

WORKING WITH TEXT

BEYOND WORD CLOUDS

PRESENTED BY ANNA NICANOROVA

bigdata

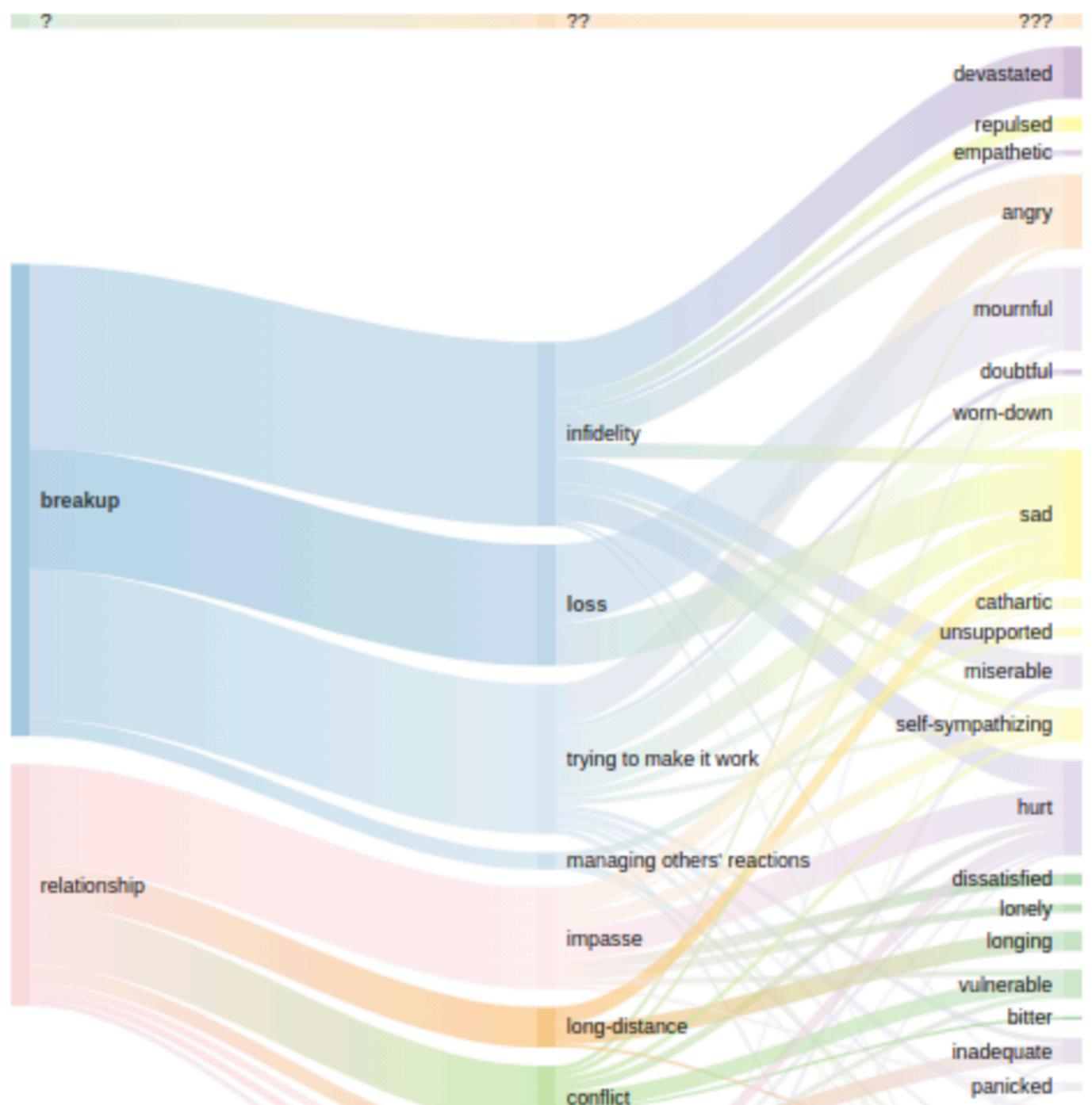
training showcase
microsoft day enterprisethecube
enter time come hortonworks keynotes
new tomorrow database
amitwalia amitwalia
including free json
using boost jimmchugh
world ready find need schedule w help
become going need livedemo^o help
algorithms datakind learn azure social
today via kdnuggets 337 merv insight check
streaming building win meetlike session apachespark
cloudera wish spark great million bigdatany^o
jakeporway good analytics makes insights
thoughtspot talk mongodb mapr
data us hadoop things
booth available challenge



a a a
correlation strength

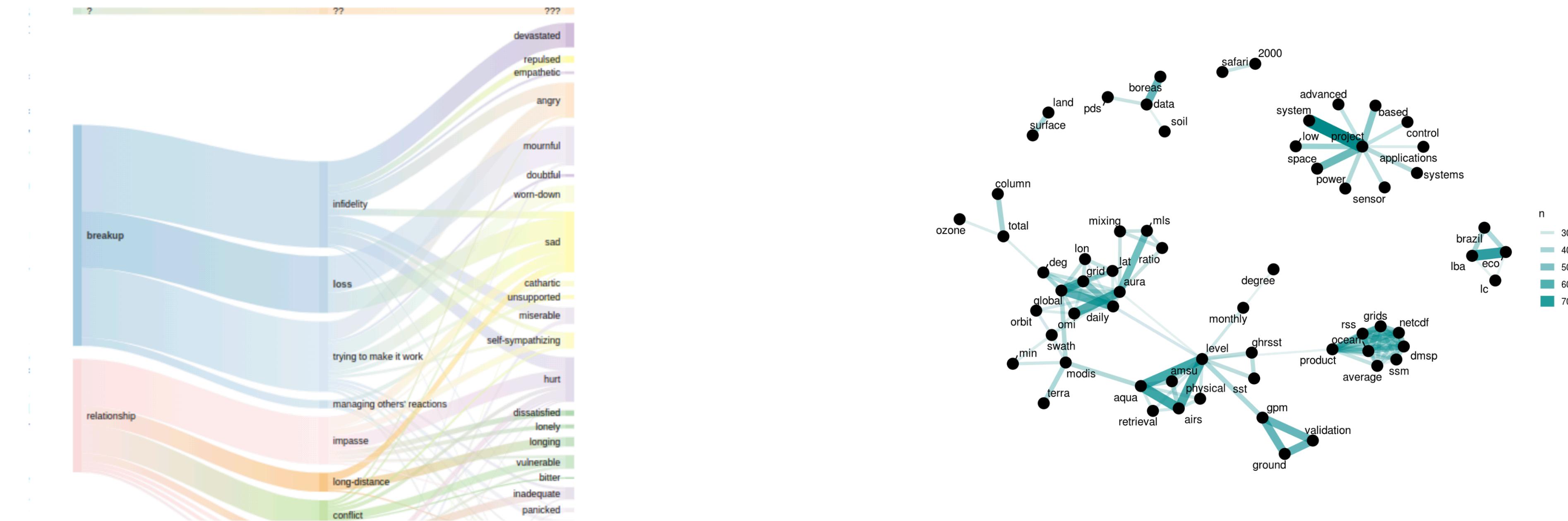
b **b** **b**
prevalence in topic

ALTERNATIVES TO WORD CLOUDS



SANKEY DIAGRAM

ALTERNATIVES TO WORD CLOUDS



SANKEY DIAGRAM

WORD NETWORKS

See also: <http://texttexture.com/>

ALTERNATIVES TO WORD CLOUDS



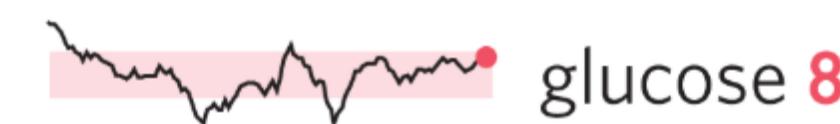
TEXT SUMMARIZATION

ALTERNATIVES TO WORD CLOUDS



TEXT SUMMARIZATION

INFO HIGHLIGHT
(SPARKLINE/CALLOUT)



