara. 2015. Relation between Word Order Parameters and Suicide / Homicide Rates, *Journal of Yamanashi Eiwa College*, Vol.13, pages 9-29.
- Terumasa Ehara. 2016. Structural Distance, Lexical Distance and Speaker's Genetic Distance of Languages, *Proceedings of The 22th Annual Meeting of The Association for Natural Language Processing*, P5-1, pages 123-126.
- Terumasa Ehara. 2017. Relation between the Word Order Characteristics and Suicide/Homicide Rates (6), *Proceedings of The 23th Annual Meeting of The Association for Natural Language Processing*, P4-4, pp.190-193.
- Tatiana M. Karafet, Fernando L. Mendez, Monica B. Meilerman, Peter A. Underhill, Stephen L. Zegura, and Michael F. Hammer. 2008. New Binary Polymorphisms Reshape and Increase Resolution of the Hu-

man Y Chromosomal Haplogroup Tree, *Genome Research*, Vol. 18, pages 830-838.

Mitsuru Sakitani. 2008. A long journey of Japanese populations revealed by DNA analysis, *Showado, Kyoto*, pages 37-39 (in Japanese).

崎谷満. 2008. DNAでたどる日本人10万年の旅, 昭和

堂, 京都, pages 37-39.

Wikimedia Foundation. 2015. Y-chromosome haplogroups by populations, *Wikipedia*. ([https://en.wikipedia.org/wiki/Y-chromosome\\_haplogroups\\_by\\_populations](https://en.wikipedia.org/wiki/Y-chromosome_haplogroups_by_populations)). Accessed on 2015-11-20.

## **Appendix 1 Base data used in the research**

Feature values for O and V are 1:OV, 2:VO, 3>No dominant order, blank: No data.

Feature values for A and N are 1:AN, 2:NA, 3:No dominant order, blank: No data.

GCD means great circular distance of languages originated by Amharic position.

n means Number of sample.

A to T means Y-chromosome haplogroup name (cell values are relative frequency of YHg in the sample).

