

COIN7 CoolPeople

Final presentation

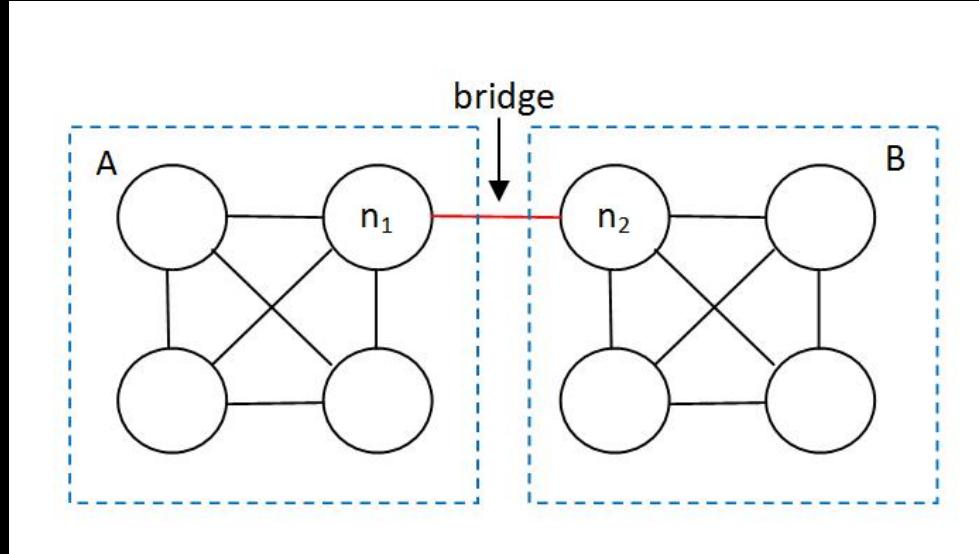
David Huser
Igor Bosnjak

Goals

Extract **names of persons** in **text sources** and generate a **network file** compatible with e.g. Gephi network visualization.

Two persons are linked if they appear in the same document.

If a person is in more than a document it becomes a bridge.



Methods used

1. Apache openNLP

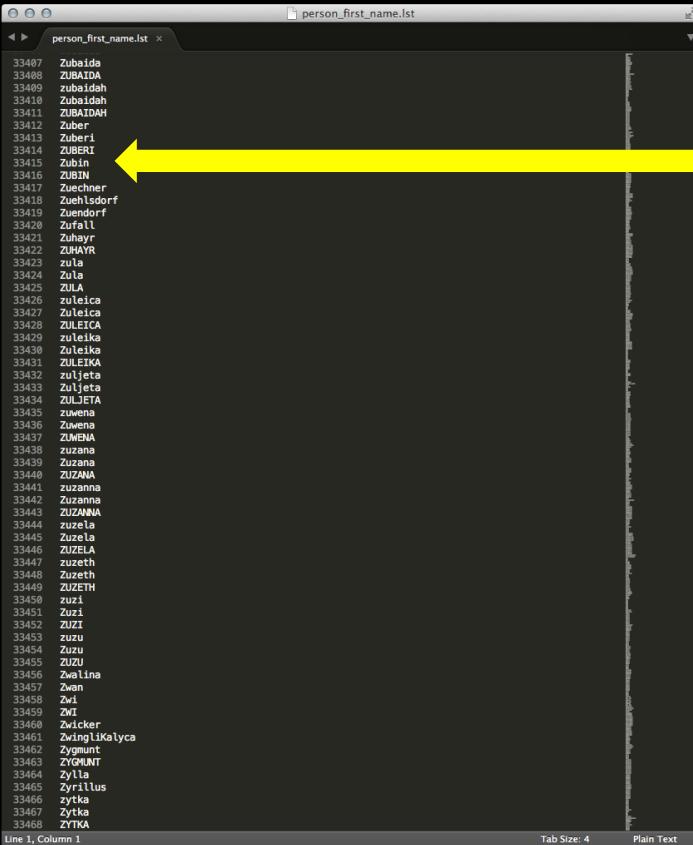
Natural Language Processing

Machine learning based toolkit for the processing of natural language text.