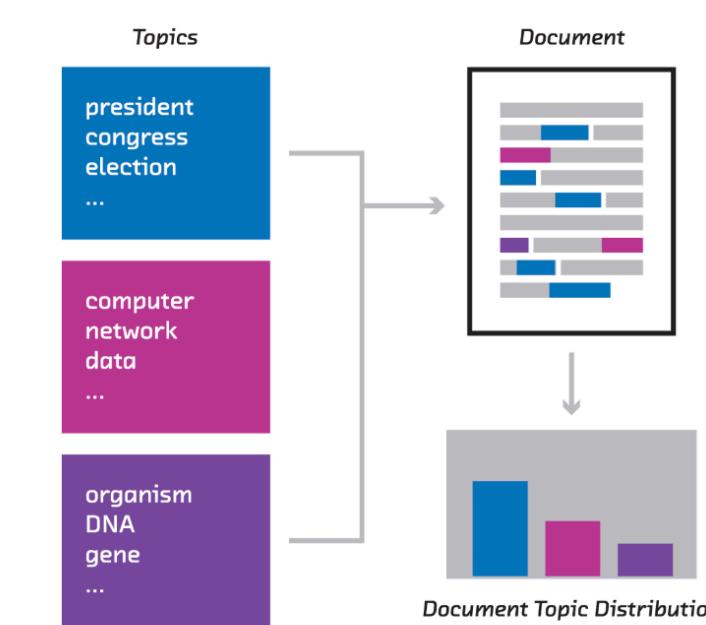
ALTERNATIVES TO WORD CLOUDS



TEXT SUMMARIZATION



INFO HIGHLIGHT
(SPARKLINE/CALLOUT)



TOPIC MODELING

ELEPHANT IN THE ROOM



ELEPHANT IN THE ROOM



AI

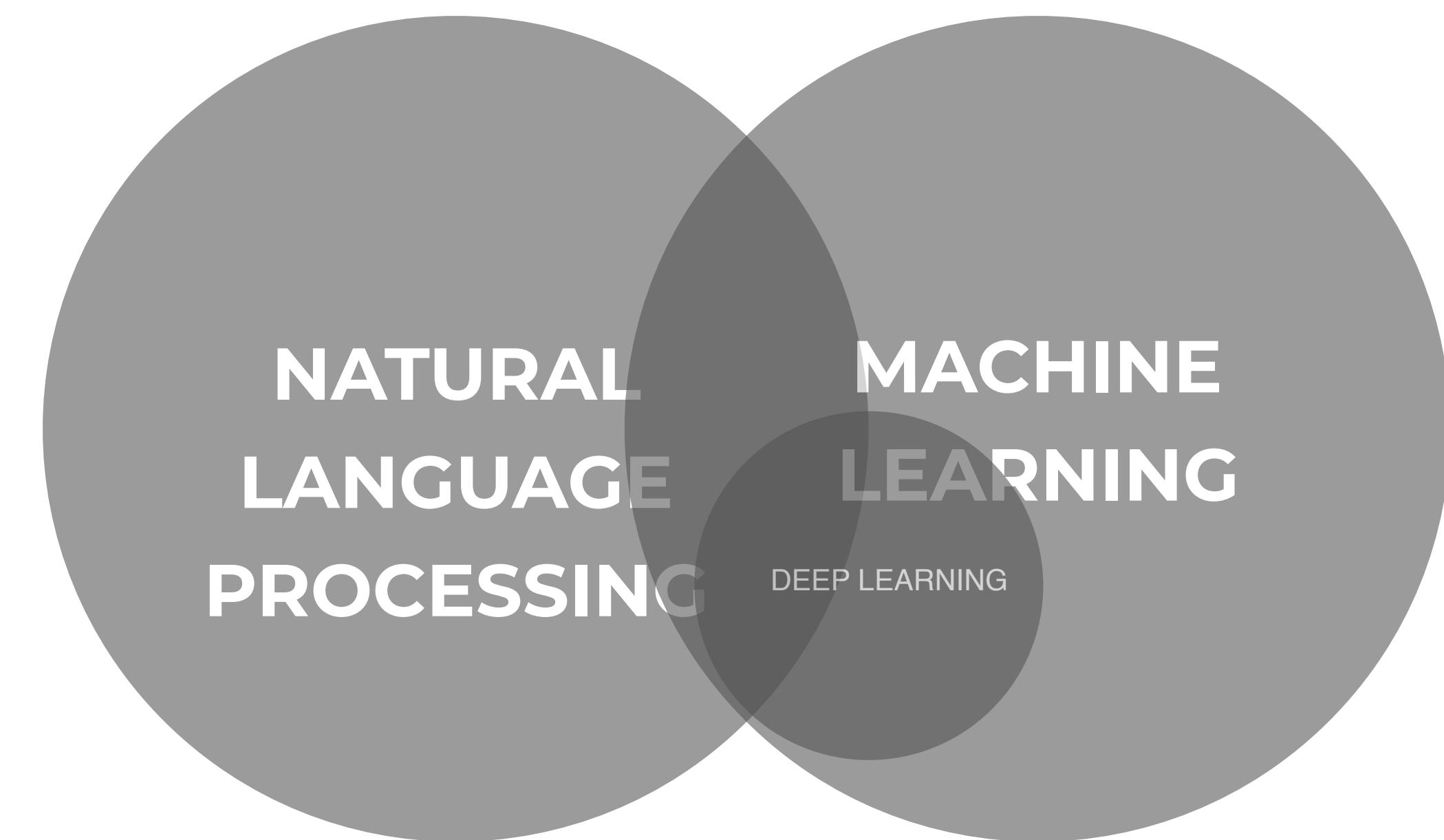
TEXT IS HARD

Consider robots

- They can walk (now run, and maybe even climb stairs)
- They can play chess
- They can play football
- They can recognize people by face
- ⊗ They can't hold a conversation
- Language turns out to be very difficult

