CRICOS PROVIDER 00123M

**Students:** Peter Roush and Bryce Shi

**Supervisors:** Derek Abbott, Maryam Ebrahimpour and Brian Ng

# Project 44: Cracking the Voynich Code

Final Seminar

# Outline

- ## The Voynich Manuscript
  - Objectives, Background and Motivation
- ## Analysing the Manuscript
  - Techniques, Research, and Testing
- ## The Information Learnt
  - Results and Analysis
- ## Project Management
  - Team Roles, Milestones, and Budgeting
- ## Conclusions

Background, Motivation, and Objectives

# The Voynich Manuscript

# A Brief History

- Voynich Manuscript
  - Found in an Italian Castle by Wilfred Voynich, a book collector
  - Pages and some references dated to the 15th Century
  - Author or authors unknown
  - Language unknown
  - Pictures have been inconclusively matched to plants in Europe and South America

- Electronic Transcriptions
  - At least two different languages or dialects
  - Hard to separate letters into a fixed alphabet
  - Interlinear Transcription File

# Current Theories

- Early Language or Writing System
  - Early Welsh (Tim Ackerson)
  - Romanised Manchu Chinese (Zbigniew Banasik)
- Code
  - Fake cipher related to Arabic numerals (D'Imperio)
  - Cipher by Roger Bacon (William Newbold)
  - Cipher by Antonio Averlino (Nick Pelling)
  - Certain pages are key to unlocking the mystery (Mark Sullivan)
- Hoax
  - Written to scam money out of Rudolf II (Raphael Mnishovsky)
  - Written by Voynich for money and fame

# Voynich Manuscript

- **Part 1** (Herbal)
  *129 pages*
- **Part 2** (Astronomical)
  *12 pages*
- **Part 3** (Biological)
  *20 pages*
- **Part 4** (Cosmological)
  *20 pages*
- **Part 5** (Pharmaceutical)
  *18 pages*
- **Part 6** (Recipes)
  *25 pages*

Detailed chemical analysis can be found at Yale:
http://beinecke.library.yale.edu/sites/default/files/voynich_analysis.pdf



Pictures reproduced from Beinecke Library under the free public domain licence

# Characters (EVA Alphabet)



| A | B | C | D | E | F | G | H | I | J | K | L | M |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| N | O | P | Q | R | S | T | U | V | X | Y | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|

Superis habeo gratiam
quorum maiestate sug
gerente mihi fauorum
opperfici· djksvwxyzi

Humanist miniscule writing (left)

# Objectives

- Develop Data Mining Techniques for the unknown language/code in the Voynich Manuscript.

- Compare linguistic features of the Voynich Manuscript and other languages.

- Determine whether the language in the Voynich manuscript is real, a code, or a hoax.

- Develop a code base and documentation to aid future projects.

Research, Methods and Tools
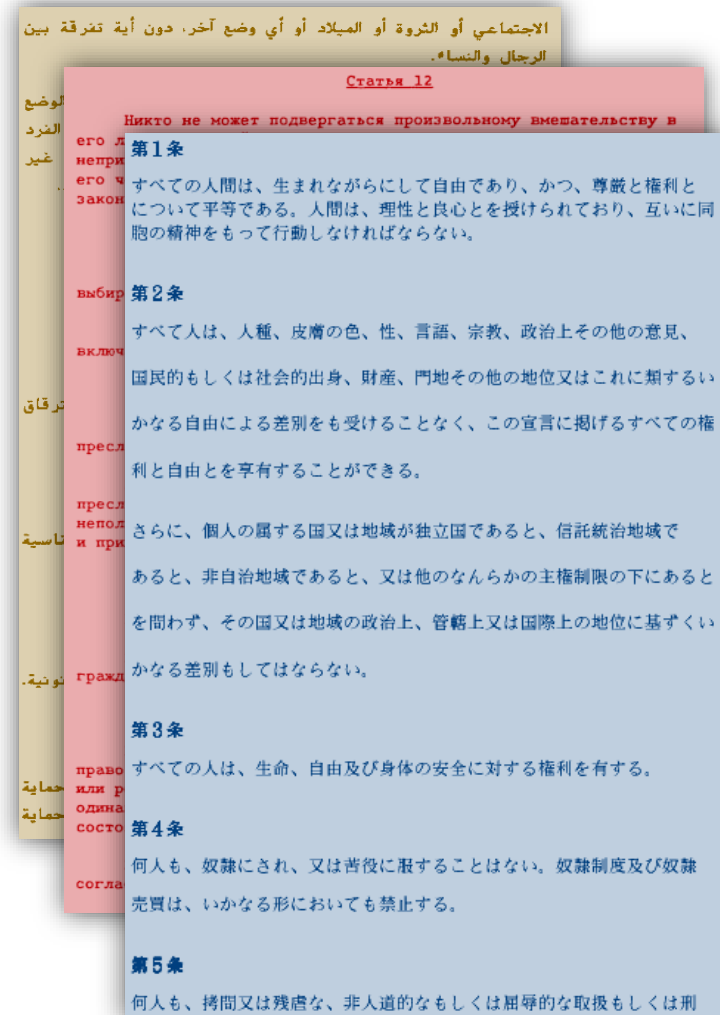
# Analysing The Manuscript

