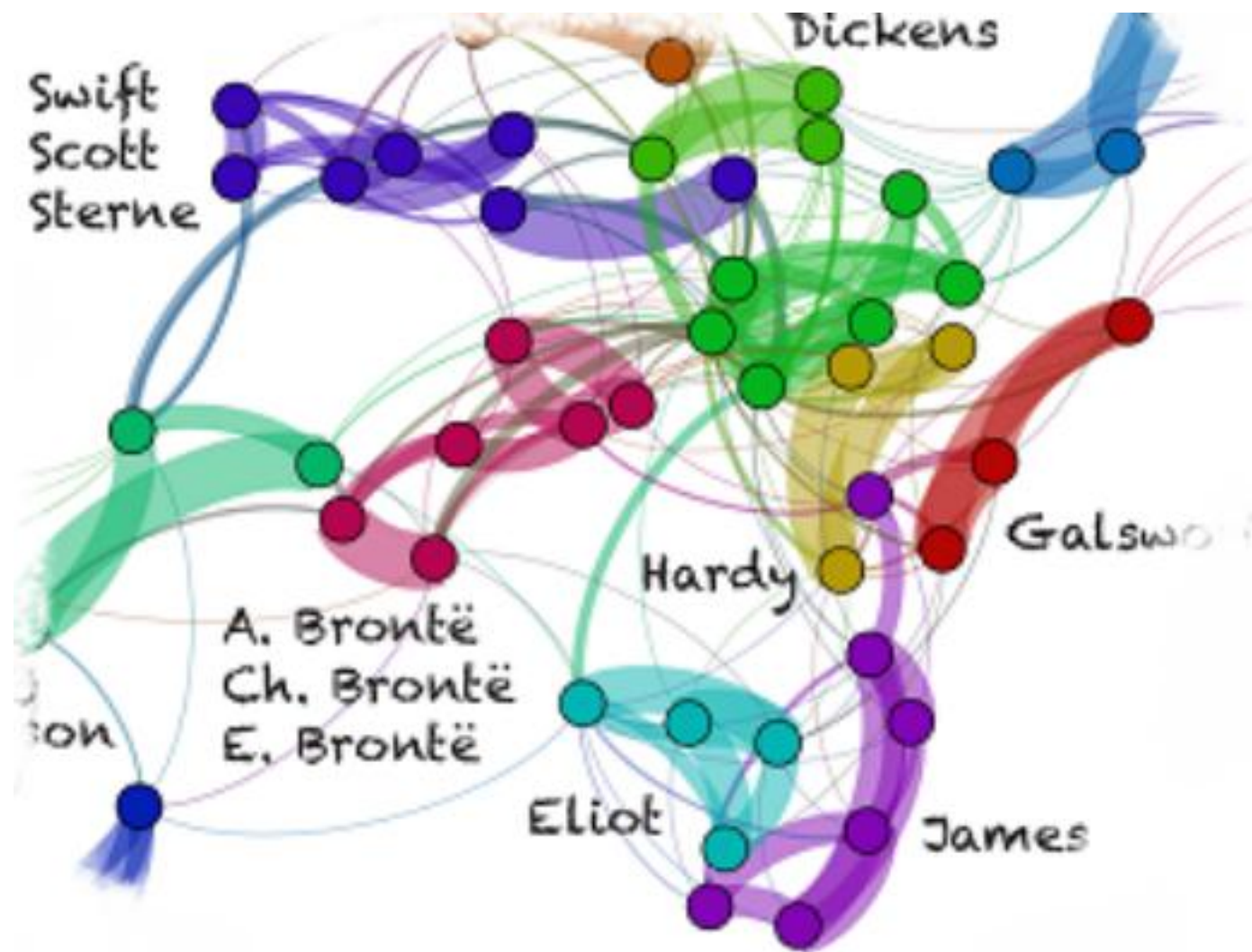


Stylometria z pakietem Stylo

Stylometria jest metodą analizy dzieł sztuki. Pozwala ona na określenie stylu danego dzieła, i na podstawie tego – kto jest jego potencjalnym twórcą, albo czy dana osoba naprawdę jest autorem.

W ramach niniejszego projektu w celu wykonania takiej analizy skorzystano z pakietu „Stylo”.



