

Written Language vs. Non-Linguistic Symbol

Systems: A Statistical Comparison

Supplementary Materials

Richard Sproat

Google, Inc.

76 9th Avenue

New York, NY 10011, USA

1 Introduction

2 In this supplement, I describe the methodology used to arrive at the results discussed in the
3 main paper. This includes a description of our corpora, as well as the statistical techniques and
4 the results of applying them to the corpora. I compare our results to those of previous work, in
5 particular (1, 2) and (3).

6 Since the notion of a non-linguistic system will be unfamiliar to many readers — even
7 though everyone uses such systems every day — I start with a brief taxonomy of such systems

1 as a segue into the more technical discussion to follow.

2 **Taxonomy**

3 Linguistic symbol systems — writing — have been classified in various ways by scholars of
4 writing systems for the past half century (4–9). While the proposed classifications differ, they
5 all share the property that they differentiate systems on the basis of the kind of information
6 encoded in the system, whether it be phonological segments, syllables, morphemes or words.

7 Systematic taxonomies of *non-linguistic* systems have certainly not been as prevalent and
8 in fact may not exist. Where such systems are discussed — e.g. (7, 10) — it is usually just a few
9 systems that are presented in juxtaposition to writing systems. In a forthcoming publication
10 (11), I outline and justify a preliminary taxonomy of non-linguistic systems. Here I merely
11 summarize the taxonomy, with examples, to give some context for the ensuing discussion.

12 But before we delve into the taxonomy, I should perhaps start by asking: Is the distinction
13 between linguistic and non-linguistic systems that I laid out in the main paper really that clear?
14 Lee and his colleagues (12) in their response to my critique (13) of their claim that Pictish
15 symbols were writing, claimed that they adhere to a much looser notion of what it means to be a
16 writing system, accepting Powell’s (14) definition, namely that “writing is a system of markings
17 with a conventional reference that communicates information”. Notice that Powell’s definition
18 makes no reference to (natural) language. Indeed from Powell’s point of view practically every
19 symbol system we are discussing here would be considered writing. If Lee and his colleagues

1 really accept such a broad definition, then their paper, published in the *Proceedings of the*
2 *Royal Society*, would seem to have simply derived the obvious fact that Pictish symbols were
3 a “system of markings” that had a “conventional reference that communicates information”,
4 something that nobody would have doubted: Presumably the Picts would not have expended
5 energy carving symbols on stones if they were not intended to communicate something, and the
6 fact that the symbols repeat on different stones suggests that they were conventional. In contrast,
7 by claiming that their method “reveals” Pictish symbols to be language, it seems that Lee and his
8 colleagues intended that readers accept the conclusion that the symbols served a similar function
9 to the marks that you are currently reading. Indeed, Powell aside, most students of writing
10 systems would accept the conclusion that writing is a special form of symbol system, one that
11 encodes linguistic information that is mostly or completely lacking in non-linguistic systems.
12 Without accepting such a distinction, work such as that of (1–3) would have no purpose. With
13 this distinction in mind I turn to a taxonomy.

14 Non-linguistic systems may be classified into at least the following types:

15 • **Simple informative systems.** In simple informative systems, the symbols by and large
16 convey a single piece of information. The helical red, white and blue barber pole, for
17 example, indicates the presence of barber shop. In weather reports, icons are used to rep-
18 resent various states of the weather, such as whether it is sunny, partly cloudy, raining,
19 etc. Further examples: other symbols of guild such as three balls for a pawnbroker; insti-
20 tutional logos; traffic information signs; ownership signs, such as brands, e.g. Mongolian
21 horse brands (15) or house marks (http://en.wikipedia.org/wiki/House_

1 mark).

2 • **Emblematic systems.** Closely related to simple informative systems are emblematic
3 systems where the symbols represent some special distinction earned by the bearer. Ex-
4 amples: symbols of military rank and distinction, scouting merit badges, Phi Beta Kappa
5 keys, and symbols of other scholarly fraternities; letter grades on academic assignments
6 or in courses.

7 • **Religious Iconography.** Religious iconography could also be characterized as a simple
8 informative system, except that here though here the notion of “single piece of informa-
9 tion” is much less clear, since such symbols are intentionally often highly multivocal in
10 the meanings they evoke. Examples: Christian cross, star of David, star and crescent,
11 dharmachakra, swastika (in its original Buddhist usage) .

12 • **Heraldic systems.** Heraldic systems are similar to emblematic system in that they usually
13 represent a particular set of features of the bearer, including possibly marks of distinction.
14 They differ, however, in that heraldic systems are frequently highly combinatoric, involv-
15 ing “texts” built of many symbols, often with a quite rigid syntax. Examples: European
16 heraldry, *kudurrus*, some functions of Totem poles (16) .

17 • **Formal systems.** In a formal system, the individual symbols have well-defined meanings,
18 and there are generally strict rules on how the symbols may be combined. Examples:
19 mathematical symbols, alchemical symbols, chemical notation, Feynman diagrams (17),
20 programming flowcharts and Systems Biology Graphical Notation (18).

- 1 • **Performative systems.** Performative systems indicate a sequence of actions to be taken
2 to perform a particular task. Perhaps the most familiar to many people today is are the
3 wordless assembly instructions that come with furniture from Ikea. Other examples: Silas
4 John’s system for notating Apache prayers (19), musical notation (20), dance notation
5 and other movement notation systems (21), chess notation, systems that can be used to
6 indicate the sequence of plays in a game, knitting patterns (10).
- 7 • **Narrative systems or “prompt” texts.** Narrative systems are used to recount stories and
8 as such are the most language-like of the non-linguistic systems. In narrative systems,
9 the symbols typically represent actors or events in the story in an iconic way. Examples:
10 Dakota winter counts (22), (probably) Nazi symbology (23).
- 11 • **Purely decorative systems.** Some systems that involve what are commonly thought of
12 as symbols seem nonetheless to be purely decorative. In such systems, the symbols may
13 derive historically from symbols that had meanings or ranges of meanings, but where
14 those meanings are quite irrelevant in their current use. A clear modern example is the use
15 of Chinese characters in body tattoos, worn by people who may be completely unaware
16 of the character’s original meaning, and use them solely because they look “cool”. Other
17 examples are Pennsylvania German barn stars (24), and Asian emoticons (e.g. (25)) where
18 in the latter case symbols from various scripts are combined into a “text” that represents
19 an image, usually a face. The face itself may convey some sort of emotion (e.g. sadness,
20 via depiction of crying), but other than that has no real meaning and mostly functions to

1 decorate the surrounding text.

2 For the purposes of the discussion here, I am particularly interested in systems where the
3 symbols may be combined into texts. As we have seen already in the discussion above, some
4 systems — narrative systems, performative systems, heraldic systems, formal systems, for ex-
5 ample — are particularly to be found in texts, but in fact systems of almost any kind may be
6 used combinatorically, depending upon what kinds of things the system denotes. Clearly a
7 mathematical system that only allowed single symbols, or a performative system for furniture
8 assembly that was similarly simplex, would not be of much utility. On the other hand, in reli-
9 gious symbology, a single symbol on its own may suffice to communicate the intended message:
10 to those devoted to the faith, the Christian cross can be a highly evocative symbol.

11 **3 Materials and Data Preparation**

12 For this work I selected seven non-linguistic symbol systems and fourteen linguistic symbol
13 systems for analysis. Note that in addition to the non-linguistic symbol systems described here,
14 a further five are under development and will be released along with the corpora described here.

15 **3.1 Non-linguistic corpora**

16 Two kinds of non-linguistic corpora were collected. The first is ancient or traditional systems
17 which serve as exemplars of what a complex ancient non-linguistic system can be like. If a set
18 of “texts” in a previously unknown symbol system is discovered among the ruins of an ancient

1 civilization, these systems could serve as models of what that symbol system might have been,
2 assuming it was not a form of written language. Our current set of such systems comprises
3 Mesopotamian deity symbols (*kudurrus*) (26), Vinča symbols (27), Pictish symbols (28–32),
4 Totem poles (16, 33–43), and Pennsylvania German barn stars (24, 44, 45). The development of
5 these ancient or traditional corpora is described in (46). Two of these systems are of interest
6 for other reasons: Mesopotamian deity symbols and Vinča systems were mentioned by (1) as
7 models for low- and high-entropy systems, respectively. As we shall see, Rao and his colleagues
8 were completely wrong about this, a point that could only have been known for certain by
9 actually doing the work of collecting the corpora. A third — Pictish symbols — were, of course,
10 already discussed by (3), and argued to be writing: their inclusion here as non-linguistic accords
11 with what is probably still the most widely accepted view, but in addition to that, the tests
12 presented below — including a retrained version of one of Lee and colleagues’ own measures
13 — are at least consistent with it being non-linguistic.

14 The second kind are modern systems, which can easily be collected electronically. These,
15 presumably, are poorer models of ancient symbol systems, but have the advantage that they are
16 easy to collect from online sources, which was mostly not the case for the ancient or traditional
17 systems, which had to be transcribed from print sources.

18 After discussing the various corpora, we present in Section 3.1.8 statistics on the size of the
19 various corpora. **S2** contains examples of the first five symbol systems, along with the corre-
20 sponding transcription in the XML markup scheme that was introduced in (46). One feature
21 of our XML markup is that it allows one some flexibility in how one extracts the elements of

1 the text. For example, if there are clear lines in the text, then we preserve line information. If
2 the text is circular (so that it is not defined where the beginning or end of the text occurs) that
3 information is also marked. Furthermore, one of the questions that arises in the analysis of a
4 symbol system is whether to treat elements that seem to be composed of more atomic symbols
5 as single symbols or as a composite of the individual symbols. The `symbolUnit` preserves
6 the grouping structure so that one can make the choice of level of analysis later on.

7 At the head of each section I propose a classification for the symbol system in question in
8 terms of the taxonomy developed in Section 2. The breakdown of the types is as follows. Note
9 that totem poles account for two of the types — heraldic and narrative, since poles can have
10 either of these functions:

Type	# of systems represented
Heraldic	3
Narrative	1
Decorative	2
Unknown	2
Simple informative	1
Formal	1

11
12 **3.1.1 Vinča symbols**

13 **Classification: Unknown, possibly religious**

14 The Vinča were a late Neolithic people of Southeastern Europe between the sixth to the third

1 millenium BCE who left behind pottery inscribed with symbols. Often the symbols occurred
2 singly, but on about 120 items the symbols occur in “texts” of two or more symbols; see **S2**
3 for a sample Vinča text. In such texts, the symbols are sometimes, but not always arranged
4 linearly as in a typical script. The most authoritative compilation of the Vinča materials is (27),
5 who classifies the symbols into about 200 types, and provides a corpus of the multisymbol
6 texts from 123 sources, comprising more than 800 tokens. In addition to their sometimes linear
7 arrangement, the Vinča symbols had other script-like properties: in Winn’s analysis, some signs
8 seemed to be comprised of two signs ligatured together. Though Winn characterizes the system
9 as “pre-writing”, it must be stressed that there is no reason to think this was a linguistic writing
10 system. Very likely, the system had religious significance. As Winn states (p. 255): “In the final
11 analysis, the religious system remains the principle source of motivation for the use of signs.”
12 However the meaning of the symbols remains unknown.

13 Developing Winn’s materials into an electronic corpus was relatively straightforward, since
14 he very nicely lays out the texts in his catalog, giving both the form and the type number for
15 each sign. We followed his linearization of the texts, but since he also provides line drawings
16 of the artifacts themselves, it was possible in some cases to indicate the true spatial relationship
17 between the symbols.

