

# Written Language vs. Non-Linguistic Symbol Systems: A Statistical Comparison

Richard Sproat

Google, Inc.

76 9th Avenue

New York, NY 10011, USA

1 **Are statistical methods useful in distinguishing written language from *non-***  
2 ***linguistic* symbol systems? Some recent papers in *Science* (1) and the *Proceed-***  
3 ***ings of the Royal Society* (2) have claimed so.**

4 **Using a larger set of non-linguistic and comparison linguistic corpora than**  
5 **were used in these and other previous papers, I show that none of the previous**  
6 **proposed methods are useful as published. However, one of the measures pro-**  
7 **posed in (2) (with a different cutoff value), as well as a novel measure based on**  
8 **repetition, turn out to be good measures for classifying symbol systems into the**  
9 **two categories. For the two ancient symbol systems of interest to (1) and (2) —**

1        **the Indus Valley “script” and Pictish symbols, respectively, both of these mea-**  
2        **asures classify them as non-linguistic, contradicting the findings of that previous**  
3        **work.**

## 4    **Introduction**

5    One of the defining habits of human beings is the use of visual marks to convey information.  
6    People of all cultures have been communicating with symbols etched on stone, pottery, brick,  
7    wood, preparations of leaves, and many other materials for thousands of years. Many different  
8    types of information can be represented with symbols. Familiar systems such as traffic signs,  
9    written music or mathematical symbols represent information relevant to those domains. In  
10   cases like these the information conveyed does not depend upon reference to any specific lan-  
11   guage: for example most traffic signs are intended to be “readable” in any language. Thus in  
12   such cases the information conveyed is *non-linguistic*. One special type of symbol system rep-  
13   resents linguistic information — phonological segments, syllables, morphemes or in some cases  
14   words: this latter type of *linguistic* symbol system we call *writing*, and the regular use of such a  
15   system defines a *literate* culture (3). The first literate cultures for which we have clear evidence  
16   arose about 5,000 years ago in Mesopotamia and at roughly the same time in Egypt (4).

17     Given this definition, the distinction between linguistic and non-linguistic systems may  
18     seem a clear one, but of course the definition presumes one *knows* what the symbols denote.  
19     In the case of an ancient culture that left behind “texts” consisting of otherwise unknown sym-

1 bols, the default situation is that one does *not* know what the symbols mean. Two examples of  
2 such cultures are the 3rd millenium BCE Indus Valley culture (5, 6), and the Pictish culture of  
3 Iron Age Scotland (7). In both cultures one finds texts consisting of one or more pictographic  
4 symbols. The fact that the symbols are pictographic is of no importance to their status, since all  
5 ancient writing systems are also pictographic. The traditional assumptions about these systems  
6 — that the Indus symbols were writing, and that the Pictish were probably not — are also of no  
7 importance.

8 One might therefore be tempted to ask if the statistical distribution of symbols in linguistic  
9 systems has some sort of “signature” that can distinguish these systems from non-linguistic  
10 systems. Indeed, just such a claim has been made in recent papers by Rajesh Rao and colleagues  
11 (1, 8) for the Indus Valley symbols, and Rob Lee and colleagues (2) for Pictish symbols. Both  
12 of these strands of work make use of the information-theoretic measure of *entropy* (9): in the  
13 case of Rao et al’s work either conditional entropy or what Rao terms “block entropy”; in the  
14 case of Lee et al’s work a more complex set of measures which include conditional entropy as  
15 one of the terms. In both cases the results seem to suggest that the symbol systems of interest  
16 behave more like linguistic systems than they do like the handful of non-linguistic systems that  
17 the authors compared them to. Lee et al interpret their method as *discriminative* in that they  
18 claim to have found a technique that can classify an unknown system as either non-linguistic  
19 or as one of a set of defined types of linguistic system. Rao and colleagues (1, 8) on the other  
20 hand interpret their method as “inductive”: the entropic values for the Indus symbols do not  
21 demonstrate that the system was writing, but if you assume that the system *was* writing, then it

1 follows that the entropic values should look like those of known writing systems, which is what  
2 they claim to have found.

