# Knowledge Node and Relation Detection
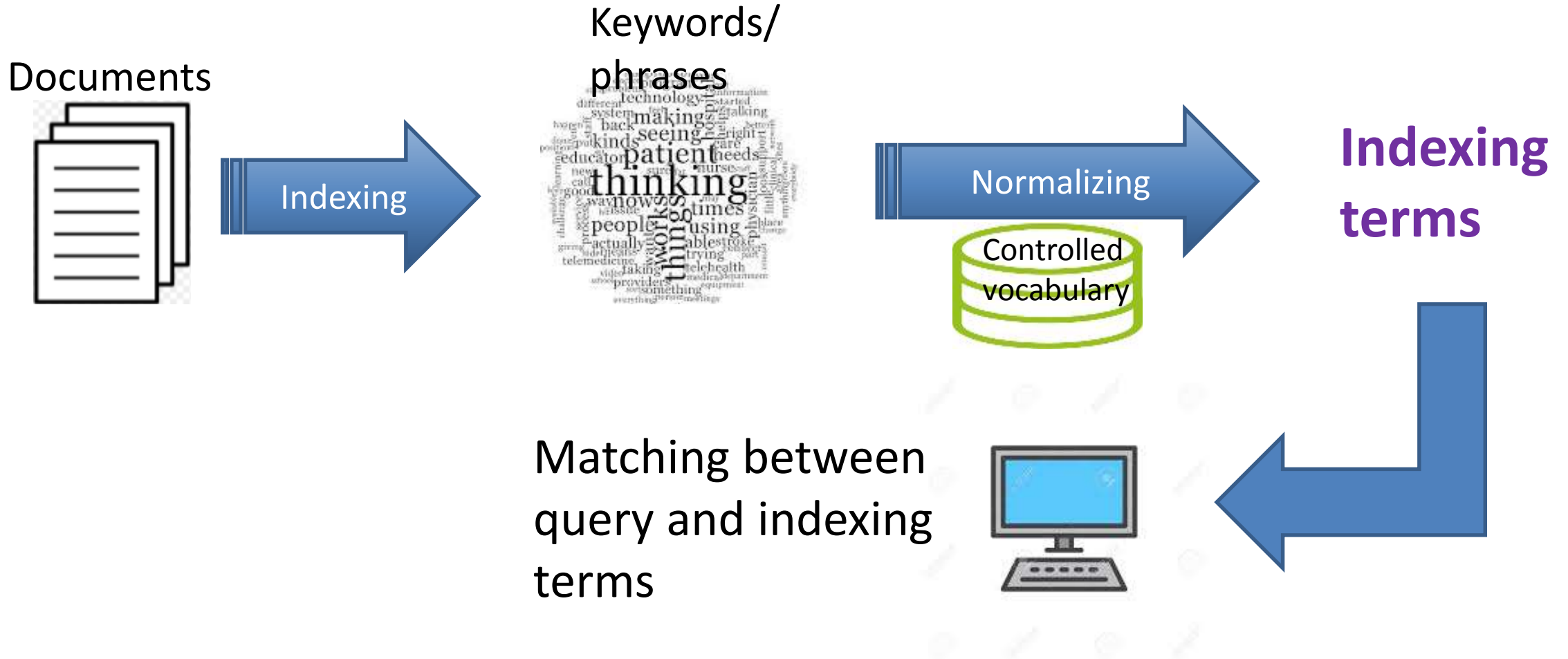
Jian Qin

School of Information Studies

Syracuse University

NKOS Workshop at DCMI/TPDL, Porto, Portugal, September 14, 2018

# It is all about subject content representation

Documents

Keywords/ phrases

Indexing →

Normalizing

Controlled vocabulary

**Indexing terms**

Matching between query and indexing terms

# Increase interactivity by transforming the way knowledge is represented



https://fusiontables.google.com/DataSource?docid=1Gs7wXxBI5TeiUrsV3MBoJNHouJEqjMk-ZSxmrsOC#chartnew:id=4

