



Part IV: Weakly-Supervised Text Classification: Embeddings with Less Human Effort

AAAI 2022 Tutorial


Pre-Trained Language Representations for Text Mining

Yu Meng, Jiaxin Huang, Yu Zhang, Jiawei Han

Computer Science, University of Illinois at Urbana-Champaign

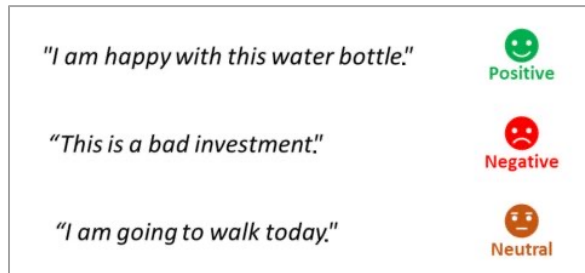
February 23, 2022

Outline

- ❑ What Weakly-Supervised Text Classification Is, and Why It Matters 
- ❑ Flat Text Classification
 - ❑ Embedding: WeSTClass [CIKM'18]
 - ❑ Pre-trained LM: ConWea [ACL'20], LOTClass [EMNLP'20], X-Class [NAACL'21]
- ❑ Text Classification with Taxonomy Information
 - ❑ Embedding: WeSHClass [AAAI'19]
 - ❑ Pre-trained LM: TaxoClass [NAACL'21]
- ❑ Text Classification with Metadata Information
 - ❑ Embedding: MetaCat [SIGIR'20]
 - ❑ Pre-trained LM: MICoL [WWW'22]

Text Classification

- Given a set of text units (e.g., documents, sentences) and a set of categories, the task is to assign relevant category/categories to each text unit
- Text Classification has a lot of downstream applications

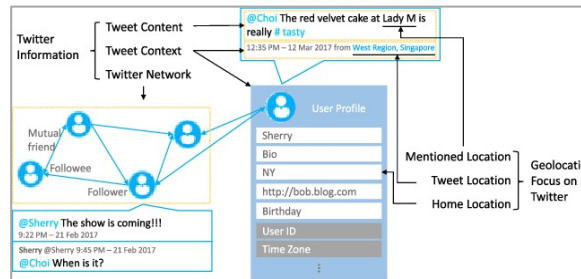


"I am happy with this water bottle." 😊 Positive

"This is a bad investment." 😞 Negative

"I am going to walk today." 😐 Neutral

Sentiment Analysis



Twitter Information: Tweet Content, Tweet Context, Twitter Network

Mutual friend, Follower, Followee

User Profile: Sherry, Bio, NY, http://bob.blog.com, Birthday, User ID, Time Zone

Geolocation Focus on Twitter: Mentioned Location, Tweet Location, Home Location

Location Prediction



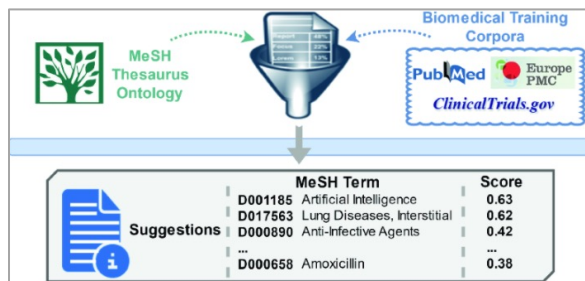
Raw Text Data

Technology

Sports

Fashion

News Topic Classification



Suggestions	MeSH Term	Score
	D001185 Artificial Intelligence	0.63
	D017563 Lung Diseases, Interstitial	0.62
	D000890 Anti-Infective Agents	0.42

	D000658 Amoxicillin	0.38

Paper Topic Classification

From: Jack
To: Alice
Subject: instruction copy
Monday 11:12 AM

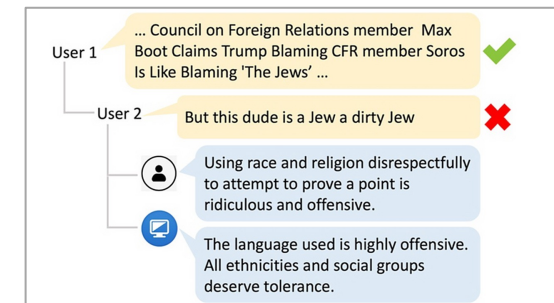
Hi Alice,

Thanks for the great tutorial about the new development tools. It was great catching up with you. We should get together sometime to continue our discussion soon. I am writing to check if you still have the document you mentioned about how to submit jobs to the new cluster? **If so, could you please share it with me?** I need to run some tests for the new feature.

Thank you!

Best,
Jack

Email Intent Identification



User 1: ... Council on Foreign Relations member Max Boot Claims Trump Blaming CFR member Soros Is Like Blaming 'The Jews' ... ✓

User 2: But this dude is a Jew a dirty Jew ✗

Using race and religion disrespectfully to attempt to prove a point is ridiculous and offensive.

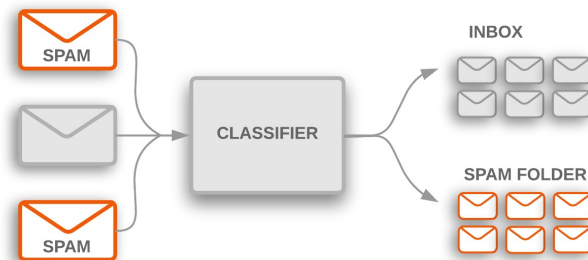
The language used is highly offensive. All ethnicities and social groups deserve tolerance.

Hate Speech Detection

Different Text Classification Settings: Single-Label vs. Multi-Label

❑ **Single-label:** Each document belongs to one category.

❑ Ex. Spam Detection



❑ **Multi-label:** Each document has multiple relevant labels.

❑ Ex. Paper Topic Classification

📄 BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Abstract

We introduce a new language representation model called BERT, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications. BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing the GLUE score to 80.5 (7.7 point absolute improvement), MultiNLI accuracy to 86.7% (4.6% absolute improvement), SQuAD v1.1 question answering Test F1 to 93.2 (1.5 point absolute improvement) and SQuAD v2.0 Test F1 to 83.1 (5.1 point absolute improvement).

Related Topics ⓘ

Question answering Language model Natural language understanding Named-entity recognition SemEval Inference Winograd Schema Challenge Sequence labeling

Artificial intelligence Computer science Transformer (machine learning model) View Less ^

<https://academic.microsoft.com/paper/2963341956/>

Different Text Classification Settings: Flat vs. Hierarchical

❑ **Flat:** All labels are at the same granularity level.

❑ Ex. Sentiment Analysis of E-Commerce Reviews (1-5 stars)

★★★★★ It works, it's nice, comfortable, and easy to type on. Not loud (unless you're a key pounder)

This keyboard works. It's comfortable, sensitive enough for touch typers, very quiet by comparison to other mechanicals (unless, of course, you're a 'key pounder'), and the lit keys are excellent for people like me who tend to prefer to work in a cave-like environment.

<https://www.amazon.com/gp/product/B089YFHYY5/>

❑ **Hierarchical:** Labels are organized into a hierarchy representing their parent-child relationship.

❑ Ex. Paper Topic Classification (the arXiv category taxonomy)

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

We introduce a new language representation model called BERT, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models, BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications. BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing the GLUE score to 80.5% (7.7% point absolute improvement), MultiNLI accuracy to 86.7% (4.6% absolute improvement), SQuAD v1.1 question answering Test F1 to 93.2 (1.5 point absolute improvement) and SQuAD v2.0 Test F1 to 83.1 (5.1 point absolute improvement).