# Electronic Transcriptions

# Testing Methodology

- Used Takahashi Transcription and EVA alphabet for all tests
- Handwritten text files for basic verification
- 10 Comparison Texts of similar length in selected languages
  - English (3 Texts)
  - Latin
  - Italian
  - Hungarian
  - Hebrew (Without vowel accents)
  - Chinese (Simplified Characters)
  - Chinese (Pinyin)

# The UN Declaration of Human Rights

- 382 translated languages

- Allows greater selection of comparison languages.

- Translations contain an average of 1800 word tokens.



Picture Reproduced From: www.boes.org (Public Domain)

# Collocations

- A collocation is a word combination that occurs more often than would be expected by chance:
  - "Strong Tea"
  - "Friendly Footing"
  - "Saucer of Milk"
  - "Scotland Yard"
- Collocations indicate names and expressions in a language, and don't translate well into other languages.

$$\text{pmi}(x; y) \equiv \log \frac{p(x,y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)}.$$

# TF-IDF

- TF: Term Frequency
  - Proportional to the number of times a word is used in a document or section

- IDF: Inverse Document Frequency
  - Inversely Proportional to the number of documents or sections in which a word appears

- TF-IDF scores provide a way to find words relevant to a section, while ignoring words that are common across all sections.

# Word Recurrence Interval (WRI)

- WRI is defined as the number of words in between successive occurrences of a keyword

- Keyword being: I

  <span style="color:purple">1     2     3     4     5     6     7   8   9    10      11</span>

  I have six locks on my door all in a row. When I go out, I lock every other one. I figure no matter how long somebody stands there picking the locks, they are always locking three.

- Word Recurrence interval is: {0, 11, 2, 4}

# Support Vector Machine (SVM)

Unknown Data



Class 1

- SVM is a binary classifier
- Defines a decision point from a set of training data which is split into two distinct classes
- Assigns new testing data into one of those classes based on the decision point.
- Can be used for authorship detection

# Language Investigations (Herbal Book)

- **Language and grammar was lax at times**
  - Repeated letters skipped
  - Words abbreviated with symbols
- **Position dependent letters**
  - Two different interchangeable versions of letter 's'
- **Different authors, different substitutions**
  - Separate authors would substitute words with own symbols
- **Penmanship questionable**
  - Words sometimes written as one word sometimes split apart
- **Words continued on different lines**
  - Occasionally would have an indicator to show word had been split

Results and Analysis

# Information Learnt

| Section | Currier Language | Pages | Tokens | Words | Words per Page | Full Alphabet Length | Common Alphabet Length |
|---|---|---|---|---|---|---|---|
| Cosmological | Unknown | 20 | 3008 | 1521 | 150 | 27 | 24 |
| Biological | B | 20 | 6917 | 1549 | 346 | 21 | 18 |
| Herbal A | A | 97 | 7956 | 2492 | 82 | 32 | 21 |
| Herbal B | B | 32 | 3442 | 1349 | 108 | 23 | 20 |
| Recipes | B | 25 | 11417 | 3328 | 457 | 29 | 19 |
| Pharma | A | 18 | 2573 | 1139 | 143 | 21 | 19 |
| Zodiac | Unknown | 12 | 1331 | 808 | 111 | 20 | 19 |
| Unclassified | Unknown | 12 | 1276 | 708 | 106 | 28 | 24 |
| Missing | | 20 | 0 | 0 | 0 | 0 | 0 |
| **Full Manuscript** | | **256** | **37945** | **8105** | **161** | **47** | **21** |