Supports the most common NLP tasks:

tokenization, sentence segmentation, part-of-speech tagging, named entity extraction, chunking, parsing, coreference resolution

2. Whitelist



```
person_first_name.lst
33407 Zubaida
33408 ZUBAIDA
33409 zubaida
33410 Zubaida
33411 ZUBAIDA
33412 Zuber
33413 Zuberi
33414 ZUBERI
33415 Zubin
33416 ZUBIN
33417 Zuleika
33418 Zuchsdorf
33419 Zuendorf
33420 Zufall
33421 Zubayr
33422 ZUHAYR
33423 zula
33424 Zula
33425 ZULA
33426 zuleica
33427 Zuleica
33428 ZULEICA
33429 zuvana
33430 Zuleika
33431 ZULEIKA
33432 zuljeta
33433 Zuljeta
33434 ZULJETA
33435 zuwena
33436 Zuwena
33437 ZUMENA
33438 zuzana
33439 Zuzana
33440 ZUZANA
33441 zuzanna
33442 Zuzanna
33443 ZUZANNA
33444 zuzela
33445 Zuzela
33446 ZUZELA
33447 zuzeth
33448 Zuzeth
33449 ZUZETH
33450 zuri
33451 Zuri
33452 ZUZI
33453 zuzu
33454 Zuzu
33455 ZUZU
33456 Zwolina
33457 Zwan
33458 Zwei
33459 ZKI
33460 Zwickler
33461 ZwingleiKalyca
33462 Zygmunt
33463 ZYGMUNT
33464 Zylla
33465 Zyrillus
33466 zytka
33467 Zytka
33468 ZYTKA
```

If a word in the input document contains one of the names:

Add the **following word** as a new person,
if the prename and the lastname are
Uppercase.

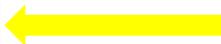
3. GEXF4J

A java library for the GEXF file format

Compare

The following table will help if you want to choose in which format you want to encode your data. If you plan to work only with Gephi, we recommend to use GEXF, for many reasons. The table criteria don't mention all features of formats but concentrate on these supported by Gephi.

	Edge List/Matrix Structure	XML Struture	Edge Weight	Attributes	Visualization Attributes	Attribute Default Value	Hierarchical Graphs	Dynamics
CSV	■							
DL Ucinet	■							
DOT Graphviz		■		■				
GDF								
GEXF	■	■	■	■	■	■	■	■
GML		■	■	■				
GraphML	■	■	■	■	■	■	■	■
NET Pajek	■	■	■	■				
TLP Tulip								
VNA Netdraw		■	■					
Spreadsheet*						■		