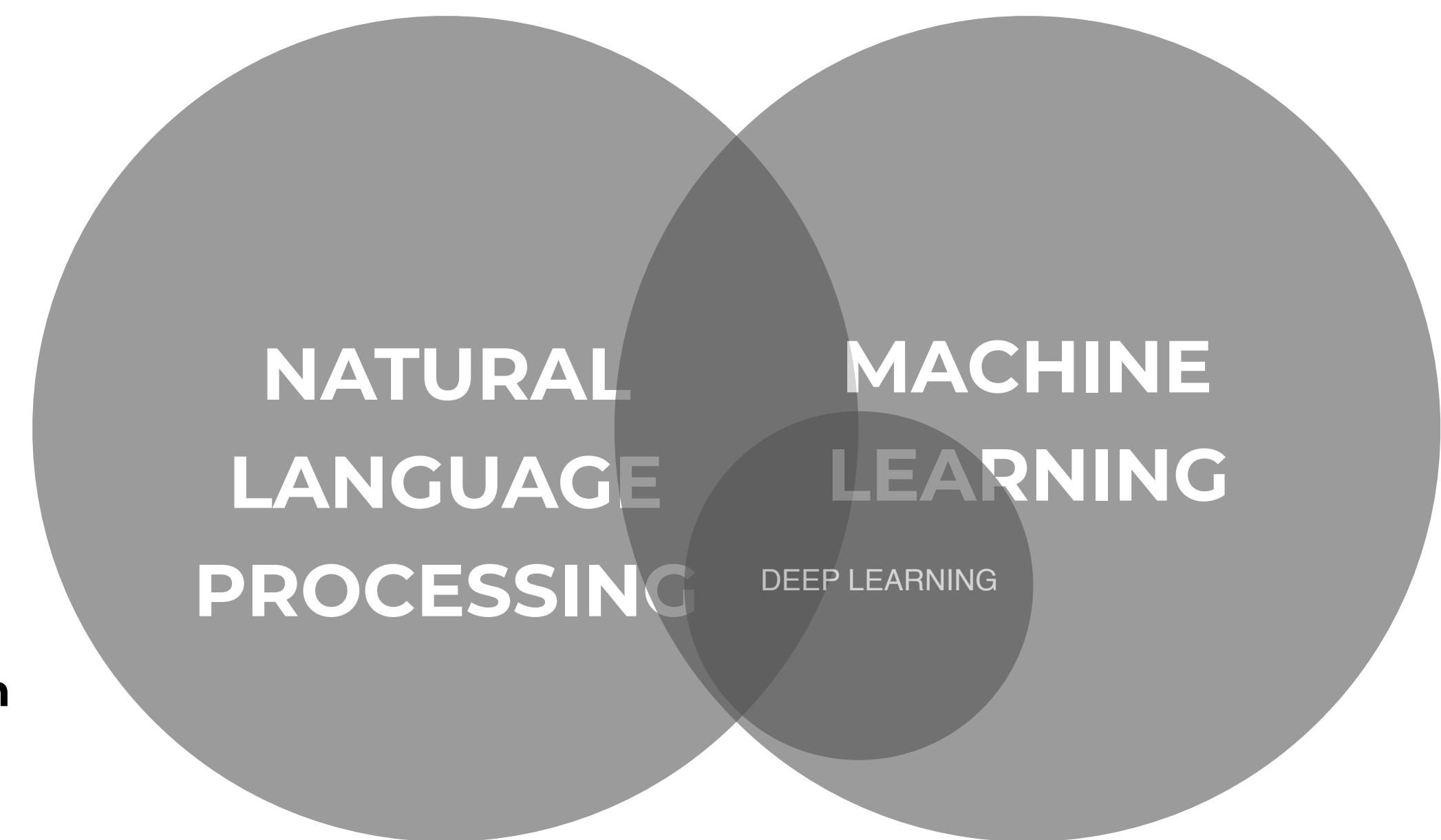


ML/NLP/AI...



ML/NLP/AI...

- Syntax – What part of given text is grammatically true.
- Semantics – What is the meaning of given text?
- Pragmatics – What is the purpose of the text?
- Phonology – It is systematic organization of sounds in language.
- Morphology – It is a study of words formation and their relationship with each other.



- Text classification and clustering
- Information retrieval and extraction
- Machine translation(one language to another)
- Question and answering system
- spelling and grammar checking
- Topic modeling and sentiment analysis
- Speech recognition

LANGUAGE TECHNOLOGY

MOPSTLY SOLVED

SPAM DETECTION

- Let's go to Agra!
- Buy VIAGRA

PART OF SPEECH TAGGING

- Colorless (ADJ) green (ADJ) ideas (NOUN)
sleep (VERB) furiously (ADV)

NAMED ENTITY RECOGNITION

- Einstein (Person) met with UN(ORG)
officials in Princeton (Location)

LANGUAGE TECHNOLOGY

MOSTLY SOLVED

SPAM DETECTION

- Let's go to Agra!
- Buy VIAGRA

PART OF SPEECH TAGGING

- Colorless (ADJ) green (ADJ) ideas (NOUN)
sleep (VERB) furiously (ADV)

NAMED ENTITY RECOGNITION

- Einstein (Person) met with UN(ORG)
officials in Princeton (Location)

MAKING GOOD PROGRESS

SENTIMENT ANALYSIS

- Best pizza in New York
- We waited for 30 minutes to get seated

MACHINE TRANSLATION

- I like learning Chinese
- 我喜欢学中文
- Wǒ xǐhuān xué zhōngwén

INFORMATION EXTRACTION

- You're all booked for your flight NYC to SF on
May 6th
- Fly May 6th. Add to calendar?

WORD SENSE DISAMBIGUATION

- I need batteries for my mouse

LANGUAGE TECHNOLOGY

MOSTLY SOLVED

SPAM DETECTION

- Let's go to Agra!
- Buy VIAGRA

PART OF SPEECH TAGGING

- Colorless (ADJ) green (ADJ) ideas (NOUN)
sleep (VERB) furiously (ADV)

NAMED ENTITY RECOGNITION

- Einstein (Person) met with UN(ORG)
officials in Princeton (Location)

MAKING GOOD PROGRESS

SENTIMENT ANALYSIS

- Best pizza in New York
- We waited for 30 minutes to get seated

MACHINE TRANSLATION

- I like learning Chinese
- 我喜欢学中文
- Wǒ xǐhuān xué zhōngwén

INFORMATION EXTRACTION

- You're all booked for your flight NYC to SF on
May 6th
- Fly May 6th. Add to calendar?

WORD SENSE DISAMBIGUATION

- I need batteries for my mouse

STILL REALLY HARD

QUESTION/ANSWER

- How effective is Ibuprofen in reducing fever
in patients with acute febrile illness?

PARAPHRASING

- XYZ acquired ABC yesterday
- ABC has been taken over by XYZ

SUMMARIZATION

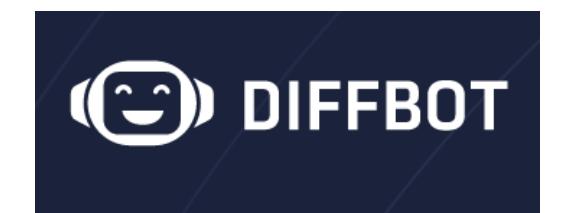
- The Dow Jones is up
 - The S&P500 jumped
 - Housing prices rose
- Economy is good

DIALOGUE

- Where is Citizen Kane playing in NYC?
- Angelica theater at 7:30. Do you want a ticket?

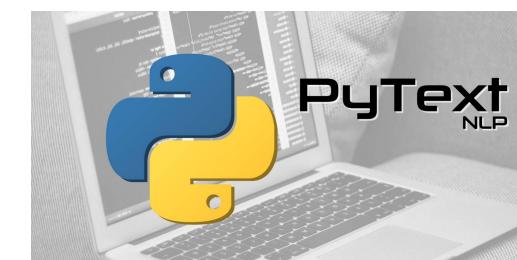
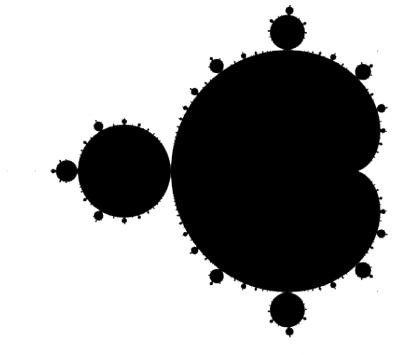
PAID NLP API'S

You don't really need to write your own code - there are plenty providers, who for a certain charge will process text for you



NLP IN PYTHON

Python is by far the most popular language for doing Natural Language Processing. This list is constantly updated as new libraries come into existence.



