

ConTEXTual Scripts: A visual exploration of scripts in their context

Tamar Rucham

A Thesis in the Field of Information Technology
for the Degree of Master of Liberal Arts in Extension Studies

Harvard University

May 2014

Abstract

The visual shapes of characters today stem from a long history of evolution and is deeply rooted in the historical context in which the script evolved. Exploring the visual connections reveals those underlying ties between scripts and gives a different insight into those visual encodings on language.

We wish to enable exploration of the visual themes of scripts in their historical and geographical context. In order to create a platform for comparison between scripts we have analyzed the visual data of 11 different scripts. Calculating on an aggregation of lines and curves to create average glyph information for each script, we have defined a distance method to measure similarity between scripts. To validate our results we have randomly split each character set in two and measured the calculated distance between the parts. Combined with historical and geographical information, we used this data to create an interactive web-based visualization. The visualization allows users to assess the similarity of different scripts and explore the contextual information that may be leading to this relative similarity or disparity. It provides insight into the evolution of the scripts, their geographic context and lays out some of the visual themes these scripts posses.

Acknowledgments

(It's customary to acknowledge the Thesis Director here) This line is four lines below the title, and the title is 1.5 inches from the top.

Table of Contents

Table of Contents	v
List of Figures	vii
List of Tables	x
Chapter 1 – Introduction	1
1.1 Background	1
1.2 Problem and Task Definition	2
1.3 Result	3
1.4 Organization of this document	4
Chapter 2 – Related Work	6
2.1 Scripts visualizations	6
2.2 Linguistics and Language Visualizations	8
2.3 Writing Systems Text Visualization	10
2.4 Art Related Works	12
2.5 Heatmaps	14
2.6 Previous projects	15
Chapter 3 – Implementation	22
3.1 Data Collection and Analysis	22
3.2 Data Visualization	24
Chapter 4 – Methods	26
4.1 Writing Systems vs. Scripts	26

4.2 The Free Type Library.....	27
4.3 Defining the Dataset.....	29
4.4 Evolution Tree	30
4.5 Distance Formulas and Normalization	31
4.6 Validation	33
4.7 First Visualization as a Connected Graph.....	34
4.8 Heatmap.....	36
4.9 Comparison Section.....	37
4.10 Letters Scatter Plots and Clustering.....	39
Chapter 5 – Summary and Conclusions	41
5.1 Discussion.....	41
5.2 Conclusions.....	42
5.3 Future Work	43
References.....	46

List of Figures

Figure 1: The full visualization with the comparison between Arabic and Latin selected.....	4
Figure 2: The result of the animation done by Fradkin on the evolution of the Latin alphabet from its Phoenician origin.....	8
Figure 3: The color coded map and the tree representing sentence structure relation and evolution based on Gell-Mann paper.....	9
Figure 4: A visualization by Teresa Elms mapping the lexical difference between European languages.....	10
Figure 5: An image from Abdul-Rahman et al rule-based visual mappings showing three different layouts of a poem.....	11
Figure 6: A DocBurst analysis of a science text book rooted at the word ‘idea’ with a search on ‘pl’	12
Figure 7: An image out of the wordcollider video.....	13
Figure 8: The Wikipedia article about JPEG converted into an image using the null sets application for representing text in an image.....	14
Figure 9: An image describing Rufiange et al suggestion for combining tree and matrix displays into one. The different color coding in the heat map itself are to identify the tree (green) hierarchy.....	15
Figure 10: Hindi’s space distribution.....	18
Figure 11: Space distributions of 10 selected writing systems.....	18

Figure 12: The space distribution and characters of Thai with a specific pixel and corresponding letters highlighted.....	19
Figure 13: The space distribution and characters of Chinese with a specific pixel and corresponding letters highlighted.....	19
Figure 14: The Turkic origin evolution tree.....	20
Figure 15: The language Bashkir is selected and the corresponding alphabet is displayed.....	20
Figure 16: Additional information from Wikipedia on the selected language and script.....	21
Figure 17: A diagram of the data flow in python code for data collection and analysis.....	24
Figure 18: Sample characters from Chinese, English, Hebrew and Tibetan, with control lines of Bezier curves provided by the freetype library.	29
Figure 19: Details from Wikipedia for the Hebrew script showing the evolution from Egyptian Hieroglyphs to the script under the Parent systems section. This section was used to construct the evolution trees of the scripts.....	31
Figure 20: A connected graph based on the Arial font representing similarity scores of Latin, Greek, Cyrillic, Hebrew, Arabic, Thai and Devanagari. The length of the edge corresponds to the similarity score.	35
Figure 21: A connected graphs based on the Courier New font. The top cluster, which contains Thai and Devanagari, does not exist with the Arial font analysis and the distances between Latin, Greek, Cyrillic and Hebrew are greater.....	35

Figure 22: The heatmap showing the similarity score between the 11 scripts. This image shows the heatmap without the overlaying scatterp plots along the diagonal that exist in the visualization..... 37

Figure 23: The comparison section of the visualization showing the comparison of Latin and Arabic. The scripts names at the top are links to the Wikipedia articles. The scatter plots show the distribution and clusters of the letters as well as the representative characters. The scale shows the similarity score between Latin and Arabic – 0.75. The world map shows the countries which use a Latin derived alphabet (pink), the countries that use the Arabic script (orange) and countries that use both (brown). 38

Figure 24: Scatter plots of 3 scripts – Latin, Arabic, and Telugu. This small multiple shows the difference in distribution which translates into a higher similarity score between Arabic and Latin vs. Telugu. 40

List of Tables

Table 1: Data gathered on the 4 sample characters from different scripts using the freetype library. This displays number of closed contours, total segments, lines and curves for each of the character.....	29
Table 2: The normalized average similarity score of each script randomly split in two over 150 iterations. A score of 1 means identity.	34

Chapter 1 – Introduction

The focus of this thesis is to explore visual themes of different scripts. The word script is used to refer to the character set used by different writing systems to encode spoken language; e.g. the Latin alphabet, is a script used by many different writing systems to express languages such as English, German, French and Malay. In the field of linguistics some research exists in the development of specific scripts (see Section 2). However the goal of this project is not to explore the evolution from past to present, but to compare different current scripts in light of their historical connections.

1.1 Background

The visual aspect of written language played an important role in the history of writing. Before written language as we know it today was *proto-writing*, a picture writing system. Proto-writing uses ideograms or pictograms - graphic symbols that represent ideas or objects respectively - and does not directly translate to specific words from spoken language. The shape of the symbols conveyed the information and a reader would not necessarily need to know the spoken language of the writer in order to gather the information contained in the symbols. When most *true-writing* evolved (*true-writing* being a system in which the entire content of spoken language can be encoded), symbols that represent whole words were used to encode information (logograms). Again the shape of

the symbol was related to the information, but unlike in picture-writing systems, they represented specific words of the writer's language. From there, the phonetic system stemmed, in which symbols represent sounds that are combined to phonetically construct words from spoken language. The shape of the symbol no longer has a specific meaning; rather it can be combined to create many different words. Some scripts, such as the Chinese characters or Egyptian hieroglyphs, preserved both systems such that symbols can function both as logograms and as phonemes. In general, most scripts lost the connection that used to exist between the visual shape and the meaning. However the visual aspect of the scripts we use today stemmed from a long history of evolution.

1.2 Problem and Task Definition

Though we use scripts today as a functional tool, they encompass a rich and intriguing visual aspects and allow a fascinating and unique point of view into the world's history. The shapes of the different glyphs, which began as a representation of animals and tools, have long forsaken their visual origin. However the path that led to their current form and their ties to other scripts, both visually and historically, still exist. While there are several visualizations for spoken languages analysis on the one hand, and some for specific scripts and alphabets evolutions on the other, we have yet to find a visualization that analyzes and allows comparison between different scripts. We feel that such a comparison exposes the visual aspects, which in turn reveal the ties between the scripts and encourage further investigation and exploration of the scripts.

The goal of this project is to expose different scripts and allow a visual exploration of them as a first step. Second to enable an exploration of historical and geographical data in light of the visual connection.

1.3 Result

The end result of this project is a web based interactive visualization that compares 11 “live” scripts (scripts that are in use today), all of which stemmed from the Egyptian Hieroglyph origin. The visualization contains several parts that are linked to each other. On the left there is an evolutionary tree of the scripts leading up to a heatmap mapping the similarity between all scripts. An information area that displays specific information when comparing to scripts follows this.

In the heatmap color intensity maps to similarity – light colored rectangles represent two scripts that are less similar than those with a stronger shade. The heatmap diagonal is the intersection of scripts with themselves (identity) and since the similarity of a script to itself is 1, it contains no added information. Therefore we chose to overlay these rectangles and use this space in order to show a small scatter plot of the letters distribution for the corresponding script. This scatter plot shows the distribution of the letters within the script with respect to the number of lines and curves. All scripts are identically scaled so this small multiple can be used to compare scripts by their letters distribution.

Selecting one of the rectangles representing the similarity between two scripts highlights the two scripts, displays extra information on the right and shows a larger version of the scatter plots (see Figure 1). The comparison section shows representative characters, an enlarged version of the scatter plot with the character set as the markers, and a world map with current script distribution. The name of the script at the top of the information section is a link to its Wikipedia article for further investigation. Each of the two compared scripts is represented by a distinct color throughout the visualization; countries that use both scripts are in a darker shade.

