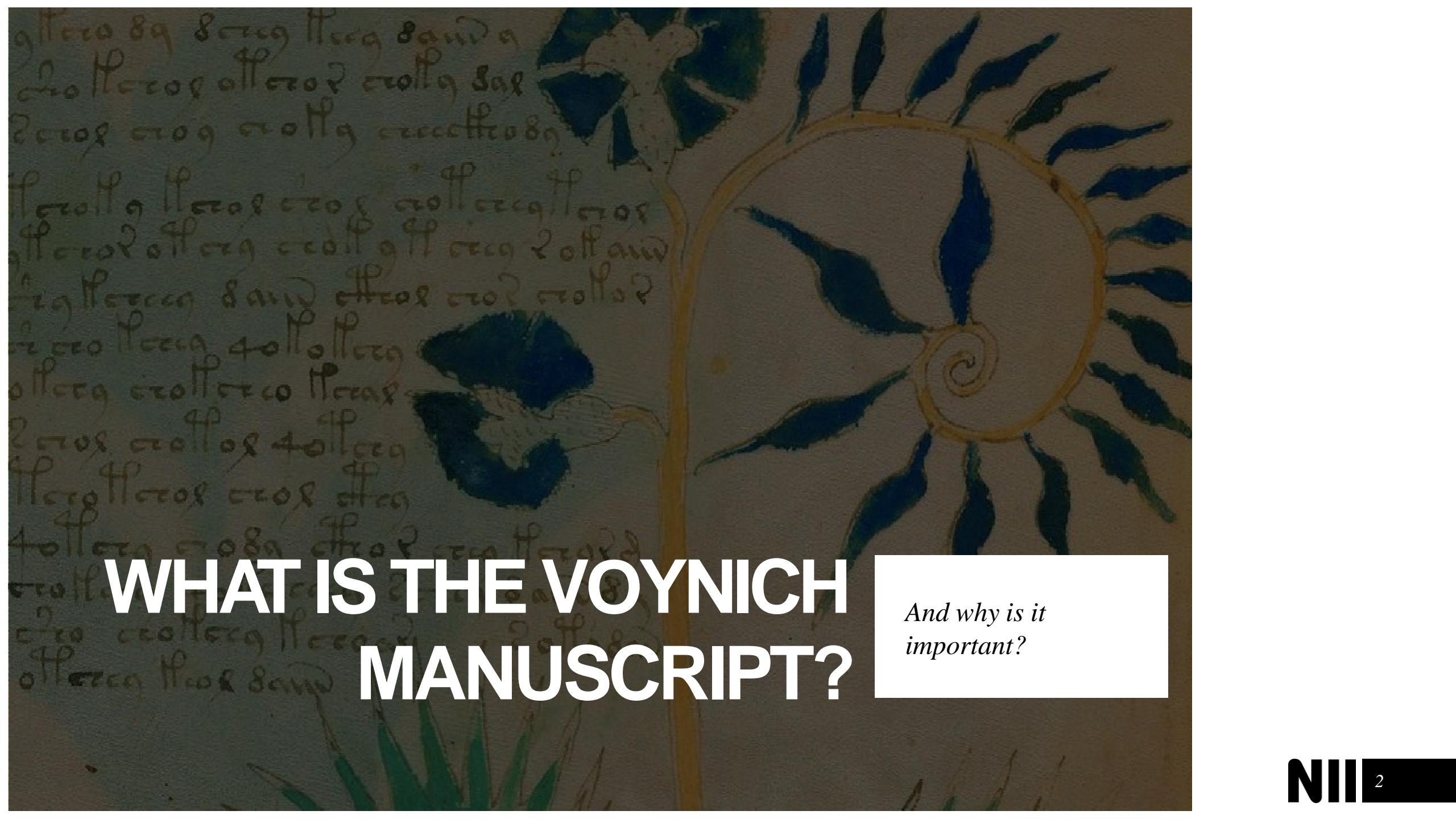


ON IDENTIFYING NATURAL LANGUAGES

*A case study on the
Voynich Manuscript*

João Figueira





WHAT IS THE VOYNICH MANUSCRIPT?

*And why is it
important?*

WHAT IS THE VOYNICH MANUSCRIPT?

A 15th Century Illustrated Book

- Acquired by Polish book seller Wilfrid Voynich (Michał Habdank-Wojnicz) in Italy, in 1912
- 200 pages of vellum (goat skin)
- Text on every page and big illustrations on most of them
- About late medieval herbology, astrology, pharmacology, and health
- Currently at the Yale library of rare books



WHY IS IT INTERESTING?



A Real-Life Indiana Jones Mystery

- Unknown Origin
- Unknown Author
- Unknown Writing System
- Unknown Language or Language Family

WHY IS IT INTERESTING?

The manuscript has been systematically studied by professional and amateur linguists and cryptographers alike

Many have made decipherment claims but the book has never demonstrably decoded

It is still entirely unclear if the book is a Hoax, is written in a plain Natural language, or is Enciphered!



“DECODED”



Scientists claim to crack an elusive centuries-old code – and it's Hebrew

Hebrew

Voynich Manuscript SOLVED:
World's most mysterious book deciphered after 600 years

Latin

Has the Voynich Manuscript Finally Been Decoded?:
Researchers Claim That the Mysterious Text Was Written in
Phonetic Old Turkish

in [Books, History](#) | February 21st, 2019 [19 Comments](#)

21 FEBRUARY, 2014 - 21:30 APRILHOLLOWAY

First words in mysterious Voynich Manuscript decoded

Turkish

Linguists often use “circular” and self-reaffirming thought processes to try to translate words in the manuscript

Arabic

HOAX, CODE, OR NATURAL LANGUAGE?

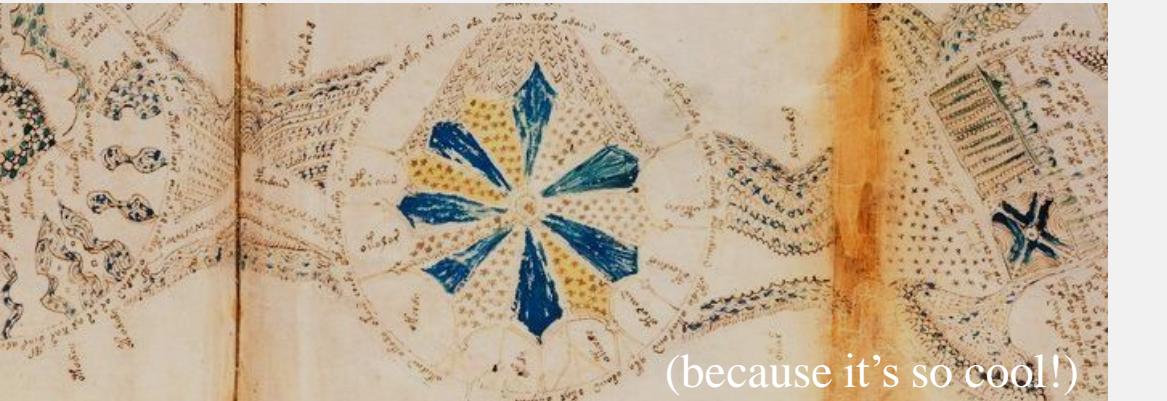
- What tools to use?
 - How can we identify a natural language?

