

Natural language samples

Version 1.01. Last updated 2015-03-28.

Intro

Summary: This pack provides 52 natural language text samples in one convenient download. They are all computer friendly and contain the major languages you're looking for, so don't sweat anywhere else. ☺

It is designed to help with statistical comparisons of natural language and the Voynich Manuscript "text". This includes both types - those attempting to identify if it is written in a natural language at all, and those attempting to identify it with specific languages and features. However you can use it for anything else you want.

This pack is free to download, share and use but please read the information in section 4.

Before you do any experiments, please read [this article](#) (and its comments) for an interesting discussion on the matter. It points out potential problems you might not have thought about.

All future updates of this pack, and other resources, will be hosted at <https://briancham1994.wordpress.com/voynich-resources/> .

Contents

1. Rationale of selection

Lists the principles behind the selection 2

2. Overview of selection

Summarises the selection of languages..... 3

3. Full list

Includes details of all languages and samples included, ordered by code 5

4. Legal and authorship information

Explains the copyright status of this pack, who made it, and contact details 14

5. Works cited

References 15

1. Rationale of selection

In light of the shortcomings of previous studies, I have compiled a large selection of linguistic samples according to the following principles:

- Select a variety of writing systems to capture the peculiarities of transcription.
- Select a variety of writing styles and genres.
- Select a large variety of languages to represent the range of statistical properties in natural languages:
 - Span many geographic areas
 - Span many time periods
 - Span many major [language families](#)
 - Span many language features
- Select languages possibly related to the Voynich Manuscript:
 - Select common languages in 15th century European manuscripts.
 - Select languages directly proposed in previous studies, marginalia readings and decipherment attempts.

2. Overview of selection

Spatial and temporal summaries of the language selection are below.

2a. Geographic spread

The languages selected cover every inhabited continent, but with a focus on Western Eurasia. The map below shows the approximate location of each language as a red dot.



2b. Historical spread

The languages range from the 4th millennium BCE to the present day. The samples range from the late 3rd millennium BCE to 2014 CE - a span of over four thousand years.

2c. Writing systems

Statistics of text may be affected by both the peculiarities of the underlying information and the writing system used to record it.

This selection includes a range of types of different writing systems, rather than having one standardised representation for every language. This includes things like [alphabets](#) (e.g. Czech), [abjads](#) (e.g. Arabic), [abugidas](#) (e.g. Thai), high [orthographic depth](#) (e.g. Middle French), low orthographic depth (e.g. romanised Japanese), academic reconstructions (e.g. transliteration of Old Egyptian) and an artificially constructed system for personal use (e.g. [ASLSJ transcription](#) of American Sign Language).

2d. Technical details

All samples are saved as .txt files with [ANSI or UTF-8 encoding](#). Punctuation, letter cases and some chapter/paragraph numbers have been retained, with the justification that computer code can easily be designed to ignore these if desired.

Efforts have been taken to obtain samples close to or over the length of the Voynich Manuscript (approx. 35,000 words), but please note that this is very difficult for rare or extinct languages.

Samples that used [logographic](#) or [gestural systems](#) were transcribed into an alphabetic system for ease of computer processing and comparison. The original versions of these samples have been included wherever possible.

3. Full list

Time and location information is approximate.

