

# **Exploration and application of complex graph properties**

by  
**Zi Yuan Chen**

# Contents

<b>List of Figures</b>	<b>iii</b>
<b>List of Tables</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Overview . . . . .	1
1.2 History . . . . .	1
1.3 Basic notation and terminology . . . . .	4
<b>2 Graph Properties</b>	<b>8</b>
2.1 Trophic Coherence . . . . .	8
2.1.1 Inclusion of weights and equation improvement . . . . .	9
2.2 Clustering Coefficient . . . . .	10
2.2.1 Clustering for simple graphs . . . . .	11
2.2.2 Clustering for weighted graphs . . . . .	12
2.3 Centrality . . . . .	16
2.3.1 Betweenness centrality . . . . .	16
2.3.2 Generlisation to directed graphs . . . . .	18
2.3.3 Other centrality values . . . . .	20
2.4 Webpages . . . . .	20
2.4.1 Hyperlink-Induced Topic Search . . . . .	21
2.4.2 Page Rank . . . . .	22
<b>3 Application of Graph Properties</b>	<b>26</b>
3.1 Early Experimentations . . . . .	26
3.2 Zipf's Law . . . . .	31
<b>4 Analysing Languages</b>	<b>33</b>
4.1 Linguistics . . . . .	33
4.2 Text Corpus . . . . .	34
4.3 Indo-European Language Family . . . . .	34
4.3.1 English . . . . .	35
4.3.2 German . . . . .	40
4.3.3 French . . . . .	44

4.4	Japonic . . . . .	49
4.4.1	Japanese . . . . .	49
4.5	Sino-Tibetan . . . . .	54
4.5.1	Chinese . . . . .	54
<b>5</b>	<b>Conclusion</b>	<b>59</b>
5.1	Final Correlations . . . . .	59
5.1.1	Local Clustering Coefficient . . . . .	59
5.1.2	Betweenness Centrality . . . . .	59
5.1.3	Closeness Centrality . . . . .	59
5.1.4	Page Rank . . . . .	59
5.1.5	Trophic Level . . . . .	59
5.2	Further works and Applications . . . . .	59
<b>Bibliography</b>		<b>60</b>
<b>Appendices</b>		<b>65</b>
<b>Appendix A Karate Club</b>		<b>65</b>
A.1	Karate Club Adjacency Matrix . . . . .	65
<b>Appendix B Langauages</b>		<b>67</b>
B.1	Text Corpus . . . . .	67
B.2	English Language Table . . . . .	67

# List of Figures

1.1	The original seven bridges of Königsberg and it's graph representation that only focusses on the key details and disregards the irrelevant information. Achieved by the use of vertices and edges. . . . .	2
1.2	The solution for the four colour map problem with regards to the counties in the UK sourced from Robin Wilson [58]. . . . .	3
1.3	The 2 simple types of graphs within graph theory that are the basic building blocks for more complex structures and theorems. . . . .	5
1.4	A simple network flow with a source node and a sink node, the current flow is 0 and the edge capacities are shown. The image is sourced from Wikipedia [38]. . . . .	6
1.5	A simple graph along with it's adjacency matrix. . . . .	7
2.2	A simple graph of six vertices and seven edges. Vertices $v_1$ and $v_3$ have two geodesics meaning that the central vertices are $v_4$ and $v_5$ who share partial control but $v_0$ has full control and is most central between $v_1$ and $v_3$ . Note that $v_3$ and $v_6$ cannot be central as they only have one edge each and they only have one geodesics because if $v_5$ or $v_1$ is included then it is no longer the shortest path between $v_3$ and $v_6$ . . . . .	17
2.3	a . . . . .	19
2.4	A graph representation of webpages and their in and out links. Algorithms such as Page Rank and HITS uses this graph format to help demonstrate their calculations. This image was sourced from the paper by Devi, Gupta and Dixit[16] whom describes and compares Page Rank algorithm to the HITS algorithm. . . . .	23
3.1	The initial graph based on the karate club dataset pre-existing within the library that was generated through the python program. Contains 34 clubs and their connections to one another with no important meaning with the positions of the vertices. . . . .	27
3.2	Centrality values as the x-axis . . . . .	28

3.3	Graph generated similarly to the betweenness graph for the karate dataset but the vertices are plotted with local clustering coefficient as the x-axis and trophic levels as the y-axis. . . . .	29
3.4	The plot of Zipf's law containing 30 different language corpus generated from the first 10 million words in each language from Wikipedias. The image was sourced from Wikipedia[30]. . . . .	32
4.1	Graph created from the corpus labelled with words or in integers. . . . .	35
4.2	Betweenness and closeness centrality values displayed on the x-axis in graphical form. . . . .	38
4.3	Similarly as before with betweenness and closeness but with local clustering and page rank instead. . . . .	39
4.4	The German word graph and numbered equivalent of the word graph generated from the German translation of the "Sleeping Beauty" corpus. . . . .	41
4.5	Betweenness and closeness centrality values displayed on the x-axis based on the German numbered word graph. . . . .	43
4.6	Displays the local clustering and page rank on the x-axis instead of the centrality values. . . . .	44
4.7	The French word graph and numbered equivalent of the word graph generated from the French translation of the "Sleeping Beauty" corpus. . . . .	45
4.8	Betweenness and closeness centrality values displayed on the x-axis based on the French numbered word graph. . . . .	47
4.9	Displays the local clustering and page rank on the x-axis instead of the centrality values. . . . .	48
4.10	The Japanese word graph and numbered equivalent of the word graph generated from the Japanese translation of the "Sleeping Beauty" corpus. . . . .	50
4.11	Betweenness and closeness centrality values displayed on the x-axis based on the Japanese numbered word graph. . . . .	52
4.12	Displays the local clustering and page rank on the x-axis instead of the centrality values. . . . .	53
4.13	The Chinese word graph and numbered equivalent of the word graph generated from the Chinese translation of the "Sleeping Beauty" corpus. . . . .	55
4.14	Betweenness and closeness centrality values displayed on the x-axis based on the Chinese numbered word graph. . . . .	57
4.15	Displays the local clustering and page rank on the x-axis instead of the centrality values. . . . .	58

# List of Tables

3.1	Table containing all the values calculated for each vertex of the graph.. The order is in the club number depicted in the first column. . . . .	30
4.1	Partial extracts of the table data for graphical properties of the English Story Corpus. . . . .	36
4.2	Partial extracts of the table data for graphical properties of the German Story Corpus. . . . .	42
4.3	Partial extracts of the table data for graphical properties of the French Story Corpus. . . . .	46
4.4	Partial extracts of the table data for graphical properties of the Japanese Story Corpus. . . . .	51
4.5	Partial extracts of the table data for graphical properties of the Chinese Story Corpus. . . . .	56

# Chapter 1

## Introduction

### 1.1 Overview

The application of mathematics have played a key role in the development of technological advancements and breakthroughs in sciences. Throughout history, mathematics has provided us with an increasing amount of various sub branches such as discrete maths, applied maths, cartesian geometry, algebra, calculus and many more. All of which can be applied to the real world, whether it is for construction, physics or day to day life.

A key branch under discrete mathematics is graph theory where models can be developed that represents relationships between different objects. Graphs has a range of uses both in the mathematical world and the real world. They can be used to visually represent large sets of numerical data so that different properties can be derived from the graph such as clustering of certain areas, the connectivity between vertices or edges and the correlations that they may have. They are also widely known as networks with some examples such as a friendship network, business networks and even a food chain levels. These are the areas that I will explore along with a way to visualise them by using python. Analysing specific properties that may influence the visualisation of the graphs and thus the outcome of the relationships between the vertices.

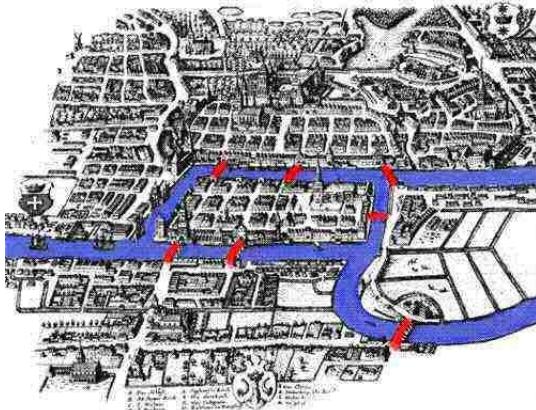
### 1.2 History

Initially, graph theory was introduced in 1735 as a form of solution to the seven bridges of Königsberg problem which solved by Leonard Euler[46]. This famous problem involved an island within Königsberg that had a river, Pregel, surrounding the island via a fork, there were seven bridges that crossed this river from the island to other major landmasses of Königsberg, Prussia. The island had four bridges, two north, two south onto the mainland, one bridge to a neighbouring island and the neighbouring island itself had two bridges which totals to the seven bridges stated

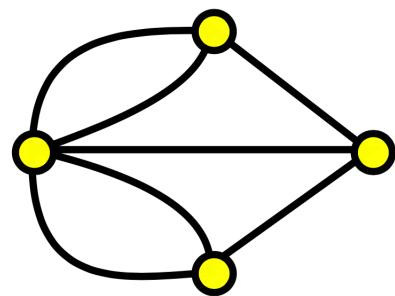
by the problem, see Figure 1.1a.

Due to the location in which the island was situated, the problem was to determine whether a route exists that manoeuvres through all the bridges exactly once and must return to the starting location. Leonard Euler proceeded to analyse the problem by evaluating the only the key areas, this was the land masses and the bridges. Other information such as the sizes of the island, bridge type or length were irrelevant. Consequently, the problem could be portrayed via dots and lines to give a simplistic view (Figure 1.1b), once developed, this was known as vertices which denoted the key interest and the edges that are incident to the vertices to denote the connections/relations between them.

By removing the irrelevant information, Euler was able to simplify and visualise the problem thus discovering the fact that for there to be a solution, each vertex must have an even number of edges incident to the vertex (even degree) as you require one edge for entering and another for exiting otherwise not all edges will be part of the final path. All vertices in the Königsberg problem have odd degree, thus what's known as an *Eulerian path* (a path that traverses all edges exactly once) doesn't exist for this problem and since a Eulerian path doesn't exist then a *Eulerian circuit* cannot exist either (a Eulerian path that returns to the starting vertex). Therefore Euler's solution proves that there is no solution to the seven bridges of Königsberg and is regarded as the first proof in relation to graphs. This led to the birth of graph theory.



(a) Königsberg and the seven bridges, from MacTutor Archive [51]



(b) Graph representation to Königsberg problem from online source [54]

Figure 1.1: The original seven bridges of Königsberg and it's graph representation that only focusses on the key details and disregards the irrelevant information. Achieved by the use of vertices and edges.

After Eulerian paths and cycles were introduced, the next famous puzzle in relation to graph theory was invented in 1857 and was known as the Icosian Game[10] by William Rowan Hamilton. The objective of the puzzle was to find a cycle that

visits all vertices exactly once and returns to the starting vertex. This type of cycle is later called a *Hamilton cycle* along with the definition of a *Hamilton path* which doesn't have the requirement to return to the starting vertex.

In the mathematical world, Leonard Euler is known for the *Euler's identity* within complex mathematics which states that for a real number  $x$ ,  $e^{ix} = \cos(x) + i \sin(x)$ . This has been crucial in many subject areas such as in physics and engineering. Additionally, in 1850, Euler uncovered another formula to be known as *Euler's polyhedra formula* which states that  $F + V - E = 2$  where  $F$  denotes the number of faces,  $V$  as the number of vertices and the number of edges as  $E$  of this graph model. As polyhedrons can be depicted as graphs, algebraic topology benefited from Euler's polyhedron formula where more complex surfaces could be studied such as the surface of a torus. Based upon this formula, the *Euler characteristic* was formalised to describe the topological characteristic of various complex surfaces with it's formula as  $F + V - E = 2 - 2g$  where  $g$  is denotes the number of "holes" the surface has (formally known as the genus).

Furthermore, graph theory has assisted in problems such as the four-colour map problem in the 1850s that states whether all the countries can be coloured with only four colours on a map such that no adjacent countries were coloured with the same colour. In which the solution wasn't found until 1972 by Kenneth Appel and Wolfgang Haken[42] through the assistance of a computer. An example of a four colour problem is colouring the counties in the UK with only four colours and the solution for this is shown in Figure 1.2.

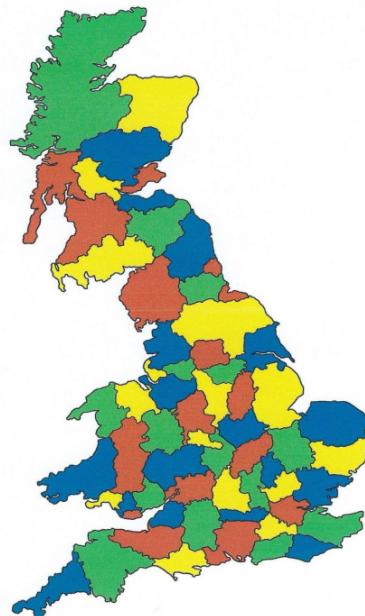


Figure 1.2: The solution for the four colour map problem with regards to the counties in the UK sourced from Robin Wilson [58].

The Chinese postman problem is another such graph theory problem where you must find the shortest path that uses all the edges in the graph at least once. A similar alternate version is called the travelling salesman problem where you find a shortest path that uses all edges exactly once in the graph and must end at the starting vertex. These such problems are used in Linear programming to find optimal solution in routing or pathing between locations. Variants of the CPP (Chinese postman problem) includes undirected Chinese postman problem (UCPP) which contains different restrictions depending on the subset of edges (see the paper [29]) and Chinese Postman Problem with load-dependent costs (CPP-LC[14]) that includes weights upon the edges so that an optimal route can be generated. This can be widely used and applied to problems such as the best route to lower vehicle CO<sub>2</sub> emissions. So these problems can be applied to the modern day scenarios and are still valuable to companies now.

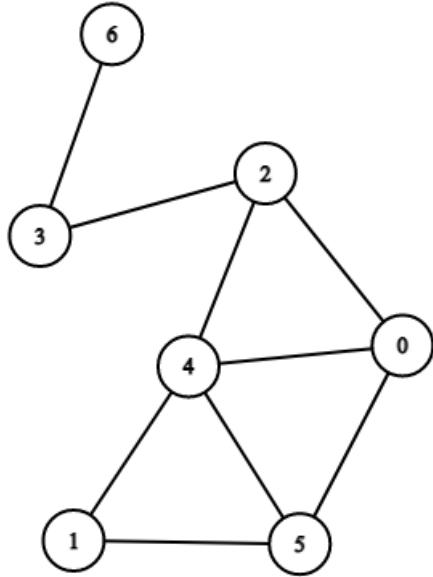
Therefore graph theory studies have been researched extensively since 1735 with many beneficial factors brought into the real world. This is possible due to the versatile nature that data structures in the real world can be represented as a mathematical structure as graphs/networks. Examples of what they can represent ranges from simple relationships between people to the complex structure of the brain by studying it's anatomic structure and assigning vertices according to the sections of the brains and edges as the links between them. The links are typically representations of the neurons in the brain. Further detail of brain mapping into graphs can be read by the paper [25]. Therefore by using graph models, various patterns and correlations can then be derived to generate graph properties that can be studied to develop useful information and possible improvements to the whole graph.

### 1.3 Basic notation and terminology

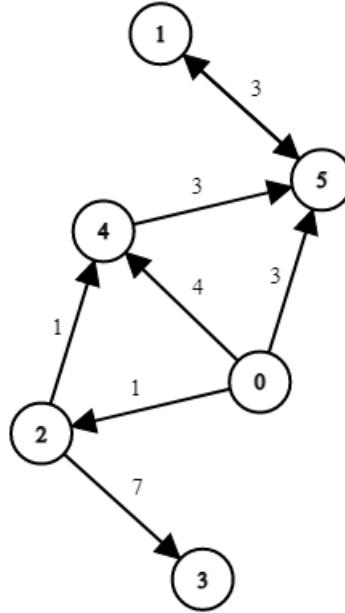
Graphs or networks are mathematical constructs that are formed by a collection of vertices and edges. Vertices  $V$  represents individual objects such as land masses, companies, houses, people etc. Edges  $E$  represents the connections between the vertices such as their relationship, flow of water, supply chain etc. Sets  $V(G)$  or  $V$  and  $E(G)$  or  $E$  forms the graph  $G$  and can be written as  $G = (V, E)$  with  $E$  being a subset of  $V \times V$  so an edge  $e \in E$  can be written as  $v_1v_2$  if  $e$  connects  $v_1$  and  $v_2$  where  $v_1, v_2 \in V$ . There's variation among the graphs as they can be directed or undirected, the edges may carry weights and they may contain self loops. Figure 1.3a and 1.3b shows simple graphs, one of which is undirected and another which is directed. A graph  $H = (V', E')$  is a *subgraph* of  $G = (V, E)$  if  $V' \subset V$  and  $E' \subset E$ .

*Order* is the mathematical term that represents the number of the vertices in the vertex set  $V(G)$ . *Size* is number of vertices in the edge set  $E(G)$ . The *degree* of a vertex, denoted by  $\deg(v)$  is the number of edges that are connected (otherwise known as *incident*) to the vertex, discussed previously in Euler's solution

to the seven bridges of Königsberg problem. Additionally,  $\delta(G)$  and  $\Delta(G)$  represents the minimum and maximum degree in  $G$  respectively.  $G$  is a *regular* graph if  $\delta(G) = \Delta(G)$ . Vertices  $v_1$  and  $v_2$  are *adjacent* if there exists an edge  $e \in E$  that connects them. Vertex  $v_2$  is also known as a *neighbour* to  $v_1$  and is part of the set  $N(v_1)$  which denotes the neighbours of  $v_1$ .



(a) An undirected graph with 7 vertices, 9 edges and a average vertex degree of  $18/7$



(b) A directed graph with 7 vertices and 7 directed edges that contain weights

Figure 1.3: The 2 simple types of graphs within graph theory that are the basic building blocks for more complex structures and theorems.

There are multiple ways to traverse a graph. A *walk*  $w = v_1v_2v_3v_4v_5...v_n$  is a sequence of vertices such that  $E(w) = (v_1v_2, \dots, v_{n-1}v_n)$  where vertices and edges can be repeated. They can be *open* or *closed* depending if the final vertex is equal to the starting vertex. A *trail* is an open walk where no edges are repeated but vertices may be repeated. When all the edges are traversed exactly once, it's known as a *Eulerian trail* (mentioned previously as a Eulerian path) and the graph is called *semi-eulerian* or *traversable*. Similarly for a trail that's closed (returns to the start), then it is known as a *circuit* and if all edges are used then it's called a *Eulerian circuit* (or Eulerian cycle) and the graph is defined to be *Eulerian*. A *path* is a trail but with no vertex repetitions and if the size of the path is equal to the size of the graph then it's a *Hamilton path*. Finally, a *cycle* is a path that ends at the starting vertex and if all vertices are visited exactly once then it's known as a *Hamilton cycle* meaning that the graph is *Hamiltonian*.

When considering network with flows, vertices can be known as nodes and may have capacities that limit the overall flow through the network shown in Figure 1.4. These networks are especially used when considering plumbing, water pipes and even evacuation routes in a building. Within a room there's a capacity that is represented by it's node capacity and the weights of the edges can demonstrate the rate of flow along with it's maximal flow, i.e. when people are evacuating a building, the corridors have a limit to the amount of people that may pass through. Networks can be used to model social network processes through the study of small corporate groups to generate a communications network. This network representation will have a flow of sentiments based on social network theory[60] that is constrained in three ways. Firstly by any existing direct relationships within the group, that will be denoted by vertices. Secondly, the frequency of communication of the relationships defined in the first point. Lastly, the breadth of the existing relationships in the network. Thus, by using graphs and networks, social behaviours with groups or companies can be studied on giving ways to more phycological information through numerical data.

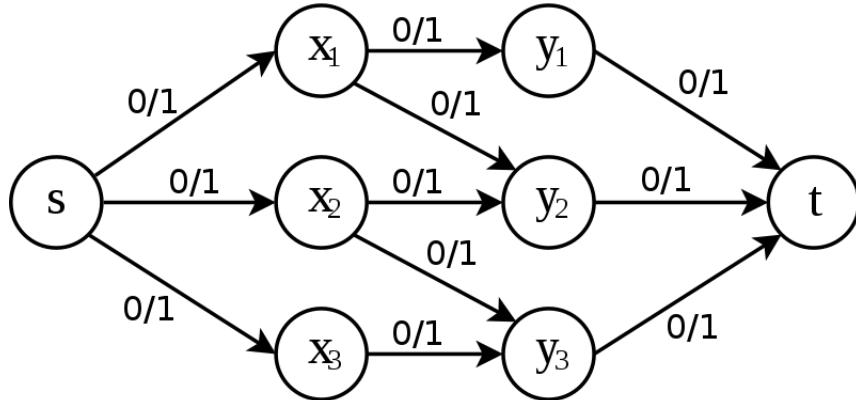


Figure 1.4: A simple network flow with a source node and a sink node, the current flow is 0 and the edge capacities are shown. The image is sourced from Wikipedia [38].

Additionally they can be represented by the use of matrices to enable the use of matrix calculations on the data sets. They're known as *adjacency matrices*[36] and within this matrix holds the number of edges incident to each vertex and it's connection based on the location of this value as the rows and columns represents the vertices. A *weighted adjacency matrix* will instead hold the weights of each edge in the matrix. An  $n \times n$  adjacency matrix  $A = (a_{ij})$  for  $i, j = 1, \dots, n$  is defined by  $a_{ij} = 1$  if there exists an edge from vertex  $i$  to vertex  $j$ . A matrix is always symmetric when considering an undirected simple graph as an edge will contribute to both sides of the matrix. Examples of an adjacency matrix along with it's graph representation is shown in Figure 1.5.