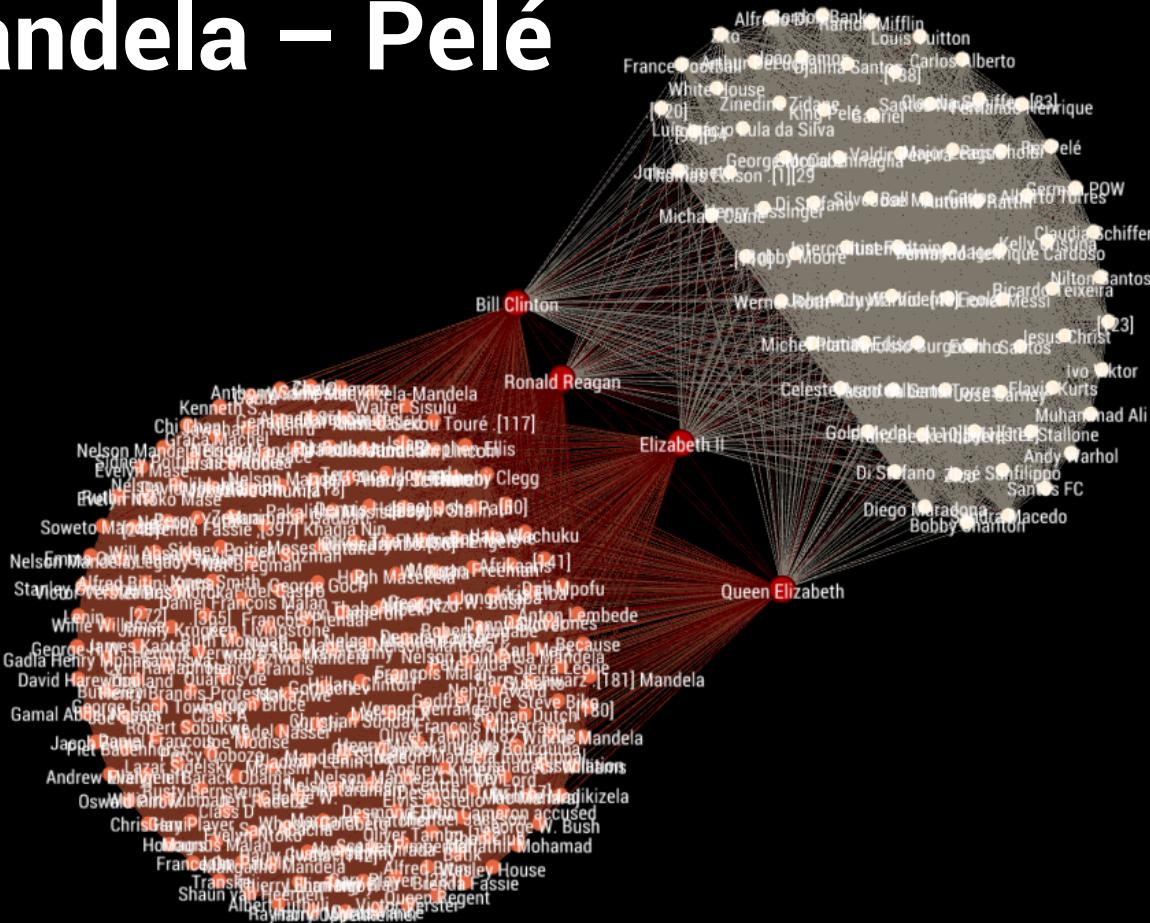


Main results

Nelson Mandela – Pelé



Nelson Mandela – Pelé



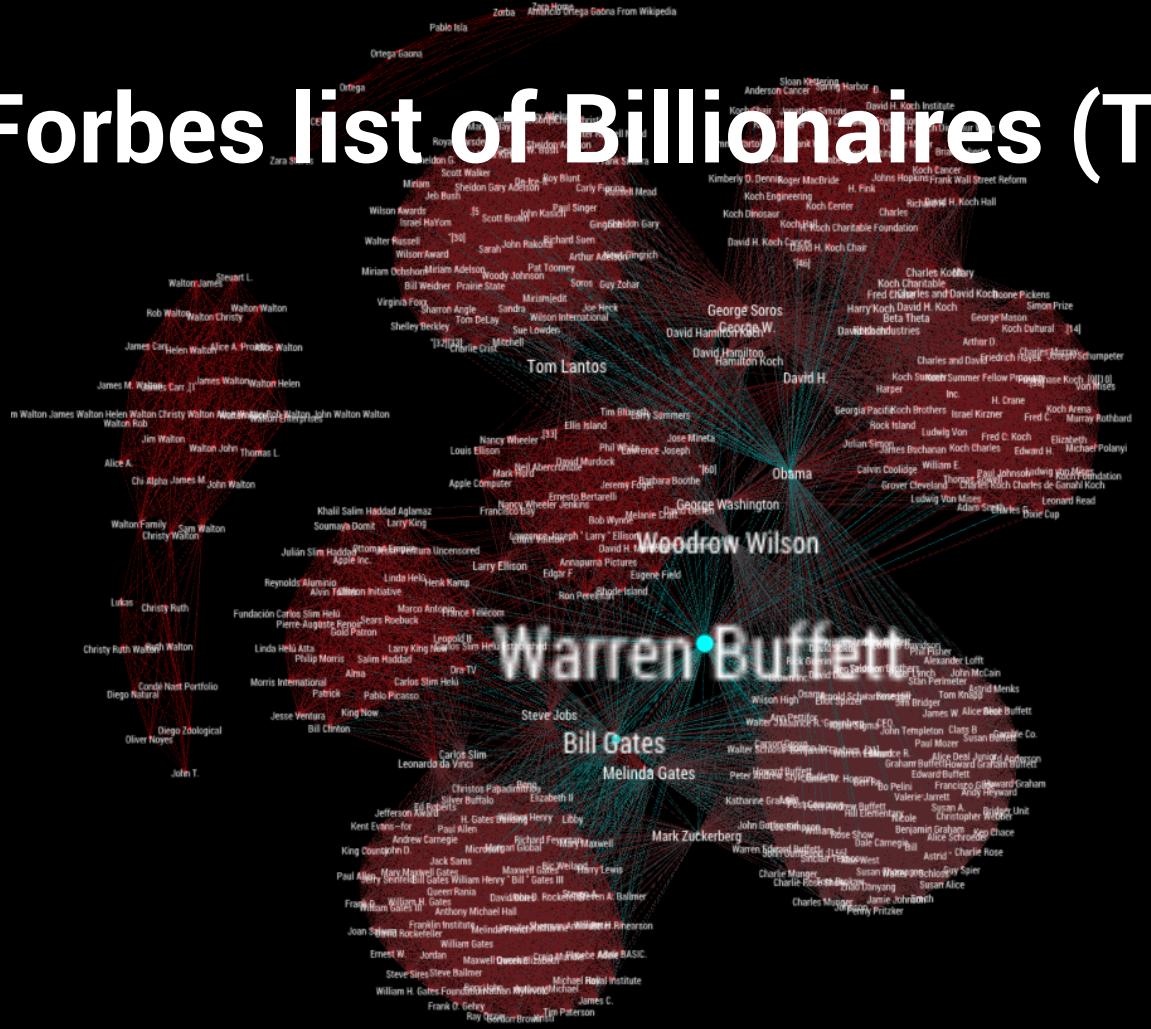