Subjects: Computation and Language (cs.CL)

Cite as: arXiv:1810.04805 [cs.CL]

(or arXiv:1810.04805v2 [cs.CL] for this version)

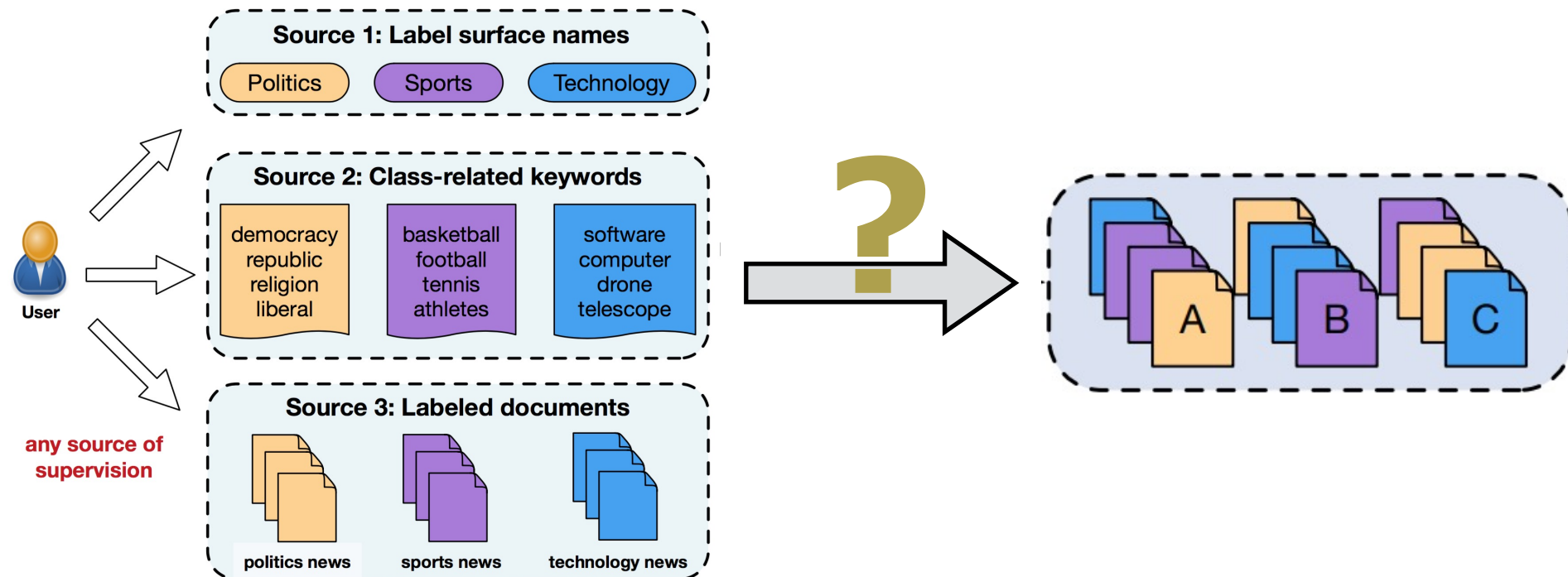
<https://arxiv.org/abs/1810.04805>

Weakly-Supervised Text Classification: Motivation

- ❑ Supervised text classification models (especially recent deep neural models) rely on a significant number of manually labeled training documents to achieve good performance.
- ❑ Collecting such training data is usually expensive and time-consuming. In some domains (e.g., scientific papers), annotations must be acquired from domain experts, which incurs additional cost.
- ❑ While users cannot afford to label sufficient documents for training a deep neural classifier, they can provide a small amount of seed information:
 - ❑ Category names or category-related keywords
 - ❑ A small number of labeled documents

Weakly-Supervised Text Classification: Definition


- Text classification without massive human-annotated training data
 - **Keyword-level weak supervision:** category names or a few relevant keywords
 - **Document-level weak supervision:** a small set of labeled docs



General Ideas to Perform Weakly-Supervised Text Classification

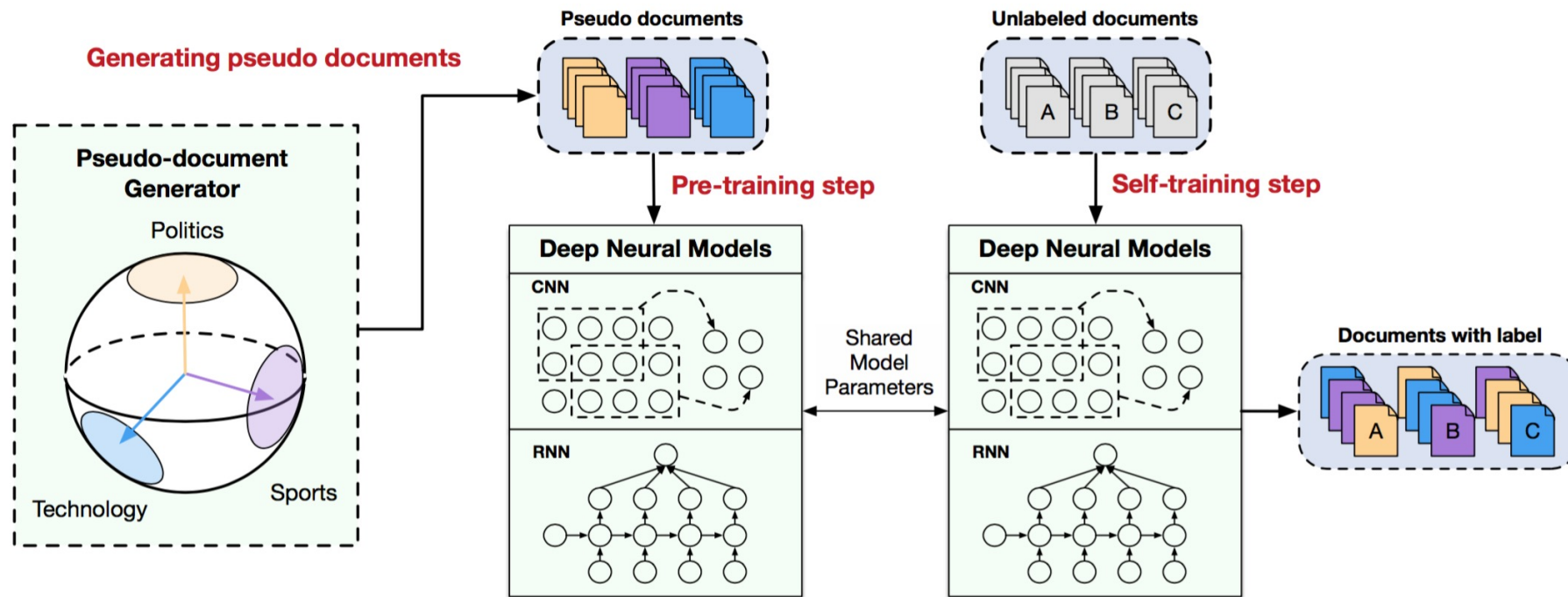
- ❑ Joint representation learning
 - ❑ Put words, labels, and/or documents into the same latent space using **embedding learning** or **pre-trained language models**
- ❑ Pseudo training data generation
 - ❑ Retrieve some unlabeled documents or synthesize some artificial documents using **text embeddings** or **contextualized representations**
 - ❑ Give them pseudo labels to train a text classifier
- ❑ Transfer the knowledge of **pre-trained language models** to classification tasks

Outline

- ❑ What Weakly-Supervised Text Classification Is, and Why It Matters
- ❑ Flat Text Classification
 - ❑ Embedding: WeSTClass [CIKM'18] 
 - ❑ Pre-trained LM: ConWea [ACL'20], LOTClass [EMNLP'20], X-Class [NAACL'21]
- ❑ Text Classification with Taxonomy Information
 - ❑ Embedding: WeSHClass [AAAI'19]
 - ❑ Pre-trained LM: TaxoClass [NAACL'21]
- ❑ Text Classification with Metadata Information
 - ❑ Embedding: MetaCat [SIGIR'20]
 - ❑ Pre-trained LM: MICoL [WWW'22]