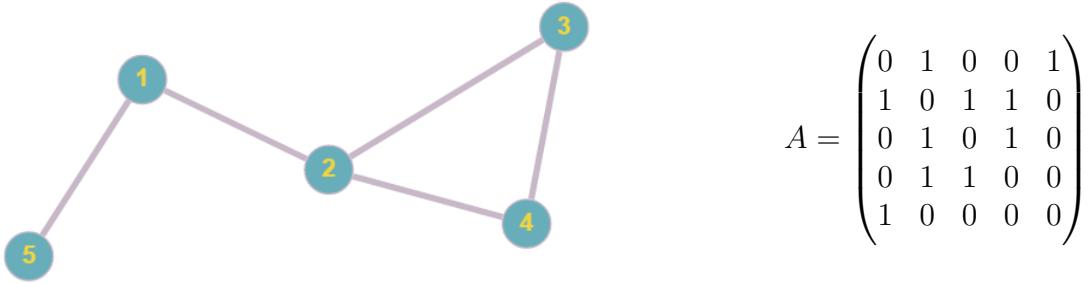


Figure 1.5: A simple graph along with it's adjacency matrix.

The main focus for future chapters will be on weighted directed graphs and the data that can be perceived by them. Chapter 2 will describe and outline specific graph properties that will be applied to datasets. This means that by analysing the graphs, numerical values can be generated to represent certain factors of the graph. These factors can then be applied to the graph to rearrange it's shape so that further correlation can be identified between all the vertices and edges.

# Chapter 2

## Graph Properties

There exists many different graph properties that can be applied according to specific graph types. Hence we will focus on properties that will be applied to connected weighted graphs to provide a collection of views targeted for the same layout structure and to enable a wider selection of properties. Hence, unless stated otherwise, the properties described in this chapter will be of the form of connected weighted graphs.

### 2.1 Trophic Coherence

Firstly the graph property of *trophic coherence*[32], this property determines the stability of the graph and provides trophic levels to the graph and it's individual vertices. Trophic levels are taken from ecology and applied to graphs to generate a height based format for the vertices of a graph. Thus, for a graph  $G = (V, E)$  (can be directed), we have  $V$  which is the set of all vertices in  $G$  and  $E$  which is the set of all the edges within  $G$ . The graph can be represented with an adjacency matrix  $A$  where  $a_{ij}$  denotes the elements within the matrix. The standard trophic level definition on vertex level uses the in degree and the out degree of the vertex  $v_i$  given by Equation 2.1.

$$k_i^{\text{in}} = \sum_j a_{ij}, \quad k_i^{\text{out}} = \sum_j a_{ji} \quad (2.1)$$

So, the standard trophic levels for vertices  $v_i \in V$  is formulated as:

$$t_i = 1 + \frac{1}{k_i^{\text{in}}} \sum_j a_{ij} t_j \quad (2.2)$$

By ecological convention,  $t_i = 1$ , if the vertex  $v_i$  is *basal*. A vertex is known to be basal if it has no edged directed to it, i.e.  $k_i^{\text{in}} = 0$ . The trophic level equation can also simply be written in matrix form by 2.3 with  $\mathbf{z}$  be defined as  $z_i = \max(k_i^{\text{in}}, 1)$  and  $\Delta = \text{diag}(\mathbf{z}) - A$ .

$$\Delta \mathbf{t} = \mathbf{z} \quad (2.3)$$

All vertices will receive a trophic level if and only if the laplacian  $\Delta$  is invertible as the row sum of elements in  $\Delta$  for a vertex that isn't basal is zero. If there are no basal vertices in the graph then  $\Delta$  will be equal to zero and thus be singular (i.e. no inverse). This is discussed by Samuel Johnson [31] in the investigation of stability and such dynamical features within graphs. So for the standard trophic levels equation to work for the entire graph, there must exist at least one basal vertex meaning that this is a limitation to the definition of trophic levels.

### 2.1.1 Inclusion of weights and equation improvement

In a weighted graph, edges carry a weight between a vertex  $u$  to a vertex  $v$ . Weighted matrix  $W$  are used as a representation of the entire graph along with incorporating the direction of each edge and if there exists self loops. The elements of the weighted matrix are  $w_{uv}$ . If the graph isn't weighted then the edge is valued as 1 if the edge exists and 0 otherwise, i.e. the adjacency matrix of  $G$ . The total weight (also known as the strength) for each vertex is defined by the weights into a vertex  $u$  and the weights out of vertex  $u$ , which is shown by  $s_v$  in Equation 2.4. This is essentially the same as the in and out degrees of Equation 2.1 mentioned previously but instead of ones and zeros, the weight values are taken into account. The imbalance of the vertex  $u$  is defined by the weights in of a vertex minus the weights out of the vertex shown by  $i_v$  in Equation 2.5.

$$s_v = \sum(w_{uv}) + \sum(w_{vu}) \quad (2.4)$$

$$i_v = \sum(w_{uv}) - \sum(w_{vu}) \quad (2.5)$$

Vectors  $\mathbf{s}$  and  $\mathbf{i}$  holds all the values of the strength and imbalance for all vertices respectively. We let  $\mathbf{h}$  be a vector, then the graph Laplacian operator in matrix form is defined by Equation 2.6.

$$\Delta = \text{diag}(u) - W - W^T \quad (2.6)$$

Therefore to get the trophic levels for each vertex with consideration for the additional weights, we solve the system of equations for vector  $\mathbf{h}$  as shown by Equation 2.7.

$$\Delta \mathbf{h} = \mathbf{v} \quad (2.7)$$

The values within this vector  $h$  corresponds to the trophic levels for the relative vertices. The values are used to illustrate the various trophic levels within a graph to give a hierarchical format visualisation. Trophic levels aren't unique solutions due to the fact that an arbitrary constant can be added to each component of the graph

if there are multiple components to generate new levels that would be correct. The benefit that this is that Equation 2.7 can use an arbitrary constant  $c$  for a vertex in  $\mathbf{h}$ . If there are multiple components to the graph then a vertex in each component. A unique solution can be found this way in which the trophic level values can be shifted so that a better graphical display can be generated. For example, have the lowest level be zero.

Trophic level values can also be used to equate the overall *trophic incoherence* of the graph as a whole rather than just looking at the graph on a vertex level. By using the trophic levels from  $\mathbf{h}$ , the equation for the trophic incoherence is defined by R. S. MacKay, S. Johnson, B. Sansom[31] as:

$$F_0 = \frac{\sum_{uv} w_{uv}(h_v - h_u - 1)^2}{\sum_{uv} w_{uv}} \quad (2.8)$$

The possible shifting of the trophic levels doesn't affect the trophic incoherence hence it is independent. The incoherence is strictly ranged from zero to one. If  $F_0 = 0$  then the graph is *maximally coherent* as it would mean that all levels in the graph has a difference of exactly one which means that the graph is perfectly separated into levels. Whereas if  $F_0 = 1$  then the graph is *maximally incoherent* and levels are harder to decipher. As  $F_0$  measure the incoherence, by taking  $1 - F_0$ , this would instead measure the coherence of the graph as they're each others converses.

In conclusion, trophic levels can be applied to weighted directed graph and any subcategories to achieve a hierarchical view of the graph. This eases the visualisation of many datasets and is used to decipher valuable information that may be of use. Through the combination of other graph properties which are described in this chapter, various combinations of these properties will yield different visualisations.

## 2.2 Clustering Coefficient

The *transitivity* [50] of a graph, also known as *clustering*, is a property of a graph that measured the density of triangles within the graph, where three vertices are connected together. Used to quantify the graph's connectivity strength as it determines the fraction of triangles over the possible triangles that could be formed within the graph. Another perspective is that the coefficient quantifies the probability of a vertex  $a$  having an edge to vertex  $c$  if  $ab, bc \in E(G)$ . Thus, the *clustering coefficient* determines how complete the graph is with a value of 1 meaning it is complete. There are two popular introductions of clustering coefficient, the *local clustering* and the *global clustering*. The global clustering coefficient essentially measures the completeness of the graph by measuring the number of existing triangles divided by the number of possible triangles in the graph. The local clustering coefficient measures the clustering coefficient for each vertex rather than the whole graph, the measurement is taken by the number of triangles that has a connection to the vertex

over the number of triples centred on this vertex. In other words, the local value demonstrates how close the neighbours of this vertex is to being a complete graph (a *clique*).

### 2.2.1 Clustering for simple graphs

To determine the values of the clustering coefficient for simple connected graphs that are unweighted and undirected, the global clustering coefficient is defined by equation 2.9 where  $\sum T$  denotes the number of triangles (closed triplets) and  $\sum \tau$  denotes the number of connected triplets in the graph.

$$C = \frac{3(\text{Number of total triangles})}{\text{Number of total connected triples}} = \frac{\sum T}{\sum \tau} \quad (2.9)$$

An alternative equation which was demonstrated by M.E.J. Newman [41] through the studies of complex networks in terms of social networks where the clustering coefficient determines the likelihood that a friend of your friend is also your friend. So, the alternative equation is written in the form of equation 2.10 where  $\sum P_2$  denotes the number of paths with length two within the graph.

$$C = \frac{6(\text{Number of total triangles})}{\text{Number of paths with length 2}} = \frac{\sum T}{\sum P_2} \quad (2.10)$$

By considering the vertices of the graph, the local clustering coefficient can be defined to give such a value to each vertex  $v \in V(G)$  and is given by the equation 2.11 where  $i$  is the index of the vertex. This definition is from [41] and proposed by Watts and Strogatz [56] where they analysed small world networks in relation to various real world systems by the use of clustering coefficients and random graphs to formulate certain similarities. Note that if the degree of a vertex is one then the coefficient can be determined as 0, otherwise the equation will lead to 0/0.

$$C_i = \frac{\text{Number of triangles connected to } i}{\text{Number of triples centred on vertex } i} \quad (2.11)$$

Another representation of the global clustering coefficient is to take the averages of all the local coefficients[47]. When the vertices have a degree of 0 or 1 then  $C_i = 0$  so clustering coefficient is defined by equation 2.12.

$$C = \frac{1}{n} \sum_i C_i \quad (2.12)$$

Later the clustering coefficients will be used as assistance to model graphs generated off of a dataset. However to provide more accurate values, then weights and directions would need to be taken into account. Thus, the definitions of clustering coefficients must be developed further.

### 2.2.2 Clustering for weighted graphs

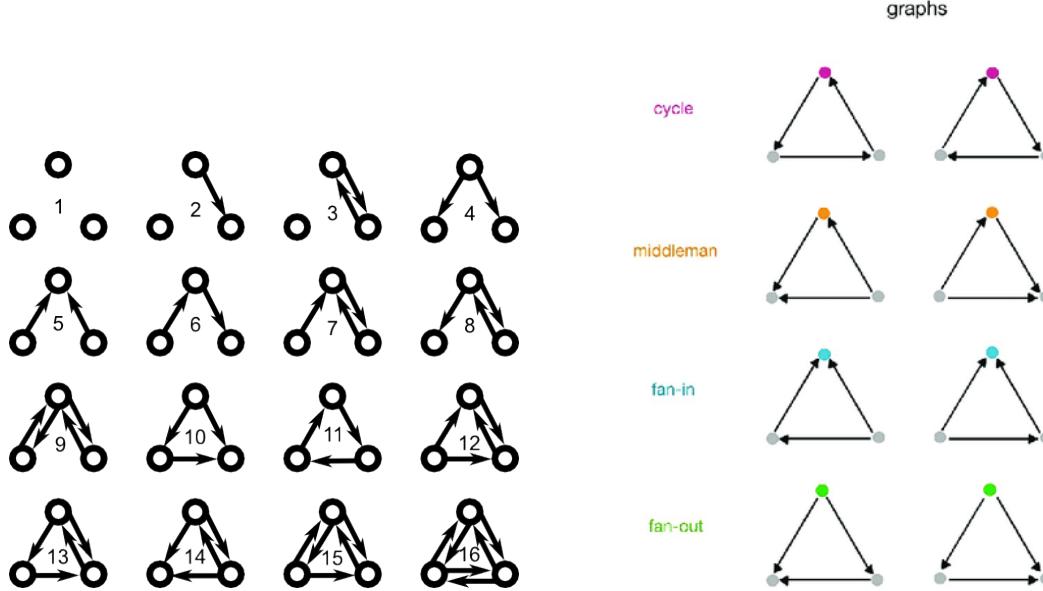
Now by considering graphs as before but weighted, the equations undergo changes. For the instance of weighted graphs, there are multiple different definitions of clustering coefficients, each with slight variation in values depending on the type of graph. This section will summarise a couple of the different definitions for weighted and directed graphs and further detail can be analysed from Tanguy and Anna Levina on weighted directed clustering [20]. In this paper, four different definitions are reviewed which are the Barrat definition, Onnela's definition, Zhang & Horvath and their own continuous definition for weighted graphs. Zhang & Horvath[61] have used their definition of weighted clustering coefficients to analyse gene co-expression networks to review their functionality. Additionally, by soft or hard thresholding, it enables them to determine relationships between the clustering coefficient and gene networks within biology.

A simple idea to associate the clustering coefficient with regards to the edge weights is to define a value  $w$  that represents the value of the triplet.  $w$  can be the summation of the triplet, the mean of the triplet or another suitable method depending on the purpose. Then equation 2.13 calculates the weighted clustering coefficient[44] where  $T$  denote the triangles in the graph and  $\tau$ , the triples.

$$C = \frac{\text{Total of closed } w}{\text{Total of } w} = \frac{\sum_T w}{\sum_\tau w} \quad (2.13)$$

Weights can be added trivially through this way, in which it won't affect the equations formalised beforehand. So now considering the addition of directions as well as weights which causes further complexities in the values of clustering coefficients due to the various number of different motifs used to describe the nature of the triangles. For instance, there are sixteen possible motifs for directed graphs of three vertices shown in Figure 2.1a. However if we consider only connected triangles, they can be organised into four types of motif groups known as a cycles, Middleman, Fan-in and Fan-out as demonstrated in Figure 2.1b which are used in the study of higher order motifs and synaptic integration by Bojanek, Zhu and Maclean[4]. An interesting result used in this paper is the isomorphisms between the middleman, fan-in and fan-out motifs.

Consideration of the edge's directions yields better accuracy in the coefficient values. One of the versions mentioned in the paper [20] was Fagios where he introduces the clustering coefficient to binary directed networks which are equivalent to simple directed connected graphs. Firstly the equation for the directed version without the consideration of weights is defined by the ratio of all directed triangles centred on a vertex  $i \in V(G)$  and the number of all possible triangles that could be formed with vertex  $i$ . These are called  $t_i$  and  $T_i$  respectively. Before this equation, prior properties of the graph are necessary so that the equation can be easily formulated. Thus, consider a graph  $G = (V, E)$  with its matrix representation as the adjacency matrix  $A$  along with  $V_1$  as the column vector, dimension  $n$  of the



(a) All sixteen motifs of possible connections between three nodes and their directions shown on the edges. Ranges from zero edges to the maximum of six edges. Image sourced from Math Insight[28].

(b) The four named categories of motifs in which all directed connected triangles fall under. Each pair of motifs is listed as the corresponding types which are cycles, middleman, fan-in and fan-out. Image displaced from [43].

graph, of only 1s.  $A_i$  is the  $i$ -th row of the adjacency matrix. The in-degrees and out-degrees of a graph are the total number of edges going in or out of a vertex  $i \in V(G)$  respectively, the total degree is the sum of the in and out degrees shown by equations 2.14.

$$\begin{aligned} \text{in}(d_i) &= \sum_{i \neq j} a_{ji} = (A^T)_i V_1 & (2.14) \\ \text{out}(d_i) &= \sum_{i \neq j} a_{ij} = (A)_i V_1 \\ \text{tot}(d_i) &= \text{in}(d_i) + \text{out}(d_i) \sum_{i \neq j} a_{ji} + \sum_{i \neq j} a_{ij} = (A^T + A)_i V_1 \end{aligned}$$

When the edge is directed both ways, this degree is the summation of the products of all the edges of vertex  $v$  that are bidirectional. Formally can be shown as equation 2.15 with  $A_{ii}$  as the  $i$ -th element of the diagonal for the matrix product of  $A$ .

$$\text{bi}(d_i) = \sum_{i \neq j} a_{ij} a_{ji} = A_{ii}^2$$

This equation is demonstrated by Fagio [19] that measures the clustering coefficient for each vertex with directed edges with the consideration of the eight triangles that this vertex could form shown previously in figure 2.1b.

So this equation can be demonstrated with vertex  $i$  and pairs of neighbours  $j$  and  $k$  that essentially shows that the equation 2.15 calculates the triangles formed by  $v$  over the possible triangles with the deduction of  $2\text{bi}(d_i)$  as otherwise if vertex  $i$  and  $j$  had edges directed to each other, this causes a count of two additional triangles.

$$\begin{aligned} C_i &= \frac{\frac{1}{2} \sum_j \sum_k (a_{ij} + a_{ji})(a_{ik} + a_{ki})(a_{jk} + a_{kj})}{\text{tot}(d_i)(\text{tot}(d_i) - 1) - 2\text{bi}(d_i)} \\ &= \frac{(A + A^T)_{ii}^3}{2(\text{tot}(d_i)(\text{tot}(d_i) - 1) - 2\text{bi}(d_i))} = \frac{t_i}{T_i} \end{aligned} \quad (2.15)$$

Weights of the edges can be implemented into Fagio's equation of clustering coefficients by simply using a weighted adjacency matrix  $W$  instead of the adjacency matrix  $A$  for graph  $G$ . Mentioned before, the weights of the triangles can be considered differently however in Fagio's generalisation of the clustering coefficient equation, the mean weights of the triangles are utilised. This is shown by taking a cube root all elements of the weighted matrix  $W$  which can be denoted as  $W^{[1/3]}$ . Therefore by subbing  $W^{[1/3]}$  in the place of  $A$  in equation 2.15, the weighted directed version can be achieved, formally as equation 2.16.

$$C_i = \frac{(W^{[1/3]} + (W^{[1/3]})^T)_{ii}^3}{2(\text{tot}(d_i)(\text{tot}(d_i) - 1) - 2\text{bi}(d_i))} = \frac{t_i}{T_i} \quad (2.16)$$

However with this generalisation of the formula, it considers any triangle formed by a vertex  $i$  as equal values. This means that the directions in the triangles are meaningless however due to the nature of directed graphs, their directions are what gives the graph its flow of information and different directions will lead to different interpretations of the graph. Thus an improvement to properly include the directions of edges is to consider each motif separately or in Fagio's case, treat them by considering the 4 types of categories the motifs can fall under mentioned before in Figure 2.1b.

These were cycles, middleman, fan-in and fan-out and by measuring each specific category then we can get more accurate coefficients for the relative pattern. The definition of number of all possible triangles was defined in equations 2.15 and 2.16 as  $T_i$ . Similarly with the directed triangles that are actually formed by a vertex  $i$  by  $t_i$ . These can be both decomposed into the 4 types of motifs which can then be used to create the clustering coefficient for each specific motif. Note that the sum of the clustering coefficient for specific motifs will be the general clustering coefficients for all triangles. The equations are decomposed as follows:

$$\begin{aligned}
T_i &= \text{tot}(d_i)(\text{tot}(d_i) - 1) - 2\text{bi}(d_i) \\
&= \text{in}(d_i)\text{out}(d_i) - \text{bi}(d_i) + \text{in}(d_i)\text{out}(d_i) - \text{bi}(d_i) \\
&\quad + \text{in}(d_i)(\text{in}(d_i) - 1) + \text{out}(d_i)(\text{out}(d_i) - 1) \\
&= \text{cyc}(T_1) + \text{mid}(T_2) + \text{fan-in}(T_3) + \text{fan-out}(T_4)
\end{aligned} \tag{2.17}$$

$$\begin{aligned}
t_i &= (W^{[\frac{1}{3}]} + W^T)_{ii} \\
&= (W^{[\frac{1}{3}]})_{ii}^3 + (W^{[\frac{1}{3}]}W^{[\frac{1}{3}]}{}^TW^{[\frac{1}{3}]}))_{ii} + (W^{[\frac{1}{3}]}{}^TW^{[\frac{1}{3}]}{}^2)_{ii} + (W^{[\frac{1}{3}]}{}^2W^{[\frac{1}{3}]}{}^T)_{ii} \\
&= \text{cyc}(t_1) + \text{mid}(t_2) + \text{fan-in}(t_3) + \text{fan-out}(t_4)
\end{aligned} \tag{2.18}$$

So both equations can be split according to their different motifs through logical and algebraic reasoning. For  $T_i$ , the maximum number of directed triangles for cycles, middleman, fan-ins and fan-outs that can be formed equates to the total number of triangles for a vertex  $i$ . For example, when calculating cycles, the maximum number of directed cycles that can be formed with vertex  $i$  is it's degree in multiplied by it's degree out, hence  $\text{in}(d_i)\text{out}(d_i)$ . However if a neighbour of  $i$  has an edge to and from the same vertex then this would account to an additional triangle counted hence the subtraction of the bidirectional edges by  $-\text{bi}(d_i)$ . Notice that middleman and cycles only differ according to the direction of the pairs of neighbours connected to vertex  $i$  that forms the triangle hence the reason why  $\text{cyc}(T_1) = \text{mid}(T_2)$ . Then for the calculations of fan-in motif, it's the multiplications of the in degrees of  $i$  and in degrees - 1 (as an edge is already considered) which gives the maximum fan-in styled triangles. Similarly for the fan-out motif.