Czym jest Stylo

- Pakiet stylo został opracowany w 2016 przez Macieja Edera. Jest to potężne narzędzie do przeprowadzenia różnorodnych testów i analiz na zbiorach tekstowych. Wybór algorytmów i innych ustawień odbywa się przez interfejs graficzny, ale możliwe jest też ręcznie ustawianie parametrów w kodzie programu. Pakiet również przedstawia liczne metody wizualizacji wyników, co znacznie upraszcza ich odczyt i zrozumienie.

Stylometry with R | stylo | set parameters

JT & LANGUAGE FEATURES STATISTICS SAMPLING

INPUT: plain text xml xml (plays) xml (no titles) html

LANGUAGE: English English (contr.) English (ALL) Latin Latin (u/v > u)
Polish Hungarian French Italian Spanish
Dutch German CJK Other UTF-8

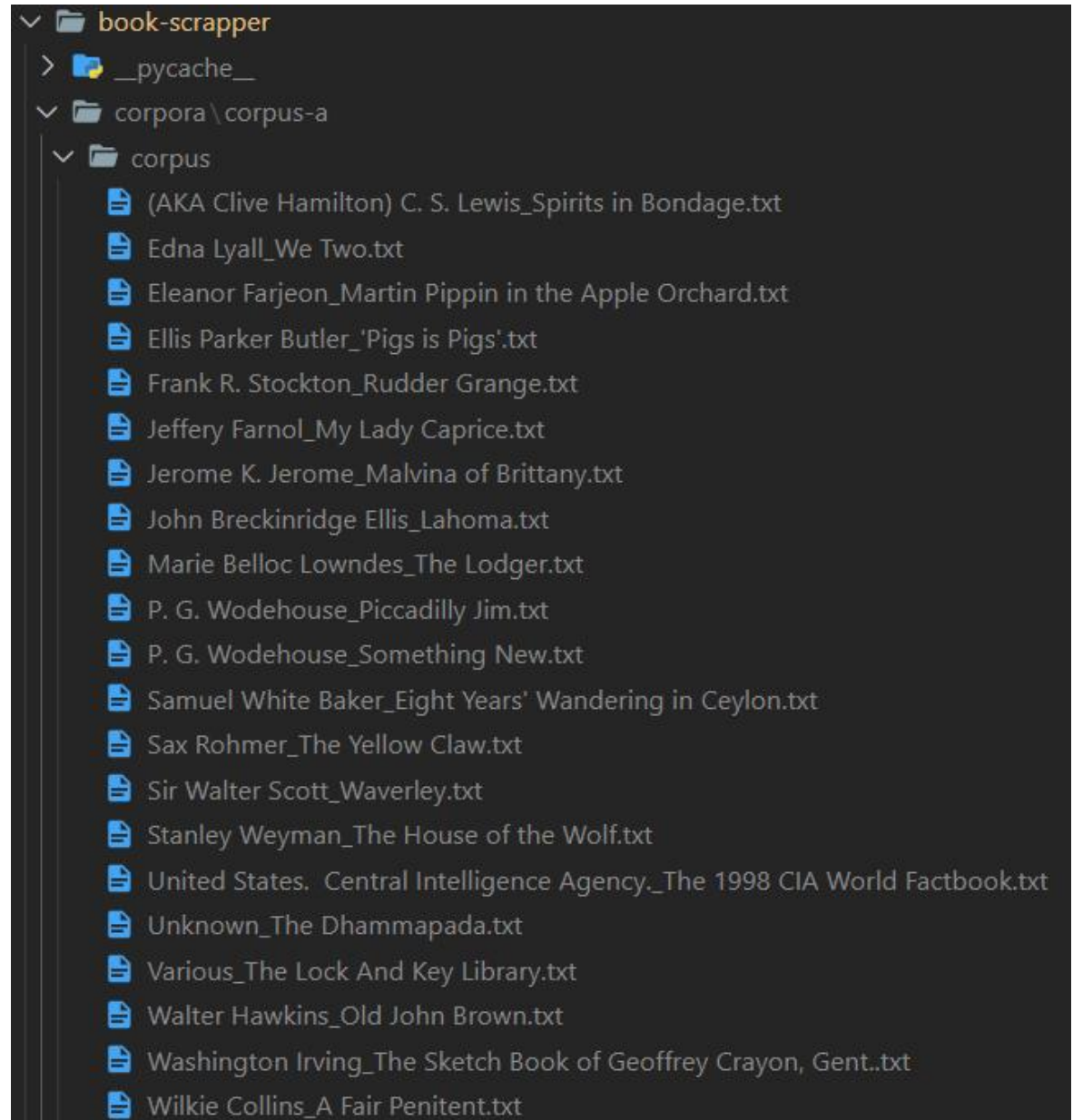
OK

Praca ze Stylo



- Ważnym etapem pracy ze Stylo jest przygotowanie danych. Pakiet zawiera kilka przykładowych zbiorów tekstów i ich cech na potrzeby testowania, ale wraz z tym narzuca pewną „architekturę” organizacji projektu.

W naszym projekcie stworzyliśmy web scrapper, przy pomocy którego zostały pobrane teksty różnych autorów. Ograniczyliśmy się do 21 tekstu.



Następnym etapem jest uruchomienie Stylo albo z poziomu R, albo z poziomu Python.

```
library(stylo)
stylo(analysis.type = "CA",
write.png.file = TRUE, custom.graph.title = "CA Analysis",
gui = FALSE)
```

```
from rpy2 import robjests

R = robjests.r
R.library("stylo")
R.stylo(gui = True)
```

Ustawienia: język i dane wejściowe

Stylometry with R | stylo | set parameters

INPUT & LANGUAGE	FEATURES	STATISTICS	SAMPLING	OUTPUT	
INPUT:	plain text	xml	xml (plays)	xml (no titles)	html
	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
LANGUAGE:	English	English (contr.)	English (ALL)	Latin	Latin (u/v > u)
	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Polish	Hungarian	French	Italian	Spanish
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Dutch	German	CJK	Other	Native encoding
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="checkbox"/>

OK

Ustawienia: cechy

Stylometry with R | stylo | set parameters

INPUT & LANGUAGE | FEATURES | STATISTICS | SAMPLING | OUTPUT

FEATURES: words chars ngram size preserve case

MFWS SETTINGS: Minimum Maximum Increment Start at freq. rank

CULLING: Minimum Maximum Increment List Cutoff Delete pronouns

VARIOUS: Existing frequencies Existing wordlist Select files manually List of files