3 In this paper we show that none of these approaches work when one compares a larger  
4 and more balanced set of non-linguistic systems drawn from a range of cultures and eras. Our  
5 non-linguistic corpora, detailed in **S1**, **S2**, comprise: Vinča symbols (10), Pictish symbols (7,  
6 11–14), Mesopotamian deity symbols (*kudurrus* (15–17)), Totem Poles (18–29), Pennsylvania  
7 German barn stars (“hex signs” (30–32)), sequences of five-day forecast icons from the Weather  
8 Underground, and individual Unicode code points in a corpus of Asian emoticons. We also  
9 included a small corpus of Indus Valley bar seal texts — a subset of the texts that Rao and  
10 colleagues had at their disposal — developed as part of earlier work.

11 We compared these seven systems with fourteen writing systems coded as indicated in  
12 parentheses: Amharic (syllables), Arabic (letters), Chinese (characters), Ancient Chinese Ora-  
13 cle Bone texts (characters), Egyptian (glyphs), English (letters), Hindi (letters), Korean (coded  
14 two ways: syllables and letters), Mycenaean Greek (Linear B – syllables), Malayalam (letters),  
15 Oriya (letters), Sumerian (glyphs), Tamil (letters). Again, details of these linguistic corpora can  
16 be found in **S1**.

## 17 **Results**

18 We computed both Rao et al’s (1) bigram conditional entropy and Rao’s (8) block entropy for  
19 our corpora (**S1**). For the block entropy, since Rao gave a reference to the exact method he

1 used (33), complete with a pointer to the MATLAB software and the parameter settings, it is  
 2 possible to produce plots that are directly comparable to his. These, side by side with Rao's  
 3 plot, are given in Figures 1 and 2. Note that unlike Rao's clean separation of linguistic and  
 4 non-linguistic systems 1, the true situation evidenced in 2 is much more messy. Linguistic (red)  
 5 and non-linguistic (blue) systems are interleaved, with no clear pattern. A few points are worth  
 6 noting. First of all note that *kudurrus* have relatively high entropy, which is surprising given  
 7 that Rao in various publications (1, 34) has insisted that the entropy for *kudurrus* must be low.  
 8 Second, our bar seal corpus (green) matches Rao's Indus corpus (dark blue in his plot) very  
 9 well; while they are not expected to be identical, given that the bar seals are a subset of Rao's  
 10 corpus, it would be surprising if they showed radically different behavior. Finally note that  
 11 Ancient Chinese has a very low entropy, like Rao's estimates for Fortran: this is not surprising  
 12 given the repetitive nature of Oracle Bone texts, and it underscores the point that measures like  
 13 entropy are highly sensitive to the *kind* of text being analyzed.

14 Lee and colleagues develop two measures,  $U_r$  and  $C_r$ , defined as follows. First,  $U_r$  is  
 15 defined as  $U_r = \frac{F_2}{\log_2(N_d/N_u)}$ , where  $F_2$  is the bigram entropy,  $N_d$  is the number of bigram types  
 16 and  $N_u$  is the number of unigram types.  $C_r$  is defined as  $C_r = \frac{N_d}{N_u} + a\frac{S_d}{T_d}$ , where  $N_d$  and  $N_u$   
 17 are as above,  $a$  is a constant (for which, in their experiments, they derive a value of 7, using  
 18 cross-validation),  $S_d$  is the number of bigrams that occur once and  $T_d$  is the total number of bi-  
 19 gram tokens. They use these two features to train a decision tree to classify symbol systems into  
 20 four bins — non-linguistic, and three types of linguistic writing systems: if  $C_r \geq 4.89$ , the sys-  
 21 tem is linguistic, and then depending on the value of  $U_r$ , the system is segmental ( $U_r < 1.09$ ),

1 syllabic ( $U_r < 1.37$ ) or else logographic. On the basis of that tree they classify Pictish as a  
2 linguistic system. The problem is that when we apply their system with their parameters to our  
3 non-linguistic corpora, it almost uniformly classifies them as linguistic. The classifications are  
4 as follows: Asian emoticons (linguistic: letters), Barn stars (linguistic: letters), Mesopotamian  
5 deity symbols (linguistic: syllables), Pictish symbols (linguistic: words), Totem poles (linguis-  
6 tic: words), Weather icon sequences (linguistic: letters), Vinča symbols (non-linguistic). Note  
7 that the Pictish symbols are classified as words, which replicates Lee et al's conclusion.