18 **3.1.2 Mesopotamian deity symbols on *kudurrus***

19 **Classification: Heraldic**

20 **S2** gives an example of a *kudurru* boundary stone with deity symbols. *Kudurrus* were

1 legal documents that specified property rights. The stones often also included actual cuneiform
2 Babylonian writing. The deity symbols had an essentially heraldic function, to indicate favored
3 deities of the owner of the property. Our source for the deity symbol corpus was Seidl (26). We
4 picked texts from all stones where the depictions in (26) were clear enough to read. Fortunately
5 (26) includes a chart that lists, for each stone, the symbols that appear on that stone (though not
6 in the order they appear); that chart proved useful for checking that the reading of the symbols
7 was correct.

8 **3.1.3 Pictish stones**

9 **Classification: Unknown, possibly heraldic**

10 The Picts were an Iron Age people (or possibly several peoples) of Scotland who, among
11 other things, left a few hundred standing stones inscribed with symbols, with “texts” ranging
12 from one to a few symbols in length. The meaning of the symbols is unknown, but until recently
13 it would never have occurred to anyone to assume that they are some form of writing. If nothing
14 else, the Picts evidently had at least some literacy in the (segmental) Ogham script (47). An
15 example of Pictish symbols can be found in **S2**.

16 Fortunately for our work, a corpus of images, comprising 340 stones, along with information
17 on the symbols on each stone is available from the University of Strathclyde <http://www.stams.strath.ac.uk/research/pictish/database.php>. The stones are cross-
18 referenced to standard texts such as (28–32). The main work that was needed was to correct
19 the ordering of symbols in some cases since the ordering of the symbols in the Strathclyde
20

1 transcription does not always correspond to the symbols' ordering on the stone.

2 **3.1.4 Totem Poles**

3 **Classification: Heraldic, narrative**

4 Totem poles are the product of a wide range of Native American cultures of the Pacific
5 Northwest, dating from the 19th and 20th centuries. The texts, which consist of anthropo-
6 morphic and zoomorphic symbols, carved vertically on cedar or other tree trunks, represent a
7 variety of kinds of information from legends, genealogical information, or depictions of impor-
8 tant events. Types of poles include memorial poles, house frontal poles, mortuary pole, and
9 heraldic poles. Different tribes had different styles so that, for example, on Haida poles the
10 figures are carved in bas relief (39). The “texts” also served different purposes among different
11 groups: Thus totem poles are a good example of how a symbol system can vary across different
12 regions and cultures without that variation implying that the system encoded language; cf. (48).

13 We used a number of published resources in developing our data. In addition to Malin (39),
14 we transcribed texts from (16, 33–38, 40–43).

15 There are many recurring themes on totem poles. For example, certain combinations of
16 symbols always allude to the same folk story. A whale with a man and a woman always refers
17 to the Nanasimget story (42). In principle then this could be represented by a single symbol
18 (e.g. NANASIMGET), but instead we elected to use the `symbolUnit` tag (46) to represent the
19 grouping:

20 `<symbolUnit>`

1 <symbol><title>Whale</title></symbol>
2 <symbol><title>Man</title></symbol>
3 <symbol><title>Woman</title></symbol>
4 </symbolUnit>

5 **3.1.5 Pennsylvania German Barn Stars**

6 **Classification: Decorative**

7 Barn stars, commonly known as “hex signs”, are a traditional decorative art among Penn-
8 sylvania German (“Dutch”) communities. They are found widely in northeastern Pennsylvania,
9 particularly Berks County, as well as in scattered communities in Ohio and elsewhere. They
10 are mostly used to decorate barns, but can also be found on other structures, such as porches.
11 Traditional barn symbols consist mostly of stars, rosettes, wheels-of-fortune, and swastikas.

12 There is a common belief that “hex signs” had a magical function, such as to ward off evil
13 spirits, and indeed this belief was sometimes expressed by the German farmers who used the
14 symbols on their barns (44). And the barn symbols themselves can be traced back in many
15 cases to symbols that probably had definite sets of meanings at one time (24, 44, 45). That
16 said, the consensus of much of the small scholarly literature on this topic is that barn stars
17 probably had no particular meaning at all for the people who used them and were used rather
18 for decoration (24, 45).¹

¹If barn stars are purely decorative, why include them here? There are two reasons. First of all, the symbols can occur in “texts” of several symbols in length. While the texts look rather unlike written texts — for one thing, they

1 Our barn star corpus is based upon the slide collection of W. Farrell, who toured areas
2 around Berks County in the 1940's and photographed many barns that were decorated with

are almost always symmetric — they are nonetheless interesting from the point of view of developing statistical
techniques that might possibly distinguish linguistic from non-linguistic systems.

Second, much of the critique of (49) has depended on a fundamental confusion of what is meant by the term “non-linguistic symbol system,” the most common misconception being that we were referring to randomly arranged meaningless symbols. A good example of this misconception is expressed by Vidale in his critique of our paper (50). Using the obviously decorative pottery designs of Shahr-e Sukhteh and elsewhere as his examples he states (page 344):

Together with coupling and opposition of selected symbols, systematic, large-scale redundancy (constant repetition of the same designs or symbols) is a distinctive feature shared by the more evolved and formally elaborated non-linguistic symbolic systems considered (highly repetitive patterns on the pottery of Shahr-i Sokhtai, endless repetition of icons such as scorpions, men-scorpions, temple facades, water-like patterns and interwoven snakes at Jiroft, and redundant specular doubling of most major symbols in the Dilmunite seals). While positional regularities might be detected in part of the Jiroft figuration, redundancy in all these systems dismiss one of the basic assumption of Farmer & others, who take the rarity of repeating signs as a proof of the non-linguistic character of the Indus script.

Repetition is indeed a characteristic of decorative art: consider the pineapple motif popular in American Colonial interior decoration and stenciled in repeated patterns on walls. And, not surprisingly, high symbol repetition rates do show up as a strong feature of barn stars. But high rates of repetition are not particularly characteristic of non-linguistic systems, though as we shall see high rates of local *reduplicative* repetition relative to the total repetition rate *do* seem to be characteristic of many non-linguistic systems. Systems that have massively higher than expected symbol repetition, such as barn stars, tend to be decorative. Vidale is simply confused on this point.

1 barn stars. His slides, in five boxes of about 100 slides each, are now housed in the archives of
2 the Berks County Historical Society in Reading, PA. The most useful boxes, in terms of amount
3 of material, are boxes 1 and 2: in these boxes the photos show the whole barn with associated
4 decorations. Boxes 3 and 4 contain a little more material, but in many cases the photos only
5 show a portion of the barn and its decorations, and there are a number of duplicates of barns
6 already seen in Boxes 1 and 2. Box 5 seems for the most part to be useless from the point of
7 view of collecting material on barn stars.

8 The designs of barn stars are quite numerous, but break down into a relatively small set
9 of categories, following the classification of Farrell himself, and Graves (24). The main ones,
10 in order of descending frequency of occurrence, are: 8-POINT-STAR, 5-POINT-STAR, 4-POINT-
11 STAR, 6-POINT-STAR, ROSETTE, SWIRLING-SWASTIKA, 12-POINT-STAR, WHEEL-OF-FORTUNE,
12 14-POINT-STAR, 7-POINT-STAR, 10-POINT-STAR, 15-POINT-STAR, 16-POINT-STAR. In many
13 cases the names of the slides in Farrell's system give a clear indication of what category the
14 (main) symbol on the barn belongs to: thus Box01/F.106St6 depicts a barn with three
15 six-pointed stars (St6). Farrell however does not distinguish some cases that Graves does dis-
16 tinguish: the main example of this is rosette, which Farrell classifies under six-pointed star; see
17 Figure 1. In such cases I followed Graves' classification, where I was able to clearly assign
18 the symbol to the rosette category. In one notable case I invented a new category, what I term
19 4-SLICE-PIE, which Farrell classifies as 4-POINT-STAR; see Figure 2. The "four-slice pie" is,
20 however, so distinct that it seems to deserve its own category.

1 **3.1.6 Weather Icons**

2 **Classification: Simple Informative**

3 The Weather Underground (www.wunderground.com) provides weather forecasts for
4 many parts of the world. The forecast includes icons that represent the predominant weather ex-
5 pectation for a given day. For example, [http://icons-pe.wxug.com/i/c/a/rain.](http://icons-pe.wxug.com/i/c/a/rain.gif)
6 [gif](http://icons-pe.wxug.com/i/c/a/rain.gif) is the icon for rain during the day. There are about 20 distinct icons.² These icons,
7 taken in series, form a “text” that corresponds to the weather predictions for a five-day period,
8 one icon per day. In this case we are dealing with a human-designed symbol system, but one
9 where the distribution of the symbols is determined by natural phenomena (or more properly a
10 computational model thereof).

11 We stored weather icon “texts” by collecting weekly forecasts from the Weather Under-
12 ground site for a selection of 161 cities throughout the world for 72 days, giving us a corpus of
13 50,710 symbols.³ See Figure 3 for an example of a weather “text”.

14 **3.1.7 Asian Emoticons**

15 **Classification: Decorative**

16 As part of a project on normalization of Twitter messages, my colleagues and I developed
17 an analysis system to detect and parse Asian emoticons (*kaomoji*) (25). Unlike the familiar

²There are actually double this number since there are separate icons for day and night: for example alongside the “partly cloudy” icon for daytime, there is a nighttime version. In this work we only included the daytime icons.

³Note that on occasion the scraping of the page resulted in no data returned.

1 90-degree flipped ASCII “smileys” — :-), ;-), :- (, 8-) ...— Asian emoticons (so-called
2 because they were popularized by Japanese and other East Asian users) are oriented horizon-
3 tally, and make use of a much wider range of characters. Some examples can be seen in Figure 4.
4 Traditional ASCII smileys are relatively limited, comprising perhaps a few tens of examples.
5 Asian smileys, in contrast, are productive and open ended: our collection includes thousands of
6 examples.

7 Of interest from the point of view of this project are the individual characters used in the
8 emoticons. Asian emoticons tend to be somewhat (though often not perfectly) symmetric. How-
9 ever unlike in the symmetric “texts” found with Pennsylvania barn stars, the mate characters
10 found in Asian emoticons are different symbols, chosen because they are visually close mirror
11 images. A statistical analysis of the symbol distributions would easily miss the fact that the
12 texts are symmetric.

13 **3.1.8 Statistics on Corpus Sizes**

14 The sizes comprising the numbers of texts, the number of tokens, the number of types and the
15 mean text length of the above-discussed corpora is given in Table 1.

16 In addition, we also included in our analysis a corpus of Indus bar seals from Harappa and
17 Mohejo Daro, which we developed as part of the work reported in (49). The details on this latter
18 corpus are given below:⁴

⁴The reason for not using the same Indus data as that used in (1, 2) is that none of the Indus corpora that have been collected by various groups over the years have been made available to other researchers, which has in turn

Corpus	# Texts	# Tokens	# Types	Mean text length
Indus bar seals	206	1,265	209	6.14

1
 2 Note that the bar seals are relatively late instances of Indus inscriptions, and that the mean
 3 length of the texts, 6.14, is substantially longer than the mean length of all Indus inscriptions
 4 taken together, which is 4.5 (49, 51).