Anne Frank – Pablo Picasso



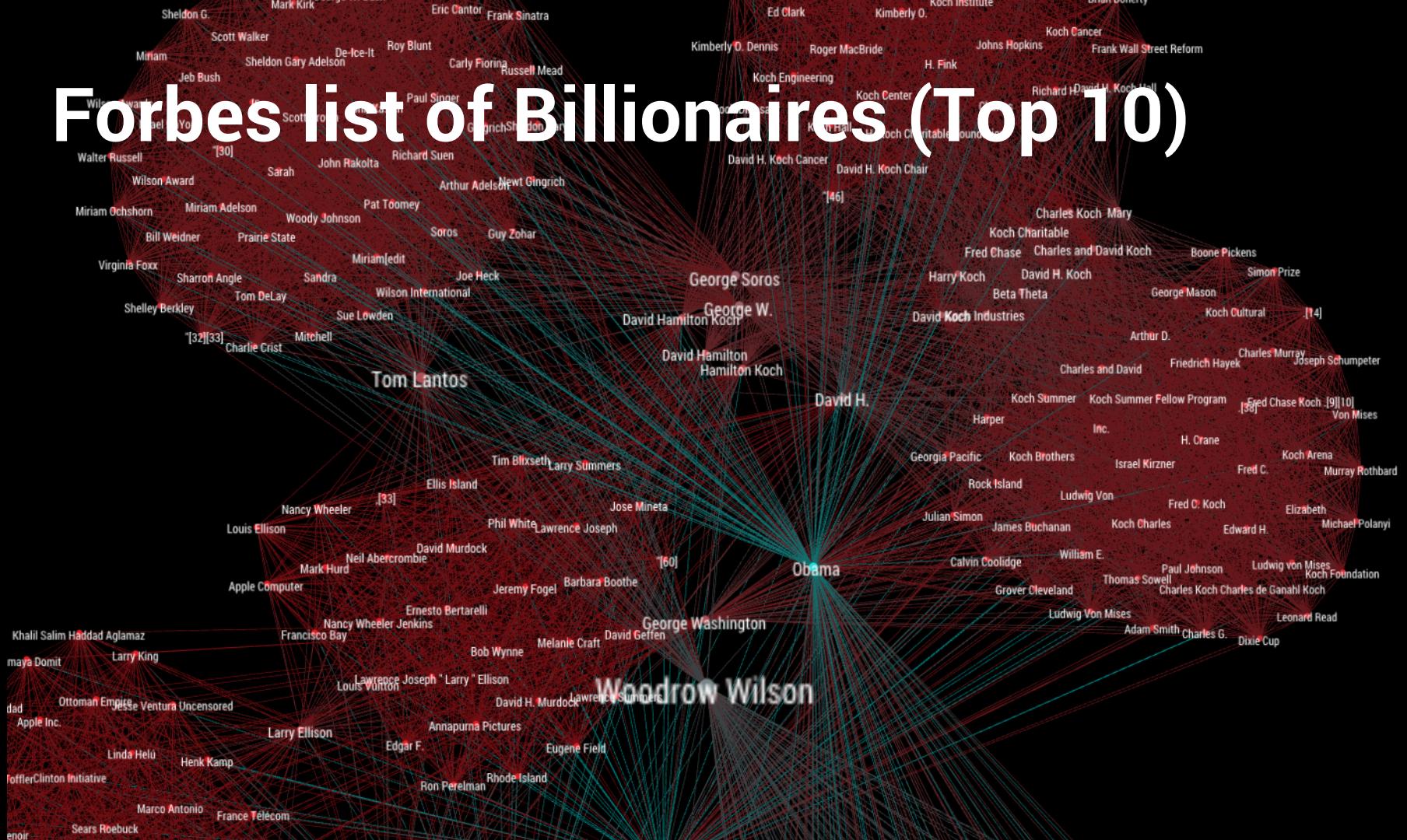
Forbes list of Billionaires (Top 10)