Code	Language		Sample
	Name	Details	
N-AIN	Ainu	Location: Northern Japan Time: Unknown to present Family: Isolate Notes: Alphabetic.	<i>The Book of Common Prayer</i> by John Batchelor (translated 1896), with non-Ainu phrases removed. 4,540 words. Public domain.
N-ALB	Old Albanian	Location: Albania Time: 14 th century to unknown Family: Indo-European, Albanian Notes: Has a lot of consonants. Alphabetic.	First 104 chapters of <i>Meshari</i> by Gjon Buzuku (1555) 62,972 words. Public domain.
N-ARA	Arabic	Location: Middle East and North Africa Time: 6 th century to present Family: Afro-Asiatic Notes: Arabic script is an abjad, which the Voynich Manuscript's text may statistically resemble according to Reddy & Knight. Arabic's glyphs have a strong ordering system.	<i>Commentary on Anatomy in Avicenna's Canon</i> (نون اقلد احي رشت حرش) (سي فن ل ا ن ب ا) by Ibn al-Nafis (انيس ن ب ا) (13 th century) ¹ 99,065 words. Public domain.
N-ARM	Armenian (modern)	Location: Armenia Time: 18 th century to present Notes: Alphabetic. Family: Indo-European, Armenian	<i>Don Quixote</i> (<i>Ղնւ Կիիսուն</i>) by Miguel de Cervantes (1605), translated by Paul Makintyan (<i>Պողոս Մակինցյան</i>) (1982) 54,180 words. Public domain.
N-ARW	Arawak	Location: Guianas and Caribbean Time: Unknown to present Family: Arawakan Notes: The language spoken in Cuba, which is where Christopher Columbus spent some time. Notes: Alphabetic.	<i>Acts of the Apostles</i> , translated by Theodore Shultz (1850) 22,807 words. Public domain.
N-ASL	American Sign Language	Location: Anglo-America Time: 19 th century – present Notes: Signed language. Represents an example of a very artificial and "engineered" notation system.	Extract from <i>Alice in Wonderland</i> by Lewis Carroll (1865), transcribed by Thomas Stone in ASLSJ notation (2014). 461 words. © Thomas Stone . Used with permission.

¹ Sorry, my version of Microsoft Office 2013 can't combine Arabic alphabet text correctly.

N-BAQ	Basque	Location: Pyrenees Time: Unknown to present Family: Isolate Notes: Alphabetic.	Gero by Axular (1643) 78,344 words. Public domain.
N-CUR	Eurasian Curlew Birdsong	Location: Eurasia Time: Unknown to present Family: Animal language Notes: Represents an example of non-human communication, and a very artificial and “engineered” notation system.	Song of a Eurasian Curlew (no date), transcribed in ABC notation. 166 words. Public domain.
N-CZE	Czech	Location: Czech Republic Time: c. 14 th century to present Family: Indo-European, Slavic Notes: Language spoken in Prague, the first known location of the Voynich Manuscript. Alphabetic.	“Book of Genesis” and “Book of Exodus” from the <i>Kralice Bible</i> by the Unity of the Brethren (translated 1613) 54,018 words. Public domain.
N-EGY	Old Egyptian	Location: Ancient Egypt Time: 3 rd millennium BCE Family: Afro-Asiatic Notes: Represents an example of an academically constructed notation system. Originally logographic and syllabic.	<i>Autobiography of Weni</i> by Weni the Elder (late 3 rd millennium BCE), transliterated by Geoffrey Graham 1,011 words. Public domain.
N-EMY	Ch’olti’	Location: Guatemala Time: Unknown to c. 17 th century Family: Mayan Notes: Logographic or alphabetic.	<i>Bocabulario Grande</i> by Fray Francisco Morán (1695) 9,047 words. Public domain. Notes: The text is a dictionary from Spanish to Ch’olti’. Represents an example of a bilingual text.
N-ENM	Middle English	Location: England Time: 12 th to 15 th century Family: Indo-European, Germanic Notes: Writing system has high orthographic depth and a lot of inconsistency. Influenced by French. Alphabetic.	<i>The Canterbury Tales</i> by Geoffrey Chaucer (late 14 th century) 150,502 words. Public domain. Notes: Represents an example of a text with inconsistent spellings.
N-FIN	Finnish	Location: Finland Time: 15 th century to present Family: Uralic Notes: Has a lot of vowels. Alphabetic.	<i>Kalevala</i> by multiple, collated by Elias Lönnrot (1835) 67,443 words. Public domain. Note: Represents an example of a poetic writing style.