HOAX, CODE, OR NATURAL LANGUAGE?

- What tools to use?
- How can we identify a natural language?

New options might be open now with advances in NLP

My main point is to share my big interest for this book with other NLP researchers so that more ideas might surge for exploring the properties of this book



(because it's so cool!)

offered offload code off
and as code gotcode off
as offload gotcode code
code recode gotcode
code query code
2 offload offload off
as offload code
as gotcode Hello So
as code code gotcode
Hello as Hello code
as offload code
as code offcode gotcode
Hello offload gotcode
as offload offcode off
as saw as offload re
offload gotcode code
as code gotcode code
as code code
as code as offload off
as code gotcode offload
as code code as off
and code Hello as
as offload offcode off
as code code as off

WHAT DO WE ALREADY KNOW?

*About the Voynich
Manuscript*

ORIGINS OF THE MANUSCRIPT

- The vellum was carbon dated to 1404-1438
- Made with techniques that were common all over Europe
- Uses pigments that were cheap and common in Europe
- Earliest letter about its existence suggests it was sold to Rudolph II, emperor of the Holy Roman Empire, in the late 16th century
- Nothing is known about the almost 200 years between its creation and sale
- Sold to him for 600 ducats (pure gold coins) (about 70 thousand euro worth of gold today)



CONTENTS OF THE MANUSCRIPT

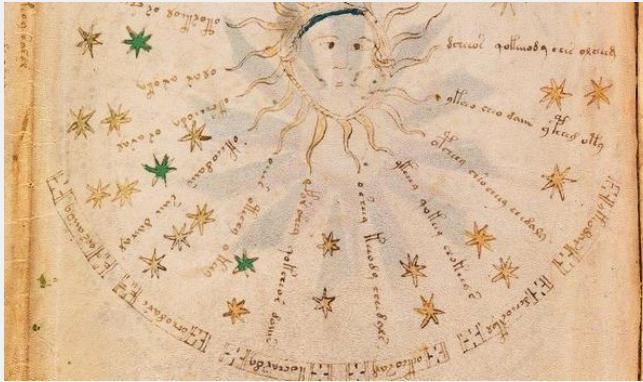
Herbal

(drawings of plants)



Astrological

(zodiac star maps)

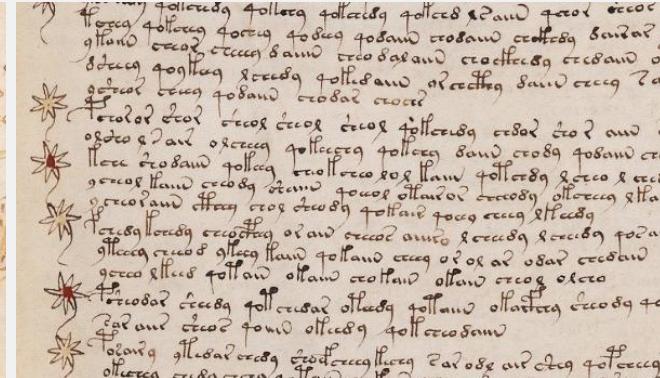


Biological

(bathing women)



Based on illustrations, the book can be divided into 6 different sections



Cosmological
(circular illustrations)

Pharmaceutical
(parts of plants)

Recipes
(itemized text)

THE TEXT

The language is often called “Voynichese” and is visually similar to many European scripts

Transcriptions of the text use different numbers of letters
(depends on how some “connected” symbols are
interpreted)

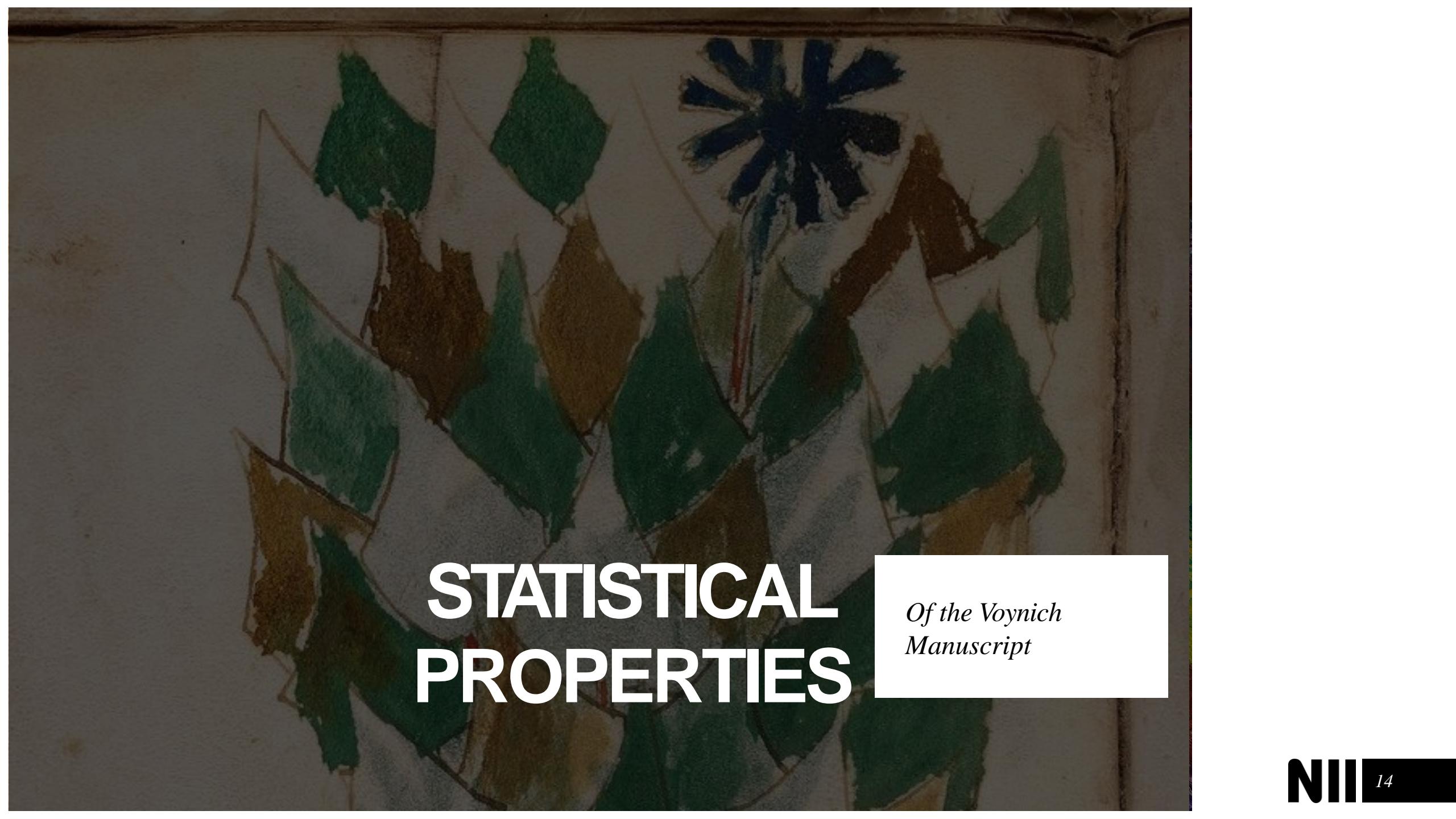
About 35 common and frequent letters with as many as 200 rare letters