Then for the actual triangles formed, they can be broken down by algebra shown by equation 2.19 which equates to the weighted matrix component seen in the original equation. Hence by algebraic manipulation,  $t_i$  can be separated into the 4 motifs shown before.

$$\text{cyc}(t_i) = \frac{1}{2} \sum_j \sum_h w_{ij}w_{jh}w_{hi} + w_{ih}w_{hj}w_{ji} \tag{2.19}$$

$$\begin{aligned}
&= \frac{1}{2}((W^{[\frac{1}{3}]})_{(i)}W^{[\frac{1}{3}]}(W^{[\frac{1}{3}]}{}^{(i)}) + (W^{[\frac{1}{3}]}{}^T)_{(i)}(W^{[\frac{1}{3}]}{}^T(W^{[\frac{1}{3}]})(W^{[\frac{1}{3}]}{}^{(i)}) \\
&= (W^{[\frac{1}{3}]})_{(i)}W^{[\frac{1}{3}]}(W^{[\frac{1}{3}]}{}^{(i)}) = (W^{[\frac{1}{3}]})_{ii}^3
\end{aligned} \tag{2.20}$$

Finally, the clustering can be calculated according to the specific motifs of cycles, middleman, fan-in and fan out. Which can be formally defined as:

$$x(C_i) = \frac{x(t_i)}{x(T_i)} \tag{2.21}$$

where  $x = \{\text{cyc, mid, fan-in, fan-out}\}$  .

By demonstrating one specific formulation of the clustering coefficient applied to weighted directed graphs, we notice that depending on variations of properties in relation to the triangles, different values may occur for the same graph. Such properties include the various different motifs of the triangles and the consideration of the edge weights. This is what gave way to the multiple different versions as each has their benefit depending on what the graph represents which are thoroughly examined in the paper [20], all versions can also be analysed against each other to determine which has the best performance depending on the ideal requirements of the task[12].

## 2.3 Centrality

Positioning of the graph is a vital area in displaying and extracting important information which can then be assigned to the vertices and edges. The benefit of this is that it may be used to identify key areas within a graph or network such as the links between companies and which one has the most influence. This can also be applied to brain networks, the spreading patterns of disease and is useful to characterise specific areas to give a new interpretation of the graph. One of the major centrality values is the betweenness of vertices in a graph which will be described in further detail in the next section.

### 2.3.1 Betweenness centrality

Betweenness centrality measures the centrality of a graph by using the shortest paths of pairs of vertices. The introduction of this idea was through the view of a communications network. Where a point(or vertex) of the communications network is deemed to be central if it lies on the shortest path between another pair of points in the network. Alex Bavelas[2] formalised this idea of centrality where he suggests that a person in a group is in the central position if that person lies on the shortest path between other connecting pairs. With the implication that this person then holds the power or responsibility for the others due to information exchange that must go through that point. A simplistic way of viewing the betweenness value for a vertex is that the larger the value, the greater number of shortest path connections it has to other vertices. In other words, if the value is large for a vertex, then the travel time from this vertex to other vertices is shorter.

The betweenness centrality will be discussed based on Freeman's interpretation of betweenness centrality measure[21]. For a simple unconnected and undirected graph  $G = (V, E)$ , consider all the unordered pairs of vertices  $v_i, v_j \in V$  with  $i \neq j$ . This pair must either be disconnected or has at least one path connecting them with its path length based on the number of edges contained within. The path or paths

connecting  $v_i$  to  $v_j$  with the shortest length is known to be the *geodesics*. If the path is larger than one edge (the vertices were adjacent) then vertices in this geodesic are considered to be central. Depending if the vertex is the only central vertex between the pair of vertices determines whether it has complete or partial control of their link. The intuition of control in betweenness was expressed by Shimbel[52] where work sites are considered as the vertices. Meaning that the vertices with control will be the connecting sites between other sites. Which in terms means that the connecting sites holds responsibility to the other sites and must relay information and resources to them. Figure 2.2 expresses the idea of central vertices and the geodesics between them. Vertices  $v_1$  and  $v_3$  have two geodesics meaning that they share power. Also vertices  $v_3$  and  $v_6$  only has one geodesics through  $v_4$  and  $v_0$  so either  $v_4$  or  $v_0$  can have complete control.

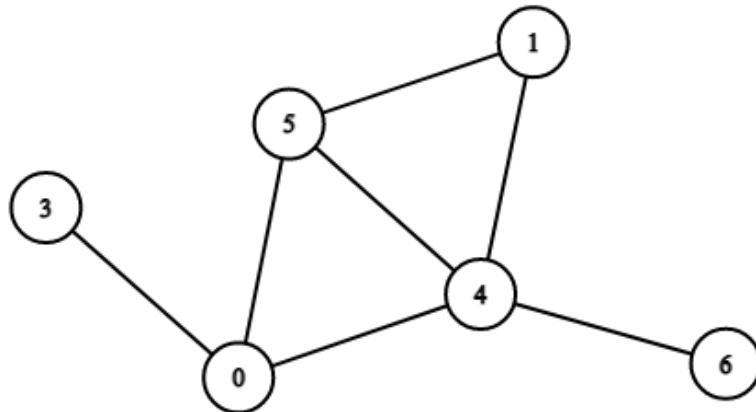


Figure 2.2: A simple graph of six vertices and seven edges. Vertices  $v_1$  and  $v_3$  have two geodesics meaning that the central vertices are  $v_4$  and  $v_5$  who share partial control but  $v_0$  has full control and is most central between  $v_1$  and  $v_3$ . Note that  $v_3$  and  $v_6$  cannot be central as they only have one edge each and they only have one geodesics because if  $v_5$  or  $v_1$  is included then it is no longer the shortest path between  $v_3$  and  $v_6$ .

To use the geodesics for a pair of vertices, if there are more than one geodesics then they are considered to have equal probability in deciding which one to be used. This is just simply given by  $\frac{1}{g_{ij}}$  where  $g_{ij}$  are the number of geodesics between vertices  $v_i$  and  $v_j$ . The formulation of partial betweenness  $b_{ij}(v_k)$  is defined using the idea of geodesics, so for a vertex  $v_k$  in  $G$  and a vertex pair  $v_i$  and  $v_j$ , the partial betweenness for  $v_k$  can be calculated based on the vertex pair and is given by the equation

$$b_{ij}(v_k) = \frac{1}{g_{ij}}(g_{ij}(v_k)) = \frac{g_{ij}(v_k)}{g_{ij}}(i \neq j \neq k) \quad (2.22)$$

where  $g_{ij}(v_k)$  is the amount of geodesics connecting vertices  $v_i$  and  $v_j$  that include

$v_k$  in it's path. This essentially translate to the probability that the geodesic chosen for  $v_i$  and  $v_j$  contains  $v_k$ . Additionally, notice that  $b_{ij}(v_k) = 0$  if there doesn't exist a path between the vertex pair,  $v_i$  and  $v_j$ . This can be extended to calculate the centrality of each vertex. So, for a graph of size  $n$ , the centrality value for a vertex  $v_k \in V$  can be defined by

$$C_B(v_k) = \sum_i^n \sum_j^n b_{ij}(v_k) \quad (2.23)$$

where  $i < j$ . So this defines the betweenness centrality for  $v_k$  and it's value increases depending on the amount of geodesics that  $v_k$  is a part of. The maximum value[22] was proved by Freeman through the use of a star with the central point as  $v_k$  as all vertices are rachable through this central one. Furthermore, this means there are  $n(n - 1)/2$  paths between all the unordered pairs of the star graph  $S$ . With  $n - 1$  of these connected to the central vertex  $v_k$ . So the betweenness centrality for  $S$  is

$$C_B(v_k) = \frac{n(n - 1)}{2} - (n - 1) = \frac{n^2 - 3n + 2}{2} \quad (2.24)$$

And if any new edge is added that doesnt increase the branches of the star, then a new geodesic would form without  $v_k$  meaning that the value will fall. Therefore Equation 2.26 expresses the betweenness centrality of any vertex in a graph  $G$  by it's representation though a ratio with the maximal value.

$$C'_B(v_k) = \frac{2C_B(v_k)}{n^2 - 3n + 2} \quad (2.25)$$

### 2.3.2 Generlisation to directed graphs

The key idea of betweenness centrality is to evaluate the graph and produce values according to the shortest paths of all the possible pairs of the graph. The higher the betweenness value, the more the vertex is likely to lie on the shortest paths of any two vertices. As this concept investigates the vertices and shortest paths, weighted edges wouldn't benefit this graph evaluation as weights could cause a pair of vertices to be seen as further apart than they actually are within the graph. In the experimentations of social networks[23], the centrality values are used on undirected graphs however an idea to generlise the betweenness to weighted graphs is to take the weights as indication of the distance of the vertices. Meaning that the geodesics of any pair will be defined on the smallest total value of paths between them rather than the shortest path length as shown on Figure 2.3.

Hence, we discuss the generalisation in this section only to directed graphs and if the graph has weights, then they are not impleted to ensure betweenness values will not be influenced. The geodesic proportions of paths from  $v_i$  and  $v_j$  was defined earlier in Equation 2.22. Consequently, this eqation can be utilised to define pair-dependency of vertices  $v_i$  to  $v_k$  where the vertext  $v_i$  has to depend on vertex  $v_k$  in

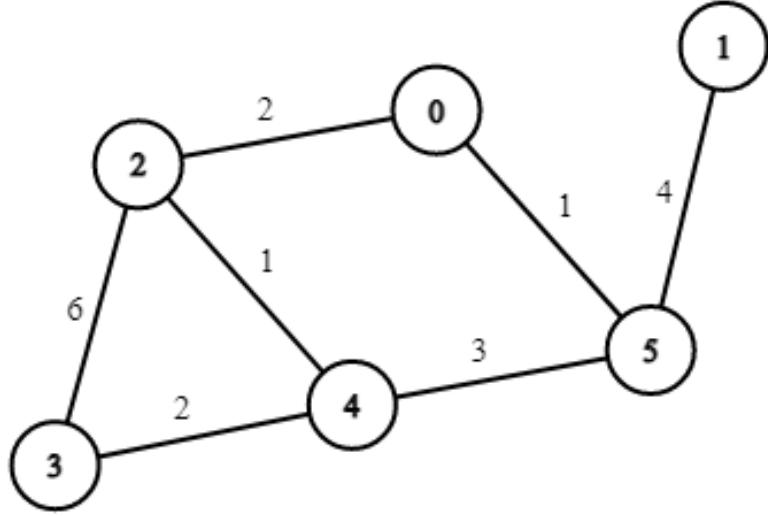


Figure 2.3: a

order to get to other vertices such as  $v_j$  on it's geodesics. In other words,  $v_k$  acts like a gatekeeper to  $v_i$ . Therefore, for a graph with  $n$  vertices, the pair dependency is defined to be

$$d_{ik}^* = \sum_{j=1}^n b_{ij}(v_k) \quad (2.26)$$

where  $i \neq j \neq k$ . Matrices can be used to store the pair-dependency values to provide an ease of use and is a better representation for all the values. The results can be arranged into the matrix  $D$  defined as  $D = (d_{ik}^*)$ . The elements of the matrix measures how much the vertex corresponding to the row number depends on vertex corresponding to the column number to be able to connect to other vertices in the graph. Additionally the betweenness centrality can also be calculated based on this matrix  $D$  through the summation of the columns in which the sum will give the betweenness centrality for the column number that represents the vertex. Otherwise shown as

$$\sum_{i=1}^n d_{ik}^* = 2C_B(v_k) \quad (2.27)$$

Which means that the betweenness centrality of  $v_k$  is double of the pair dependency column sum[57]. This is because for an undirected graph, the upper and lower diagonal matrix are equal due to the symmetry of graph  $G$ . The generalisation for directed graph can then be shown to be

$$C_B(v_k) = \sum_{i=1}^n d_{ik}^* \quad (2.28)$$

### 2.3.3 Other centrality values

Centrality in itself is a larger area within graph theory and other than the betweenness values, there are other similar attributes such as the closeness centrality, eigenvalue centrality, Katz centrality[34] and the Hyperlink-Induced Topic Search (HITS) centrality. All of which calculate values for the vertices or edges given their positions within the graph by various diffent methods. Interestingly the closeness centrality can be seen as the duality of betweenness centrality as they can be obtained from row and column summatiuons of the depency relation defined in the paper by Brandes, Borgatti and Freeman[5]. Betweenness studies the vertices that acts as bridghes between other geodesics whereas the closeness centraliuty measures the average distance of the shortest paths between any pairs of vertices. A vertex with high closeness value means that the distance to any other vertex is short on average. Closeness centrality[6] can simply be calculated as the inverse of total shortest distance from a vertex  $i$  to all of vertices, demonstrated by equation 2.29.

$$C_C(v_i) = \frac{1}{\sum_j^n d(v_i, v)} \quad (2.29)$$

where  $i \neq j$ . The calculation is for simply connected graphs however can be easily expanded to directed and weighted graphs by modying the calculation of the mesaure of distances for the graph in question. I.e. Take the total weights of the path length rather than the path length for weightedness and to consider only correct paths(travelling along an edge in the permited direction) when investigating directed graphs.

Therefore by taking these properties, the graph can be rearranged in accordance to their values. This is accomplished by having their property values as their position vector and is done so similarly in the paper by Juan, Alvarez, Villasante and Ruiz-Frau[15] where they relabelled graphs with their normalised centrality values. This ranges from the betweenness centrality to eignevector centrality. Thus a similar idea can be applied along with the other properties explained and studied in this chapter.

## 2.4 Webpages

One of the largest graphs in the modern world is the network of web-pages, especially Google's search engine. To help navigate through the vast quantity of pages, Brin and Page[7] developed an algorithm known as the Page Rank algorithm. The Page rank gives a quantified meaning to the importance of the webpages and their links to other web-pages. Additionally the Page Rank is known as a centrality value which

was discussed briefly in the last section along with the Hyperlink-Induced Topic Search (HITS) which was also created with a similar goal to Page Rank. Both these values are based upon digraphs (directed graphs) and was generated to help with the navigation of the world wide web and provide user's with the highest quality webpages that were also most relevant to their search criteria.

### 2.4.1 Hyperlink-Induced Topic Search

HITS calculate the ranks of *authorities* and *hubs* in relation to their in-links and out-links[37], i.e. the edges pointing in and the edges going out of the vertices in the graph/network. Authorities and Hubs are assigned to webpages(vertices) depending on their number of in-links and out-links. For the HITS algorithm, the webpages who has lots of in-links pointing to it are denoted as authorities and the webpages who has lots of out-links pointing to other webpages are denoted as the hubs. These can be identified through the use of the HITS algorithm as the calculations are an iterative process where it enforeces the authorities and hubs by bringing the authorities to the surface of the graph. Thus isolating them from other the other webpages. This is achieved by generating the hub and authority values by mutual reinforcement. Based on vertices  $i \in V$  from a graph of webpages  $G = (V, E)$ , the hub value  $h_i$  and authority value  $a_i$  are first set to a value of 1. Then by Kleinberg's HITS algorithm, they are updated iteratively through the formulas:

$$a_i^{(k)} = \sum_{j:j \rightarrow i \& j \neq i} h_j^{(k-1)}, \quad h_i^{(k)} = \sum_{j:i \rightarrow j \& j \neq i} a_j^{(k)} \quad (2.30)$$

where  $j$  are denoted as the links from and to the webpages. The  $k$  depicts the  $k^{\text{th}}$  iteration of the algorithm, so the authority values depend on the previous iteration of hub values and the hubs are calculated based on the current authority values. Which gives the mutual reinforcement of both values. As the graph  $G$  can be represented by an adjacency matrix  $A = [a_{ij}]$ , the formulas can be expressed through matrices[11] and vectors instead as:

$$\mathbf{a}^{(k)} = \mathbf{A}^T \mathbf{h}^{(k-1)}, \quad \mathbf{h}^{(k)} = \mathbf{A}^T \mathbf{a}^{(k)} \quad (2.31)$$

where the hub and authority values are adapted into vectors  $\mathbf{a}^{(k)}$  and  $\mathbf{h}^{(k)}$ . The vectors are shown as

$$\mathbf{a}^{(k)} = \begin{bmatrix} a_1^{(k)} \\ a_2^{(k)} \\ \vdots \\ a_n^{(k)} \end{bmatrix}, \quad \mathbf{h}^{(k)} = \begin{bmatrix} h_1^{(k)} \\ h_2^{(k)} \\ \vdots \\ h_n^{(k)} \end{bmatrix} \quad (2.32)$$

Thus accomplishing the goals of generating the values which depict the hubs and authoritities more clearly within the graph. Although this algorithm is currently

defined only for directed weightless graphs so we extend the algorithm to include a weighted edge version.

### Normalisation and weights

As the number of iterations increases, the values generated may increase. This means that at some point, the values will become too large to be used for any calculations so they can be normalised such that this doesn't occur. After  $k$  iterations, the normalised formulas are demonstrated by Equations (2.34) and the general outline for the HITs algorithm is demonstrated by Agosti and Pretto[1] as the following:

```

 $\mathbf{a}^{(0)} := \mathbf{u}$  ,  $\mathbf{h}^{(0)} := \mathbf{u}$ ;
for  $k := 1$  to  $K$  do
     $\mathbf{a}^{(k)} = \mathbf{A}^T \mathbf{h}^{(k-1)}$ ;
     $\mathbf{h}^{(k)} = \mathbf{A}^T \mathbf{a}^{(k)}$ ;
    normalise  $\mathbf{a}^{(k)}$  such that  $\|\mathbf{a}^{(k)}\| = 1$ ;
    normalise  $\mathbf{h}^{(k)}$  such that  $\|\mathbf{h}^{(k)}\| = 1$ ;
end for
 $\mathbf{a} := \mathbf{a}^{(K)}$  ,  $\mathbf{h} := \mathbf{h}^{(K)}$ ;
```

where  $K$  denotes the maximum number of iterations and  $\mathbf{u}$  be the vector for the first iteration of hub and authority values, also known as the base case. The base case for  $\mathbf{u}$  will just be the vector of ones known as  $\mathbf{1}$  or  $\mathbf{e}$  in linear algebra. So then the normalised values after  $k$  iterations is given as follows:

$$\mathbf{a}^{(k)} = (\mathbf{A}^T \mathbf{A})^{k-1} \mathbf{A}^T \mathbf{u}, \quad \mathbf{h}^{(k)} = (\mathbf{A} \mathbf{A}^T)^k \mathbf{u} \quad (2.33)$$

To incorporate weights of the edges into the algorithm, the weighted matrix  $W$  for the graph  $G$  can be used in place of the adjacency matrix. Simply by replacing the  $A$  in the formulas, the weighted version can be generated as the formulas:

$$\mathbf{a}^{(k)} = (\mathbf{W}^T \mathbf{W})^{k-1} \mathbf{W}^T \mathbf{u}, \quad \mathbf{h}^{(k)} = (\mathbf{W} \mathbf{W}^T)^k \mathbf{u} \quad (2.34)$$

#### 2.4.2 Page Rank

A widely known algorithm that contributes to the internet of navigation within Google is the Page Rank. The Page Ranks are needed because many different search query's can be entered onto the search engine to produce lots of results that contain the same or similar words to the search enquiry so a method to help organise or prioritise is necessary. The Page Rank is used to rank the web-pages according to the amount of backlinks a page may have and the number citations that reference a particular page. An example of a webpages and links as a graph is shown in Figure 2.4.

For the graph  $G = (V, E)$  the webpages can be seen as vertices  $v \in V$  with their links as the edges directed edges between pages, the same way to how they were

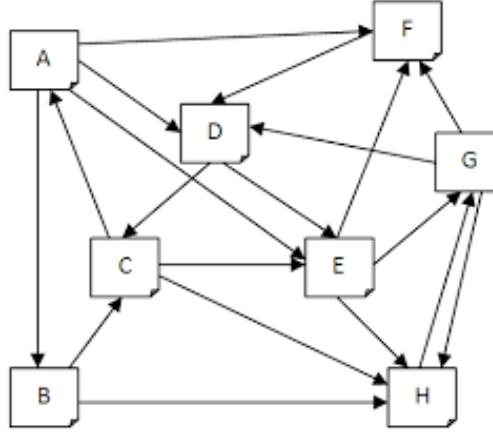


Figure 2.4: A graph representation of webpages and their in and out links. Algorithms such as Page Rank and HITS uses this graph format to help demonstrate their calculations. This image was sourced from the paper by Devi, Gupta and Dixit[16] whom describes and compares Page Rank algorithm to the HITS algorithm.

represented in the HITS formulas. We can let  $F_v$  be the set of webpages that  $v$  points to, i.e. the forward links. Similarly the backwards links as  $B_v$  which is the set of webpages that points to  $v$ . The simplified version of Page Rank[45] can then be defined with  $N_v = |F_v|$  as the following:

$$PR(v) = c \sum_{w \in B_v} \frac{PR(w)}{N_w} \quad (2.35)$$

where  $c$  is a variable used for normalisation purposes so can be modified accordingly. So Page Rank is calculated by the even distribution of the Page Rank for webpage  $v$  among webpages that  $v$  points to. The values that links to other webpages from  $v$  are then used to calculate their Page Rank. Hence giving an iterative approach to Page Rank as the algorithm travels along the links of the webpages. However if a webpage doesn't have outlinks and only inlinks from other webpages then their Page Rank is never distributed to others causing it's value to accumulate. This situation is known as a *rank sink*.