8 We then considered other measures that might potentially distinguish linguistic from non-  
9 linguistic systems. One promising measure is the ratio of the number of symbols that repeat in a  
10 text and are *adjacent* to the symbol they repeat ( $r$ ), to the number of total repetitions in the text  
11 ( $R$ ) (**S1**). This measure is by far the cleanest separator of our data into linguistic versus non-  
12 linguistic, with higher values for  $\frac{r}{R}$  (e.g. 0.85 for barn stars, 0.79 for weather icons, and 0.63  
13 for totem poles), being nearly always associated with non-linguistic systems, and lower values  
14 (e.g. 0.048 for Ancient Chinese, 0.018 for Amharic or 0.0075 for Oriya) being associated  
15 with linguistic systems. If we set a value of  $\frac{r}{R} = 0.10$  as the boundary, only Sumerian and  
16 Mesopotamian deity symbols are clearly misclassified, while Asian emoticons and Egyptian  
17 are ambiguously on the border. Both Pictish ( $\frac{r}{R} = 0.26$ ) and Indus symbols ( $\frac{r}{R} = 0.58$ ) are  
18 solidly non-linguistic by this measure. The latter is particularly noteworthy in light of the  
19 discussion in (35) that identified odd repetition patterns in the Indus corpus as problematic  
20 for the linguistic hypothesis. Note that the repetition measure correlates negatively with mean  
21 text length (Pearson's  $r = -0.49$ ). However the correlation is clearly not perfect and  $\frac{r}{R}$  still

1 contributes to a reasonable split between linguistic and non-linguistic corpora, even when we  
2 artificially truncate the data to make the texts and corpus sizes more comparable as detailed in  
3 **S1**.

4 Finally we trained classification and regression trees (36) with the above-discussed and other  
5 measures in a series of experiments, holding out our Indus corpus, and then in a second set of  
6 experiments the Pictish corpus. In both cases the vast majority of experimental runs — 98% and  
7 97% respectively — classified these corpora as *non-linguistic*. See **S1** for further details. The  
8 features that were most often picked by the trees as discriminative were (not surprisingly) our  $\frac{r}{R}$   
9 repetition measure — and Lee and colleagues'  $C_r$ . As we discussed above, Lee et al's measures  
10 perform very poorly on our non-linguistic corpora when used out of the box with their published  
11 cutoff values. But  $C_r$ , with cutoff values retrained by the CART algorithm, actually turns out to  
12 be a reasonable discriminator (though not as good as  $\frac{r}{R}$ ). On the other hand, interestingly,  $C_r$   
13 also correlates positively with mean text length (Pearson's  $r = 0.4$ , or  $r = 0.72$  if we exclude  
14 the Amharic data, which for some reason are an outlier).

15 It is clear from the above discussion that while  $\frac{r}{R}$  and  $C_r$  turn out to be useful all of the  
16 methods as proposed in the previous literature fail as *discriminative* methods: for example, Lee  
17 et al's published decision tree misclassifies almost all of our non-linguistic systems as writing.  
18 But what about Rao et al's "inductive" interpretation? Again, in Rao's view, the fact that the  
19 entropy growth curves for the Indus Valley symbols fall into the same narrow band as linguistic  
20 systems, serves as sort of "sanity check" on the hypothesis that the Indus signs were writing.  
21 But what we have shown is that this narrow band of entropic growth curves is densely populated

1 by all sorts of symbol systems, both linguistic and non-linguistic; and among the non-linguistic  
2 systems many that are clearly meaning-bearing. So the middle band of the entropy curves serves  
3 equally well as a sanity check that one is dealing with a meaning bearing — but not necessarily  
4 linguistic — symbol system. Since nobody has disputed that the Indus signs must have meant  
5 something to the people that used them, this tells us nothing new.