OTHER QUICK OBSERVATIONS



- There is not a single instance of a word being erased or corrected (which suggests the book might be a copy) but some words have been retouched
- Serious cryptography systems (that were not simple letter reordering or substitution ciphers) did not reach Europe until years after the manuscript was written
- Some pages appear to have different handwritings. As many as 6 scribes are said to have been involved in writing the book



STATISTICAL PROPERTIES

*Of the Voynich
Manuscript*

WORD TOKENS

2 and saw oxseed one
4 Hauw offaw creting &g of Ham saw &g saw creting
offaw offaw creting offeed qotting Hauw saw offaw
offaw saw offaw creting offeed &g 2 and offaw saw &g
follow offaw offaw oxseed &g 2 and offaw saw &g
or or of creting saw qotting saw creting &g for
qotting of creting 4 offaw saw offaw creting offeed
and seeing Hauw offaw saw olof creting saw &g

ample text with
ost common word
the book labelled

170 000 characters (approximately)

37 919 word tokens

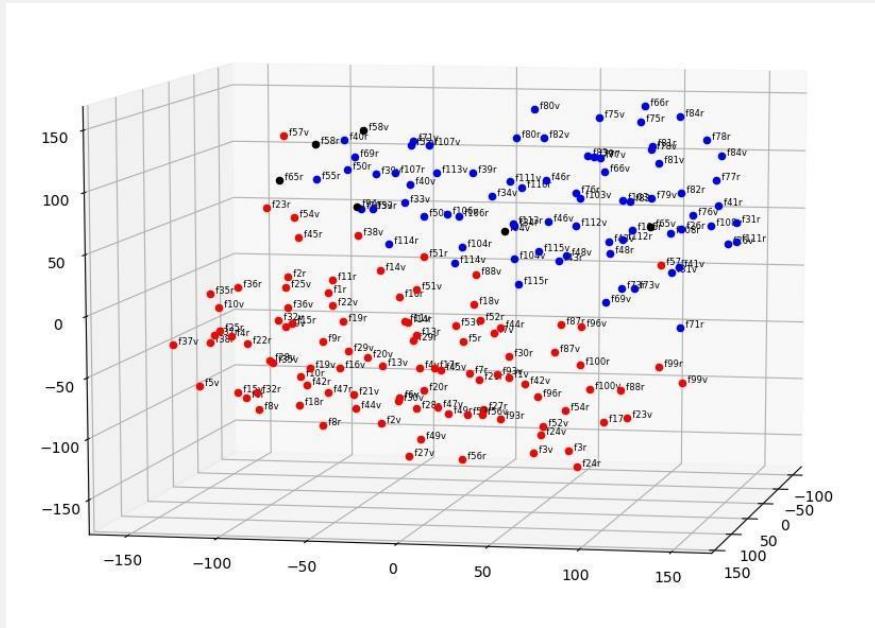
8 114 unique word tokens

POSSIBLE DIFFERENT LANGUAGES

One of the first things that where first noticed about the first and second order statistics of the text was that they where sometimes very different between two pages

As indicated by P. Currier, the pages are written in one of two languages (Currier A, or Currier B) with distinct statistical properties and rules for word formation

The two languages have also been associated with different handwritings



The two languages might also be different dialects, different ciphers, language domains

For NLP purposes this means that our already small dataset of Voynichese is actually two smaller datasets!

PROBABLY DIFFERENT LANGUAGES

Words common in Language A, but rare in Language B

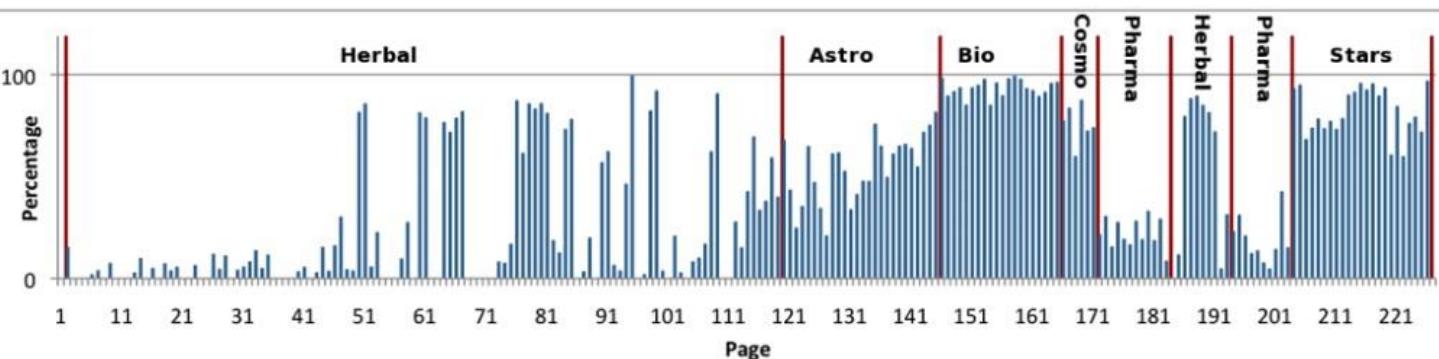
Word	# in Language A	# in Language B
ଶ୍ରୀ	61	1
ଶ୍ରୀମତୀ	20	1
ବେଳାନ୍ତି	14	1
ଶ୍ରୀମତୀ	28	2
ଶ୍ରୀମତୀ	13	1
ବେଳାନ୍ତି	13	1
ଶ୍ରୀମତୀ	11	1
ଶ୍ରୀମତୀ	21	2

Most research focuses on Currier B due to the larger number of pages

Words common in Language B, but rare in Language A

Word	# in Language B	# in Language A
ଫଲାଙ୍କ	159	1
ଫୋ	112	1
ଫଲାଙ୍କ	107	1
ଫଲାଙ୍କ	50	1
କଟିଙ୍ଗ	28	1
ଗୋ	21	1
ଫଲାଙ୍କ	81	4
ଫଲାଙ୍କ	40	2
କଟିଙ୍ଗ	18	1
କାନ୍ଦ	31	2

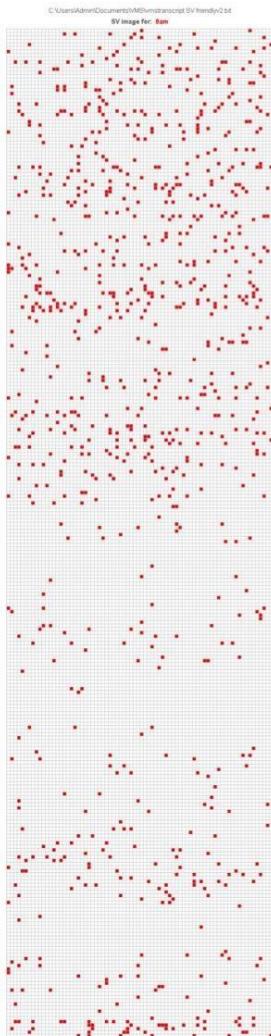
Figure 2: VMS sections, and percentage of word tokens in each page that are tagged as language B by the HMM.