School of Information Studies
SYRACUSE UNIVERSITY

# Limitations

- Lack of rich relations between concepts (or other types of things) beyond scope relations
  - Between publications and data
  - Between datasets in different data repositories
  - Inside data and/or publications:
    - Between different types of entities
    - Between different topics
    - …
- Discrete terms from indexing process that must rely on relations defined in controlled vocabularies to show relations

# In linked data age...

**Knowledge Organization (KO)**

KO systems or structures

Codified in some formats and structures

**Knowledge Representation (KR)**

Knowledge nodes and relations

Codified in triples or structures that can facilitate computational processing and analysis

# How knowledge is represented?

- Currently two practices:

  - From natural language in full-text documents
    - Traditional indexing
    - Natural language processing and machine learning

  - From existing KOS through remodeling and restructuring
    - Converting existing KOS into linked data service (LCSH, MeSH, AAT)
    - Transform legacy data into linked data (e.g., library linked data)

School of Information Studies
SYRACUSE UNIVERSITY

# Paradigm shift from term-based representation to node-relation representation

N1H1 — causes → Influenza

N1H1 — is-a → Virus

Influenza — is-a → Disease

Virus — causes → Disease

Concept detection is relatively straightforward with help of KOS. Relation detection, however, has become the jewel in the crawn (or bottleneck problem) for representation knowledge in linked data age.

# Experiment

- 30 documents from PubMed
- Hand annotated 150 sentences to identify knowledge nodes (k-nodes) and relations in format [k-node(A), relation, k-node(B)]
- Indexing software used: MetMap and SemRep
  - Both support concept detection backed by UMLS
  - SemRep supports relation extraction
- Evaluation of results used Bilingual Evaluation Understudy (BLEU) and cosine similarity algorithm

School of Information Studies
SYRACUSE UNIVERSITY

# Research questions

- To what extent manually annotated and automatically generated k-nodes and relations are similar or dissimilar?

- What are some of the patterns of agreement and/or disagreement between the two sets of results?

- How can human intelligence (human-intervened k-node and relation recognition) be translated into machine intelligence for more accurate knowledge representation?

School of Information Studies
SYRACUSE UNIVERSITY

# Findings: degree of abstraction

| Sentence | Manually annotated k-nodes | MetaMap extracted k-nodes | SemRep extracted k-nodes |
|---|---|---|---|
| **Unlike most pathologic testing, which serves as an adjunct to establishing a diagnosis, the results of HER2 testing stand alone in determining which patients are likely to respond to trastuzumab, a monoclonal antibody against HER2.** | • pathologic testing<br>• HER2 testing<br>• monoclonal antibody<br>• trastuzumab<br>• HER2<br>• diagnosis | • pathologic testing<br>• results of her2 testing<br>• respond to trastuzumab<br>• results of her2 testing<br>• a monoclonal antibody against her2<br>• diagnosis | • pathologic<br>• testing<br>• HER2<br>• testing<br>• trastuzumab<br>• monoclonal antibody<br>• diagnosis |

# Findings: degree of abstraction

| Sentence | Manually annotated k-nodes | MetaMap extracted k-nodes | SemRep extracted k-nodes |
|---|---|---|---|
| At present, several preanalytic factors, including the time from tissue removal to tissue fixation, are underappreciated as important variables that have the potential to negatively impact the consistency and reliability of HER2 testing. | time from tissue removal to tissue fixation preanalytic factor HER2 testing preanalytic factor consistency reliability | time from tissue removal tissue fixation several preanalytic factors reliability of her2 testing several preanalytic factors consistency | time removal tissue fixation factors HER2 testing consistency |

# Findings: degree of abstraction

| Relations detected by SemRep | Relations from manual annotation | | |
|---|---|---|---|
| | Exact match | Similar/Partial match | No match |
| AFFECTS | affects | allows, improves, impacts, promotes provides, controls | is against is essential to documents enumerates confirms assesses assays begins with demonstrates establishes harbors has identifies includes is approved by is performed by predicts responds to |
| IS-A | is-a | is a kind of, exists, is equivalent to, is a prototype for, is given as | |
| ASSOCIATED_WITH | is associated with | is-for correlates | |
| AUGMENTS | expands | | |
| CAUSES | Causes, makes, determines | leads to, promotes, drives, improves | |
| COMPARED_WITH | | is measured by, is-tested-by, measures, is-in-context | |
| LOCATION_OF | | | |
| METHOD_OF | is-method-for | | |
| PART_OF | is-part-of | is a factor of, has-attribute, has condition of | |

