# Overview of Natural Language Processing (NLP) in biomedical and cancer research

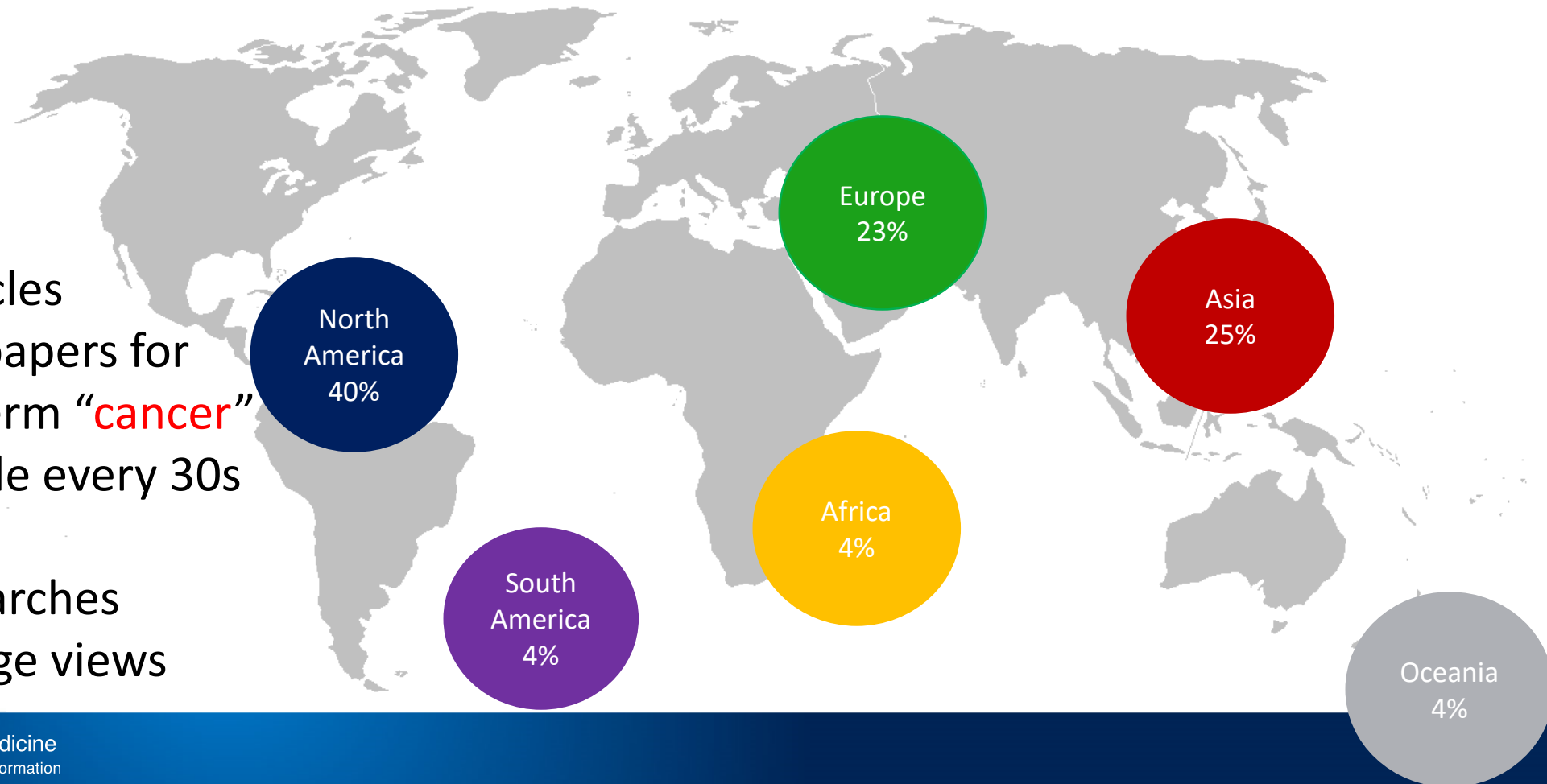**Yifan Peng, Qingyu Chen**

**NCBI/NLM/NIH**

U.S. National Library of Medicine
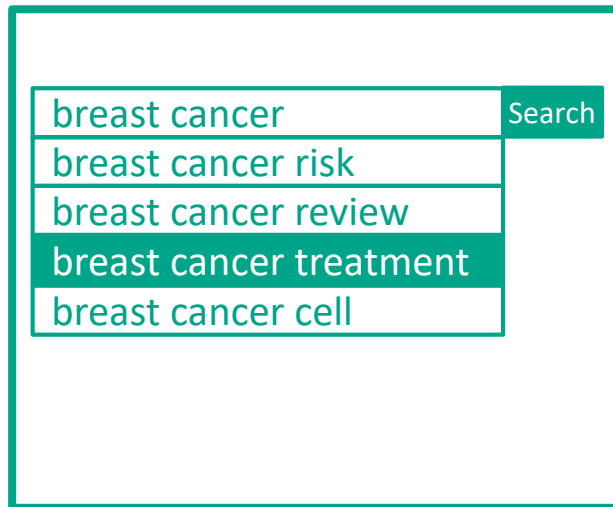National Center for Biotechnology Information

# PubMed: biomedical literature search engine



- ➢ 28+ million articles
  - ➢ 3.8 million papers for the query term "cancer"
  - ➢ A new article every 30s
- ➢ Daily usage
  - ➢ 3 million searches
  - ➢ 9 million page views