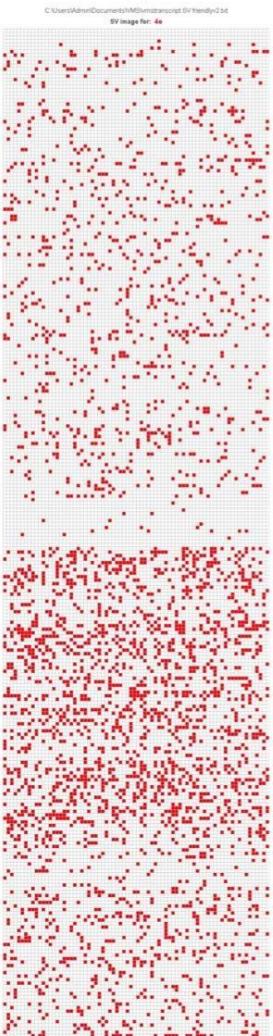
PROBABLY DIFFERENT LANGUAGES

Common syllables in the Voynich Manuscript, showing entire manuscript

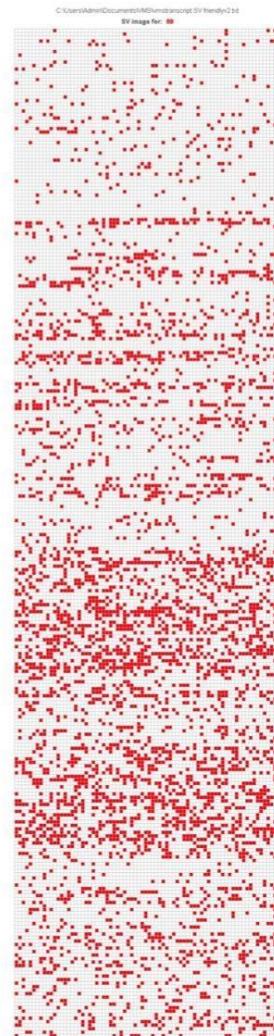
8AM



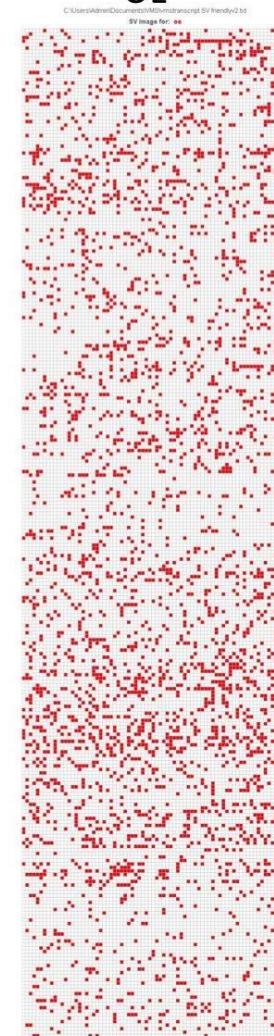
4o



89

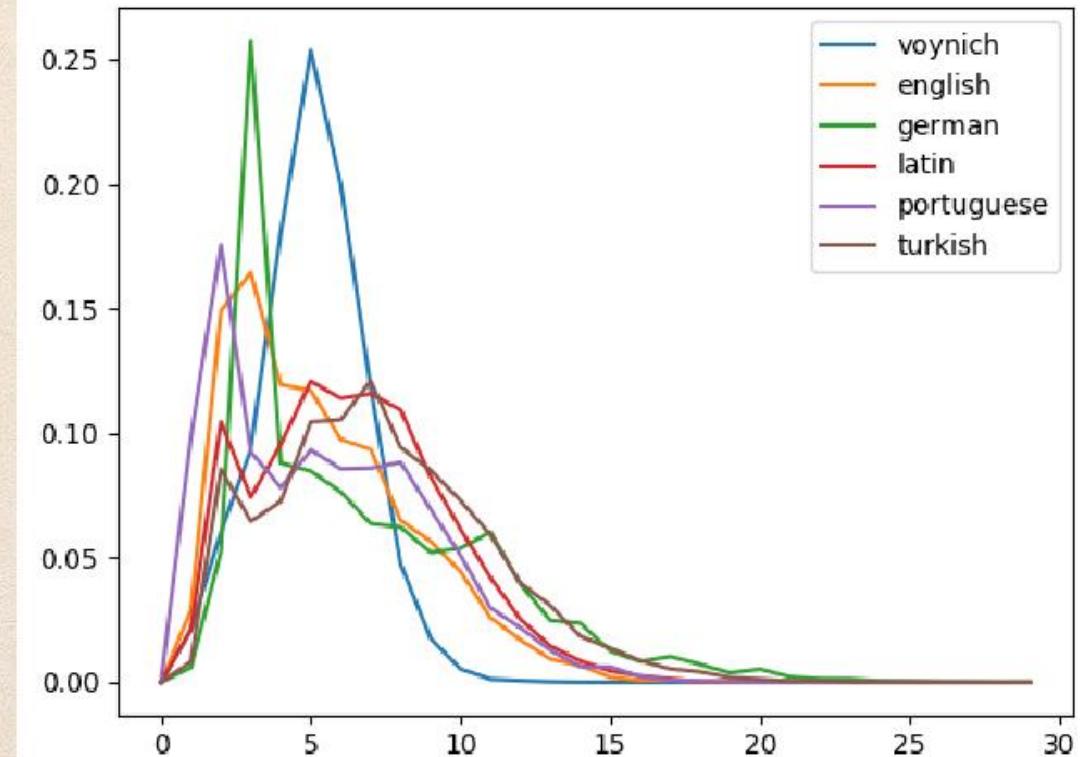


OE

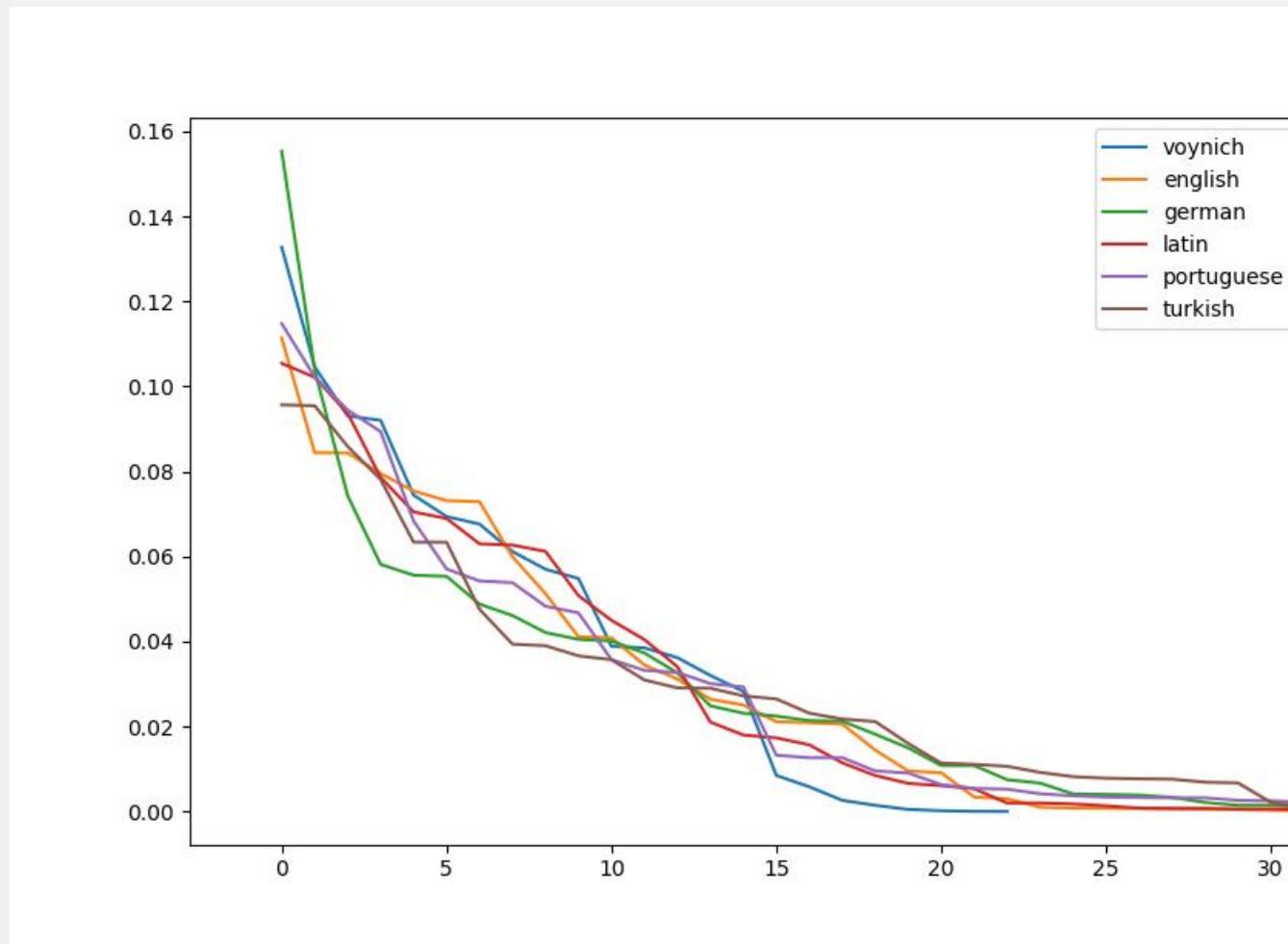


Certain syllables
are very frequent
