

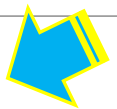


# **Towards Automatic Construction of Text-Rich Information Networks from Text**

**JIAWEI HAN  
COMPUTER SCIENCE  
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN  
AUGUST 18, 2022**

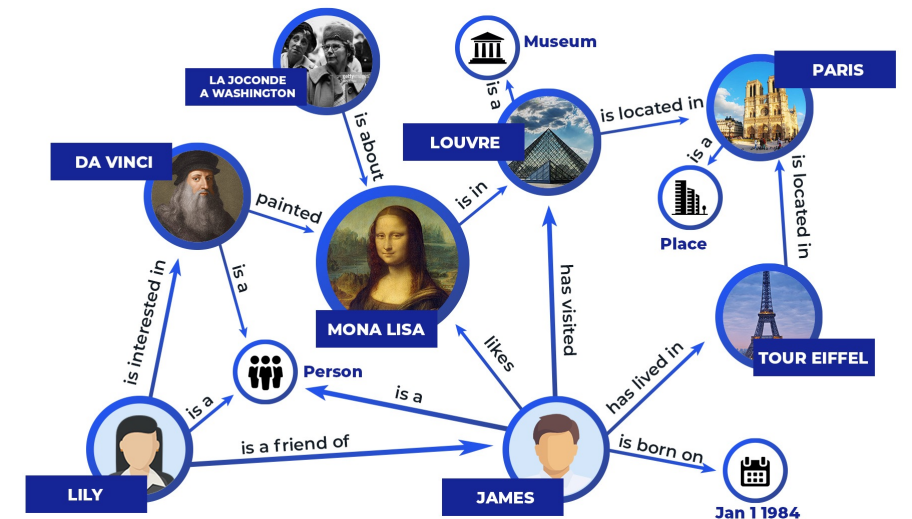
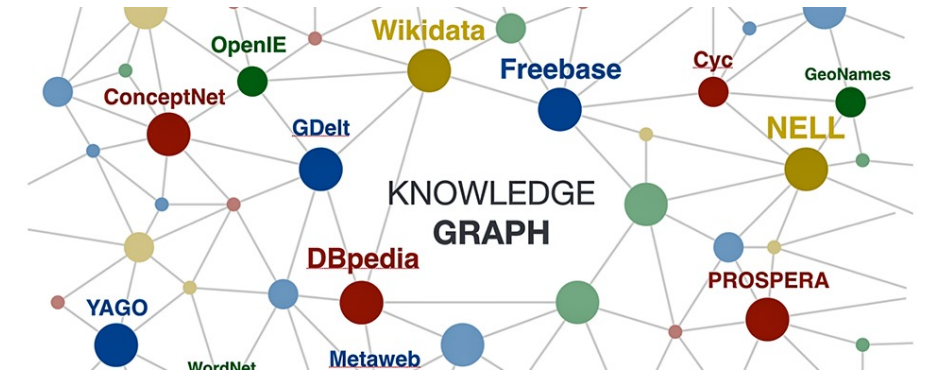
# Outline

---

- ❑ **What Kinds of Text-Rich Information Networks Do We Really Need?** 
- ❑ **Key Issue: Construction of Theme-/Corpus-Based Information Networks**
- ❑ **The Role of Embedding and PLM in Information Network Construction**
- ❑ **Data Preparation: Taxonomy-Guided Text Classification**
- ❑ **Identifying Information Network Primitives: Entities, Properties and Relations**
- ❑ **Conclusion: Towards Theme/Corpus-Based Information Network Construction**

# Info. Networks Are Used to Solved Real-World Problems

- ❑ Current information networks used in our research
  - ❑ Knowledge graphs: one gigantic graph for real-world?
  - ❑ Citation graphs
  - ❑ DBLP: authors, venue, keywords, citations & affiliations are very different types of links
  - ❑ Network repository contains ~40 different kinds of graphs (<https://networkrepository.com/network-data.php>)
  - ❑ Are we really using our network mining studies solving our real-world problems?
  - ❑ What are the burning problems we are solving — the real-life problem in scale?



Ack. Figures are from Google images


# Most Real-Life Info. Networks Need to be Constructed

---

- ❑ Most daily life data or the problems to be solved are essentially info. networks
  - ❑ News events: essentially information networks to be constructed
  - ❑ Tweet networks need to be understood from structured text analysis
  - ❑ University Web pages (departments, professors, courses, students, ...) are also information networks
  - ❑ Research literatures are also information networks
    - ❑ Types, entities, relations of many different types
    - ❑ Not just meta-data: authors, venues, keywords, citations, .....
- ❑ If we want our research or technology to be relevant
  - ❑ We have to solve real-world network problems
- ❑ If we want to solve real-world network problems
  - ❑ We have to study how construct real-world networks from unstructured data

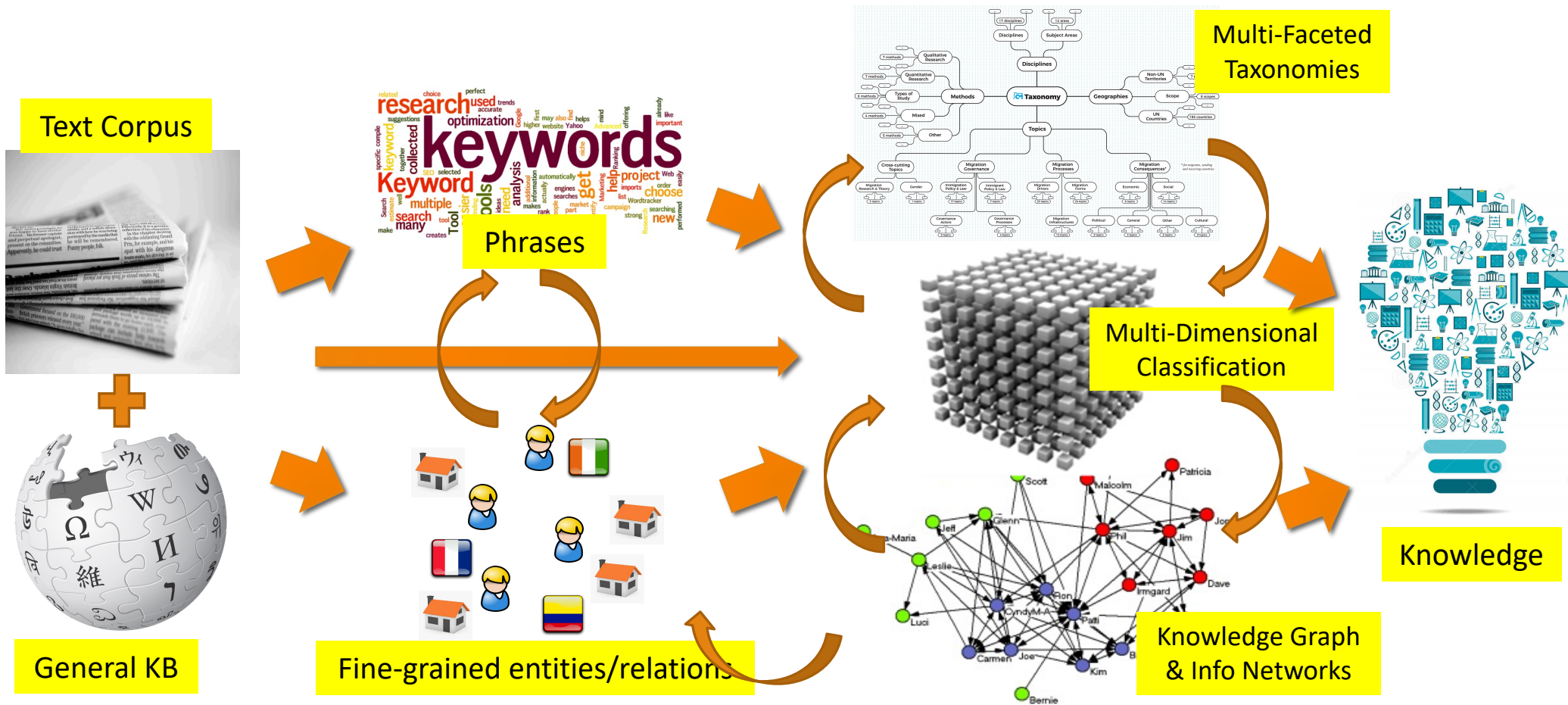
# Outline

---

- ❑ **What Kinds of Text-Rich Information Networks Do We Really Need?**
- ❑ **Key Issue: Construction of Theme-/Corpus-Based Information Networks** 
- ❑ **The Role of Embedding and PLM in Information Network Construction**
- ❑ **Data Preparation: Taxonomy-Guided Text Classification**
- ❑ **Identifying Information Network Primitives: Entities, Properties and Relations**
- ❑ **Conclusion: Towards Theme/Corpus-Based Information Network Construction**


# Automated, Local Information Network Construction

- Our General Roadmap: Mining structuring from unstructured text
  - One gigantic knowledge graph vs. many small structured, type networks
  - Automated construction vs. human annotated construction