A remedy to having a rank sink is to use a damping factor to illustrate the probability that the user follows the links on the webpage. This takes into the consideration that users could skip pages or go directly to another webpage that wasn't linked through the url. Hence  $(1 - d)$  is considered as the distribution of the Page Rank from webpages that wasn't directly linked to it, i.e. no direct edges between them. Thus for the page  $v$  with  $b_i \in B_v$ , the Page Rank[7] is defined through Equation 2.36.

$$PR(v) = (1 - d) + d(PR(b_1)/F(b_1) + \dots + PR(b_n)/F(b_n)) \quad (2.36)$$

where  $F(v)$  was retrieved from the set  $F_v$  for webpage  $v$ . Page Rank is applied

to each page and it is repeated on further pages until the equation converges. The damping factor adds randomness into the network of webpages so that the rank sink doesn't occur and ensures the convergence isn't reached too quickly. Usually the damping factor is taken to be 0.85.

By using summation, the Page Rank formula can be simply reduced to

$$PR(v) = (1 - d) + d \sum_{w \in B_v} \frac{PR(w)}{N_w} \quad (2.37)$$

## Personlisation

The algorithm for Page Rank can be modified depending on the use and aims. Such modifications include adjusting the vertex and edge values, modifying the damping factor or to introduce a new variable into the algorithm. An example of a personlisation is the Weighted Page Rank[59] where larger values are assigned to more popular/important webpages rather than the even distribution that occurred beforehand. Webpages that are outlinked will instead receive a value proportional to the pages popularity which are based off of their number of in links and out links. These popularity values of the in links and outlinks are represented as  $W_{v,w}^{in}$  and  $W_{v,w}^{out}$  accordingly. So  $W_{v,w}^{in}$  is calculated by the inlinks of webpage  $v$  shown to be  $I_v$  and the inlinks of all the webpages that the webpage  $w$  references, call this set of pages  $R(w)$ , as  $I_p$  where  $p \in R(w)$ . This can then be formulated into the equation

$$W_{v,w}^{in} = \frac{I_v}{\sum_{p \in R(w)}} I_p \quad (2.38)$$

Similarly, the same can be done for the outlinks

$$W_{v,w}^{out} = \frac{O_v}{\sum_{p \in R(w)}} O_p \quad (2.39)$$

Where  $O_v$ ,  $O_p$  is defined the same as the inlinks previously but with the use of the outlinks as replacement. Therefore, the Page Rank formula can be modified to include the webpages importance giving the formula

$$PR(v) = (1 - d) + d \sum_{w \in B_v} PR(w) W_{v,w}^{in} W_{v,w}^{out} \quad (2.40)$$

So, this formula is focussed more upon the webpages that are visited more frequently by users and ensures they end up with a higher Page Rank.

However, for the purpose of general directed weighted graphs, the edge weights are lost in the formula as the algorithm updates the Page Rank of every vertex in each iteration meaning that the weights can be disregarded and replaced with the Page Ranks instead. To ensure this doesn't happen, the weights of all the edges must be included in the ranks calculation. This is accomplished by summing up

the weight values of the in and out edges and incorporating it into the formula as achieved by Equation 2.41.

$$PR(v) = (1 - d) + d \sum_{w \in B_v} \frac{PR(w)w_{(w \rightarrow v)}}{N_w} \quad (2.41)$$

where  $N_w = \sum_y A_{w,y}w_{(w \rightarrow y)}$  is redefined as the sum of weights of the out linked edges in relation to vertex  $w$ .

Whilst HITS and Page Ranks are designed for the search engine, the network of web-pages essentially is a large directed graph that can contain weights hence why the Page Rank can be used on any derivatives of a weighted directed graph. Thus, there are additional graphical property that can be used to help analyse the structure and linkage of a graph.

# Chapter 3

## Application of Graph Properties

As various graph properties have been discussed and defined, we now apply these properties onto connected graphs to demonstrate their use. As well as this, they are used to outline and modify the graph such that an alternative layout may be given, revealing correlations between the vertices or edges within the graph. Through the use of Python, I have coded a program to display a graph either generated from a weighted matrix, a adjacency matrix or a graph data set that's pre-existing. To help accomplish this, I have used Tiago's Graph Tool library for python which contains useful documentation and functions to achieve the graph generations as well as other mathematical libraries for complex arithmetic. The general idea is to compare various graph properties by modifying their positions according to the values of their graph properties. For simplicity and the goal of being comparable, we choose the y-axis of the graph to be based upon the trophic coherence and the x-axis to vary between the different properties discussed in the last chapter.

### 3.1 Early Experiments

By programming and using tools from the various libraries within python, I was able to generate graphs and define each property value based upon the generated graphs. Then the positioning of each vertex can be modified in relation to the values calculated.

Initially, out of the many pre-existing graphical datasets from Tiago's library, I experimented on a smaller dataset that demonstrates the relationships between karate clubs in a city so that I can test and generate a visualisation of this dataset. This dataset involved 34 karate clubs where the initial graph can be shown by Figure 3.1 with the adjacency matrix in Appendix A.1.

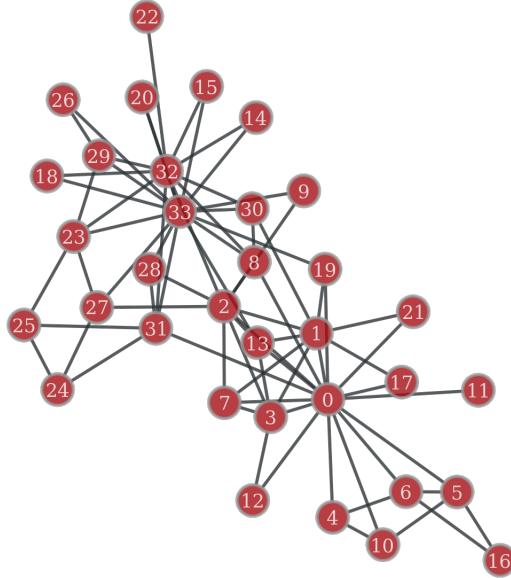
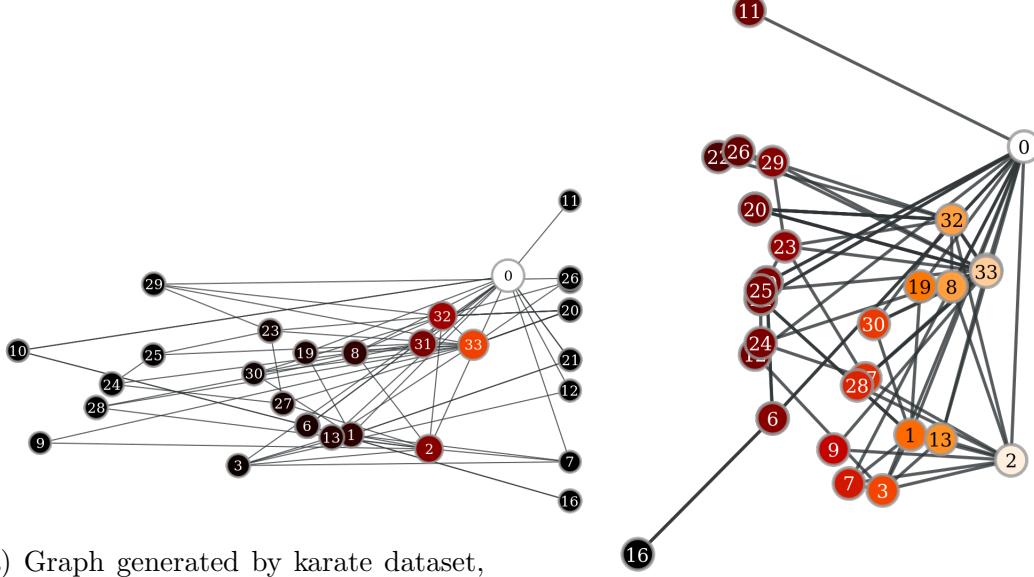


Figure 3.1: The initial graph based on the karate club dataset pre-existing within the library that was generated through the python program. Contains 34 clubs and their connections to one another with no important meaning with the positions of the vertices.

Current positioning of the graph are determined based on the idea that the vertex do not overlap and the connections are all easily visible. So the dataset's positioning has no real benefits other than having good visibility. Any correlations or vital information can not be derived from the initial graph as the positions do not represent anything. The only information that can be derived easily are the certain outliers who only have one edge, vertices 22 and 11. The vertices that have more edges can also be seen such as vertices 32, 33 and 0 under closer examination. But these vertices are difficult to distinguish at a first glance. Thus, we now include various different graph properties discussed in the last chapter to ensure that more can be derived. This can be seen in Figures 3.2 and 3.3 by using the trophic coherence values for each vertex as their y-value and another property value for their x-value.

For karate clubs, as the graph is not directed, the trophic levels do not represent a clear hierarchical format since there is no distinction between "upstream" or "downstream" of information with the edges of the graph. Consequently, the vertices in the karate graph do not have information that passes in one direction. So for vertices such as 1 or 16, they could be recognised as the start or end of the network hence the reason they have the largest and smallest trophic levels. Even if the flow of information is bidirectional, trophic levels are still useful in implementing some structure into the karate dataset. Further datasets involving languages that will be explored will be directional so that trophic levels can be used optimally.



(a) Graph generated by karate dataset, the y-axis represents the trophic coherence values and the x-axis represents the betweenness centrality values for all the vertices. Additionally added colour changes between the x-axis to give a clearer visualisation of the separations.

(b) Graph generated similarly to the betweenness graph for the karate dataset but the vertices are plotted with local clustering coefficient as the x-axis and trophic levels as the y-axis.

Figure 3.2: Centrality values as the x-axis

Figure 3.2a shows the karate graph with the x-axis representing the betweenness value. The betweenness value are scaled by a factor of 10 to give a clearer visualisation with larger betweenness value further on the x-axis. Which also means the vertex is more frequently involved among short paths of the connections. In this case, the connections of the clubs. We can see that vertex 0 is the furthest right vertex so is involved in the most short paths between all the vertices of the karate graph. On the other hand, vertices such as 11, 7, 16 and more are clubs who are on the outskirts with no proper connections to other clubs. Can be seen on the left side on the figure. Therefore, in relation to this dataset, the clubs with larger betweenness are the clubs who are more centralised in a city and have more meaningful links to others.

We compare this to another centrality value, the closeness values who also has been scaled by a factor of 10. Figure 3.2b represents this in place of the betweenness value. Through comparison of both, vertices are positioned similarly to betweenness. This is because betweenness are values when considering all vertices within the graph whereas closeness considers all neighbours of a specific vertex. In other words, betweenness measured the control a vertex has over the flow of information through the entire graph, whereas closeness measures the control over the flow of information

with vertices in close proximity (i.e. neighbours). Therefore the vertices on the right of both the betweenness and closeness graph would be the most important/largest clubs. Additionally, clubs such as 19 who have a larger closeness compared to betweenness means that it is important to the clubs in close proximity of itself, in other words, the club is the most important/largest within its local area.

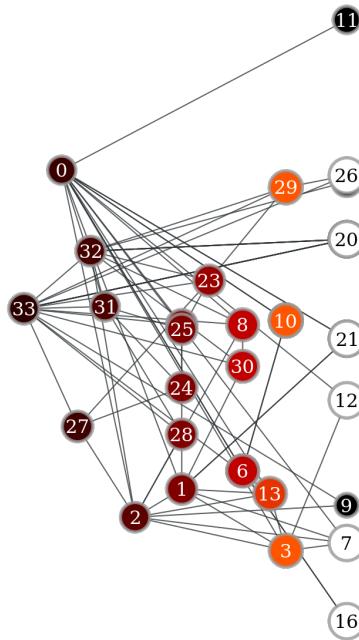


Figure 3.3: Graph generated similarly to the betweenness graph for the karate dataset but the vertices are plotted with local clustering coefficient as the x-axis and trophic levels as the y-axis.

Figure 3.3 is shown with local clustering coefficient as the x-axis instead of betweenness. After manipulation the vertices positions, notice that the clubs with less connections to the major clubs (vertices of high degree) are visually seen to the left in both graphs as these vertices would have a low clustering coefficient as well as a low betweenness. On Figure 3.3, the vertices with the best connections to major clubs are seen further to the right which we see are the karate clubs 26, 20 and 12 etc. However vertices with high betweenness seen before such as vertex 33 instead have a smaller local clustering as their club has connections to smaller clubs, decreasing the overall value of it's own connections. So this means that clubs on the right have quicker access/communication with better clubs and are the closest to them.

All the values for each club can be shown in table format, Table 3.1, where apart from the initial column, each column is one of the different graph property value.

Table 3.1: Table containing all the values calculated for each vertex of the graph.. The order is in the club number depicted in the first column.

<b>Club</b>	<b>Trophic Levels</b>	<b>Betweenness</b>	<b>Closeness</b>	<b>Local Clustering</b>
0	0	0.437645803	0.568965517	0.15
1	2.12094015	0.053873557	0.485294118	0.333333333
2	2.30804964	0.152263709	0.559322034	0.244444444
3	2.53039414	0.011961881	0.464788732	0.666666667
4	1	0.000631313	0.379310345	0.666666667
5	2	0.029987374	0.38372093	0.5
6	2	0.029987374	0.38372093	0.5
7	2.47969196	0	0.44	1
8	1.02851103	0.056737013	0.515625	0.5
9	2.22893206	0.000847763	0.428571429	0
10	1	0.000631313	0.379310345	0.666666667
11	-1	0	0.366666667	0
12	1.53039414	0	0.370786517	1
13	2.15210654	0.04159151	0.507692308	0.6
14	0.45900156	0	0.370786517	1
15	0.45900156	0	0.370786517	1
16	3	0	0.284482759	1
17	1.12094015	0	0.375	1
18	0.45900156	0	0.370786517	1
19	1.02788172	0.02936057	0.492537313	0.333333333
20	0.45900156	0	0.370786517	1
21	1.12094015	0	0.375	1
22	0.07623827	0	0.34375	0
23	0.73415591	0.018244949	0.392857143	0.4
24	1.44596889	0.002209596	0.375	0.333333333
25	1.05781988	0.003840488	0.375	0.333333333
26	0.03492233	0	0.358695652	1
27	1.70452843	0.02170214	0.452054795	0.166666667
28	1.75702472	0.001794733	0.445945946	0.333333333
29	0.1140399	0.003869048	0.38372093	0.666666667
30	1.30422637	0.014727633	0.458333333	0.5
31	0.90660502	0.140348425	0.540983607	0.2
32	0.53811913	0.182157888	0.515625	0.181818182
33	0.92088243	0.275055616	0.540983607	0.116666667

This is the early experimentations of the positioning of vertices within the graph so instead of using simple datasets, I will generate the datasets based upon various different languages as well as their sentence structure. To understand this further, words in languages must be given a rank which can be demonstrated through Zipf's Law discussed in the next section.

## 3.2 Zipf's Law

Zipf's law analyses the natural languages and the frequency of words that appear in them. Alternatively, Zipf's Law[27] is generally seen as the frequencies of specific events are inversely proportional to their rank that is determined through this law. The law was proposed by George Kinbgsley Zipf when researching the various frequencies of words within the English language. The law states that the  $r^{\text{th}}$  most frequent word in the language has a frequency of  $f(r)$  that has a relation with the inverse of  $r$  where  $r$  is the *frequency rank* for the word and  $f(r)$  as the frequency of the word in the corpus examined (The *corpus* means the collection of written text).

$$f(r) \propto \frac{1}{r^\alpha} \quad (3.1)$$

This is the scale for  $\alpha \approx 1$  and means that the most frequent word in the examined text which is  $r = 1$  has its frequency of appearance to 1, the next most frequent word which is  $r = 2$  has a frequency of  $\frac{1}{2^\alpha}$  and so on. This Zipf's law can be drawn on a graph to show a relation and when  $\log(f)$  is drawn against  $\log(r)$ , the graph generates a curve that closely resembles a straight line with a slope of -1. This is known as Zipf's curve and later in the 1960s, this was reinforced by the law being correct for smaller corpora[53]. However the curve varies depending on the corpora as expected and the higher ranking words deviated more from the straight line. Therefore, Mandelbrot derived a generalisation for Zipf's law to adjust to the frequency distributions within the different languages. Mandelbrot proposed to adjust the rank by a constant  $\beta$ , demonstrated by

$$f(r) \propto \frac{1}{(r + \beta)^\alpha} \quad (3.2)$$

Generalisation of Zipf's law can then be applied to various different corpus of languages so that a frequency distribution can be viewed for the corpus. An example of this can be seen in Figure 3.4.

So words in a corpus has a systematic relationship between their rank in their occurrence table and their frequency that they appear in their corpus. Meaning that there are words within languages that are used more commonly such as "the", "or", "of" that account for most of the word occurrences and other words such as "xylophone" and "accordion". A larger corpus is studied by Bentz, Kiela, Hill and Butterly [3] where they study zipfs law for Old English and Modern English. They study the frequency and ranks of each word and then compare them between the old and new English. In doing so, for Old English, the words "and" is ranked first with a frequency of 1731 whereas in Modern English, "the" is ranked first with a frequency of 1775 and "and" is second instead with a frequency of 1024. By looking at more words and the comparisons between them, Old English has a larger number of distinct types whilst Modern English had less but this meant that they had higher frequencies within their first 100 words. In conclusion Zipf's Law is a useful tool as

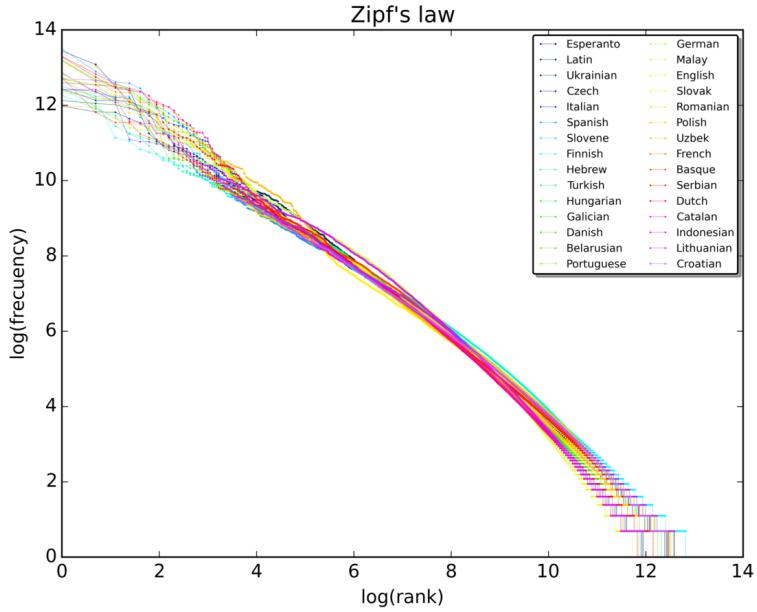


Figure 3.4: The plot of Zipf's law containing 30 different language corpus generated from the first 10 million words in each language from Wikipedias. The image was sourced from Wikipedia[30].

languages tend to follow Zipf's curve in terms of their frequencies and rank meaning that by following an existing corpus of language, the data can be extrapolated and used in other corpus's to determine similarities between the known and the unknown texts. As well as to use it to determine the types of words that are deemed to be most common in a language. We will use Zipf's Law when comparing texts of different languages in further experiments.

# Chapter 4

## Analysing Languages

Through the earlier experimentations with the graph generations based on specific properties, improvements have been made. This includes normalisation of the values into the range of 0-10 rather than using a scale factor. Which ensures that all the graphs will be similarly shaped and have the same axis size for an easier visual comparison. Also I've included page rank as well as the other graph properties used before. Finally the datasets I will use will be generated through my program by an input of text. Afterwards will be converted into graphs where the words represent the vertices and the edges are directed to the next word in the order of the text. A factor that affected the dataset was the punctuation so to achieve congruent data, the punctuation are stripped unless they are sentence enders such as full stops, question marks, exclamation marks etc.

Therefore, the graph we experiment on are directed graphs generated from various language families.

### 4.1 Linguistics

Modern languages are descendants of ancestral languages through evolution of linguistics. Through the different ages of the world, language has always been a key part in communication between societies. They are developed and taught to newer generations to reach the stages in the current world. The history of languages can be viewed as a family tree where modern languages are nearer the bottom. Along this trees, there are groups of languages that share a common ancestor. These are the *language family* of those languages.

Estimations of around 500 language families exist and Campbell [9] has reported that there are exactly 406 independent language families including dead languages and *language isolates* (where the language does not fit into a language family). According to Ethnologue[18], who are the research centre for language intelligence, there are 142 different living language families. Of the living families, six are considered to be the major families and are Indo-European, Afro-Asiatic, Niger-Congo,

Austronesian, Sino-Tibetan and Trans-New Guinea.

The aim of my research is to study modern languages that fall under the Indo-European language family and the Sino-Tibetan language family. The main families are known as the *proto-language* as they are the parent language where many languages are derived from[49]. The languages being English, German and Dutch that are Germanic language under the Indo-European family. Russian and Polish which are Balto-Slavic language under Indo-European family. French and Spanish which are Latin languages under the Indo-European language. Finally Chinese which falls under the Sino-Tibetan family. Furthermore, I will also look at Japanese which part of the Japonic language family but it would be considered a language isolate if the Ryūkyūan languages were not distinct from Japanese[8]. Therefore these are all the languages in which I will translate the text extract into datasets.

## 4.2 Text Corpus

The best way in comparing the results to other languages is to have a dataset that is based on the same text extract. Thus, the text extract that is chosen should be simple and well known, in my case, I have chosen to use the popular story in all languages, "Sleeping Beauty". To ensure that the same version is used, the Grimm Brothers version is utilised where the original was in German and extracted from the book of children stories "Kinder- und Hausmärchen"[24]. So using translations of this story, graphs are generated based on each language studied. Additionally, instead of using the entirety of the story, the first two paragraphs are used so that the graphs are not overwhelmingly dense which tells the story as follows:

In times past there lived a king and queen, who said to each other every day of their lives, "Would that we had a child!" and yet they had none. . . There were thirteen of them in his kingdom, but as he had only provided twelve golden plates for them to eat from, one of them had to be left out.