WeSTClass: Pseudo Training Data + Self-Training

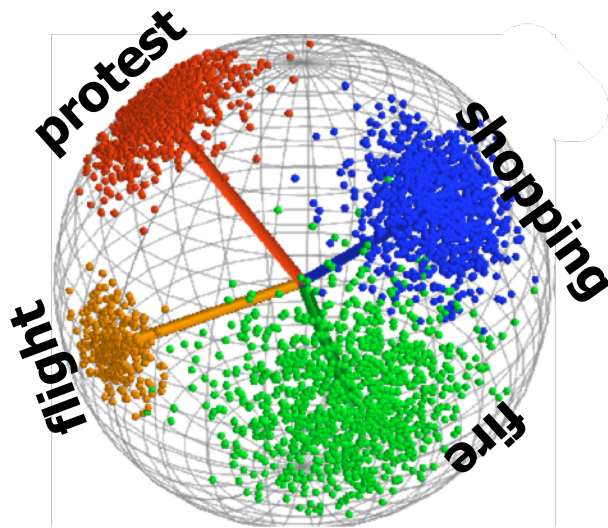
- ❑ Embed all words (including label names and keywords) into the same space
- ❑ Pseudo document generation: generate pseudo documents from seeds
- ❑ Self-training: train deep neural nets (CNN, RNN) with bootstrapping



Meng, Y., Shen, J., Zhang, C., & Han, J. "Weakly-supervised neural text classification", CIKM'18.
Applicable to both keyword-level and document-level supervision.

WeSTClass: Pseudo Document Generation

- Fit a von-Mises Fisher distribution for each category according to the keywords
 - Category name as supervision? Find nearest words as keywords
 - A few documents as supervision? Retrieve words with high TF-IDF scores
- Sample bag-of-keywords as pseudo documents for each class



Mean
direction

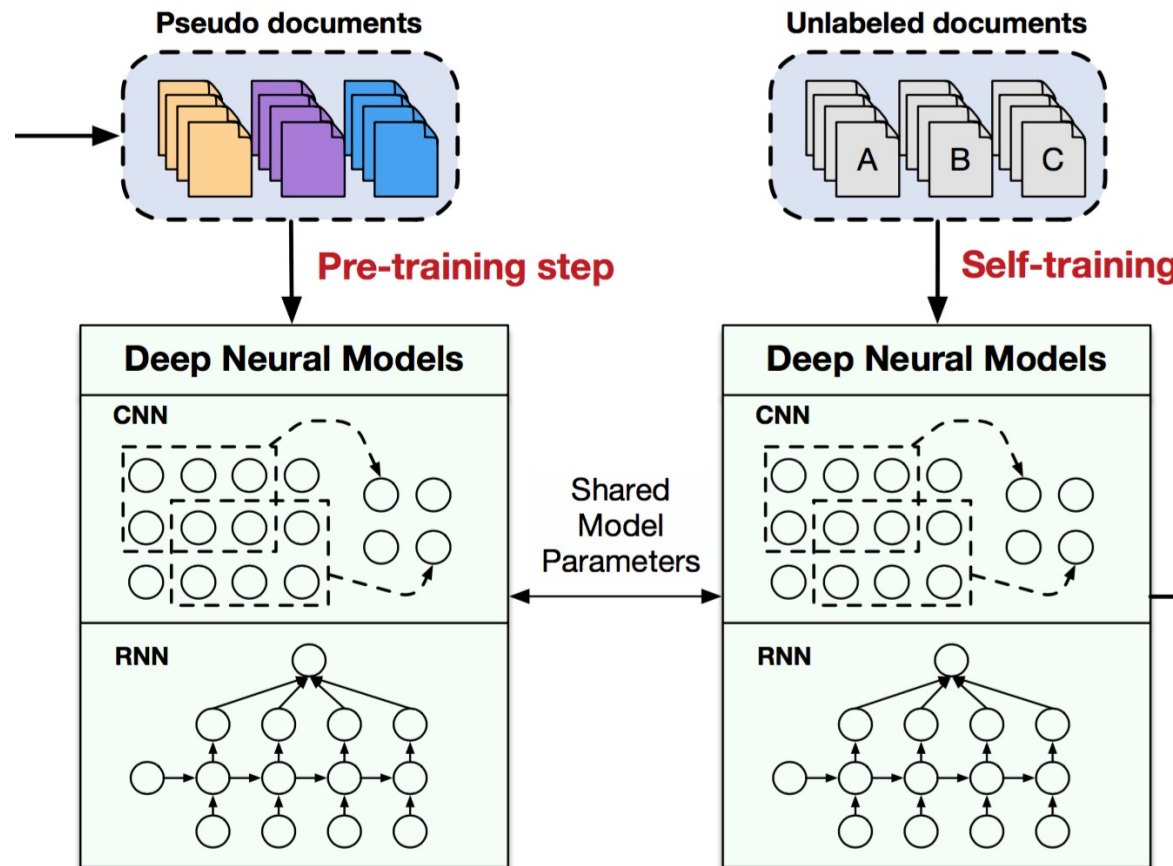
Concentration
parameter

$$p(\mathbf{x}|\mu, \kappa) = C_D(\kappa) \exp(\kappa \mu^T \mathbf{x})$$

$$C_D(\kappa) = \frac{\kappa^{D/2-1}}{I_{D/2-1}(\kappa)}$$

WeSTClass: Self-Training Deep Neural Nets

- ❑ **Pre-training:** Use pseudo documents to initialize DNNs (e.g., CNN, RNN)
- ❑ **Self-training:** Iteratively refine DNNs in a self-boosting fashion



new score of
label j for document i

$$l_{ij} = \frac{y_{ij}^2 / f_j}{\sum_{j'} y_{ij'}^2 / f_{j'}}$$

WeSTClass: Experiment Results

- Datasets: (1) NYT, (2) AG’s News, (3) Yelp
- Evaluation: use different types of weak supervision and measure accuracies

Macro-F1 scores:

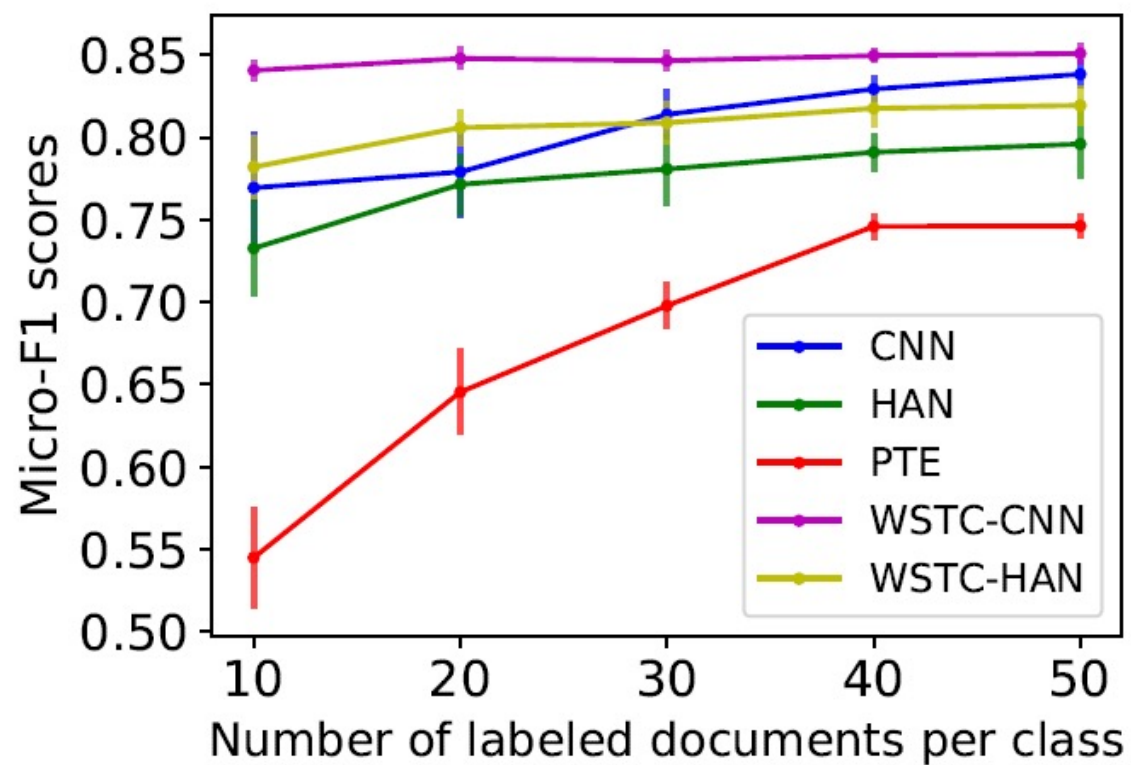
Methods	The New York Times			AG’s News			Yelp Review		
	LABELS	KEYWORDS	DOCS	LABELS	KEYWORDS	DOCS	LABELS	KEYWORDS	DOCS
IR with tf-idf	0.319	0.509	-	0.187	0.258	-	0.533	0.638	-
Topic Model	0.301	0.253	-	0.496	0.723	-	0.333	0.333	-
Dataless	0.484	-	-	0.688	-	-	0.337	-	-
UNEC	0.690	-	-	0.659	-	-	0.602	-	-
PTE	-	-	0.834 (0.024)	-	-	0.542 (0.029)	-	-	0.658 (0.042)
HAN	0.348	0.534	0.740 (0.059)	0.498	0.621	0.731 (0.029)	0.519	0.631	0.686 (0.046)
CNN	0.338	0.632	0.702 (0.059)	0.758	0.770	0.766 (0.035)	0.523	0.633	0.634 (0.096)
NoST-HAN	0.515	0.213	0.823 (0.035)	0.590	0.727	0.745 (0.038)	0.731	0.338	0.682 (0.090)
NoST-CNN	0.701	0.702	0.833 (0.013)	0.534	0.759	0.759 (0.032)	0.639	0.740	0.717 (0.058)
WESTCLASS-HAN	0.754	0.640	0.832 (0.028)	0.816	0.820	0.782 (0.028)	0.769	0.736	0.729 (0.040)
WESTCLASS-CNN	0.830	0.837	0.835 (0.010)	0.822	0.821	0.839 (0.007)	0.735	0.816	0.775 (0.037)

Micro-F1 scores:


IR with tf-idf	0.240	0.346	-	0.292	0.333	-	0.548	0.652	-
Topic Model	0.666	0.623	-	0.584	0.735	-	0.500	0.500	-
Dataless	0.710	-	-	0.699	-	-	0.500	-	-
UNEC	0.810	-	-	0.668	-	-	0.603	-	-
PTE	-	-	0.906 (0.020)	-	-	0.544 (0.031)	-	-	0.674 (0.029)
HAN	0.251	0.595	0.849 (0.038)	0.500	0.619	0.733 (0.029)	0.530	0.643	0.690 (0.042)
CNN	0.246	0.620	0.798 (0.085)	0.759	0.771	0.769 (0.034)	0.534	0.646	0.662 (0.062)
NoST-HAN	0.788	0.676	0.906 (0.021)	0.619	0.736	0.747 (0.037)	0.740	0.502	0.698 (0.066)
NoST-CNN	0.767	0.780	0.908 (0.013)	0.553	0.766	0.765 (0.031)	0.671	0.750	0.725 (0.050)
WESTCLASS-HAN	0.901	0.859	0.908 (0.019)	0.816	0.822	0.782 (0.028)	0.771	0.737	0.729 (0.040)
WESTCLASS-CNN	0.916	0.912	0.911 (0.007)	0.823	0.823	0.841 (0.007)	0.741	0.816	0.776 (0.037)

WeSTClass: Effect of # Labeled Documents

- Compare the performances of five methods on the AG's News dataset by varying the number of labeled documents per class and



Outline

- ❑ What Weakly-Supervised Text Classification Is, and Why It Matters
- ❑ Flat Text Classification
 - ❑ Embedding: WeSTClass [CIKM'18]
 - ❑ Pre-trained LM: ConWea [ACL'20], LOTClass [EMNLP'20], X-Class [NAACL'21] 
- ❑ Text Classification with Taxonomy Information
 - ❑ Embedding: WeSHClass [AAAI'19]
 - ❑ Pre-trained LM: TaxoClass [NAACL'21]
- ❑ Text Classification with Metadata Information
 - ❑ Embedding: MetaCat [SIGIR'20]
 - ❑ Pre-trained LM: MICoL [WWW'22]

Language Models for Weakly-Supervised Classification

- ❑ The previous approaches only use the local corpus
- ❑ Fail to take advantage of the general knowledge source (e.g., Wikipedia)
- ❑ Why general knowledge?
 - ❑ Humans can classify texts with general knowledge
 - ❑ Common linguistic features to understand texts better
 - ❑ Compensate for potential data scarcity of the local corpus
- ❑ How to use general knowledge?
 - ❑ Neural language models (e.g., BERT) are pre-trained on large-scale general knowledge texts
 - ❑ Their learned semantic/syntactic features can be transferred to downstream tasks

ConWea: Disambiguating User-Provided Keywords

- ❑ User-provided seed words may be ambiguous.

- ❑ Example:

Class	Seed words
Soccer	soccer, goal, penalty
Law	law, judge, court

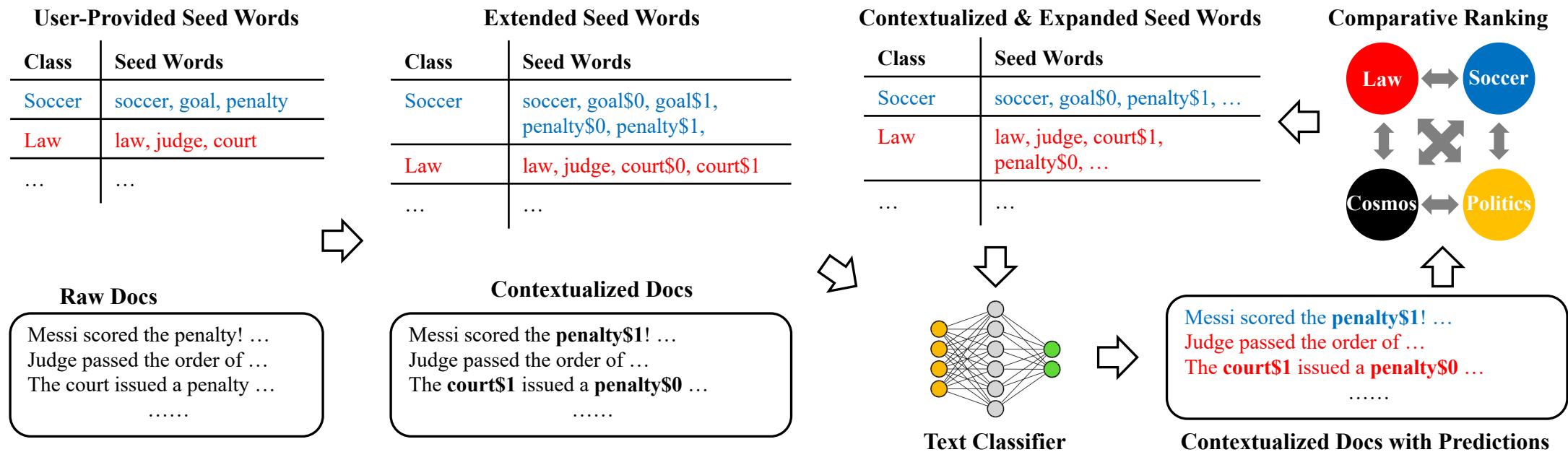
- ❑ Classify the following sentences:
 - ❑ Messi scored the penalty.
 - ❑ John was issued a death penalty.
- ❑ Disambiguate the “senses” based on contextualized representations

Mekala, D. & Shang, J. “Contextualized Weak Supervision for Text Classification”, ACL’20. [Keywords as supervision.](#)

ConWea-related slides credit to Jingbo Shang

ConWea: Clustering for Disambiguation

- For each word, find all its occurrences in the input corpus
- Run BERT to get their contextualized representations
- Run a clustering method (e.g., K-Means) to obtain clusters for different “senses”



ConWea: Experiment Results

□ Ablations:

- ConWea-NoCon: Variant of ConWea trained without contextualization.
- ConWea-NoExpan: Variant of ConWea trained without seed expansion.
- ConWea-WSD: Variant of ConWea with contextualization replaced by a word sense disambiguation algorithm.