# Findings: types of k-nodes and relations

*Simple k-node relations*: two simple k-nodes are connected by a direct relation in the form of a single verb: A→B

(amplification_of_HER2_gene,   promotes,   receptor_activation)
(tumor,                                        harbors,               HER2_molecular_alteration)

*Compound k-node relations*: refers to situation where one k-node is related to more than one k-node that has the same or different relations: A→(B₁…Bₙ)

(overexpression_of_receptor,   mediates,       biology_behavior_of_HER2-positive_tumor_cells)
(overexpression_of_receptor,   mediates,       clinical_behavior_HER2-positive_tumor_cells)
(overexpression_of_receptor,   drives,         proliferation_of_tumor_cells)
(overexpression_of_receptor,   drives,         survival_of_tumor_cells)

(overexpression_of_receptor,   mediates,       (biology_behavior_of_HER2-positive_tumor_cells, clinical_behavior_HER2-positive_tumor_cells))
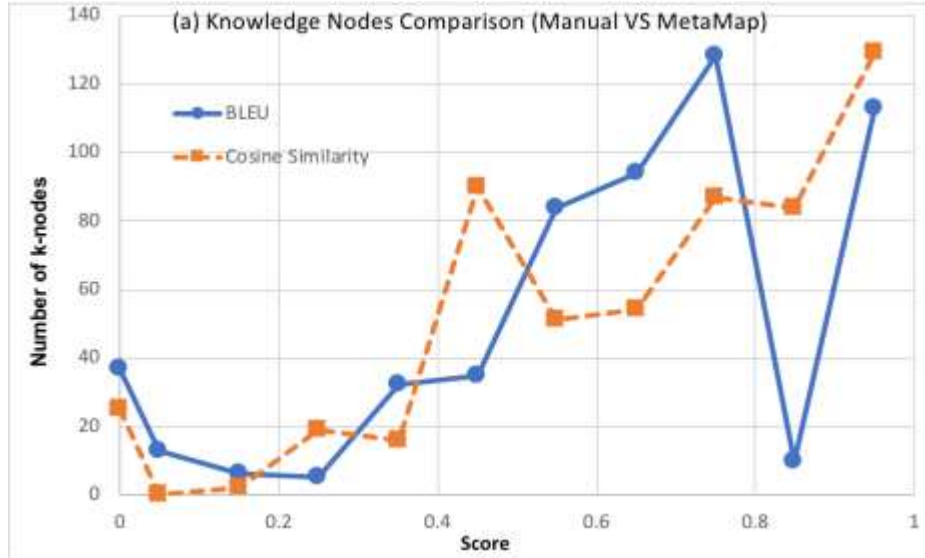(overexpression_of_receptor,   drives,         (proliferation_of_tumor_cells, survival_of_tumor_cells))

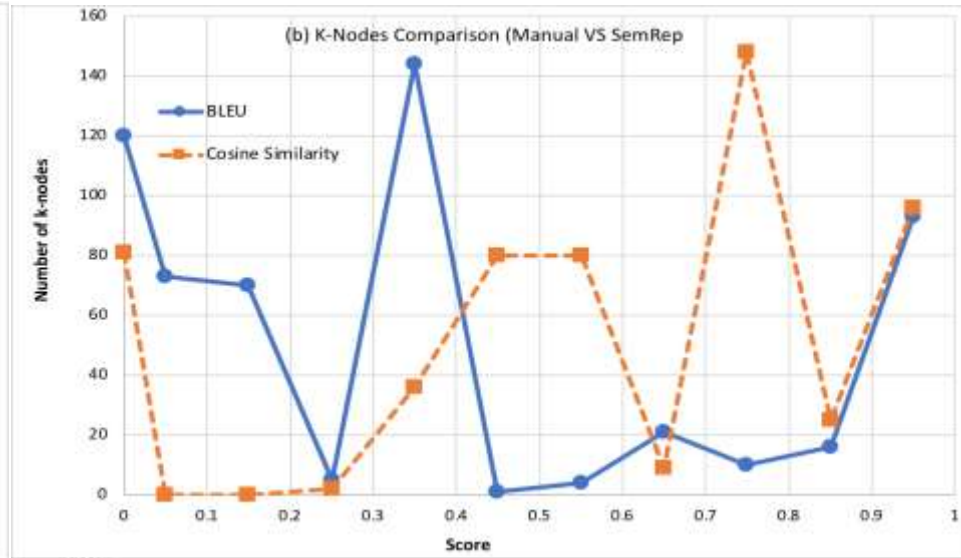# Findings: types of k-nodes and relations (cont'd)

