**Intro and Hypothesis**

The world wide web contains a network of web pages that are interconnected in many different ways. Wikipedia itself is its own network of a vast collection of different topics, ideas, and other entries. Our goal was to find an interesting hypothesis to analyze that dealt with this network of pages within Wikipedia. An initial path of thought was that almost all Wikipedia pages are connected in some way. This means that if each Wikipedia page was a node on a graph, almost the graph would be connected if the edges between nodes were links from one Wikipedia page to another. An even more specific idea was that all pages on Wikipedia formed a connected graph if one just uses edges only for the first accessible in-text link on each Wikipedia page.

After some further searching, we discovered a hypothesis centered around the idea that the Philosophy page is reachable by almost any given starting Wikipedia page by only traversing through pages by clicking the first in-text link on each page. It seems intuitive that the Philosophy page could potentially be reached by almost any starting page from any link on the page, but only accessing the first link to do the same does not seem as easily possible. When researching this further we came up with the following hypothesis: In Wikipedia it is theoretically possible to reach the Philosophy page by just clicking on the first link on every page from **97%** of all Wikipedia pages. We also found estimates suggesting that it takes a **median of 23 clicks** to get from a starting page to the Philosophy page. We decided to test this hypothesis by going through a representative sample of the (approximately) 29 million English Wikipedia pages and storing their paths to Philosophy in graph form. We then wanted to check the median path size after starting at a random selection of pages. We also wanted to look at what we have

designated as "funnels", which are nodes that more than one inward edge passing through them towards Philosophy.

We decided to further develop this project in the implementation side by building a basic UI that would allow users to choose any exit node (Philosophy in our study) and then find paths from any source node to this exit node as well as path lengths. The UI would also let users find funnels after running multiple paths and also generate a graph (with visuals) of all paths generated. A final feature of the UI would be the ability for a user to generate a random graph in a similar way to the analysis we were performing for this project.

**Implementation**

In order to run the analysis, we first developed a "Graph" class which uses nodes from a "Node" class to represent a graph of nodes in an adjacency list form. We then built a "Scraper" class which would go through each Wikipedia page starting from a source node. At each page the scraper finds the first link to another Wikipedia page that is in the text body on the page (not the sidebar information). The scraper then accesses the next page through that link and adds it to the graph as a new node. This process continues until reaching the Philosophy page and the graph is populated with all the nodes found along the path from the source node to Philosophy. We also incorporated checks to stop when cycles are found if a source does not lead to Philosophy. We also wanted to be able to see the graphs produced by our analysis, so we implemented a graph visualization tool using a Java library called GraphStream. A feature was also included that would allow us to output the adjacency list representation of the graph, as well as the list of funnels, as a text file. Another thing we wanted to implement was a way to find

distances between any node and philosophy, not just source nodes from created paths. To achieve this we implemented a BFS algorithm that would work with our graph.

For the UI implementation we first created a window where users could specify the end node by providing a link to a Wikipedia page. Then, the user is able to submit queries of source Wikipedia page links. When running the tool, the user UI will display all previous queries as well as the length of their respective paths to the end node. During each query, the entire path is displayed until the next query is run. We also implemented a way for the user to find the "funnels" after running as many queries as desired by opening a new window that displays the funnels and the number of inward edges each funnel has. The graph display feature was also added that will display a visual of the graph containing the paths from all queries to the exit node. The random graph button on the UI was implemented to mimic our analysis that we conducted by generating a random graph from a random selection of Wikipedia pages. This was achieved by accessing the built in Wikipedia link on the website that redirects to a random page.

**Analysis Setup**

We ran some early trial runs with the code we implemented to make sure everything was working the way we wanted. Initially, we found that every page was following the same path to Philosophy with each path length around 28-29 nodes traversed. After further investigation we realized that every page was getting directed through a path that contained a language of origin such as Greek, Latin, or French. This was not ideal for our analysis as each source node ended up having essentially the same path through a page like "Greek" or "Ancient Greek Language". This problem made sense because we were using only pages that contained English words, and

most words from the English language have an origin in languages like Greek or Latin. This problem was further emphasized by the Wikipedia page structure in which any pages with terms that were not proper nouns had the language of origin for the word or terms in parentheses at the beginning of the text body. Due to the prominence of this issue and how it made the analysis uninteresting, we decided to ignore these links for languages of origin and they were removed from our scraper. This meant the scraper would instead go to the next link outside of the initial parentheses containing language of origin in each Wikipedia page.