No.	WALS language name	O and V	A and N	GCD	n	Entropy	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
100	Lakhota	1	2	2.00711	106	0.384	0	0	5.5	0	0	0	0	0	0	0	0	0	0	0	0	0	39.62	42.45	0	0
46	Fula (Cameronian)	2	2	0.42179	63	0.385	3.17	0	0	0	60.27	0	0	0	0	0	0	0	0	0	0	0	0	22.2	0	4.76
158	Tongan	2	2	2.55251	55	0.385	0	0	23	0	0	0	0	0	0	0	1	0	8	0	60	0	0	0	0	
29	Cree (Plains)	2	1	1.92881	359	0.387	0	0	6.39	0	0	0	0	0	0	0	0	0	0	0	0	0	29.17	45.97	0	0
69	Jul hoan	2	2	0.5846	64	0.388	36	8	0	0	56	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
127	Polish	2	1	0.77636	1006	0.389	0	0	0	4.08	0	0	0	17.19	2.27	0	0	0	3.65	0	0	0	68.66	0	0	
63	Indonesian	2	2	1.19295	1536	0.394	0	0	15.48	0	0	0.33	0	1.22	0	0	3.66	0	1.21	0	66.44	0	0.14	0	2.61	
135	Samoan	2	2	2.628	87	0.395	0	0	60.93	0	0	0	0	0	0	3.43	0	2.28	0	28.73	0	0	1.14	0		
38	Evenki	1	1	1.39714	96	0.398	0	0	67.7	0	0	0	0	0	0	5.2	0	0	0	0	0	0	4.2	1	0	
31	Dagur	1	1	1.39488	39	0.403	0	0	30.8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
49	Goro	1	2	0.90768	71	0.408	0	0	8.5	0	0	4.2	0	1.4	0	0	4.2	0	0	0	0	0	0	0		
107	Mari (Meadow)	1	1	0.83139	159	0.414	0	0	0	0	0	0	0	5.65	1.9	0	0	0	44	0	0	0	47.14	0	0	
47	Fur	1	2	0.20333	32	0.415	31.3	3.1	0	0	59.4	0	0	0	6.3	0	0	0	0	0	0	0	0	0		
92	Koryak	1	1	1.71994	27	0.415	0	0	59.2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
6	Arrernte (Mparntwe)	1	2	1.76791	108	0.416	0	0	66.2	0	0	3	0	0	0	22.2	0	0	0	0.9	0	0	8	0		
30	Catalan	2	2	0.78125	250	0.425	0	0	0	0	9.55	0	5.41	0	6.14	0.4	0	0	0	0	0	0	0	67.16		
34	Dovayo	2	2	0.42965	72	0.43	1.4	12.5	0	0	58.4	0	0	0	0	0	0	0	0	0	0	0	0	0		
68	Japanese	1	1	1.62985	493	0.43	0	0	6.88	36.8	0	0	0	0	0	0	0	0	0	2.05	53.19	0	0.4	0		
140	Slovene	2	1	0.7157	75	0.431	0	0	0	0	2.7	0	2.7	0	30.7	4	0	0	0	0	0	0	0	60	0	
14	Beja	3	0	1.04367	42	0.433	4.8	0	0	52.4	0	0	0	38.1	0	0	0	0	0	0	0	0	4.8	0		
94	Kwao	2	2	2.16138	32	0.433	0	0	0	0	0	0	0	59.4	0	9.4	0	28.1	0	0	0	3.1	0			
111	Moro	2	2	0.13555	28	0.435	46.4	14.3	0	0	39.3	0	0	0	0	0	0	0	0	0	0	0	0			
44	French	2	2	0.8471	23	0.438	0	0	0	0	8.7	0	0	0	17.4	4.3	0	0	0	0	0	0	52.2	0		
33	Dirka	3	2	0.17423	107	0.44	53.37	24.29	0	0	22.34	0	0	0	0	0	0	0	0	0	0	0	0			
134	Romaní (Welsh)	2	1	0.94294	57	0.442	0	0	0	0	29.8	0	0	0	5.6	5.3	1.8	0	0	0	0	0	0	3.6		
51	German	3	1	0.83365	1231	0.448	0	0	0	0	6.2	0	0	0	23.58	3.95	0	0	0	1.58	0	0	0	56.92		
77	Khoekhoe	1	1	0.7067	413	0.449	42.13	12.11	0	0	43.1	0	0	0	0.24	0	0	0	0	0	0	0	0.71	0		
15	Batak (Karo)	2	2	1.04981	57	0.45	0	0	5.3	1.8	0	14	0	0	0	3.5	0	0	0	61.3	0	0	0	3		
104	Maybrat	2	2	1.65217	245	0.461	0	0	16.29	0	0	0	0	0	0	6.97	0	64.92	0	1.99	0	0	0	9.38		