in only one of the
languages

WORD SIZES SEEM VERY NATURAL



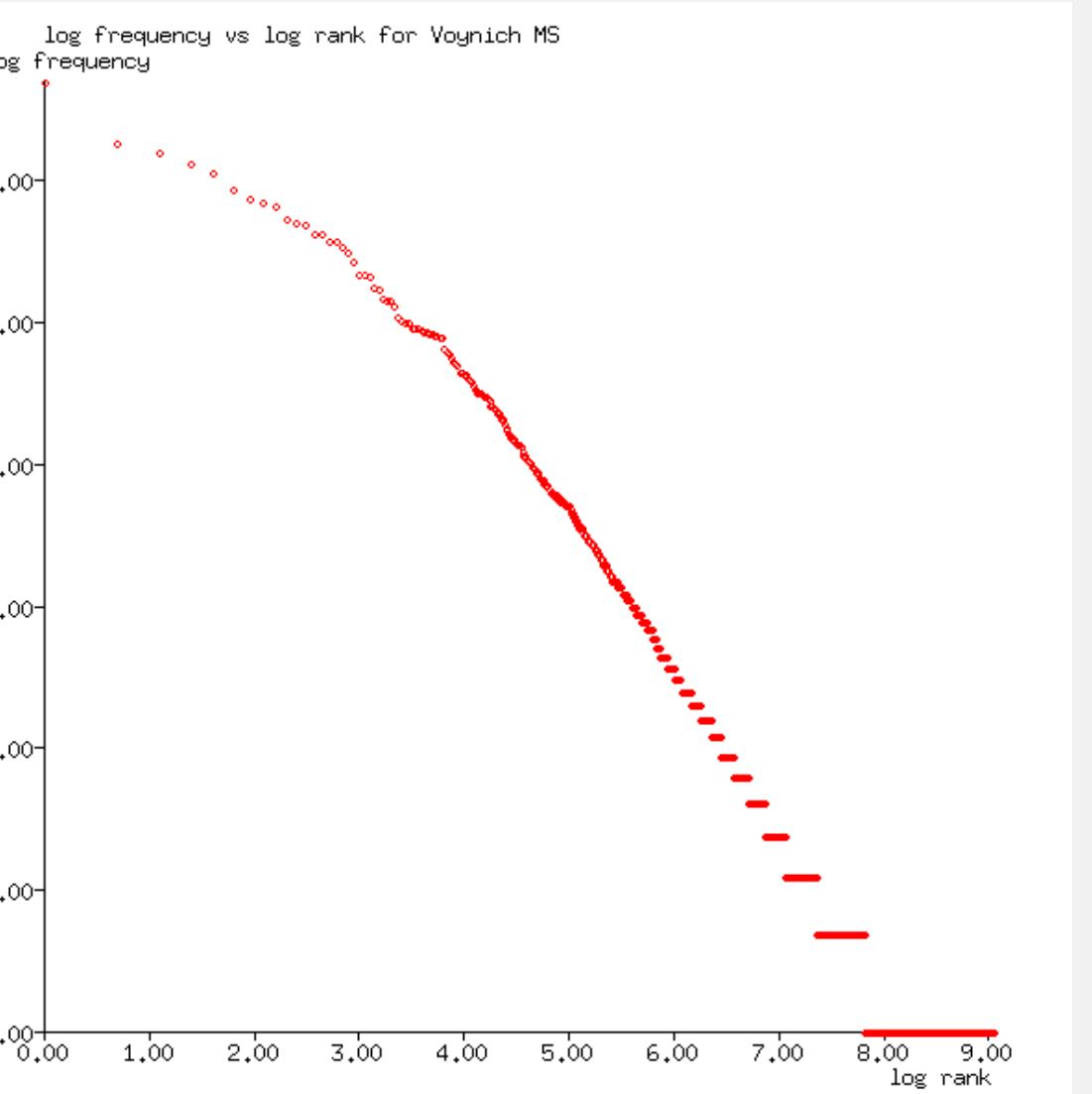
LETTER FREQUENCIES ARE ALSO NATURAL



WORD FREQUENCIES FOLLOW ZIPLF'S LAW

Zipf's law for word length is also followed by Voynichese (shorter words are more frequent)

This is expected
of a real natural
language



ARE THERE VOWELS AND CONSONANTS?