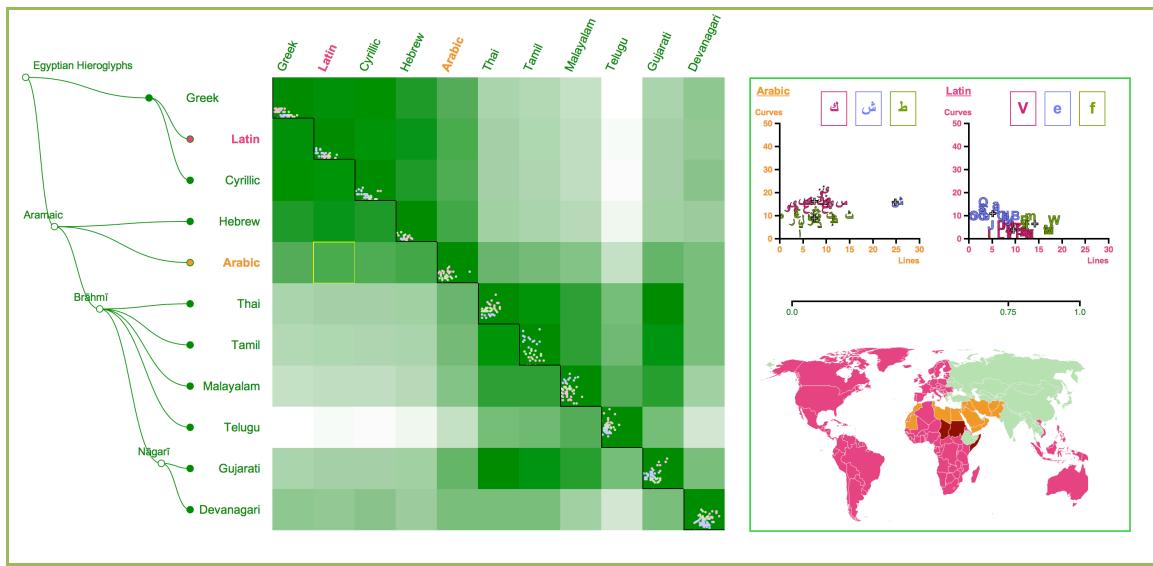


Figure 1: The full visualization with the comparison between Arabic and Latin selected.

1.4 Organization of this document

The related work in chapter 2 covers work done both in the field of linguistics and in the field of data visualization. Chapter 3 contains the implementation details relating technical details of the project. Chapter 4

describes the methods used to implement this project and the reasoning behind them. Chapter 5 is conclusions and summary.

Chapter 2 – Related Work

In this section we review related work, from the domains of linguistics, visualizations and art. We review existing visualizations on scripts, visualizations on languages (evolution etc.) and text visualizations related to writing systems in specific as well as some art projects. We also review heatmaps, which is the main visualization technic used in this project. While some of these projects analyze text visually they do not provide insight into visual themes of the scripts themselves. Lastly we will present previous projects, which formed the groundwork to this thesis project.

2.1 Scripts visualizations

In the field of linguistics some visualizations have been made to illustrate the evolutions of different scripts evolution. These visualizations, while they provide insight into the historical context of the scripts, do not investigate and discuss the visual themes of the scripts. Moreover, these visualizations usually only provide a depth exploration, from past to present, but are less focused on a breadth view to compare parallel scripts.

Fradkin, for his class on the History of Alphabets at the University of Maryland¹, has made several animations (executed by Charlie Seljos) that demonstrate the evolution of different alphabets, for example the evolution of the Latin alphabet from Phoenician displayed in Figure 2. It visually shows how letters gain their visual form and displays the change of style over the course of history. It provides an interesting insight into how the characters transformed into what we know and use today. While it is a fascinating exploration of the history and evolution it offers no analysis of main visual themes, and although the evolution of several alphabets is available, the focus is in exploring each script path separately, making it difficult to compare the visual themes of different modern scripts.

Ada Yardeni, in her Hebrew book HarpatkaOT, lays down the transformation of each character in the Hebrew alphabet from the Canaitic origin to their modern form. The book also visually displays the parallel path to the “sibling” characters in Latin and Arabic scripts, which stemmed from the same origin (XXX add figure). Thus this book allows an in depth, per character, analysis of these scripts (Yardeni, 1993).

¹ <http://terpconnect.umd.edu/~rfradkin/alphapage.html>

Phoenician -- c. 900 B.C.	𐤀 𐤁 𐤂 𐤄 𐤅 𐤆 𐤈 𐤉 𐤊 𐤋 𐤌 𐤍 𐤎 𐤏 ߱ ߲ ߳ ߴ ߵ ߶ ߷ ߸ ߹ ߺ ߻ ߻ ߻ ߻ ߻ ߻
Earliest Greek -- c. 750 B.C. (Western Variant)	Α Β Κ Δ Ε Ζ Η Θ Ι Κ Λ Μ Ν Ξ Ο Π Μ Ο Ρ Σ Τ Υ Φ Χ Ψ
Etruscan -- c. 650 B.C.	𐌈 𐌉 𐌊 𐌋 𐌌 𐌍 𐌏 𐌐 𐌑 𐌒 𐌓 𐌔 𐌕 𐌖 𐌗 𐌘 𐌙 𐌚 𐌚 𐌚 𐌚 𐌚 𐌚
Latin -- c. 500 B.C.	A B C D E F G H I K L M N O P Q R S T V X
C to G -- 3rd cent. B.C.	A B C D E F G H I K L M N O P Q R S T V X Y Z
Latin -- 1st cent. B.C.	A B C D E F G H I K L M N O P Q R S T V X Y Z
Latin -- Middle Ages	A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
Some European Additions	À Ç Ð É Ì Ñ Ø Š Ü

Figure 2: The result of the animation done by Fradkin on the evolution of the Latin alphabet from its Phoenician origin.

XXX Add image from Yardeni's book.

2.2 Linguistics and Language Visualizations

In this subsection we will review visualizations that illustrate relationships between languages and writing systems. The datasets that are visualized in these works are related and in some senses are parallel to the data presented and visualized in this research. They attempt to represent their data in the historical and geographical context and sometimes use methods similar to the ones used in this thesis.

An article by Yasin in Stanford University² describing Gell-Mann & Ruhlen research on sentence structure relation between spoken languages, uses visualization technics similar to our. The said article uses a tree to express the evolution of languages with color-coding to indicate the sentence structure used in each. A map is displayed on the side to allow for the comparison to geographical data (see Figure 3). However the color-coding used in the map

² <http://humanexperience.stanford.edu/languagetree>

does not correlate to the color encoding of the tree, making the visual comparison challenging (Gell-Mann & Ruhlen, 2011).

Teresa Elms created a visualization of the lexical distances between European languages³ based on linguistics research by Tyshchenko (Tyshchenko, 2000). The color and size is used in the visualization to relate historical information about those languages, while the distance and lines are used to relate the lexical distance between them. This allows the user to visually compare the relation, historically and lexically between the languages (see Figure 4).

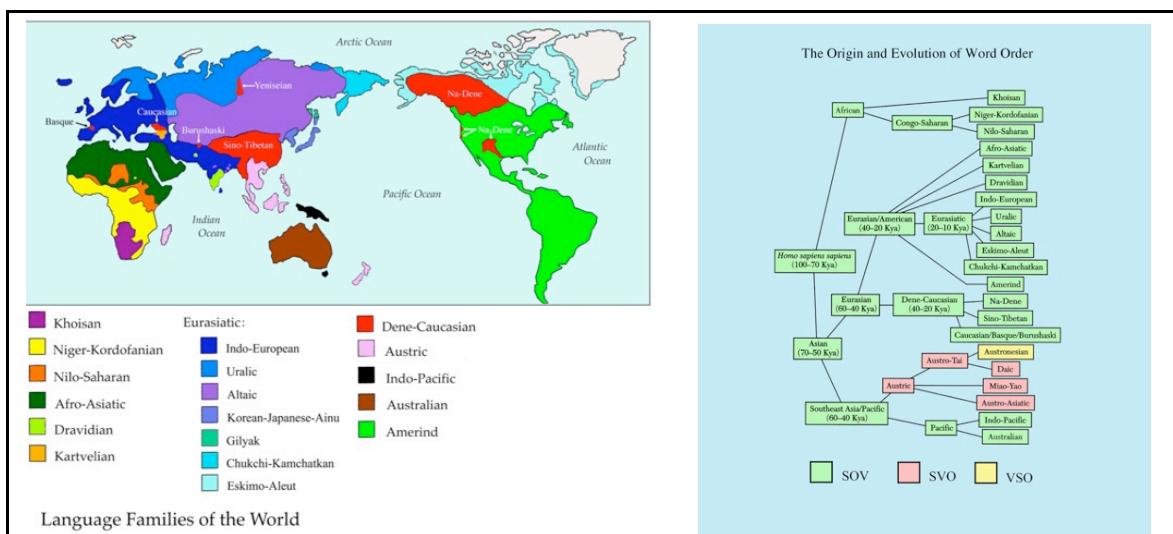


Figure 3: The color coded map and the tree representing sentence structure relation and evolution based on Gell-Mann paper.