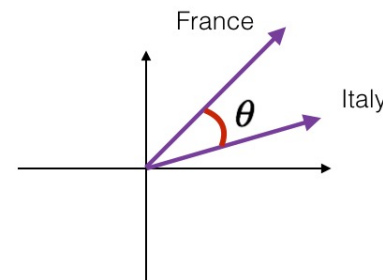
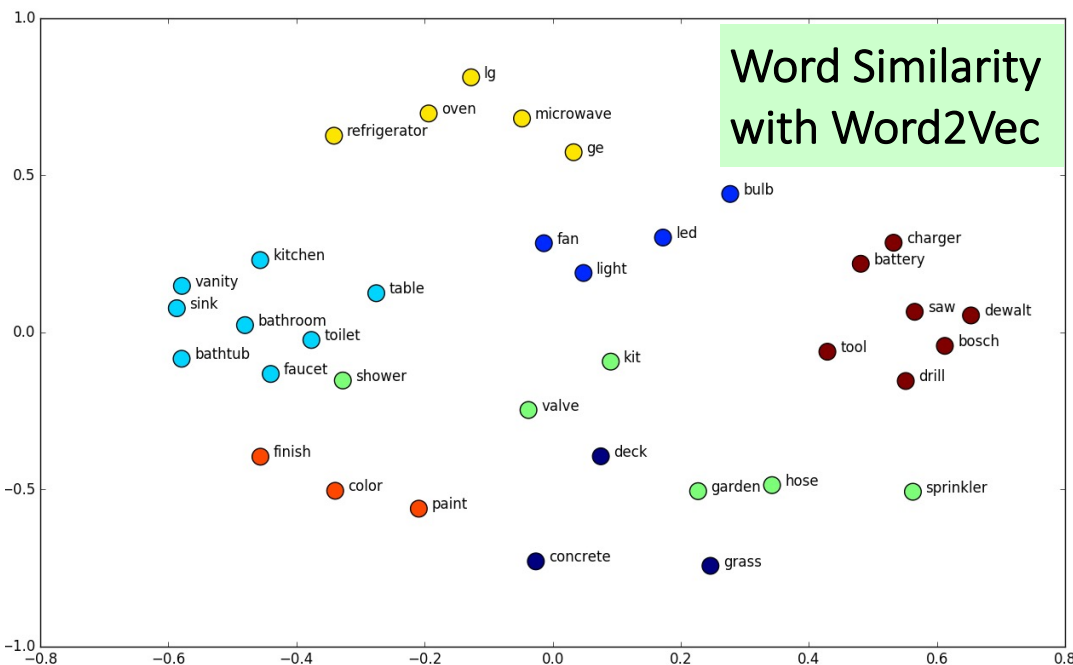
# Outline

---

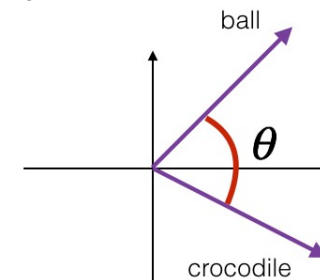
- ❑ **What Kinds of Text-Rich Information Networks Do We Really Need?**
- ❑ **Key Issue: Construction of Theme-/Corpus-Based Information Networks**
- ❑ **The Role of Embedding and PLM in Information Network Construction** 
- ❑ **Data Preparation: Taxonomy-Guided Text Classification**
- ❑ **Identifying Information Network Primitives: Entities, Properties and Relations**
- ❑ **Conclusion: Towards Theme/Corpus-Based Information Network Construction**

# Representation Learning in Text: Text Embedding

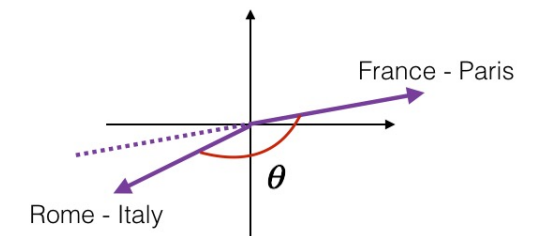
- Distributive representation: Embedding words in lower-dimension space
  - Word2Vec (Google), GloVe (Stanford), fastText (Facebook)
  - Handling sparsity & high dimensionality: Similar words are embedded closer
- Most text embeddings are trained in the Euclidean space but used on spherical space (i.e., cosine similarity)



France and Italy are quite similar  
 $\theta$  is close to  $0^\circ$   
 $\cos(\theta) \approx 1$

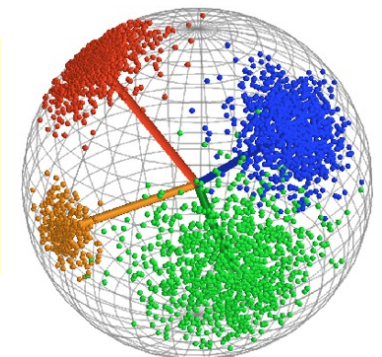


ball and crocodile are not similar  
 $\theta$  is close to  $90^\circ$   
 $\cos(\theta) \approx 0$



the two vectors are similar but opposite  
the first one encodes (city - country)  
while the second one encodes (country - city)  
 $\theta$  is close to  $180^\circ$   
 $\cos(\theta) \approx -1$

Spherical Text Embedding [NeurIPS'19]:  
embeddings are normalized, and  
spherical clustering algorithms are used





# Joint Embedding: Integrating Local and Global Contexts

- Local contexts can only partly define word semantics in unsupervised word embedding learning

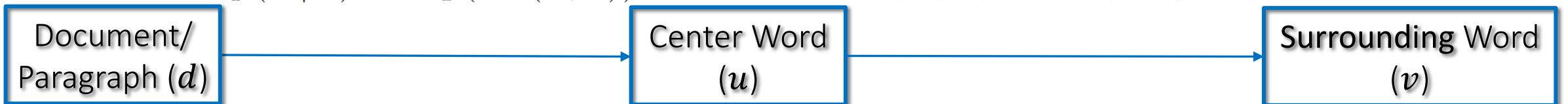
Local contexts of  
"harmful"

If I hear someone screwing with my car (ie, setting off the **alarm**) and **taunting** me to come out, you can be very sure that my Colt Delta Elite will also be coming with me. It is not the screwing with the car that would get them **shot**, it is the potential physical **danger**. If they are **taunting** like that, it's very possible that they also intend to **rob** me and or do other physically **harmful** things. Here in Houston last year a woman heard the sound of someone ...

- Design a generative model on the sphere that follows how humans write articles:
  - First a general idea of the paragraph/doc, then start to write down each word in consistent with not only the paragraph/doc, but also the surrounding words

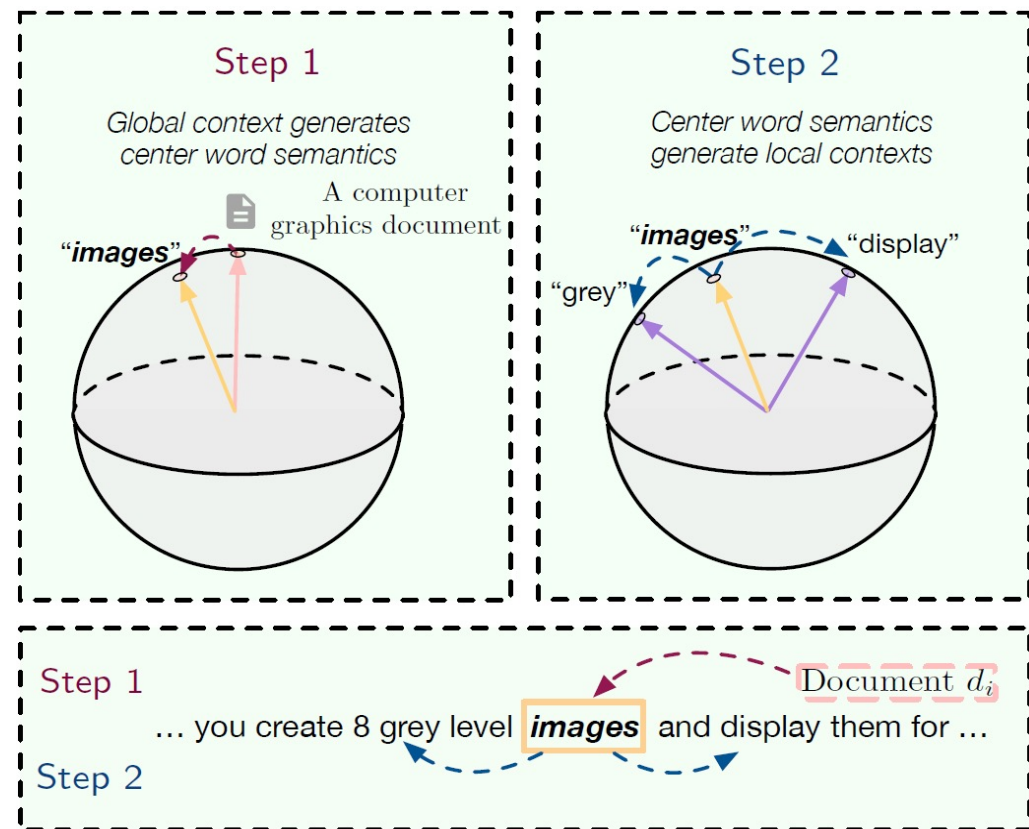
$$p(u | d) \propto \exp(\cos(\mathbf{u}, \mathbf{d}))$$

$$p(v | u) \propto \exp(\cos(\mathbf{v}, \mathbf{u}))$$



# Joint Spherical Embedding: Performance Comparison

## Understanding the Spherical Generative Model



Global Context Helps Interpreting Acronyms

## Word Similarity: Performance Comparison

Table 1: Spearman rank correlation on word similarity evaluation.