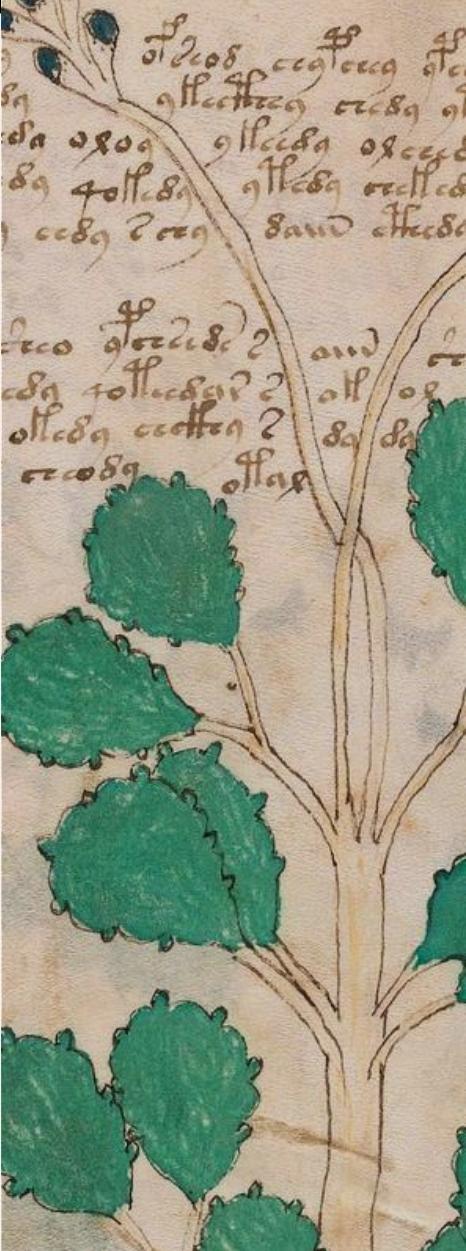
Sukhotin's algorithm for vowel identification (works perfectly on modern languages) separates 4 symbols as vowels but with a low degree of separation, and many words do not contain any of these symbols

Expectation Maximization algorithms for a two-state Hidden Markov Model yields similar results

Further research has shown that the letters can be grouped into two classes: letters that appear at the end of words, and all others

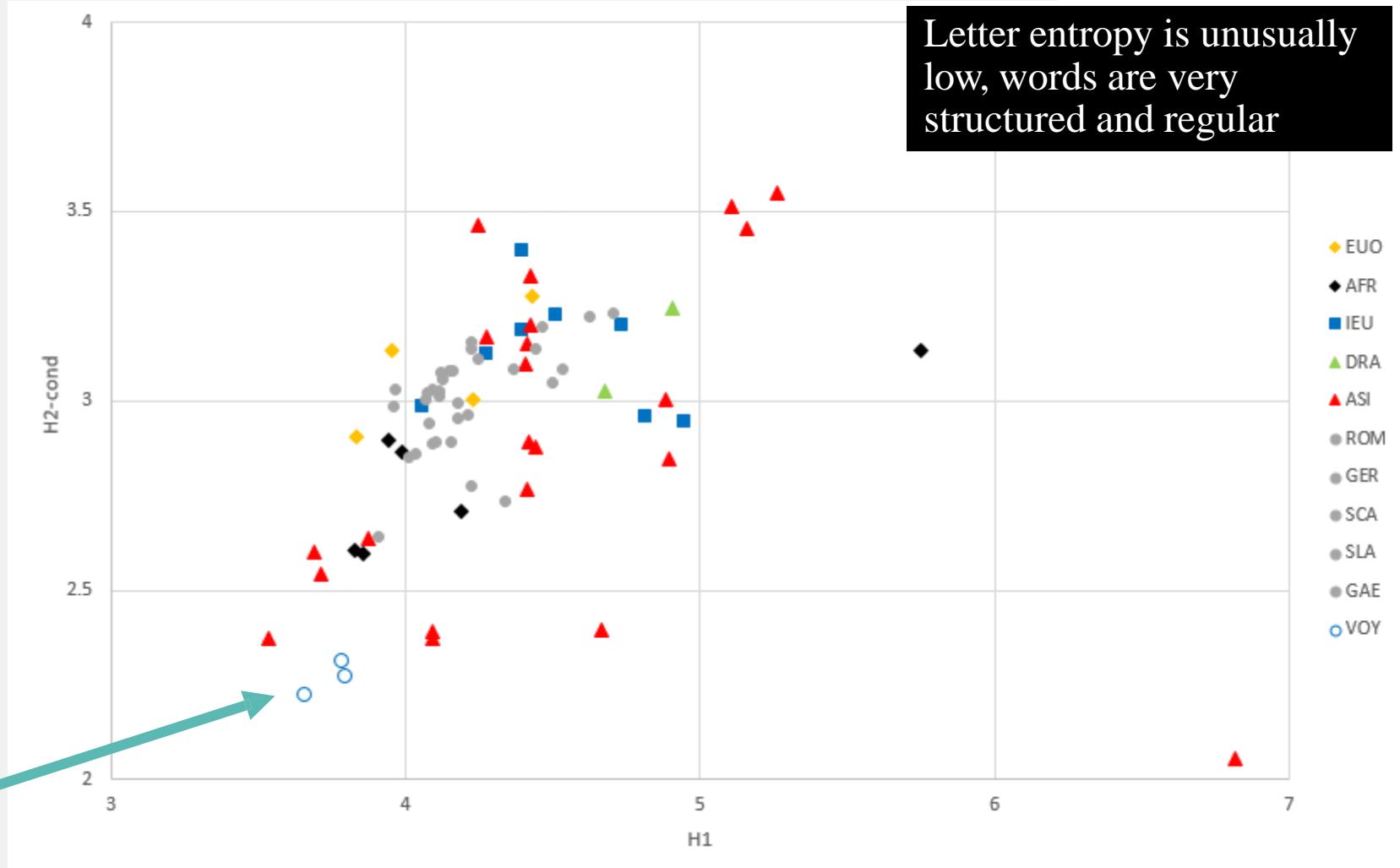
This is a weird property not shared with modern languages

RegEx grammar for word construction is basically a^*b



WORD STRUCTURE AND LETTER ENTROPY

Code	Meaning
ROM	Romance Languages
GER	Germanic Languages
SCA	Nordic Languages
SLA	Slavic Languages
GAE	Gaelic Languages
EUO	Other European
AFR	African
IEU	Other Indo-European
DRA	Dravidian
ASI	Asian
VOY	Voynichese



GRAMMATICAL STRUCTURE AND WORD ENTROPY

The text has famously low quantities of repeated word pairs and trios

The text feels very random in word order, the current word shows little to no information about the next

Table 4: Predictability of words (over 10-fold cross-validation) with bigram contexts, compared to unigrams.

