# ARElight: Context Sampling of Large Texts for Deep Learning Relation Extraction

Nicolay Rusnachenko, Huizhi Liang, Maxim Kolomeets, Lei Shi

{name.surname}@newcastle.ac.uk

**Newcastle University**

## Outline

1. Information Retrieval and Large Texts
2. Existed Systems
3. Processing and Analysis of Large Texts

Outline
Introduction
Systems
ARElight

**Information Retrieval**
Large Documents
Analysis

## Information Retrieval (IR)

*Information Extraction* – one of the direction in Natural Language Processing (NLP) aimed on retrieving content from structuring textual information:

- **Objects** (entities, events)
- Establishing **relations** between objects (semantic, <u>sentiment</u>)[1]

[1] Iris Hendrickx et al. "SemEval-2010 Task 8: Multi-Way Classification of Semantic Relations between Pairs of Nominals". In: *Proceedings of the 5th International Workshop on Semantic Evaluation.* Uppsala, Sweden: Association for Computational Linguistics, July 2010, pp. 33–38. URL: https://aclanthology.org/S10-1006.

Outline
Introduction
Systems
AORElight

Information Retrieval
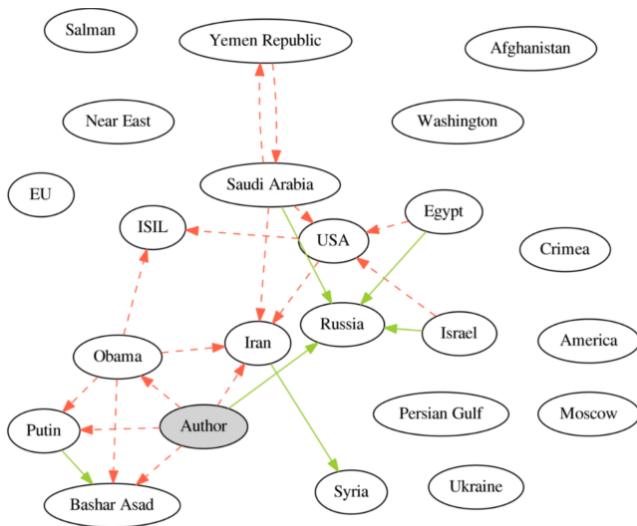**Large Documents**
Analysis

# Large Document[2]



Russia criticized Belarus for permitting Georgian President Mikheil Saakhashvili to appear on Belorussian television. "The appearance was an unfriendly step towards Russia," the speaker of Russian parliament Boris Gryzlov said. ... Saakhashvili announced Thursday that he did not understand Russia's claims. Moscow refused to have any business with Georgia's president after the armed conflict in 2008 ...
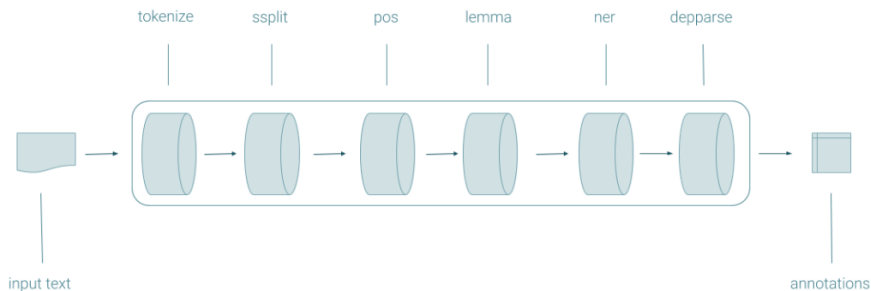
[2] Eunsol Choi et al. "Document-level Sentiment Inference with Social, Faction, and Discourse Context". In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics, Aug. 2016, pp. 333–343. DOI: 10.18653/v1/P16-1032. URL: https://aclanthology.org/P16-1032.

Outline
Introduction
Systems
ARElight

Information Retrieval
Large Documents
Analysis

# Representation and ways to Analyse

# Existed Systems

Outline
Introduction
Systems
ARElight

Pipeline-based
Target-Oriented Systems
Attention
Summary of Limitations

# Pipeline-based Concept[3]



[3] Christopher Manning et al. "The Stanford CoreNLP Natural Language Processing Toolkit". In: *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Baltimore, Maryland: Association for Computational Linguistics, June 2014, pp. 55–60. DOI: 10.3115/v1/P14-5010. URL: https://aclanthology.org/P14-5010.