Embedding Space	Model	WordSim353	MEN	SimLex999
Euclidean	Word2Vec	0.711	0.726	0.311
	GloVe	0.598	0.690	0.321
	fastText	0.697	0.722	0.303
	BERT	0.477	0.594	0.287
Poincaré	Poincaré GloVe	0.623	0.652	0.321
Spherical	<b>JoSE</b>	<b>0.739</b>	<b>0.748</b>	<b>0.339</b>

Table 5: Effect of Global Context on Interpreting Acronyms.

Acronyms	Global ( $\lambda = \infty$ )	Local ( $\lambda = 0$ )
CMU	<b>mellon, carnegie,</b> andrew, pa, pittsburgh	andrew, kfnjyea00uh, am2x, mr47, devineni
UIUC	<b>urbana, illinois, uxa,</b> <b>univ, uchicago</b>	uxa, ux4, ux1, mrcnext, cka52397
UNC	<b>chapel, carolina, astro,</b> images, usc	launchpad, gibbs, umr, lambada, jge
Caltech	<b>california, gap, institute,</b> keith, <b>technology</b>	juliet, jafoust, lmh, henling, bdunn
JHU	<b>johns, camp, hopkins,</b> nation, grand	pablo, hasch, iglesias, davidk, atlantis

# Discriminative Topic Mining via Category Name-Guided Embedding

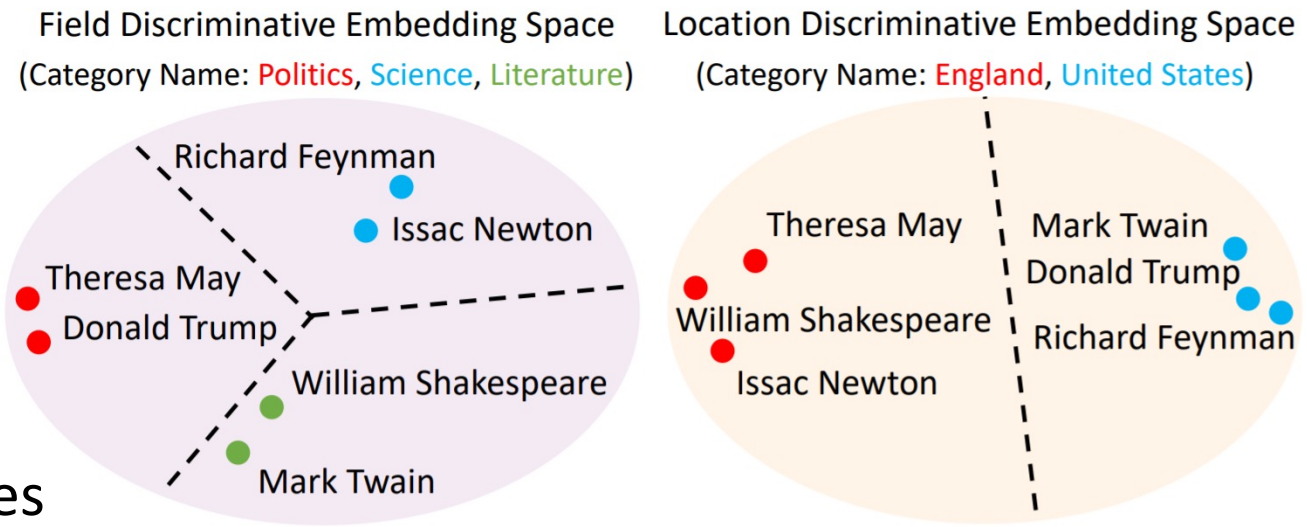
- Traditional text embedding (e.g., Word2Vec, GloVe, fastText, JoSE)
  - Mapping words with similar local contexts closer in the embedding space
  - Not imposing particular assumptions on the type of data distributions
- CatE: Category Name-guided Embedding [WWW'20]**
  - Weak guidance: leverages *category names* to learn word embeddings with discriminative power over the specific set of categories

### CatE: Inputs

- Category names + Corpus

### CatE: Outputs

- The same set of celebrities are embedded differently given different sets of category names



# Method of CatE: Category-name guided text Embedding

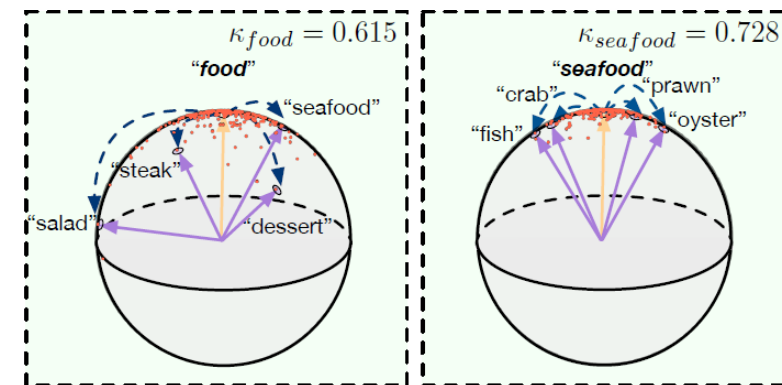
- A category-name guided text embedding learning module (E):
  - Takes a set of category names to learn category distinctive word embeddings by modeling the text generative process conditioned on the user provided categories

$$P(\mathcal{D} | \mathcal{C}) = \prod_{d \in \mathcal{D}} p(d | c_d) \prod_{w_i \in d} p(w_i | d) \prod_{\substack{w_{i+j} \in d \\ -h \leq j \leq h, j \neq 0}} p(w_{i+j} | w_i)$$

- A category representative words retrieval module (R):
  - Selects category representative words based on both word embedding similarity and word distributional specificity

The two modules (E + R) collaborate in an iterative way:

- E refines word embeddings and category embeddings
- R selects representative words that will be used by E in the next iteration



# Performance Study on Discriminative Topic Mining

## Quantitative comparison

TC: topic coherence

MACC: Mean accuracy

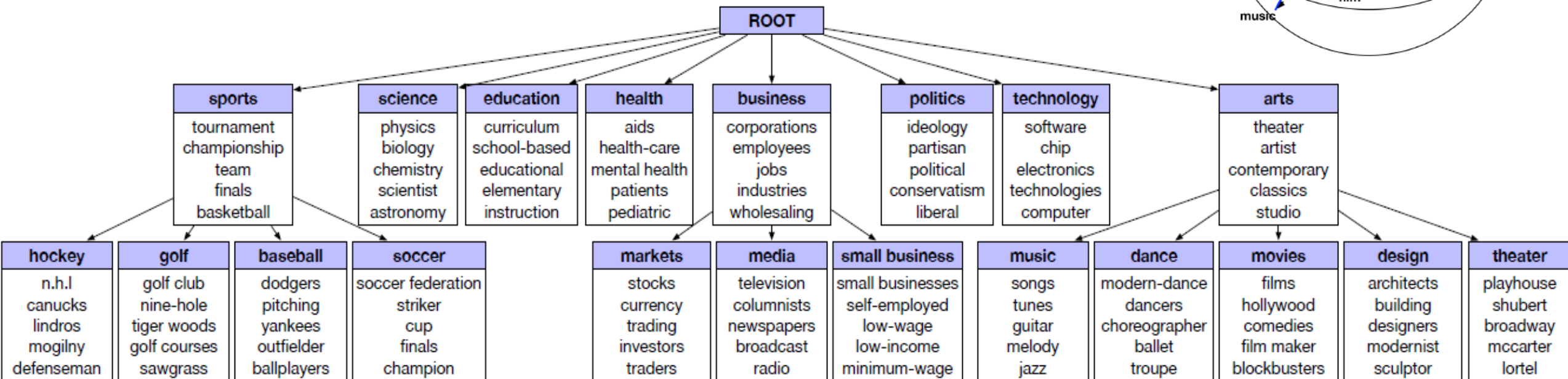
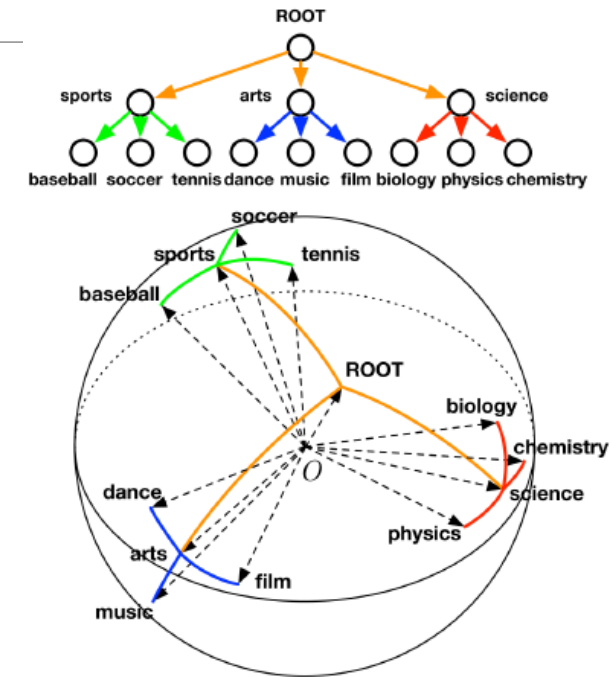
## Qualitative Comparison of Discriminative Topic Mining