TUTORIAL PLAN

STEP 1:

CLEANING DATA

STEP 2:

TEXT EXPLORATION

STEP 3:

ENTITY EXTRACTION

STEP 4:

SUMMARIZATION

STEP 5:

TOPING MODELING

STEP 6:

GETTING FANCY &
CREATIVE APPROACH TO
SHOWING DATA RESULTS

ART BOOKS FROM GUGGENHEIM MUSEUM

GUGGENHEIM the Solomon R. Guggenheim Museum

Share
 Favorite
 Play All

MORE

ABOUT

COLLECTION

204 RESULTS

Search this Collection

Metadata
 Text contents

Media Type
 texts 204

Year

- 2003 3
- 2002 3
- 2001 2
- 2000 4
- 1999 7
- 1998 8

[More ▶](#)

Topics & Subjects [Aa](#)

- Art, Modern 28
- Solomon R. Guggenheim 28
- Museum 28
- Painting 19
- Art 13
- Painting, Modern 12
- Kandinsky, Wassily, 1866-1944 11

[More ▶](#)

Collection

- American Libraries 204
- the Solomon R. Guggenheim Museum 204
- georgia.gallup Favorites 202
- francescaruffo Fa... 190

ON THE SPIRITUAL IN ART
BY WASSILY KANDINSKY

POINT AND LINE TO PLANE
BY WASSILY KANDINSKY

Point and line to plane : contribution to the
by Kandinsky, Wassily, 1866-1944; Dearstyne, Howard, tr;
60,014 | 70 | ★ | 0

On the spiritual in art :
First complete English
by Kandinsky, Wassily, 1866-1944; Rebay, Hilla, 1899-1967;
29,250 | 52 | ★ | 0

Gustav Klimt and Egon Schiele

KLIMT
SCHIELE

The Italian metamorphosis, 1943-1945

27,689 | 44 | ★ | 0

Kazimir Malevich : suprematism

Kazimir Malevich, 1878-1935; Drutt, Severinovich, 1988

23,428 | 24 | ★ | 1

CHINA 5000 YEARS

Collection

China, 5000 years : innovation and
by Lee, Sherman E.

The Great utopia : the
Russian and Soviet avant-gardes

The Great Utopia : the
Russia and Soviet avant-gardes

16,649 | 40 | ★ | 1

PICASSO AND THE WAR YEARS 1937-1945

Picasso and the war years, 1937-1945

Roy Lichtenstein

16,437 | 37 | ★ | 0

max ernst

Max Ernst : a retrospective

16,437 | 37 | ★ | 0

A NOTE ON SCRAPING

```
def get_book_list():
    '''get book list available from archives
    # there is infinite scroll on the page, so use this to load more: https://www.accordbox.com/blog/how-crawl-infinite-scrolling-pages-using-python/
    '''
    global alldata
    npages = 3
    alldata = pd.DataFrame(columns=[])

    for pagenum in range(1,npages):
        url = "https://archive.org/details/guggenheimmuseum?and%5B%5D=mediatype%3A%22texts%22&sort=titleSorter&page=" + str(pagenum)
        page = requests.get(url)
        soup = BeautifulSoup(page.text, 'html.parser')
        books = soup.find_all("div", {"class": "item-ia"})
        data = pd.DataFrame(columns[])
        data['html'] = books
        alldata = alldata.append(data)
```

STEP 0

SET-UP

EXECUTABLE NOTEBOOK (IF YOU WANT TO RUN REMOTE):

<https://notebooks.azure.com/anna-nicanorova/projects/wcaiconf-2019>

GITHUB CODE (IF YOU WANT TO RUN LOCALLY): https://github.com/AnnaNican/wcaiconf_2019

TEXT FILES : https://github.com/AnnaNican/wcaiconf_2019/tree/master/data/books

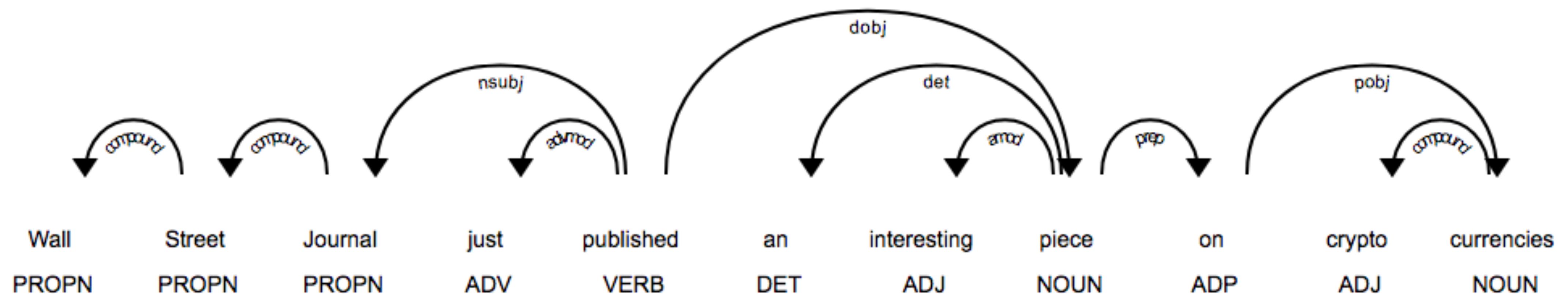
STEP 1

LOAD AND CLEAN DATA

	ASCII	ISO-8859-15 (latin-9)	CP-1252 (Windows 1252)	UTF-8
a	01100001	01100001	01100001	01100001
€	NA	10100100	10000000	11100010 10000010 10101100
¤	NA	NA	10100100	11000010 10100100