N-FRM	Middle French	<p>Location: France</p> <p>Time: 14th to 17th century</p> <p>Family: Indo-European, Romance</p> <p>Notes: Postulated to be the language of the month name marginalia in the zodiac section. Writing system has high orthographic depth. Alphabetic.</p>	<p><i>Pantagruel</i> by François Rabelais (1530)</p> <p>37,017 words. Public domain.</p>
N-GEO	Georgian	<p>Location: Georgia</p> <p>Time: 6th century to present</p> <p>Family: Kartvelian</p> <p>Notes: Alphabetic.</p>	<p><i>The Knight in the Panther's Skin</i> (<i>Vepkhistqaosani/ვეფხისტყაოსანი</i>) by Shota Rustaveli (შოთა რუსთაველი) (12th century)</p> <p>46,557 words. Public domain.</p>
N-GLE	Irish (modern) , a.k.a. Irish Gaelic	<p>Location: Ireland</p> <p>Time: 13th century to present</p> <p>Family: Indo-European, Celtic</p> <p>Notes: Alphabetic.</p>	<p><i>Constitution of the Irish State</i> by multiple (1937)</p> <p>17,157 words. Public domain.</p> <p>Note: Represents an example of a legal/formal writing style.</p>
N-GMH	Early New High German	<p>Location: Southern Germany</p> <p>Time: 14th to 17th century</p> <p>Family: Indo-European, Germanic</p> <p>Notes: A few art history parallels may link the Voynich Manuscript to Germany. Postulated to be one of the languages in the f66r marginalia and f116v marginalia.</p> <p>Notes: Alphabetic.</p>	<p>"Gospel of Matthew" and "Epistle to the Romans" from the <i>Luther Bible</i> by Martin Luther (translated 1522)</p> <p>37,537 words. Public domain.</p>
N-GRC	Ancient Greek	<p>Location: Greece</p> <p>Time: 9th to 4th century BCE</p> <p>Family: Indo-European, Hellenic</p> <p>Notes: A few have claimed parallels between the Voynich Manuscript and Greek manuscripts.</p> <p>Notes: Alphabetic.</p>	<p><i>History of Animals</i> (<i>Των περί τα ζώα ιστοριών</i>) by Aristotle (Ἀριστοτέλης) (350 BCE)</p> <p>94,864 words. Public domain.</p>
N-GRN	Guarani	<p>Location: Paraguay, Bolivia</p> <p>Time: Unknown to present</p> <p>Family: Tupani</p> <p>Notes: Alphabetic.</p>	<p><i>Explanation of Catechism in the Guarani Language</i> by Nicolas Yapuguai (1724)</p> <p>106,106 words. Public domain.</p> <p>Note: This sample uses poor quality OCR scanned text. This has been deliberately chosen to emulate a text sample with wildly inconsistent orthography.</p>