5 **3.2 Linguistic Corpora for Comparison**

6 For comparison with the non-linguistic systems described above, I gathered a set of linguistic
 7 corpora from various existing sources. Just as I have tried to collect non-linguistic corpora
 8 of a variety of types, so I have also tried to gather linguistic corpora that are varied along a
 9 couple of dimensions. The first dimension is age: we have both extremely ancient systems
 10 — Sumerian, Egyptian, Ancient Chinese and Linear B (Mycenaean Greek), as well as various
 11 modern systems. The second dimension was type of writing system: we have examples of
 12 morphosyllabic writing, syllabaries, alphasyllabic writing, abjads and segmental writing. See,
 13 e.g., (7) for descriptions of what these various terms mean. In most cases each individual “text”
 14 in our corpus corresponded to an individual line in the encoding in the source corpus from
 15 which we extracted our sample.

16 The following sections give brief descriptions of each of the corpora. Summary statistics
 17 can be found in Table 2.

made it difficult to verify results.

1 **3.2.1 Amharic**

2 Collection of texts from <http://www.waltainfo.com>. The text was tokenized at the
3 syllable level, with spaces represented as a separate token (“#”).

4 **3.2.2 Modern Standard Arabic**

5 Collection of the first 10,000 headlines from the Linguistic Data Consortium (LDC) Arabic Gi-
6 gaword corpus <http://www ldc .upenn .edu/Catalog/catalogEntry.jsp?catalogId=LDC2003T12>. The text was tokenized at the letter level, with spaces represented as a separate
7 token (“#”).

9 **3.2.3 Modern Chinese**

10 Collection of the first 10,000 headlines from the LDC Chinese Gigaword corpus [http://](http://www ldc .upenn .edu/Catalog/catalogEntry.jsp?catalogId=LDC2003T09)
11 www ldc .upenn .edu/Catalog/catalogEntry.jsp?catalogId=LDC2003T09.
12 The text was tokenized at the character level.

13 **3.2.4 Ancient Chinese**

14 Oracle bone texts from Chenggong University, Taiwan. We kept only lines for each document
15 that contain no omissions. The text was tokenized at the character level.

1 **3.2.5 Egyptian**

2 Collection of texts transcribed using the JSesh hieroglyph editor (<http://jsesh.qenherkhopeshef.org/>) downloaded from <http://webperso.iut.univ-paris8.fr/~rosbord/hieroglyphes>.
3
4 The text was tokenized at the glyph (individual symbol) level.⁵

5 **3.2.6 English**

6 Collection of the first 10,000 headlines from the LDC English Gigaword corpus <http://www ldc.upenn.edu/Catalog/catalogEntry.jsp?catalogId=LDC2003T05>.
7
8 The text was tokenized at the letter level, with spaces represented as a separate token (“#”).

9 **3.2.7 Hindi**

10 Collection of the first 1,000 lines of text from the Hindi corpus developed at the Central Institute
11 of Indian Languages. The text was tokenized at the letter level, with spaces represented as a
12 separate token (“#”).

13 **3.2.8 Korean**

14 Collection of the first 10,000 headlines from the LDC Korean Newswire corpus <http://www ldc.upenn.edu/Catalog/catalogEntry.jsp?catalogId=LDC2000T45>.
15

⁵The following texts were used, all in <http://webperso.iut.univ-paris8.fr/~rosbord/hieroglyphes/>: CT160_S2P.hie, DoomedPrince.hie, HAtra.hie, HetS.hie, L2.hie, LC26.hie, Pachereniset.hie, Prisse.hie, gurob.hie, amenemope/*.gly, ikhernofret.hie, ineni.hie, kagemni.hie, lebensmuede.hie, mery.hie, naufrage.hie, sethnakht.hie, twobro.hie, year400.hie.

1 Following the standard Unicode encoding of Korean, the text was tokenized at the Hangul syl-
2 lable level, with spaces represented as a separate token (“#”).

3 **3.2.9 Korean Jamo**

4 Same source as above, but here the text was tokenized at the level of the (*jamo*) — the individual
5 letters that make up the syllable composites, with spaces represented as a separate token (“#”).

6 **3.2.10 Linear B**

7 Transcription by the author of 300 tablets from (52). Omitted lines with uncertainties and/or
8 omissions. Tokenized at the glyph level.

9 **3.2.11 Malayalam**

10 The Malayalam corpus (937 lines) developed at the Central Institute of Indian Languages. The
11 text was tokenized at the letter level, with spaces represented as a separate token (“#”).

12 **3.2.12 Oriya**

13 Collection of the first 1,000 lines of text from the Oriya corpus developed at the Central Institute
14 of Indian Languages. The text was tokenized at the letter level, with spaces represented as a
15 separate token (“#”).

1 3.2.13 Sumerian

2 Tablets from Gudea ruler of Lagah (c. 2100 BCE), extracted from the Cuneiform Digital Library
3 Initiative (<http://cdli.ucla.edu/>) via a transliteration search for *gu3-de2-a*, with “rime
4 3” in the primary publication field. Omitting lines with “#”, “<” or “[”, since these mark
5 uncertainty/reconstructions.⁶ “Texts” consisted of individual lines in the CDLI transcription,
6 meaning that the Sumerian “texts” are rather short.⁷

7 3.2.14 Tamil

8 Collection of the first 1,000 lines of text from the Tamil corpus developed at the Central Institute
9 of Indian Languages. The text was tokenized at the letter level, with spaces represented as a
10 separate token (“#”).

⁶Thanks to Chris Woods for suggesting appropriate search terms.

⁷One issue with Sumerian that will need to be resolved in future work is that the symbols are transcribed, following standard Sumerological practice, with romanized spellings of either the pronunciation of the symbol (if in lower case) or of the morpheme denoted by the symbol (if in upper case), followed by a subscript. Thus *aya*₂ is a transcription of 𒀝, which is one of the ways of writing the two-syllable sequence *a-ya*. The problem is that the system is not many-to-one, but many to many: 𒀝 could also be *duru*₅, *eš*₁₀, and several others. Standard Sumerological transcriptions thus give an overestimate of the actual number of glyph types appearing in a document. This issue can only be resolved by knowing which transcribed elements map back to a single glyph.

1 **3.3 Preparation**

2 Non-linguistic corpora were processed by extracting each text from the XML markup, and rep-
3 resenting the corpus with one text per line. In cases where there was a `SymbolUnit` (see
4 above) we extracted the individual symbols making up that unit. Linguistic corpora were repre-
5 sented with one text per line, tokenized as described above.

6 All probability-based entropy statistics discussed below were computed using tools based
7 on the OpenGrm NGram library (53), available from <http://www.opengrm.org>; these
8 tools will be made available along with the corpora. We used Kneser-Ney smoothing, including
9 the probability and thus entropy estimates for unseen cases. Start and end symbols to pad the
10 text are implicitly computed by the software.

11 **4 Methods and Results**

12 In the ensuing sections we describe the various statistical measures we applied to our corpora,
13 along with the results we derived.

14 **4.1 Bigram Conditional Entropy**

One measure that has been widely used to measure the information structure of language is
conditional entropy (54), for example, bigram conditional entropy, defined as follows:

$$H(Y|X) = - \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log p(y|x) \quad (1)$$

1 Entropy is a measure of information content: in the case of bigram entropy, if for any symbol x
2 there is a unique y that can follow x , then the message carries no information, and the entropy
3 will be minimal. If on the other hand, for any symbol x any other symbol can follow, and
4 with equal probability, then the entropy will be maximal. Most real symbol systems, including
5 language and a whole range of non-linguistic systems, fall between these two extremes.

6 Bigram conditional entropy was used by (1) to argue for the linguistic status of the Indus
7 Valley symbols. One can compute entropy over various sized units of granularity — symbols
8 in the script or symbol system, words, etc. — and for various portions of the corpus. In Rao
9 and colleagues' approach they computed the entropy between the n most frequent units, then
10 the $2n$ most frequent, the $3n$ most frequent, and so forth. This procedure derives an entropy
11 growth curve that starts relatively low and grows until it reaches the full entropy estimate for
12 the corpus, given a particular granularity.

13 Figure 5 shows these curves, taking $n = 20$, for a variety of actual writing systems, Ma-
14 hadevan's corpus of Indus inscriptions, and two artificial systems with maximum entropy, and
15 minimum entropy, which Rao and colleagues label as "Type 1" and "Type 2", respectively.
16 Figure 6 shows the *relative conditional entropy* — "the conditional entropy for the full corpus
17 relative to a uniformly random sequence with the same number of tokens" — for the various sys-
18 tems represented in Figure 5, as well as the real non-linguistic systems, DNA, protein ("prot")
19 and Fortran ("prog lang").

20 While the maximum ("Type 1") and minimum ("Type 2") curves are generated from com-
21 pletely artificial systems, Rao and colleagues argue that they are in fact reasonable models for

1 actual symbol systems: Vinča symbols and Mesopotamian deity symbols, respectively. Their
2 beliefs are based, respectively, on Winn’s (55) description of a certain subset of the Vinča cor-
3 pus as having no discernible order; and on the fact that there is a certain hierarchy observed in
4 the ordering of the symbols on *kudurru* stones. In (13) and (56), I argued that Rao’s interpreta-
5 tion of how these symbol systems work is a misinterpretation, and indeed the results we present
6 below support that conclusion. Rao and colleagues also present results for DNA and protein in
7 Figure 6 (though not in Figure 5): it is clear from that figure that the curves of these biological
8 symbol systems are very close to the maximum entropy curves.⁸⁹

9 Turning to our own data, in Figure 7 I show conditional entropy growth curves for all of
10 our non-linguistic systems (in blue), all of our linguistic samples (in red), and for our corpus of
11 Indus bar seals (in green).

12 As can be seen in the plot, the results are pretty much across the map. The *kudurru* inscrip-
13 tions turn out, contrary to what Rao and colleagues claim — and what Rao continued to insist
14 in (2) — to have the highest entropy of all the systems. Chinese shows the lowest. Between

⁸At least in the case of DNA, though, this is because Rao and colleagues take the “alphabet” of DNA to be the base pairs, an alphabet of size 4. But these are surely not the right units of information: the base pairs are effectively the “bits” of DNA. The equivalent for natural language would be taking the basic units of English to be the ones and zeroes of the ASCII encoding — and throwing in large amounts of random-looking “non-coding regions” to mimic the non-coding regions of DNA — surely a meaningless comparison for these purposes.

⁹Note that while Rao reports that they used Kneser-Ney smoothing (57), a number of details are not clear: for example, it is not clear whether they included estimates for unseen bigrams in their analysis. Discrepancies such as this and others may help explain why the entropy estimates I present below are on bulk lower than those presented in Rao and colleagues’ work.

1 these two extremes the linguistic and non-linguistic systems are mixed together. The Indus
2 bar seal corpus falls in the middle of this distribution, but that is quite uninformative since the
3 distribution is a mix of both types of system.

4 One concern is that the corpora are of many different sizes, which would definitely affect
5 the probability estimates. How serious is this effect? Could the (from Rao and colleagues'
6 point of view) unexpectedly high entropy estimate for the *kudurru* corpus be due simply to
7 undersampling? In Figure 8 I address that concern. This plots the bigram conditional entropy
8 growth curves for three samples from our Arabic newswire headline corpus containing 100
9 lines, 1,000 lines and 10,000 lines of text. Note that Arabic has roughly the same number of
10 symbol types as the *kudurru* corpus; also, the smallest sample of 100 lines is about 4 times
11 the size of the *kudurru* corpus, and thus is within the same order of magnitude. The entropy
12 does indeed change as the corpus size changes, with the smaller corpora having higher overall
13 entropy than the larger corpora. This is not surprising, since with the smaller corpora, the
14 estimates for unseen cases are poorer and, one would expect, more uniform, leading in turn to
15 higher entropy estimates. But the difference is not huge, suggesting that one cannot attribute
16 the relatively high entropy of the *kudurru* corpus wholly to its small size. Given a larger corpus,
17 the estimates would surely be lower. But there is clearly no reason to presume, as Rao does,
18 that a large *kudurru* corpus would show a conditional entropy near zero.

