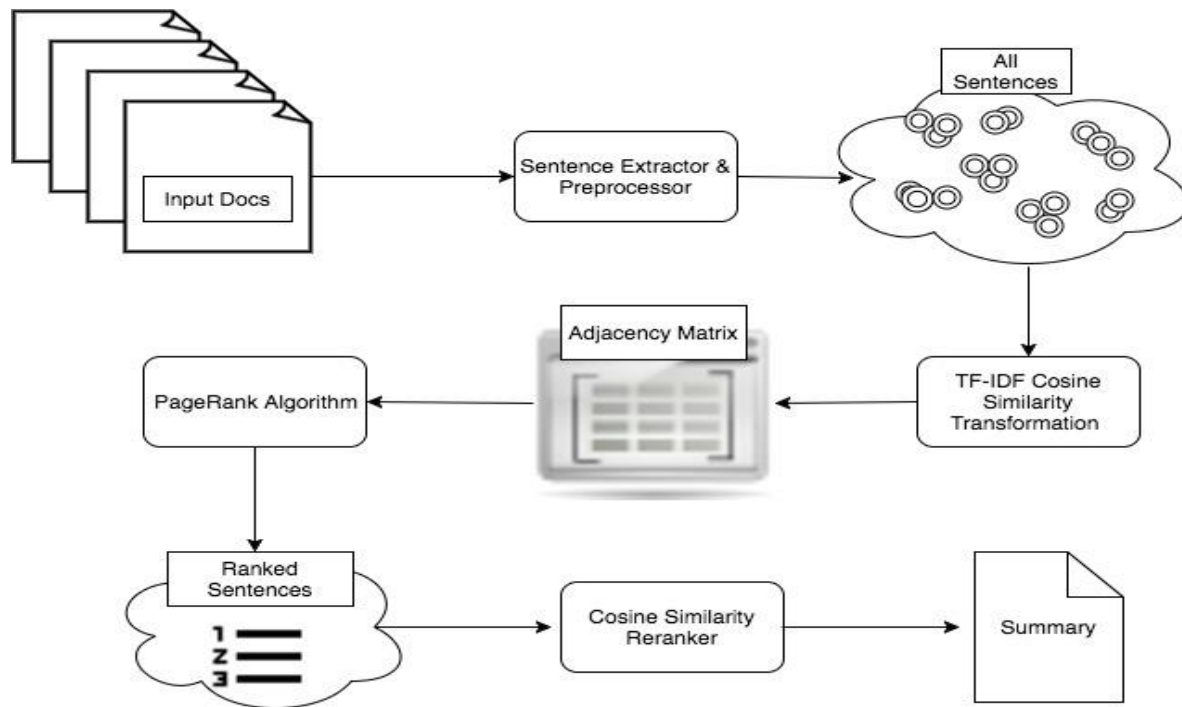# Multi-Document Summarization

Eslam Elsawy, Audrey Holmes, Masha Ivenskaya

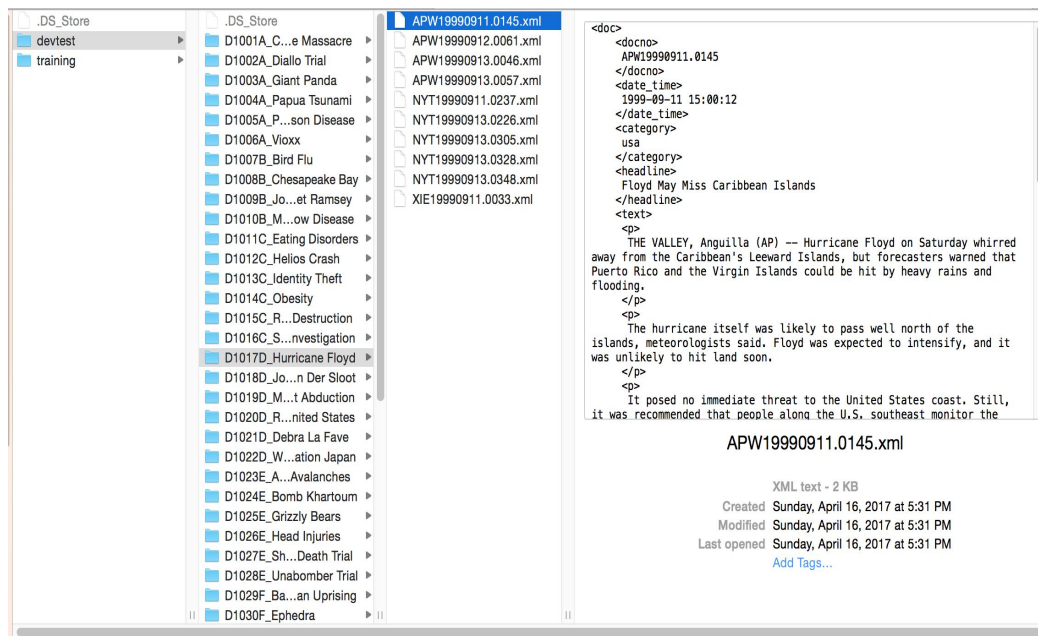# System Architecture

# Dataset Pre-Processing

Issues:

- Non-uniform file naming scheme
- Encoded characters (e.g: &amp )
- Not rooted xml files

What we did:

- Separate cleaning module

Gain:

- Easier to explore the dataset
- Better running times

# Content Selection - LexRank - Sentence Similarity

- All sentences ⟶ adjacency matrix by similarity score

```
[[ 1.          0.17963573  0.27465779  0.06461521  0.21442665]
 [ 0.17963573  1.          0.19960539  0.10907585  0.23165034]
 [ 0.27465779  0.19960539  1.          0.14474842  0.26945008]
 [ 0.06461521  0.10907585  0.14474842  1.          0.05448286]
 [ 0.21442665  0.23165034  0.26945008  0.05448286  1.        ]]
```

- Approaches to measure pairwise similarity:

  1.Cosine Similarity

  **2.Cosine Similarity with TF-IDF (highest Rouge scores)**

  3.Doc2Vec model trained on NLTK Reuters Corpus

# Content Selection - LexRank - Algorithm

- Binarize similarity matrix **M** with 0.15 threshold.
- Implement LexRank using Power Method:

1. Initialize $p_0$ vector with uniform distribution.
2. Iteratively update $p_t$ such that $p_t = M^T p_{t-1}$ until $||p_t - p_{t-1}||$ is sufficiently small.
3. Return $p_t$.

# Information Ordering & Content Realization

- Input: sentences sorted by score
- TF-IDF cosine similarity ordering [1]
- One Sentence per line
- Max 100 words

Threshold?

- At 1.0 => R-1 R = 0.239
- At 0.2 => R-1 R = 0.258 (+0.019)

# Results

- Runtime ~= 5 - 10 seconds
- Best results:
  - Content Selection: TF-IDF cos sim
  - Ordering: threshold = 0.2