N-HAW	Hawaiian	<p>Location: Hawaii</p> <p>Time: Unknown to present</p> <p>Family: Austronesian, Polynesian</p> <p>Notes: Stolfi once briefly mentioned that the entropy of Voynich Manuscript text is most similar to Hawaiian. [unfortunately I cannot find the original source again] Very small alphabet.</p> <p>Notes: Alphabetic.</p>	<p><i>Ka Hoku o ka Pakipika</i> volume 1, numbers 1 and 2 by multiple (1861)</p> <p>51,987 words. Public domain.</p>
N-HIT	Hittite	<p>Location: Hittite Empire (present day Turkey)</p> <p>Time: 16th to 13th century BCE</p> <p>Family: Indo-European, Anatolian</p> <p>Notes: Represents an example of an academically constructed notation system. Originally logographic and syllabic.</p>	<p><i>Illuyanka</i> (CTH 321) by unknown (c. 2nd millennium BCE)</p> <p>1,781 words. Public domain.</p> <p>Notes: Represents an example of a copy of a badly degraded text. Many consecutive lines are multiple interpretations or sources of the same original line. Line and source codes have been removed.</p>
N-IND	Indonesian	<p>Location: Insular Southeast Asia</p> <p>Time: Unknown to present</p> <p>Family: Austronesian, Nuclear Malayo-Polynesian</p> <p>Notes: Alphabetic.</p>	<p><i>Penalty and Passion</i> by Merari Siregar (1920)</p> <p>40,296 words. Public domain.</p>
N-ITA	Florentine , a.k.a. Florentine Italian or Tuscan Italian	<p>Location: Florence</p> <p>Time: 14th to 19th century</p> <p>Family: Indo-European, Romance</p> <p>Notes: Related to present day Italian. A few art history parallels may link the Voynich Manuscript to northern Italy. Alphabetic.</p>	<p><i>The Divine Comedy</i> by Dante Alighieri (early 14th century)</p> <p>96,316 words. Public domain.</p>
N-JPN	Late Middle Japanese	<p>Location: Japan</p> <p>Time: 12th to 16th century</p> <p>Family: Japonic/Isolate</p> <p>Notes: Limited set of phonemes. Extremely shallow orthographic depth. Strong CVC structure. Originally logographic and syllabic script.</p>	<p><i>Chronicles of the Authentic Lineages of the Divine Emperors</i> (<i>Jinnō Shōtōki</i>/神皇正統記) by Kitabatake Chikafusa (北畠親房) (early 14th century)</p> <p>69,165 words. Public domain.</p> <p>Notes: Romanised with the Hepburn system, without capitalisation.</p>

N-KAA	Karakalpak , a.k.a. Qaraqalpaq	Location: Uzbekistan Time: 13 th century to present Family: Turkic Notes: Has similar letter frequency distribution according to Jaśkiewicz. Features vowel harmony.	Collected works of Berdaq G'arg'abay ulı (19 th century) 1,127 words. Public domain. Notes: Uses the 1991-2009 Latin Qaraqalpaq alphabet.
N-KAL	Greenlandic	Location: Greenland Time: 13 th century to present Family: Eskimo-Aleut Notes: Polysynthetic language; words can be very long. Alphabetic.	<i>Universal Declaration of Human Rights</i> by the United Nations, translated by unknown (1998) 1,046 words. Public domain.
N-KAN	Middle Kannada	Location: Southwest India Time: 13 th to 18 th century Family: Dravidian Notes: Has similar letter frequency distribution according to Jaśkiewicz. The script is made of half-letters that combine into compounds.	Collected <i>Vachana Sitya</i> (ವಚನ ಸಾಹಿತ್ಯ) of Sarvajna (ಸರ್ವಜ್ಞ) (16 th century) 12,212 words. Public domain.
N-KBD	Kabardian Circassian , a.k.a. Kabardian, Kabardino- Cherkess or East Circassian	Location: North Caucasus Time: Unknown to present Family: Northwest Caucasian Notes: Has similar letter frequency distribution according to Jaśkiewicz. The script has few vowels and many consonants, which require large multi-glyph combinations.	Various <i>Narts</i> (Нартхымэ акъыбарыхэ) by multiple (date unknown), compiled by the Kabardian Science and Research Institute (1951) 25,775 words. Public domain. Note: This sample uses poor quality OCR scanned text. This has been deliberately chosen to emulate a text sample with wildly inconsistent orthography.
N-KOR	Middle Korean	Location: Korea Time: 10 th to 16 th century Family: Koreanic/Isolate Notes: Alphabetic.	<i>Tale of Hong Gildong</i> (홍길동전) by Heo Gyun (허균) (c. 16 th century) 7,158 words. Public domain.
N-LAT	Classical Latin	Location: Roman Empire Time: 75 BCE to 3 rd century Family: Indo-European, Italic Notes: Common scientific language in Medieval Europe. Features heavy inflection. Alphabetic.	<i>Natural History</i> by Pliny the Elder (c. 77-79) 402,702 words. Public domain.
N-LIT	Lithuanian	Location: Lithuania Time: c. 11 th century to present Family: Indo-European, Baltic Notes: Alphabetic.	<i>Ways of the Ancient Lithuanian</i> by Simonas Daukantas (1845) 58,674 words. Public domain.