		NYT				20 Newsgroup			
		5-Class (Coarse)		25-Class (Fine)		6-Class (Coarse)		20-Class (Fine)	
Methods		Micro-F ₁	Macro-F ₁	Micro-F ₁	Macro-F ₁	Micro-F ₁	Macro-F ₁	Micro-F ₁	Macro-F ₁
Baselines	IR-TF-IDF	0.65	0.58	0.56	0.54	0.49	0.48	0.53	0.52
	Dataless	0.71	0.48	0.59	0.37	0.50	0.47	0.61	0.53
	Word2Vec	0.92	0.83	0.69	0.47	0.51	0.45	0.33	0.33
	Doc2Cube	0.71	0.38	0.67	0.34	0.40	0.35	0.23	0.23
	WeSTClass	0.91	0.84	0.50	0.36	0.53	0.43	0.49	0.46
	ConWea	0.95	0.89	0.91	0.79	0.62	0.57	0.65	0.64
Ablations	ConWea-NoCon	0.91	0.83	0.89	0.74	0.53	0.50	0.58	0.57
	ConWea-NoExpan	0.92	0.85	0.76	0.66	0.58	0.53	0.58	0.57
	ConWea-WSD	0.83	0.78	0.72	0.64	0.52	0.46	0.49	0.47
Upper bound	HAN-Supervised	0.96	0.92	0.94	0.82	0.90	0.88	0.83	0.83

LOTClass: Find Similar Meaning Words with Label Names

- Find topic words based on label names
 - Overcome the low semantic coverage of label names
- Use language models to predict what words can replace the label names
 - Interchangeable words are likely to have similar meanings

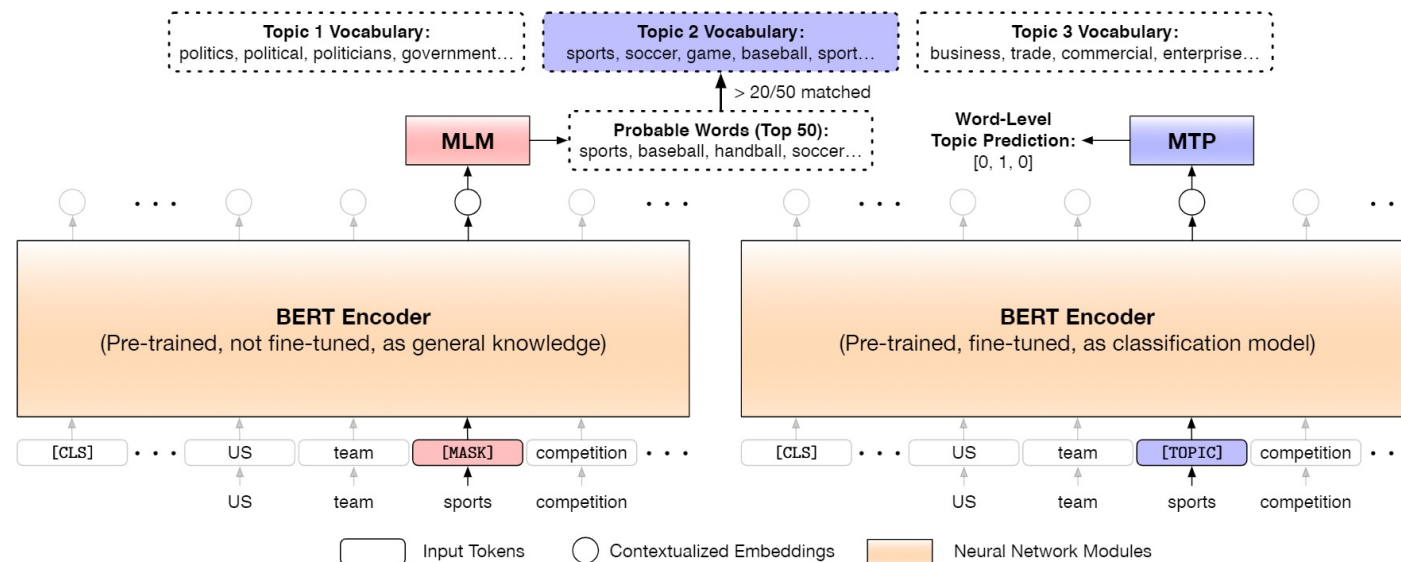
Sentence	Language Model Prediction
The oldest annual US team sports competition that includes professionals is not in baseball, or football or basketball or hockey. It's in soccer.	sports, baseball, handball, soccer, basketball, football, tennis, sport, championship, hockey, ...
Samsung's new SPH-V5400 mobile phone sports 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, ...

Table 1: BERT language model prediction (sorted by probability) for the word to appear at the position of “sports” under different contexts. The two sentences are from *AG News* corpus.

Meng, Y., Zhang, Y., Huang, J., Xiong, C., Ji, H., Zhang, C., & Han, J. “Text Classification Using Label Names Only: A Language Model Self-Training Approach”, EMNLP’20. [Category names as supervision.](#)

LOTClass: Contextualized Word-Level Topic Prediction

- ❑ Context-free matching of topic words is inaccurate
 - ❑ “Sports” does not always imply the topic “sports”
- ❑ Contextualized topic prediction:
 - ❑ Predict a word’s implied topic under specific contexts
 - ❑ We regard a word as “topic indicative” only when its top replacing words have enough overlap with the topic vocabulary.



LOTClass: Experiment Results

- Achieve around 90% accuracy on four benchmark datasets by only using at most 3 words (1 in most cases) per class as the label name
- Outperforming previous weakly-supervised approaches significantly
- Comparable to state-of-the-art semi-supervised models

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
	Ours w/o. self train	0.822	0.850	0.844	0.781
	Ours	0.864	0.889	0.894	0.906
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	0.937	0.972

How Powerful Are Vanilla BERT Representations in Category Prediction?

- An average of BERT representations of all tokens in a sentence/document preserves domain information well

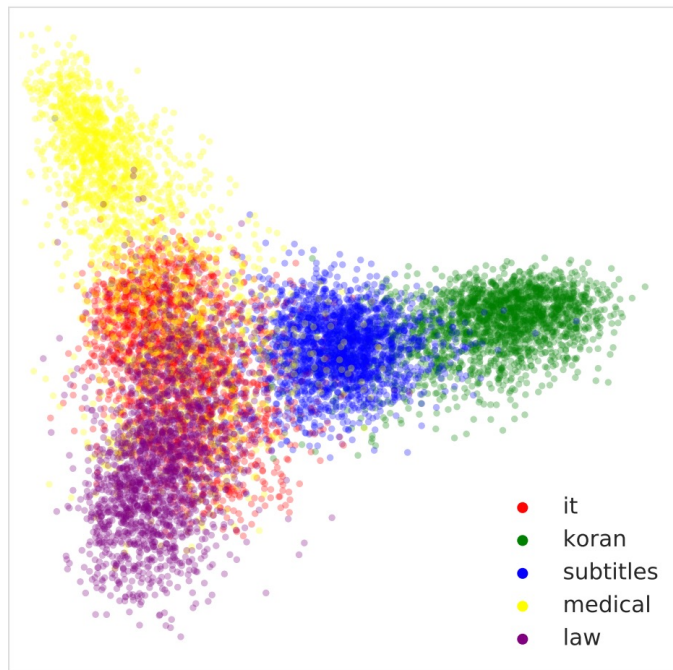


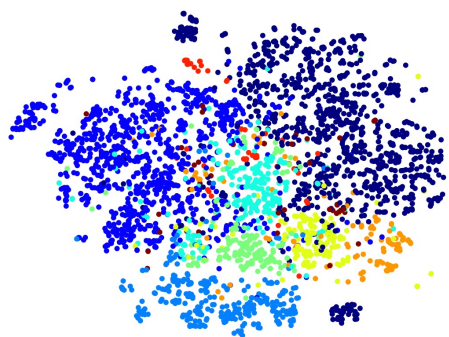
Figure 1: A 2D visualization of average-pooled BERT hidden-state sentence representations using PCA. The colors represent the domain for each sentence.