³ <http://elms.wordpress.com/2008/03/04/lexical-distance-among-languages-of-europe/>

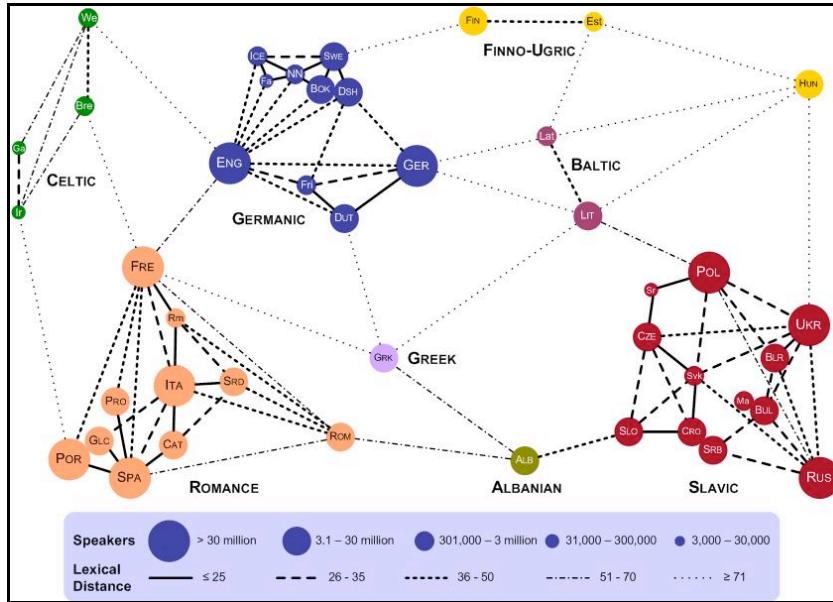


Figure 4: A visualization by Teresa Elms mapping the lexical difference between European languages.

2.3 Writing Systems Text Visualization

In this section we will survey some text visualization relating to writing systems. There is a significant body of work on visualizing text in many different ways, here we will present some that explore the text with an emphasis on the context of writing systems and language structure.

Abdul-Rahman et al created a tool that analyses text in general and poetry in specific in a visual way (see Figure 5). The text is analyzed both phonetically and semantically including the connections and features of the different blocks. The approach also uses a visual tool in order to analyze text, and is an example of collaboration between two remote domains in order to create a new tool for analyzing poetry. However this tool is used to investigate the inner structure of the text rather than its visual form and does not currently promote the

comparison between different scripts of even different text segments (Abdul-Rahman, et al., 2013).

Collins et al in their project DocBurst visualize text based on language structure. They extend the basic idea of visualizing word frequencies in text and give the lexical language context. The user selects a word and based on their database the visualization is populated with all of its hyponyms, coloring and highlighting the ones that occur in text by frequency (see Figure 6). This provides an interesting insight into text from a lexical perspective (Collins, Carpendale, & Penn, 2009).

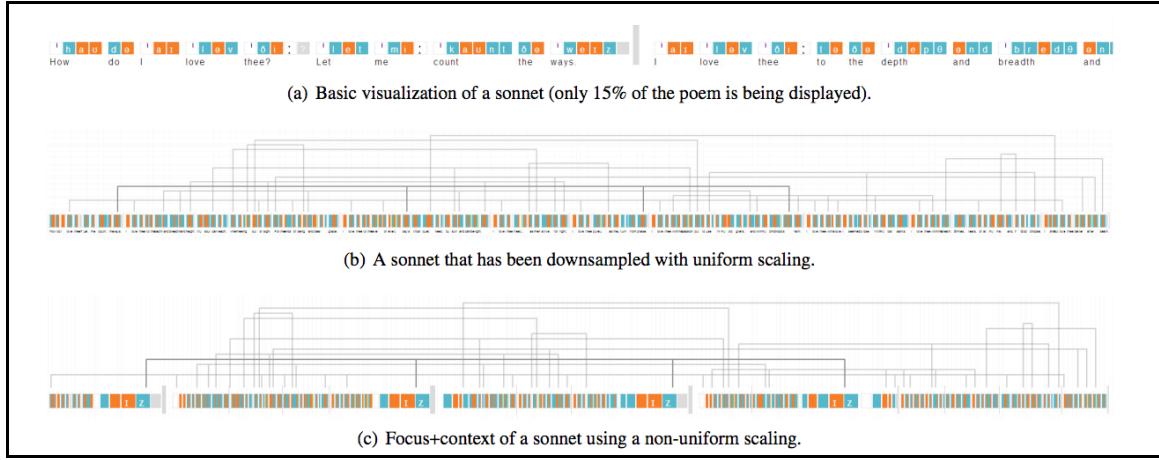


Figure 5: An image from Abdul-Rahman et al rule-based visual mappings showing three different layouts of a poem.

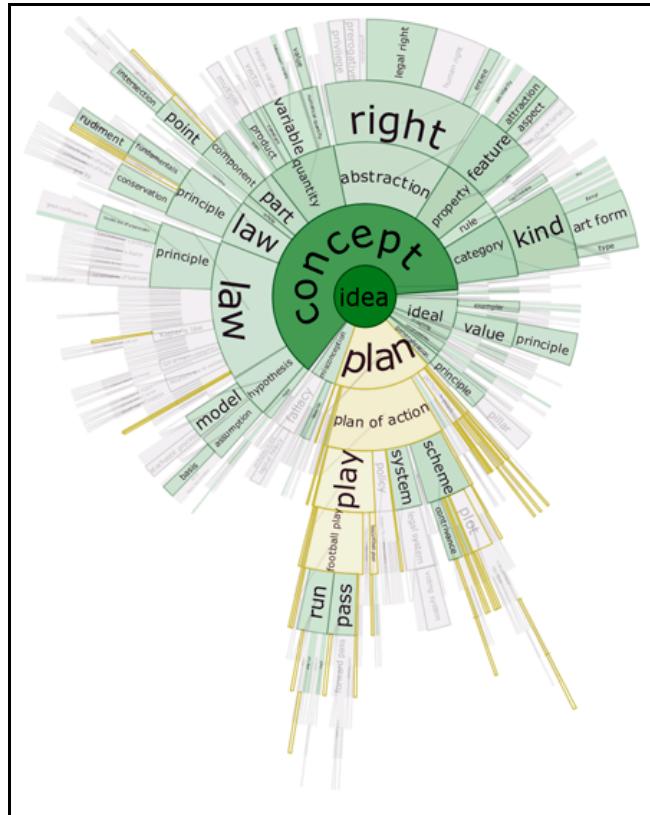


Figure 6: A DocBurst analysis of a science text book rooted at the word ‘idea’ with a search on ‘pl’.

2.4 Art Related Works

*Wordcollider*⁴ is an artistic visualization done by Moritz Heller inspired by particles collision that accelerates two phrases into each other, giving a different visual theme to each letter’s phonetic characteristics (Figure 7). The end result is a visual representation of the two phrases phonetically. Though the approach is interesting and creates an intriguing visual “footprint” of the text, it does not allow for exploration or comparison and is in essence an attempt to visualize sounds (phonetic characteristics). Moreover this approach “visualizes” phonetics, pr in

other words, attempts to give a visual form to sound, while this thesis focuses on analyzing the character sets themselves (Heller, 2012).

Null sets is an artwork that visualizes text files as images, representing their size and structure (see Figure 8). The result is an abstract image of varying color that creates a visual rhythm that is based on the text and therefore represents the structure. Although this project takes a similar approach in the concept of visualizing text and despite the fact it allows for a comparison between different text segments, the end result is not derived by and does not indicate of the visual themes of the original text (Szczepanski & Meaney).

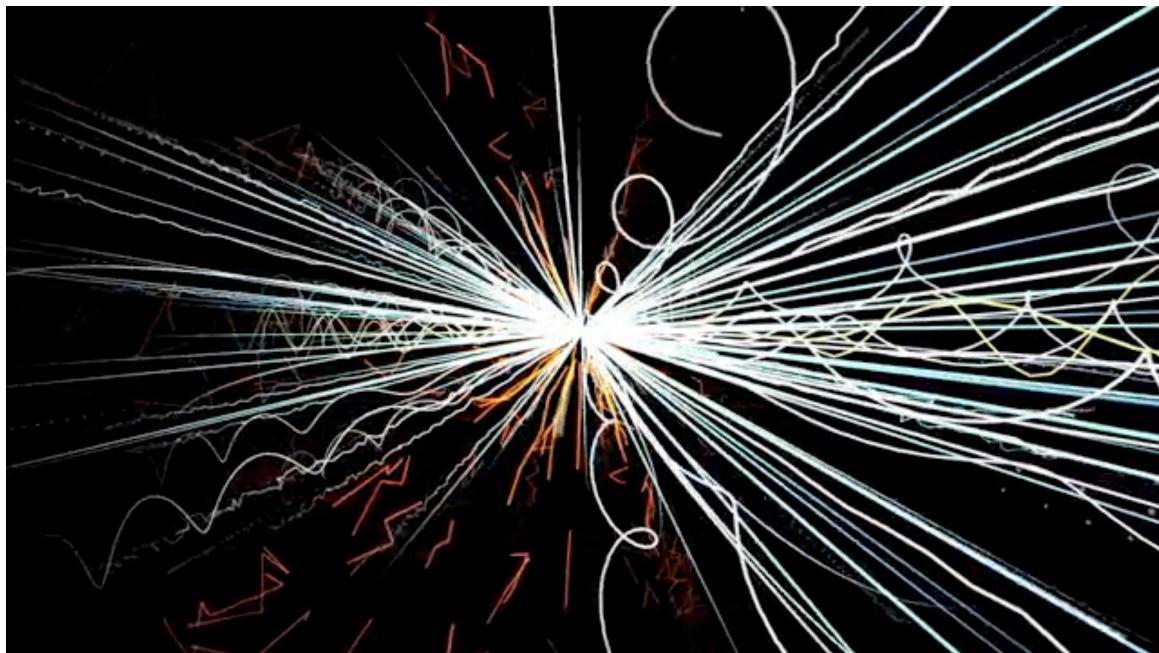


Figure 7: An image out of the wordcollider video.

⁴ <http://vimeo.com/37015401>



Figure 8: The Wikipedia article about JPEG converted into an image using the null sets application for representing text in an image.

2.5 Heatmaps

The main visualization technic used in this thesis is the heatmap; Here we will briefly survey this and related comparable technics.