| ROUGE-L | Recall |
|---------|--------|
| ROUGE-1 | 0.25785 |
| ROUGE-2 | 0.07108 |
| ROUGE-3 | 0.02438 |
| ROUGE-4 | 0.00847 |

Sample output
Topic: JonBenet Ramsey Murder

```
1   Hunter took the JonBenet case to the grand jury
    shortly after a former Boulder police detective on
    the case and three former friends of the Ramseys
    publicly demanded that Colorado's governor, Roy
    Romer, replace Hunter on the case with a special
    prosecutor.
2   Although the police chief and district attorney both
    have said that the Ramseys fall under ``the umbrella
    of suspicion,'' they have not formally named any
    suspects.
3   Burke was in the family's Boulder home when 6-year-
    old JonBenet was found beaten and strangled Dec. 26,
    1996.
4   Police say her parents, John and Patsy Ramsey,
    remain under suspicion.
```

# Issues and Successes

Success: Implementing tf-idf raised ROUGE-1 recall from 0.13 to 0.26.

Problems that need fixing:

- Meta-info like:
  - LITTLETON, Colo. (AP) --
  - NEW YORK _ The parents of Amadou Diallo plan to meet with the Bronx district
- Quotes without attributions
- Pronouns without reference
- Ordering of the sentences within summary
- Long sentences

# References

[1] Radev, Dragomir R., et al. "MEAD-A Platform for Multidocument Multilingual Text Summarization." *LREC*. 2004.

[2] Erkan, Günes, and Dragomir R. Radev. "Lexrank: Graph-based lexical centrality as salience in text summarization." Journal of Artificial Intelligence Research 22 (2004): 457-479.

[3] Lin, Chin-Yew. "Rouge: A package for automatic evaluation of summaries." Text summarization branches out: Proceedings of the ACL-04 workshop. Vol. 8. 2004.

# Questions ?