OK

Ustawienia: statystyki i algorytmy

Stylometry with R | stylo | set parameters

INPUT & LANGUAGE | FEATURES | STATISTICS | SAMPLING | OUTPUT

STATISTICS: Cluster Analysis MDS PCA (cov.) PCA (corr.) tSNE

Consensus Tree Consensus strength

DELTA DISTANCE: Classic Delta Cosine Delta Eder's Delta Eder's Simple Entropy

Manhattan Canberra Eudidean Cosine Min-Max

OK

Ustawienia: pobieranie próbek

Stylometry with R | stylo | set parameters

INPUT & LANGUAGE | FEATURES | STATISTICS | SAMPLING | OUTPUT

No sampling Normal sampling Random sampling

Sample size Random samples

OK

Ustawienia: zapis wyników

Stylometry with R | stylo | set parameters

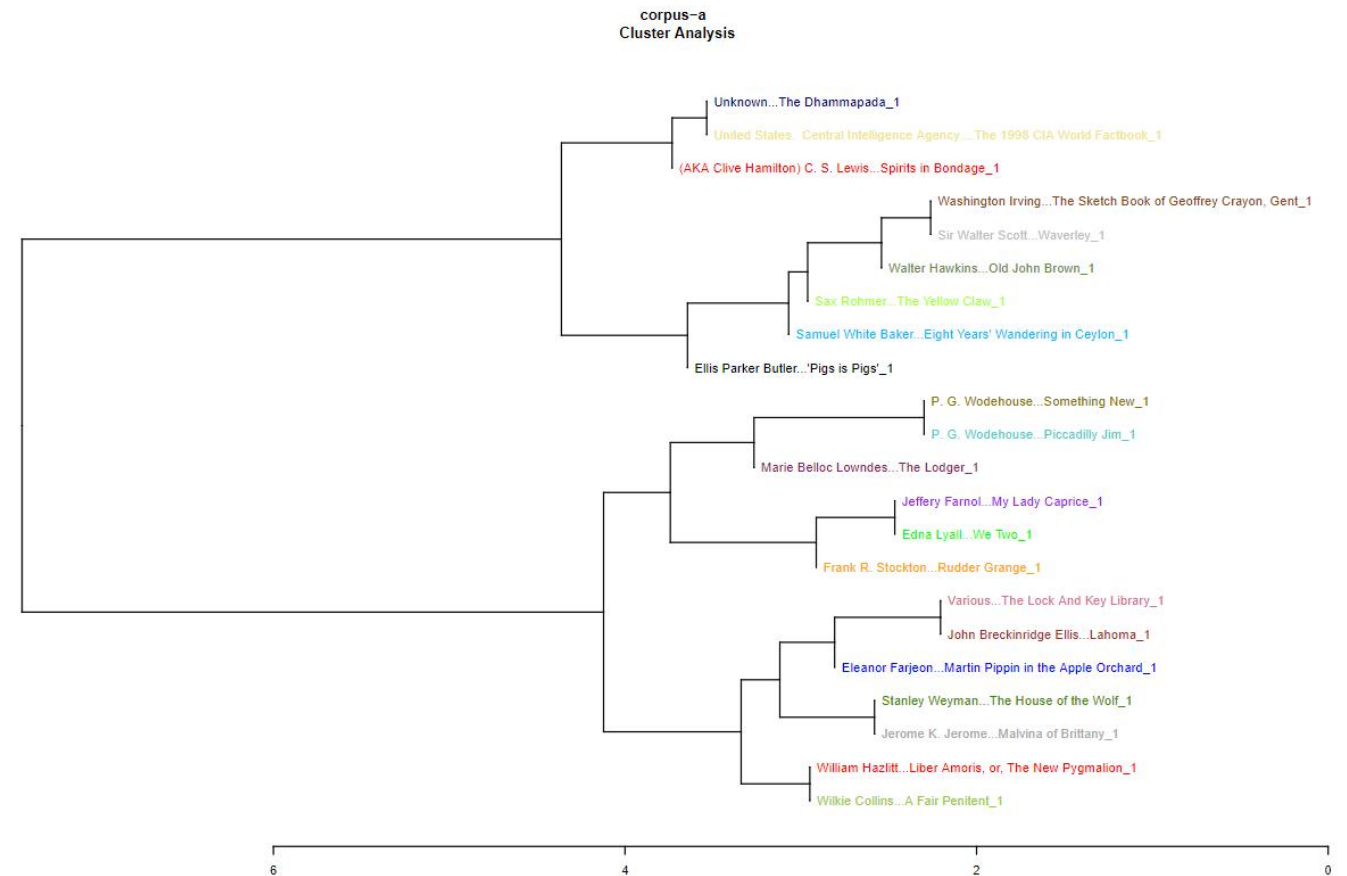
INPUT & LANGUAGE	FEATURES	STATISTICS	SAMPLING	OUTPUT	
GRAPHS:	Onscreen <input checked="" type="checkbox"/>	PDF <input type="checkbox"/>	JPG <input type="checkbox"/>	SVG <input type="checkbox"/>	PNG <input checked="" type="checkbox"/>
PLOT AREA:	Set default <input type="checkbox"/>	Plot height <input type="text" value="7"/>	Plot width <input type="text" value="7"/>	Font size <input type="text" value="10"/>	Line width <input type="text" value="1"/>
		Colors <input checked="" type="radio"/>	Grayscale <input type="radio"/>	Black <input type="radio"/>	Titles <input checked="" type="checkbox"/>
PCA/MDS:	Labels <input checked="" type="radio"/>	Points <input type="radio"/>	Both <input type="radio"/>	Margins <input type="text" value="2"/>	Label offset <input type="text" value="0"/>
PCA FLAVOUR:	Classic <input checked="" type="radio"/>	Loadings <input type="radio"/>	Technical <input type="radio"/>	Symbols <input type="radio"/>	
VARIOUS:	Horizontal CA tree <input checked="" type="checkbox"/>	Save distance table <input type="checkbox"/>	Save features <input type="checkbox"/>	Save frequencies <input type="checkbox"/>	Dump samples <input type="checkbox"/>