1 4.2 Block Entropy

We turn now to a second entropic measure, what Rao (2) terms “block entropy”, defined as follows:

$$H_N = - \sum_i p_i^{(N)} \log_{|V|} p_i^{(N)} \quad (2)$$

2 Here N is the length of an ngram of tokens, say 3 for a trigram, so all we are doing is computing
3 the entropy of the set of ngrams of a given length. One can vary the length being considered,
4 say from 1 to 6, and thus get estimates for the entropy of all the unigrams, bigrams, etc., all the
5 way up to hexagrams. In order to compare across symbol systems with different numbers of
6 symbols, Rao uses the log to the base of the number of symbols in each system ($\log_{|V|}$ in the
7 equation). Thus for DNA the log is taken to base 4. The maximum entropy thus normalized is
8 then just N , the length of the ngrams being considered, for any system.

9 Figure 1a in the main paper shows the block entropy estimates from (2) for a variety of
10 linguistic systems, some non-linguistic systems and the Indus corpus from (1). For this work,
11 Rao used an estimation method proposed in (58), which comes with an associated MATLAB
12 package. In this case, therefore, it is possible to replicate the method used by Rao and produce
13 results that are directly comparable with his results. These are shown in Figure 1b in the main
14 paper. As with the conditional entropy, and unlike Rao’s clear-looking results, the systems are
15 quite mixed up, with linguistic and non-linguistic systems interspersed. The Indus system is
16 right in the middle of the range. Note that the curve for our Indus corpus is rather close to the
17 curve for the Mahadevan corpus seen in Figure 1a in the main paper, suggesting that even though

1 the bar seals constitute a smaller subset of the whole corpus with rather different average text
 2 length, there is some similarity in the distributions of ngrams. Note on the other hand that our
 3 sample of Sumerian is radically different from Rao's in its behavior, and in fact looks more like
 4 Fortran. The Ancient Chinese Oracle Bone texts are also similar to Fortran, which is perhaps
 5 not surprising given that one finds quite a few repeating formulae in this corpus.

6 **4.3 Measures Derived in Part from Entropy**

Lee and colleagues (3), in their study of Pictish symbols, develop two measures, U_r and C_r ,
 defined as follows. First, U_r is defined as:

$$U_r = \frac{F_2}{\log_2(N_d/N_u)} \quad (3)$$

where F_2 is the bigram entropy, N_d is the number of bigram types and N_u is the number of
 unigram types. C_r is defined as:

$$C_r = \frac{N_d}{N_u} + a \frac{S_d}{T_d} \quad (4)$$

7 where N_d and N_u are as above, a is a constant (for which, in their experiments, they derive a
 8 value of 7, using cross-validation), S_d is the number of bigrams that occur once (*hapax legom-*
 9 *ena*) and T_d is the total number of bigram tokens; this latter measure will be familiar as $\frac{n_1}{N}$, the
 10 Good-Turing estimate of the probability mass for unseen events.

11 Lee et al. use C_r and U_r to train a decision tree to classify symbol systems. The resulting
 12 tree is shown in Figure 9. Thus if $C_r \geq 4.89$, the system is linguistic. Subsequent refinements
 13 use values of U_r to classify the system as segmental ($U_r < 1.09$), syllabic ($U_r < 1.37$) or else

1 logographic. As I note in the main paper, when one applies Lee et al’s tree, “out of the box” to
2 our data, all but one of our systems are classified as linguistic; see Table 3.

3 **4.4 Models of Association**

All of the measures thus far discussed involve entropy. Another property of language is that some symbols tend to be strongly associated with each other in the sense that they occur together more often than one would expect by chance. There is by now a large literature on association measures, but all of the approaches compare how often two symbols (e.g., words), occur together, compared with some measure of their expected cooccurrence. One popular measure is Shannon’s (pointwise) mutual information (54, 59) for two terms t_i and t_j :

$$\text{PMI}(t_i t_j) = \log \frac{p(t_i, t_j)}{p(t_i)p(t_j)} \quad (5)$$

4 This is known to produce reasonable results for word association problems, though there are
5 also problems with sensitivity to sample size as pointed out in (60).

6 In this study I use pointwise mutual information between adjacent symbols to estimate how
7 strongly associated the most frequent symbols are with each other. For words at least, it is often
8 the case that the most frequent words — function words such as articles or prepositions —
9 are *not* strongly associated with one another: consider that *the the* or *the is* are not very likely
10 sequences in English.

I therefore look at the mean association of the n most frequent symbols, normalized by the

number of association measures computed,

$$\frac{\sum_{j=0}^n \sum_{i=0}^n \text{PMI}(t_i t_j)}{n^2} \quad (6)$$

1 and let n range from 10, 20, . . . up to the k such that the first k symbols account for 25% of the
2 corpus. As with the entropy computations, I estimated probabilities of bigrams and unigrams
3 using OpenGrm.

4 As we can see, and as with the entropic measures, there is no clear separation of linguistic
5 and non-linguistic systems. Both kinds of systems have symbols that are more or less strongly
6 associated with each other (as well as symbols pairs that are not so strongly associated). For
7 example, in the case of Mesopotamian deity symbols there is a particularly strong association
8 between the HORNEDCROWN and the SYMBOLBASE symbols.

9 **4.5 Models of Repetition**

10 To the extent that they involve the statistical properties of the relations between symbols, the
11 measures that we have discussed so far are all local, involving adjacent pairs of symbols. But
12 despite the fact that local entropic models have a distinguished history in information-theoretic
13 analysis of language dating back to Shannon ((54)), language also has many non-local proper-
14 ties.

15 One is that symbols tend to repeat in texts. Suppose you had a corpus of boy scout merit
16 badge sashes and you considered each badge to be a symbol. The corpus would have many
17 language-like properties. Since some merit badges are awarded far more than others, one

1 would find that the individual symbols follow a roughly Zipfian distribution, as we see in
2 Figure 11. Furthermore, some badges tend to be earned before others, and while there is no
3 requirement that merit badges be applied to the sash in any particular order (<http://www.scoutinsignia.com/sash.htm>), nonetheless one would expect that the earlier earned
4 badges would tend to come earlier in the text. Thus we would expect a hierarchy rather like
5 what we find in the *kudurru* texts. All of this implies that entropic measures of the kinds we
6 have been considering would find “structure” in these “texts”. Yet such measures would fail
7 to capture what is the most salient feature of merit badge “texts”, namely that a merit badge
8 is never earned twice, and therefore no symbol repeats. In this feature, more than any other, a
9 corpus of merit badge “texts” would differ from linguistic texts.

11 Symbols in writing systems, whether they represent segments, syllables, morphemes or
12 other linguistic information, repeat for the simple reason that the linguistic items that they rep-
13 resent tend to repeat. It is hard to utter a sentence in any natural language without at least
14 some segments repeating. Obviously as the units represented by the symbols in the writing
15 system become larger (e.g. whole morphemes versus phonological segments), the probability
16 of repetition decreases, but we would still expect some repetition.

17 It is not only the amount of repetition in a corpus, but how the repetition manifests itself, that
18 is important. In (49) we argued that the Indus inscriptions not only showed rather low repetition
19 rates but that when one does find repetition of individual symbols within an inscription, the rep-
20 etition patterns often involve repetitions of adjacent symbols as well as symmetric arrangement
21 (see (49), Figure 6).

1 We can capture part of this oddity with a measure of repetition that seeks to quantify the *rate*
2 *of iterated repeats*, relative to the *total number of repetitions*. Specifically, we count the number
3 of tokens in each text that are repeats of a token that has already occurred in that text, and sum
4 that number over the entire corpus. Call this number R . We then count the number of tokens
5 that are repeats in each text, and which are furthermore adjacent to the token they repeat. Sum
6 that over the entire corpus, and call this number r . So for example for a single text:

7
8

9 A A B A C D B

10

11 R would be 3 (two *As* are repeats, and one *B*), and r would be 1 (since only one of the re-
12 peated *As* is adjacent to a previous *A*). Thus the repetition rate $\frac{r}{R}$ would be 0.33. Table 4 shows
13 our corpora ranked by $\frac{r}{R}$. This measure is by far the cleanest separator of our data into lin-
14 guistic versus non-linguistic. If one sets a value of 0.10 as the boundary, only Sumerian and
15 Mesopotamian deity symbols are clearly misclassified (while Asian emoticons and Egyptian
16 are ambiguously on the border). Not surprisingly, this is the feature most often picked by the
17 decision tree classifier we will discuss in the next section.

18 One important caveat is that the repetition measure is also negatively correlated with mean
19 text length: the Pearson's correlation for mean length and $\frac{r}{R}$ is -0.49 . This makes sense given
20 that the shorter the text, the less chance there is for repetition, whereas at the same time, the more
21 chance that if there *is* repetition, the repetition will involve adjacent symbols. Of course the cor-

1 relation is not perfect, meaning that $\frac{r}{R}$ probably also reflects intrinsic properties of the symbol
2 systems. How much is length a factor? One-way analyses of variance with class (linguistic/non-
3 linguistic) as the independent variable and mean text length or $\frac{r}{R}$ as dependent variables yielded
4 the following results:

5

$$\text{Mean text length } F = 5.75 \quad p = 0.027$$

6

$$\frac{r}{R} \quad F = 15.83 \quad p = 0.0008$$

7

8 Mean text length is thus well predicted by class, with non-linguistic systems having shorter
9 mean lengths. But the repetition rate is even better predicted, suggesting that the results above
10 cannot be simply reduced to length differences.

11 To see this in another way, consider the results of computing the same repetition measures
12 over our corpora where we have artificially limited the texts to length no more than six (by
13 simply trimming each text to the first six symbols). Also, in order to make the corpora more
14 comparable, we limit the shortened corpus to the first 500 “texts” from the original corpus. As
15 we can see in Table 5, the separation is not as clean as in the case of Table 4. Nevertheless,
16 the five corpora with the highest $\frac{r}{R}$ are non-linguistic (if one counts the Indus system as non-
17 linguistic) and the nine corpora with the lowest values are all linguistic. Again this measure is a
18 far better separator than any of the models that we have discussed previously, which are based
19 ultimately on probabilities.

4.6 Two-Way Classification Using Decision Trees

This study has collected a variety of features that might in principle be useful for determining whether an unknown symbol system is linguistic or non-linguistic. Some of these features have been explored in previous literature, others are novel here. A common technique in computational linguistics is to combine features, each of which may on its own not be a very good predictor, into a discriminative model. To that end I used the above features, along with the Classification and Regression Tree (CART) algorithm (61) to train (binary branching) classification trees for a binary classifier (linguistic vs. non-linguistic). The features used were:

- The PMI association measure over the symbol set comprising 25% of each corpus
- Block entropy (2) calculations for $N = 1$ to $N = 6$ (block 1 ... block 6).
- Maximum conditional entropy (I) calculated for Figure 7.
- F_2 , the \log_2 bigram conditional entropy for the whole corpus, from (3).
- U_r from (3).
- C_r from (3).
- The repetition measure $\frac{r}{R}$.

Since the set of corpora is small — only 21 in total, not counting the Indus bar seals — I tested the system by randomly dividing the corpora into 16 training and 5 test corpora, building, pruning and testing the generated trees, and then iterating this process 100 times. The mean

1 accuracy of the trees on the held out data sets was 0.8, which is above the baseline accuracy
2 0.66 of always predicting a system to be linguistic (i.e., picking the most frequent choice, which
3 in every case was to classify the system as linguistic).