*Complex k-node relations*: multiple k-nodes and the relations chained together by "bridge" k-nodes: A→(B→C)

(HER2_testing,   determines,   (patient,   responds-to,   trastuzumab))
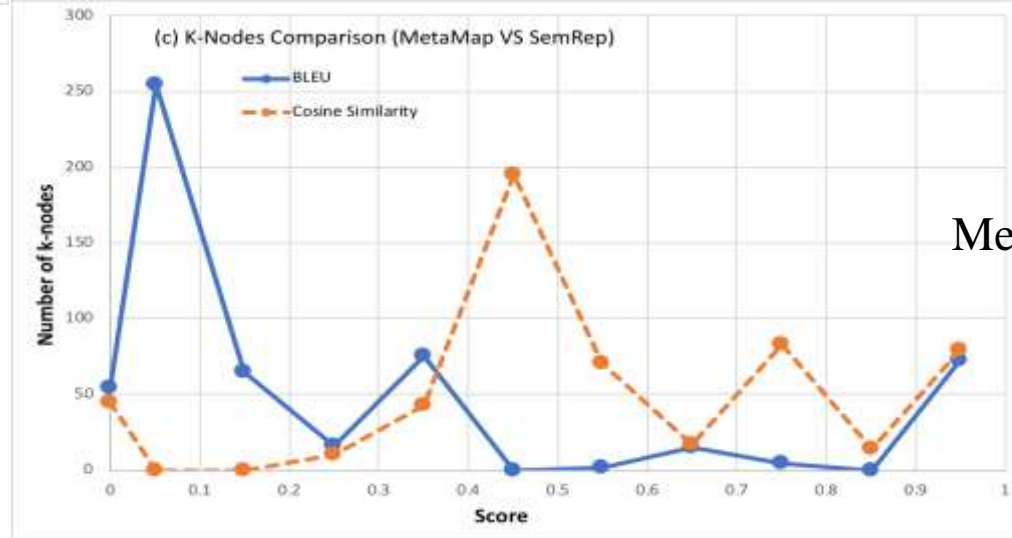
(trastuzumab,   is-a,   monoclonal antibody against HER2)

School of Information Studies
SYRACUSE UNIVERSITY

# Evaluation scores for three k-node detection methods



(a) Knowledge Nodes Comparison (Manual VS MetaMap)

Manual vs. MetaMap

(b) K-Nodes Comparison (Manual VS SemRep

Manual vs. SemRep

(c) K-Nodes Comparison (MetaMap VS SemRep)

MetaMap vs. SemRep

# Evaluation scores for UMLS term matching


(a) UMLS Terms Comparison (Manual VS MetaMap)

Manual vs. MetaMap


UMLS Terms Comparison (Manual VS SemRep)

Manual vs. SemRep


UMLS Terms Comparison (MetaMap VS SemRep)

MetaMap vs. SemRep

# Discussion and conclusion

- knowledge node and relation recognition from full-text documents is highly challenging, yet critically important in the big data era.

- Each of k-node and relation detection methods has different areas of strengths and limitations.

- Automatic tools have a long way to go

- K-node and relation representation facilitates knowledge network generation

School of Information Studies
SYRACUSE UNIVERSITY