Wilkinson and Friendly describe the history of heat map matrices while describing their most common use cases and applications (Wilkinson & Friendly , 2008). Rufiange et al suggests a new method of combining trees and heatmap matrices when the data contains multiple dimensions by combining color-coding and other methods of distinguishing the hierarchical data (. While we do use a combined tree and heat map, this approach is more appropriate when dealing with more complex datasets (Rufiange, McGuffin, & Fuhrman, 2012).

Ghoniem et al compare using a heatmap to a linked node when representing rich data. They show that when presenting large and complex data sets, a heat map matrix representation can outperform in some tasks of conveying data to users compared to a linked node. While our data set is small

and can be represented as a connected graph (see more details in section 4.8), the heat map representation selected allows an easier exploration of the data and the possibility of further expanding the data set without compromising the readability and accessibility of the visualization (Ghoniem , Fekete , & Castagliola, 2005).

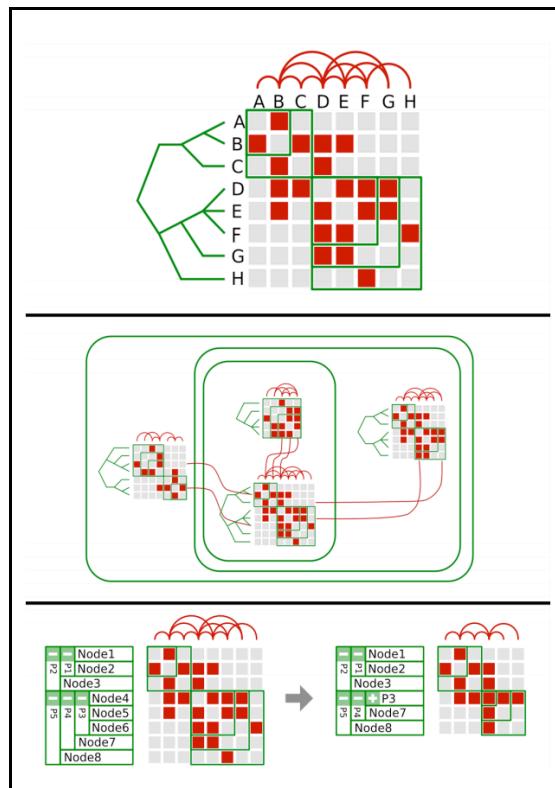


Figure 9: An image describing Rufiange et al suggestion for combining tree and matrix displays into one. The different color coding in the heat map itself are to identify the tree (green) hierarchy.

2.6 Previous projects

The next section describes two projects, *letters space distribution* and *the origin of languages and their scripts*. These projects were done as part of the visualization class with professor Hanspeter Pfister and together they form the groundwork that led to this thesis project.

Letters space distribution is a project that aims to explore the distribution in space of glyphs in text segments of varying writing systems. It allows for the exploration of the space distribution for a writing system used by a specific system and a comparison between the different writing systems. The space distribution is obtained based on text segments which provide sets of characters from each writing system. The bitmap, the image of the character as a 2D array of pixels, is retrieved for each character. Then the bitmaps of all the characters are overlayed, counting the number of times each pixel is visited. This generates a gray scale “heat map” showing how often a pixel is used by the characters in the given text segment and provides a visual insight into how the characters of that writing system are distributed in space. Note that all characters in given text segments were used, therefore if a character repeated multiple times it had a stronger impact on the resulting space distribution. For example Hindi’s space distribution (Figure 10) shows almost all characters have a line across the top, since these pixels are most strongly colored.

The first part of this visualization (Figure 11) contains a small multiple showing the space distribution derived from 10 different languages, on the first row - English, Portuguese, German, French and Malay which all use a latin derived alphabet; on the second row - Hebrew, Arabic, Hindi, Thai and Chinese which each use their own distinct writing system.

Selecting a specific writing system brings up a page that displays a larger interactive version of the space distribution for that system on the left, all the characters that appear in the text as a bubble chart on the right, and the original

text at the bottom of the page (Figure 12). The size of the bubble corresponds to the number of time that character appears in the text which is displayed at the bottom. Hovering over a pixel in the space distribution section on the left highlights that pixel in orange as well as the corresponding letters (i.e. letters that occupy the highlighted pixel) in the bubble chart on the right (Figure 12).

Since the space distribution is based on the frequency and use of specific characters, we expected to find a greater difference between the writing systems that use a latin derived alphabet, however this visualization shows that the space distribution of these systems is remarkably similar. The Malay text is slightly different as it seems to use the letter *a* more heavily than *e*, unlike the other latin derived writing systems. This slight difference as well as the high similarity of the other latin derived scripts may be explained by the fact that the spoken languages English, Portuguese, German and French all have a common Indo-European origin and share a closer history than Malay. Chinese (Figure 13) was probably the most interesting to note as it was the most distinctly different - it showed very dense characters with even distribution over a square, and a significantly higher amount of characters used with a smaller occurrences rate.

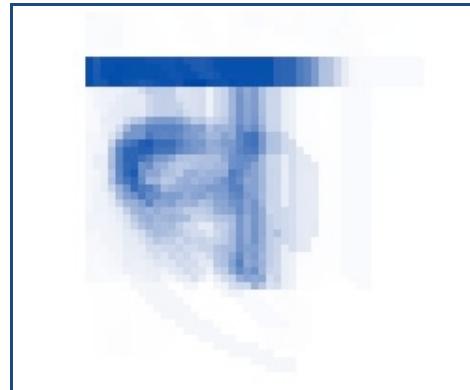


Figure 10: Hindi's space distribution.

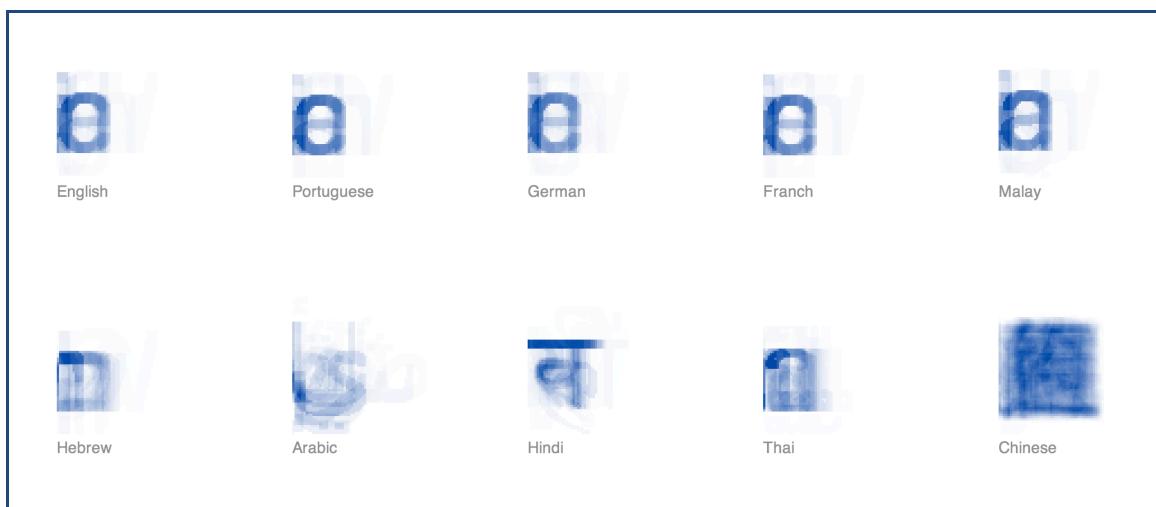


Figure 11: Space distributions of 10 selected writing systems.



Figure 12: The space distribution and characters of Thai with a specific pixel and corresponding letters highlighted.

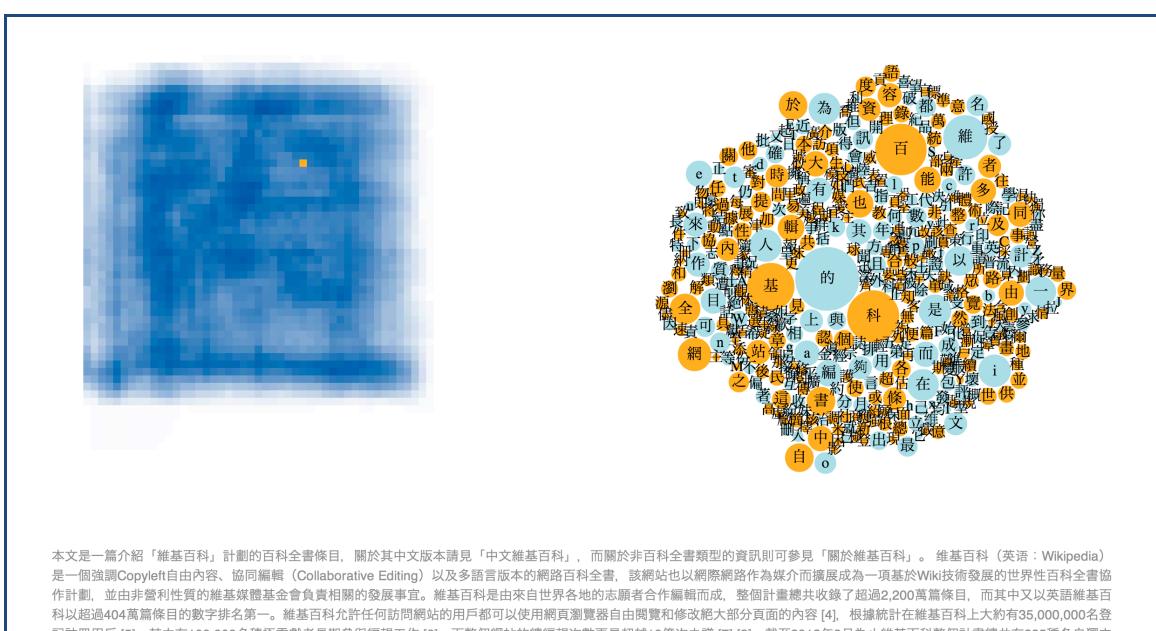


Figure 13: The space distribution and characters of Chinese with a specific pixel and corresponding letters highlighted

The origin of languages and their scripts is a project that aims to collect and visualize the evolution tree of different languages and their writing systems. Origins of spoken languages are displayed on the left and origins of writing

systems are displayed on the right. Clicking on a language family brings up the tree for that family (Figure 14). Languages that have an alphabet connected to it are larger and in full color. Selecting one of those nodes displays the alphabet node at the same level, connected to its origin. The selected node is highlighted while the rest are grayed out (Figure 15). At the bottom there are details on demand - a frame displaying the wikipedia article about the language or the alphabet (Figure 16). Selecting the alphabet will open the writing system tree, with that node already selected. When a script tree is open, selecting one of the selectable scripts will present a list of all languages that use that alphabet.