In conclusion the original nine dataset created in this way can be used for graph property calculations. Next section will study English version of the story which is part of the I

## 4.3 Indo-European Language Family

Initial the languages of the Indo-European family are studied and a few of them are described in detail later including English and German. Throughout the analysis, words are referred to as vertices, the two paragraph extract of "Sleeping Beauty" in the relative language will be referred as the story corpus and vice versa. Each word in a graph has edges which are directed to the next word relative to the two paragraph extract.

### 4.3.1 English

English words can be organised into eight different parts of speech; Nouns, Pronouns, Adjectives, Adverbs, Verbs, Prepositions, Determiners and Conjunctions. Linguistic researchers focus on the use of these categories in different situations such as through speaking or through magazines[35]. We will study the appearances of these categories in our story corpus. To achieve this, the story corpus is received as an input to my program so it may generate a usable dataset. The dataset is then converted into a directed word graph as shown by figure 4.1a. Additionally in replacement of having each vertex labelled by the corresponding word, each vertex will be labelled with an integer. The relative integer for each word will be shown on the table of values for the graph. The table will be shown later in Table ???. So we achieve the initial directed word graph both in original form and alternated form shown in figure 4.1b.

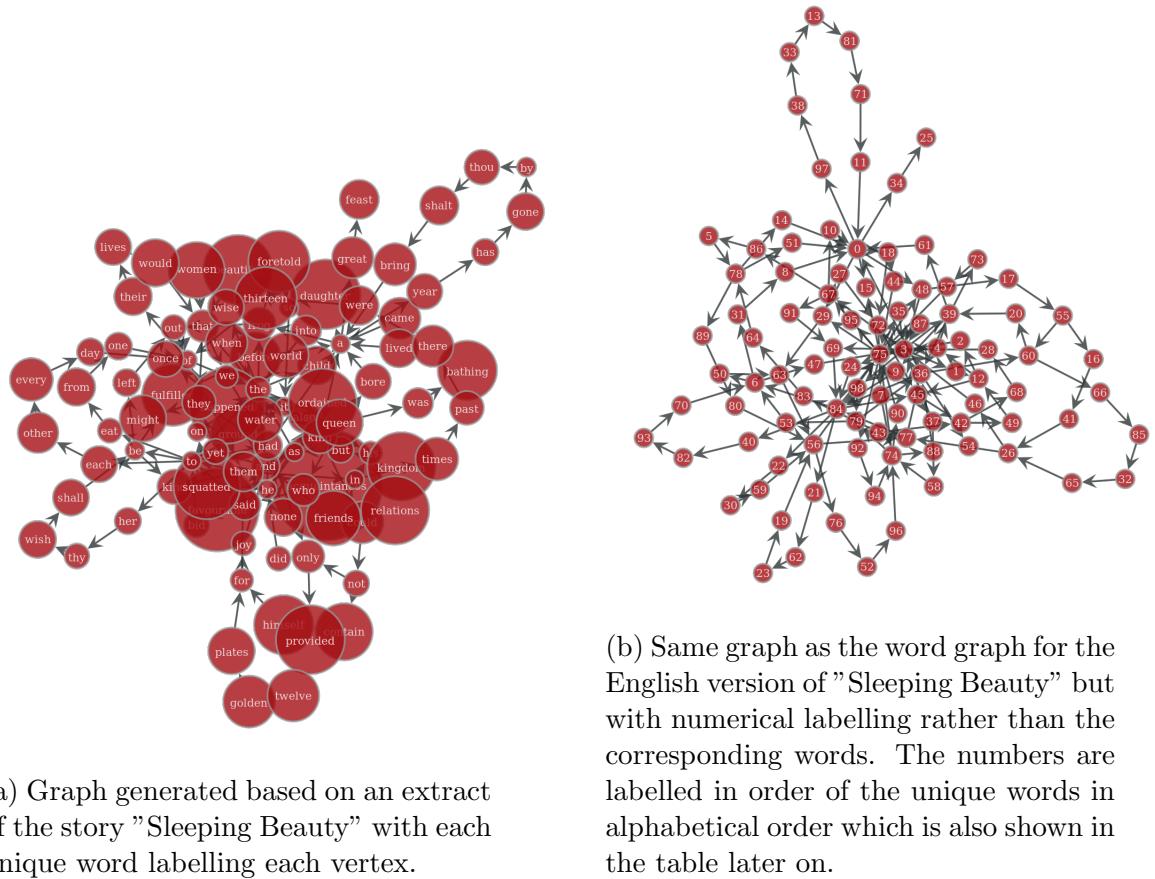


Figure 4.1: Graph created from the corpus labelled with words or in integers.

As done in the Early Experiments of the karate club dataset, we calculate the values of the various graph properties explained in Chapter 2. These are graph

properties such as local clustering coefficient, betweenness centrality, closeness centrality, trophic levels and page rank. Values are organised into their corresponding columns and presented as a table with the ten most frequent words shown in Table 4.1a below. Table also includes the number of appearances the word has denoted as count and the relative vertex number in the numbered English graph. The entire table can be seen in Appendix B.2.

Vertex	Words	Counts	TC	BC	CC	LC	PR	Vertex	Words	Counts	TC
3	and	9	3.00	10.00	2.49	0.57	8.13	25	feast	1	10.00
75	the	9	4.34	9.19	2.20	0.44	10.00	95	world	1	8.35
0	a	7	1.98	9.26	2.22	0.00	7.65	54	none	1	7.45
84	to	6	3.80	7.48	2.28	0.18	5.54	15	child	2	7.17
36	had	4	3.44	5.32	2.48	0.36	3.82	63	out	2	6.01
39	he	4	2.42	3.68	2.07	0.48	2.38	34	great	1	5.99
56	of	4	5.19	6.21	2.37	0.00	6.03	69	said	2	5.62
74	that	4	3.73	3.49	2.19	0.36	5.12	56	of	4	5.19
6	be	3	3.78	3.15	1.91	0.00	3.24	19	day	1	4.91
12	but	3	2.05	1.29	1.95	0.00	1.73	50	left	1	4.90

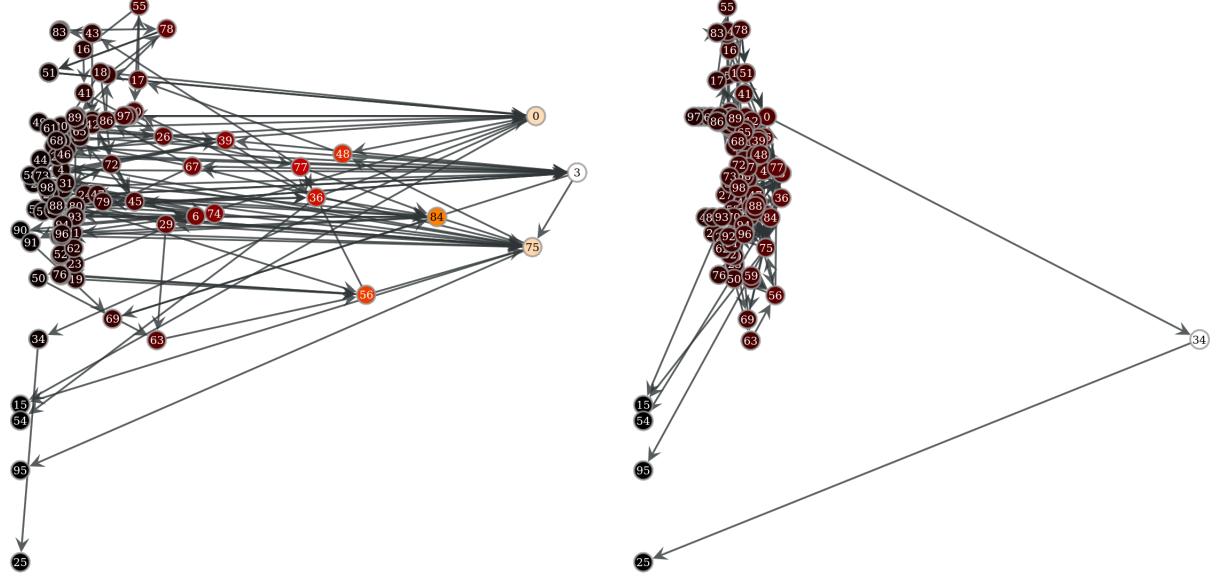
(a) The first 10 most common words of the dataset. Generated from the English version of "Sleeping Beauty" in a table format.

(b) Top 10 words ranked by their trophic levels based on the English Story Corpus.

Table 4.1: Partial extracts of the table data for graphical properties of the English Story Corpus.

We begin by analysing words with the most recurrences which are, in the order of most frequent to least, "and", "the", "a", "to", "had", "he", "of", "that", "be" and "but". Note that nouns, adverbs and adjectives do not appear in the most recurring words. These are the words that are deemed more vital in the creation of structure within a sentence. Any sentence in "Sleeping Beauty" will have a high chance of containing at least one of these words. Since this is a small dataset, we can compare these words to a larger dataset for word reoccurrences as the Zipf curve for language are similar as discussed in a previous section. We choose to compare the story corpus to the British National Corpus (BNC)[13] which is a 100 million word collections that includes both written and spoken language. The benefits in choosing BNC is that it contains older English so may provide clearer correlations to the story of "Sleeping Beauty" as the Brothers Grimm version began in the late 18th century. Hence for the BNC, the top ten frequent words[39] in order of frequency are; "the", "of", "and", "a", "in", "to", "it", "is", "to" and "was". Comparing the most frequent words in both corpora, correlations are achieved such as the repetitions of words "the", "to" etc. Such similarities reinforce the fact that the English language has a structured form that requires the use of these such words, as demonstrated with a much smaller corpus compared to the BNC.

Now onto the analysis of Trophic levels and coherence (see Table 4.1b). When applying trophic level calculations on directed word graphs, the levels represent the position of the word in a sentence similarly shown in the analysis of network data in empirically-derived directed networks[33]. For the "Sleeping Beauty" English word graph, the lower trophic levels tends to be sentence starters and the higher levels are the ends of sentences. This is supported by the data because the top five words with largest trophic levels values (ranging from 10.00-6.01) are all sentence enders. Along with the bottom five being words (trophic values ranging from 0.49-0.00) nearer the start of sentences such as "there" and "times" in relation to the corpus. However trophic incoherence is calculated to be 0.969 which means that the levels in the graph can not be distinguished and are not clear. Since the trophic coherence =  $1 - \text{trophic incoherence} = 0.03$  which is due to the fact of the vast difference of sentence lengths in the corpus. Which varies from the shortest sentence of five words and the longest of forty seven words. Consequently the for the extract used, it may not provide a clear hierarchical layout but still provides a good layout of sentence flow from the lowest level to the highest level. Demonstrated by further graphs with the trophic levels as the y-axis ranging from 0 at the top to 10 at the bottom to provide a normal cascade of words. Includes other graphical properties starting with betweenness and closeness centrality (Figures 4.2a and 4.2b) shown next.



(a) Numbered English graph positioned with trophic levels (y-axis) against the betweenness centrality (x-axis).

(b) Numbered English graph positioned with trophic levels (y-axis) against the closeness centrality (x-axis).

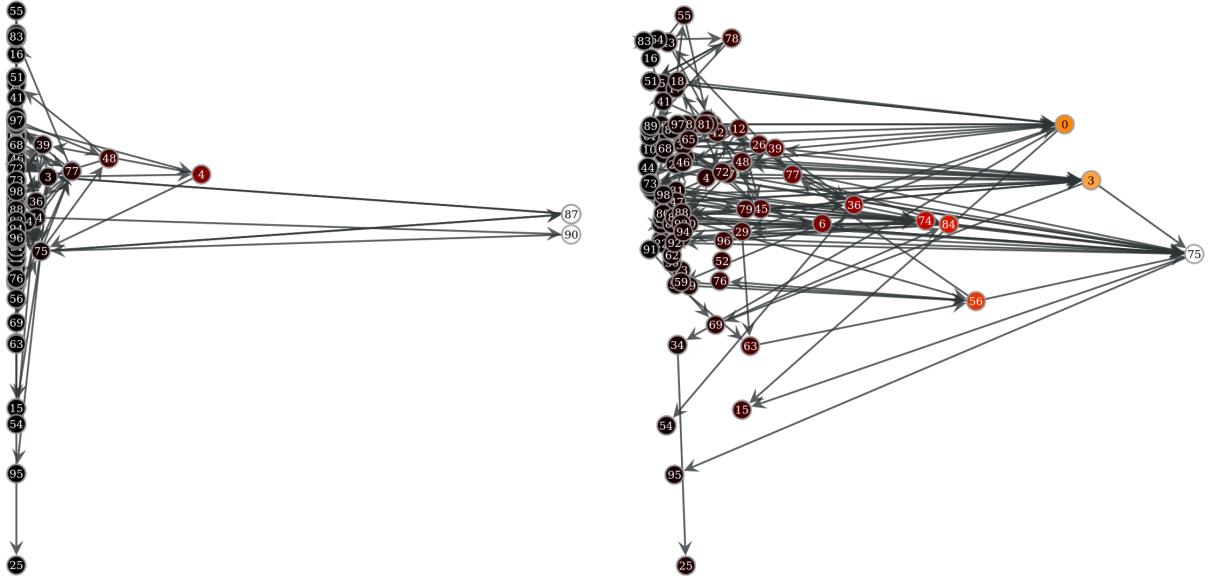
Figure 4.2: Betweenness and closeness centrality values displayed on the x-axis in graphical form.

As visually demonstrated on Figures 4.2a and 4.2b, the centrality values for each word is plotted against their trophic levels. Both have been normalised to a range of 0-10 with Figure 4.2a showing the betweenness centrality on the x-axis and Figure 4.2b showing the closeness centrality. Key vertices identified in relation to their betweenness values are vertex 3, 0 and 75 which are the words "and", "a" and "the" respectively. These are conjunctions and determiners of the English language and they have the largest frequency of appearance in the corpus relating to the word counts discussed before. So, there is a strong link between the betweenness centrality of vertices to the word counts in the text. Furthermore, these words are common in forming correct structure of an English sentence meaning that high betweenness associates the words as key bridges within a sentence.

When considering closeness centrality, the graph shows that almost all vertices have a larger closeness value in comparison to their betweenness. This is because the closeness analyses the importance of the words as within their clusters rather than the graph as a whole. So words with high closeness are key connectives in their relative clusters, in other words the sentences they are part of. However vertex 34 (the word "great") is an outlier and by further analysis, vertex 34 is the only connect between vertex 0 and 25. Vertex 25 having the highest trophic level and vertex 0 has a low trophic level and a higher degree. Consequently, the closeness

value for 34 is much higher due to the fact that it is the only predecessor of vertex 34 which in itself only has one predecessor. Thus meaning that vertex 34 is the only local bridge in its sentence giving it an extreme closeness.

In conclusion, based on the story corpus, betweenness finds the words most commonly used as connectors in sentences given the whole extract and closeness finds the words as connections within a close range of one another.



(a) Numbered English graph positioned with trophic levels (y-axis) against the local clustering coefficient (x-axis).

(b) Numbered English graph positioned with trophic levels (y-axis) against the page rank (x-axis).

Figure 4.3: Similarly as before with betweenness and closeness but with local clustering and page rank instead.

Finally Local Clustering and Page rank of the vertices in the graph are presented similarly as before with its local clustering in Figure 4.3a and page rank as Figure 4.3b. Immediately, the local clustering graph shows very few vertices who has a high local clustering, the main ones being vertices 87 and 90 (words "water" and "when" respectively). However, these do no give a clear relation other than the words connected to and from these vertices have a high degree or importance in the graph. These are words such as "and" and "a". On the other hand this is only for the English version of the story so other languages may lead to different results. The page rank of each vertex shows the importance of each word beyond their direct contact. Essentially it has elements of both closeness and betweenness centrality which the graph reinforces.

In conclusion, when analysing the English language, trophic levels have provided a naturally flow of data presentation from top to bottom in the various graphs gen-

erated. Some sections are of the graph are also grammatically correct. Betweenness Centrality and Page rank both identifies the words of most importance with respects to the words that connect them. Closeness centrality also identifies words of key importance but at a local level. As this is a smaller corpus, these words are similar. Finally the Local Clustering does not provide sufficient benefit when visualising the dataset in the English language. Therefore the results generated based on the English version can be expanded to represent the English language and can even be used to predict and analyse unknown texts or missing words with a given text by representing the missing words as vertices. The position of the vertices will determine it's importance and use within a sentence. This is of course for the English language or language with similar structure.

Now we move onto the analyse of a different translation of the corpus, German.

#### 4.3.2 German

German is a Germanic language under the Indo-European language family, same as English. However whilst Modern English no longer use inflectional case system, the German language still does[17]. So as well as the parts of speech like in English, German words can be divided into two groups, the ones who are *inflectable* and *uninflectable*. If a word is inflectable then their form changes based on the context that the word is used in. These include the three genders for the words, the four cases (nominative, accusative, genitive and dative) and number (singular or plural). Uninflectable words are known as modal particles and mainly used to highlight an emotion of the sentence in spoken language. Therefore German language considerably more varied than the English language so the dataset in theory would contain more unique words overall. This is proven to be true as the dataset for German contains 109 words whilst English had 99.

Expectations of the graph properties for the German language is that there would be more unique words of higher importance based on the different genders for words. On average, sentence lengths should also be longer due to the generated words which will be seen in trophic levels if it is the case. Word graphs for the German version of the corpus are generated and shown in figure 4.4. Similarly as before, each number references the same word in it's position.

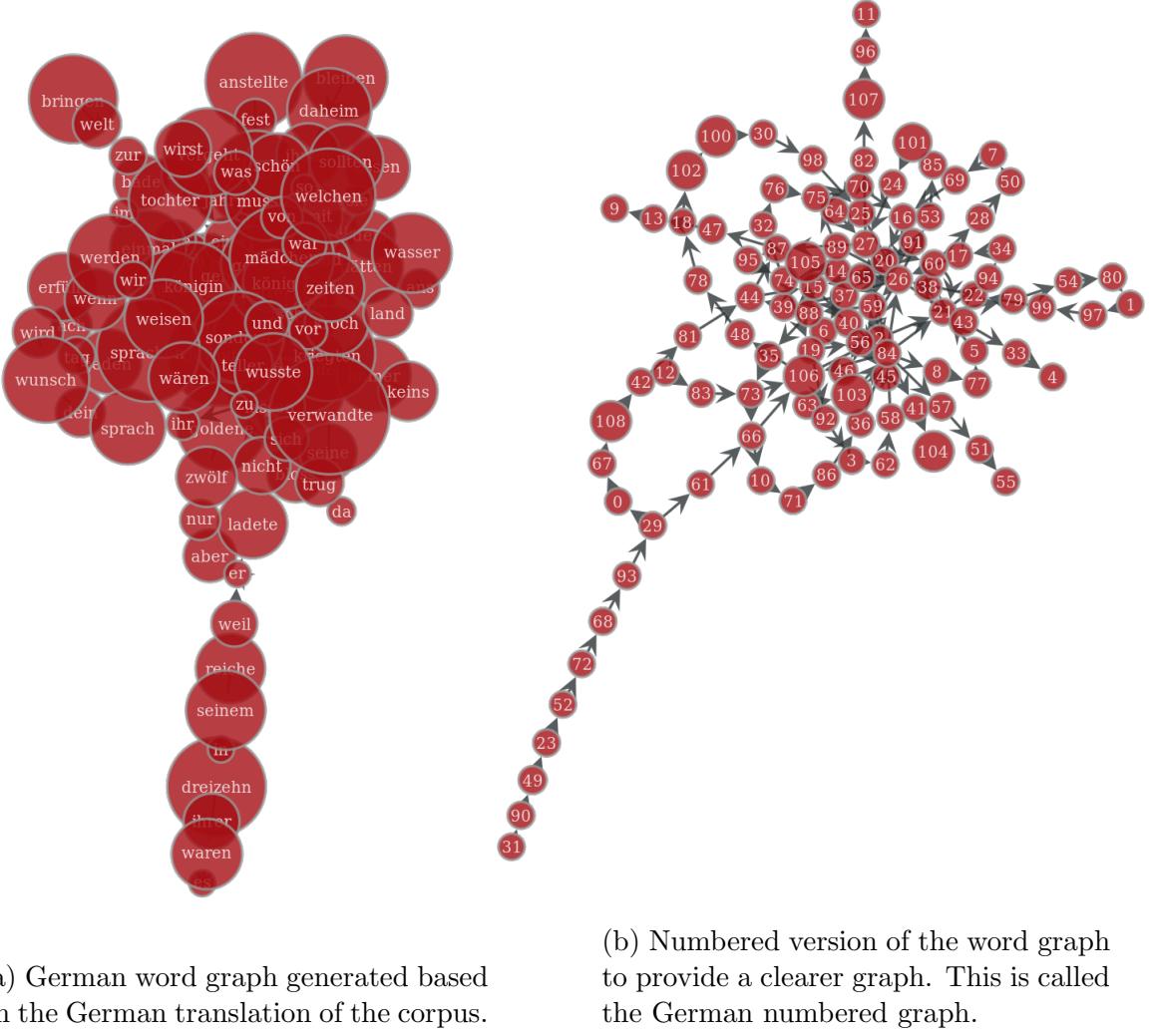


Figure 4.4: The German word graph and numbered equivalent of the word graph generated from the German translation of the "Sleeping Beauty" corpus.