6 The fact that mean text length seems to be a correlate of the two most discriminative mea-  
7 sures we have found —  $\frac{r}{R}$  and  $C_r$  — is noteworthy insofar as our non-linguistic corpora do  
8 tend to be shorter. This is not surprising since while a true writing system is expected to have  
9 long texts (since language itself is theoretically unbounded in the length of utterances that can  
10 be produced), there is no such *a priori* expectation for non-linguistic systems. One of the prob-  
11 lems for the thesis that the Indus symbols were a true writing system is the fact that all extant  
12 texts are very short (35, 37): the longest text on a single surface is 17 glyphs long, which is  
13 quite a bit shorter than our longest *kudurru* text (39 glyphs, and indeed 19 out of our 69 *kudurru*  
14 texts are longer than the longest Indus text). As has been argued elsewhere (35), the so-called  
15 “lost manuscript” hypothesis that states that longer texts were written on perishable materials,  
16 all of which have been lost, seems questionable given the absence of other “markers” of literate  
17 civilization (for example pens, styluses or ink pots). In any case the belief in the existence of  
18 a large trove of now lost literary material in the Indus “script” must be taken as a mere act of  
19 faith, in the absence of any substantive evidence for it.

20 The status of any ancient symbol system as a writing system must be supported by good em-  
21 pirical evidence. As argued in (35) for the Indus Valley symbols in particular, good arguments



1 for the linguistic status would be a decipherment into one or more languages that succeeds in  
2 convincing a wide body of scholars; the discovery of artifacts indicating an active culture of  
3 literacy; or the discovery of a long text or a text that is bilingual with a known contemporaneous  
4 writing system. These ought to count as minimal requirements for accepting the thesis that *any*  
5 unknown ancient symbol system is writing.

## 6 **References**

- 7 1. R. Rao, *et al.*, *Science* **324**, 1165 (2009).
- 8 2. R. Lee, P. Jonathan, P. Ziman, *Proceedings of the Royal Society A: Mathematical, Physical*  
9 *& Engineering Sciences* pp. 1–16 (2010).
- 10 3. J. Goody, I. Watt, *Literacy in Traditional Societies*, J. Goody, ed. (Cambridge University  
11 Press, New York, 1968), pp. 27–68.
- 12 4. P. Daniels, W. Bright, eds., *The World's Writing Systems* (Oxford, New York, 1996).
- 13 5. G. Possehl, *The Indus Age: The Writing System* (University of Philadelphia Press, Philadel-  
14 phia, 1996).
- 15 6. G. Possehl, *The Indus Age: The Beginnings* (University of Philadelphia Press, Philadelphia,  
16 1999).
- 17 7. A. Jackson, *The Pictish Trail* (Orkney Press, 1990).

- 1 8. R. Rao, *IEEE Computer* **43**, 76 (2010).
- 2 9. C. Shannon, *Bell System Technical Journal* **27**, 379 (1948).
- 3 10. S. M. Winn, *The Life of Symbols*, M. L. Foster, ed. (Westview Press, Boulder, 1990), pp.  
4 269–271.
- 5 11. A. Jackson, *The Symbol Stones of Scotland: a Social Anthropological Resolution of the*  
6 *Problem of the Picts* (Orkney Press, 1984).
- 7 12. Royal Commission on the Ancient and Historical Monuments of Scotland, Pictish symbol  
8 stones: A handlist (1994).
- 9 13. A. Mack, *Field Guide to the Pictish Symbol Stones* (The Pinkfoot Press, Balgavies, Angus,  
10 1997). Updated 2006.
- 11 14. E. Sutherland, *The Pictish Guide* (Birlinn Ltd, Edinburgh, 1997).
- 12 15. L. King, *Babylonian Boundary Stones and Memorial Tablets in the British Museum* (British  
13 Museum, London, 1912).
- 14 16. U. Seidl, *Die babylonischen Kudurru- Reliefs. Symbole mesopotamischer Gottheiten* (Uni-  
15 versitätsverlag Freiburg, 1989).
- 16 17. J. Black, A. Green, T. Rickards, *Gods, Demons and Symbols of Ancient Mesopotamia* (Uni-  
17 versity of Texas Press, Austin, TX, 1992).