	Unigram	Bigram	Improvement
VMS B	2.30%	2.50%	8.85%
English	4.72%	11.9%	151%
Arabic	3.81%	14.2%	252%
Chinese	16.5%	19.8%	19.7%
Hungarian	5.84%	13.0%	123%

Less Predictability = More Entropy

None of the tested natural languages had such a high word bigram entropy

PAGE TOPICALITY

If pages have specific topics, it is expected that there is an unusually higher frequency of some words

Table 6: Strength of page topics in VMS and other texts, cropped to be of comparable length to the VMS.

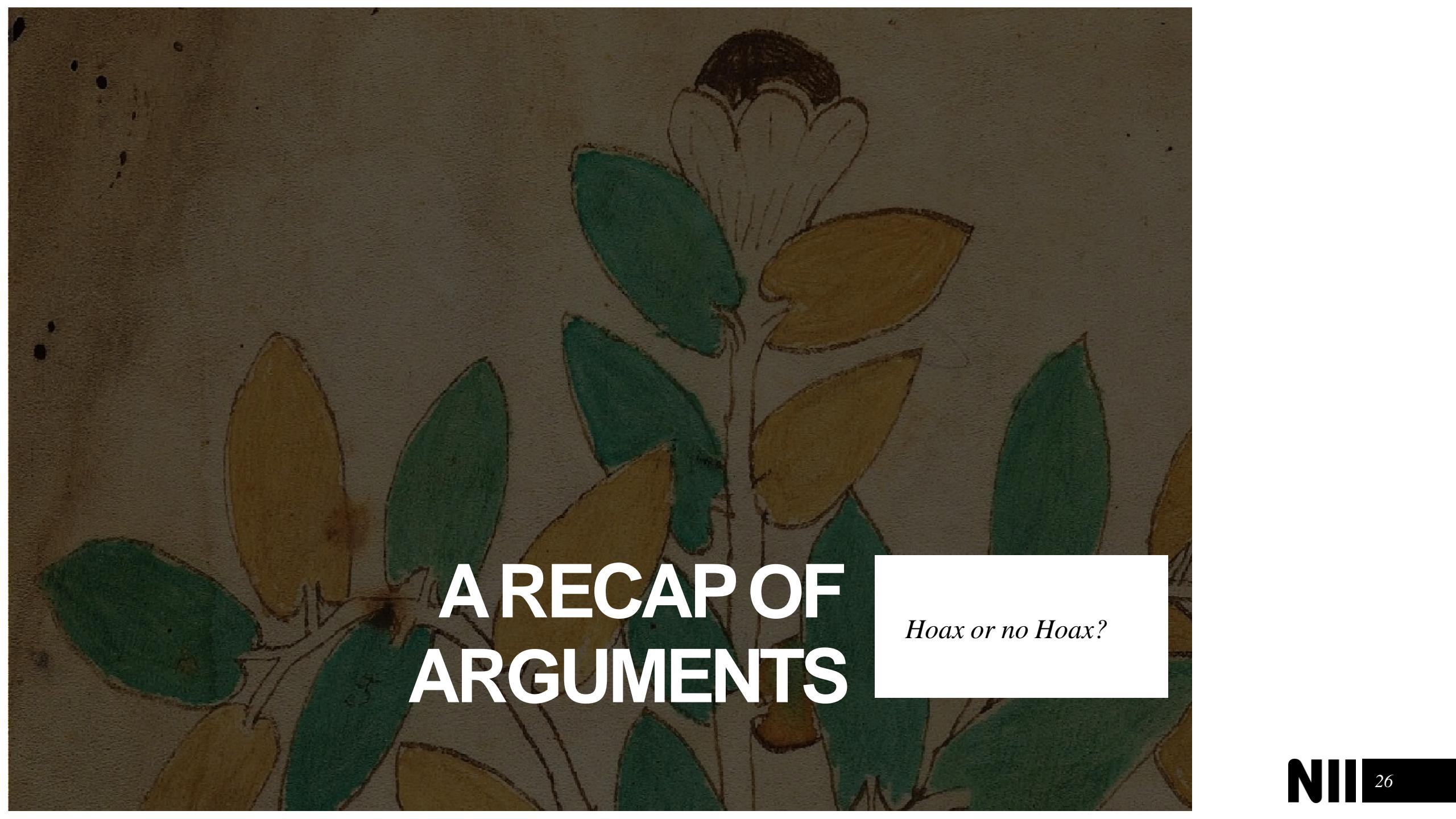
	VMS B	English WSJ	English Genesis	Arabic Quran
T	7.5	6.3	6.6	7.7
T_{rand}	7.7	6.5	7.1	7.9
$1 - T/T_{rand}$	0.033	0.037	0.069	0.025

Adjacent pages share a lot of common vocabulary, as expected of a real book with ordered pages (calculated here with bag of words cosine similarity)

For reference, the topicality of pages (calculated with word entropy) is compared to the topicality of the scrambled document

Table 7: ADJPAGESIM for VMS and other texts.

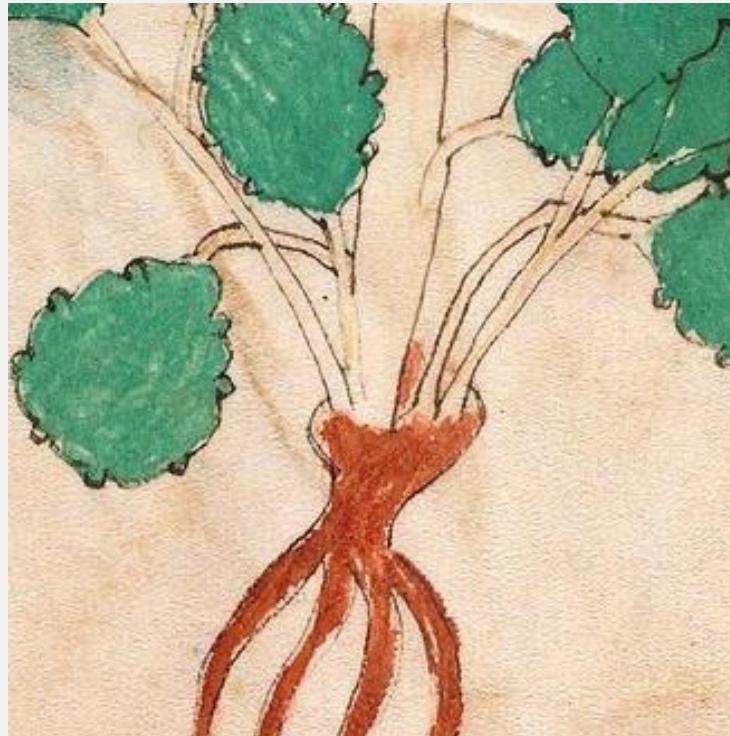
VMS B	38.8%
VMS All	15.6%
VMS B pages scrambled	0%
VMS All pages scrambled	0.444%
WSJ	1.34%
English Genesis	25.0%
Arabic Quran	27.5%

The background of the slide is a traditional East Asian ink wash painting of a magnolia flower. The flower has a dark, textured center and several large, overlapping petals in shades of white, cream, and light green. It is surrounded by several long, narrow, lanceolate leaves in various stages of growth, some green and some yellowish-green. The overall style is minimalist and expressive.