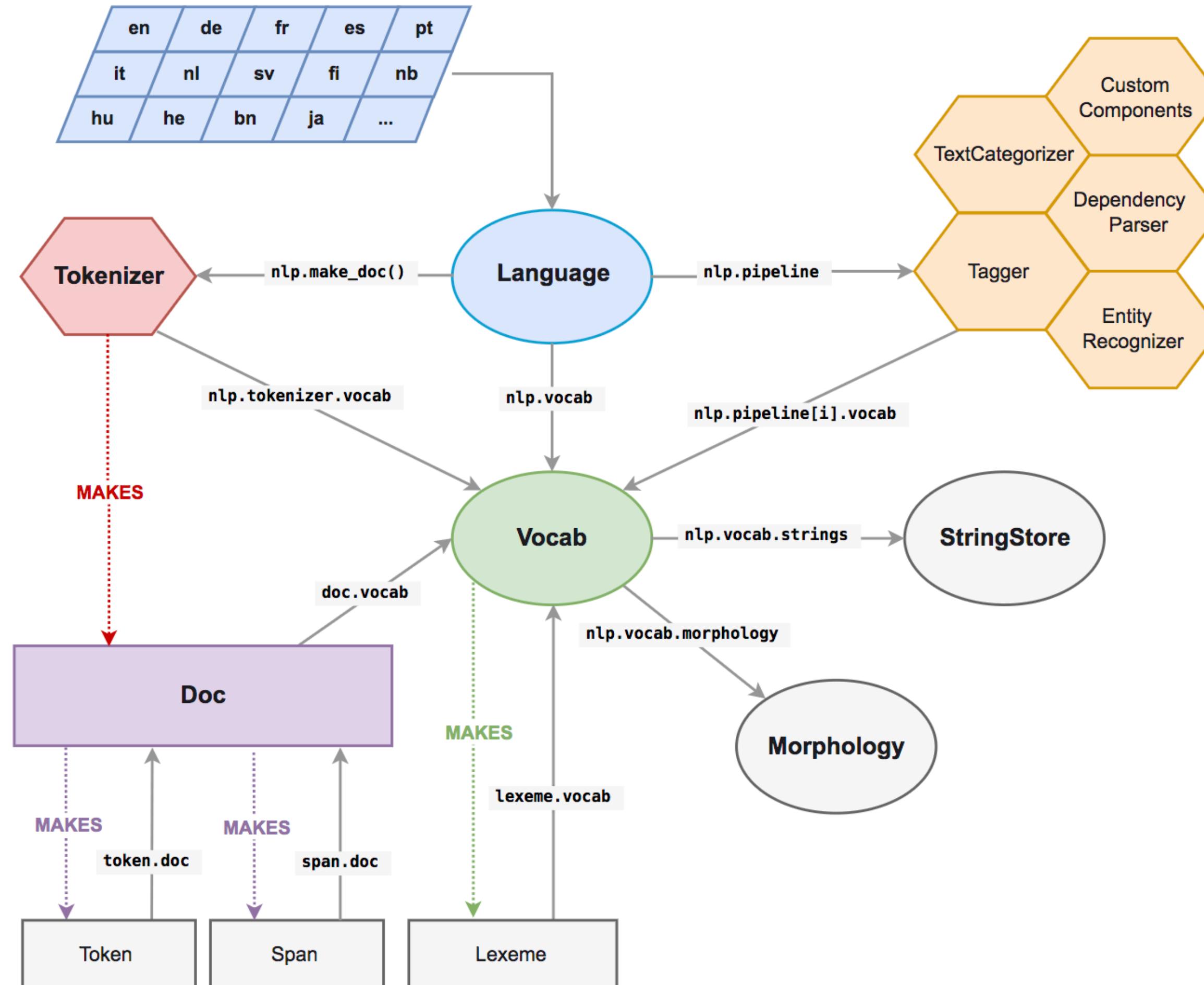
STEP 2

TEXT EXPLORATION



STEP 2

TEXT EXPLORATION



STEP 3

ENTITY EXTRACTION

The screenshot shows the displaCy Named Entity Visualizer interface. At the top, there's a navigation bar with links for About, Software, Demos (which is underlined), and Blog. A circular logo on the left contains the text "EXPLO NOISE". On the right, there's a Twitter icon.

The main title is "displaCy Named Entity Visualizer". Below it, a text input area contains a quote from Sebastian Thrun:

When Sebastian Thrun started working on self-driving cars at Google in 2007, few people outside of the company took him seriously. “I can tell you very senior CEOs of major American car companies would shake my hand and turn away because I wasn’t worth talking to,” said Thrun, now the co-founder and CEO of online higher education startup Udacity, in an interview with Recode earlier this week.

To the right of the text input is a section titled "Entity labels (select all)" with a grid of checkboxes for various entity types. Some checkboxes are checked, while others are unchecked:

- PERSON (checked)
- NORP (checked)
- FACILITY (unchecked)
- ORG (checked)
- GPE (checked)
- LOC (checked)
- PRODUCT (checked)
- EVENT (unchecked)
- WORK OF ART (checked)
- LANGUAGE (checked)
- DATE (checked)
- TIME (unchecked)
- PERCENT (checked)
- MONEY (checked)
- QUANTITY (checked)
- ORDINAL (unchecked)
- CARDINAL (unchecked)

Below the text input, there's a "Model ?" dropdown menu set to "English - en_core_web_sm (v2.0.0)".

At the bottom, the processed text is shown with entities highlighted and labeled:

When **Sebastian Thrun** PERSON started working on self-driving cars at **Google** ORG in **2007** DATE, few people outside of the company took him seriously. “I can tell you very senior CEOs of major **American** NORP car companies would shake my hand and turn away because I wasn’t worth talking to,” said **Thrun** PERSON, now the co-founder and CEO of online higher education startup Udacity, in an interview with **Recode** ORG earlier this week DATE.

A little less than a decade later DATE, dozens of self-driving startups have cropped up while automakers around the world clamor, wallet in hand, to secure their place in the fast-moving world of fully automated transportation.

STEP 4

TEXT SUMMARIZATION

Extractive Summarization

Select parts (typically sentences) of the original text to form a summary.



- Easier
- Too restrictive (no paraphrasing)
- Most past work is extractive

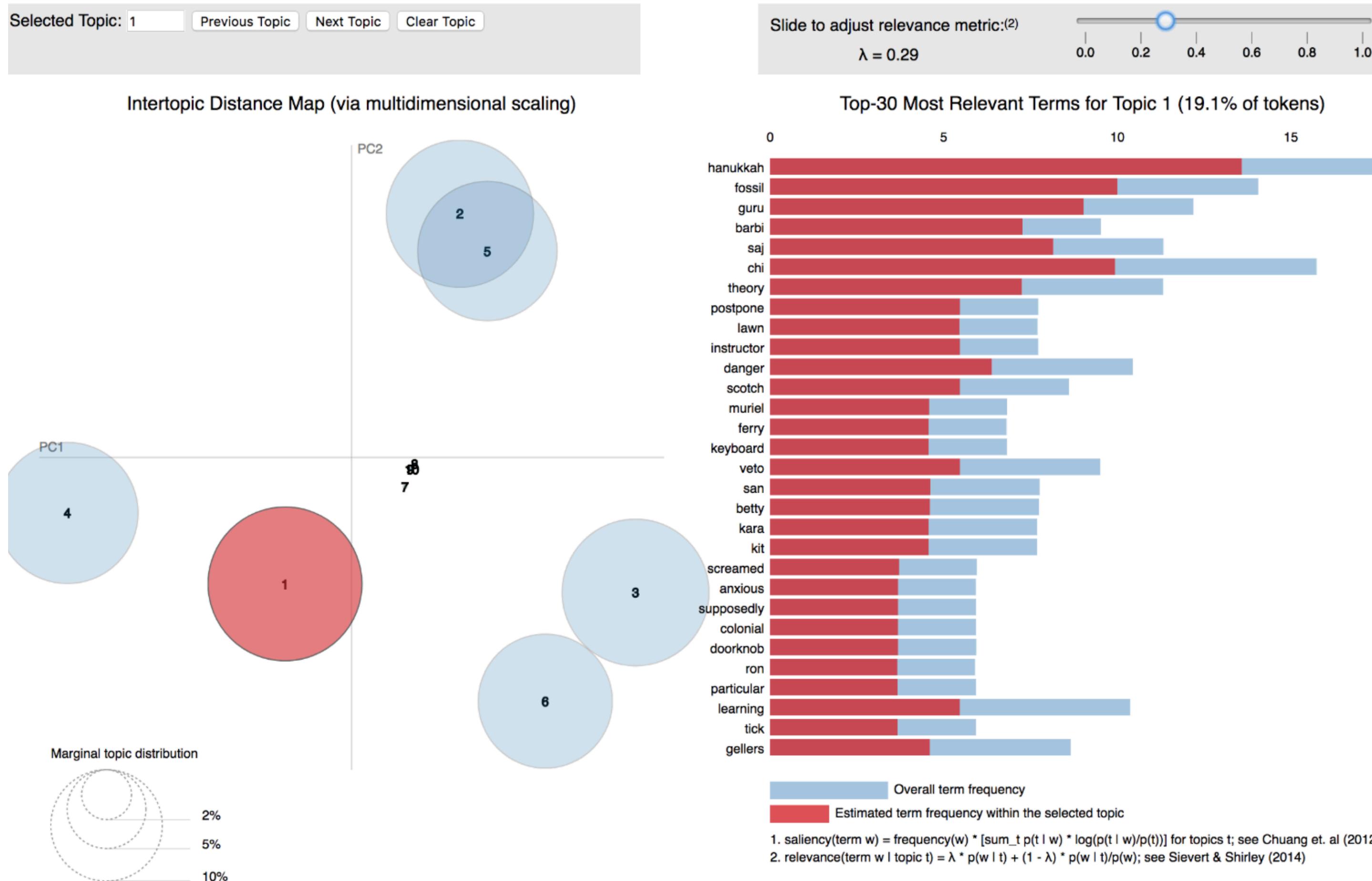
Abstractive Summarization

Generate novel sentences using natural language generation techniques.