North America 40%

South America 4%

Europe 23%
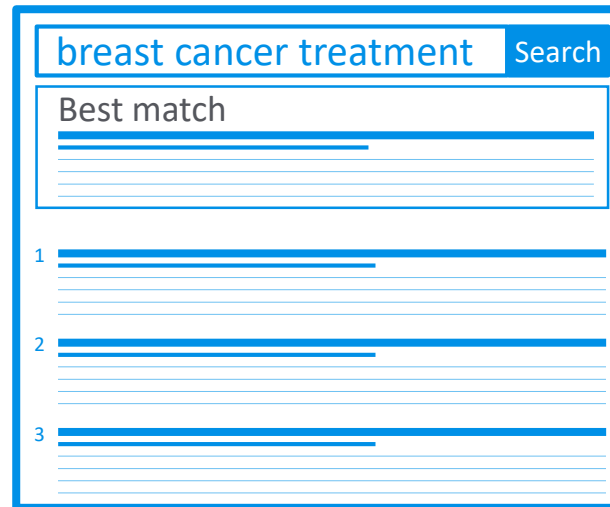
Africa 4%

Asia 25%

Oceania 4%

# What happens when you click the "Search" button on PubMed

**Search page**



- Spelling correction
- Query expansion
- Query suggestion

**Result page**



- Navigational searches
- Relevance match

**Article page**



- Related articles suggestions
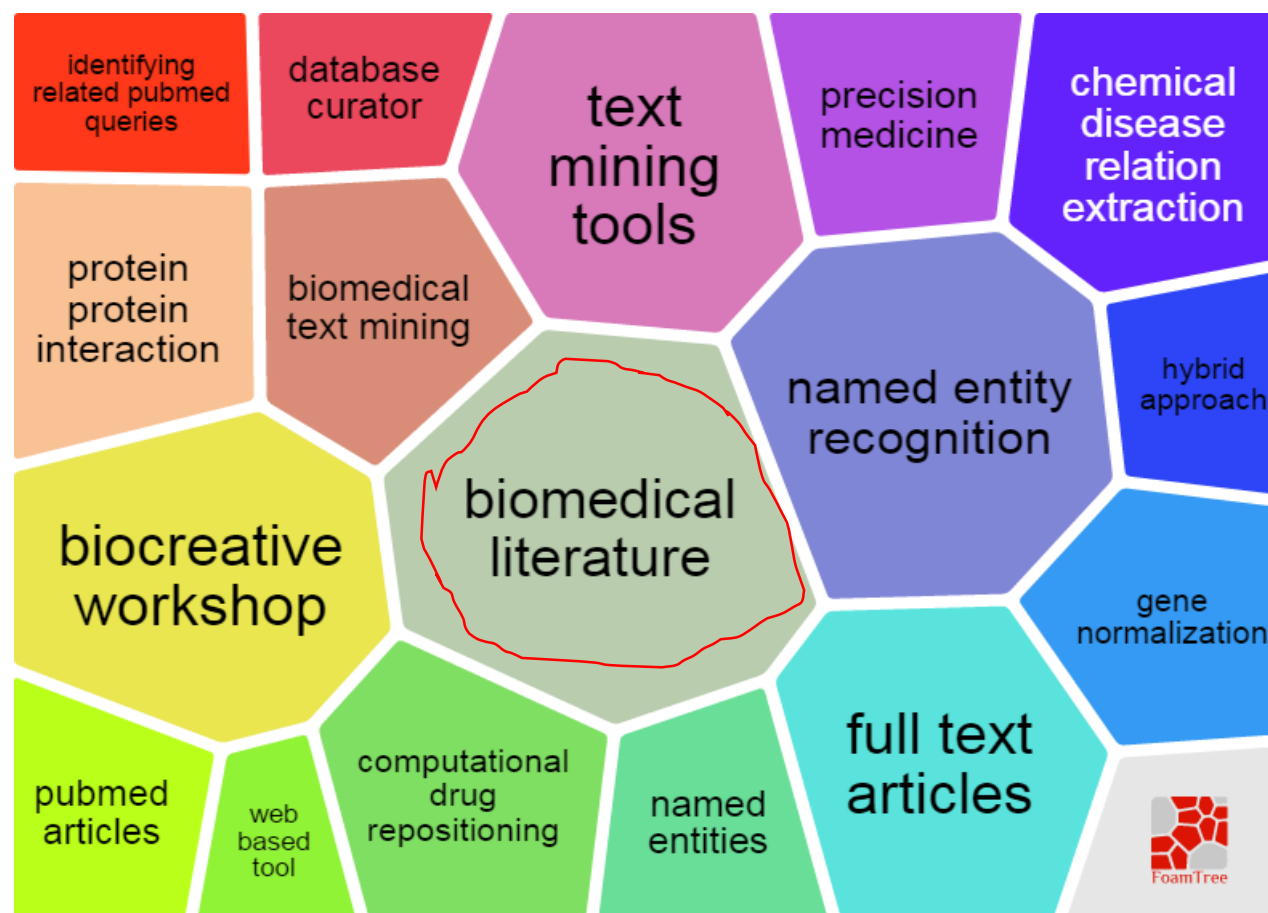- Author name disambiguation
- Citation sensor

https://www.ncbi.nlm.nih.gov/labs/pubmed/
Fiorini, et al., Best Match: New relevance search for PubMed, PLoS Biology, 2018

# What is Natural Language Processing (NLP)?

- **Natural language processing** is a field at the intersection of
  - Computer science
  - Artificial intelligence
  - Linguistics

- **Goal**: for computers to "understand" natural language in order to perform tasks that are useful

# NLP goes beyond the biomedical literature

- Biomedical Literature
- Clinical notes, EMRs
  - Chest X-ray & retinal images

# NLP is important for cancer research

- Finding relevant literature
- Extracting important entities such as <span style="color:red">cancers</span> and <span style="color:red">treatment</span> mentioned in literature
- Understanding the semantics of language
- Classifying related documents for manual curation
- Helping image analysis

# NLP helps extract information

## LitVar: Extracting mutation information from articles



Allot et al., LitVar: a semantic search engine for linking genomic variant data in PubMed and PMC. Nucleic Acids Research. 2018

https://www.ncbi.nlm.nih.gov/CBBresearch/Lu/Demo/LitVar/

# LitVar is supported by PubTator

Named entity recognition tool

https://www.ncbi.nlm.nih.gov/CBBresearch/Lu/Demo/PubTator/
Wei et. al., PubTator: a Web-based text mining tool for assisting Biocuration, Nucleic acids research, 2013.