- 1 18. British Columbia Provincial Museum, *British Columbia totem poles*, Printed by Charles F.  
2 Banfield printer to the King's most excellent majesty (1931).
- 3 19. V. E. Garfield, *The Seattle Totem Pole* (University of Washington Press, Seattle, 1940).
- 4 20. M. Barbeau, *Totem Poles*, vol. 1–2 of *Anthropology Series 30, National Museum of Canada*  
5 *Bulletin 119* (National Museum of Canada, Ottawa, 1950).
- 6 21. S. W. Gunn, *The Totem Poles in Stanley Park, Vancouver, B.C.* (Macdonald, Vancouver,  
7 BC, 1965), second edn.
- 8 22. S. W. Gunn, *Kwakiutl House and Totem Poles at Alert Bay* (Whiterocks Publications, Van-  
9 couver, BC, 1966).
- 10 23. S. W. Gunn, *Haida Totems in Wood and Argillite* (Whiterocks Publications, Vancouver, BC,  
11 1967).
- 12 24. F. Drew, *Totem Poles of Prince Rupert* (F.W.M. Drew, Prince Rupert, BC, 1969).
- 13 25. E. Malin, *Totem Poles of the Pacific Northwest Coast* (Timber Press, Portland, OR, 1986).
- 14 26. City of Duncan, *Duncan: City of Totems*. (1990). Duncan, BC.
- 15 27. H. Stewart, *Totem Poles* (University of Washington, Seattle, WA, 1990).
- 16 28. H. Stewart, *Looking at Totem Poles* (University of Washington, Seattle, WA, 1993).
- 17 29. R. D. Feldman, *Home before the Raven Caws: the Mystery of the Totem Pole*, Cincinnati,  
18 OH (Emmis Books, 2003).

- 1 30. A. C. Mahr, *Ohio History: The Scholarly Journal of the Ohio Historical Society* **54**, 1  
2 (1945).
- 3 31. T. E. Graves, *The Pennsylvania German hex sign: A study in folk process*, Ph.D. thesis,  
4 University of Pennsylvania, Philadelphia, PA (1984).
- 5 32. D. Yoder, T. Graves, *Hex Signs: Pennsylvania Dutch Barn Symbols and their Meaning*  
6 (Stackpole, Mechanicsburg, PA, 2000).
- 7 33. I. Nemenman, F. Shafee, W. Bialek, *Advances in Neural Information Processing Systems*  
8 (MIT Press, Cambridge, MA, 2002), vol. 14.
- 9 34. R. Rao, *et al.*, *Computational Linguistics* **36**, 795 (2010).
- 10 35. S. Farmer, R. Sproat, M. Witzel, *Electronic Journal of Vedic Studies* **11** (2004).
- 11 36. L. Breiman, J. Friedman, R. Olshen, C. Stone, *Classification and Regression Trees*  
12 (Wadsworth and Brooks, Pacific Grove, CA, 1984).
- 13 37. A. Parpola, *Deciphering the Indus Script* (Cambridge University Press, New York, 1994).