N-LZH	Classical Chinese	<p>Location: China</p> <p>Time: 5th century BCE to 20th century CE</p> <p>Family: Sino-Tibetan</p> <p>Notes: Tonal language. Mandarin and Voynich Manuscript words have a binomial distribution, according to Stolfi. Originally logographic.</p>	<p><i>Records of the Three Kingdoms</i> (<i>Sanguozhi/三國志</i>) volumes 1 and 2 by Chen Shou (陳壽) (3rd century) 45,478 words. Public domain.</p> <p>Notes: Romanised with Pinyin, which captures tones.</p>
N-MNC	Manchu	<p>Location: Northeast China</p> <p>Time: unknown to 19th century</p> <p>Family: Tungusic</p> <p>Notes: Alphabetic.</p>	<p><i>The Fable of Mr. Dunggwo and the Ungrateful Wolf</i> (<i>dungg'o siyanxeng jai niohe/东郭先生和狼</i>) by Ma Zhongzi (馬中錫) (c. 15th or 16th century) 586 words. Public domain.</p> <p>Notes: Represented in the Latin alphabet since the native Manchu alphabet is difficult for computers to deal with.</p>
N-MON	Classical Mongolian	<p>Location: Mongolia</p> <p>Time: c. 12th to 17th century</p> <p>Family: Mongolic</p> <p>Notes: Alphabetic.</p>	<p><i>The Secret History of the Mongols</i> by unknown (13th century), transcribed by Sergei Kozin into the Latin alphabet 28,289 words. Public domain.</p>
N-MWP	Kalaw Lagaw Ya , a.k.a. Western Torres Strait Language	<p>Location: Australia and Torres Strait Islands</p> <p>Time: Unknown to present</p> <p>Family: Pama-Nyungan</p> <p>Notes: Alphabetic.</p>	<p><i>Book of Genesis</i>, translated by Bruce E. Waters (1981) 1,174 words. © Bruce E. Waters. Free to use and share for non-commercial purposes.</p>
N-NCI	Classical Nahuatl , a.k.a. Aztec	<p>Location: Central Mexico</p> <p>Time: c. 7th to 16th century</p> <p>Family: Uto-Aztecan</p> <p>Notes: Language identified in Tucker & Talbert's decipherment attempt. Alphabetic.</p>	<p><i>Gospel of Luke</i>, translated by Mariano Paz y Sanchez (1833) 28,389 words. Public domain.</p>
N-OCI	Old Occitan	<p>Location: France, Spain, Italy, Monaco</p> <p>Time: 8th to 14th century</p> <p>Family: Indo-European, Romance</p> <p>Notes: Postulated to be the language of the f17r marginalia or the month names in the zodiac section. Alphabetic.</p>	<p><i>Libre de Memorias</i> by Jacme Mascaron (1390) 18,819 words. Public domain.</p> <p>Notes: Represents an example of a text with some missing and out-of-order pages.</p>