# Common Letter Combinations



N-Grams Ranked by Normalied Frequency (Takahashi VMS)

| N-Gram | Percentage |
|--------|-----------|
| ii | 79.4 |
| iin | 71.0 |
| ch | 70.6 |
| in | 67.0 |
| ee | 51.4 |
| ai | 51.3 |
| aii | 48.3 |
| aiin | 48.0 |
| dy | 44.7 |
| he | 43.1 |
| sh | 35.6 |
| qo | 34.2 |
| ok | 33.6 |
| ol | 31.6 |
| ed | 30.3 |
| ar | 29.9 |
| da | 29.9 |
| che | 29.4 |
| ke | 25.3 |
| al | 24.9 |

Percentage

# Word and Illustration Relationships

# Words and Illustration Relationships

| Astrological | Biological | Cosmological | Pharma | Recipes | Herbal |
|---|---|---|---|---|---|
| osar | qol | v | daiin | qokeedy | Daiin |
| oteody | qolkeedy | ytaiin | okeol | qokaiin | chor |
| oteotey | qokedy | k | ctheol | lchedy | cthor |
| eody | qokain | {&169} | olchor | lkaiin | ctho |
| okalar | shedy | {&171} | qoor | lkain | qotchor |
| okeodaly | lchedy | x | shockhey | qokain | qotchy |

# Word Lengths and Frequency

# UDHR and Word Lengths

| Text | Tolerance | Match | UDHR Match | Peak Length |
|---|---|---|---|---|
| Voynich | 10% | 45.45% | Arabic, Standard | 2 |
| Voynich | 15% | 54.54% | Arabic, Standard | 2 |
| Voynich | 25% | 63.63% | Malay (Arabic) | 4 |
| Voynich | 40% | 72.72% | Hebrew, Malay (Arabic), Guarayu, Arabic (Standard) | 4, 4 5 2 |
| Voynich | 50% | 81.81% | Arabic (Standard), Hausa (Niger), Hausa (Nigeria) | 2 2 2 |

**Voynich:**

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|
| 4.13% | 8.52% | 9.45% | 17.01% | 23.95% | 18.84% | 11.12% | 4.49% | 1.68% | 0.52% | 0.14% |

# WRI and Rank Plot

# UDHR and WRI

| Name | Tolerance | Match | UDHR Match | Comments |
|---|---|---|---|---|
| Voynich Herbal A | 10% | 17% | Bosnian (Latin) | **f15r** - **f22v** |
| Voynich Herbal A | 10% | 12% | Jola-Fonyi | **f3r** - **f10v** |
| Voynich Biology B | 10% | 3% | Hmong (Southern Qiandong), Aceh | **f83r** - **f85r1** |
| Voynich Recipe B | 10% | 22% | Bosnian (Latin), Mapudungun | **f113r** - **f114r** |
| Herbal Book | 10% | 8% | Hmong, Southern Qiandong | **16th Century** |