U.S. National Library of Medicine
National Center for Biotechnology Information

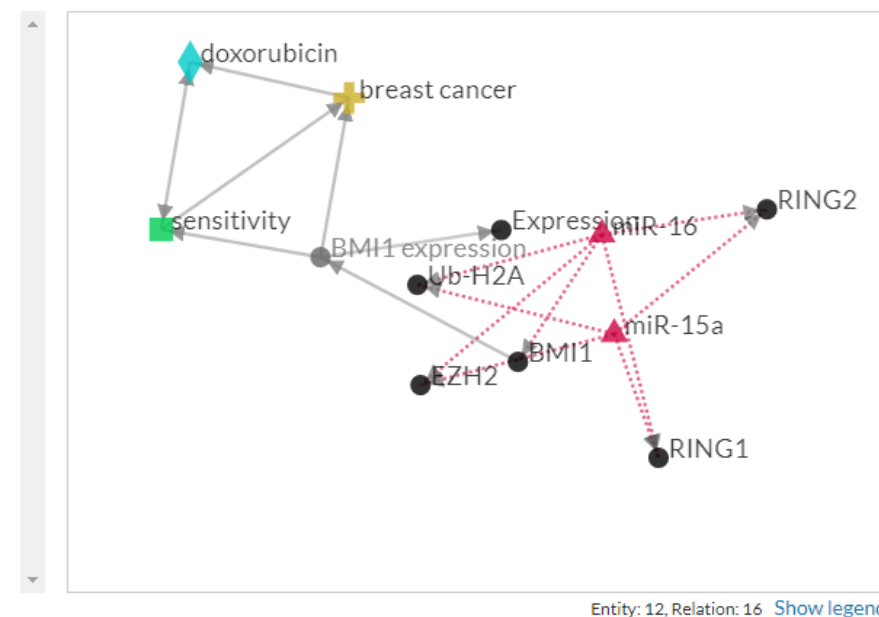# From named entities to relations

PMID: 28655885  RLIMS-P 0  eFIP 0  miRTex 9  eGARD 1                                        Issue Report

Abstract

1. miR-15a/miR-16 down-regulates BMI1, impacting Ub-H2A mediated DNA repair and breast cancer cell sensitivity to doxorubicin.

2. The B-lymphoma Moloney murine leukemia virus insertion region-1 protein (BMI1) acts as an oncogene in various cancers, including breast cancer.

3. Recent evidence suggests that BMI1 is rapidly recruited to sites of DNA double strand breaks where it facilitates histone H2A ubiquitination and DNA double strand break repair by homologous recombination.

4. Here we show that miR-15a and miR-16 expression is decreased during the initial period after DNA damage where it would otherwise down-regulate BMI1, impairing DNA repair.

5. Elevated miR-15a and miR-16 levels down-regulated BMI1 and other polycomb group proteins like RING1A, RING1B, EZH2 and also altered the expression of proteins associated with the BMI1 dependent ubiquitination pathway.

6. Antagonizing the expression of miR-15a and miR-16, enhanced BMI1 protein levels and increased DNA repair.

7. Further, overexpression of miR-15a and miR-16 sensitized breast cancer cells to DNA damage induced by the chemotherapeutic drug doxorubicin.

8. Our results suggest that miR-15a and miR-16 mediate the down-regulation of BMI1, which impedes DNA repair while elevated levels can sensitize breast cancer cells to doxorubicin leading to apoptotic cell death.

9. This data identifies a new target for manipulating DNA damage response that could impact the development of improved therapeutics for breast cancer.

Entity: 12, Relation: 16  Show legend

https://research.bioinformatics.udel.edu/itextmine/

# LitSense: sentence-level retrieval



https://www.ncbi.nlm.nih.gov/research/litsense/

# NLP helps scale up manual curation



Lee et al., Scaling up data curation using deep learning: An application to literature triage in genomic variation resources, PLoS Comp Biol. 2018.

# NLP helps biomedical image analysis

**DeepLesion**: Lesion annotation, detection, and retrieval



Unchanged **large** **nodule** bilaterally for example
**right lower lobe** [OTHER BOOKMARK] and
**right middle lobe** [BOOKMARK]

Size: large
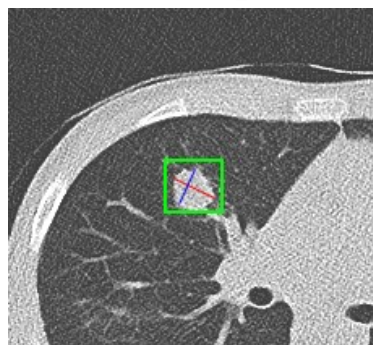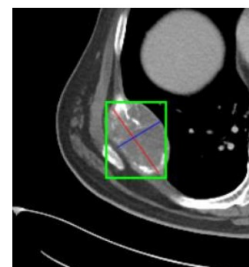Type: nodule
Body part: right mid lobe
Unrelated: right lower lobe