Figure 14: The Turkic origin evolution tree.



Figure 15: The language Bashkir is selected and the corresponding alphabet is displayed.

Bashkir language	
From Wikipedia, the free encyclopedia	
<p>The Bashkir language (Башкорт телे <i>bashqort tele</i>, pronounced [baʂ. qʊrt.tɬ], ^[i] (listen)) is a Turkic language, and is the language of the Bashkirs. It is co-official with</p>	
Bashkir	
башкорт теле, баʂqort tele	
Native to	Bashkortostan, Russia, Kazakhstan
Ethnicity	Bashkirs
Native speakers	1.45 million (2002)
Language family	Turkic <ul style="list-style-type: none">▪ Kipchak▪ North Kipchak

Figure 16: Additional information from Wikipedia on the selected language and script.

Chapter 3 – Implementation

The chapter describes the implementation details from a technical perspective. The project has two implementation parts – data gathering and analysis and web-based visualization.

3.1 Data Collection and Analysis

The data was collected and analyzed using python with the help of the Numpy library for matrix manipulation and the FreeType_py library for glyph information retrieval. The FreeType_py library uses ttf files (true type font) to provide character information. Many operating systems today come with the most common fonts already installed, however many others can be downloaded. Jack Kilmon⁵ and Ecological Linguistics⁶ both offer font files for ancient scripts. This opens the possibility of a visual exploration of scripts that are no longer in use as a possible continuation of this project.

Using a specific ttf file and Unicode indices per scripts, the code analyzes the characters and outputs the general data to a json file. Secondary processing files then use this output in two ways. One code path generates the required

⁵ <http://www.historian.net/newindex.html>

information for the heatmap – calculating and storing the similarity score between every combination of scripts indexed properly for easy heatmap access. The other code path uses the data to generate chars relations information per script, each saved into a separate json file under the chars folder. This information is further processed to create the clusters of each script (see Figure 17).

The information for the evolution tree was scrapped from Wikipedia using the Pattern library Wikipedia API (*add reference*). The file created in code was then manually adjusted to include only the scripts that are compared in this thesis. More details on scripts selection can be found in the methods section of this document.

⁶ <http://www.lingfonts.com/>

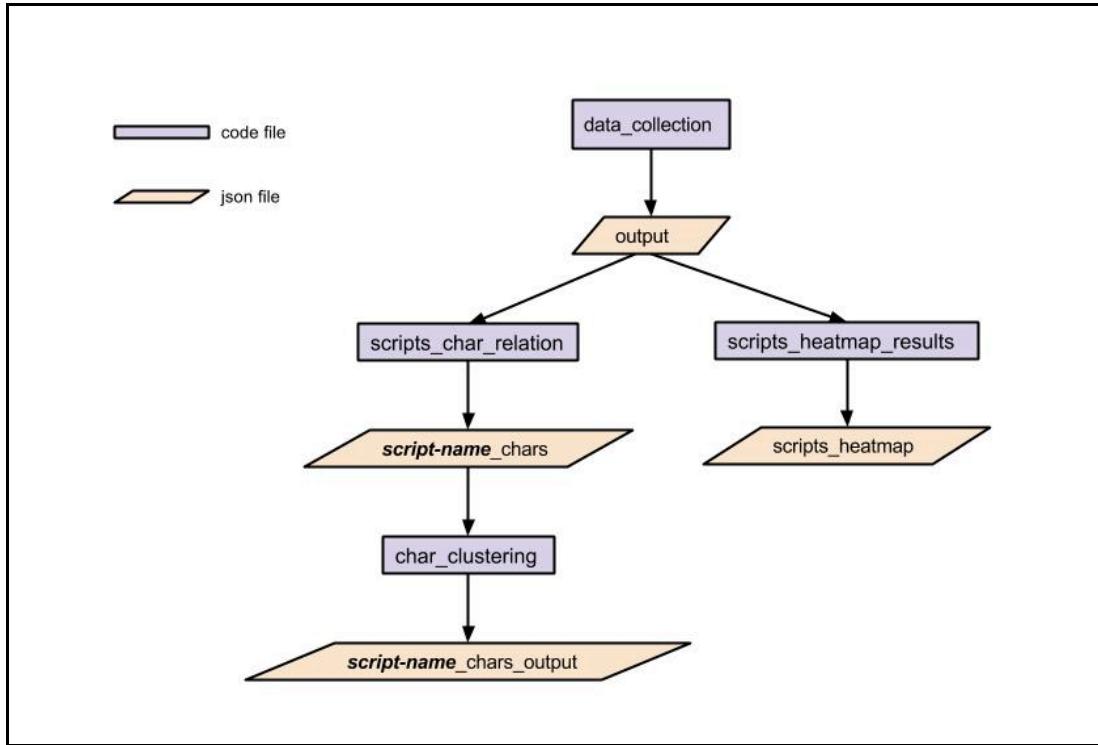


Figure 17: A diagram of the data flow in python code for data collection and analysis.

3.2 Data Visualization

The visualization is a web-based application built with html, css and javascript, using the D3 javascript library (*add reference*). The different parts of the visualization are implemented in separate files and are generated as multiple individual SVG rather than a single element in order to provide encapsulation of each element. All these visualization parts have a javascript file that is responsible for creation and setup and a css file for individual styles.

The **tree** section implements the evolution tree as... (*add reference to technical resource used*)

The **heatmap** section implements the main heatmap comparing the scripts along with the top labels (add reference to technical resource used:

<http://blog.nextgenetics.net/?e=44>)

The **data** section implements the comparison section on the right and includes the world map based on the topojson library (*add reference*). This section uses a static data file, which contains information such as links to their Wikipedia articles and geographic information by country codes. In addition it uses a json file containing information on country boundaries to generate and manipulate the world's map. This file⁷ indexes countries by their ISO 3166-1 numeric codes. Since the list of countries per script in the static data file uses the countries ISO 3166-1 alpha-3 codes a file mapping the two was downloaded⁸ and the relevant information was extracted using a simple python script.

Finally, the **char relation** section implements the per-script scatter plots (individual SVG elements), both the large version, which includes the letters, and the simple smaller version displayed inside the heatmap. This file uses the clusters of each script generated by the python script, in order to mark the clusters and the representative characters.

⁷ Coordinates downloaded from <http://bl.ocks.org/mbostock/raw/4090846/world-50m.json>

⁸ Country information downloaded from
<https://github.com/mledoze/countries/blob/master/countries.json>

Chapter 4 – Methods

The section describes the methods, algorithms and approaches taken in this thesis, from data collection to analysis and finally to visualization.

4.1 Writing Systems vs. Scripts

First we had to define the scope of the project, whether we focus on writing systems, as in the project Letters Space distribution or character sets (scripts), as in The Origin of Languages and Their Scripts. A visual aspect of a writing system is far more complex than just its script, as it also includes amongst other things, length of words and frequency of character use. While considering a comparison of the visual aspect of writing systems, we thought of using large bodies of text in different writing systems to run the analysis on. The bible was an obvious selection for its length and availability in many writing systems. This would have required some sort of validation that the length and language variety of this text is sufficient to describe the writing system. Another option is to use dictionaries, however dictionaries are not a representation of a writing system in text. Dictionaries being essentially lists of words don't seem sufficient to represent a writing system which also includes sentence structure which effects the visual aspect of a writing script. Eventually we came to the conclusion the writing systems are more vaguely and broadly defined, and

though it may be an interesting future project, for the scope of this thesis we will focus on scripts.

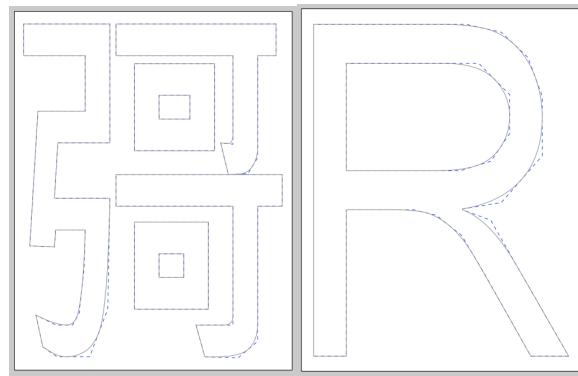
The basis for which scripts to compare was the scripts used in the previous project *The Origin of Languages and Their Scripts*, seeing we have already obtained their evolution tree in a process described later. In addition, in order to keep the scope of this thesis from getting too broad, we decided to focus on the Egyptian Hieroglyphs originated scripts only. Scripts coming from the Oracle Bone script origin such as the Chinese character set, tend to be extremely large (over 3000 characters) and sometimes not easily defined - some writing systems use a mixture of character sets, or derived characters (*reference*).