4 It is interesting to see how these 100 trees classify the Indus bar seals, which were held
5 out from the training sets. 98 of the trees classified it as *non-linguistic*, and only 2 of the trees
6 classified it as linguistic. Furthermore, the 98 that classified it as non-linguistic had a higher
7 mean accuracy — 0.81 — than the 2 that classified it as non-linguistic (0.3). Thus, more and
8 better classifiers tend to classify the Indus symbols as non-linguistic than those that classify it
9 as linguistic.¹⁰

10 Since the various runs involve different divisions into training and test corpora, an obvious
11 question is whether particular corpora are associated with better performance. If corpus X is in
12 the training and thus not in the test set, does this result in better performance, either because
13 its features are more informative for a classifier and thus helpful for training, or because it is
14 harder to classify and thus helps performance if it is not part of the test set? Tables 6 and 7
15 provide some results on this. Asian emoticons and Sumerian result in higher classification rates
16 when they occur in the training and not in the testing, perhaps because both are hard to classify
17 in testing: Sumerian is misclassified by the repetition measure, as we saw above, and emoticons
18 on that measure are close to the border with linguistic systems, as is Egyptian, which also results

¹⁰To show that this result is no artifact of the particular setup we are using, I ran the same experiment, this time dropping the Indus bar seals entirely, and holding out Oriya from the training. 100% of the trees classified Oriya as linguistic.

1 in better performance if it is not in the test data. The next two on the list, weather icons and
2 Pictish symbols, have a repetition value that places them well outside the range of the linguistic
3 corpora, so in these cases it may be that they are useful as part of training to provide better
4 classifiers. Mesopotamian deity symbols are also misclassified by the repetition measure, and
5 thus could lead to better performance when they are not part of the test data.

6 In what was just described, I included Pictish among the non-linguistic systems. Of course,
7 with the publication of (3), this classification has become somewhat controversial. What if one
8 does not include the Pictish data? In this case I had training sets of 15 corpora, and again held
9 out 5 corpora for testing. Here the overall mean accuracy is 0.84. What do these classifiers
10 make of Pictish? The results are strongly in favor of the *non-linguistic* hypothesis: 97 classified
11 Pictish as non-linguistic, with a mean accuracy of 0.81; 3 classified it as linguistic, with a lower
12 mean accuracy of 0.4.

13 It is also interesting to consider the features that are used by the trees. These are presented
14 in Table 8 for the two experimental conditions (with versus without Pictish). In both cases,
15 the most prominent feature was repetition, and the second most common is C_r . The latter is
16 interesting since if one compares Lee et al.'s tree in Figure 9, it is C_r which was involved in
17 the top-level linguistic/non-linguistic decision. This may seem surprising confirmation of Lee
18 et al.'s results, given that I showed in Section 4.3 that C_r and U_r , with Lee et al.'s settings,
19 systematically misclassified nearly all of our non-linguistic corpora. What this means, however
20 is that while their own settings — the ones they used to arrive at the linguistic hypothesis for
21 Pictish symbols — are probably not useful, the measure itself is somewhat useful, if one is

1 allowed to retrain the boundary between linguistic and non-linguistic for C_r .

2 To see the importance of the various measures in another way, Table 9 shows the results
3 of the Wilcoxon signed rank test for each feature comparing for the two populations of lin-
4 guistic and non-linguistic corpora. Only for the repetition measure, and C_r are the population
5 means different according to this test, suggesting that the other features are largely useless for
6 determining the linguistic status of a symbol system.

7 Since the repetition measure is well correlated with mean length of texts in the corpus, and
8 the non-linguistic corpora in general have shorter mean lengths than the linguistic corpora, it is
9 instructive to consider what happens when we remove that feature and train models as before.
10 As we can see in Table 10, a wider variety of trees is produced, and Lee and colleagues' C_r
11 is the favored feature, being used in 82 of the trees. The results for the held out Indus bar
12 seal corpus are not as dramatic as when I included the repetition feature, but they still highly
13 favor the non-linguistic analysis. 88 of the trees classified the system as non-linguistic (mean
14 accuracy 0.75), and 12 as linguistic (mean accuracy 0.37).

15 Similarly lower results were obtained when I replaced the repetition rate $\frac{r}{R}$ from Table 4
16 with that computed over the artificially shortened corpora seen above in Table 5, 85 of the
17 trees classified the Indus symbols as non-linguistic (mean accuracy 0.63), and 15 classified it
18 as linguistic (mean accuracy 0.36). The features used by these trees are shown in Table 11.
19 Repetition is far less dominant than it is in the original tree set shown in Table 8, but it is still
20 the second most used feature, after C_r , occurring in 35 out of 100 trees.

1 The most useful features thus seem to be the measure of repetition $\frac{r}{R}$ and C_r ; again this is
2 further confirmed by the Wilcoxon signed rank test reported above for which these were the
3 only two features that showed a significant correlation with corpus type.

4 What do these two features have in common? We know that $\frac{r}{R}$ is correlated with text length:
5 could it be that C_r is also correlated? As Figure 12 shows, this is indeed the case. Pearson's
6 r for C_r and text length is 0.39, and this increases dramatically to 0.71 when the one obvious
7 outlier — Amharic — is not considered.

As we argued above, there is a plausible story for why $\frac{r}{R}$ should correlate with text length,
but why should C_r correlate? Recall the formula for C_r , repeated here:

$$C_r = \frac{N_d}{N_u} + a \frac{S_d}{T_d} \quad (7)$$

8 where, again, N_d is the number of bigram types, N_u the number of unigram types, S_d is the
9 number of bigram hapax legomena, and T_d is the total number of bigram tokens. Now, Lee et
10 al. pad their texts with beginning and end of text delimiters (as the OpenGrm software we use
11 does implicitly). For corpora consisting of shorter texts, this means that bigrams that include
12 either the beginning or the ending tag will make up a larger portion of the types, and that the
13 type richness will be reduced. The unigrams, on the other hand, will be unaffected by this,
14 though they will of course be affected by the overall corpus size. This predicts that the term $\frac{N_d}{N_u}$
15 alone should correlate well with mean text length, and indeed it does, with $r = 0.4$ ($r = 0.72$
16 excluding Amharic).

1 U_r , which is derived from $\frac{N_d}{N_u}$ is also somewhat correlated, but negatively, and not so
2 strongly: $r = -0.29$. Thus a higher mean text length corresponds to a lower value for U_r .
3 Recall again that it is this feature that is used in Lee and colleagues' (3) decision tree for clas-
4 sifying the type of linguistic system: the lowest values of U_r are segmental systems, next syl-
5 labaries, next "logographic" systems. This makes sense in light of the (weak) correlation with
6 text length: *ceteris paribus*, it takes more symbols to convey the same message in a segmen-
7 tal system than in a syllabary, and more symbols in a syllabary than in a logographic system.
8 Thus segmental systems should show a lower U_r , syllabic systems a higher U_r and logographic
9 systems the highest U_r values.

10 Returning to C_r , non-linguistic systems do tend to have shorter text lengths than linguistic
11 systems, so one reason why C_r seems to be such a good measure for distinguishing the two
12 types, apparently, is because it correlates with text length. Of course where one draws the
13 boundary between linguistic and non-linguistic on this scale will determine which systems get
14 classified in which ways. For Lee and colleagues, Pictish comes out as linguistic only because
15 they set the value of C_r relatively low.¹¹

16 While there are obviously other factors at play—for one thing, something must explain why
17 Amharic is such an outlier — it does seem that the key insight at work here comes down to a

¹¹We have already seen that repetition on its own misclassifies Sumerian as non-linguistic, due to the short
"texts" — an artifact of the way we divided the corpus. Given what we have just discussed, we would therefore
expect that the decision trees would also misclassify it. Indeed they do: using all features, 92 trees classified it as
non-linguistic (accuracy 0.72) and 8 trees as linguistic (0.23).

1 rather simple measure: how long on average are the extant texts in the symbol system? If they
2 are short, then they are probably non-linguistic; if they are long they are probably linguistic.

3 **5 Conclusions**

4 This work has used a set of corpora of non-linguistic symbol systems that we have developed
5 and compared them using a variety of statistical measures to a large set of corpora of written
6 language. The non-linguistic corpus set was designed to include types of systems that are likely
7 to be relevant to assessing whether ancient systems such as the Indus Valley signs, or Pictish
8 symbols, are linguistic or not. I also included one system (barn stars) that is almost certainly
9 purely decorative (thus addressing a common misconception about the claims of (49) to the
10 effect that we were claiming the Indus symbols were decorative); as well as some modern
11 systems that can easily be harvested from the Web. Finally, I picked a wide variety of linguistic
12 symbol systems, ranging across both time and script type.

13 I have examined a variety of statistical measures and shown that two of them — Lee and
14 colleagues' (3) C_r and the repetition measure $\frac{r}{R}$ — are useful for distinguishing linguistic and
15 non-linguistic systems. However, the fact that they are useful seems to be in part because both
16 of them are also correlated with a measure that is far more basic and trivial: mean text length.

17 In both depth and breadth of coverage, this work arguably goes far beyond previous work
18 in this area, but the apparent conclusion — that a core distinction between linguistic and at
19 least some non-linguistic systems — comes down to text length, seems troublesome. Could the

1 answer be that trivial?

2 I believe the answer to that question is yes, for a simple reason: if you have a real writing
3 system, it allows you to write anything you want to say, meaning that you can generate very long
4 texts. Non-linguistic systems can be expected to be much more varied in this regard, depending
5 upon what kind of information they represent. Mathematical formulae can, indeed, be quite
6 long, as can deity symbol strings on *kudurru* stones. However Pictish “texts” are all very short,
7 presumably because whatever kind of information was represented by those “texts” was by
8 its nature concise. Totem pole texts are short for similar reasons (and probably also because
9 of the labor involved in carving whole tree trunks). These characterizations are informal and
10 somewhat circular to be sure, but the point is that with non-linguistic systems, unlike linguistic
11 systems, there is no *a priori* expectation that the texts should be long.

12 One of the main problems all along for the theory that the Indus Valley symbols were a
13 written language has been the extremely short length of the “texts”. While few Indus researchers
14 would perhaps admit this, the brevity of the inscriptions is an embarrassment. Why else would
15 researchers since Marshall (62) have appealed to a “lost manuscript” hypothesis if it were not
16 the sense that a symbol system that was only used for very short inscriptions did not look much
17 like a writing system? The “lost manuscript” hypothesis and its ramifications were discussed
18 extensively in (49), and those arguments will not be repeated here. But one of the conclusions
19 of that work was that one discovery that would be most needed to support the script hypothesis
20 would be a genuine long text in the symbols. This claim is supported by the statistical results
21 we have obtained here.

1 At its core this just seems like common sense, and there is a clear analogy in child language
2 development. The most basic measure of language development is *mean length of utterance*
3 (MLU), measured in words, or in morphemes. There are of course other measures, but often
4 even much more sophisticated measures correlate very strongly with MLU: for example the
5 *Index of Productive Syntax* — IPSyn (63), a complex set of syntactic and morphological features
6 used in assessing language development for English-speaking children, has an extremely high
7 correlation with MLU in Scarborough’s original data. That MLU should be so basic seems
8 sensible: if an individual only utters four-word sentences, we would not generally think of that
9 individual as having mastered his or her native language.