## Target-Oriented Systems

$$\text{TEXT} \rightarrow \text{OBJECTS}$$
$$[\text{TEXT, OBJECTS}] \rightarrow \text{RELATIONS}$$

Outline
Introduction
Systems
ARElight

Pipeline-based
**Target-Oriented Systems**
Attention
Summary of Limitations

# OpenNRE[4]



**Sentence-level RE**

*Ernest Hemingway* was raised in *Oak Park, Illinois* ⟹ *[Ernest Hemingway]* —place of birth→ *[Oak Park, Illinois]*

**Bag-level RE**

In 1921, *Ernest Hemingway* married *Hadley Richardson*, the first of his four wives

*Hadley Richardson* was the first wife of American author *Ernest Hemingway*

... ... ...

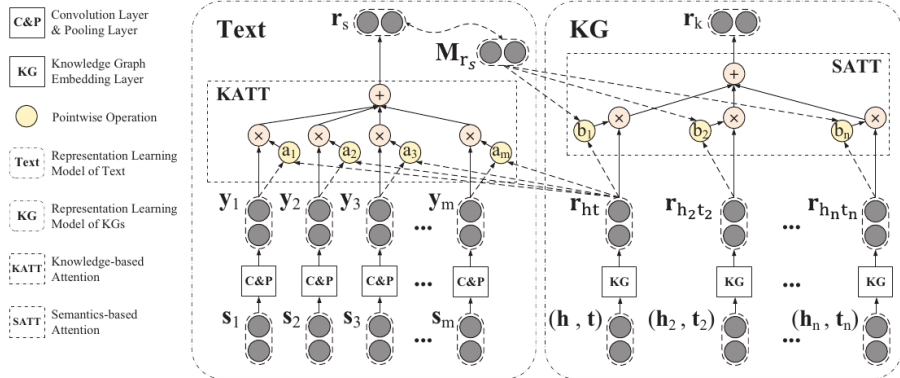⟹ *[Ernest Hemingway]* —spouse→ *[Hadley Richardson]*

**Document-level RE**

*Mark Twain* and *Olivia Langdon* corresponded throughout 1868. She rejected his first marriage proposal, but they were married in Elmira, New York in February 1870. Then, Twain owned a stake in the Buffalo Express newspaper and worked as an *editor* and *writer*. While they were living in *Buffalo*, their son *Langdon* died of diphtheria at the age of 19 months. They had three daughters: *Susy Clemens*, *Clara Clemens*, and *Jean Clemens*.
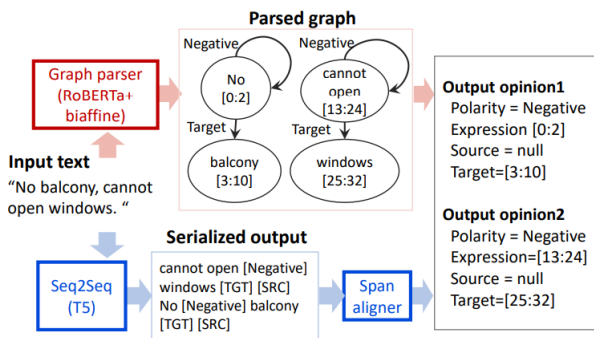
sibling of   spouse of   occupation
child of   place of death

*[Olivia Langdon]*   *[Susy Clemens]*
*[Mark Twain]*   *[Clara Clemens]*
*[editor]*   *[Langdon]*   *[Jean Clemens]*
*[writer]*   *[Buffalo]*

[4] Xu Han et al. "OpenNRE: An open and extensible toolkit for neural relation extraction". In: *arXiv preprint arXiv:1909.13078* (2019).

Outline    Pipeline-based
Introduction    **Target-Oriented Systems**
**Systems**    Attention
ARElight    Summary of Limitations

# JointNRE[5]



[5] Xu Han, Zhiyuan Liu, and Maosong Sun. "Neural Knowledge Acquisition via Mutual Attention between Knowledge Graph and Text". In: *Proceedings of AAAI*. 2018.

Outline
Introduction
Systems
ARElight

Pipeline-based
**Target-Oriented Systems**
Attention
Summary of Limitations

# T5 graph-based transformer by Hitachi[6]

[6] Gaku Morio et al. "Hitachi at SemEval-2022 Task 10: Comparing Graph- and Seq2Seq-based Models Highlights Difficulty in Structured Sentiment Analysis". In: *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*. Seattle, United States: Association for Computational Linguistics, July 2022, pp. 1349–1359. DOI: 10.18653/v1/2022.semeval-1.188. URL: https://aclanthology.org/2022.semeval-1.188.