4.2 The Free Type Library

The analysis of the scripts is done using the freetype-py library for python. The library provides both bitmap and vector information on glyphs based on a provided ttf file (true type font – *reference?*). We started by looking at the glyph-vector example provided with the library. Note that this example uses matplotlib to plot the glyph. The glyph information is initially divided into closed contours - closed lines, for example the letter **O** has two contours, the outer circle and the inner circle while the letter **S** has only one contour. The closed contours are in turn divided into segments that can be either a straight line or a quadratic Bezier curve (or curves). *Add clear and concise description of the outlines and how the*

control points define the line in TrueType based on <http://www.truetype-typography.com/ttoutln.htm> ??

Figure 18 shows a few examples of characters from different writing systems with their control lines (dashed). These are single characters from a few selected scripts we have gathered in order to explore how to access and analyze the data available with this library. These are characters from the Chinese, Latin, Hebrew and Tibetan character sets (from left to right respectively). The statistics provided by the library for each character is as presented in Table 1. This experiment demonstrates the great difference in number and type of segments that can exist. First there is a significant difference between the number of closed contours and segments between the Chinese and Tibetan selected characters on the one hand and Hebrew and Latin characters on the other. In addition, the Chinese and Tibetan samples differ in the type of segments as well – the Chinese sample contains many straight-line segments while the Tibetan one has many curves. Indeed if we examine the characters ourselves, the variance in complexity and nature that is reflected in the data is clear to see.



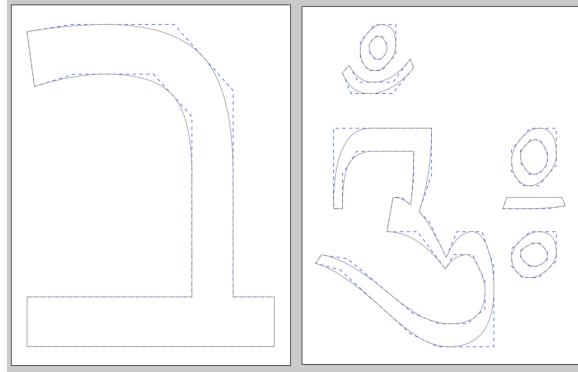


Figure 18: Sample characters from Chinese, English, Hebrew and Tibetan, with control lines of Bezier curves provided by the freetype library.

	Chinese char	English char	Hebrew char	Tibetan char
Closed contours	7	2	1	9
Total segments	57	17	12	58
Straight lines	47	11	8	9
Curves	9	6	4	49

Table 1: Data gathered on the 4 sample characters from different scripts using the freetype library. This displays number of closed contours, total segments, lines and curves for each of the character.

4.3 Defining the Dataset

Therefore we knew we wanted to use scripts to which we have font files.

At this point the question of font came up. There is great variance in visual aspect of characters between different fonts, mainly between serif and san serif fonts (*what is serif / san serif*). Even within those families we wanted to choose a font that can serve as a baseline for comparison. Therefore the Arial Unicode font was selected for representing all the scripts. The Arial Unicode font covered the most scripts of the ones we have defined in the project *The origin of languages and their scripts* and more importantly this font was created to be **as generic and blend as possible** (XXX reference quote). Out of the scripts we started with, the Arial Unicode ttf file supported 11 scripts – Greek, Latin, Cyrillic,

Hebrew, Arabic, Thai, Tamil, Malayalam, Telugu, Gujarati, and Devanagari.

Therefore we selected these scripts as our dataset to explore. For each of these scripts we manually defined the ranges of letters based on their Unicode charts (*reference*).

4.4 Evolution Tree

Before diving into the glyph analysis and script comparison, I'd like to take a moment to describe the method used to obtain the evolution tree of the scripts. This work was done as part of the project *The Origin of Languages and Their Scripts*. The codes for all the languages that have a wikipedia article about Wikipedia were scraped using the Pattern library Wikipedia API. Then the language name was matched using an external source. That name was used to search the Wikipedia article about that language, which was in turn scraped and the linked writing system was followed. Articles about languages and scripts have specific structure for the origins as can be seen in Figure 19. Languages that did not follow this structure were ignored. The trees were constructed by the python script, starting for each origin going down, checking if the node existed for each level. The data was stored in json files. After the generation of the data by the script, it was manually cleaned. There were some inconsistencies with trees starting points (for example some scripts ended their parent systems in the proto-sinaitic origin, though that is a child of the Egyptian hieroglyphs) as well as references that use slightly different names (such as Latin script vs. Latin alphabet). Also, to reduce the overwhelming amount of data, some of the

esoteric language trees and branches were removed. To narrow down the tree size, nodes that have a single child were ignored (single line of connection).

Hebrew alphabet	
	אַלְפָבֵיַת עֲבָרִי
Type	Abjad (for Hebrew, Aramaic, and Judeo-Arabic) True Alphabet (for Yiddish)
Languages	Hebrew, Yiddish, Ladino, and Judeo-Arabic (see Jewish languages)
Time period	3rd century BCE to present
Parent systems	Egyptian hieroglyphs <ul style="list-style-type: none"> • Proto-Sinaitic • Phoenician alphabet • Aramaic alphabet • Hebrew alphabet
Sister systems	Nabataean Syriac Palmyrenean Mandaic Brāhmī ¹ Pahlavi Sogdian
ISO 15924	Hebr, 125
Direction	Right-to-left
Unicode alias	Hebrew
Unicode range	U+0590 to U+05FF ↴ U+FB1D to U+FB4F ↴

Figure 19: Details from Wikipedia for the Hebrew script showing the evolution from Egyptian Hieroglyphs to the script under the Parent systems section. This section was used to construct the evolution trees of the scripts.

4.5 Distance Formulas and Normalization

For every character in each script we calculated the number of lines and the number of curves (we also got the number of closed contours but that information is not used in this project). Based on these results we calculated the average number of lines and the average number of curves per script. Next we

needed to define a distance function, how do we measure a similarity, or disparity, between scripts. This formula is used both to measure distance between characters based on lines and curves for clustering and for the similarity score between scripts based on averages.

The first formula we used was simply an accumulation of the differences between the number of lines and number of curves. Therefore this produces a value starting from 0 (identity) going up to the greatest difference (which for this formula was around 25 – comparing Greek and Telugu):

```
absolute_value(script_1_lines - script_2_lines)
+ absolute_value(script_1_curves - script_2_curves)
```

Later when we came to create clusters of the letters (described in section XXX) we started considering the number of lines and number of curves as our coordinates in a two dimensional plane. Thinking of clustering we found that a small tweak to the formula – making it a straightforward Euclidean distance formula worked better. Besides making the clustering algorithm simpler, it gave slightly better results in the validation process described in Section 4.6. Therefore the formula that is now used to calculate the similarity score is:

```
absolute_value(script_1_lines - script_2_lines)
/ absolute_value(script_1_curves - script_2_curves)
```

Finally we normalized the results of script similarity scores to a value from 0 (greatest disparity) to 1 (identity). This makes comparison, further calculation and representation of the data simpler and easier to read. To normalize the data

we take the calculated score minus the maximum found difference and multiply negative one in order to reverse the direction – the distance similarity score gets higher as the disparity increases. The condition preceding the formula below is of course to avoid division by 0 when the value equals the maximum data.

```
0 if value == maxData else (value - maxData) * (-1) / maxData
```

4.6 Validation

In order to measure the performance and validate the algorithm described in the previous section, we created a validation script. The code randomly splits the character set in two and evaluates the similarity score. The assumption is that the two subsets of the same script should come out with a good similarity score (0 being identity - same character set). This was done 150 times for each of the scripts and outputs the normalized average score received. The results were surprisingly good for such a simple algorithm, at worst getting a score of ~0.81 (see Table 2). It seems to do better with the scripts that weigh more on straight lines, such as Latin, Greek and Hebrew, than with the more curved scripts such as Tamil, Gujarati and Malayalam.

After changing the formula from the accumulation of differences into an Euclidean distance formula, we ran the validation again. The results shows that the second formula gives slightly better results; the average score of the first algorithm being ~0.909 and the second algorithm providing an average score of 0.915. The detailed results are shown in Table 2.

Script	Validation I	Validation II	Closest Script	Closest score
Telugu	0.91	0.92	Malayalam	0.73
Cyrillic	0.93	0.94	Greek	0.98
Greek	0.94	0.95	Cyrillic	0.98
Malayalam	0.89	0.89	Gujarati	0.85
Thai	0.92	0.92	Gujarati	0.99
Latin	0.93	0.94	Greek	0.96
Gujarati	0.89	0.90	Thai	0.99
Hebrew	0.94	0.94	Latin	0.92
Devanagari	0.93	0.93	Gujarati	0.58
Arabic	0.91	0.91	Hebrew	0.78
Tamil	0.81	0.83	Thai	0.95

Table 2: The normalized average similarity score of each script randomly split in two over 150 iterations. A score of 1 means identity.

4.7 First Visualization as a Connected Graph

As an experiment we calculated the similarity score for 7 of scripts – Latin, Greek, Cyrillic, Hebrew, Arabic, Thai and Devanagari – and plotted them as a connected graph (Figure 20). The length of the edges, which defines the distance between the nodes, uses the similarity score as the weight. The tight group of 4 near the bottom contains Latin, Greek, Cyrillic and Hebrew, the more remote outliers are Thai and Devanagari, and Arabic is somewhere in the middle.