No.	Name	Net worth (USD)	Age	Citizenship	Source(s) of wealth
1	Bill Gates	\$78.0 billion	58	United States	Microsoft
2	Carlos Slim & family	\$71.2 billion	74	Mexico	Telmex, América Móvil, Grupo Carso
3	Warren Buffett	\$65.5 billion	83	United States	Berkshire Hathaway
4	Amancio Ortega	\$62.3 billion	78	Spain	Inditex Group
5	Larry Ellison	\$51.4 billion	69	United States	Oracle Corporation
6	Charles Koch	\$41.4 billion	78	United States	Koch Industries
6	David Koch	\$41.4 billion	74	United States	Koch Industries
8	Christy Walton & family	\$37.8 billion	58-59	United States	Wal-Mart
9	Sheldon Adelson	\$37.7 billion	80	United States	Las Vegas Sands
10	Bernard Arnault & family	\$36.3 billion	65	France	LVMH Moët Hennessy • Louis Vuitton

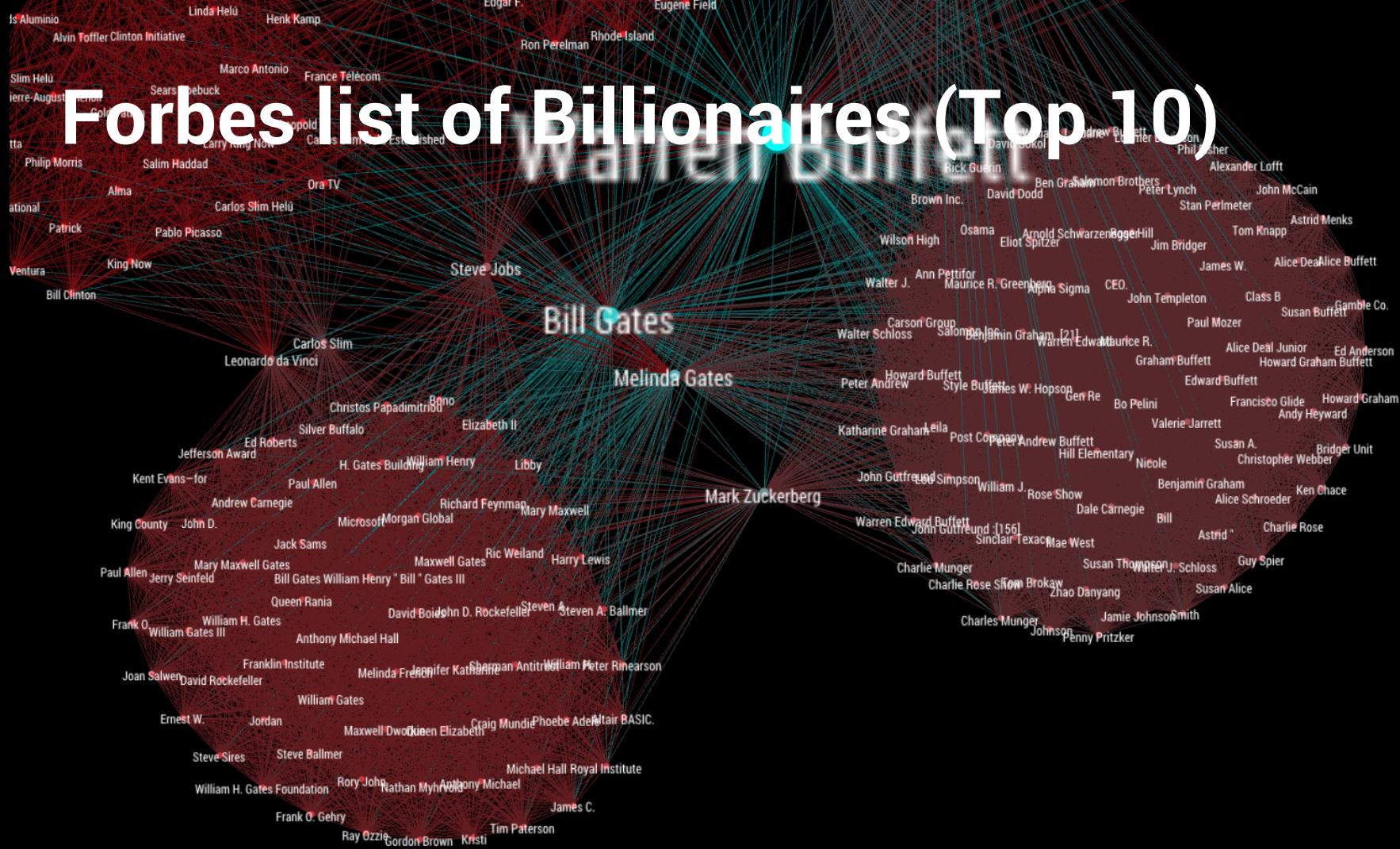
Forbes list of Billionaires (Top 10)