OK

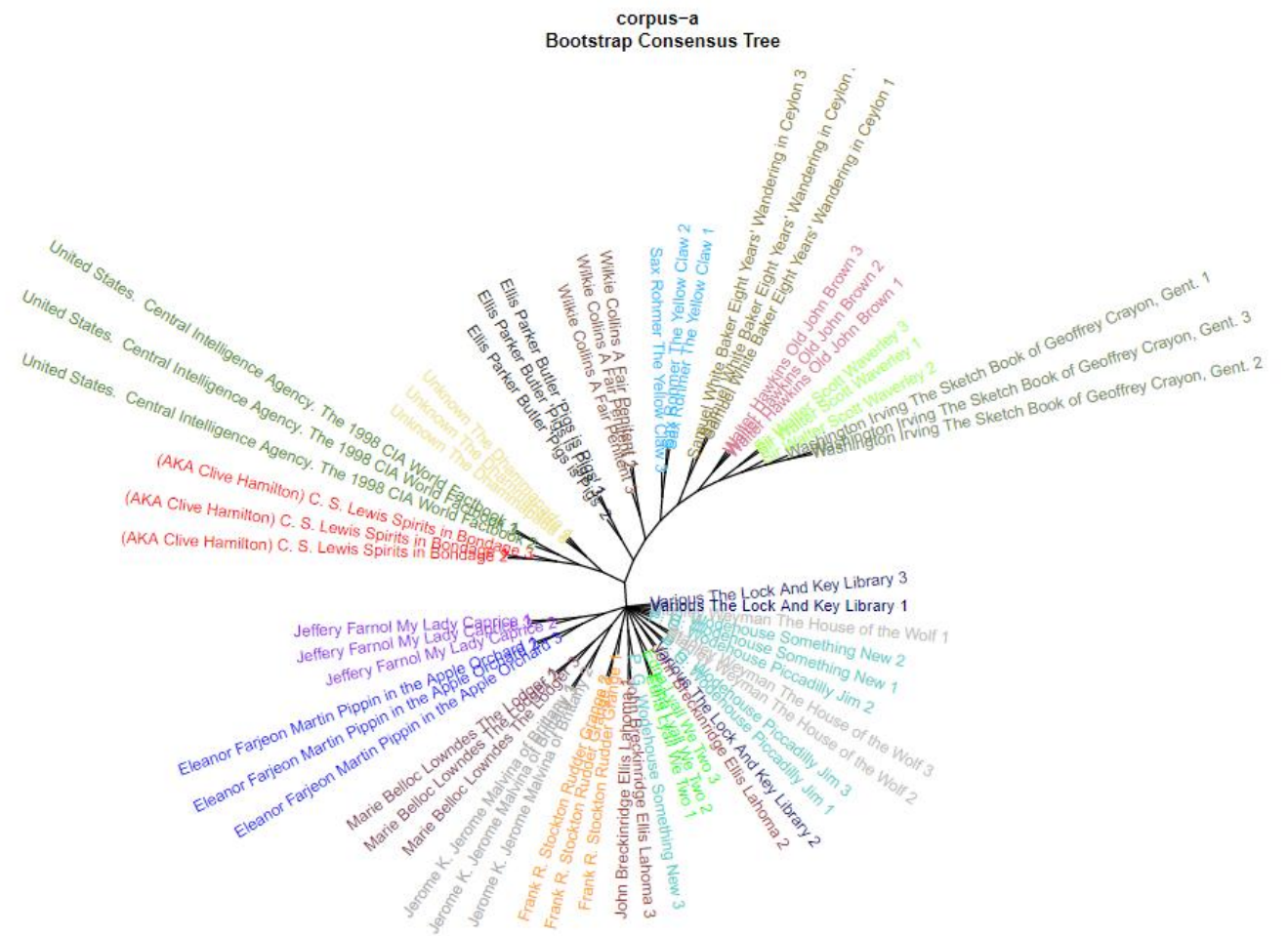
Co potrafi Stylo?