True label \ Predicted label	it	koran	subtitles	medical	law
it	1927	0	55	16	2
koran	4	1767	225	0	4
subtitles	47	21	1918	9	5
medical	340	0	82	1413	165
law	206	0	10	58	1726

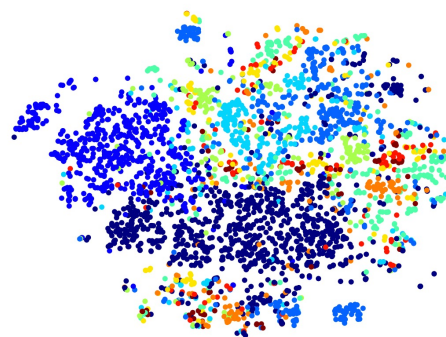
Figure 2: A confusion matrix for clustering with k=5 using BERT-base.

X-Class: Class-Oriented BERT Representations

- A simple idea for text classification
 - Learn representations for documents
 - Set the number of clusters as the number of classes
 - Hope their clustering results are almost the same as the desired classification
- However, the same corpus could be classified differently



(a) NYT-Topics



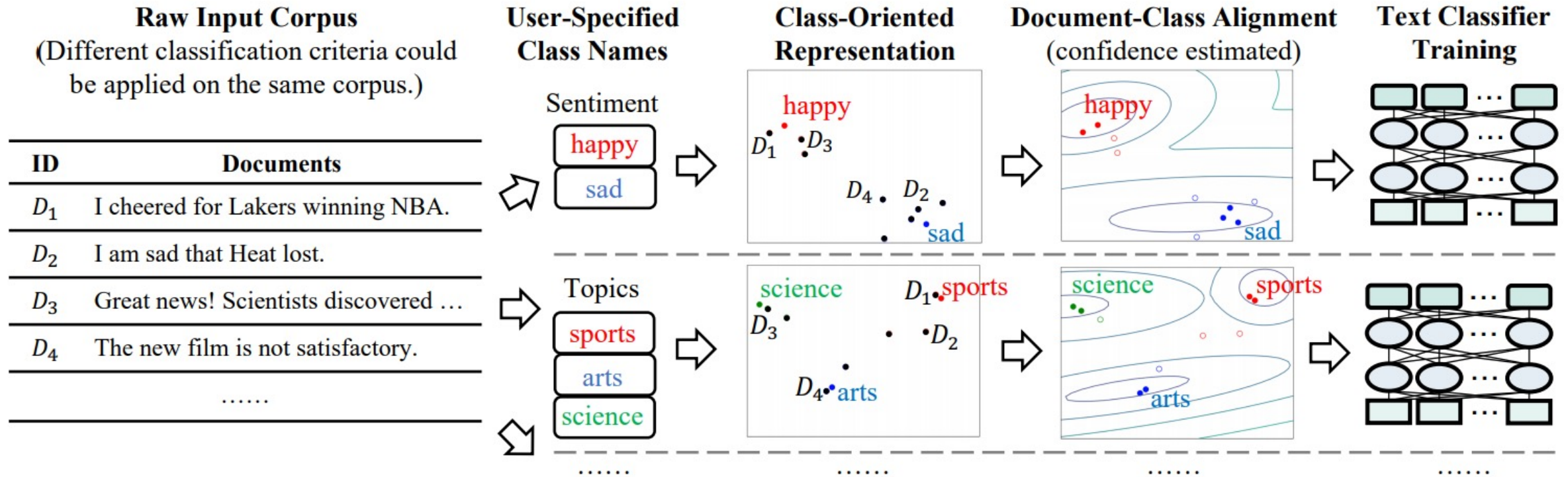
(b) NYT-Locations

Figure 1: Visualizations of News using Average BERT Representations. Colors denote different classes.

Wang, Z., Mekala, D., & Shang, J. "X-Class: Text Classification with Extremely Weak Supervision", NAACL'21. [Category Names as supervision.](#)
X-Class-related slides credit to Jingbo Shang

X-Class: Class-Oriented BERT Representations

- Clustering for classification based on class-oriented representations




X-Class: Experiment Results

- WeSTClass & ConWea consume at least 3 seed words per class
- LOTClass & X-Class use category names only

	AGNews	20News	NYT-Small	NYT-Topic	NYT-Location	Yelp	DBpedia
Corpus Domain	News	News	News	News	News	Reviews	Wikipedia
Class Criterion	Topics	Topics	Topics	Topics	Locations	Sentiment	Ontology
# of Classes	4	5	5	9	10	2	14
# of Documents	120,000	17,871	13,081	31,997	31,997	38,000	560,000
Imbalance	1.0	2.02	16.65	27.09	15.84	1.0	1.0

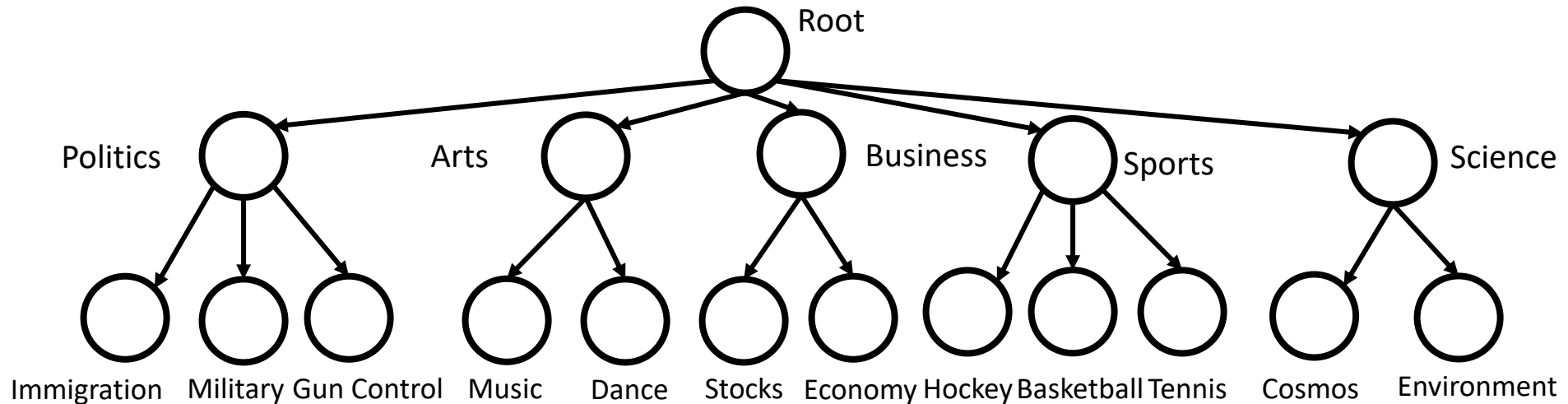
Model	AGNews	20News	NYT-Small	NYT-Topic	NYT-Location	Yelp	DBpedia
Supervised	93.99/93.99	96.45/96.42	97.95/95.46	94.29/89.90	95.99/94.99	95.7/95.7	98.96/98.96
WeSTClass	82.3/82.1	71.28/69.90	91.2/83.7	68.26/57.02	63.15/53.22	81.6/81.6	81.1/ N/A
ConWea	74.6/74.2	75.73/73.26	95.23/90.79	81.67/71.54	85.31/83.81	71.4/71.2	N/A
LOTClass	86.89/86.82	73.78/72.53	78.12/56.05	67.11/43.58	58.49/58.96	87.75/87.68	86.66/85.98
X-Class	84.8/84.65	81.36/80.6	96.67/92.98	80.6/69.92	90.5/89.81	88.36/88.32	91.33/91.14
X-Class-Rep	77.92/77.03	75.14/73.24	92.13/83.94	77.85/65.38	86.7/87.36	77.87/77.05	74.06/71.75
X-Class-Align	83.1/83.05	79.28/78.62	96.34/92.08	79.64/67.85	88.58/88.02	87.16/87.1	87.37/87.28