Curious to see the importance of the font selection, we calculated the similarity score using the Courier New font, and re-plotted (see Figure 21). Some elements were preserved, but there is a clear difference in the similarity score for Thai and Devanagari, which together form the tight cluster at the top of this graph.

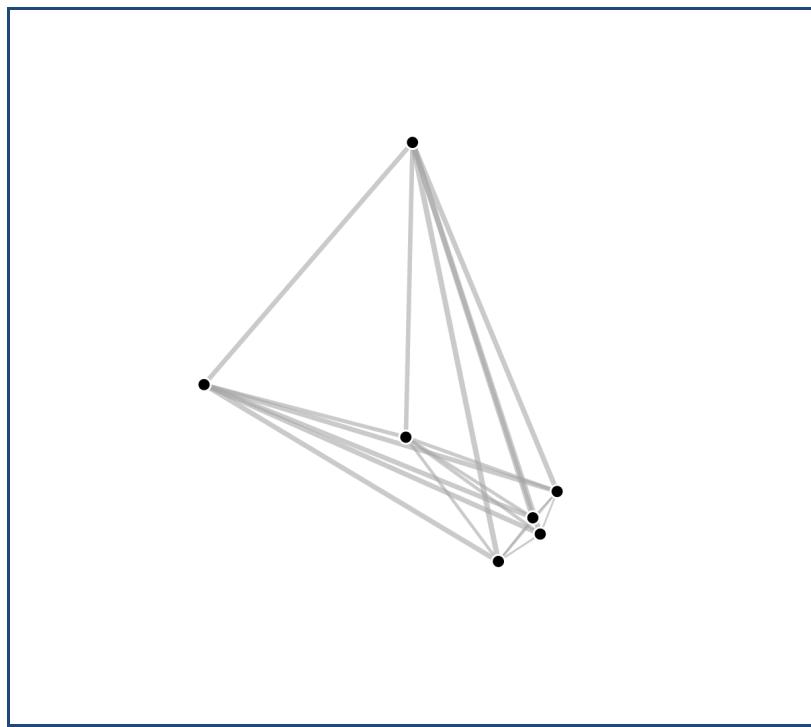


Figure 20: A connected graph based on the Arial font representing similarity scores of Latin, Greek, Cyrillic, Hebrew, Arabic, Thai and Devanagari. The length of the edge corresponds to the similarity score.

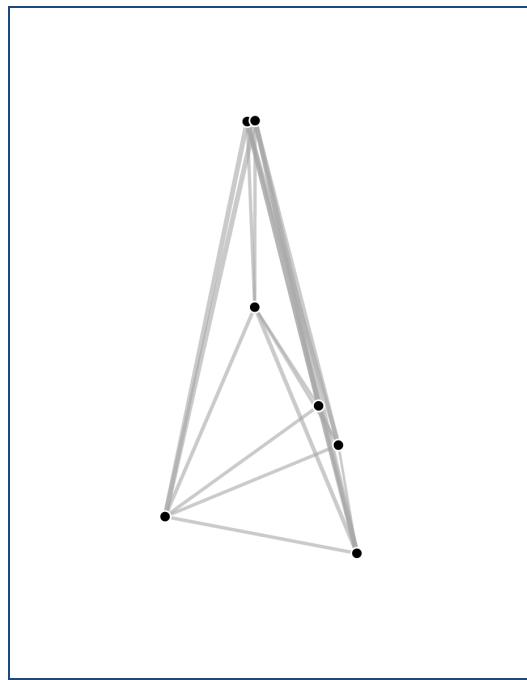


Figure 21: A connected graphs based on the Courier New font. The top cluster, which contains Thai and Devanagari, does not exist with the Arial font analysis and the distances between Latin, Greek, Cyrillic and Hebrew are greater.

4.8 Heatmap

Since our focus is on the comparison between the scripts and this is not a graph in a traditional sense we decided to use heatmaps instead (see Figure 22). In order to display our data in a heatmap we generated the similarity score for every combination of the scripts – 55 values – to be the base of our heatmap. The heatmap was created based on the example Damian Kap [resented in his blog post⁹. The intensity of the color corresponds to stronger similarity between the scripts that intersect on a specific rectangle – white representing the greatest disparity and the strongest green on the diagonal shows the identity of a script with itself. These rectangles are linked to the Comparison section described in Section 4.9. The ordering of a heatmap holds a great importance on the readability (XXX reference). Since the number of scripts is fairly low and there was an immediate correlation with the evolution tree, the scripts were manually sorted based on their similarity score. It should be interesting to run a hierarchical clustering algorithm on the scripts themselves to see how well such an algorithm corresponds to the evolution tree.

⁹ <http://blog.nextgenetics.net/?e=44>

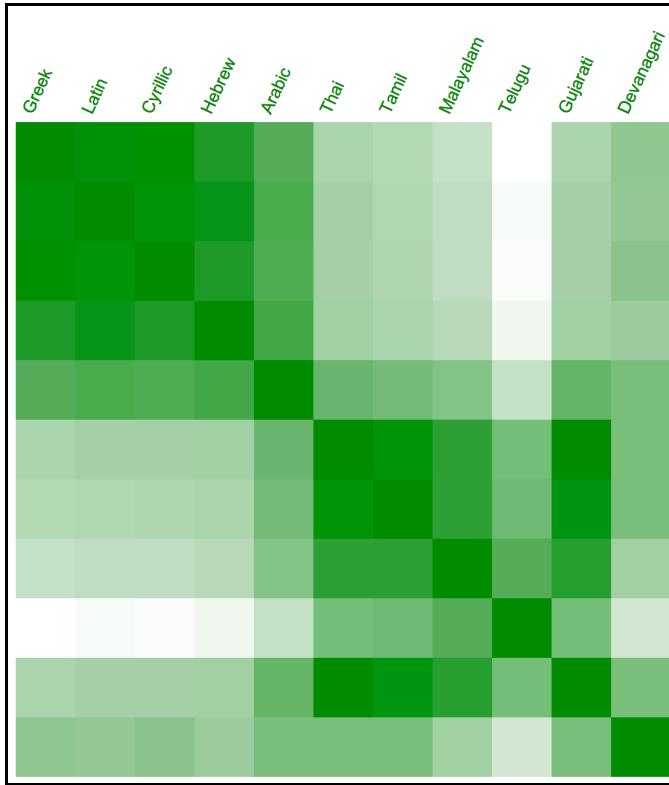


Figure 22: The heatmap showing the similarity score between the 11 scripts. This image shows the heatmap without the overlaying scatter plots along the diagonal that exist in the visualization.

4.9 Comparison Section

The comparison section is meant to provide extra data when comparing two scripts. It is linked to the heatmap so selecting one of the rectangles in the heatmap, brings up the information on the corresponding scripts (which intersect in that rectangle). The idea is to provide insight into the two scripts and what contributes to their similarity score (which is also provided as a number on a scale as part of this section). Therefore a significant part of the exploration is done in this section. The name of the script at the top is a link to its Wikipedia page to allow users to further investigate a script. Below that, the letters scatter plot which its clusters and the representative characters are displayed, more on

that in Section 4.10. At the bottom of this section a map of the world is presented with the countries in which those scripts are used are highlighted. The information for this map was manually compiled based on the Wikipedia article on Writing systems¹⁰, which was created and updated based on articles and books on the topic (Coulmas, 1999), (Daniels & Bright, 1996), (Sampson, 1985). Though the different scripts of India are mostly used in different subregions of the country, we decided to mark the entire country for all script since their proximity is great when viewing on a global scale, and the development overhead for impleneting subregions seemed unnecessary.

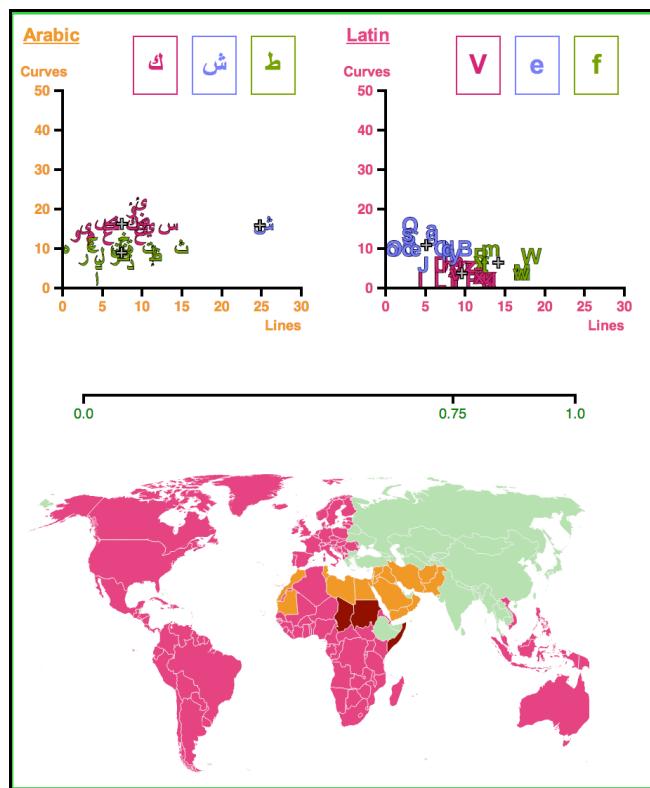


Figure 23: The comparison section of the visualization showing the comparison of Latin and Arabic. The scripts names at the top are links to the Wikipedia articles.

¹⁰ http://en.wikipedia.org/wiki/Writing_system

The scatter plots show the distribution and clusters of the letters as well as the representative characters. The scale shows the similarity score between Latin and Arabic – 0.75. The world map shows the countries which use a Latin derived alphabet (pink), the countries that use the Arabic script (orange) and countries that use both (brown).