- Comparison text of ~1500 words

- Average UDHR text length is ~1800 words

- Top 100 data points

# Word Frequency and Zipf's Law



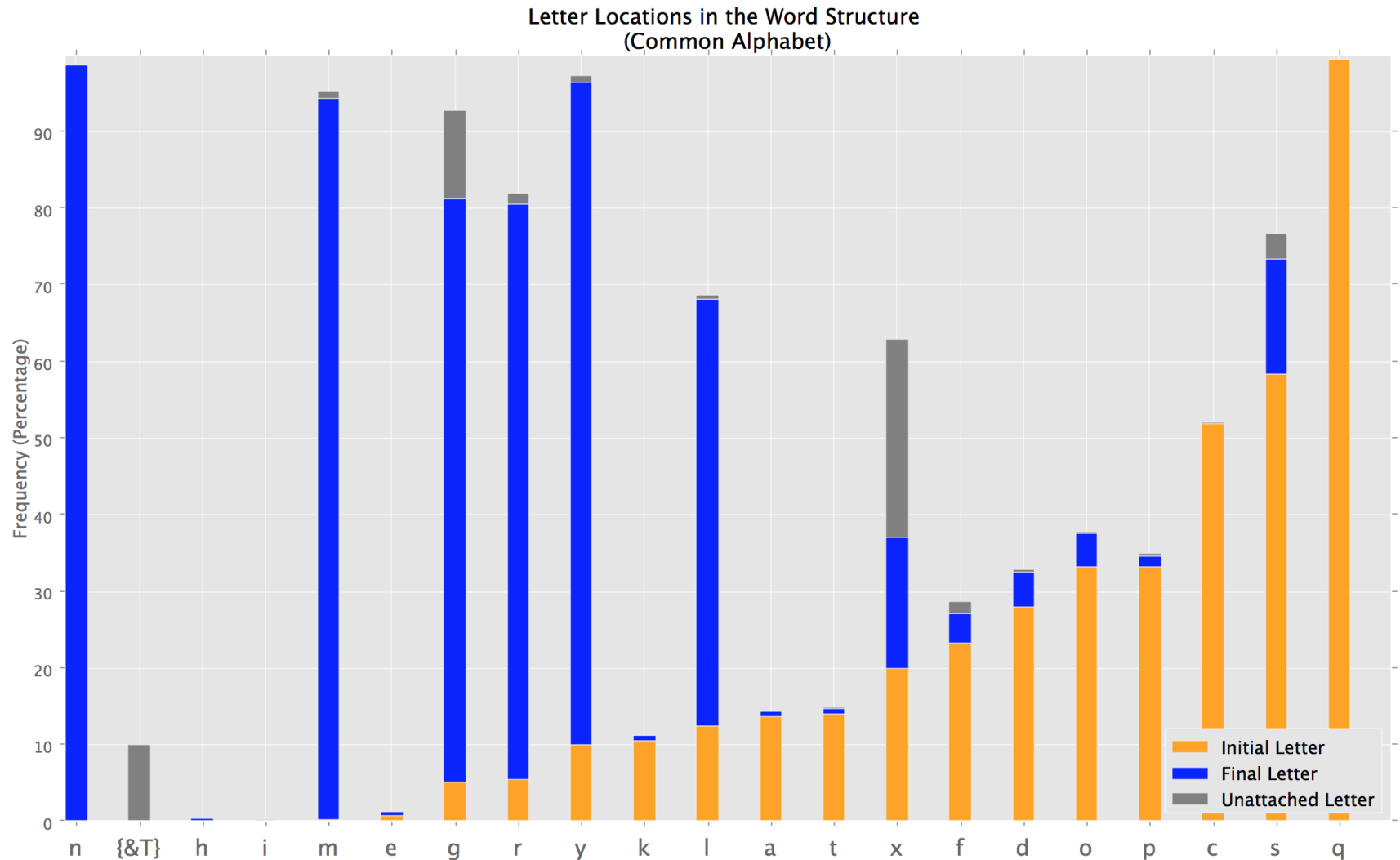Russian Frequent Words Normalised

# Word Entropy

# Collocations

# Word Structure
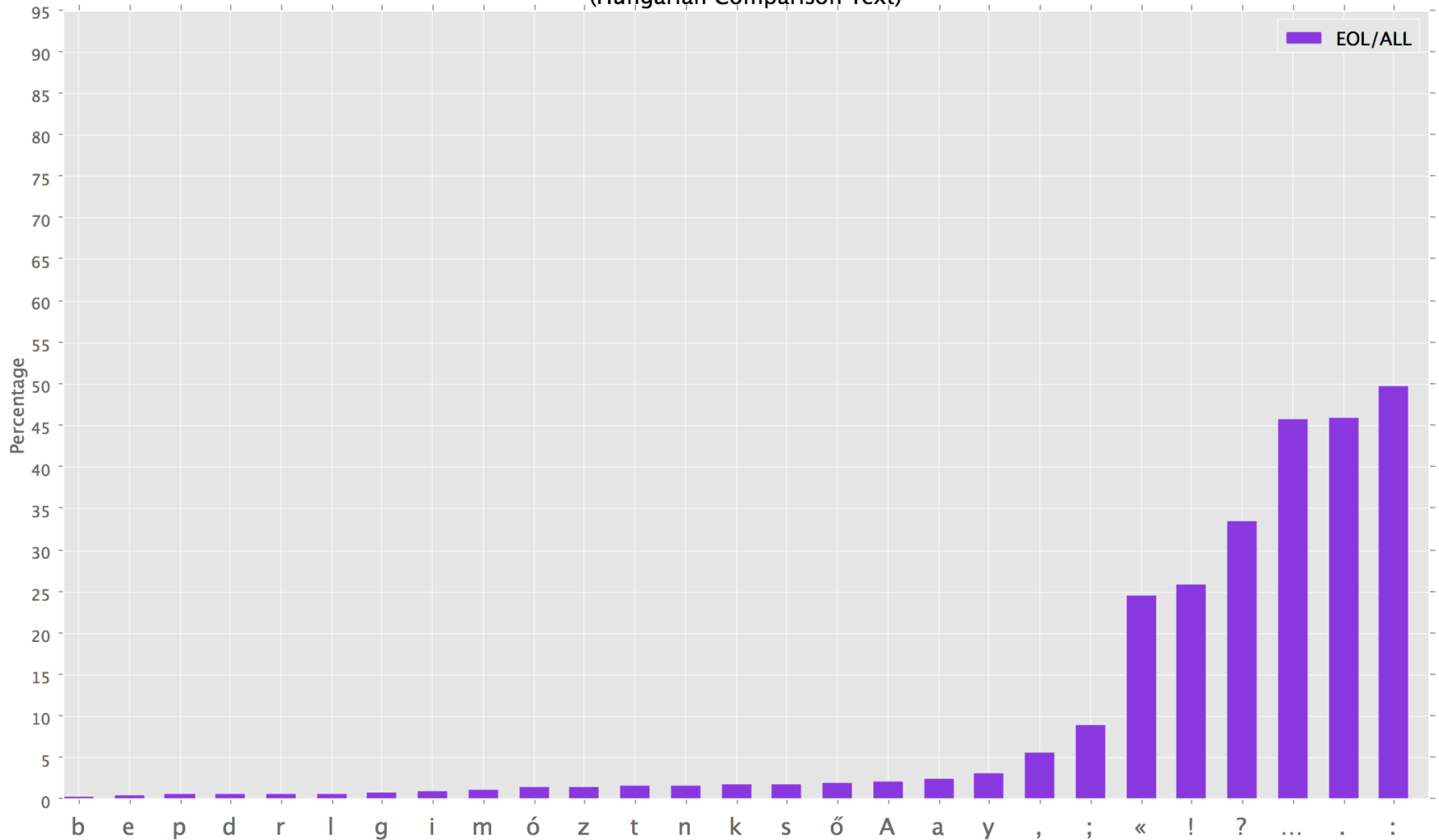


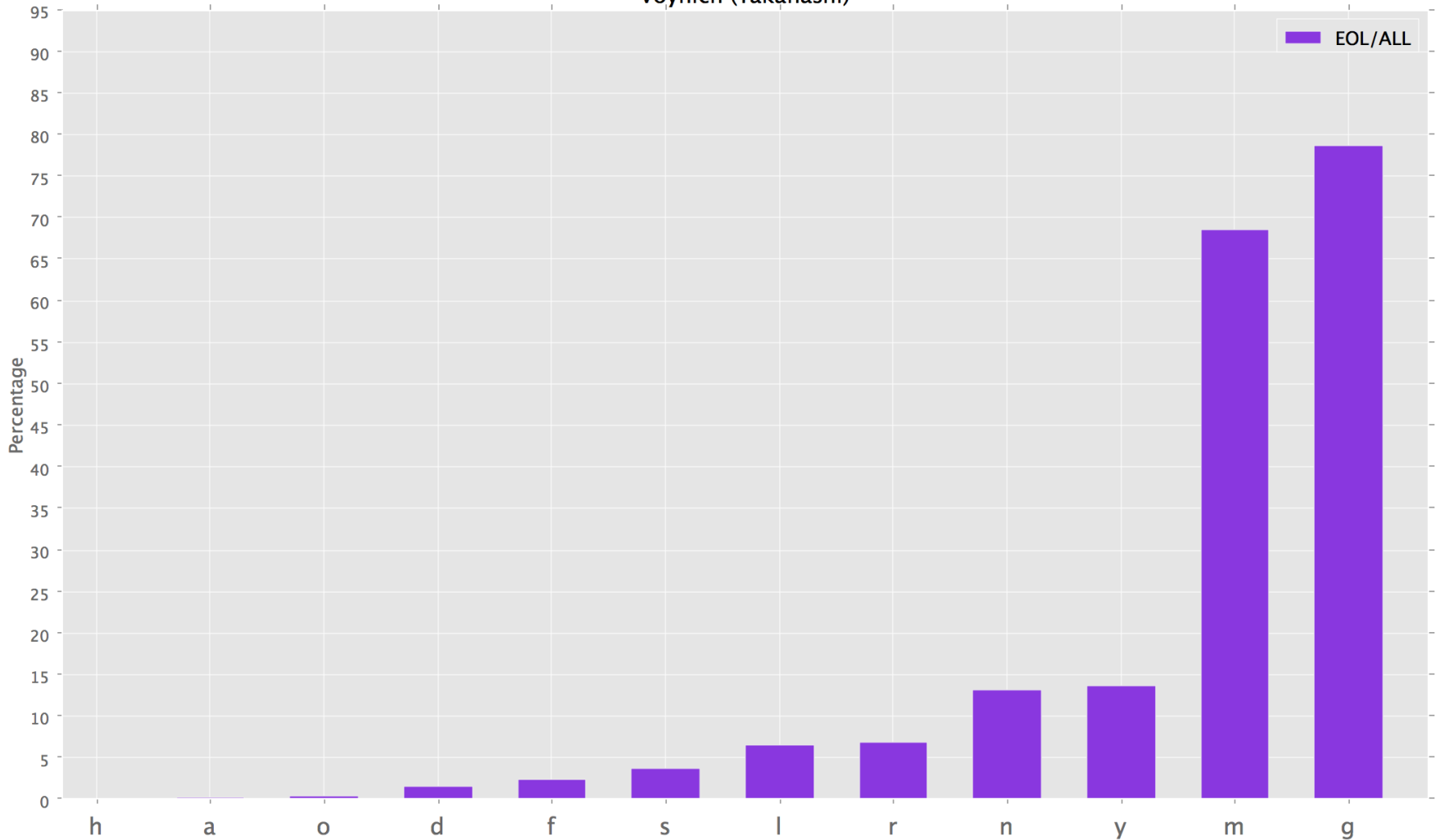Letter Locations in the Word Structure
(Common Alphabet)

# Punctuation



Significant End of Line Characters
(Hungarian Comparison Text)

# Punctuation



Significant End of Line Characters
Voynich (Takahashi)

# Support Vector Machine (SVM)

| Language | Comparisons | Group |
|---|---|---|
| Voynich $_{Takahashi}$ Normalised | Frequency | **Hebrew** |
| Voynich $_{Takahashi}$ σ | WRI | **Russian** |
| **Normal Languages Compared** | | |
| Chinese | English $_{Sherlock\ Holmes}$ | Hebrew | Hungarian |
| Italian | Latin | PinYin | Russian |

| Language | Comparisons | Group |
|---|---|---|
| Voynich $_{Takahashi}$ Herbal A | Frequency | **Zodiac** |
| Voynich $_{Takahashi}$ Herbal A | WRI | **Pharmaceutical** |
| **Voynich Languages Compared** | | |
| Biological | Cosmological | Herbal A | Herbal B |
| Pharmaceutical | Recipes | Unknown | Zodiacs |