- Obliczanie podobieństw tekstów z wykorzystaniem różnych metryk (odległość kosinusowa, entropia ...)
- Testowanie, kto jest autorem tekstu
- Wykrywanie plagiatu na podstawie analizy segmentów tekstu i uczenia maszynowego

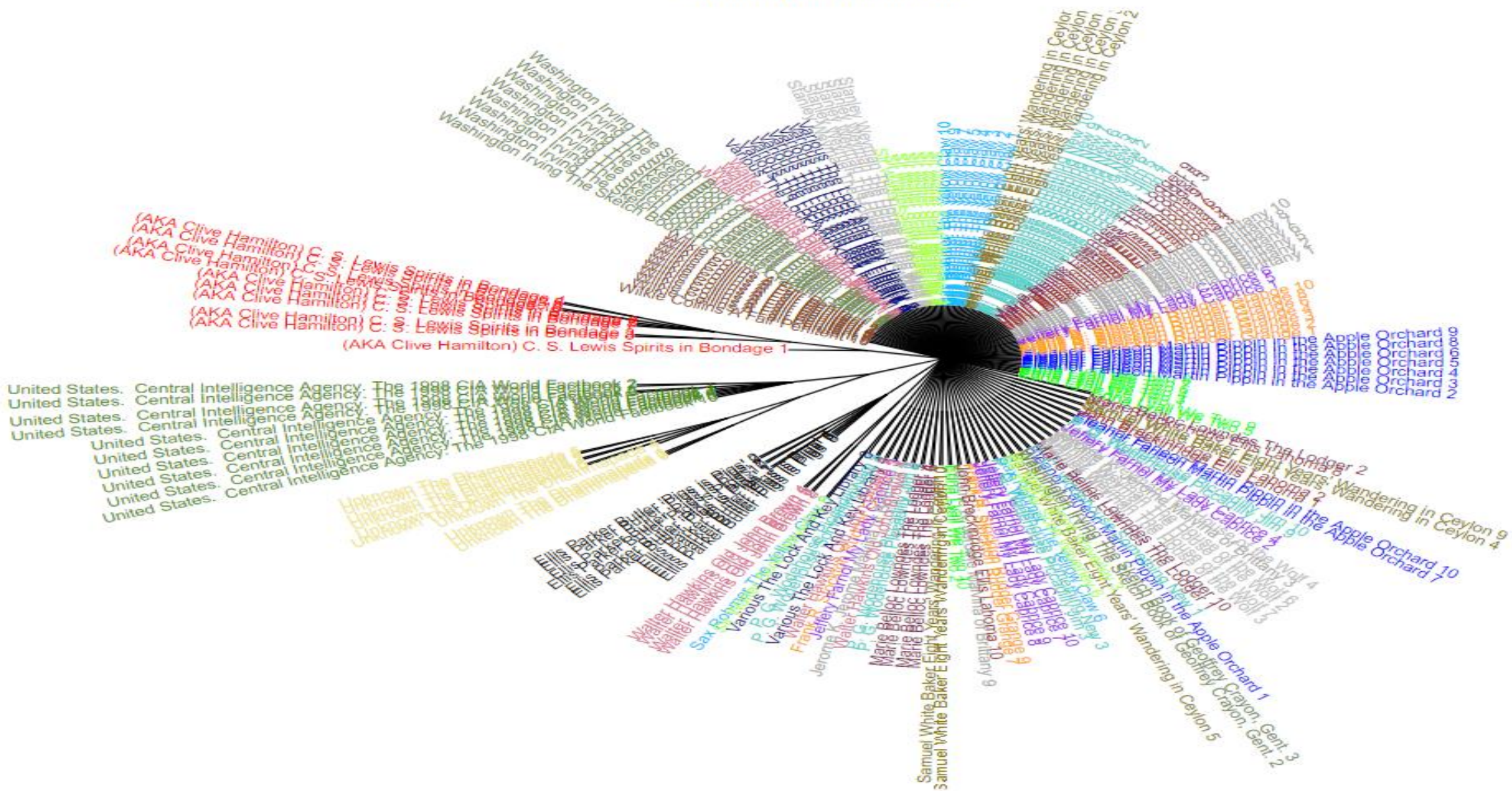
Przykłady działania



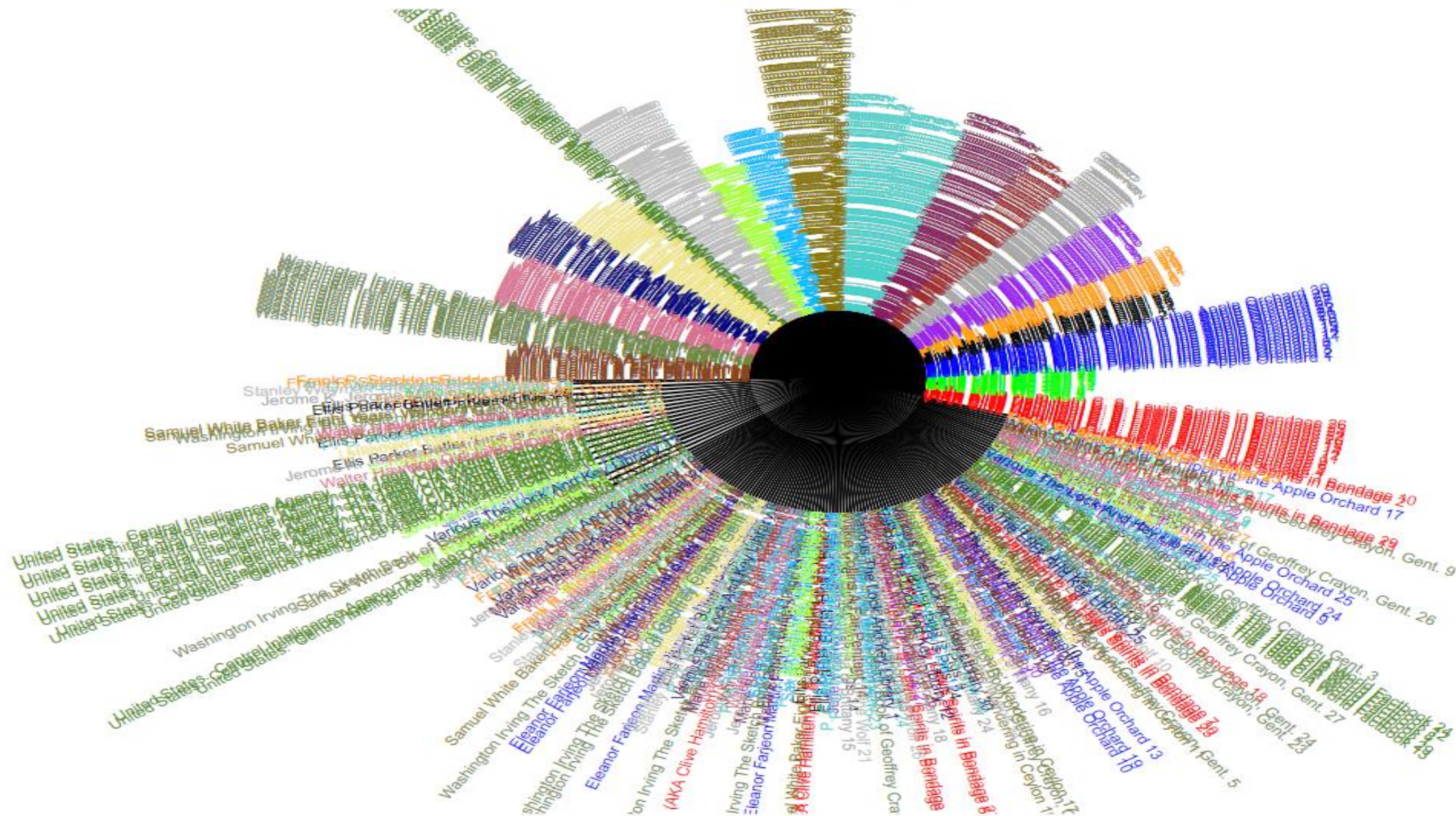
Drzewo zgodności



corpus-a Bootstrap Consensus Tree



corpus-a
Bootstrap Consensus Tree



PCA