As stated above, we utilized the Wikipedia feature of a random page redirect link to find random pages to use as source nodes. We decided to run 10 tests of 50 pages each that would create a new graph each time. This gave us 500 randomized source node samples and also made sure the program would not take too long to run or send too many requests in a short period of time to Wikipedia, also ensuring that Wikipedia would not block our computers from accessing the website.

**Analysis**

After running the analysis as specified above, the data was put into an excel file (see attached "analysis data" excel file). First, we took a look at all the distances from the source node to Philosophy and found that some had distances of zero. These Wikipedia links were checked to verify that these zeros were not due to errors. Nodes with distances of less than 5 were also checked at their respective Wikipedia links and verified, changed to zero, or removed. After cleaning the distance data, only 2 source pages were removed, resulting in a total of 498 samples. The funnel data is recorded in a different tab in the same excel sheet. Each of the 10 graphs

produced a list of funnels as well as their number of inward edges towards Philosophy. The 10 graphs from the analysis trials are included at the end of the report.

The final analysis results had a mean distance of about 17 nodes, **median distance of 17 nodes**, a max distance of 41, standard deviation of around 6.2, and a resulting **97.8%** of source nodes having a successful path to Philosophy.

**Discussion**

When looking at the pages with zero distances, it was interesting to note that they usually fell into two categories: surnames or some type of list. Both of these types of pages are represented on Wikipedia in such a way that the initial paragraph of text that we counted for the first link either does not really exist or does not have any links contained in it. The other pages that followed this result for a zero distance usually were names of people that had short pages with little to no text in the initial paragraph. One of the more interesting sources with a zero distance was the page of Billie Jean King's career statistics (https://en.wikipedia.org/wiki/Billie_Jean_King_career_statistics) which ended up producing a cycle given our method of taking the first link. The reason this happens is that Billie Jean King's page leads to the world number one women's tennis rankings page, leading to the Women's Tennis Association page. Coincidentally enough, Billie Jean King happens to be the founder of the Women's Tennis Association where the first link on the page redirects back to Billie Jean King, creating the cycle. This was the only cycle we found in our data, indicating that these such cycles are rare (estimated at about 1/500 or around 0.2% given our data).

The next point of interest are the funnels found in each graph. The trend in the funnels are usually general words of subject such as Science, Mathematics, Physics, Knowledge, etc. This makes sense because pages that would have common paths would have a higher chance of being similar if they had similar general underlying concepts, such as Basketball leading to Game Theory, Probability or Math, and something like Financial Markets also leading to Probability or Math. There are a few more specific funnels that appeared, such as Sport, Association Football, Rock Music, and Canada, but these are not as important as they are due to the random page search producing two very similar pages in the same category. These funnels also tend to exist much farther apart from the rest of the funnels and also farther from Philosophy when looking at the graphs. What is most significant is that Physics, Science, Psychology, Language, Linguistics, and Knowledge are in almost every graph as a funnel. The explanation behind this is intuitive, as eventual any source page will get to some sort of topic that boils down to a these general topics that are pillars of subject matter. These funnels are also more concentrated together in each graph and tend to appear closer to Philosophy, indicating that they are more significant than the other funnels.

Something that is noticeable in almost every graph is that there is a path through Philosophy that doesn't usually contain any paths branching away. On every graph, Philosophy appears to partition the graph into two distinct parts if it were removed. A further look reveals that one of the parts is consistently much larger than the other. If you path is traversed outward from Philosophy towards the larger part, you always end up at the Knowledge funnel. This means that Knowledge could also be used as an exit node in an analysis and would most likely produce similar results to using Philosophy as the exit node, although it may be slightly less

successful if the source nodes cut off from removing Philosophy would not have a path to Knowledge. Knowledge being an important node/funnel also is plausible in this graph given that the subjects of Knowledge and Philosophy are both closely related.

Finally, the hypothesis. Based on the results of our analysis the hypothesis is supported. We found a slightly higher percentage than the hypothesis for the amount of pages that would have a successful path. The median is also actually lower, which is actually better than the hypothesis, so it is safe to say the median from the hypothesis is at least supported in the worst case. While we only ran this with 500 samples, we would expect running our analysis on more pages would also support the hypothesis.