- More difficult
- More flexible and human
- Necessary for future progress

STEP 5 TOPIC MODELING



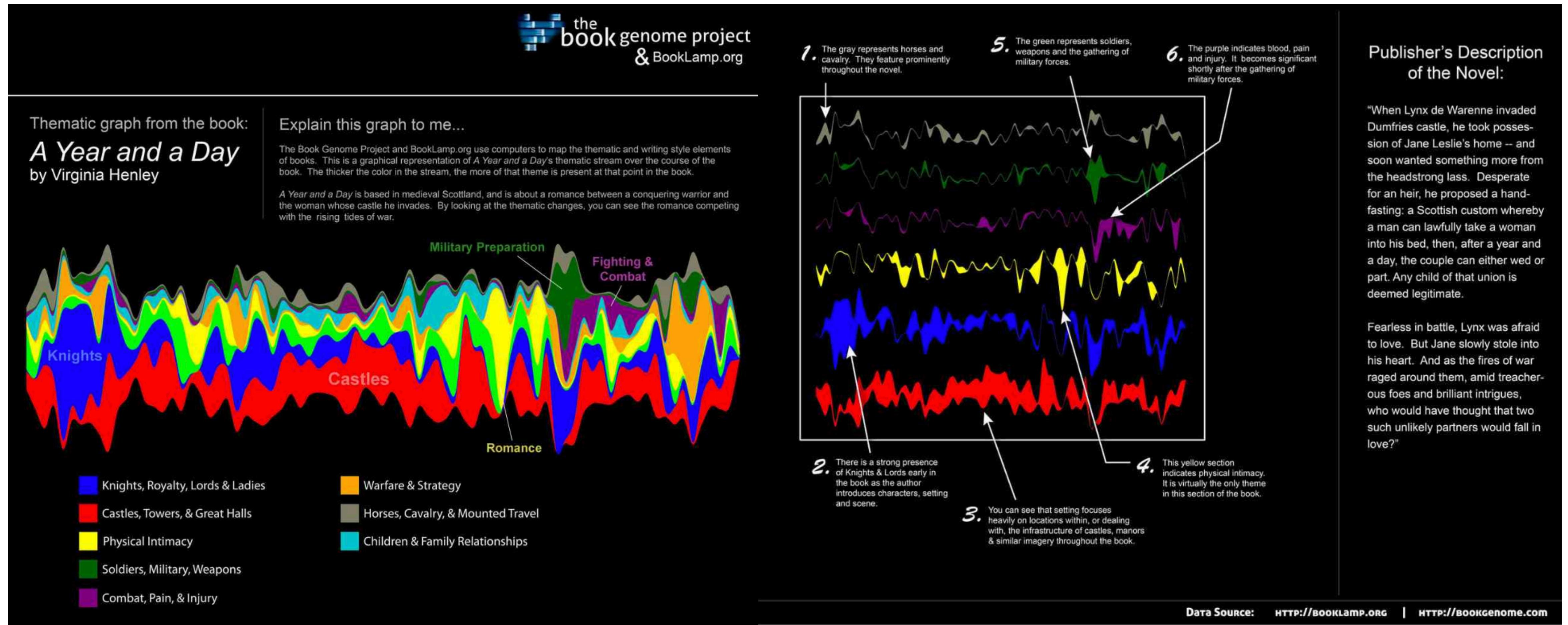
STEP 6

BEYOND WORD CLOUDS

STEP 6

BEYOND WORD CLOUDS