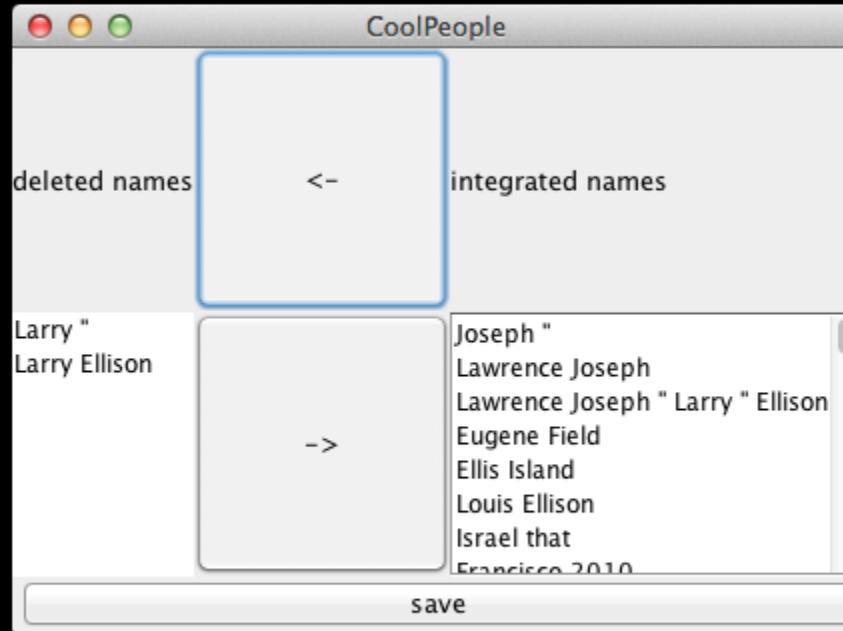
Forbes list of Billionaires (Top 10)



Forbes list of Billionaires (Top 10)



Filter GUI



Open issues

Performance

- Multithreading
- Forbes Top 10: 30 seconds @ Intel i5
- GUI improvements: Filechooser
- Save button

Future work

Extendable architecture

A document object can be extended... all we need is a String to extract names.

- Web scraping
- MySQL Adapter
- ...

Extendable export file

The Gephi export format can be used for
directed graphs, timelines, ...

Maven

Apache Maven build tool for dependencies,
testing and resources

Discussion Q&A

Links

Apache openNLP – <https://opennlp.apache.org/>

openNLP models – <http://opennlp.sourceforge.net/models-1.5/>

Gephi – <https://gephi.org/>

GEXF – <http://gexf.net/format/index.html>

Apache Maven – <http://maven.apache.org/>