10 And so it is with symbol systems. If a symbol system is a true writing system, then it is
11 able to represent anything that can be said in the language of its users. Language is in principle
12 unbounded in that there are no limitations, apart from obvious physical ones, on how long a
13 linguistic message can go on. Linguistic texts can range from short texts such as a couple of
14 names to lengthy stories that can extend into the hundreds of thousands of words or more. Texts
15 in a true writing system should thus show a wide range of lengths *because of what it is they*
16 *encode*. Indeed, in societies that make extensive use of writing, written texts can often be much
17 longer than any given spoken message, precisely because their creation is not limited by the
18 same factors, such as breathing or fatigue, that limit spoken utterances. If the Indus Valley
19 system was a full writing system as some have claimed, it would be quite anomalous if the
20 Indus Valley scribes had not figured out, over a span of roughly 700 years, that they could do a
21 bit more with the system than write short cryptic texts (49).

1 Needless to say, text length, while a plausible predictor of symbol-system status, is nonethe-
2 less a crude feature. Non-linguistic systems *can* still have relatively long texts, as in some of
3 the corpora we have collected. And while our useful features seem to reduce ultimately to text
4 length, it is not out of the question that in future work we will discover other features that are
5 useful in classifying symbol systems, and that are not correlated with length.

6 For contentious systems — in particular the Indus Valley symbols — we do not expect the
7 debate to end here. For the Indus Valley symbols in particular, too many people have a lot
8 invested in the idea that the symbols were writing to give that idea up. I hope however that I
9 have provided an analysis of statistical approaches to the question of symbol system type that
10 is better informed than what has been prominently displayed in the past few years.

11 **Acknowledgments**

12 A number of people have given input of one kind or another to this work.

13 First and foremost, I would like to thank my research assistants Katherine Wu, Jennifer
14 Solman and Ruth Linehan, who worked on the development of many of the corpora and markup
15 scheme. Katherine Wu developed the corpora for Totem poles, Mesopotamian deity symbols,
16 Vinča and Pictish, and Ruth Linehan the Heraldry corpus. Jennifer Solman developed the XML
17 markup system used in most of our corpora. An initial report on our joint work was presented
18 in (46).

19 I also thank audiences at the Linguistic Society of America 2012 annual meeting in Port-

1 land, the Center for Spoken Language Understanding, the Oriental Institute at the University of
2 Chicago, Kings College London, Johns Hopkins University and the University of Maryland, for
3 feedback and suggestions on presentations of this work.

4 Chris Woods gave me useful suggestions for Sumerian texts to use. William Boltz and Ken
5 Takashima were instrumental in pointing me to the corpus of transcribed Oracle bones.

6 I would like to acknowledge the staff of the Berks County Historical Society, in particular
7 Ms. Kimberly Richards-Brown, for their help in locating W. Farrell's slide collection of barn
8 stars.

9 I thank Brian Roark, Shalom Lappin, Suyoun Yoon and especially Steve Farmer for lots of
10 useful discussion.

11 This material is based upon work supported by the National Science Foundation under grant
12 number BCS-1049308. Any opinions, findings, and conclusions or recommendations expressed
13 in this material are those of the author and do not necessarily reflect the views of the National
14 Science Foundation.

15 **References**

- 16 1. R. Rao, *et al.*, *Science* **324**, 1165 (2009).
- 17 2. R. Rao, *IEEE Computer* **43**, 76 (2010).
- 18 3. R. Lee, P. Jonathan, P. Ziman, *Proceedings of the Royal Society A: Mathematical, Physical*
19 *& Engineering Sciences* pp. 1–16 (2010).

- 1 4. I. Gelb, *A Study of Writing: The Foundations of Grammatology* (University of Chicago,
2 Chicago, 1952).
- 3 5. G. Sampson, *Writing Systems* (Stanford University Press, Stanford, CA, 1985).
- 4 6. J. DeFrancis, *Visible Speech: The Diverse Oneness of Writing Systems* (University of
5 Hawaii Press, Honolulu, HI, 1989).
- 6 7. P. Daniels, W. Bright, eds., *The World's Writing Systems* (Oxford, New York, 1996).
- 7 8. R. Sproat, *A Computational Theory of Writing Systems* (Cambridge University Press, Cam-
8 bridge, 2000).
- 9 9. H. Rogers, *Writing Systems: A Linguistic Approach* (Blackwell, Malden, MA, 2005).
- 10 10. R. Harris, *Signs of Writing* (Routledge, London, 1995).
- 11 11. R. Sproat, *Non-Linguistic Symbol Systems: A Taxonomy and Analysis of their Statistical*
12 *Properties* (TBD, 2014).
- 13 12. R. Lee, P. Jonathan, P. Ziman, *Computational Linguistics* **36**, 791 (2010).
- 14 13. R. Sproat, *Computational Linguistics* **36** (2010).
- 15 14. B. Powell, *Writing: Theory and History of the Technology of Civilization* (Wiley-Blackwell,
16 Chichester, 2009).
- 17 15. C. Waddington, *American Ethnologist* **1**, 471 (1974).

- 1 16. M. Barbeau, *Totem Poles*, vol. 1–2 of *Anthropology Series 30*, National Museum of Canada
2 *Bulletin 119* (National Museum of Canada, Ottawa, 1950).
- 3 17. D. Kaiser, *American Scientist* **93**, 156 (2005).
- 4 18. N. Le Novère, et al, *Nature Biotechnology* **27**, 735 (2009).
- 5 19. K. Basso, N. Anderson, *Science* **180** (1973).
- 6 20. J. McCawley, *The World's Writing Systems*, P. Daniels, W. Bright, eds. (Oxford, New York,
7 1996), pp. 847–854.
- 8 21. B. Farnell, *The World's Writing Systems*, P. Daniels, W. Bright, eds. (Oxford, New York,
9 1996), pp. 855–879.
- 10 22. G. Mallery, *Pictographs of the North American Indians: A Preliminary Paper*, no. 4 in An-
11 nual Report of the Bureau of Ethnology (Smithsonian Institution, Washington, DC, 1883).
12 Available from Google Books.
- 13 23. L. Li, *Naxizu xiangxing biaoyin wenzi zidian* (Yunnan Minzu Chubanshe, 2001).
- 14 24. T. E. Graves, *The Pennsylvania German hex sign: A study in folk process*, Ph.D. thesis,
15 University of Pennsylvania, Philadelphia, PA (1984).
- 16 25. S. Bedrick, R. Beckley, B. Roark, R. Sproat, *Workshop on Language and Social Media*
17 (Montreal, Canada, 2012).

- 1 26. U. Seidl, *Die babylonischen Kudurru- Reliefs. Symbole mesopotamischer Gottheiten* (Uni-
2 versitätsverlag Freiburg, 1989).
- 3 27. S. M. Winn, *Pre-Writing in Southeastern Europe: The Sign System of the Vinča Culture,*
4 *ca. 4000 B.C* (Western Publishers, 1981).
- 5 28. A. Jackson, *The Symbol Stones of Scotland: a Social Anthropological Resolution of the*
6 *Problem of the Picts* (Orkney Press, 1984).
- 7 29. Royal Commission on the Ancient and Historical Monuments of Scotland, Pictish symbol
8 stones: A handlist (1994).
- 9 30. A. Jackson, *The Pictish Trail* (Orkney Press, 1990).
- 10 31. A. Mack, *Field Guide to the Pictish Symbol Stones* (The Pinkfoot Press, Balgavies, Angus,
11 1997). Updated 2006.
- 12 32. E. Sutherland, *The Pictish Guide* (Birlinn Ltd, Edinburgh, 1997).
- 13 33. British Columbia Provincial Museum, British Columbia totem poles, Printed by Charles F.
14 Banfield printer to the King's most excellent majesty (1931).
- 15 34. V. E. Garfield, *The Seattle Totem Pole* (University of Washington Press, Seattle, 1940).
- 16 35. S. W. Gunn, *The Totem Poles in Stanley Park, Vancouver, B.C.* (Macdonald, Vancouver,
17 BC, 1965), second edn.

- 1 36. S. W. Gunn, *Kwakiutl House and Totem Poles at Alert Bay* (Whiterocks Publications, Van-
2 couver, BC, 1966).
- 3 37. S. W. Gunn, *Haida Totems in Wood and Argillite* (Whiterocks Publications, Vancouver, BC,
4 1967).
- 5 38. F. Drew, *Totem Poles of Prince Rupert* (F.W.M. Drew, Prince Rupert, BC, 1969).
- 6 39. E. Malin, *Totem Poles of the Pacific Northwest Coast* (Timber Press, Portland, OR, 1986).
- 7 40. City of Duncan, *Duncan: City of Totems*. (1990). Duncan, BC.
- 8 41. H. Stewart, *Totem Poles* (University of Washington, Seattle, WA, 1990).
- 9 42. H. Stewart, *Looking at Totem Poles* (University of Washington, Seattle, WA, 1993).
- 10 43. R. D. Feldman, *Home before the Raven Caws: the Mystery of the Totem Pole*, Cincinnati,
11 OH (Emmis Books, 2003).
- 12 44. A. C. Mahr, *Ohio History: The Scholarly Journal of the Ohio Historical Society* **54**, 1
13 (1945).
- 14 45. D. Yoder, T. Graves, *Hex Signs: Pennsylvania Dutch Barn Symbols and their Meaning*
15 (Stackpole, Mechanicsburg, PA, 2000).
- 16 46. K. Wu, J. Solman, R. Linehan, R. Sproat, *Linguistic Society of America* (Portland, OR,
17 2012).

- 1 47. J. Rhys, *Proceedings of the Society of Antiquaries of Scotland* **26**, 263 (1892).
- 2 48. R. Rao, *et al.*, *Proceedings of the National Academy of Sciences* **106**, 13685 (2009).
- 3 49. S. Farmer, R. Sproat, M. Witzel, *Electronic Journal of Vedic Studies* **11** (2004).
- 4 50. M. Vidale, *East and West* **57**, 333 (2007).
- 5 51. A. Parpola, *Deciphering the Indus Script* (Cambridge University Press, New York, 1994).
- 6 52. M. Ventris, J. Chadwick, *Documents in Mycenaean Greek* (Cambridge University Press,
7 Cambridge, 1956).
- 8 53. B. Roark, *et al.*, *Proceedings of the Association for Computational Linguistics, System*
9 *Demonstrations* (Association for Computational Linguistics, Jeju Island, South Korea,
10 2012), pp. 61–66.
- 11 54. C. Shannon, *Bell System Technical Journal* **27**, 379 (1948).
- 12 55. S. M. Winn, *The Life of Symbols*, M. L. Foster, ed. (Westview Press, Boulder, 1990), pp.
13 269–271.
- 14 56. R. Sproat, *Computational Linguistics* **36**, 807 (2010).
- 15 57. R. Kneser, H. Ney, *International Conference on Speech and Signal Processing* (1995), pp.
16 181–184.
- 17 58. I. Nemenman, F. Shafee, W. Bialek, *Advances in Neural Information Processing Systems*
18 (MIT Press, Cambridge, MA, 2002), vol. 14.

- 1 59. K. W. Church, W. Gale, *Computer Speech and Language* **5**, 19 (1991).
- 2 60. C. Manning, H. Schütze, *Foundations of Statistical Natural Language Processing* (MIT
3 Press, Cambridge, MA, 1999).
- 4 61. L. Breiman, J. Friedman, R. Olshen, C. Stone, *Classification and Regression Trees*
5 (Wadsworth and Brooks, Pacific Grove, CA, 1984).
- 6 62. J. Marshall, *Illustrated London News* pp. 528–32, 548 (1924).
- 7 63. H. Scarborough, *Applied Psycholinguistics* **11**, 1 (1990).



Figure 1: 6-pointed star (left) versus rosette.