N-ORV	Old East Slavic , a.k.a. Old Russian, Old Ukrainian or Old Belarusian	Location: Eastern Europe Time: c. 10 th to 17 th century Family: Indo-European, Slavic Notes: Language identified in Stojko's decipherment attempt (as "Old Ukrainian"). Alphabetic.	<i>A Journey Beyond the Three Seas (Хождение за три моря)</i> by Afanasiy Nikitin (<i>Афанасий Никитин</i>) (1472) 7,269 words. Public domain. Notes: Represents an example of a travel journal, which some posit the Voynich Manuscript may be. Text contains details of Indian life and may contain a lot of foreign terms.
N-OSP	Early Modern Spanish , a.k.a. Classical Spanish, a.k.a. Castilian	Location: Castile Time: 15 th to 17 th century Family: Indo-European, Romance Notes: Alphabetic.	<i>A Short Account of the Destruction of the Indies</i> by Bartolomé de las Casas (1552) 30,136 words. Public domain.
N-PER	New Persian , a.k.a. Farsi	Location: Iran Time: 8 th century to present Family: Indo-European, Indo-Iranian Notes: Abjad	Extract from <i>The Book of Kings (Sháh Námeḥ/شاهنامه)</i> by Ferdowsi (<i>یس و درف</i>) (c. 1000) ² 75,653 words. Public domain. Notes: The language of the text is at the transition between Middle and New Persian. Represents an example of a poetic, column-based text.
N-PML	Sabir , a.k.a. Mediterranean Lingua Franca	Location: Mediterranean Time: 10 th to 19 th century Family: Creole Notes: Represents an example of a medieval European creole. Based on Italian. Alphabetic.	Extracts from <i>The Bourgeois Gentleman</i> by Molière (1670) 240 words. Public domain.
N-QUC	Classical K'iche'	Location: Guatemala Time: Unknown to c. 17 th Family: Mayan Notes: Logographic or alphabetic.	<i>Popol Vuh</i> by multiple (no date), compiled and published by Francisco Ximénez (18 th century) 28,525 words. Public domain. Notes: Represents an example of a two-column layout.
N-QUE	Quechua a.k.a. Runa Simi	Location: Andes Time: c. 14 th /15 th century to present Family: Quechuan Notes: Alphabetic.	<i>Inkarri</i> by unknown (c. 17 th or 18 th century) 3,451 words. Public domain.

² Sorry, my version of Microsoft Office 2013 can't combine Arabic alphabet text correctly.

N-RUM	Moldovan , a.k.a. Moldavian or Romanian	Location: Moldova Time: c. 20 th century to present (depending on definition of language) Family: Indo-European, Romance Notes: Has similar letter frequency distribution according to Jaśkiewicz. The labelling of a “Moldovan language” is derived from the study mentioned above, and does not reflect the political opinions of anybody involved. Alphabetic.	<i>Constitution of Transnistria (Конституция Републичий Молдовенешть Нистрене)</i> as of 2000 by multiple (2000) 4,052 words. Public domain.
N-SAN	Classical Sanskrit	Location: India Time: 2 nd millennium BCE to unknown Family: Indo-European, Indo-Iranian Notes: No native writing system.	<i>Eight Chapters (Aṣṭādhyāyī/अष्टाध्यायी)</i> by Pāṇini (पाणिनि) (4 th century BCE) 18,588 words. Public domain. Notes: Sample uses Devanagari script. Represents an example of a linguistic text that discusses the rules of the language it is written in.
N-SUX	Sumerian	Location: Southern Mesopotamia Time: 4 th to 3 rd millennium BCE Family: Isolate Notes: Represents an example of an academically constructed notation system. Originally logographic and syllabic script.	<i>Enki and the World Order</i> by unknown (4 th millennium) 2,767 words. Public domain. Notes: Represents an example of a copy of a badly degraded text.
N-SWA	Swahili	Location: Central and East Africa Time: 1 st millennium to present Family: Niger-Congo Notes: Alphabetic.	<i>Address to Ghanaian Parliament, 2009</i> , by Barack Obama (2009) 3,266 words. Public domain. Notes: Represents an example of a text transcribed from an oral speech.
N-THA	Thai	Location: Thailand Time: Unknown to present Family: Tai-Kadai Notes: Has similar letter frequency distribution according to Jaśkiewicz. Thai script is an abugida with 44 consonants and 15 vowels and 4 diacritics that combine. Tonal and analytic language. Abugida.	<i>Defeat of the Yuan (ลิลิตยวนพ่าย)</i> by unknown (c. 1475) 3,345 words. Public domain.