4.10 Letters Scatter Plots and Clustering

Visualizing the letters in a scatter plot provides valuable insight into the algorithm and the similarity or disparity between scripts. The axes are the number of lines and the number of curves per character. This shows how not only the difference in ratio between lines and curves, but also the total number of segments affects the similarity score. If we examine the scatter plots in Figure 24 we can easily see that Latin and Arabic has much less lines in general and less curves specifically. While a character in Latin will have less than 30 segments, characters in Telugu go up to 50 segments and more, most of which are curves. These differences translate into the similarity score, which is 0.75 between Latin and Arabic, vs. 0.023 and 0.27 between Telugu and Latin and Arabic respectively.

In order to derive representative characters for each script we have run a clustering algorithm on these data points, using the number lines and curves as our coordinates. Since the classic kmeans clustering algorithm has a high risk of falling into local minima, we have decided to use a bisecting kmeans algorithm as described by Peter Harrington (Harrington , 2012). Once we found the clusters, using the Euclidean distance function described in Section 31, we used the clusters centers in order to find the character closest to it. We found that having

three representative characters provided the most useful insight into the elements of the script. Since these clusters represent different “areas” of the letters, we perceive them as having varying characteristics; therefore we chose them to represent the script. Having only two representative clusters did not provide enough sample and presenting four characters or more was too much and introduces redundancies.

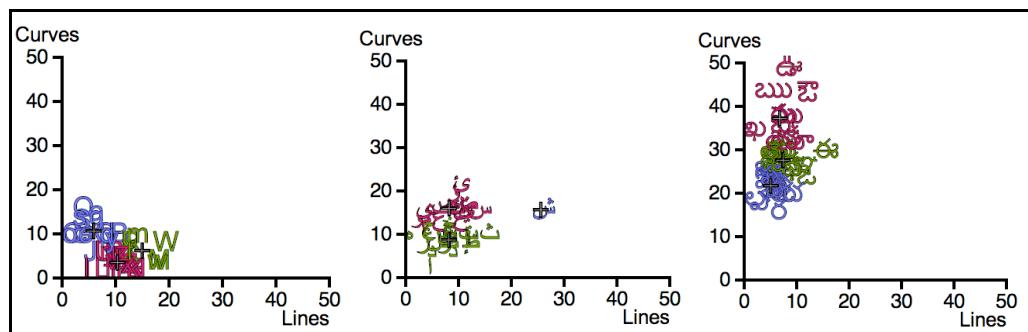


Figure 24: Scatter plots of 3 scripts – Latin, Arabic, and Telugu. This small multiple shows the difference in distribution which translates into a higher similarity score between Arabic and Latin vs. Telugu.

Chapter 5 – Summary and Conclusions

This chapter describes conclusions from the data exploration and visualization, discussion about lessons that were learned working on this thesis and future work that can be done on the topic.

5.1 Discussion

The first challenge we had to deal with was the size and fuzziness of the data. There are a great number of scripts in use today and the historical information on each is staggering and poses an interesting problem of data collection. Since this data is not numeric but is more of a collection of narratives, gathering this data is not straightforward. We had to combine web scrapping, manipulating downloaded and gathered data through scripts and manually and other methods in order to obtain the data used in this project. In addition we had to keep our focus and decide what data to explore and represent and which to put aside. This includes the decision to focus on scripts instead of writing system, the choice of scripts that derive from a single origin and focusing on one font. While all of these options pose interesting opportunities, we decided to avoid spreading too much. Another related challenge was the lack of prior work. This meant we had to gather and explore everything from scratch and constrained our scope even further. All of these challenges encouraged us to limit the scope and focus on a subset of the data and a narrower analysis.

Once we defined the data, we had to consider how best to visualize it. Instead of having multiple views into the data separated into multiple tabs or windows, we decided to condense all the information into a single linked view. This allows the user to explore, compare and investigate within the same context. While heatmaps are less familiar to users, they allowed for a simple display of the information and moreover a linking of a comparison (between two scripts) becomes intuitive.

5.2 Conclusions

The strongest most obvious conclusion was the strong correlation between the similarity clusters and the evolution tree as well as geographic information. The heatmap clearly shows two clusters, one containing the scripts closely related to Europa and the other containing the Indian and other far east scripts from the Brāhmī family (which we will refer to as the “west” and “east” cluster for ease of use). While this correlation is somewhat expected it was still exciting to see it present itself so clearly in our results and it served as some validation of our algorithm. Another interesting data point was Arabic, which serves almost as a bridge between the two clusters, with almost all of its intersections comparing above 0.5. This can be explained both geographically – since it is common in the countries between Europe and India – and by looking at historical data. In ancient times Arab traders were the main bridge between the

two worlds, transporting goods (silk and spice) from east to west (XXX reference).

Looking at the small multiple of the clusters we notice that the “west” cluster has a tendency to have less curves compared to straight lines and less segments in general. The “east” cluster is much more curved and tends to have more segments. A couple of outliers that may require more exploration are Devanagari and Telugu. Devanagari, though it is geographically and historically close to the other Brāhmī scripts, scored no more than 0.58 and not much lower in scores with the “west” cluster. The high number of segments on one hand and the high number of straight lines on the other is the cause. Telugu on the other hand, while it is relatively close in comparison with the other “east” scripts, shows the greatest difference with the “west” cluster. Again looking at the small multiple we find that Telugus characters are furthest from the origin, which translates into very complex characters. Every character in Telugu as at minimum 20 segments, and most of the segments are curves. The combination of these two features is what makes Telugu so “far” from the “west” cluster.

5.3 Future Work

We consider this project the very beginning of a journey into exploring scripts from a visual perspective. There is much work that can be done, both as a continuation of this project directly and new projects stemming from this one.

First, this project can be extended to contain more scripts. It will be particularly interesting to see a comparison with scripts that do not stem from the Egyptian Hieroglyphs origin such as the Chinese character set. In addition a comparison to scores gathered using different fonts in general and Serif scripts in particular (unlike the Arial San-Serif font used) may lead to very interesting results. Some other features that can be added include the ability to compare more than two scripts and more information upon comparison. We would have like to add the ability to brush entire sections of the heatmap to compare multiple scripts, however this will also require rethinking of the comparison area. We would have also liked to add more data in the comparison such as the estimated creation time for the given scripts. Last but not least, instead of using representative characters, a machine-learning algorithm can be written to learn the characteristics of a certain script and produce a generic character based on it. *XXX Add reference to generic numeral creation research.* This type of algorithm may extract the most crucial visual elements of a script and highlight its visual themes.

Next projects in the topic can include a good visualization of the development of different scripts from shared origins. This holds both a linguistics interest and a fascinating visualization challenge to express visual development in multiple concurrent paths. It will also make an interesting project to attempt the problem of writing systems and not scripts. This opens many questions, both from data collection and analysis as well as visualization and focus. On the other hand there is some interesting work done in the field on comparing different

European writing systems that such a project will be able to build upon. Further from this research and harder to implement, lies the idea of comparing ancient origins. The origins of modern scripts tend to be close to picture writing. We feel a compelling visual opportunity can be found in comparing Egyptian Hieroglyphs, Oracle Bone Scripts and the Mayan Hieroglyphs for example. Each of these scripts has distinct visual themes and a captivating historical context, which can form the ground for interesting explorations and visualizations.

References

- Abdul-Rahman, A., Lein, J., Coles, K., Maguire, E., Meyer, M., Wynne, M., et al. (2013). Rule-based Visual Mappings – with a Case Study on Poetry Visualization. *Eurographics Conference on Visualization (EuroVis)*. Leipzig: Computer Graphics Forum.
- Collins, C., Carpendale, S., & Penn, G. (2009). DocuBurst: Visualizing Document Content Using Language Structure. *Eurographics/IEEE-VGTC Symposium on Visualization (EuroVis)*. Berlin: Computer Graphics Forum.
- Coulmas, F. (1999). *The Blackwell Encyclopedia of Writing Systems*. Oxford: Wiley-Blackwell.
- Daniels, T. P., & Bright, W. (Eds.). (1996). *The World's Writing Systems*. Oxford University Press, USA.
- Gell-Mann, M., & Ruhlen, M. (2011). The origin and evolution of word order. *Proceedings of the National Academy of Sciences of the United States of America*.
- Ghoniem , M., Fekete , J.-D., & Castagliola, P. (2005). On the readability of graphs using node-link and matrix-based representations: a controlled experiment and statistical analysis. *Information Visualization* , 22.
- Harrington , P. (2012). Grouping unlabeled items using k-means clustering. In P. Harrington, *Machine Learning in Action*. Shelter Island, NY, USA: Manning Publications Co.
- Heller, M. (2012, February 18). *Wordcollider*. Retrieved February 24, 2014, from Vimeo: <http://vimeo.com/37015401>
- Rufiange, S., McGuffin, M. J., & Fuhrman, C. P. (2012). TreeMatrix: A Hybrid Visualization of Compound Graphs. *Computer Graphics Forum* , 31, 89-101.
- Sampson, G. (1985). *Writing Systems*. Stanford, USA: Stanford University Press.
- Szczepanski, A., & Meaney, E. (n.d.). *Info*. Retrieved February 24, 2012, from Null Sets: <http://evanmeaney.com/ns/index.html>

- Tyshchenko, K. (2000). *The Metatheory of linguistics*. Ukrain: Osnovy.
- Wilkinson, L., & Friendly , M. (2009). The History of the Cluster Heat Map. *The American Staticians* , 63 (2), 179–184.
- Yardeni, A. (1993). *HarpatkaOT*. Jerusalem, Israel: Karta.