Methods	NYT-Location		NYT-Topic		Yelp-Food		Yelp-Sentiment	
	TC	MACC	TC	MACC	TC	MACC	TC	MACC
LDA	0.007	0.489	0.027	0.744	-0.033	0.213	-0.197	0.350
Seeded LDA	0.024	0.168	0.031	0.456	0.016	0.188	0.049	0.223
TWE	0.002	0.171	-0.011	0.289	0.004	0.688	-0.077	0.748
Anchored CorEx	0.029	0.190	0.035	0.533	0.025	0.313	0.067	0.250
Labeled ETM	0.032	0.493	0.025	0.889	0.012	0.775	0.026	0.852
CatE	<b>0.049</b>	<b>0.972</b>	<b>0.048</b>	<b>0.967</b>	<b>0.034</b>	<b>0.913</b>	<b>0.086</b>	<b>1.000</b>

Methods	NYT-Location		NYT-Topic		Yelp-Food		Yelp-Sentiment	
	britain	canada	education	politics	burger	desserts	good	bad
LDA	company (×)	percent (×)	school	campaign	fatburger	ice cream	great	valet (×)
	companies (×)	economy (×)	students	clinton	dos (×)	chocolate	place (×)	peter (×)
	british	canadian	city (×)	mayor	liar (×)	gelato	love	aid (×)
	shares (×)	united states (×)	state (×)	election	cheeseburgers	tea (×)	friendly	relief (×)
	great britain	trade (×)	schools	political	bearing (×)	sweet	breakfast	rowdy
CatE	england	ontario	educational	political	burgers	dessert	delicious	sickening
	london	toronto	schools	international politics	cheeseburger	pastries	mindful	nasty
	britons	quebec	higher education	liberalism	hamburger	cheesecakes	excellent	dreadful
	scottish	montreal	secondary education	political philosophy	burger king	scones	wonderful	freaks
	great britain	ottawa	teachers	geopolitics	smash burger	ice cream	faithful	cheapskates

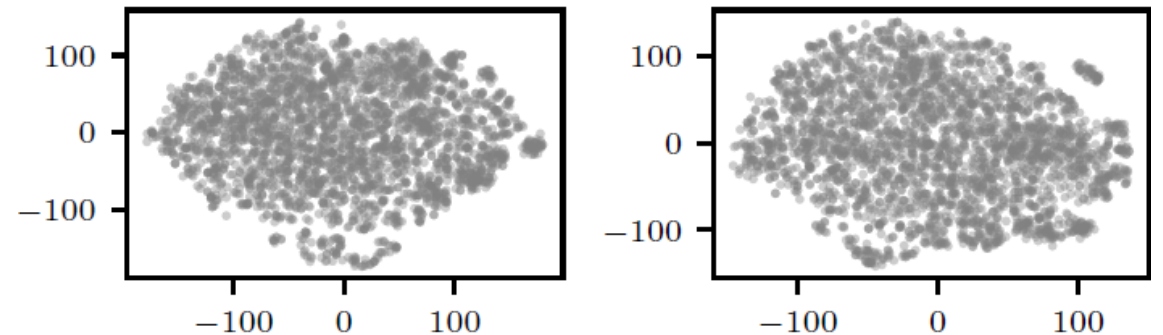
# Hierarchical Topic Mining via Joint Spherical Tree and Text Embedding [KDD'20]

- JoSH: A joint tree and text embedding method
- Simultaneous modeling of the category tree structure in the spherical space
- Effective mining of category representative, hierarchical terms
  - Ex. In PubMed literature, finding distinct terms related to *hormones, enzymes, vitamins, and vaccines*



# Topic Discovery via Latent Space Clustering of LM Embedding

- ❑ Yu Meng, Yunyi Zhang, Jiaxin Huang, Yu Zhang and Jiawei Han, “[Topic Discovery via Latent Space Clustering of Language Model Embeddings](#)”, in WWW’22
- ❑ Task: Automatic discovery of coherent and meaningful topics from text corpora
- ❑ Limitations of topic modeling (a generative process)
  - ❑ Ignoring word ordering information in text (based on the “bag-of-words” assumption)
  - ❑ cannot leverage external knowledge to learn word semantics, and
  - ❑ Inducing an intractable posterior that requires approximation algorithms
- ❑ Why not directly deploy pre-trained language models (PLMs) for topic discovery?
  - ❑ The PLM embedding space is partitioned into extremely fine-grained clusters and lacks topic structures inherently
  - ❑ PLM embeddings are high-dimensional while distance functions can become meaningless
  - ❑ Lack of good document representations from PLMs



(a) New York Times.

(b) Yelp Review.

Visualization of 3, 000 randomly sampled contextualized word embeddings of BERT: The embedding spaces do not have clearly separated clusters.

# Qualitative Evaluation of Topic Discovery


Corpus	# documents	# words/doc.	Vocabulary
NYT	31,997	690	25,903
Yelp	29,280	114	11,419

Methods	NYT					Yelp				
	Topic 1 (sports)	Topic 2 (politics)	Topic 3 (research)	Topic 4 (france)	Topic 5 (japan)	Topic 1 (positive)	Topic 2 (negative)	Topic 3 (vegetables)	Topic 4 (fruits)	Topic 5 (seafood)
LDA	olympic <u>year</u> <u>said</u> games team	<u>mr</u> bush president white house	<u>said</u> report evidence findings defense	french <u>union</u> <u>germany</u> <u>workers</u> paris	japanese tokyo <u>year</u> matsui <u>said</u>	amazing <u>really</u> <u>place</u> phenomenal pleasant	loud awful <u>sunday</u> <u>like</u> slow	spinach carrots greens salad <u>dressing</u>	mango strawberry <u>vanilla</u> banana <u>peanut</u>	fish <u>roll</u> salmon <u>fresh</u> <u>good</u>
CorEx	baseball championship playing <u>fans</u> league	house white support <u>groups</u> <u>member</u>	possibility challenge reasons <u>give</u> planned	french <u>italy</u> paris francs jacques	japanese tokyo <u>index</u> osaka <u>electronics</u>	great friendly <u>atmosphere</u> love favorite	<u>even</u> bad mean cold <u>literally</u>	garlic tomato onions <u>toppings</u> <u>slices</u>	strawberry <u>caramel</u> <u>sugar</u> fruit mango	shrimp <u>beef</u> crab <u>dishes</u> <u>salt</u>
ETM	olympic league <u>national</u> basketball athletes	government national <u>plan</u> public support	approach problems experts <u>move</u> <u>give</u>	french <u>students</u> paris <u>german</u> <u>american</u>	japanese <u>agreement</u> tokyo <u>market</u> <u>european</u>	nice worth <u>lunch</u> recommend friendly	disappointed cold <u>review</u> <u>experience</u> bad	avocado <u>greek</u> salads spinach tomatoes	strawberry mango <u>sweet</u> <u>soft</u> <u>flavors</u>	fish shrimp lobster crab <u>chips</u>
BERTopic	swimming freestyle <u>popov</u> gold olympic	bush democrats white bushs house	researchers scientists cases <u>genetic</u> study	french paris lyon <u>minister</u> <u>billion</u>	japanese tokyo ufj <u>company</u> yen	awesome <u>atmosphere</u> friendly <u>night</u> good	horrible <u>quality</u> disgusting disappointing <u>place</u>	tomatoes avocado <u>soups</u> kale cauliflower	strawberry mango <u>cup</u> lemon banana	lobster crab shrimp oysters <u>amazing</u>
TopClus	athletes medalist olympics tournaments quarterfinal	government ministry bureaucracy politicians electoral	hypothesis methodology possibility criteria assumptions	french seine toulouse marseille paris	japanese tokyo osaka hokkaido yokohama	good best friendly cozy casual	tough bad painful frustrating brutal	potatoes onions tomatoes cabbage mushrooms	strawberry lemon apples grape peach	fish octopus shrimp lobster crab