N-TPI	Tok Pisin	<p>Location: New Guinea</p> <p>Time: c. 19th century to present</p> <p>Family: Creole</p> <p>Notes: Alphabetic</p>	<p><i>Geneva Convention</i> by multiple, translated by Otto Dempwolff (1914) and by <i>Proclamation of the British Take-Over of New Guinea</i> by unknown (1914)</p> <p>552 words total. Public domain.</p>
N-TUR	Turkish (modern)	<p>Location: Turkey and Northern Cyprus</p> <p>Time: 19th century to present</p> <p>Family: Turkic</p> <p>Notes: Features heavy agglutination and vowel harmony. Alphabetic.</p>	<p><i>Treaty of Lausanne</i> by multiple (1923)</p> <p>9,699 words. Public domain.</p> <p>Note: Sample uses the modern Latin script.</p>
N-TXB	Tocharian B	<p>Location: Tarim Basin (present day Xinjiang in China)</p> <p>Time: 6th to 9th century</p> <p>Family: Indo-European, Tocharian</p> <p>Notes: Alphabetic.</p>	<p>Manuscript series <i>A, IOL Toch, M 500.1, PK AS 1-7</i> (only items exclusively in Tocharian B) by unknown (6th-9th century)</p> <p>7,903 words. Public domain.</p> <p>Note: Represents an example of a text copied from many disparate sources. Sample uses academic Latin transcription.</p>
N-VIE	Vietnamese (modern)	<p>Location: Vietnam</p> <p>Time: 19th century to present</p> <p>Family: Austro-Asiatic</p> <p>Notes: Tonal language. Alphabetic.</p>	<p><i>The Tale of Kieu</i> by Nguyễn Du (1820)</p> <p>22,848 words. Public domain.</p> <p>Note: Sample uses the modern Latin script.</p>

4. Legal and authorship information

The copyright status of each sample is stated in the right hand column of the table in section 3. In nearly all cases I have found texts that are in [public domain](#), therefore there are no restrictions on those. You are free to download, share, modify, use and abuse these files as you wish, with no permission needed.

A few cases have a different copyright status:

- N-ASL (American Sign Language) and the ASLSJ notation is copyright Thomas Stone. He has given permission to use and share this sample as long as we credit him.
- N-MWP (Kalaw Lagaw Ya) is copyright [Bruce E. Waters](#). This sample is free to use and share for non-commercial purposes.

This sample pack was selected and compiled by [Brian Cham](#) in November 2014 with help from [Thomas Stone](#) and Luman Wang. Please direct all comments, questions, suggestions and contributions to briancham1994@gmail.com.

5. Works cited

List uses MLA referencing convention.

- Jaśkiewicz, Grzegorz. "Analysis of Letter Frequency Distribution in the Voynich Manuscript." *Proceedings of the international workshop CS&P, 28-30 September 2011, Pułtusk, Poland*. Ed. M. Szczuka et al. Warsaw, Poland: Warsaw University of Technology. Print.
- Reddy, Sravana, and Knight, Kevin. "What We Know About the Voynich Manuscript." *Proceedings of the 5th ACL-HLT Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities, 2011, Portland, Oregon, USA*. Ed. Zervanou, Kalliopi, and Lendvai, Piroska. Portland, Oregon, USA: Association for Computational Linguistics, 2011. Print.
- Stojko, John. *Letters to God's Eye: The Voynich Manuscript for the first time deciphered and translated into English*. New York: Vantage Press, 1978.
- Stolfi, Jorge. "Chinese Theory Redux: Comparing the VMS and East Asian word length distributions." *Voynich Manuscript Stuff*. Stolfi, Jorge, 19 January 2002. Web. 26 February 2014. <<http://www.ic.unicamp.br/~stolfi/voynich/02-01-18-chinese-redux/>>
- Tucker, Arthur O., and Talbert, Rexford H. "A Preliminary Analysis of the Botany, Zoology, and Mineralogy of the Voynich Manuscript." *Herbalgram.org* 100 (2013): 70-85. *Herbalgram.org*. Web. 6 February 2014.