Outline

- ❑ What Weakly-Supervised Text Classification Is, and Why It Matters
- ❑ Flat Text Classification
 - ❑ Embedding: WeSTClass [CIKM'18]
 - ❑ Pre-trained LM: ConWea [ACL'20], LOTClass [EMNLP'20], X-Class [NAACL'21]
- ❑ Text Classification with Taxonomy Information
 - ❑ Embedding: WeSHClass [AAAI'19] 
 - ❑ Pre-trained LM: TaxoClass [NAACL'21]
- ❑ Text Classification with Metadata Information
 - ❑ Embedding: MetaCat [SIGIR'20]
 - ❑ Pre-trained LM: MICoL [WWW'22]

WeSHClass: Weakly-Supervised Hierarchical Text Classification

- The hierarchy has a **tree** structure. Each document is associated with **one path** starting from the root node. (E.g., the main subject of each arXiv paper.)



- Keyword-level weak supervision: The name of each node in the taxonomy, or a few keywords for each leaf category
- Document-level weak supervision: A few labeled documents for each leaf category

Meng, Y., Shen, J., Zhang, C., & Han, J. "Weakly-Supervised Hierarchical Text Classification", AAAI'19.

Applicable to both keyword-level and document-level supervision.

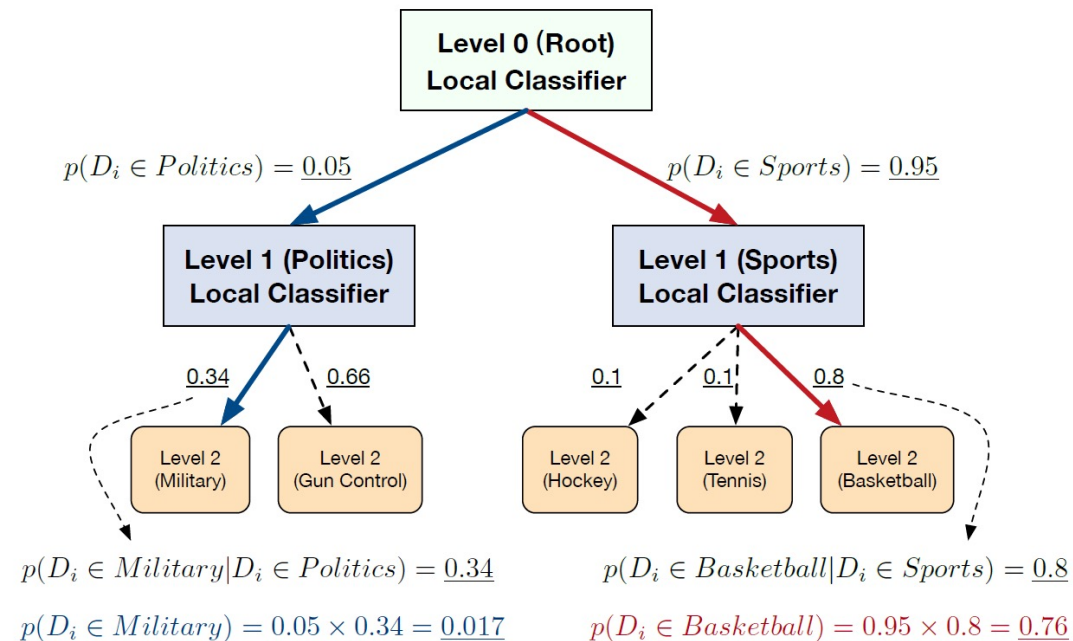
WeSHClass: Hierarchical Classification Model

□ Local Classifier Pre-training

- Generate β pseudo documents per class (recall WeSTClass) to pre-train the local classifier
- A naive way of creating the label for a pseudo document D_i^* :
 - Directly use the associated class label it is generated from; one-hot encodings;
 - Problem: classifier overfitting to pseudo documents
- Instead, use pseudo labels:
 - $$l_{ij} = \begin{cases} (1 - \alpha) + \alpha/m & D_i^* \text{ is generated from class } j \\ \alpha/m & \text{otherwise} \end{cases} .$$
 - α accounts for the “noises” in pseudo documents; it is evenly split into all m classes
- Pre-training is performed by minimizing KL divergence loss to pseudo labels

WeSHClass: Hierarchical Classification Model

- Global Classifier Per Level
 - At each level k in the class taxonomy, construct a global classifier by ensembling all local classifiers from root to level k
 - Use unlabeled documents to bootstrap the global classifier



WeSHClass: Hierarchical Classification Model

□ Global Classifier Construction

- The multiplication operation can be explained by the conditional probability formula:

$$p(D_i \in C_{child}) = p(D_i \in C_{child} | D_i \in C_{parent})p(D_i \in C_{parent})$$

- All local classifiers from root to level k are fine-tuned simultaneously via back-propagation during self-training; misclassifications at higher levels can be corrected

□ Global Classifier Self-training

- Step 1: Use the pre-trained global classifier to classify all unlabeled documents in the corpus;
- Step 2: Compute pseudo labels based on current predictions:

$$l_{ij} = \frac{y_{ij}^2 / f_j}{\sum_{j'} y_{ij'}^2 / f_{j'}} \text{ where } f_j = \sum_i y_{ij} \text{ and } y_{ij} \text{ is the current prediction}$$

- Step 3: Minimize KL divergence loss to pseudo labels
- Iterate between Steps 2 and 3 until less than $\delta\%$ of documents in the corpus have class assignment changes

WeSHClass: Experiment Results


□ Datasets

□ New York Times; arXiv; Yelp Review

□ Evaluation: Micro-F1 and Macro-F1 among all classes

Methods	NYT				arXiv				Yelp Review			
	KEYWORDS		DOCS		KEYWORDS		DOCS		KEYWORDS		DOCS	
	Macro	Micro	Macro Avg. (Std.)	Micro Avg. (Std.)	Macro	Micro	Macro Avg. (Std.)	Micro Avg. (Std.)	Macro	Micro	Macro Avg. (Std.)	Micro Avg. (Std.)
Hier-Dataless	0.593	0.811	-	-	0.374	0.594	-	-	0.284	0.312	-	-
Hier-SVM	-	-	0.142 (0.016)	0.469 (0.012)	-	-	0.049 (0.001)	0.443 (0.006)	-	-	0.220 (0.082)	0.310 (0.113)
CNN	-	-	0.165 (0.027)	0.329 (0.097)	-	-	0.124 (0.014)	0.456 (0.023)	-	-	0.306 (0.028)	0.372 (0.028)
WeSTClass	0.386	0.772	0.479 (0.027)	0.728 (0.036)	0.412	0.642	0.264 (0.016)	0.547 (0.009)	0.348	0.389	0.345 (0.027)	0.388 (0.033)
No-global	0.618	0.843	0.520 (0.065)	0.768 (0.100)	0.442	0.673	0.264 (0.020)	0.581 (0.017)	0.391	0.424	0.369 (0.022)	0.403 (0.016)
No-VMF	0.628	0.862	0.527 (0.031)	0.825 (0.032)	0.406	0.665	0.255 (0.015)	0.564 (0.012)	0.410	0.457	0.372 (0.029)	0.407 (0.015)
No-self-train	0.550	0.787	0.491 (0.036)	0.769 (0.039)	0.395	0.635	0.234 (0.013)	0.535 (0.010)	0.362	0.408	0.348 (0.030)	0.382 (0.022)
Our method	0.632	0.874	0.532 (0.015)	0.827 (0.012)	0.452	0.692	0.279 (0.010)	0.585 (0.009)	0.423	0.461	0.375 (0.021)	0.410 (0.014)