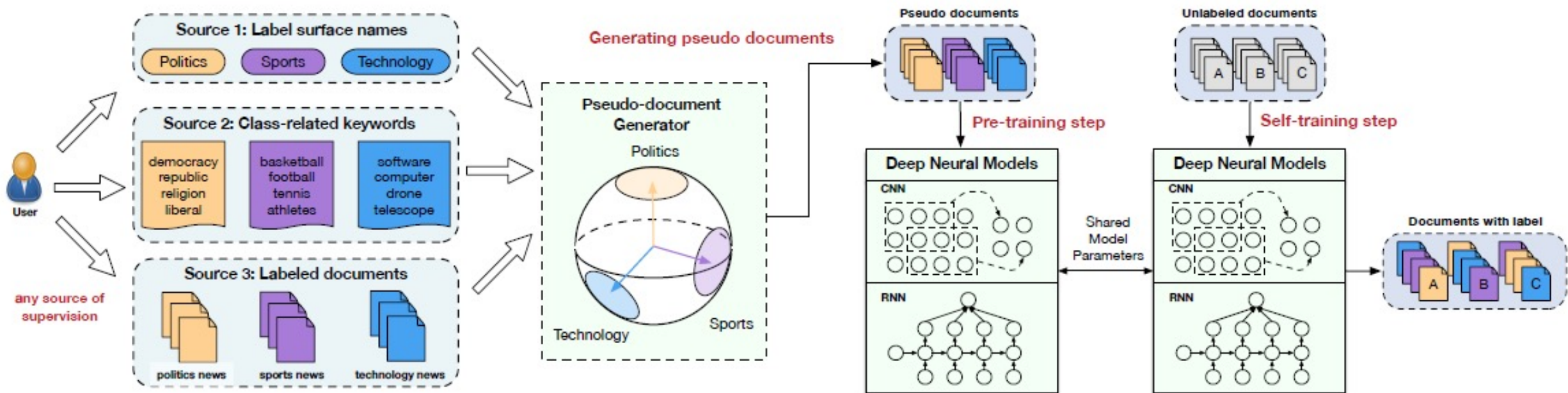
# Outline

---

- ❑ **What Kinds of Text-Rich Information Networks Do We Really Need?**
- ❑ **Key Issue: Construction of Theme-/Corpus-Based Information Networks**
- ❑ **The Role of Embedding and PLM in Information Network Construction**
- ❑ **Data Preparation: Taxonomy-Guided Text Classification** 
- ❑ **Identifying Information Network Primitives: Entities, Properties and Relations**
- ❑ **Conclusion: Towards Theme/Corpus-Based Information Network Construction**

# WeSTClass: Weakly Supervised Text Classification

- Modeling class distribution in word2vec embedding space
- Word2vec embedding captures **skip-gram (local) similarity** (i.e., words with similar local context windows are expected to have similar meanings)



WeSTClass (Weakly Supervised Text Classification): CIKM'18

WeSHClass (Weakly Supervised Hierarchical Text Classification): AAAI'19

# LOTClass: Label-Name-Only Text Classification [EMNLP'20]

- Yu Meng, et al., “Text Classification Using Label Names Only: A Language Model Self-Training Approach” [EMNLP'20]
- Inputs: A set of label names representing each class + unlabeled documents
- Method (3 steps): Make good use of pre-trained language model (e.g., BERT)
  - Step 1. Category understanding via label name replacement (learn *topic vocabulary*)
    - Ex. “sports” → {“soccer”, “basketball”, ...} (use pretrained LM to replace category name)

- Learn topic vocabulary using label name only
- Make good use of pretrained LM (e.g., BERT)
- Result from AGNews dataset

Label Name	Category Vocabulary
politics	politics, political, politicians, government, elections, politician, democracy, democratic, governing, party, leadership, state, election, politically, affairs, issues, governments, voters, debate, cabinet, congress, democrat, president, religion, ...
sports	sports, games, sporting, game, athletics, national, athletic, espn, soccer, basketball, stadium, arts, racing, baseball, tv, hockey, pro, press, team, red, home, bay, kings, city, legends, winning, miracle, olympic, ball, giants, players, champions, boxing, ...
business	business, trade, commercial, enterprise, shop, money, market, commerce, corporate, global, future, sales, general, international, group, retail, management, companies, operations, operation, store, corporation, venture, economic, division, firm, ...
technology	technology, tech, software, technological, device, equipment, hardware, devices, infrastructure, system, knowledge, technique, digital, technical, concept, systems, gear, techniques, functionality, process, material, facility, feature, method, ...

# LOTClass: Label-Name-Only Text Classification

- Step 2: Masked topic prediction: Create contextualized word-level supervisions to train the model for predicting a word's implied topic

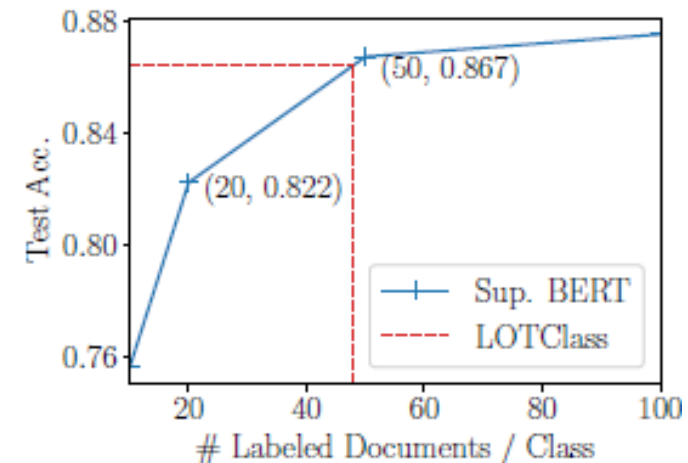
Different contexts leads to different BERT language model prediction



Sentence	Language Model Prediction
The oldest annual US team <b>sports</b> competition that includes professionals is not in baseball, or football or basketball or hockey. It's in soccer.	sports, baseball, handball, soccer, basketball, football, tennis, sport, championship, hockey, ...
Samsung's new SPH-V5400 mobile phone <b>sports</b> a built-in 1-inch, 1.5-gigabyte hard disk that can store about 15 times more data than conventional handsets, Samsung said.	has, with, features, uses, includes, had, is, contains, featured, have, incorporates, requires, offers, ...

- Step 3: Self-training: Generalize the model via self-training on abundant unlabeled data to make document-level topic prediction

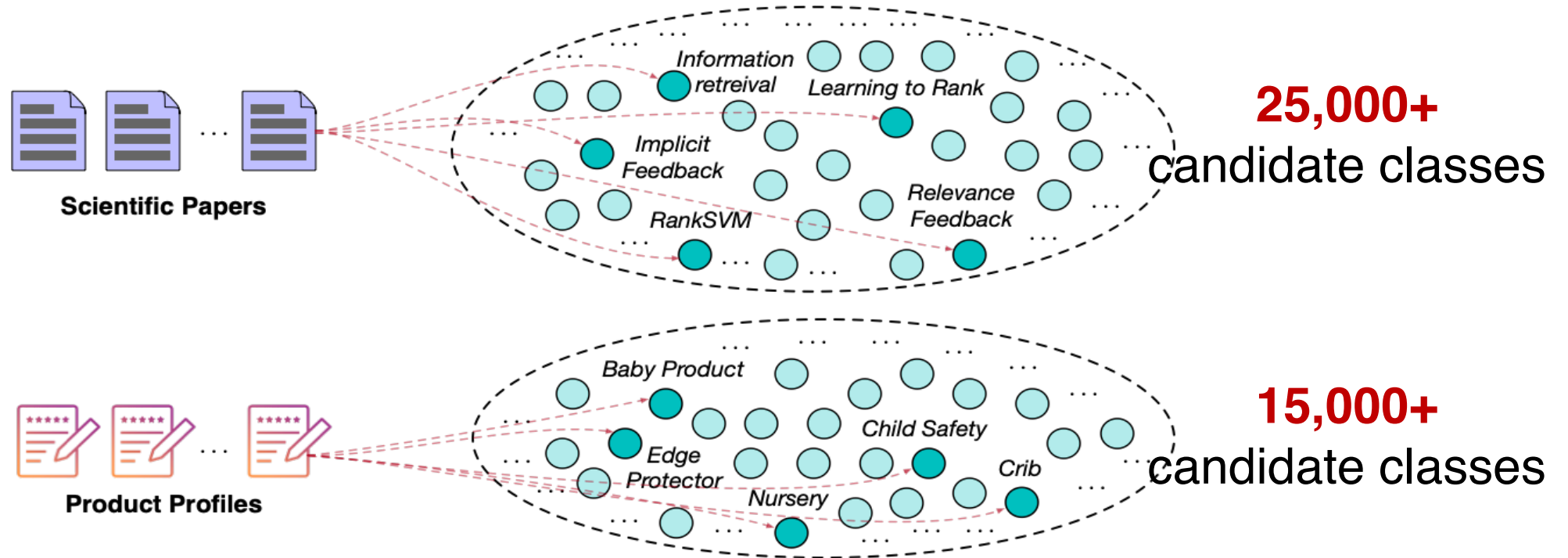
Supervision Type	Methods	AG News	DBpedia	IMDB	Amazon
Weakly-Sup.	Dataless (Chang et al., 2008)	0.696	0.634	0.505	0.501
	WeSTClass (Meng et al., 2018)	0.823	0.811	0.774	0.753
	BERT w. simple match	0.752	0.722	0.677	0.654
	LOTClass w/o. self train	0.822	0.860	0.802	0.853
	LOTClass	<b>0.864</b>	<b>0.911</b>	<b>0.865</b>	<b>0.916</b>
Semi-Sup.	UDA (Xie et al., 2019)	0.869	0.986	0.887	0.960
Supervised	char-CNN (Zhang et al., 2015)	0.872	0.983	0.853	0.945
	BERT (Devlin et al., 2019)	0.944	0.993	0.945	0.972