**Conclusion**

We conclude that our original hypothesis of 97% of Wikipedia pages having a path to Philosophy by clicking the first link, with a median path length of 23 clicks, is supported (actually better in the case of median clicks equal to 17). It is interesting that nearly all pages have some path to Philosophy in Wikipedia following these rules, and that there are similar pages that could produce the same effect, like Knowledge. We also built a simplified UI that will allow users to run our analysis on random Wikipedia pages, as well as test this concept out on any page the choose, even for different exit nodes.

# Graphs from Analysis

## (Red Node is Philosophy, Blue Nodes are Funnels)



Graph 1

Graph 2

Graph 3

Graph 4

Hijidai Station
Dosan Line
Shikoku
Japanese archipelago
Japan
Japanese language
Palau
Palauan language
Abdul Aziz bin Musaid
Saudi Arabia
August Eisenmenger
Minosov (Rokycany District)
Austrians
German language
Arabic
West Germanic languages
Central Semitic languages
Germanic languages
Semitic languages
Austric languages
Leonid Koltun
Semitic languages
NYK Virgo
Afroasiatic languages
Container ship
Language family
Cargo ship
Austroasiatic languages
Merchant ship
Watercraft
Vietnamese language
Vehicle
Vietnamese people
Sino-Tibetan language
Communication
Meaning (semiotics)
Pernille Nedergaard
Badminton
Racket (sports equipment)
Semiotics
Ball
Sphere
Geometry
Mathematics
Quantity
Counting
Element (mathematics)
Semiosis
Rudelmar Bueno de Faria
ACT Alliance
Set (mathematics)
Alliance
People
Mathematical object
Person
Abstract and concrete
Action (philosophy)
Reason
Object (philosophy)
Consciousness
Banhua
Quality (philosophy)
Printmaking
Philosophy
Pär Edwardson
Ontology
Livingston
Frölunda HC
Västra Frölunda
New York (state)
Communes of France
U.S. state
Administrative division
Red Bluff Creek
Existence
Francis Goes to the Pecos
Country
Black and white
Astronomical object
Political geography
Politics
Physical body
Reality
Continuous spectrum
List of Ultratop 40 number
Physics
Natural science
Sound recording and reproduction
Psychology
Fact
Measurement
Accounting
Branches of science
Financial accounting
Science
Biology
Asset
Finance
Linguistics
Species
Turbo magnificus
Investment fund
Data (computing)
Mass noun
Taxonomy (biology)
Knowledge
Software
Free and open
Vorbis
Class (biology)
Help:Media
Geography
Indian subcontinent
Entognatha
Hungarian language
Region
Human settlement
Hexapoda
Doba, Satu Mare
Kashmir
Exonym and endonym
Insect
Jammu and Kashmir
Moth
Ladakh
Persian language
Village
Pigritia gruis
Yapola River
Saeid Marouf
Discipline (academia)
Media studies
Sądcza, Lubusz Voivodeship
Social science
Marshall McLuhan
Political science
System
Legitimacy (political)
Social group
Authority
Mathematical model
Executive (government)
Society
Game theory
Public policy
Sociolinguistics
Zero
Public administration
Sociology
Competition
Seth McClung
Local government
Variety (linguistics)
Sport
Bat
Municipal corporation
Standard language
Social norm
Golf
Cricket
List of United States
Vulgar Latin
English
Technical standard
Contact sport
Motorsport
First
New York City
American languages
Community
International standard
Combat sport
Team sport
Charles Flood
Boroughs of New York
Eastern Romance languages
Behavior
Literary award
International standard
Mixed martial arts
Grand Prix motorcycle racing
Brooklyn
Romanian language
English
Hans Christian Andersen Award
Matt Dwyer
American football
Naomi Taniguchi
Tom Murphy (athlete)
The Unsaved
Sunčana Škrinjarić
International Organization for Standardization
Association football
Danish 3rd Division
Rule
Model (person)
Rebecca Jackson Mendoza
List of romanizations
National Football League
1971 Danish 1st Division
Pin
ISO 15919
NFL regular season
Ben
India
2010 NFL season
Karnataka
2010 New York Giants season
Adavisomapur, Haveri

Ubangi Stomp
Rockabilly
Rock and roll
Popular music
Music industry
Musical composition
Originality
Replica
Copying
Information
Uncertainty
Epistemology
Greek language
Modern Greek
Medieval Greek
Classical antiquity
History
Martin Temple
Robert Medley
George Santayana
Royal Academy of Arts
Design Council
Burlington House
Spain
Piccadilly
Diamante Cup
City of Westminster
Immigration to Great Britain
Sri Lanka
Ireland
Continent
Europe
England
Regions of England
Landmass
Norway
East Midlands
De Winterhead
South America
Ivar Aasen
Argentina
Nynorsk
Buenos Aires
Bokmål
En Vivo (Selección Nacional de Panamá)
Norwegian language
Magnhild Eia