Table 4.2a shows the most common ten words in the German dataset along with each graph property value. The translation to English for each word in the order of the table is as follows; "and", "a" (Masculine), "the" (Feminine), "to", "a" (Feminine), "Queen", "King", "not", "himself/herself/itself" (dependent on the pronoun this refers to) and "the" (Neutral). As expected when comparing to the English dataset most words appear in both translations as the top ten, in particular the top three. Which we recall was "and", "the" and "a". However rather than a cumulative count of "the" in English, the German translation has multiple versions which may counteract each others importance. By retaining the inflectable changes, the average count is lower for the German translation. Additionally the range has also decreased. For example, "the" in English was split into "die", "der" and "das"

in German. By noting these key differences, the graphical properties will be analysed starting with trophic coherence.

Vertex	Words	Counts	TC	BC	CC	LC	PR
84	und	7	6.85	9.90	2.51	0.77	8.35
26	ein	7	6.90	10.00	2.19	0.66	10.00
21	die	4	6.59	3.37	1.56	0.00	4.49
106	zu	3	6.42	4.38	1.65	4.00	2.77
27	eine	3	7.07	2.93	1.73	0.00	4.17
58	königin	3	6.41	2.84	1.78	0.00	2.67
57	könig	2	6.57	5.37	2.21	10.00	2.15
66	nicht	2	6.18	2.41	1.54	10.00	1.80
73	sich	2	5.97	1.52	1.55	10.00	1.45
15	das	2	6.91	2.42	1.94	0.00	2.75

(a) Top 10 words with the highest frequency in the German translation of the corpus. Shown in table format with other graphical properties.

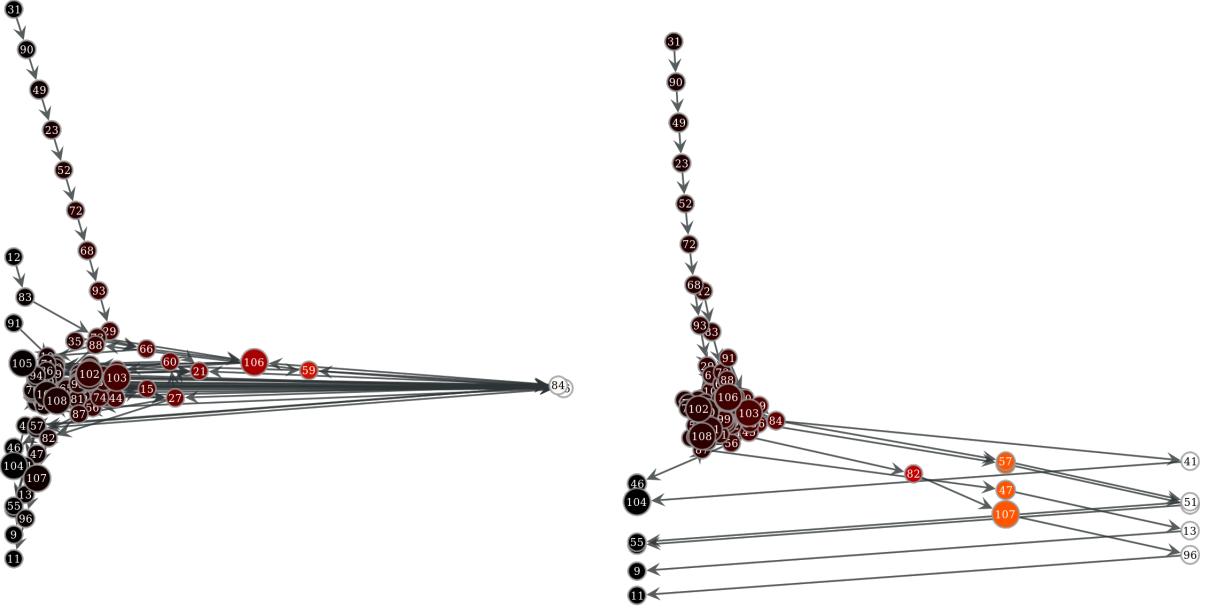
Vertex	Words	Counts	TC
11	bringen	1	10.00
9	bleiben	1	9.56
97	welt	1	9.27
4	anstellte	1	9.09
55	keins	1	9.04
13	daheim	1	8.83
107	zur	1	8.54
33	fest	1	8.36
51	immer	1	8.31
91	wären	1	8.31

(b) Top 10 words with highest trophic levels in the German translation dataset.

Table 4.2: Partial extracts of the table data for graphical properties of the German Story Corpus.

Analysing the trophic levels, top ten seen in Table 4.2b, the results show that most words congregate at around 7.1. This can be more evidently seen in visual representations on figures 4.5 and 4.6. From the same figures, there is a unique path from vertex 31 to 29 which is equivalent "Es waren ihrer dreizehn in seinem Reiche, weil er". This is the beginning of a sentence belonging to the German story corpus and holds all unique words that are not used in other sentences. Hence has not been influenced by other vertices so has a clear hierarchical structure which leads to the low trophic levels. Whereas most other sentences share words which causes the conglomeration nearer 7. So the German language has more options for word choices that causes the increased possibility of unique sentences. Meaning that their trophic levels are not evenly distributed like the English version.

Trophic coherence has also benefited from the uniqueness of word choices as for the German graph, trophic coherence is calculated to be 0.23. This is larger than the English graph which was 0.03. Therefore the German graph has a larger hierarchical structure compared to the English graph with clearer levels as demonstrated in further graphs below when analysing other properties.



(a) Positions of the German numbered graph but with trophic levels and betweenness on the y-axis and x-axis respectively.

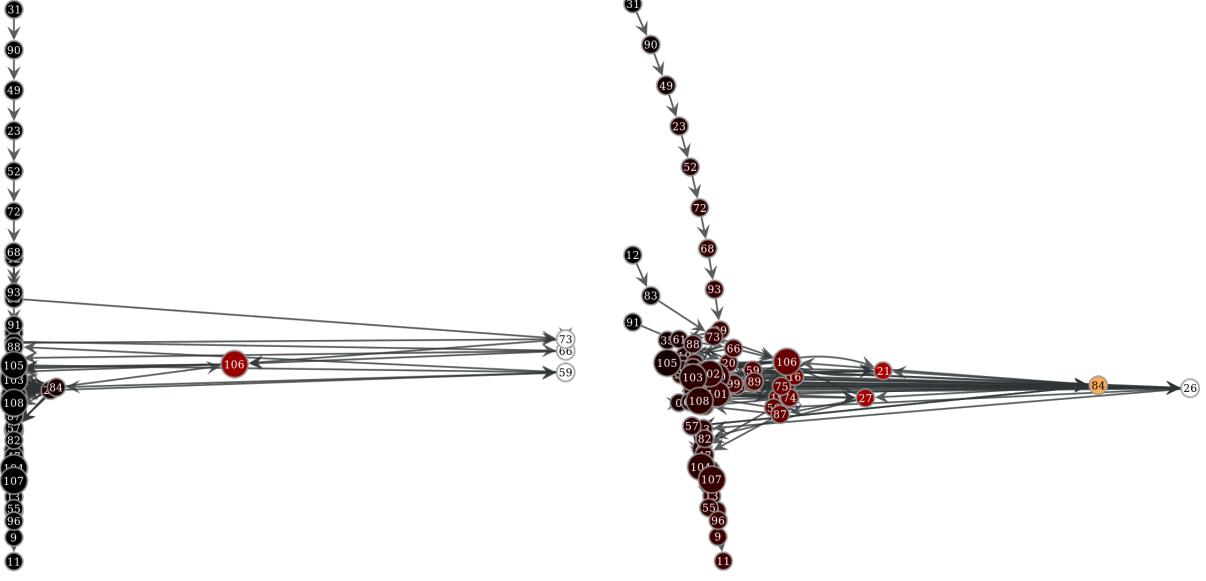
(b) Similar to the betweenness German graph but with closeness centrality values instead.

Figure 4.5: Betweenness and closeness centrality values displayed on the x-axis based on the German numbered word graph.

Assessing the graphs visually through Figures 4.5a and 4.5b. Vertices 84 and 26 have the highest betweenness which correlates to the frequencies of the words in the German Story Corpus. This was also the case of the English version as these refer to the words "ein" and "und" which are the most commonly used as bridges/links in sentences.

With closeness centrality, since it is a measure of influence on nearby words, the graph shows that the vertices 97, 13, 33, 51 and 41 has the largest values. These all correlate to the second to last word of each sentence apart from the first sentence which has the word "Kind". However "Kind" is used elsewhere meaning that the vertices mentioned before are unique and essential as a bridge in its nearby words. Otherwise the graph will be disconnected if they are removed. Note that the orange vertices in Figure 4.5b are the unique words predecessors to the vertices mentioned earlier and the last word of every sentence always has a value of zero.

Therefore betweenness identifies the words of key importance that are used most commonly as connectors. Meanwhile closeness identifies the words most likely to isolate vertices of the graph when following the sentence flow, i.e. breaks up the rest of the sentence into unique words.



(a) Displays the local clustering coefficient against the trophic levels.

(b) Displays the page rank against the trophic levels.

Figure 4.6: Displays the local clustering and page rank on the x-axis instead of the centrality values.

Nothing clear can be correlated by studying the local clustering as almost all have a local clustering of zero apart from six vertices. Three of which are vertices 57, 73 and 66 which corresponds to the words "könig", "sich" and "nicht" which are not unique words in the German story corpus. Whereas vertices with a high page rank correlate to either conjunctions, vertex 84 (und), or words that accompany other words like pronouns or articles, vertices 26 (ein), 27 (eine) etc. Therefore high page rank relates to basic building blocks for the language that are used frequently, can be extended to the German language as a whole.

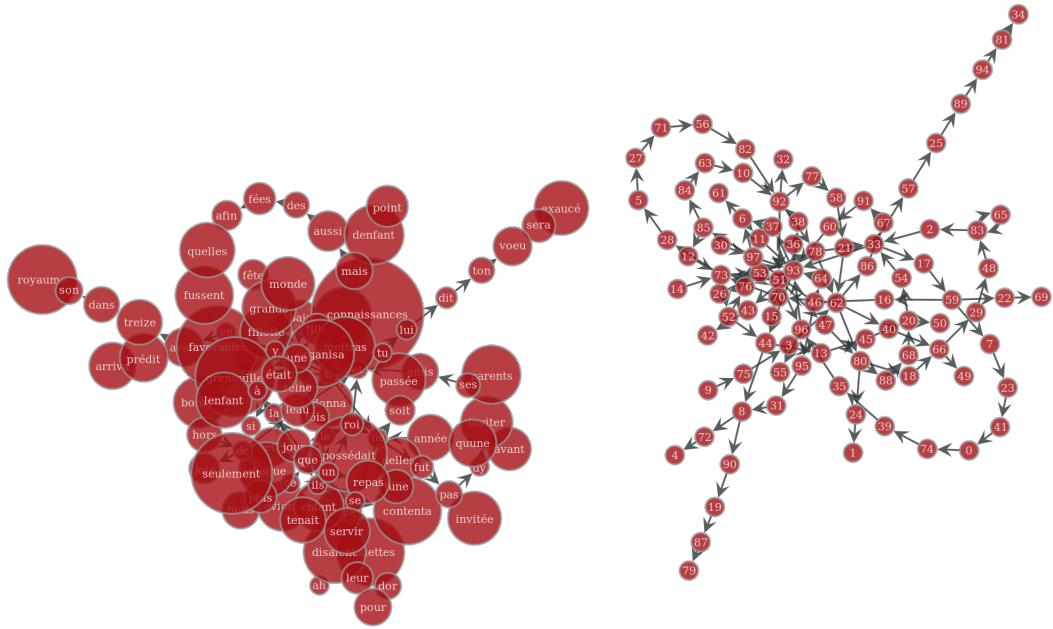
In conclusion, similar results were seen when analysing the graphical properties based on the German and English translations. With German having a clearer visualisation and stronger correlation compared to the English language due to having inflections.

### 4.3.3 French

English and German are Germanic languages under the Indo-European family. Instead of another similar language, we study a different branch under the same family. The language in question is French which lies under the Italic branch of Indo-European family. Like before with German, French contain the same parts of speech as English and also inflectional words like with German. So words can be inflected

by number (singular or plural), gender, person, case, aspect and mood. However whilst German has three genders (Masculine, Feminine and Neutral), French[26] only has the use of two (Masculine and Feminine).

The story corpus is translated into the French story corpus so graphs can be generated off this. Where the initial word graphs shown in figure 4.7.



(a) French word graph generated based on the French translation of the corpus.

(b) Numbered version of the word graph to provide a clearer graph. This is called the French numbered graph.

Figure 4.7: The French word graph and numbered equivalent of the word graph generated from the French translation of the "Sleeping Beauty" corpus.

Words of highest frequency are shown in the table extract 4.3a where the majority of words seen are inflectable and used as a way to give more information. Such as the particle of vertex 93 meaning a/an or he/it represented by vertex 46. So these would be used more frequently and it is true for the case of "Sleeping Beauty".

Vertex	Words	Counts	TC	BC	CC	LC	PR
93	une	6	4.26	10.00	1.81	0.91	10.00
62	ne	5	2.91	5.16	1.64	0.00	5.87
46	il	5	2.44	1.86	1.86	0.00	1.26
73	que	4	3.79	4.31	1.85	2.86	4.65
92	un	4	3.61	8.19	1.94	0.00	6.58
52	la	3	4.20	1.19	1.54	0.00	1.50
33	et	3	4.19	8.75	1.90	0.00	6.30
76	reine	3	4.13	1.92	1.68	10.00	2.43
51	jour	3	3.58	3.40	1.82	4.00	3.73
97	était	3	3.44	1.90	1.64	4.00	2.28

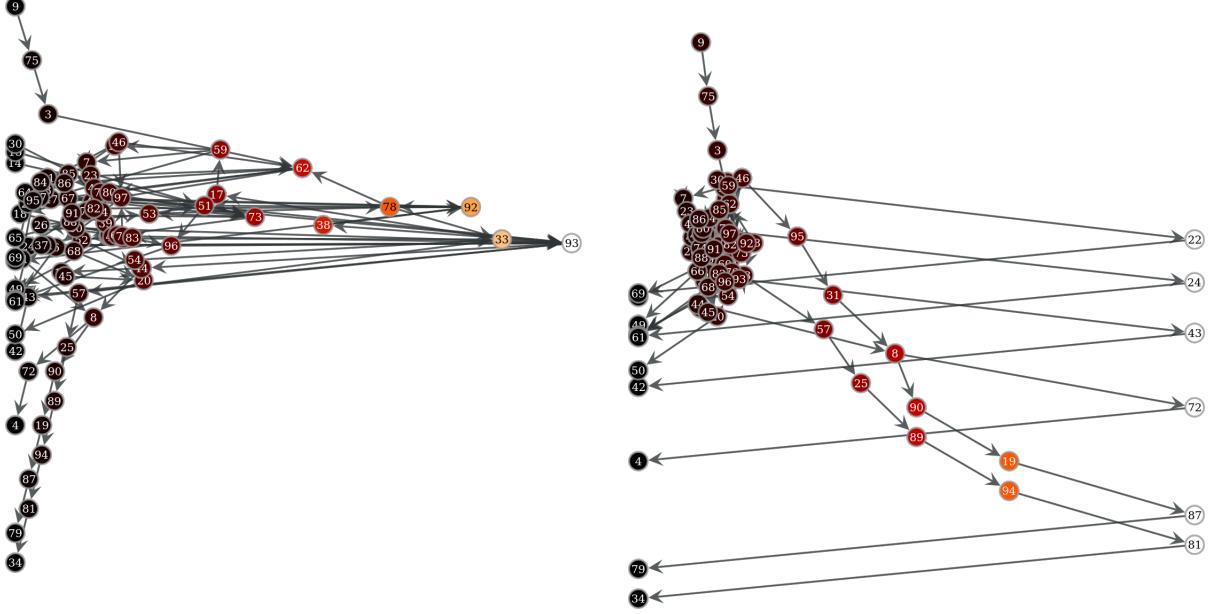
(a) Top 10 words with the highest frequency in the French translation of the corpus. Shown in table format with other graphical properties.

Vertex	Words	Counts	TC
34	exaucé	1	10.00
79	royaume	1	9.47
81	sera	1	9.03
87	son	1	8.50
94	vœu	1	8.06
4	arriva	1	7.53
19	dans	1	7.53
89	ton	1	7.09
72	prédit	1	6.56
90	treize	1	6.56

(b) Top 10 works with highest trophic levels in the French translation dataset.

Table 4.3: Partial extracts of the table data for graphical properties of the French Story Corpus.

In relation to the story corpus, French has a high trophic coherence which is calculated to be 0.35 which is higher than both English and German. Even though it may not be close to 1, compared to previous languages, a clearer level structure can be represented for this language when referencing the story corpus. Top ten trophic levels can be seen in Table ?? where the higher the value, the closer the word is to the end of a sentence as all top 10 are within the last four words in their relative sentence. However the predecessors of the words have an impact on their trophic level. If the predecessors of a word is involved in other sentences at a earlier stage then the predecessors trophic value is lower which means the original value of the word is lowered. This can be seen by vertex 49 with a trophic level of 5.08 even though it is last in it's sentence. So in general, unique sentences tend to have a higher value overall.



(a) Positions of the French numbered graph but with trophic levels and betweenness on the y-axis and x-axis respectively.

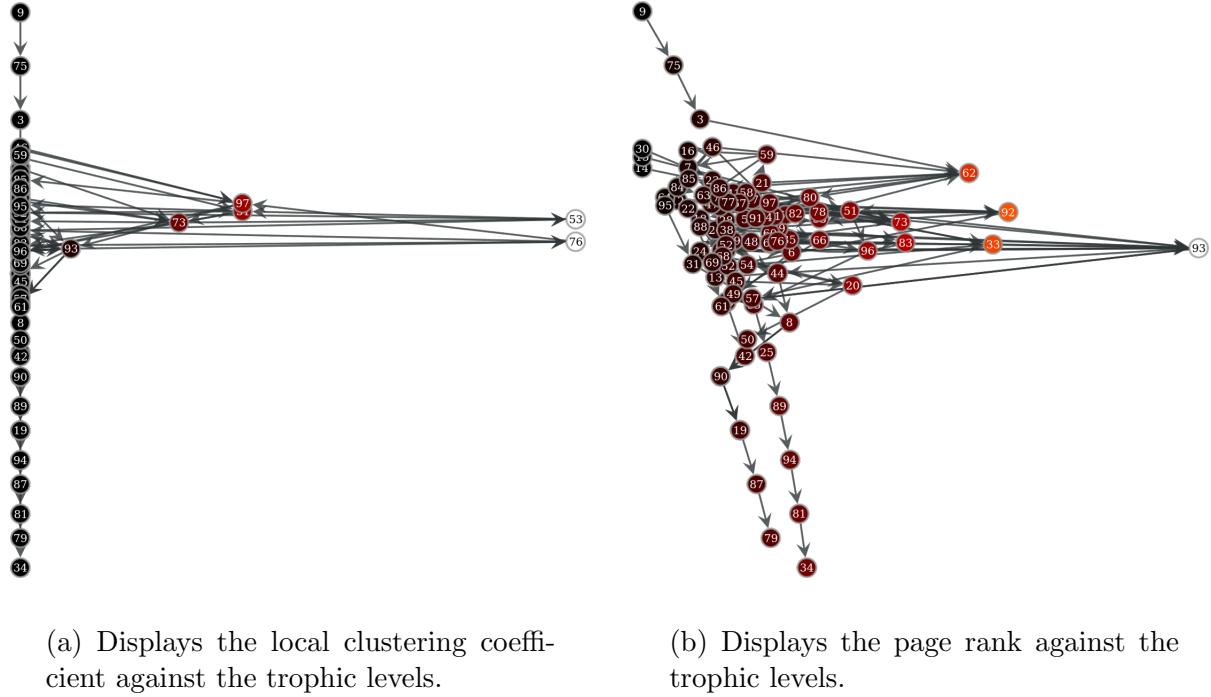
(b) Similar to the betweenness French graph but with closeness centrality values instead.

Figure 4.8: Betweenness and closeness centrality values displayed on the x-axis based on the French numbered word graph.

Vertices with high betweenness has many other vertices that must pass through it to reach other vertices within their sentence. As mentioned previously, this helps identify the key bridges in the sentences by finding the vertices which control the majority of the flow of data. Furthermore the betweenness tends to coincide with the frequency as the more frequency a word is, the more likely it has more words that require it for completion. On the other hand there are vertices such as 78 ("roi") and 38 ("fois") who have a betweenness value of 6.73 and 5.54 and a count of 2 and 1 respectively. By further analysis, the reason they have high betweenness, is that they are used with a predecessor and successor of higher betweenness value, i.e. "un roi et" and "une fois un". Meaning vertices of higher betweenness like "un", "et" and "une" must pass through "roi" or "fois" to complete their pathing. Therefore even if they have a low count, they are involved a bridges to other vertices of high betweenness.

As discovered in the German analysis, closeness identifies the vertices most likely divide the graph from the directed path, i.e. the words most likely to isolate the remaining sentence. Otherwise it means that the vertices are involved in other diverse sentences and has more in/out edges forming a conglomerate of vertices. Studying the graph in figure 4.8b, there are more vertices with high closeness and

the graph shows clearly that these vertices would cause isolation for further vertices along it's path. A side note is that French seems to be the most unique language as there are a higher number of clear divisions compared to the previous languages.



(a) Displays the local clustering coefficient against the trophic levels.

(b) Displays the page rank against the trophic levels.

Figure 4.9: Displays the local clustering and page rank on the x-axis instead of the centrality values.