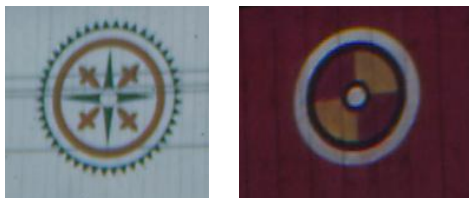


Figure 2: 4-pointed star (left) versus 4-slice pie.

PDX 2011/4/29 3:0:50



Figure 3: A sample weather icon sequence: forecast for Portland, Oregon, April 29, 2011.

| ^ _ ^ | (๑ ̇ ๑)
[๐ _ -] (̇ ̇)
\ (^ v ^) / (* ≧ m ≦ *)

Figure 4: Some representative *kaomoji* emoticons

Corpus	# Texts	# Tokens	# Types	Mean text length
Asian emoticons	10,000	59,186	334	5.9
Barn stars	310	963	32	3.1
Mesopotamian deity symbols	69	939	64	13.6
Pictish stones	283	984	104	3.5
Totem poles	325	1,798	477	5.5
Vinča	591	804	185	1.4
Weather icons	10,142	50,710	16	5.0

Table 1: Number of texts, type and token counts, and mean text length for the non-linguistic corpora.

Corpus	# Texts	# Tokens	# Types	Mean text length
Amharic	111	17,747	219	159.9
Arabic	10,000	429,463	62	42.9
Chinese	10,000	136,246	2,738	13.6
Ancient Chinese	22,359	91,010	701	4.1
Egyptian	1,259	38,804	691	30.8
English	10,000	503,309	76	50.3
Hindi	1,000	61,254	77	61.2
Korean	10,000	223,869	984	22.4
Korean (Jamo)	10,000	438,440	102	43.8
Linear B	439	5,465	220	12.4
Malayalam	937	46,497	80	49.6
Oriya	1,000	40,242	75	40.2
Sumerian	11,528	50,318	692	4.4
Tamil	1,000	38,455	61	38.5

Table 2: Linguistic corpora for comparison.

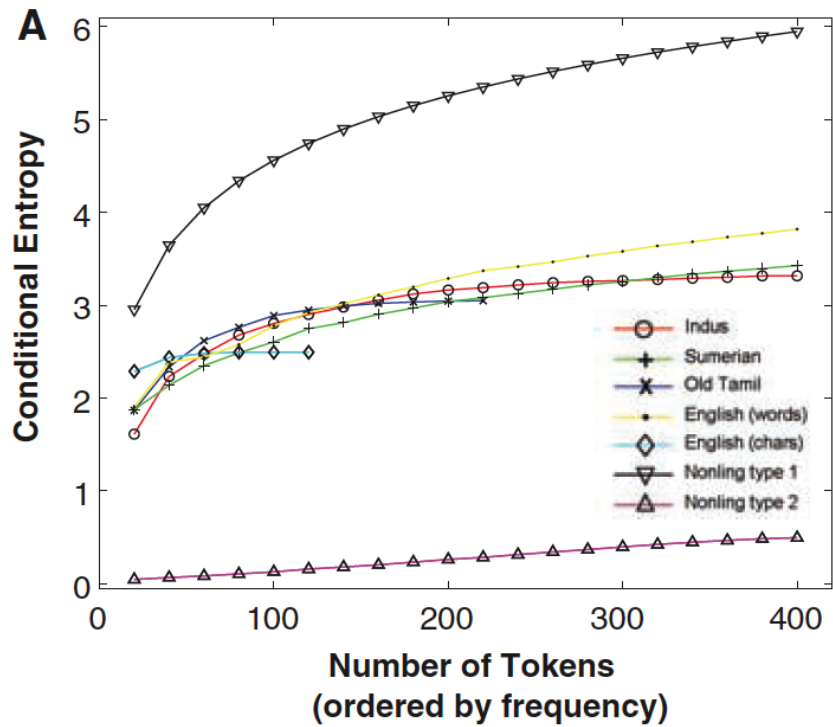


Figure 5: Bigram conditional entropy growth curves of various r linguistic symbol systems, the Indus Valley symbols, and two artificial models of non-linguistic systems, from (1). See the text for further explanation. Used with permission of the American Association for the Advancement of Science.

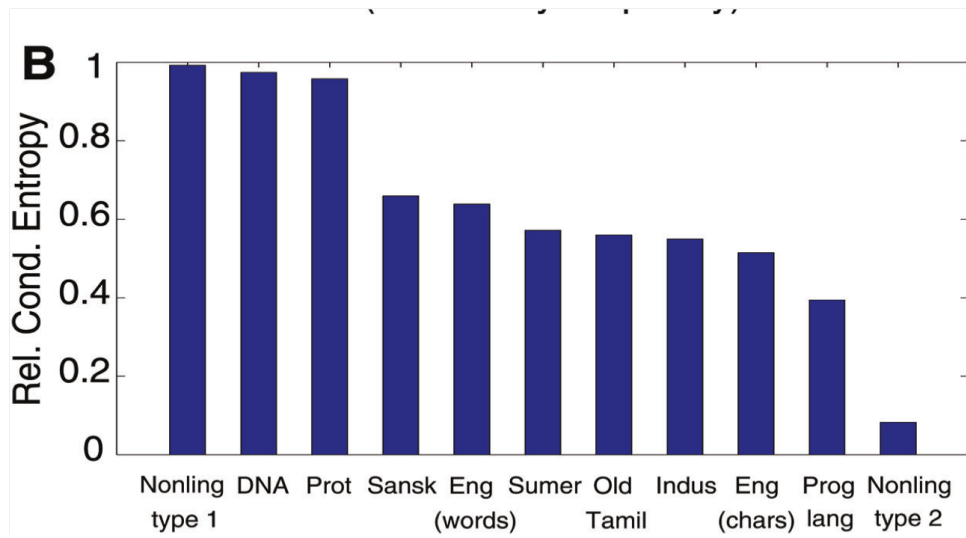


Figure 6: Relative conditional entropy (conditional entropy relative to a uniformly random sequence with the same number of tokens) of various systems from (1). Used with permission of the American Association for the Advancement of Science.

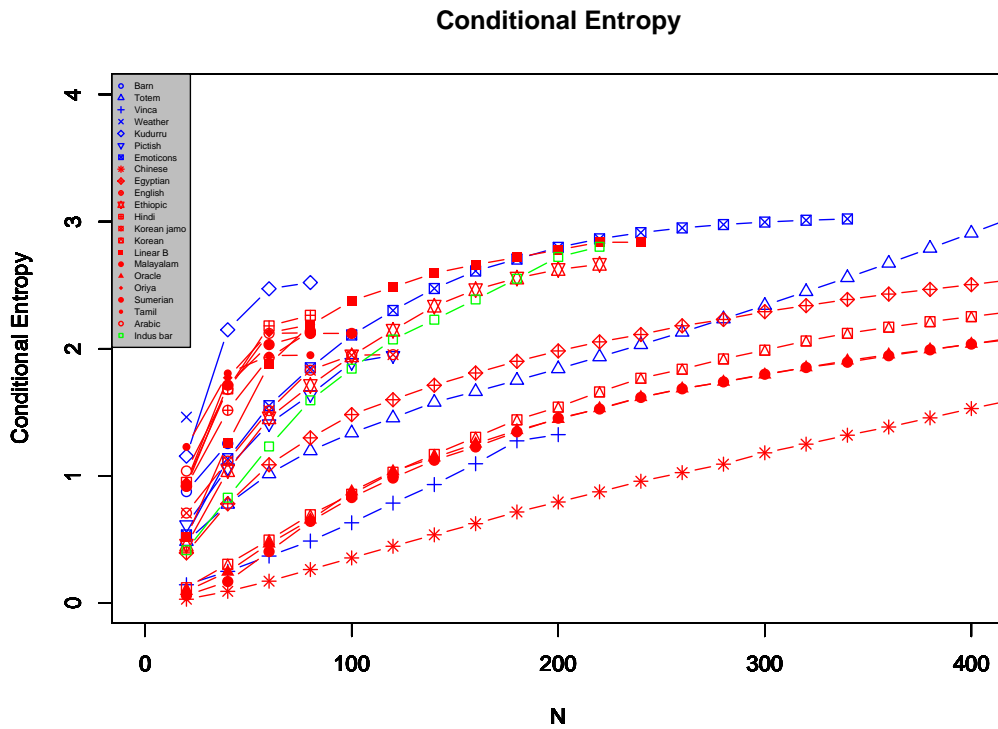


Figure 7: Bigram conditional entropy growth curves computed for our linguistic (red) and non-linguistic (blue) corpora, and for Indus bar seals (green).

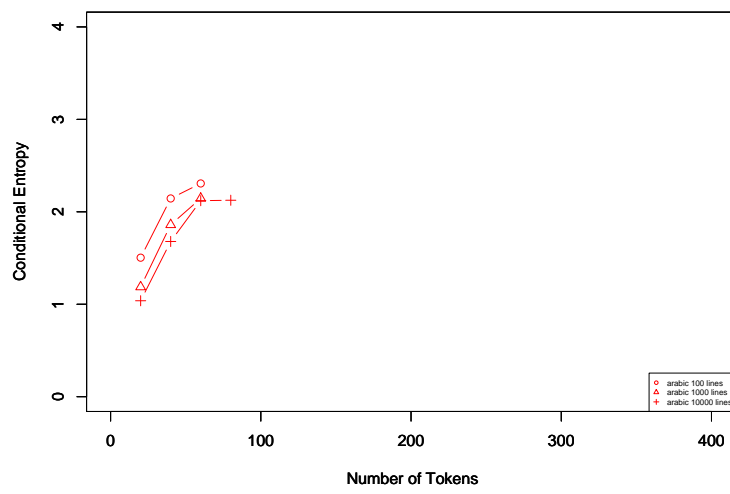


Figure 8: Bigram conditional entropy growth curves for three different sized Arabic corpora.

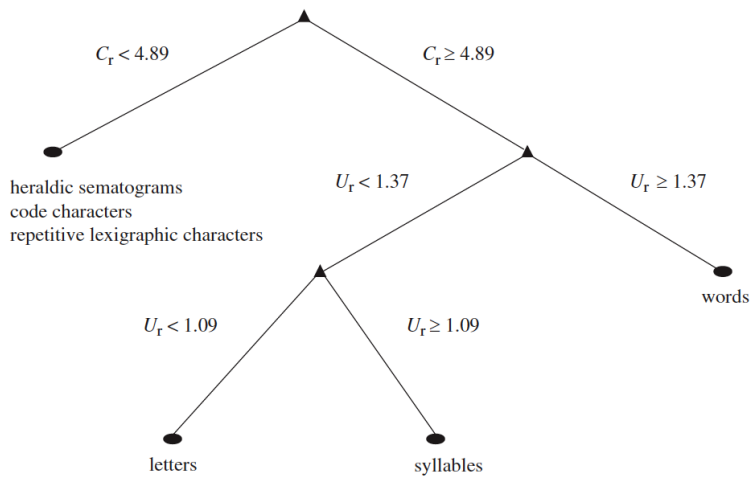


Figure 9: Reproduction of Figure 6, page 9, from: (3).

Corpus	Classification
Asian emoticons	linguistic: letters
Barn stars	linguistic: letters
Mesopotamian deity symbols	linguistic: syllables
Pictish symbols	linguistic: words
Totem poles	linguistic: words
Vinča symbols	non-linguistic
Weather icon sequences	linguistic: letters

Table 3: The results of applying Lee et al.'s (3) decision tree from Figure 9 to our data.