Graph 5

My Excuse
Rock music
Popular music
Music industry
Musical composition
Originality
Replica
Copying
Villa di Maiano
Villa
Ancient Rome
Historiography
Historian
Recorded history
History
George Santayana
Two Small Bodies
Thriller (genre)
Film genre
Narrative
Information (media)
Documentary film
Uncertainty
Epistemology
Greek language
Modern Greek
Medieval Greek
Classical antiquity
The Naked Eye (1956 film)
Spain
Europe
Finland
Continent
Thomas Heywood
Landmass
South America
Land
Argentina
Earth
Arsenal de Sarandí
Planet
Sebastián Balmaceda
The Four Prentices of London

ITV Tyne Tees
ITV (TV network)
Free
Television
Monochrome
Lightness
Colorimetry
Dihydroartemisinin/piperaquine
Combination drug
Dosage form
Medication
Drug
Warrior monk
Technology
List of art media
Monk
Material
Asceticism
Plant
Sense
Chemical substance
Multicellular organism
Physiology
Matter
Star
Astronomical object
Function (biology)
Classical physics
Sun
Biology
Physical body
Craig Miller (runner)
Fabien Bacquet
Lancaster, Pennsylvania
Soissons
South Central Pennsylvania
Communes of France
U.S. state
Dual sport
United States
Country
Political geography
Iran at the 2008 Summer Paralympics
Politics
Hypericum mutilum
Hypericum
List of Hypericum species
Type species
International Code of Zoology
Convention (norm)
Natural science
Electromagnetic radiation
Stenognatha
Physics
Natural science
Solid
Projection screen
Drama (film and television)
In the Name of the Law (1952 film)
Four Trials
The Disappearance of Aimee
Light
Ceramic
Electric light
Ceramic art
Chiengzao National Natural Park
Holy Spirit
Nature
Location
Geography
Branches of science
Chemistry
Light fixture
Style (visual arts)
Visual arts
Luxo Jr. (character)
Architectural style
Mission Revival architecture
Emeraude Toubia
List of The Mortal Instruments characters
The Mortal Instruments
Fantasy literature
Fictional universe
Psychology
Linguistics
Organic chemistry
Hydrocarbon
Petroleum reservoir
Oil megaprojects
Oil megaprojects (2020)
Two Shoes (song)
The Cat Empire
Felix Riebl
Cities (The Cat Empire album)
Consistency