ARECAP OF ARGUMENTS

Hoax or no Hoax?

WHY IT'S PROBABLY A HOAX



- No one has ever been able to decipher it
- Almost completely random word order with no grammatical structure
- The lack of any instance of a word being corrected
- Sold in ancient times for a lot of money which would justify creating it

WHY IT'S PROBABLY NOT A HOAX

- Follows Zipf's law for word frequencies
- Natural word sizes and letter frequencies
- The two identified languages/dialects
- Pages have topics, nearby pages have similar topics
- A book of this size would be very hard to fill with fake text





LEADING HYPOTHESES

For origin

IF NOT A HOAX

Cypher: Could be an existing European language intentionally masked behind some cypher

Shorthand: Could be stenography for efficient writing of an existing language. Individual symbols could replace long words

Lost Natural Language: Could be the last written evidence of some minor language that has gone extinct

Phonetic Writing: Could be a known or lost language written mostly phonetically due to lack of a consistent writing backlog. Different people might write the same words as different sounds (which would explain the two apparent languages)



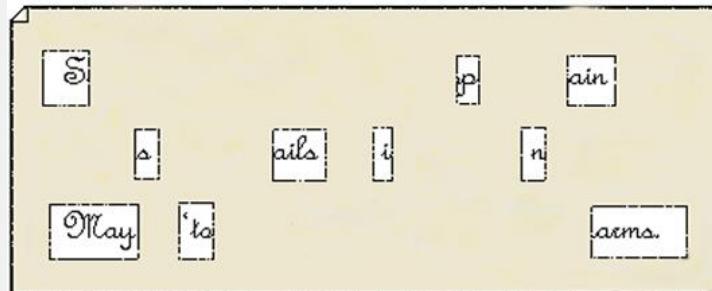
IF A HOAX



Why: The book was sold for a lot of money

How: Either by trying to write random letters, or by using dice or some other system for quickly creating random words, or by trying to shuffle a real text in a real language

Sir John regards you well and spekes again that
all as rightly 'waile him is yours now and ever.
May he 'lone for past 2'lays with many charms.

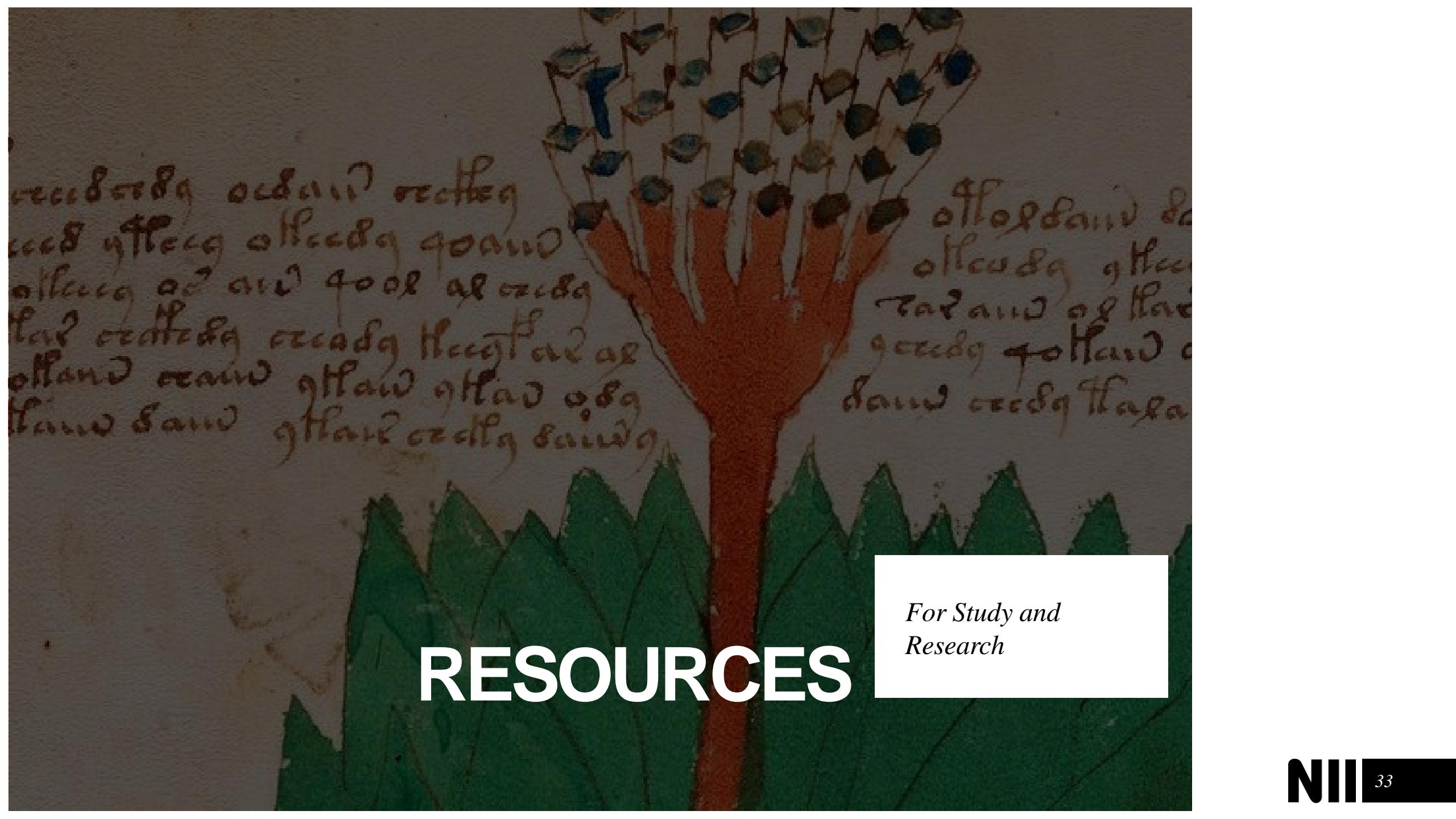


A Cardan grille, used for hiding messages, could have had been used to create the text



WHAT TO DO?

*As Computational
Linguists*



RESOURCES

*For Study and
Research*

See the book (includes mapped transliteration and word finder):

<http://www.voynichese.com>

Text descriptions of the images and comparison of all known transliterations:

<http://voynich.freie-literatur.de/index.php>

Most popular site for Voynich resources, includes download links for transcriptions:

<http://www.voynich.nu>

Article summarizing some of the known properties of the book:

<https://www.isi.edu/natural-language/people/voynich-11.pdf>

Github repository I made with one of the transcriptions and a python parser for it:

<https://github.com/joaoperfig/voynichstudies>





THANK YOU!