# Multiple Discriminant Analysis (MDA)

# Multiple Discriminant Analysis (MDA)

Risk Management, Budgeting, Timeframes and Approach

# Project Management

# Risk Management and Budget
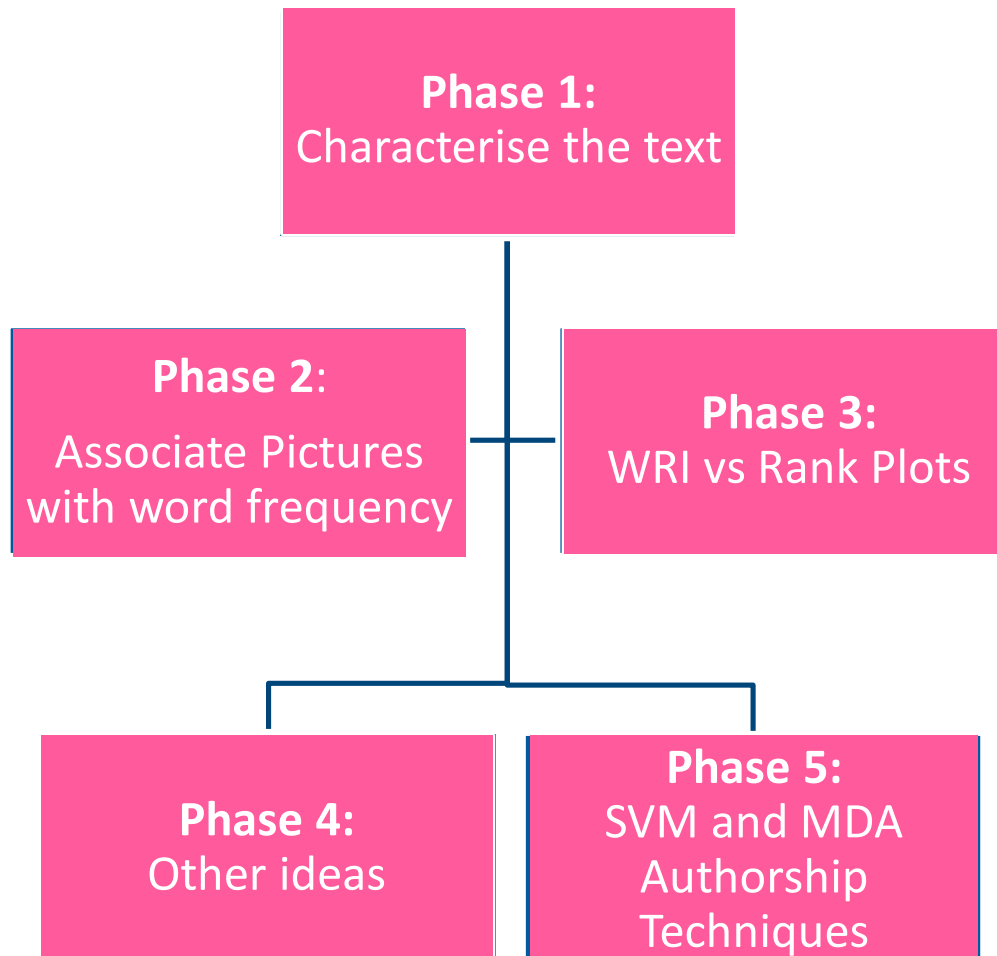
| No. | Risk | Likelihood | Consequence | Risk Level |
|---|---|---|---|---|
| 1 | Not understanding the project correctly and the processes required | Almost Certain | Moderate | Very High |
| 2 | Inaccurate allocation of time and resources to a particular area | Likely | Major | Very High |
| 3 | Health issues due to long periods of time sitting and working at a PC | Likely | Moderate | High |
| 4 | Files and working copies lost | Rare | Major | Medium |
| 5 | UofA Electrical Engineering server down for unknown reasons | Unlikely | Moderate | Medium |
| 6 | Not being able to solve the Voynich Manuscript code | Almost Certain | Negligible | Medium |

- $396.46 (Spent on 3 books, printing and lamination)