(c) Expanded **right** posterior **rib** lesion

**Posterior left rib mass**

**Right chest wall mass**

(d) Complex **retroperitoneal mass** involving the region of the **tail and body of the pancreas**

**Pancreatic tail mass**

Centrally **hypoattenuating mass** within the **pancreatic tail**

Yan et al., Fine-grained lesion annotation in CT images with knowledge mined from radiology reports., ISBI, 2019

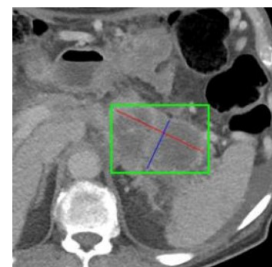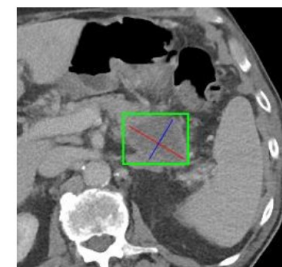# Common Thorax Disease Classification and Reporting in Chest X-rays

Chest X-ray: the largest public X-ray image dataset generated by NLP tools

- Over 100,000 frontal-view X-ray images

- 30k unique patients

- 14 common thorax diseases (e.g. pneumonia)

Wang et al., TieNet: Text-Image Embedding Network for Common Thorax Disease Classification and Reporting in Chest X-rays. CVPR 2018

# A summary of NLP methods

- Rule-based
  - For instance, if two genes co-occur in literature, it will be considered as interacted
  - Simple and efficient, but cannot tackle complex scenarios and not generalizable

- Machine learning
  - Manually derive features as inputs
  - Have been used over decades, but are limited to domain knowledge

- Recent methods: **deep learning**
  - Automatically derive features and representations
  - Have outperformed traditional machine learning methods since 2010
  - Have been widely applied in NLP applications: question answering, translation…
  - Open issues: privacy and interpretability

# Summary

- What is NLP?

- Why is NLP important?
    - NLP helps find relevant papers
    - NLP helps extract information (named entity, relation, …)
    - NLP helps image analysis

- NLP methods overview

# Text Mining Group @ NCBI/NLM



Zhiyong Lu
(Principal Investigator)

Alexis Allot

Qingyu Chen
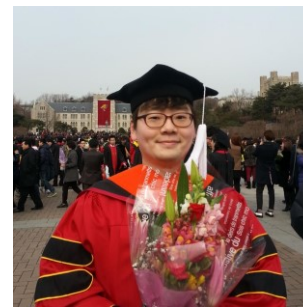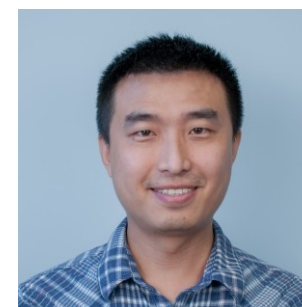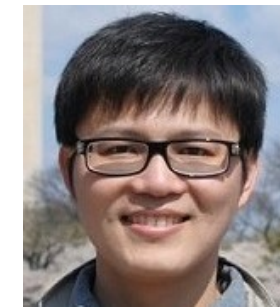
Donald Comeau

Rezarta Dogan

Alan Hsu

Sun Kim

Won Kim

Robert Leaman

Kyubum Lee

Yifan Peng

Chih-Hsuan Wei

Lana Yeganova

# Resources

- FAES course
  - BIOF395 "Introduction to Text Mining", Fall 2019 (instructed by us)

- NIH.AI workshop on NLP
  - **Who**:        Entry-level to advanced NIH researchers working with NLP
  - **When**:     April or May, 2019
  - **Mail list**: BIOINFORMATICS-SIG-L@LIST.NIH.GOV

- Reviews
  - https://www.ncbi.nlm.nih.gov/labs/pubmed/27807747
  - https://www.ncbi.nlm.nih.gov/labs/pubmed/22549152
  - https://www.ncbi.nlm.nih.gov/labs/pubmed/19649304

# Q & A

yifan.peng@nih.gov
qingyu.chen@nih.gov