Outline
Introduction
**Systems**
ARElight

Pipeline-based
Target-Oriented Systems
**Attention**
Summary of Limitations

## Attention

For input $X \in R^N$:

- $O(N^2)$ original self-attention[7] computation complexity;

How to address this problem:

1. Sparse version of Self-attention;

2. #1 with Global Attention;

3. **Structuring**[8] – limit attention on sentences, paragraphs, etc. via masking.

512 (BERT, T5) $\rightarrow$ 1K (ETC), 4K/8K/16K (LongT5), 32K (ChatGPT4)

---

[7] Ashish Vaswani et al. "Attention is all you need". In: *Advances in neural information processing systems* 30 (2017).

[8] Joshua Ainslie et al. "ETC: Encoding long and structured inputs in transformers". In: *arXiv preprint arXiv:2004.08483* (2020).

# Summary of Limitations

- Pipeline-based:
    - is considered for the **whole document**.[3]
- Target-oriented:
    - Input Size Limitations.[7] (512-32K tokens at present)

Outline
Introduction
Systems
ARElight

Concept
Sampler
Inference
Demo

# Demo

Outline
Introduction
Systems
ARElight

**Concept**
Sampler
Inference
Demo

## Overall Demo Concept

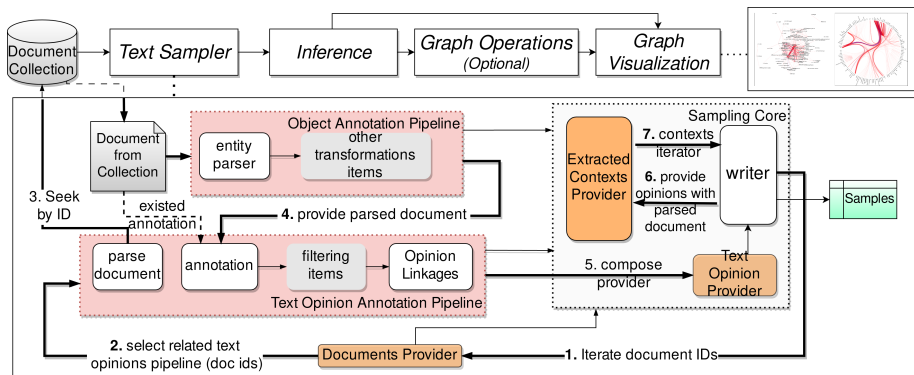**Text Sampler** – extract small portions of text (frames)[9] from (i) large document (samples) and/or (ii) collection of documents[10].



We consider sentiment analysis problem with classes: **positive**, **negative**.

---

[9] Heike Adel et al. "DERE: A task and domain-independent slot filling framework for declarative relation extraction". In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. 2018, pp. 42–47.

[10] Adam Roberts et al. "Scaling Up Models and Data with t5x and seqio". In: *arXiv preprint arXiv:2203.17189* (2022). URL: https://arxiv.org/abs/2203.17189.

Outline
Introduction
Systems
ARElight

Concept
**Sampler**
Inference
Demo

# Architecture of the Sampler and Overall Workflow



**Two declarative pipelines**[1] for separate annotation of **Objects** and **Relations**.

1 `https://github.com/nicolay-r/AREkit/wiki/Task-Schemata`

## Inference

- Sentiment Relation Extraction[11]
- Using `OpenNRE`[4] and BERT-based models as inference.

---

[11] Nicolay Rusnachenko. "Language Models Application in Sentiment Attitude Extraction Task". Russian. In: *Proceedings of the Institute for System Programming of the RAS (Proceedings of ISP RAS), vol.33.* 3. 2021, pp. 199–222.

Outline
Introduction
Systems
ARElight

Concept
Sampler
**Inference**
Demo

# Serializing Graphs

Forced[a]



Radial[a]



---

a https://observablehq.com/@d3/
force-directed-graph/2

---

a https://observablehq.com/@d3/
hierarchical-edge-bundling

Outline
Introduction
Systems
ARElight

Concept
Sampler
Inference
Demo

# System Demo[2]



(c) visualisation model selector

(a) dataset selector

(b) visualisation options

(d) force layout visualization model

(d) radial layout visualization model

# Thank you for attention!

https://nicolay-r.github.io