Finally the local clustering and page rank. Only two vertices have a local clustering value of 10 which are vertices 53 "le" and 76 "reine". Nothing concrete can be derived from this other than as before, the vertices with a local clustering value are not an end/start of a branch. Page rank values for the French story corpus and measures the influence the vertices have on one another other than for their immediate neighbours. This is why vertex 46 "il" has a low page rank value of 1.26 as it is located at the start of the sentence so does not have enough unique vertices encompassing it to be of influence.

In conclusion, French has similar results in comparison to English and German which makes sense as they are a part of the same language family. So to see a different perspective, the study of other language families are undertaken such as the Japonic language family.

## 4.4 Japonic

The Japonic language family is the protolanguage for Japanese and Ryūkyūan languages[55] so it is a small language family compared to other language families like the Indo-European family. Translating the corpus into Japanese gives us a new dataset to use for graph property analysis which will be studied in comparison to the previous language family. Japanese translation of the story will be referred to as the Japanese corpus and the same

### 4.4.1 Japanese

To give a larger variety of languages, Japanese is studied as it uses vastly different grammar compared to other languages. Japanese has only five lexical word classes which includes nouns, verbal nouns, nominal adjectives, verbs and adjectives. The order in which the words are structured are different to Indo-European languages since Japanese uses SOV (Subject-object-verb) compared to languages like German and English which uses SVO (Subject-verb-object). As well as this, Japanese does not have major reliability on grammatical number or gender and is more focused on their system of honorifics which indicates the speaker, listener or person it references. Therefore, Japanese may have words borrowed from other languages (these are referred to as loanwords[40]) but their grammar is different in comparison.

We follow the same process as before and input the Japanese Story corpus into my program to produce the basic graphs. For simplicity, we just show the numbered version of the word graph (see Figure 4.10) with the words they reference to in the table of data (See Table ??).

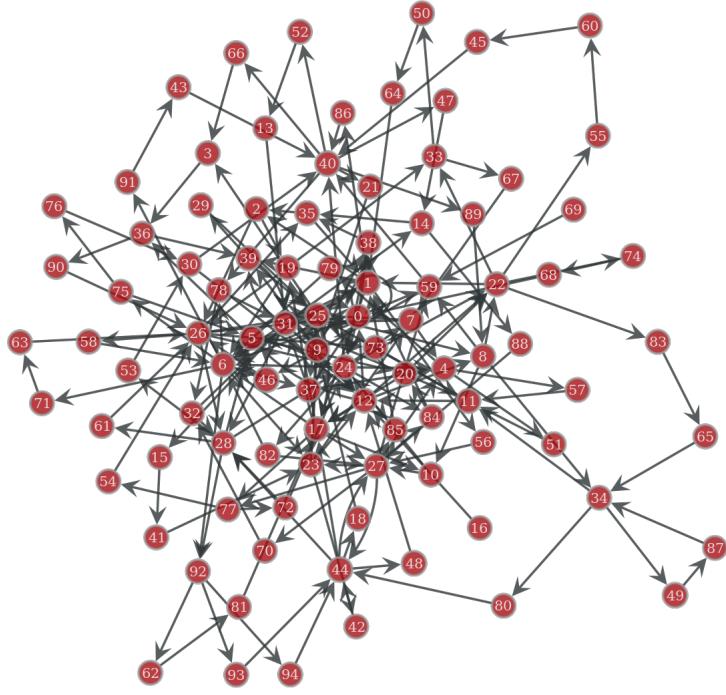


Figure 4.10: The Japanese word graph and numbered equivalent of the word graph generated from the Japanese translation of the "Sleeping Beauty" corpus.

With previous languages having a largest frequency of nine, with Japanese there are six words with eleven or more appearances (see Table 4.4a), largest being fourteen for vertex 17. This word is a form of conjugation which means that it's use depends on the inflections of it's associated words. Generally vertex 17 is most commonly used to express the past tense for the phrase in context. So rather than having a separate word for different tenses, Japanese uses accompanying words instead which leads to the higher frequency counts for said words. Furthermore, this is the case for most of the word with a large count like how vertex 12 is used to empathically add more information (an empathic and), vertex 6 is a particle so means "but" or indicated the sentence subject. Therefore, Japanese words has heavy reliability on each other to define their meaning so the graph contains a lot of edges denoting these relationships. We study other properties to see if there is other correlation.

Vertex No.	Word	Count	TL	BC	CC	LC	PR
17	たこ	14	4.82	8.12	9.50	1.06	6.12
12	し	13	3.46	8.19	8.73	0.30	6.89
31	ま	13	3.05	10.00	9.31	0.77	8.66
1	い	11	5.86	7.46	8.51	0.99	7.72
6	が	11	4.84	9.97	8.68	0.59	7.98
25	な	11	3.01	9.40	9.56	0.66	10.00
27	の	9	6.54	6.24	8.19	0.33	5.71
24	と	9	5.10	8.35	10.00	0.77	5.64
26	に	7	4.44	6.74	9.56	0.77	7.81
44	人	6	5.35	5.34	8.04	0.22	4.75

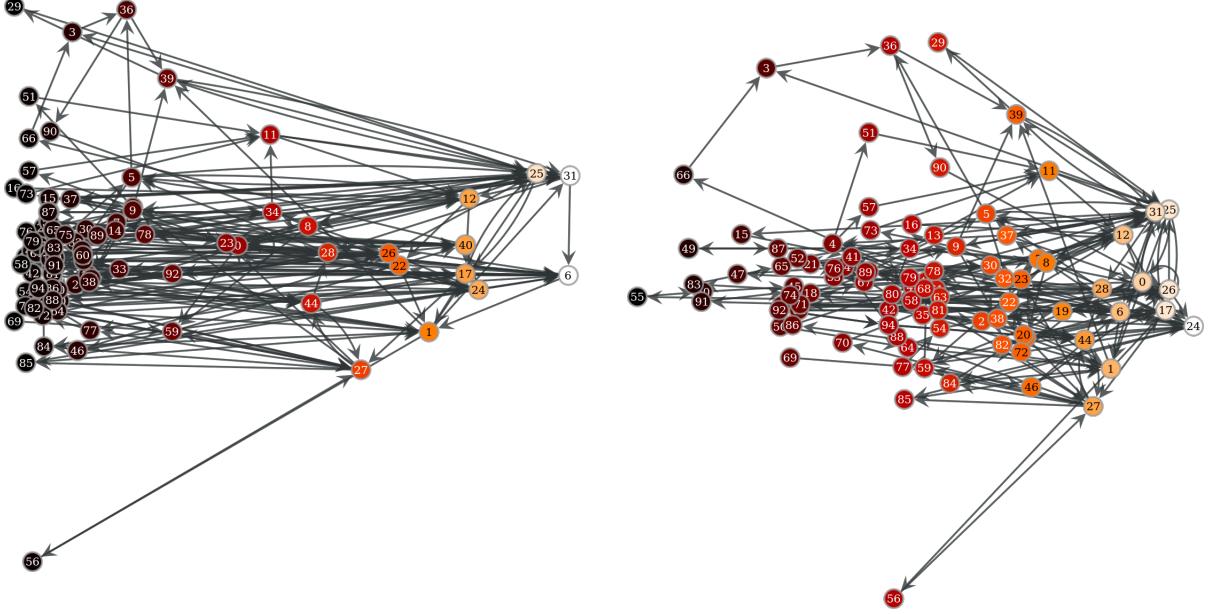
(a) Top 10 words with the highest frequency in the Japanese translation of the corpus. Shown in table format with other graphical properties.

Vertex No.	Word	Count	TL
56	女	3	10.00
27	の	9	6.54
85	賢	2	6.42
46	供	3	6.20
84	言	3	6.13
1	い	11	5.86
59	子	4	5.85
77	王	3	5.83
69	昔	1	5.67
72	様	2	5.57

(b) Top 10 works with highest trophic levels in the Japanese translation dataset.

Table 4.4: Partial extracts of the table data for graphical properties of the Japanese Story Corpus.

Trophic levels cannot be applied to this dataset like with previous languages. This is due to the graph having bidirectional edges such as the edge between vertex 25 and 5 which is part of the sentence about having only 12 plates. Also vertex 0 contains a self loop which and is part of the word "everyday", the back to back word that creates a different meaning is known as *reduplication*. So there is not a clear hierarchical structure to be inferred to. The largest trophic level is vertex 56 (refers to female) which appears in middle of sentences with the Japanese story corpus. Currently the trophic levels cannot be calculated accurately but can be if the dataset is split into specific group of words depending on it's meaning rather than character by character. For example vertex 59 and 46 means "child" together but is split in the dataset as the program doesn't recognise characters together as a singular word. However this either requires the full understanding of the Japanese language or heavy cross-referencing and checking.



(a) Positions of the Japanese numbered graph but with trophic levels and betweenness on the y-axis and x-axis respectively.

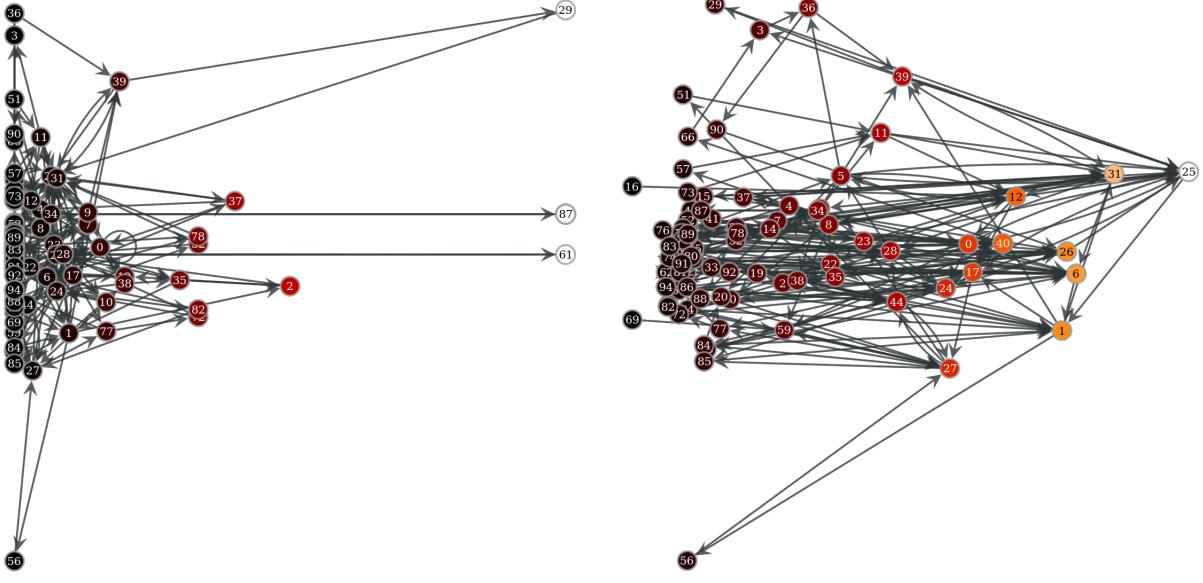
(b) Similar to the betweenness Japanese graph but with closeness centrality values instead.

Figure 4.11: Betweenness and closeness centrality values displayed on the x-axis based on the Japanese numbered word graph.

The graphs in the centrality figure 4.11 shows that the words are more evenly distributed unlike the Indo-European languages. Which correlates to the importance of the Japanese characters nature of being reliant on one another for their meaning and demonstrates the complexities of the language. However higher betweenness still correlates to the character's importance within the corpus where other characters require them to complete their meaning or structure. Although further detailed cannot be determined by the visualisation other than importance.

Japanese's closeness centrality values has a larger average than Indo-European languages. The closeness is a closer look on the betweenness centrality so all vertices in the betweenness graph is essentially shifted right if the vertices have connections to the vertices of high betweenness.

Therefore only a general correlation can be demonstrated through the centrality graphs. Clearly results can be gathered with a better reword or Japanese fluency.



(a) Displays the local clustering coefficient against the trophic levels.

(b) Displays the page rank against the trophic levels.

Figure 4.12: Displays the local clustering and page rank on the x-axis instead of the centrality values.

Local clustering (see the graph in Figure 4.12a) measures how close the vertex is part of a unique directed triangle, i.e. a directed 3 cycle. The vertices 29, 87 and 61 are vertices who are involved in such triangles where they each have a unique in and out edge to other vertices. Lower values means they are not involved in many triangles. Thus local clustering can identify certain triples in the Japanese texts that would have unique meanings.

Page rank derives the vertices who have an influence to other vertices nearby but not immediate neighbours like what closeness does. However nothing further can be correlated based on this graph (see Figure 4.12b).

In conclusion the way the Japanese story corpus is split into the database (character by character), only general correction can be determined based upon the graphs and their properties. Stronger links as to what exact words have importance cannot be obtained here without the full understanding of the language. On the other hand, we see that in general the graphs are different to the ones based on the Indo-European languages which demonstrates their differences. To see this further we analyse another language and it's language family.

## 4.5 Sino-Tibetan

Sino-Tibetan is the final language family we analyse which contains over 400 languages. The two groups of languages that are successors of this family are Tibeto-Burman language group and Chinese language group. Rather than looking at multiple, we focus on Chinese so translate the story corpus into Chinese then proceed as before. Note that Mandarin Chinese is the 2nd most spoken language after English.

### 4.5.1 Chinese

The main branches of Chinese are Mandarin or Cantonese with mandarin considered the standard and is the official language of China. Mandarin is also the language in which the corpus was translated into but will refer to this as Chinese. Chinese has a SVO (subject-verb-object) sentence structure like with English. It is made up of syllables[48] which come with three parts, an initial consonant, the tone, and a final. These three determine what syllable the word is and can be written in Pinyin, which is the standard system for romanised spelling. A simple word like "hello" has Pinyin as "nǐ hǎo" where n and h are initials, i and ao are finals and the accents on the i and a are the tone for the syllables. The tone matches the pitch of the syllable and leads to different meanings for the each word. Although each syllable in Mandarin do have meanings, they are usually combined so that words can be formed. E.g. In Pinyin, "dì" means earth but "dì tú" means map.

Similarly to Japanese, Chinese has very few inflections and use other particles which accompany the syllable to express verbal aspect. Another similarity is the use of reduplication where a syllable is repeated to give a different meaning. Parts of speech for Chinese is split into two major categories, the function words and the content words. The function words contains nouns, pronouns, verbs, auxiliary verbs, adjectives, number words, measure words, interjections and onomatopoeias. Function words contain conjunctions, prepositions and particles. Furthermore, each subgroup mentioned can also have further branches dependent on the context it is used on. For example, for nouns, there can be proper nouns, location nouns, place nouns and time nouns.

Therefore Chinese contains a larger amount of unique syllables compared to the words in the English language. Also there is a reliability factor given the context of use (like when studying Japanese). Keeping this in mind, we generate and analyse the graphs beginning with the Chinese Story Corpus's word graph seen in figure ?? (shown with numerical values which references each word in Table ??).

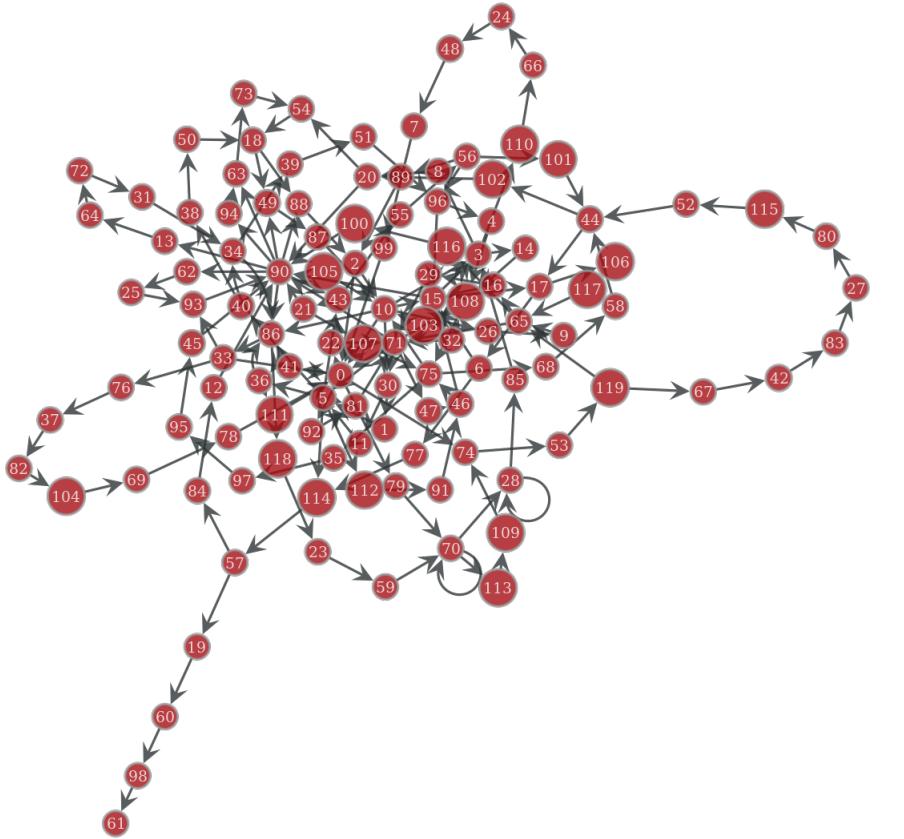


Figure 4.13: The Chinese word graph and numbered equivalent of the word graph generated from the Chinese translation of the "Sleeping Beauty" corpus.

The syllable with the higher count (see Table 4.5a) is vertex 90 which is "de" in Pinyin. This is a particle so has a grammatical meaning and different functions to denote possession or to be accompanied by other parts of speech. Other syllables with high count include vertex 71, "yǒu", vertex 5, "gè" and vertex 0, "yī". These all have various meanings depending on their use, vertex 71 can mean "have", "is" and "are", vertex 5 can mean "this", "that", party of a size or a classifier for people/objects and vertex 0 can mean "a" or part of a number. Therefore the syllables with a high count represents the words who have the most meanings and important in the sentence structure or gender related words such as vertex 15 or 43.

Vertex No.	Word	Count	TL	BC	CC	LC	PR
90	的	11	6.11	10.00	1.96	0.17	10.00
0	一	9	5.38	9.87	2.09	0.38	8.54
5	个	8	5.54	4.99	1.96	0.44	4.55
15	他	7	6.55	5.15	2.19	0.56	4.61
71	有	7	5.81	6.15	2.37	0.73	6.18
10	了	7	4.63	6.15	2.42	0.36	7.85
86	王	6	9.44	5.64	1.75	0.00	6.03
43	女	6	5.55	2.44	1.59	1.67	1.79
108	邀	4	5.24	3.02	1.76	0.67	3.90
33	后	3	10.00	3.56	1.91	0.00	1.36

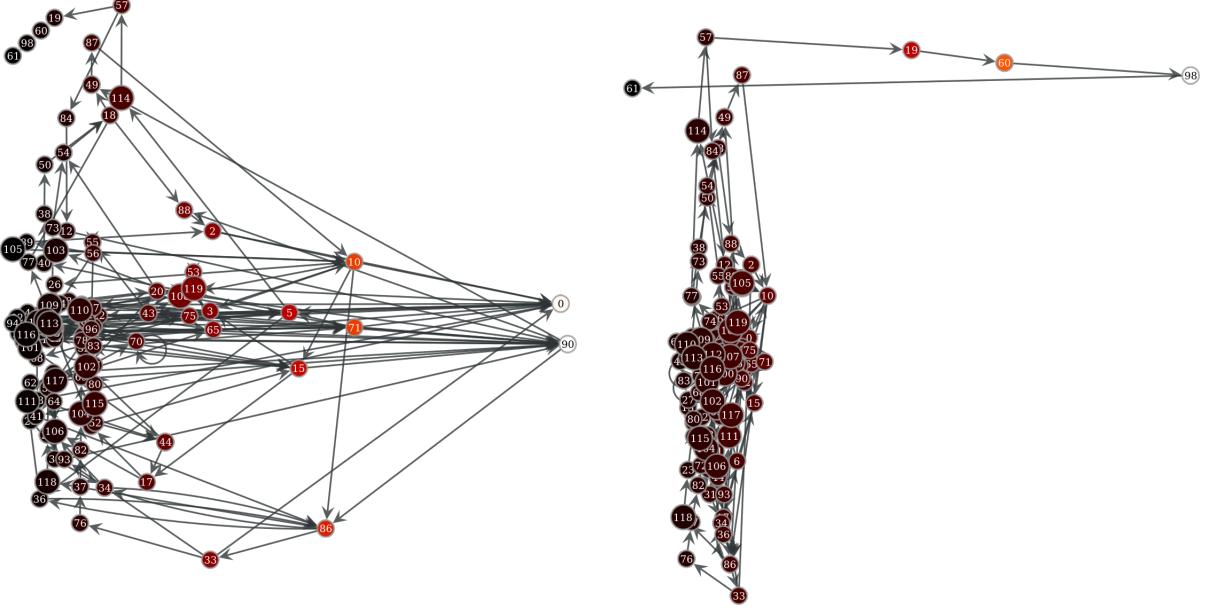
(a) Top 10 words with the highest frequency in the Chinese translation of the corpus. Shown in table format with other graphical properties.

Vertex No.	Word	Count	TL
33	后	3	10.00
86	王	6	9.44
76	正	1	9.34
36	国	3	8.90
34	和	2	8.70
37	在	1	8.68
17	们	3	8.60
118	高	1	8.59
93	真	2	8.19
31	友	1	8.18

(b) Top 10 works with highest trophic levels in the Chinese translation dataset.

Table 4.5: Partial extracts of the table data for graphical properties of the Chinese Story Corpus.

Due to bi direction edges and reduplication of syllable usage, clear hierarchical representation cannot be obtained but trophic levels are useful to give a distribution of the data based on their positioning. Thus if there are paths that continually traverse downwards, then are grammatically correct, otherwise it means that some syllables are used in various other contexts which shifts its value. For example, the path from vertex 57 to vertex 61 which can be seen in the Chinese betweenness graph in Figure 4.14a.