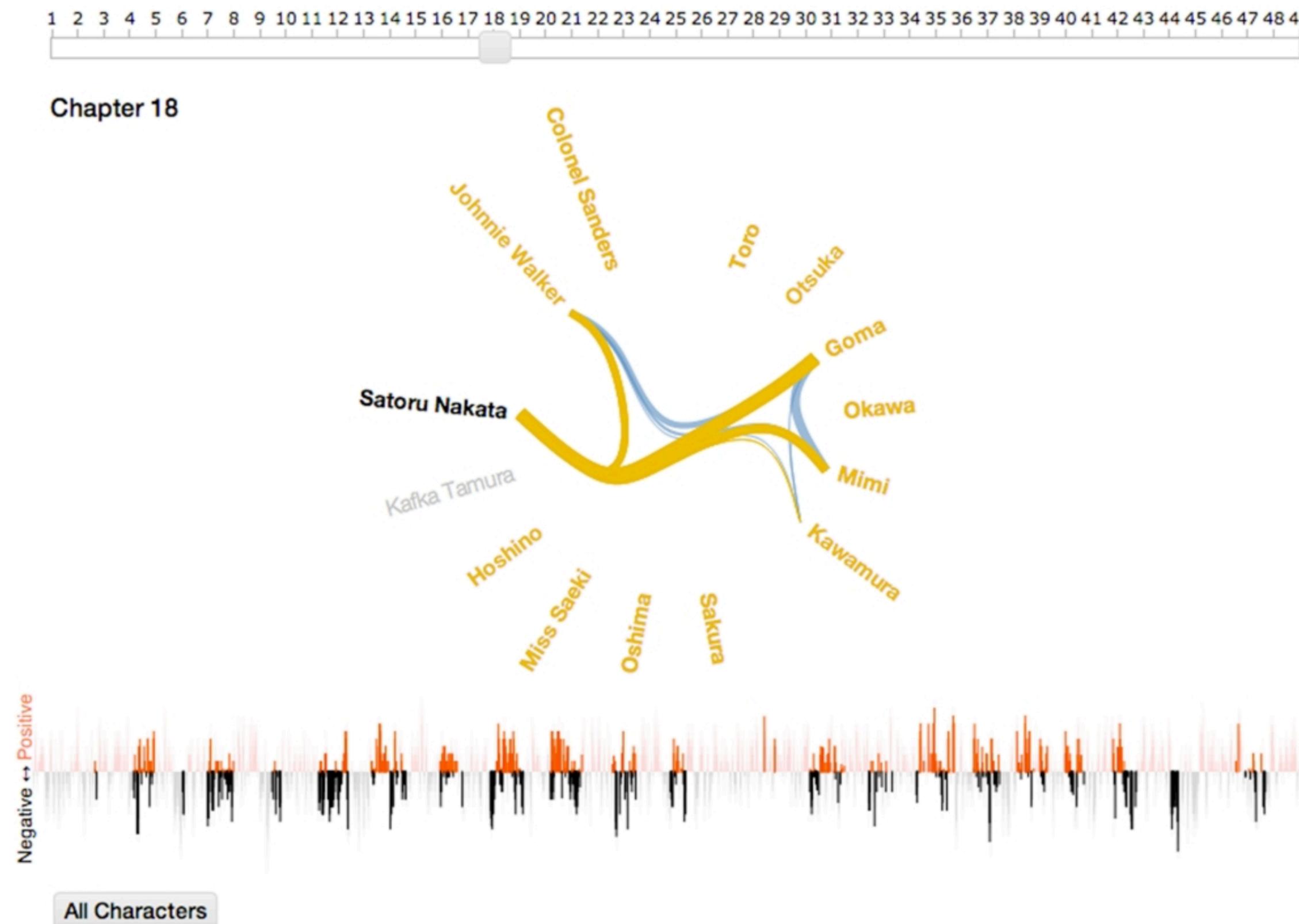
THEMATIC FLOW OF THE DOCUMENT



STEP 6

BEYOND WORD CLOUDS

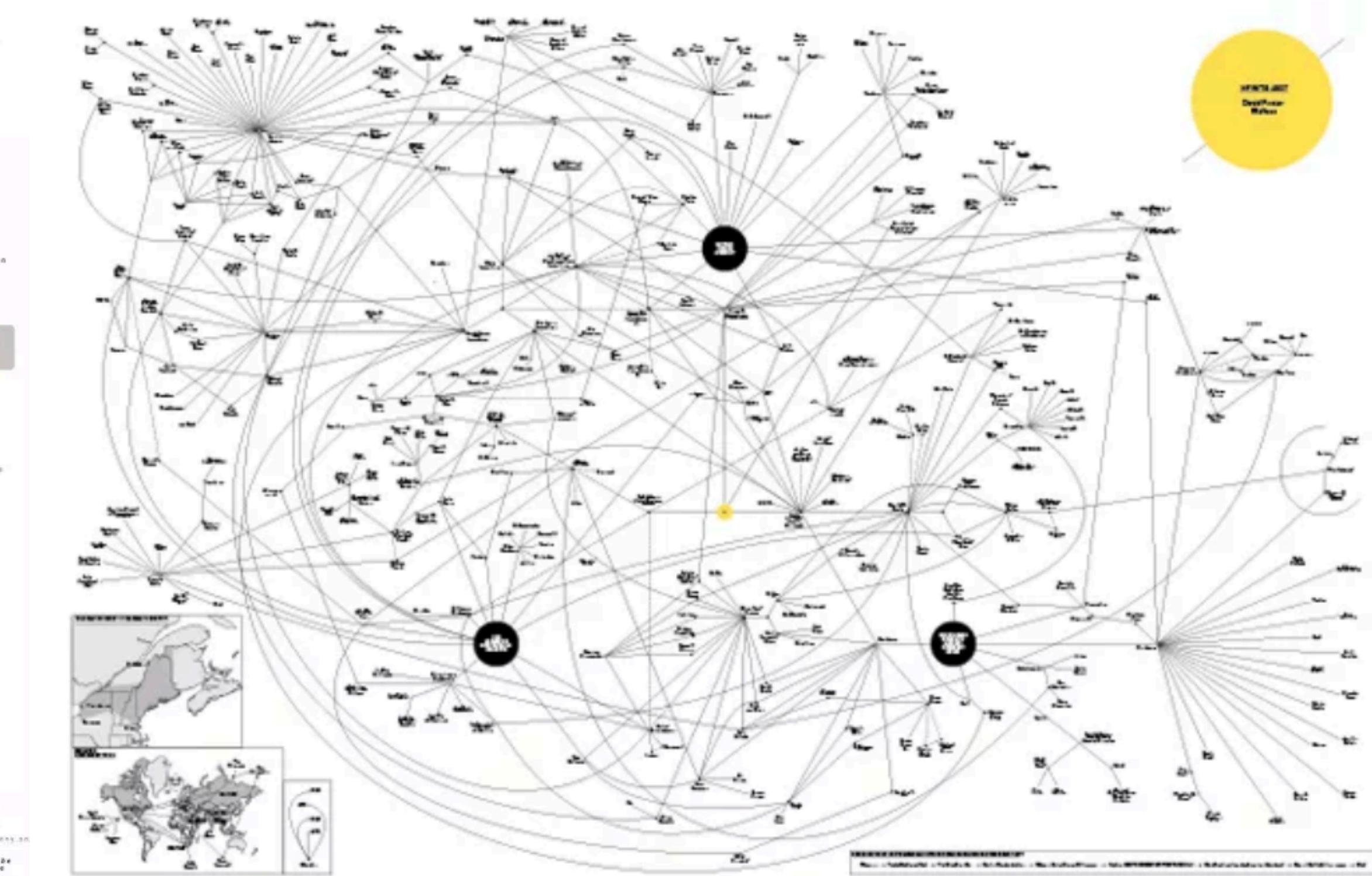
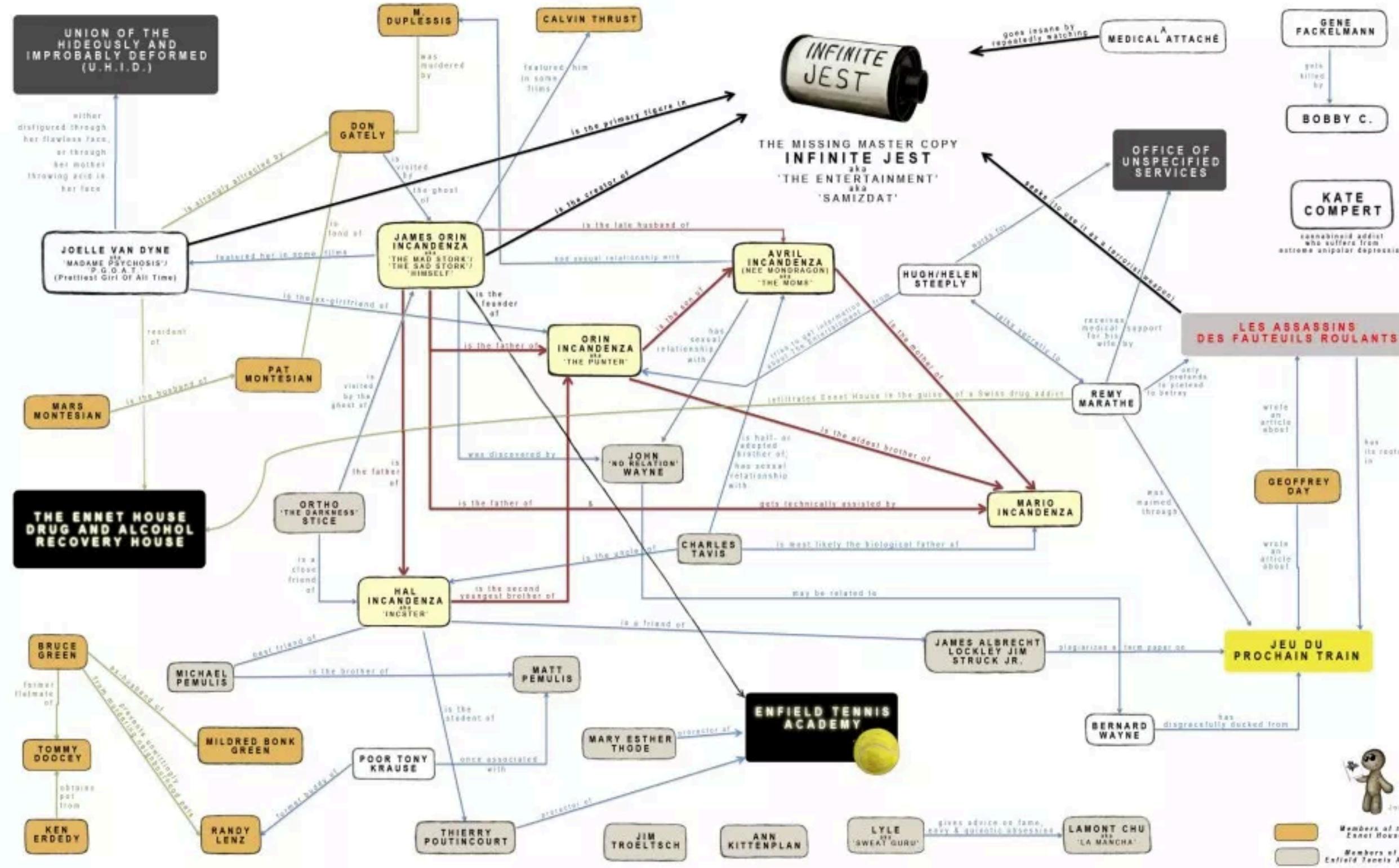
NARRATIVE LINES