Classical logic
Logic
Data compression
Dialect
Pinyin
Chinese characters
Fact
Knowledge
Reality
Existence
Engineering
Process (engineering)
Computer science
Science
Alexandre
Paul Delaroche
Romanticism
Digipak
Trademark
Ontology
Philosophy
Logical form
Action (philosophy)
Artificial Intelligence
Glossary of Artificial Intelligence
Marshall McLuhan
Sovereign (semiosis)
Semiosis
Frederick William, Duke of Courland
German language
West Germanic languages
Germanic languages
Semiotics
Message
Meaning (semiotics)
Communication
Lhota (Kladno District)
Kladno District
Czech language
West Slavic languages
Language family
System
Law
Social science
Political science
Legitimacy (political)
Authority
Executive (government)
Public administration
Local government
Society
Mathematics
Jurisprudence
Legal person
Marshall McLuhan
Natural person
Afroasiatic languages
Classical language
Jefferson Township, Maries County, Missouri
Civil township
Sociology
Sociolinguistics
Social norm
Game theory
Tax
Company
Corporation
Tax avoidance
Tax break
Private university
Southern Methodist University
Entrepreneurship
1940 SMU Mustangs football team
Ernie Ball
Lain language
Northwest Semitic languages
Hebrew language
Asteraceae
Family (biology)
Variety (linguistics)
Standard language
Community
Competition
Sport
Small business
Israel
Zalman Shazar Junior High School
Primary (church officer)
American English
British English
Nation
Social group
Solomon Islands
Team (rugby union)
Phil Bennett (rugby union)
Hierarchy
Outline (list)
Outline of Bible
Bible
Semitic people
Abrahamic religions
Christianity
List of Christian denominations by number of members
Catholic Church
Catholic Church in Colombia
Colombians
Santos Gutiérrez
Antarctic Peninsula
James Ross Island
Cape Gage
Lotus 95T
Formula One
Social behavior
Behavior
International law
Sovereign state
Alice Pollard
Culture
Music
Music genre
Jazz
Collin Walcott
American football
Basketball
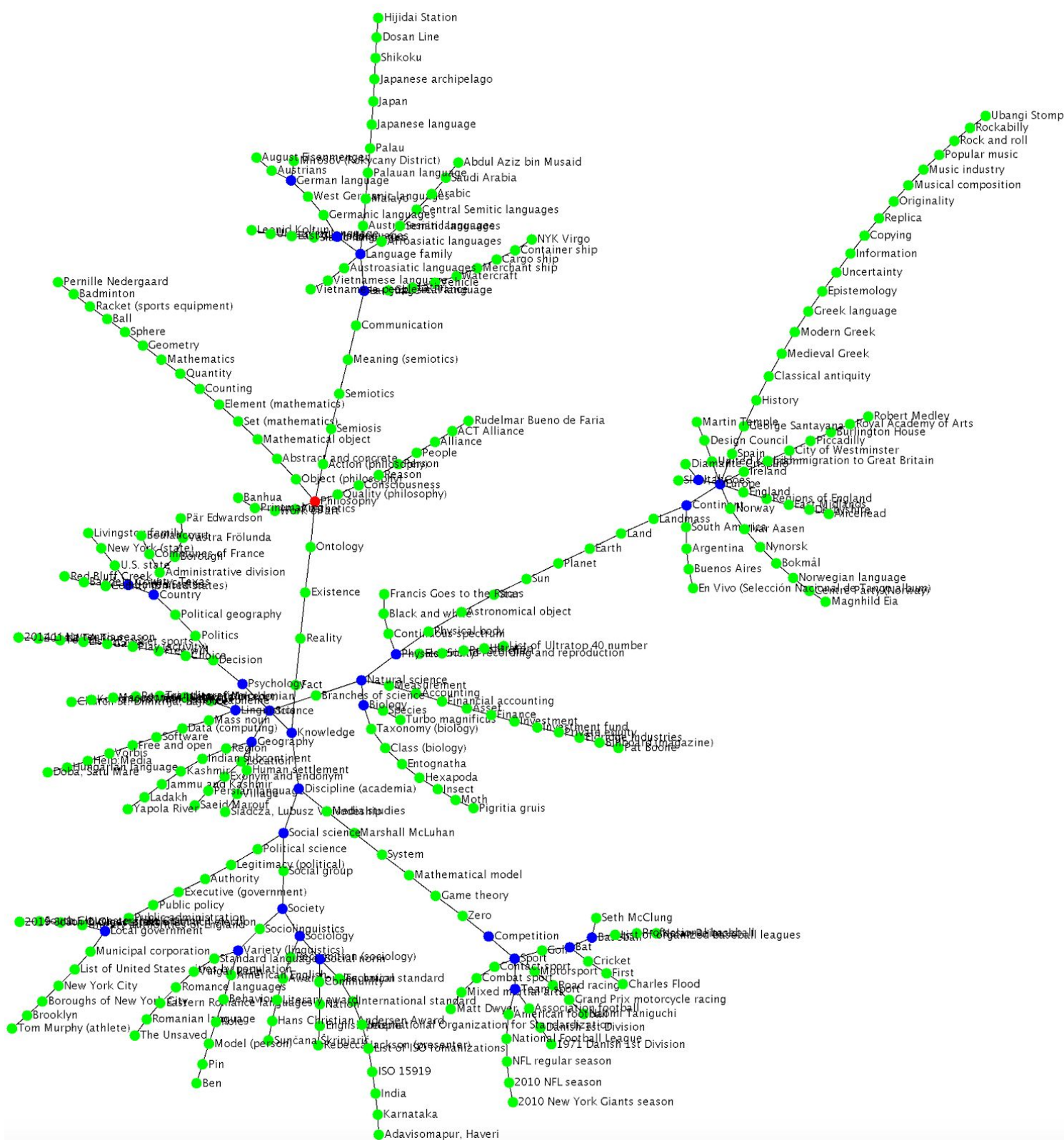Pete Latan
Brian Lynch (basketball)
Judy Cloomie
Demographics of South Africa

Graph 6

Graph 7

Graph 8

Graph 9

Graph 10