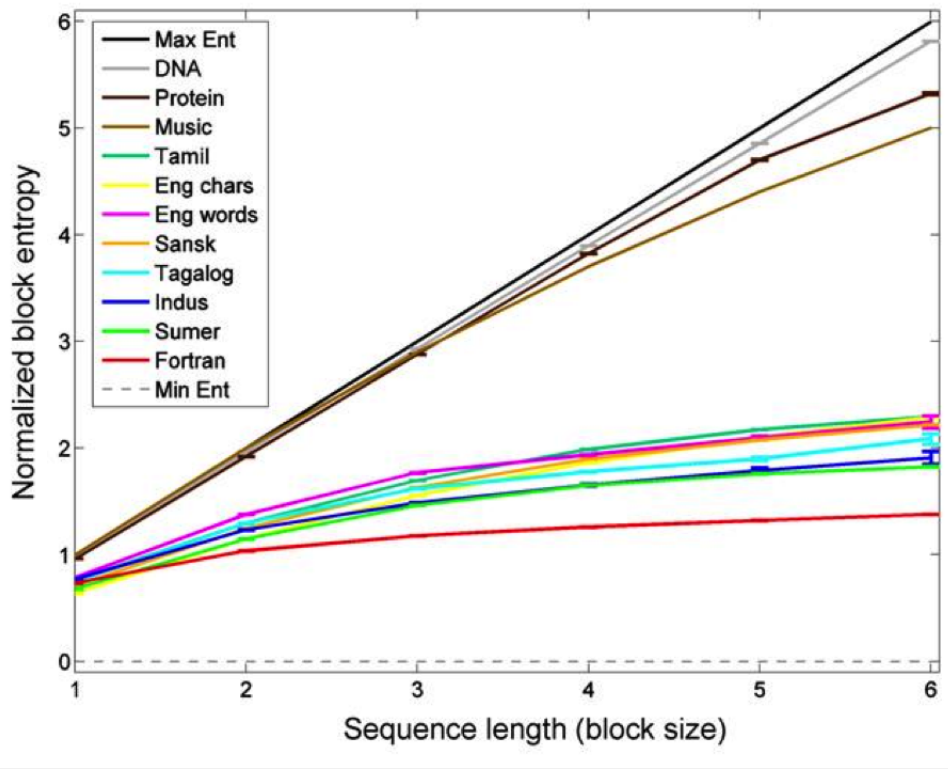


Figure 1: Block entropy estimates from (8) for a variety of linguistic systems, some non-linguistic systems, and the Indus corpus. Used with permission of the IEEE.

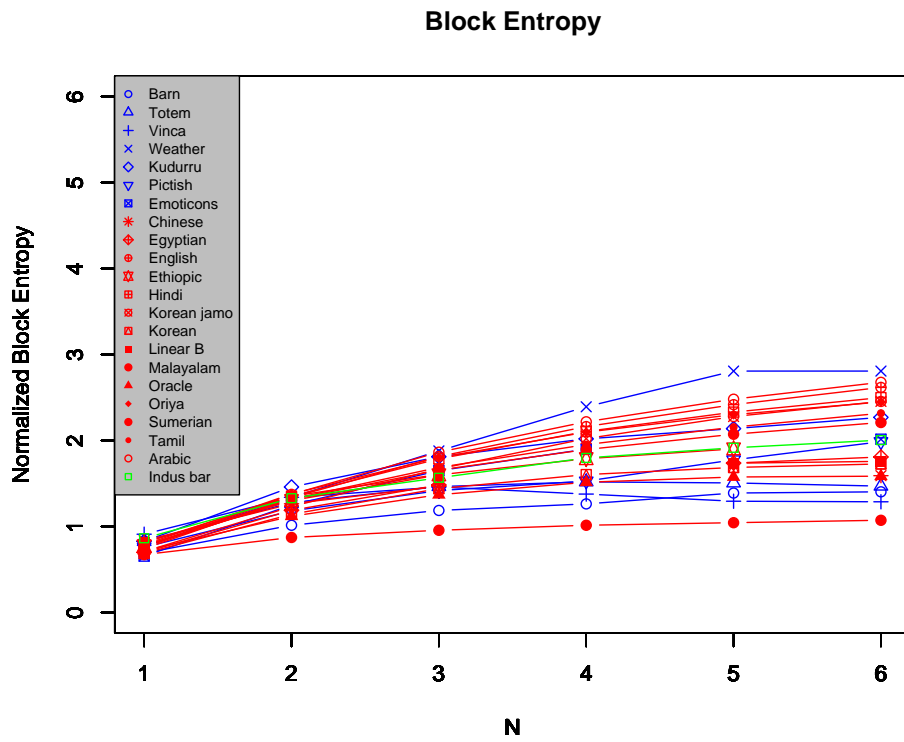


Figure 2: Block entropy values for our linguistic (red) and non-linguistic corpora (blue), and the Indus bar seal corpus (green).