STEP 6

BEYOND WORD CLOUDS

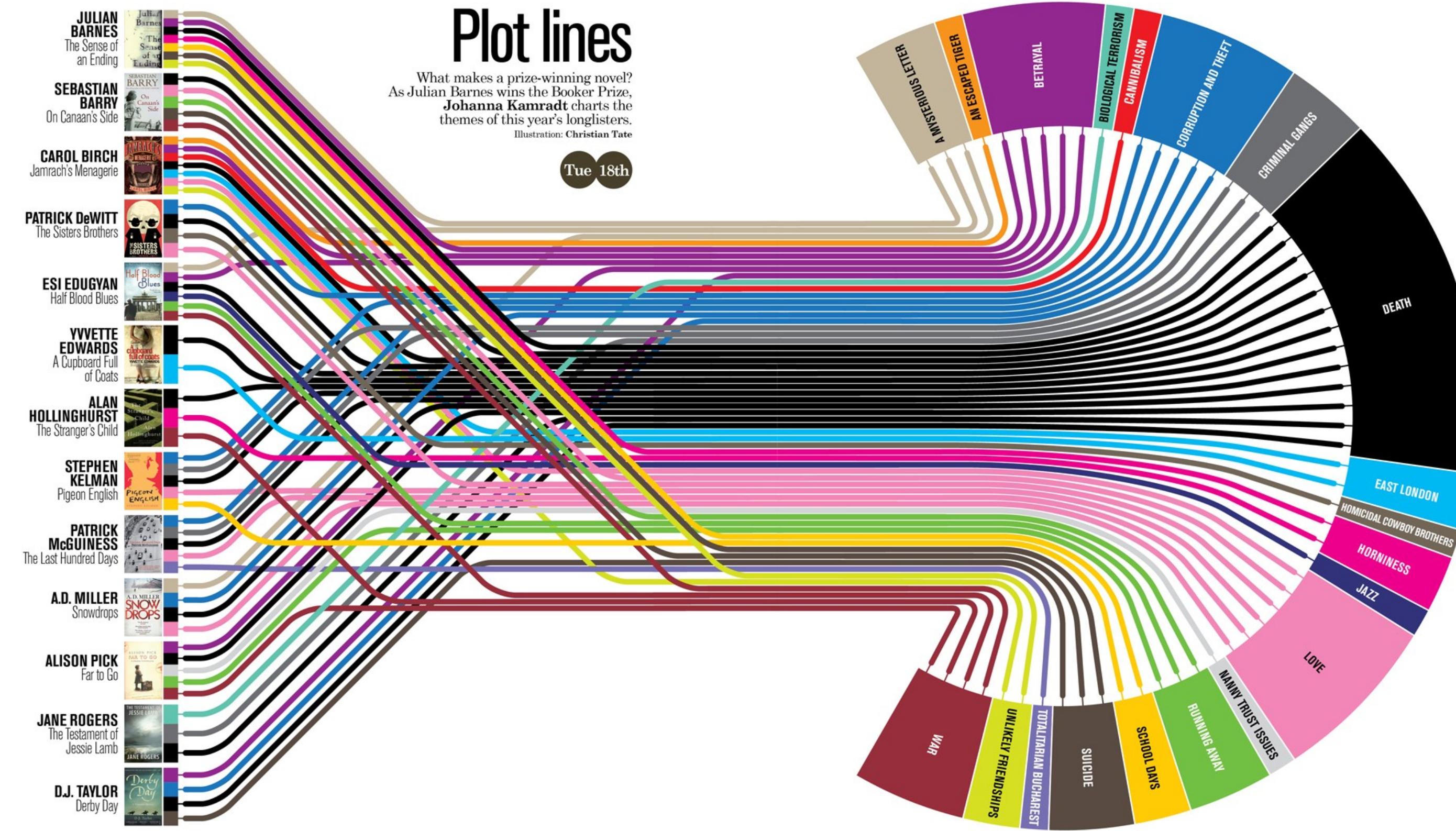
CONNECTION BETWEEN CHARACTERS



STEP 6

BEYOND WORD CLOUDS

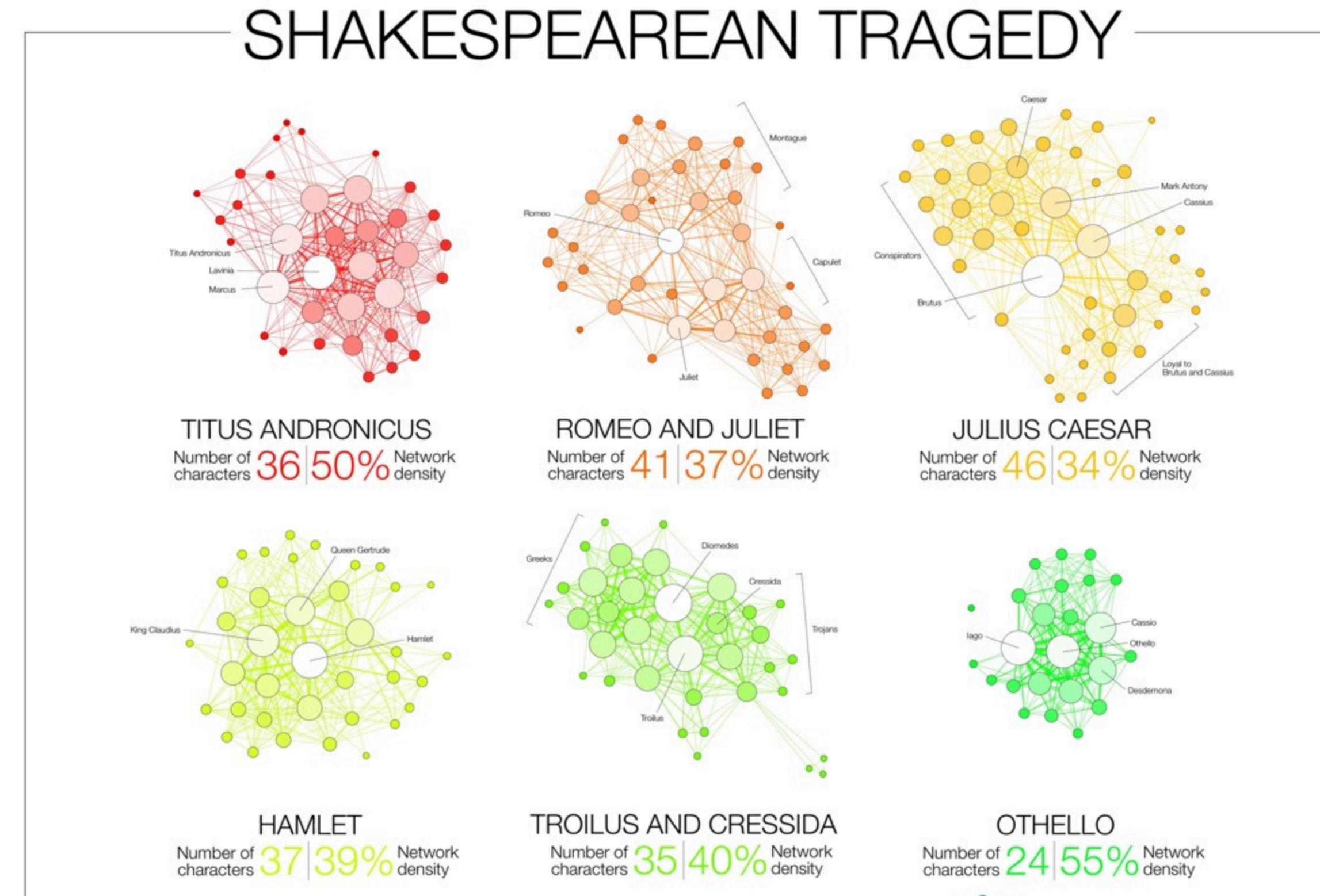
PLOT/ STORY LINES



STEP 6

BEYOND WORD CLOUDS

CHARACTER NETWORK GRAPHS



SUMMARY

Summary

- word clouds are bad
 - the visualization is hard to use
 - they mangle ideas by counting words
- alternative to counting words are methods that require additional text work
- Language work is still falling short in areas paraphrasing, summarization and
- There are plenty of paid APIs and free libraries in Python to perform basic ML tasks on text
- There are a lot of great visualization alternatives like narrative, story lines, but their require a lot of additional algorithmic processing
- With processes like Entity Extraction, Topic Modeling and summarization we can create better ways to highlight the ideas on unstructured text data