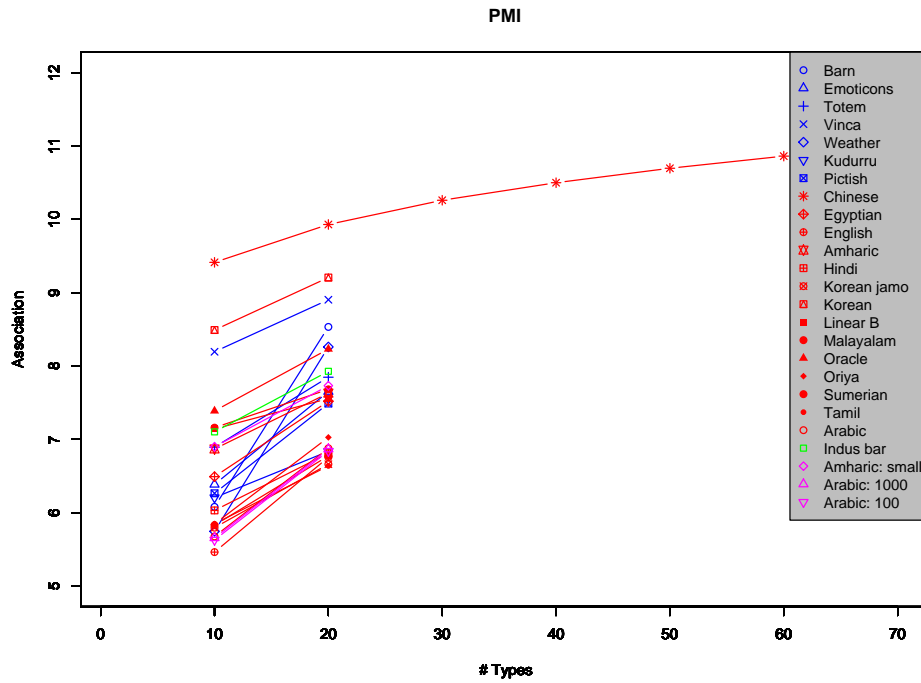


Figure 10: PMI association computations for our corpora, computed over subsets of each corpus starting with the 10 most frequent symbols, the 20 most frequent, and so forth, up to 25% of the corpus. See the text for more detailed explanation.

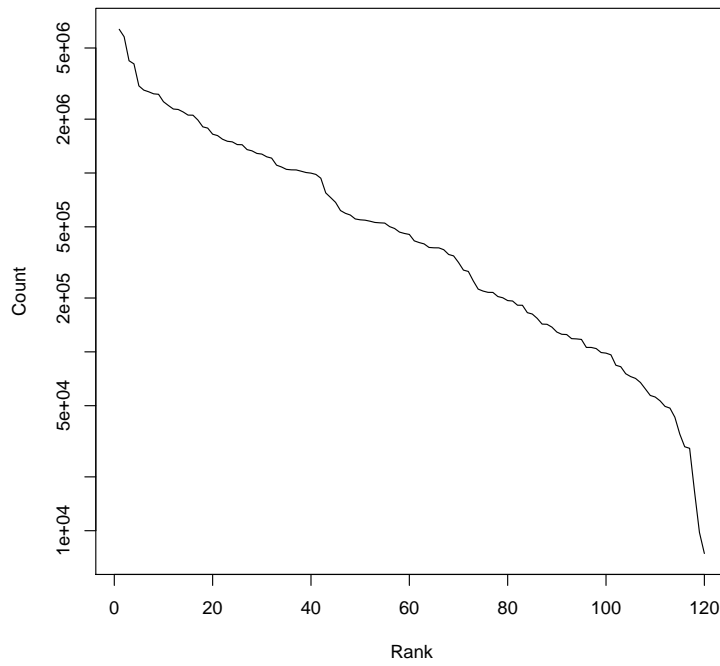


Figure 11: Unigram distribution of awarded boyscout merit badges, from http://meritbadge.org/wiki/index.php/Merit_Badges_Earned, showing the standard power-law inverse log-linear relationship between count and rank.

Corpus	$\frac{r}{R}$
Barn Stars	0.86
Weather Icons	0.79
<i>Sumerian</i>	0.67
Totem Poles	0.63
Vinča	0.59
Indus bar seals	0.58
Pictish	0.26
Asian emoticons	0.10
Egyptian	0.10
<i>Mesopotamian deity symbols</i>	0.099
Linear B	0.055
Oracle Bones	0.048
Chinese	0.048
English	0.035
Arabic	0.032
Korean jamo	0.022
Malayalam	0.022
Korean	0.020
Amharic	0.018

Oriya	0.0075
Tamil	0.0060
Hindi	0.0017

Table 4: Repetition rate $\frac{r}{R}$ for the various corpora.

Corpus	$\frac{r}{R}$
Barn Stars	0.85
Weather Icons	0.80
Indus bar seals	0.77
Totem Poles	0.71
Vinča	0.63
Egyptian	0.42
Sumerian	0.33
Chinese	0.31
Pictish	0.29
English	0.29
Mesopotamian diety symbols	0.287
Amharic	0.25
Asian emoticons	0.22
Linear B	0.21
Malayalam	0.19
Arabic	0.14
Korean jamo	0.08
Oriya	0.04
Korean	0.015

Hindi	0.015
Tamil	0.0040
Oracle Bones	0.0

Table 5: Repetition rate $\frac{r}{R}$ for versions of the corpora with artificially shortened texts of length 6 or less, and no more than 500 texts.

Corpus	Mean accuracy	# training runs
Asian Emoticons	0.85	79
Mesopotamian deity symbols	0.85	71
Sumerian	0.84	76
Egyptian	0.83	72
Weather icons	0.82	80
Pictish	0.82	79
Malayala	0.81	79
Hindi	0.80	78
Oriya	0.80	82
Korean jamo	0.80	77
Arabic	0.80	74
Linear B	0.79	70
English	0.79	76
Tamil	0.79	76
Amharic	0.79	73
Totem poles	0.79	86
Chinese	0.79	75
Korean	0.79	72
Oracle bones	0.78	69

Vinča	0.78	76
Barn stars	0.78	80

Table 6: Mean accuracy of results for each of the corpora when occurring in the training data.

Corpus	Mean accuracy	# training runs
Barn stars	0.91	20
Totem poles	0.89	14
Vinča	0.88	24
Oracle bones	0.85	31
Chinese	0.85	25
Korean	0.84	28
English	0.84	24
Tamil	0.84	24
Amharic	0.84	27
Linear B	0.83	30
Korean jamo	0.83	23
Arabic	0.82	26
Oriya	0.82	18
Hindi	0.81	22
Malayala	0.80	21
Pictish	0.75	21
Weather icons	0.74	20
Egyptian	0.74	28
Mesopotamian deity symbols	0.70	29

Sumerian	0.69	24
Asian Emoticons	0.64	21

Table 7: Mean accuracy of results for each of the corpora when occurring in the test data.

with Pictish		without Pictish	
# trees	features	# trees	features
84	repetition	71	repetition
13	C_r	26	C_r
1	block 5	1	block 5
1	block 3	1	association
1	block 1	1	

Table 8: Features used by trees under the two training conditions. The left column in each case is the number of trees using the particular combination of features in the right column. The final line for the trees trained without Pictish data, is for one tree with a single node that predicts “linguistic” in all cases.

Measure	W	p
association	29	0.15
block 1	42	0.64
block 2	49	1
block 3	56	0.64
block 4	49	1
block 5	63.5	0.30
block 6	56	0.64
maximum conditional entropy	58	0.54
F_2	63	0.32
U_r	40	0.54
C_r	84	0.0074**
repetition	6	0.00052**

Table 9: Wilcoxon signed rank test for each feature with the two populations being linguistic and non-linguistic corpora. Only repetition and C_r show significance.

# trees	features
57	C_r
25	C_r , association
6	association
4	maximum conditional entropy
2	block 5
2	
1	U_r
1	block 6
1	association, block 5
1	association, block 4

Table 10: Features used by trees when repetition is removed. The sixth row consists of two trees where there is a single leaf node with the decision “linguistic”.

# trees	features
36	C_r
30	repetition
11	C_r , association
8	
4	maximum conditional entropy
3	association
2	U_r
2	maximum conditional entropy, repetition
2	block 4, repetition
1	block 5
1	block 1, repetition

Table 11: Features used by trees when repetition is replaced by the repetition rate computed over artificially shortened corpora (Table 5. The fourth row consists of twelve trees where there is a single leaf node with the decision “linguistic”.

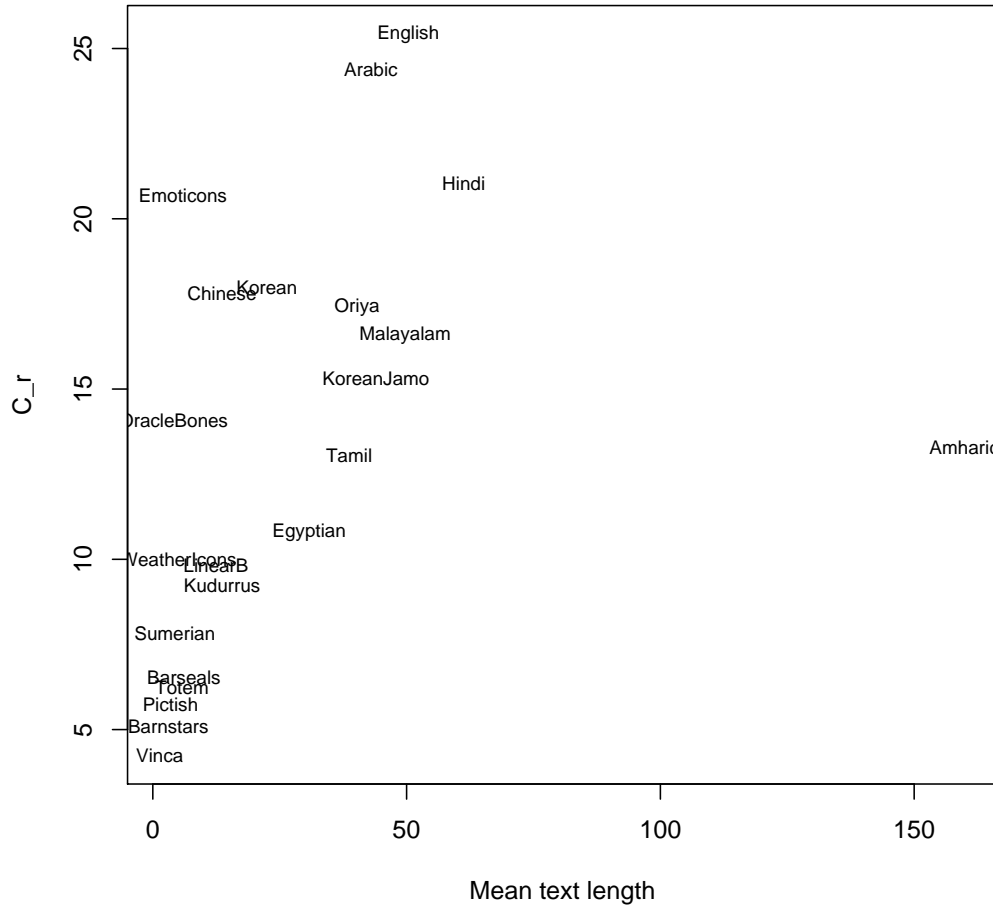


Figure 12: Correlation between C_r and text length. Pearson's r is 0.39, but when the obvious outlier Amharic is removed, the correlation jumps to 0.71.