Label-name only is equiv. to 48 labels in Supervised BERT

# Need: “Structuring”/Tagging Unstructured Documents

- Task: Tag each doc. with a set of relevant classes from a huge candidate pool

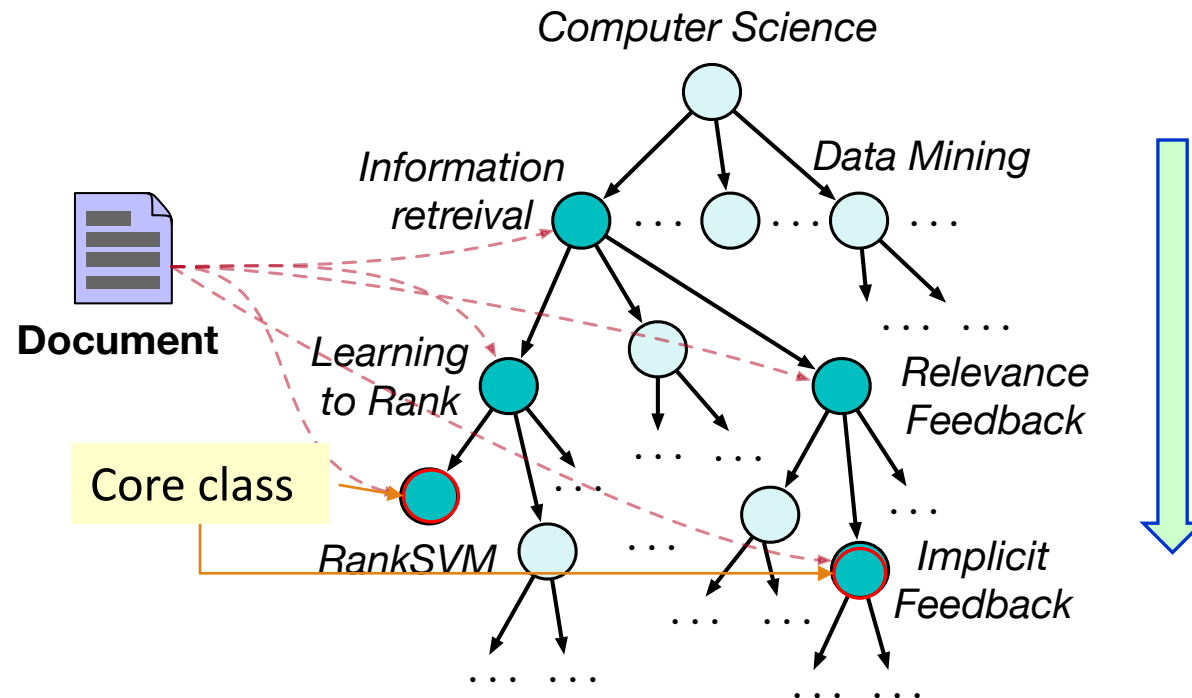


- Challenges:

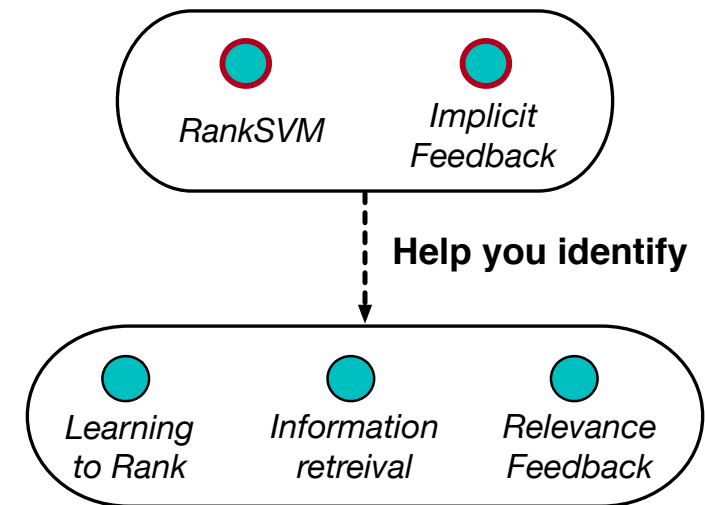
- Huge label space, multi-label tagging
- Limited labeled data— hard for supervised models

# TaxoClass [NAACL'21]: Taxonomy Comes to Rescue

- J. Shen, et al. “TaxoClass: Hierarchical Multi-Label Text Classification Using Only Class Names”, NAACL'21
- Taxonomy!— Structure the huge label space by organizing classes hierarchically
  - Enable fast label space exploration in a top-down way
- Facilitate multi-label tagging by capturing class relations

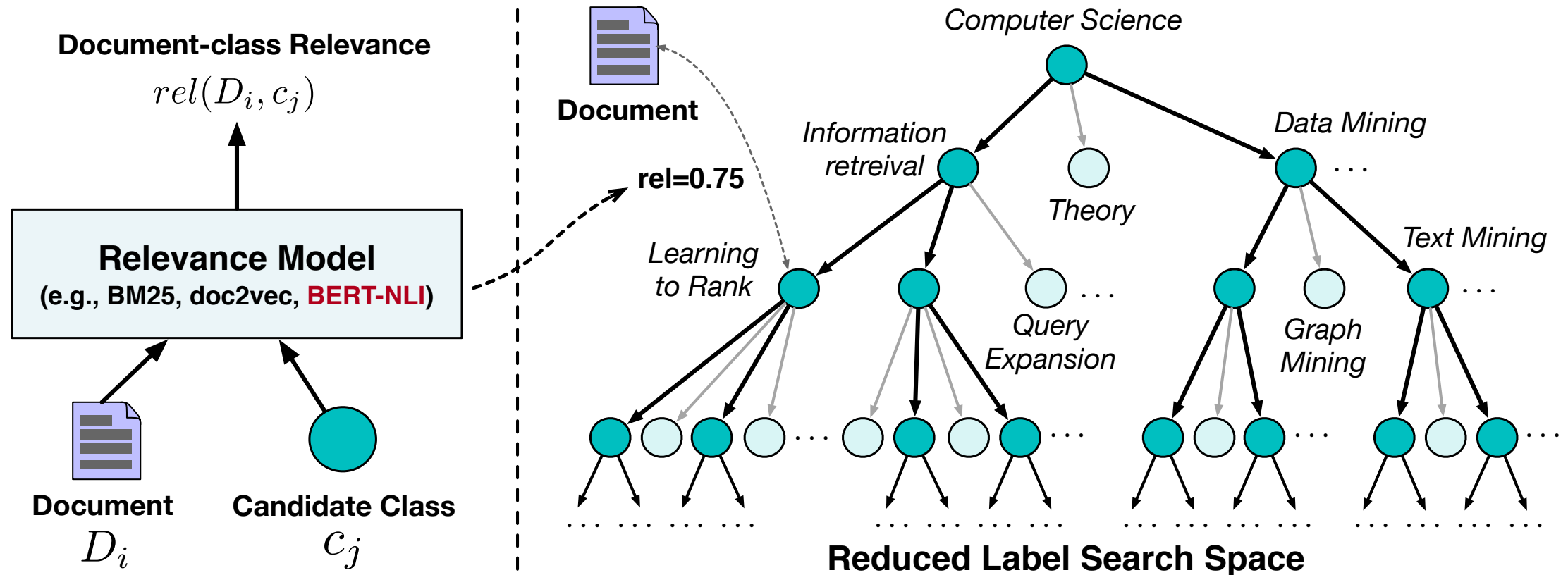


A few most essential **core classes**



# TaxoClass: A Weakly-Supervised Classification Method based on Taxonomy

- Shrink the label search space with top-down exploration
- Use a **relevance model** to filter out completely irrelevant classes for each document