143	Saami (Northern)	2	1	1.04193	38	0.463	0	0	0	0	0	0	0	31.6	0	0	0	44.7	0	0	0	0	0	0	23.7	
165	Udihe	1	1	1.53759	76	0.463	0	0	42.12	0	0	0	0	0	0	1.31	0	0	14.49	42.13	0	0	0			
147	Swedish	2	1	0.92273	160	0.47	0	0	0	0	1.3	0	0	0	37.5	0	0	0	0	0	0	0	37.5			
12	Avar	1	1	0.58191	177	0.473	0	0	0	0	1.68	0	7.39	0	11.7	66.18	0	3.94	0	11	0	0	0	13.58		
35	Dutch	3	1	0.86823	410	0.478	0	0	0	0	2.9	0	4.1	0	32.9	5.1	0	0	0.2	0	0	0	53.5	0		
23	Buriat	1	1	1.2194	238	0.481	0	0	63.9	0	0	0	0	0	0	8.8	0	0	20.2	0	1.7	1.7	2.9			
163	Tuvaluan	1	2	2.47249	100	0.481	0	0	17	0	0	2	0	0	0	36	0	0	0	45	0	0	0			
16	Belorussian	3	1	0.78052	68	0.489	0	0	0	0	4.4	0	1.5	0	25	0	0	0	8.8	0	0	0	50	0		
156	Tolai	2	3	1.99723	395	0.495	0	0	2.3	0	0	0	0	0	0	46.8	0	39	0	7.1	0	0	0	4.6		
145	Spanish	2	2	0.83355	940	0.496	0	0	0	0	13.64	0	3.46	0	7.76	9.64	0	0	0	0	0	0	0	57.97		
121	Ngada	2	2	1.47874	71	0.507	0	0	63.3	0	0	0	0	0	0	11.3	0	2.8	1.4	8.5	0	0	0	12.7		
164	Tatar	1	1	0.81811	76	0.512	0	0	0	0	2.6	0	11.8	0	0	47.4	0	0	0	1.3	0	0	0	36.8		
149	Tahitian	2	2	2.95615	810	0.513	0	0	44.68	0	0	0	0	0	0	0	0	0	0	0	0	0	0	11.0		
113	Mundari	1	1	0.80939	1681	0.516	0	0	0	0	3.67	0	24.18	0	4.13	0	0	0	0	56.17	1.66	0	7.5			
59	Anharic	1	1	0	290	0.517	17.27	1.02	0	0	46.92	0	0	0	27.96	0	0	0	0	0	0	0	3.78			
25	Chukchi	1	1	1.74162	24	0.517	0	0	4.2	0	0	0	0	0	0	0	0	0	0	58.3	0	20.8	15.5			
28	Coptic	2	2	0.29634	66	0.519	0	0	7.6	0	0	21.2	0	0	0	45.5	0	0	0	0	0	0	0	15.2		
117	Nubian (Dongolese)	1	2	0.1891	78	0.522	0	0	3.85	0	0	23.1	0	0	0	2.55	43.6	0	0	0	0	0	0	0		
167	Ukrainian	2	1	0.68458	103	0.526	0	0	0	0	1.94	0	1.94	0	21.34	7.75	0	0	5.85	0	0	0	58.26			
175	Yup'ik (Central)	3	1	1.9026	149	0.527	0	0	0.68	0	0	0	0	0	0	0	0	0	27.97	0	4.03	40.93	17.45			
76	Khalaj	1	1	0.47634	136	0.528	0	0	4.45	0	0	0.75	0	0	0	2.21	0	2.95	0	36.04	0	4.44	0	49.27		
124	Ossetic	1	1	0.58316	533	0.535	0	0	0	0	0.93	0.94	60.39	0	2.44	14.06	2.25	0.58	0	0.17	0	0.56	1.32	10.31		
150	Taijik	1	2	0.70658	38	0.535	0	0	2.6	0	0	3	0	0	0	5	0	18.4	0	8	0	0	0	0	0	
132	Russian	2	2	0.80285	995	0.539	0	0	0	0	3.33	0	0.5	0	14.39	2.5	0	0	0	26.51	0	0	0	48.46		
10	Arabic (Modern Standard)	2	2	0.27007	4564	0.541	0.18	0.11	0	0	41.91	0.64	2.88	0	1.57	36.92	0.42	1.32	0	0	0	0	0	7.49		
60	Hunzib	1	1	0.57545	104	0.571	0	0	0	0	0	0	0	7.68	0	14.4	48.07	0	0.96	0	0	0	0	0	0.94	
119	Nepali	1	1	0.83045	77	0.625	0	0	7.8	0	0	0	0	0	0	11.7	0	10.4	0	0	0	0	0	0	5.91	
160	Tetun	2	2	1.56401	39	0.625	0	0	43.5	0	0	0	0	0	0	10.3	0	10.3	0	18	0	0	0	51.3		
90	Karakhay-Balkar	1	1	0.58783	38	0.63	0	0	0	0	2.6	0	28.9	2.6	2.6	7.9	0	5.3	0	0	0	0	0	0	34.3	
75	Khalakha	1	1	1.17079	260	0.635	0	0	53.25	1.89	0	0.78														