# Final Approach

**Phase 1:**
Characterise the text

**Phase 2**:
Associate Pictures with word frequency

**Phase 3:**
WRI vs Rank Plots

**Phase 4:**
Other ideas

**Phase 5:**
SVM and MDA Authorship Techniques
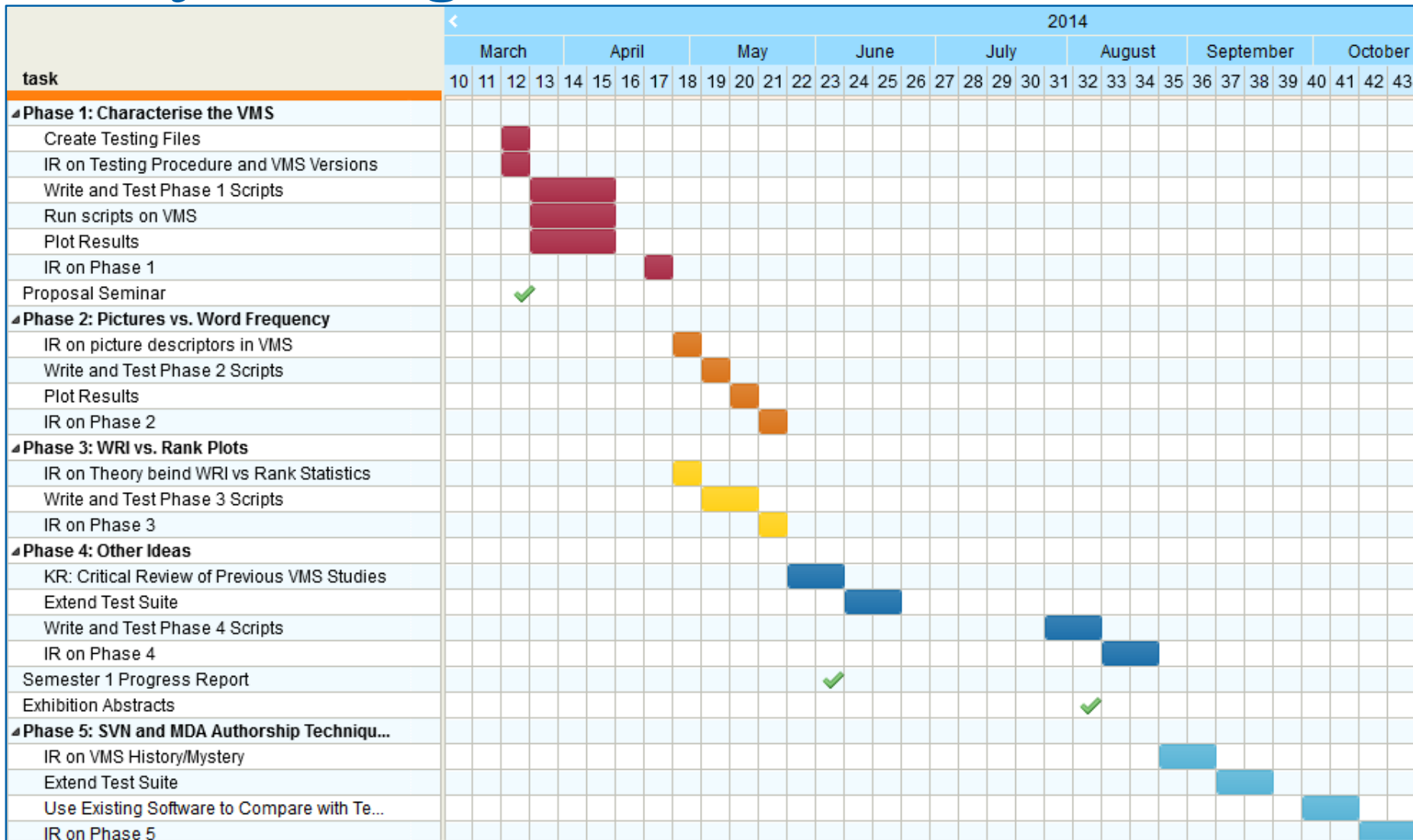
# Team Roles

- Peter
  - Python Code
  - Phase 2
  - Phase 4
  - Compilation of testing material
  - Research as necessary

- Bryce
  - MATLAB Code
  - Phase 3
  - Phase 5
  - Analysis of known 15th Century Text
  - Research as necessary

# Project Progress

Interpretation of Results

# Conclusion

# Conclusions

- The writing and language in the Voynich appears to have evolved over time, making analysis difficult.

- There is a relationship between language and section, but this may not have anything to do with illustrations

- Based on characteristics such as word length distribution and WRI, appears similar to languages such as Hebrew and Latin

- May contain punctuation, based on line characteristics.

- Weak word order, indicating lack of phrases and proper nouns, or perhaps indicating the characteristics of a code

# Future Pathways

- Expand research into word/illustration relationship

- Test the effect of modified alphabets

- Expand research into authorship if possible

- Hidden Markov Model classification of text

- Develop a rule-based grammar for the the Manuscript if possible

- Test characteristics against transcriptions of known 15[th] century codes

# Questions?



Reproduced under the Creative Commons Attribution-Non Commercial 2.5 licence from **xkcd**