# TaxoClass: Case Studies

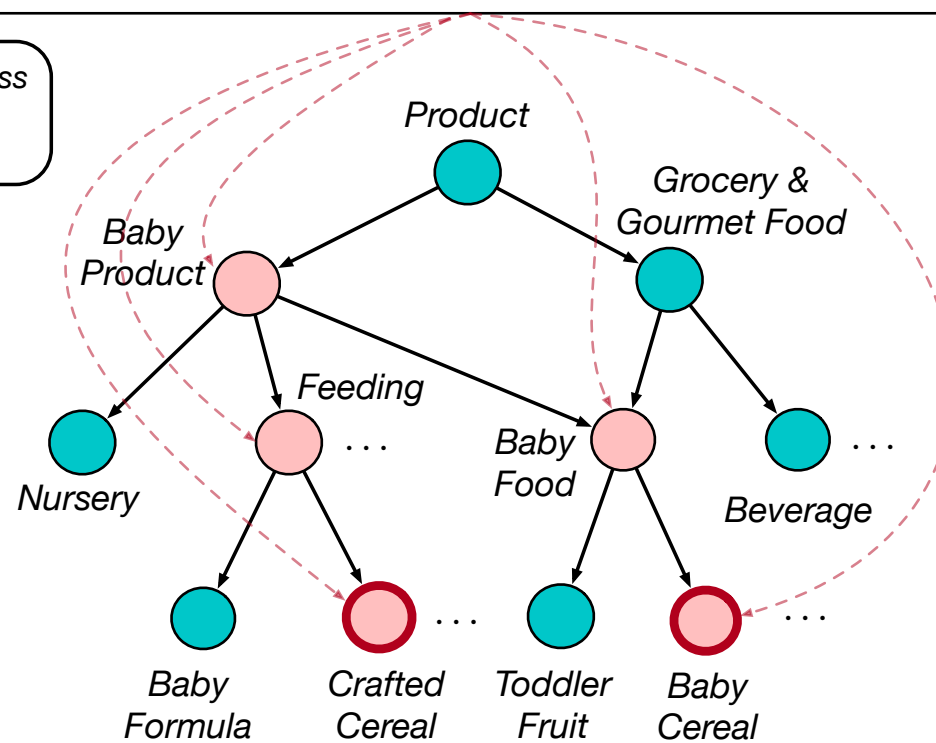
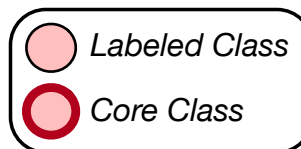
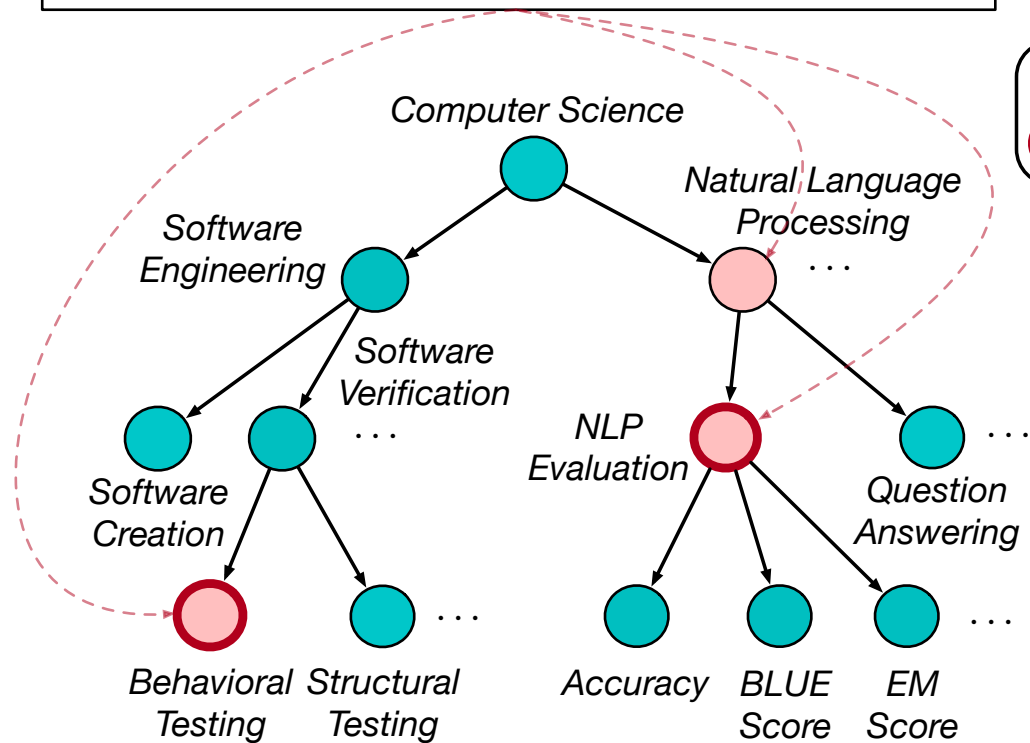


## Document

Inspired by principles of **behavioral testing** in software engineering, we introduce CheckList, a task-agnostic methodology for **testing NLP models**...

## Document

When our **son** was about **4 months old**, doctor said we could give him **crafted cereal** so we bought it. It digests well and doesn't lock up his bowels at all ...





# TaxoClass: Performance Comparison

Methods	Amazon		DBPedia		
	Example-F1	P@1	Example-F1	P@1	
Weakly-supervised multi-class classification method → WeSHClass (Meng et al., AAAI'19)	0.246	0.577	0.305	0.536	
Semi-supervised methods using 30% of training set {	SS-PCEM (Xiao et al., WebConf'19)	0.292	0.537	0.385	0.742
	Semi-BERT (Devlin et al., NAACL'19)	0.339	0.592	0.428	0.761
Zero-shot method ← Hier-0Shot-TC (Yin et al., EMNLP'19)	0.474	0.714	0.677	0.787	
TaxoClass (NAACL'21)	<b>0.593</b>	<b>0.812</b>	<b>0.816</b>	<b>0.894</b>	

$$\text{Example-F1} = \frac{1}{N} \sum_{i=1}^N \frac{2|true_i \cap pred_i|}{|true_i| + |pred_i|}, \text{P@1} = \frac{\#docs \text{ with top-1 pred correct}}{\#total docs}$$

- vs. WeSHClass: better model document-class relevance
- vs. SS-PCEM, Semi-BERT: better leverage supervision signals from taxonomy
- vs. Hier-0Shot-TC: better capture domain-specific information from core classes

Amazon: 49K product reviews (29.5K training + 19.7K testing), 531 classes  
DBPedia: 245K Wiki articles (196K training + 49K testing), 298 classes

# Outline

---

- ❑ **What Kinds of Text-Rich Information Networks Do We Really Need?**
- ❑ **Key Issue: Construction of Theme-/Corpus-Based Information Networks**
- ❑ **The Role of Embedding and PLM in Information Network Construction**
- ❑ **Data Preparation: Taxonomy-Guided Text Classification**
- ❑ **Identifying Information Network Primitives: Entities, Properties and Relations**
- ❑ **Conclusion: Towards Theme/Corpus-Based Information Network Construction**



# The ChemNER Framework

Input Corpus

Entity Span Detection

**S1:** [Methyl-14C]S-dThd was synthesized by rapid methylation of ...  
**S2:** ... Suzuki-Miyaura cross-coupling reactions were carried out ...  
**S3:** Although it was necessary to employ a stoichiometric quantity of palladium, it is noteworthy that the cross-coupling proceeded in the presence of a wide array of functional groups.  
**S4:** ... can undergo a transmetalation with either BBA or the rapidly forming boronic acid ...

Flexible KB-Matching

Knowledge Bases

**S1:** [Methyl-14C]S-dThd was synthesized by rapid methylation of ...

ORGANIC COMPOUNDS, ORGANIC POLYMERS      ORGANIC REACTIONS

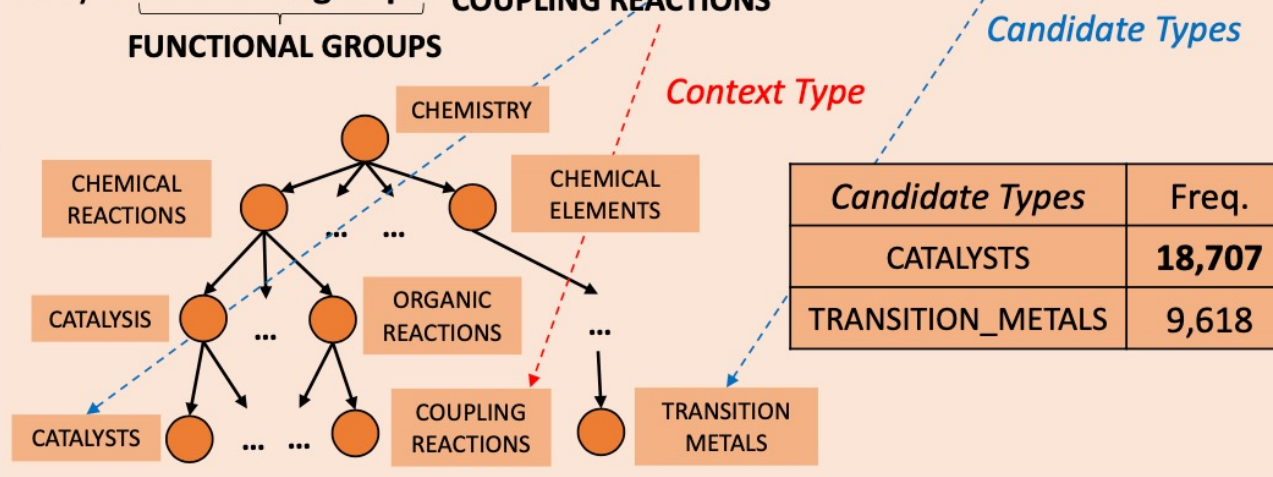
TF-IDF Scores	ORGANIC COMPOUNDS	ORGANIC POLYMERS	Biomolecules	...
methyl	0.0177	0.0139	0.0010	...
thd	0.0256	0.0115	0.0417	

**S2:** ..., Suzuki-Miyaura cross-coupling reactions were carried out ...

COUPLING REACTIONS

Ontology-guided Multi-type Disambiguation

**S3:** Although it was necessary to employ a stoichiometric quantity of palladium, it is noteworthy that the cross-coupling proceeded in the presence of a wide array of functional groups.



Sequence Labeling Models

BiLSMT-CRF, RoBERTa, ChemBERTa, ...