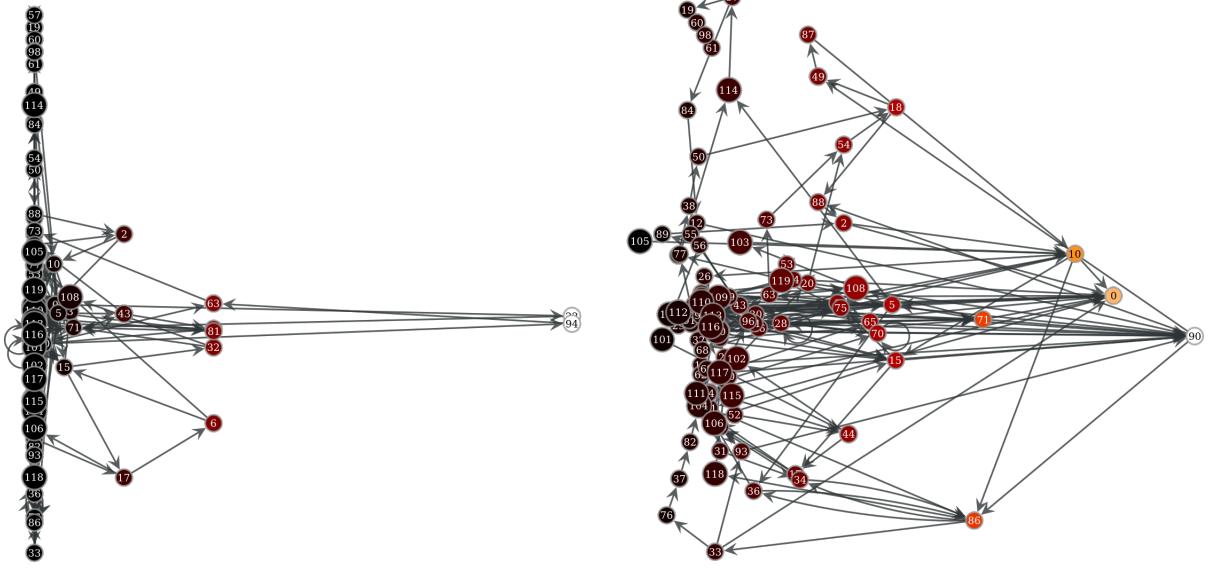
(a) Positions of the Chinese numbered graph but with trophic levels and betweenness on the y-axis and x-axis respectively.

(b) Similar to the betweenness Chinese graph but with closeness centrality values instead.

Figure 4.14: Betweenness and closeness centrality values displayed on the x-axis based on the Chinese numbered word graph.

Betweenness in the Chinese Story corpus (see figure 4.14a) are the syllables are shown to represent the particles of vertex 90 and vertex 10. Reiterating the fact that they are the syllables that act most like bridges and needed to derive further meaning. Hence betweenness relate closely to the counts of the syllables discussed earlier.

Similarly mentioned in the German analysis, closeness measures the importance of words immaterially surrounded them. Which means the vertices that are most likely to isolate further vertices in it's path, i.e. they are the vertices vital in it's path structure (controls the flow of data). Therefore, demonstrated by Figure 4.14, vertex 98 has highest closeness as without this vertex, vertex 61 is isolated. Subsequently the predecessors of vertex 98 have higher closeness. However there are no other vertices with high closeness since most other syllables have a lot more unique connections to various other vertices which infers that they have multiple meanings rather than the uniqueness of vertex 98's portion of the path.



(a) Displays the local clustering coefficient against the trophic levels.

(b) Displays the page rank against the trophic levels.

Figure 4.15: Displays the local clustering and page rank on the x-axis instead of the centrality values.

Local clustering identifies vertices involved in unique directed 3-cycles. The same idea was discovered when analysing the Japanese story corpus. For example, vertex 22 is uniquely part of the directed 3-cycle, vertex  $0 \rightarrow 22 \rightarrow 71 \rightarrow 0$  where vertex 22 is seen as the top because vertices 22 and 71 are involved with other vertices. Similarly how vertex 94 is involved in it's own directed 3-cycle. So local clustering doesn't identify the vertices of high importance but rather certain triples within the corpus.

The page rank identifies the vertices with influence to vertices other than their immediate neighbours. Evidently vertices with high betweenness will also have a high page rank as they are more commonly used in the corpus.

# **Chapter 5**

## **Conclusion**

Based upon the five languages shown in the previous chapter. The results demonstrate the similarities of languages that are part of the same family as well as their uniqueness and the possible causes for their variety in graph property value.

### **5.1 Final Correlations**

Based on each graph property that was studied, we take their correlations within the story corpus used and summarise.

#### **5.1.1 Local Clustering Coefficient**

#### **5.1.2 Betweenness Centrality**

#### **5.1.3 Closeness Centrality**

#### **5.1.4 Page Rank**

#### **5.1.5 Trophic Level**

### **5.2 Further works and Applications**

# Bibliography

- [1] Maristella Agosti and Luca Pretto. A theoretical study of a generalized version of kleinberg’s hits algorithm. *Information Retrieval*, 8(2):219–243, 2005.
- [2] Alex Bavelas. A mathematical model for group structures. *Human organization*, 7(3):16–30, 1948.
- [3] Christian Bentz, Douwe Kiela, Feli Hill, and Paula Butterly. Zipf’s law and the grammar of languages: A quantitative study of old and modern english parallel texts. *Corpus Linguistics and Linguistic Theory*, 10(2):175–211, 2014.
- [4] Kyle Bojanek, Yuqing Zhu, and Jason Maclean. Cyclic transitions between higher order motifs underlie sustained asynchronous spiking in sparse recurrent networks. *PLoS computational biology*, 16:e1007409, 09 2020.
- [5] Ulrik Brandes, Stephen P Borgatti, and Linton C Freeman. Maintaining the duality of closeness and betweenness centrality. *Social networks*, 44:153–159, 2016.
- [6] Ulrik Brandes and Christian Pich. Centrality estimation in large networks. *International Journal of Bifurcation and Chaos*, 17(07):2303–2318, 2007.
- [7] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems*, 30(1-7):107–117, 1998.
- [8] Lyle Campbell. Language isolates and their history, or, what’s weird, anyway? In *Annual Meeting of the Berkeley Linguistics Society*, volume 36, pages 16–31, 2010.
- [9] Lyle Campbell. How many language families are there in the world? *Anuario del Seminario de Filología Vasca” Julio de Urquijo”*, 52(1/2):133–152, 2018.
- [10] Stephen C. Carlson. graph theory. Dec 2022.
- [11] Alexander Chatzigeorgiou, Nikolaos Tsantalis, and George Stephanides. Application of graph theory to oo software engineering. In *Proceedings of the 2006 international workshop on Workshop on interdisciplinary software engineering research*, pages 29–36, 2006.

- [12] G.P. Clemente and R. Grassi. Directed clustering in weighted networks: A new perspective. *Chaos, Solitons & Fractals*, 107:26–38, 2018.
- [13] BNC Consortium et al. British national corpus. *Oxford Text Archive Core Collection*, 2007.
- [14] Ángel Corberán. The chinese postman problem with load-dependent costs. *Transportation Science*, 52(2):370–386, March 2018.
- [15] Silvia de Juan, Andres Ospina-Alvarez, Sebastián Villasante, and Ana Ruiz-Frau. A graph theory approach to assess nature’s contribution to people at a global scale. *Scientific Reports*, 11(1):9118, 2021.
- [16] Pooja Devi, Ashlesha Gupta, and Ashutosh Dixit. Comparative study of hits and pagerank link based ranking algorithms. *International Journal of Advanced Research in Computer and Communication Engineering*, 3(2):5749–5754, 2014.
- [17] Martin Durrell. *Hammer’s German grammar and usage*. Routledge, 2011.
- [18] David M. Eberhard, Gary F. Simons, and Charles D. Fennig. Ethnologue: Languages of the world.
- [19] Giorgio Fagiolo. Clustering in complex directed networks. *Physical Review E*, 76(2), aug 2007.
- [20] Tanguy Fardet and Anna Levina. Weighted directed clustering: Interpretations and requirements for heterogeneous, inferred, and measured networks. *Phys. Rev. Res.*, 3:043124, Nov 2021.
- [21] Linton C Freeman. A set of measures of centrality based on betweenness. *Sociometry*, pages 35–41, 1977.
- [22] Linton C Freeman et al. Centrality in social networks: Conceptual clarification. *Social network: critical concepts in sociology*. Londres: Routledge, 1:238–263, 2002.
- [23] Linton C Freeman, Douglas Roeder, and Robert R Mulholland. Centrality in social networks: II. experimental results. *Social networks*, 2(2):119–141, 1979.
- [24] J. Grimm and W. Grimm. *Kinder und haussmärchen: gesammelt durch die Brüder Grimm*. Number v. 1-2 in Kinder und haussmärchen: gesammelt durch die Brüder Grimm. Dieterich, 1857.
- [25] Michael Hart, R.J.F. Ypma, Rafael Romero-García, Stephen Price, and John Suckling. Graph theory analysis of complex brain networks: New concepts in brain mapping applied to neurosurgery. *J. Neurosurg.*, 124:1665–1678, 06 2016.

- [26] Roger Hawkins and Richard Towell. *French grammar and usage*. Routledge, 2015.
- [27] William L. Hosch. Zipf's law. 2009.
- [28] Math Insight. The 16 possible network motifs involving three nodes. [Online; accessed 8 February, 2023].
- [29] Stefan Irnich. Undirected postman problems with zigzagging option: A cutting-plane approach. *Computers & Operations Research*, 35(12):3998–4010, December 2008.
- [30] Sergio Jimenez. Zipf 30wiki en labels, 2015.
- [31] Samuel Johnson. Digraphs are different: Why directionality matters in complex systems. *Journal of Physics: Complexity*, 1(1):015003, 2020.
- [32] Samuel Johnson, Virginia Domínguez-García, Luca Donetti, and Miguel A Munoz. Trophic coherence determines food-web stability. *Proceedings of the National Academy of Sciences*, 111(50):17923–17928, 2014.
- [33] Samuel Johnson and Nick S Jones. Looplessness in networks is linked to trophic coherence. *Proceedings of the National Academy of Sciences*, 114(22):5618–5623, 2017.
- [34] Leo Katz. A new status index derived from sociometric analysis. *Psychometrika*, 18(1):39–43, 1953.
- [35] Sukany Khaisaeng and NK Dennis. A study of part of speech used in online student weekly magazine. *International Journal of Journal of Research Granthaalayah*, 5(4):43–50, 2017.
- [36] U. Knauer. *Algebraic graph theory : morphisms, monoids, and matrices / by Ulrich Knauer*. De Gruyter studies in mathematics ; 41. 2011.
- [37] Amy N Langville and Carl D Meyer. A survey of eigenvector methods for web information retrieval. *SIAM review*, 47(1):135–161, 2005.
- [38] Chin Ho Lee. Maximum bipartite matching to max flow, 2009. [Online; accessed 5 February, 2023].
- [39] Geoffrey Leech, Paul Rayson, et al. *Word frequencies in written and spoken English: Based on the British National Corpus*. Routledge, 2014.
- [40] Akira Miura. The influence of english on japanese grammar. *The Journal of the Association of Teachers of Japanese*, 14(1):3–30, 1979.

- [41] M. E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45(2):167–256, jan 2003.
- [42] Koji Ohnishi. Towards a brief proof of the four-color theorem without using a computer: theorems to be used for proving the four-color theorem. *Artificial Life and Robotics*, 14(4):551–556, Dec 2009.
- [43] Scientific Figure on ResearchGate. Cyclic transitions between higher order motifs underlie sustained asynchronous spiking in sparse recurrent networks, 2020. [Online; accessed 12 February, 2023].
- [44] Tore Opsahl and Pietro Panzarasa. Clustering in weighted networks. *Social networks*, 31(2):155–163, 2009.
- [45] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.
- [46] James Powell and Matthew Hopkins. *A Librarian’s Guide to Graphs, Data and the Semantic Web*. Chandos Information Professional Series. Chandos Publishing, 2015.
- [47] Liudmila Ostroumova Prokhorenkova and Egor Samosvat. Global clustering coefficient in scale-free networks, 2014.
- [48] Claudia Ross and Jing-heng Sheng Ma. *Modern Mandarin Chinese grammar: A practical guide*. Routledge, 2017.
- [49] Bruce M Rowe and Diane P Levine. *A concise introduction to linguistics*. Taylor & Francis, 2022.
- [50] Thomas Schank and Dorothea Wagner. Approximating clustering coefficient and transitivity. *Journal of Graph Algorithms and Applications*, 9(2):265–275, 2005.
- [51] School of Mathematics and Statistics, University of St Andrews, Scotland. Königsberg bridges, 2000. [Online; accessed 2 February, 2023].
- [52] Alfonso Shimbel. Structural parameters of communication networks. *The bulletin of mathematical biophysics*, 15:501–507, 1953.
- [53] Elvira I Sicilia-Garcia, Ji Ming, Francis J Smith, et al. Extension of zipf’s law to words and phrases. In *COLING 2002: The 19th International Conference on Computational Linguistics*, 2002.
- [54] Katie Steckles. The bridges of königsberg, 2018. [Online; accessed 2 February, 2023].

- [55] Alexander Vovin. Origins of the japanese language. In *Oxford Research Encyclopedia of Linguistics*. 2017.
- [56] Duncan J. Watts and Steven H. Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 393(6684):440–442, Jun 1998.
- [57] Douglas R White and Stephen P Borgatti. Betweenness centrality measures for directed graphs. *Social networks*, 16(4):335–346, 1994.
- [58] Robin Wilson. Four colours suffice, 2017. [Online; accessed 3 February, 2023].
- [59] Wenpu Xing and Ali Ghorbani. Weighted pagerank algorithm. pages 305–314, 2004.
- [60] Wayne W. Zachary. Modeling social network processes using constrained flow representations. *Social Networks*, 6(3):259–292, 1984.
- [61] Bin Zhang and Steve Horvath. A general framework for weighted gene co-expression network analysis. *Statistical Applications in Genetics and Molecular Biology*, 4(1), 2005.

# **Appendix A**

## **Karate Club**

### **A.1 Karate Club Adjacency Matrix**

This is the adjacency matrix generated for the karate club dataset.



# Appendix B

## Langauages

### B.1 Text Corpus

In times past there lived a king and queen, who said to each other every day of their lives, "Would that we had a child!" and yet they had none. But it happened once that when the queen was bathing, there came a frog out of the water, and he squatted on the ground, and said to her: "Thy wish shall be fulfilled; before a year has gone by, thou shalt bring a daughter into the world."

And as the frog foretold, so it happened; and the queen bore a daughter so beautiful that the king could not contain himself for joy, and he ordained a great feast. Not only did he bid to it his relations, friends, and acquaintances, but also the wise women, that they might be kind and favourable to the child. There were thirteen of them in his kingdom, but as he had only provided twelve golden plates for them to eat from, one of them had to be left out.

### B.2 English Language Table

Entire table of graph property values for the english word graph that was generated from the first two paragraphs of the story "Sleeping Beauty".

Vertex	Words	Counts	TC	BC	CC	LC	PR
0	a	7	1.98	9.26	2.22	0.00	7.65
1	acquaintances	1	2.53	0.60	1.64	0.00	0.36
2	also	1	3.20	0.24	1.82	0.00	0.16
3	and	9	3.00	10.00	2.49	0.57	8.13
4	as	2	2.95	0.73	2.16	3.33	1.13
5	bathing	1	1.24	1.56	1.51	0.00	0.33
6	be	3	3.78	3.15	1.91	0.00	3.24
7	beautiful	1	3.29	0.77	1.81	0.00	0.42
8	before	1	2.58	0.83	1.84	0.00	0.74
9	bid	1	3.11	0.38	1.88	0.00	0.11
10	bore	1	2.43	0.73	1.83	0.00	0.10
11	bring	1	1.98	1.96	1.84	0.00	1.24
12	but	3	2.05	1.29	1.95	0.00	1.73
13	by	1	1.98	1.91	1.22	0.00	1.00
14	came	1	1.20	0.51	1.84	0.00	0.12
15	child	2	7.17	0.00	0.00	0.00	1.78
16	contain	1	0.78	1.13	1.55	0.00	0.12
17	could	1	1.33	2.12	1.32	0.00	0.57
18	daughter	2	1.19	1.44	1.70	0.00	0.61
19	day	1	4.91	1.00	1.94	0.00	0.82
20	did	1	2.16	0.73	1.72	0.00	0.30
21	each	1	4.08	0.94	1.25	0.00	0.31
22	eat	1	4.15	0.70	1.42	0.00	0.31
23	every	1	4.63	0.98	1.64	0.00	0.68
24	favourable	1	3.40	1.16	1.87	0.00	0.36
25	feast	1	10.00	0.00	0.00	0.00	0.76
26	for	2	2.35	2.57	2.17	0.00	2.09
27	foretold	1	3.39	0.52	1.48	0.00	0.57
28	friends	1	2.71	0.67	2.02	0.00	0.57
29	frog	2	3.93	2.62	1.72	0.00	1.78
30	from	1	4.49	0.72	1.64	0.00	0.51
31	fulfilled	1	3.18	0.82	1.57	0.00	0.59

32	golden	1	2.17	1.05	1.55	0.00	0.67
33	gone	1	1.98	1.89	1.09	0.00	0.89
34	great	1	5.99	0.32	10.00	0.00	0.61
35	ground	1	3.67	0.60	2.00	10.00	0.55
36	had	4	3.44	5.32	2.48	0.36	3.82
37	happened	2	2.36	1.07	2.07	0.00	0.72
38	has	1	1.98	1.87	0.99	0.00	0.76
39	he	4	2.42	3.68	2.07	0.48	2.38
40	her	1	3.79	0.94	1.12	0.00	0.31
41	himself	1	1.57	1.15	1.81	0.00	0.35
42	his	2	2.13	1.29	1.52	0.00	1.33
43	in	2	0.49	1.29	1.50	0.00	0.44
44	into	1	2.77	0.36	1.82	0.00	0.07
45	it	3	3.52	2.05	1.79	0.00	2.12
46	joy	1	2.67	0.78	2.02	0.00	0.70
47	kind	1	3.39	1.38	2.01	0.00	0.59
48	king	2	2.66	5.79	2.10	1.67	1.78
49	kingdom	1	2.09	0.33	1.64	0.00	0.38
50	left	1	4.90	0.32	1.63	0.00	0.59
51	lived	1	1.20	0.51	1.84	0.00	0.12
52	lives	1	4.46	0.73	1.56	0.00	1.42
53	might	1	3.65	0.28	1.61	0.00	0.60
54	none	1	7.45	0.00	0.00	0.00	0.41
55	not	2	0.00	2.13	1.51	0.00	0.73
56	of	4	5.19	6.21	2.37	0.00	6.03
57	on	1	3.70	0.41	1.82	0.00	0.33
58	once	1	3.05	0.18	1.81	0.00	0.12
59	one	1	4.84	0.74	1.94	0.00	0.68
60	only	2	1.90	2.05	1.55	0.00	1.15
61	ordained	1	2.20	0.53	1.83	0.00	0.11
62	other	1	4.35	0.96	1.42	0.00	0.51
63	out	2	6.01	2.45	1.93	0.00	1.93
64	past	1	0.44	0.71	1.51	0.00	0.24

65	plates	1	2.26	1.06	1.81	0.00	0.81
66	provided	1	1.99	1.01	1.20	0.00	0.30
67	queen	3	2.89	3.09	1.88	1.00	1.53
68	relations	1	2.42	0.65	1.70	0.00	0.38
69	said	2	5.62	1.65	1.87	0.00	1.30
70	shall	1	3.79	1.00	1.62	0.00	0.82
71	shalt	1	1.98	1.94	1.57	0.00	1.18
72	so	2	2.85	1.63	1.72	0.00	1.42
73	squatted	1	3.06	0.39	1.55	0.00	0.11
74	that	4	3.73	3.49	2.19	0.36	5.12
75	the	9	4.34	9.19	2.20	0.44	10.00
76	their	1	4.83	0.71	1.36	0.00	1.38
77	them	3	2.90	5.03	2.40	1.00	2.70
78	there	3	0.41	2.63	1.76	0.00	1.59
79	they	2	3.52	1.48	2.05	0.00	1.85
80	thirteen	1	3.60	1.00	1.94	0.00	0.35
81	thou	1	1.98	1.93	1.37	0.00	1.10
82	thy	1	3.79	0.96	1.25	0.00	0.51
83	times	1	0.47	0.70	1.32	0.00	0.00
84	to	6	3.80	7.48	2.28	0.18	5.54
85	twelve	1	2.08	1.03	1.35	0.00	0.50
86	was	1	2.06	1.54	1.32	0.00	0.10
87	water	1	3.67	0.60	2.00	10.00	0.55
88	we	1	3.59	0.64	2.01	0.00	0.69
89	were	1	2.01	0.98	1.65	0.00	0.12
90	when	1	4.03	0.00	1.82	10.00	0.69
91	who	1	4.26	0.19	1.59	0.00	0.10
92	wise	1	4.14	0.75	1.54	0.00	0.55
93	wish	1	3.79	0.98	1.41	0.00	0.68
94	women	1	3.93	0.76	1.80	0.00	0.72
95	world	1	8.35	0.00	0.00	0.00	0.55
96	would	1	4.10	0.75	1.82	0.00	1.45
97	year	1	1.98	1.86	0.91	0.00	0.61
98	yet	1	3.26	0.48	1.72	0.00	0.36