ORGANOMETALLIC CHEMISTRY

??? [NOT IN KB] => OXOACIDS

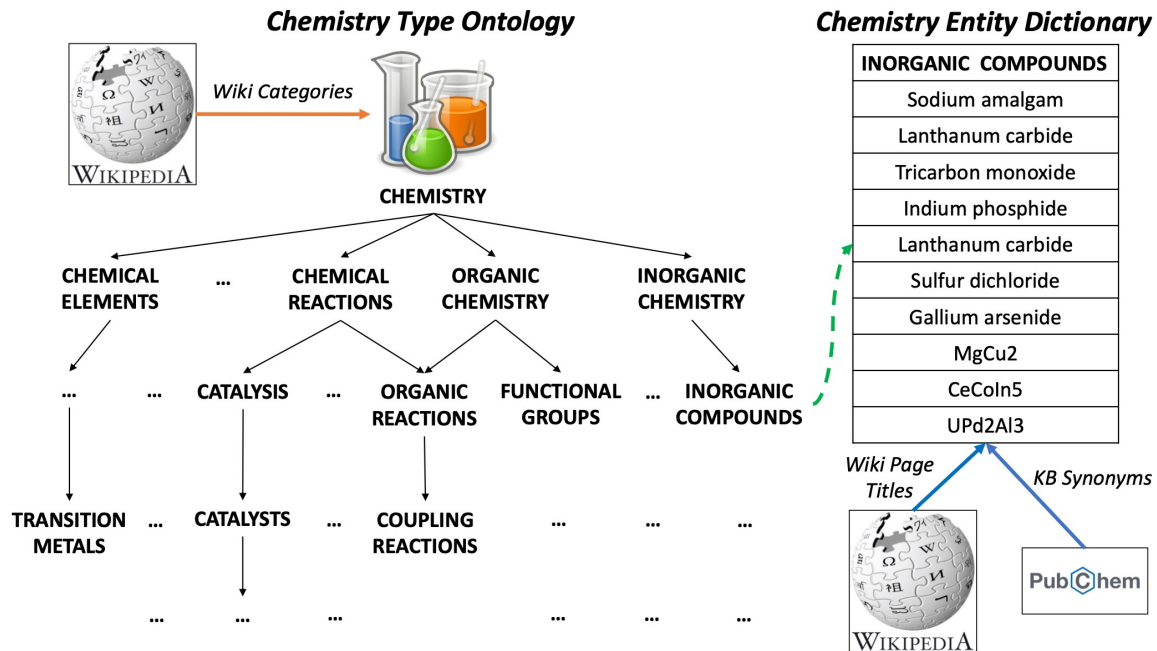
**S4:** ... can undergo a transmetalation with either BBA or the rapidly forming boronic acid ...

OXOACIDS

"either ... or ..." pattern learned by Sequence Labeling Model

# Chemistry Ontology

- ❑ Fine-grained chemistry type ontology:
  - ❑ Wikipedia categories rooted under *Chemistry*
  - ❑ Categories => Entity Types
  - ❑ Associated Page Titles => Entity Dictionaries
  - ❑ Expert proved 62 fine-grained types



## Category:Chemistry

From Wikipedia, the free encyclopedia

### Subcategories

This category has the following 73 subcategories, out of 73 total.

- ▶ [Chemistry literature](#) (3 C, 2 P)
- 
- ▶ [Chemists](#) (12 C, 3 P)
- \*
- ▶ [Chemistry set index pages](#) (1 C, 655 P)
- +
- ▶ [Chemical elements](#) (132 C, 127 P)
- A**
- ▶ [Acid–base chemistry](#) (5 C, 49 P)
- ▶ [Analytical chemistry](#) (19 C, 222 P)
- ▶ [Astrochemistry](#) (1 C, 38 P)
- ▶ [Atmospheric chemistry](#) (24 P)
- M**
- ▶ [Materials science](#) (35 C, 400 P)
- ▶ [Medicinal chemistry](#) (8 C, 77 P, 10 F)
- ▶ [Metallurgy](#) (14 C, 161 P)
- ▶ [Microwave chemistry](#) (4 P)
- ▶ [Chemical mixtures](#) (6 C, 44 P)
- ▶ [Molecular physics](#) (10 C, 79 P)
- ▶ [Molecules](#) (10 C, 20 P)
- N**
- ▶ [Chemical nomenclature](#) (4 C, 84 P)
- ▶ [Nuclear chemistry](#) (8 C, 59 P)

### Pages in category "Chemistry"

The following 132 pages are in this category, out of 132 total. This list may not reflect recent changes ([lea more](#)).

- [Chemistry](#)
- \*
- [Portal:Chemistry](#)
- 0–9**
- [2-Hexoxyethanol](#)
- A**
- [Acid–base reaction](#)
- [Actinide chemistry](#)
- [Allotropy](#)
- [Alloy](#)
- [Fluorine cycle](#)
- [Forensic chemistry](#)
- [Free element](#)
- G**
- [Geometry index](#)
- [Glossary of chemistry terms](#)
- [Gold cycle](#)
- [Green chemistry](#)
- H**
- [Harbi al-Himyari](#)
- I**
- [Iolomics](#)

# Chem NER: Performance Comparison

## Dataset:

- Training: **85,702** unlabeled sentences + **62** fine-grained chemistry types
- Test: **3,000** expert-annotated sentences

Method	Precision	Recall	F1 Score
KB-Matching	0.21	0.12	0.15
BiLSTM-CRF (2016)	0.22	0.10	0.14
RoBERTa (2019)	0.24	0.18	0.20
ChemBERTa (2020)	0.18	0.12	0.14
AutoNER (2018)	0.21	0.04	0.06
BOND (2020)	0.19	0.13	0.15
<b>ChemNER (2021)</b>	<b>0.69</b>	<b>0.34</b>	<b>0.46</b>

$$\text{Precision (P)} = \frac{\# \text{Truth Positive}}{\# \text{Prediction}}$$

$$\text{Recall (R)} = \frac{\# \text{True Positive}}{\# \text{Ground - Truth}}$$

$$\text{F1 Score} = \frac{2 \times P \times R}{P + R}$$

**+0.26 absolute F1 ↑**

# Outline

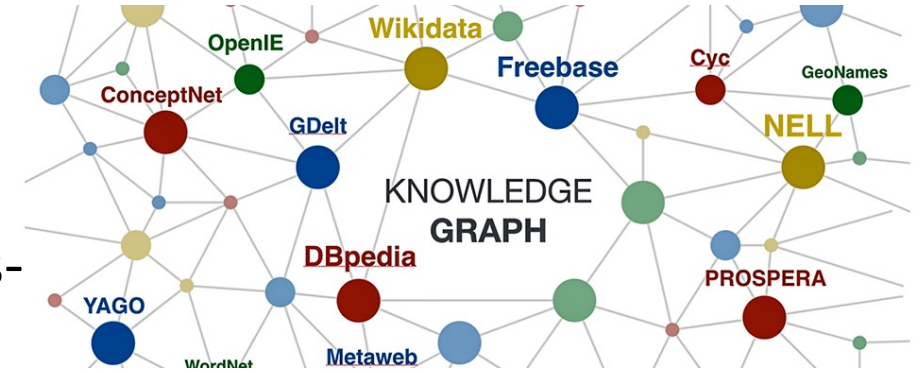
---

- ❑ **What Kinds of Text-Rich Information Networks Do We Really Need?**
- ❑ **Key Issue: Construction of Theme-/Corpus-Based Information Networks**
- ❑ **The Role of Embedding and PLM in Information Network Construction**
- ❑ **Data Preparation: Taxonomy-Guided Text Classification**
- ❑ **Identifying Information Network Primitives: Entities, Properties and Relations**
- ❑ **Conclusion: Towards Theme/Corpus-Based Information Network Construction**

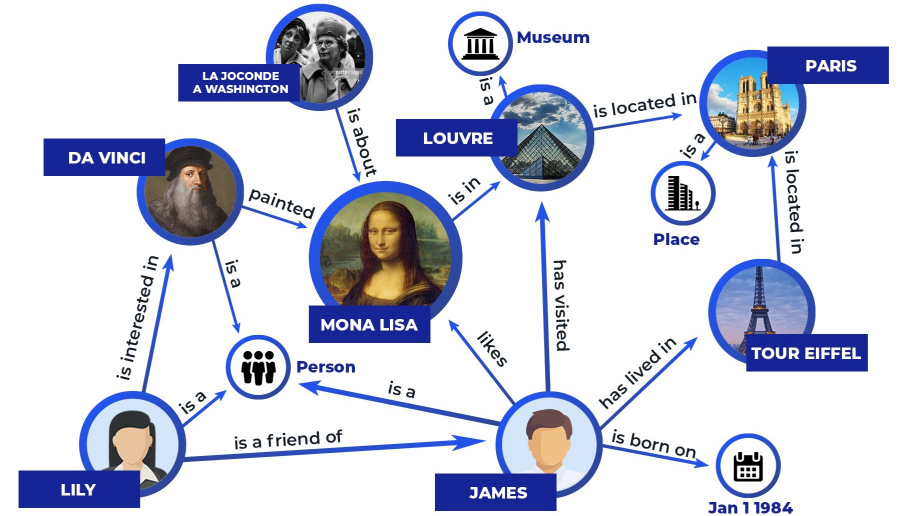


# Conclusions

- ❑ What kinds of info. networks do we really need?
  - ❑ Theme-/corpus-based info. networks
- ❑ Key issue: Automated construction of theme-/corpus-based info. networks from text
  - ❑ Exploring the power of embedding and Pre-trained Language Models (PLMs)
  - ❑ Collecting and preparing data using taxonomy-guided text classification
  - ❑ Identifying info. networks primitives: entities, properties and relations
- ❑ Towards theme/corpus-based info. networks construction



Typical KGs from Knowledge-Bases



Typed Entity-Relation-Property Graphs from Text

Ack. Figures are from Google images