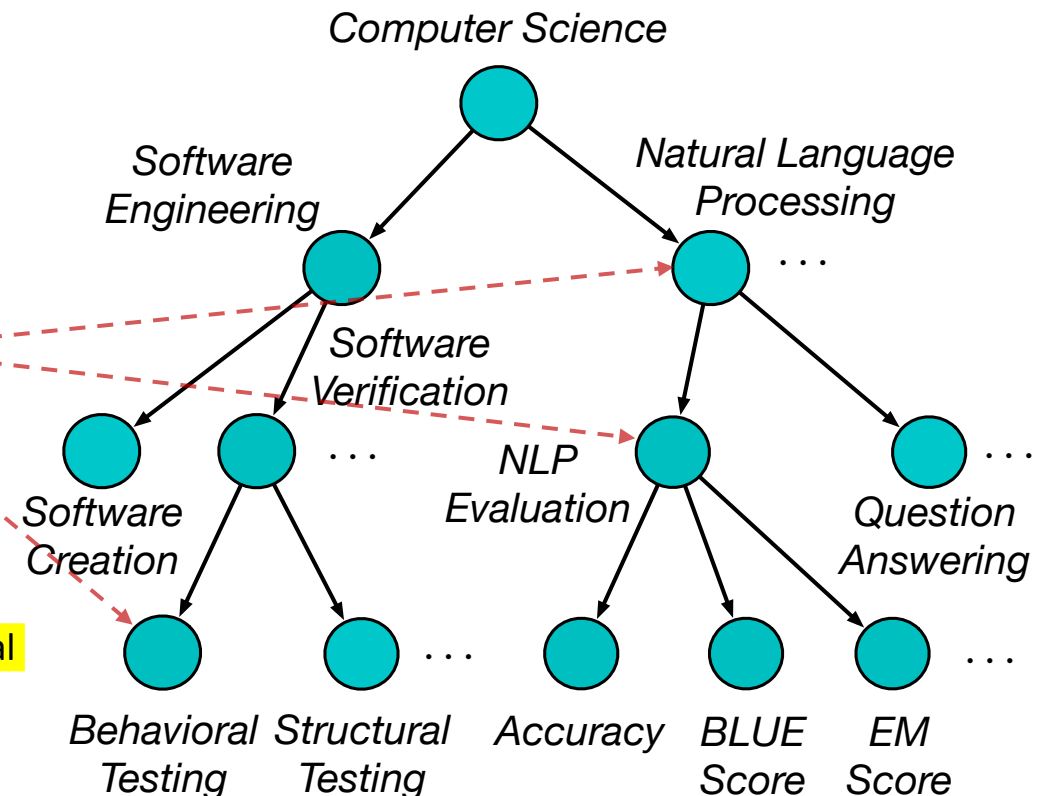
Outline

- ❑ What Weakly-Supervised Text Classification Is, and Why It Matters
- ❑ Flat Text Classification
 - ❑ Embedding: WeSTClass [CIKM'18]
 - ❑ Pre-trained LM: ConWea [ACL'20], LOTClass [EMNLP'20], X-Class [NAACL'21]
- ❑ Text Classification with Taxonomy Information
 - ❑ Embedding: WeSHClass [AAAI'19]
 - ❑ Pre-trained LM: TaxoClass [NAACL'21] 
- ❑ Text Classification with Metadata Information
 - ❑ Embedding: MetaCat [SIGIR'20]
 - ❑ Pre-trained LM: MICoL [WWW'22]

TaxoClass: Weakly-supervised Hierarchical Multi-Label Text Classification

- ❑ The taxonomy is a directed acyclic graph (DAG)
- ❑ Each paper can have multiple categories distributed on different paths
- ❑ Category names can be phrases and may not appear in the corpus

Document
Measuring held-out accuracy often overestimates the performance of <i>NLP</i> models... Inspired by principles of <i>behavioral testing</i> in software engineering, we introduce CheckList, a task-agnostic methodology for <i>testing NLP models</i> ...



Shen, J., Qiu, W., Meng, Y., Shang, J., Ren, X., & Han, J., "TaxoClass: Hierarchical Multi-Label Text Classification Using Only Class Names", NAACL'21. **Category names as supervision.**

TaxoClass: Why Category Names Only?

- ❑ Taxonomies for multi-label text classification are often big.
 - ❑ Amazon Product Catalog: $\times 10^4$ categories
 - ❑ MeSH Taxonomy (for medical papers): $\times 10^4$ categories
 - ❑ Microsoft Academic Taxonomy: $\times 10^5$ labels
- ❑ Impossible for users to provide even a small set of (e.g., 3) keywords/labeled documents for each category

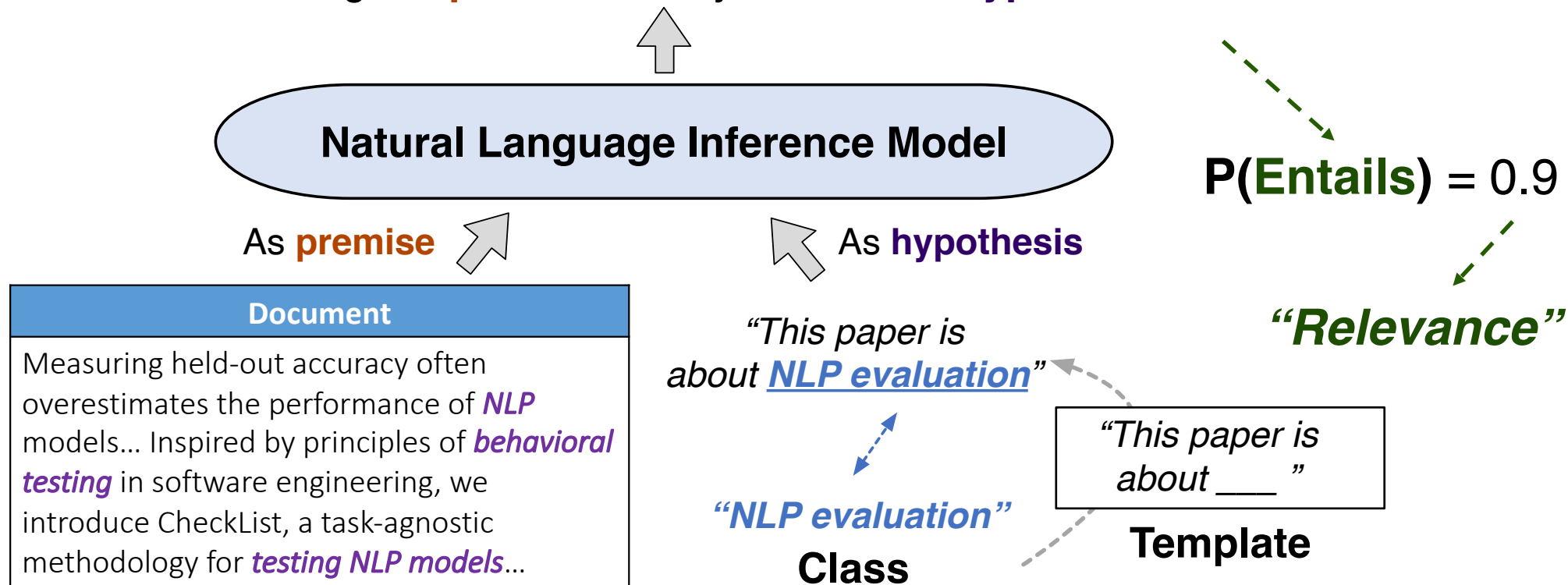


<https://academic.microsoft.com/home>

TaxoClass: Document-Class Relevance Calculation

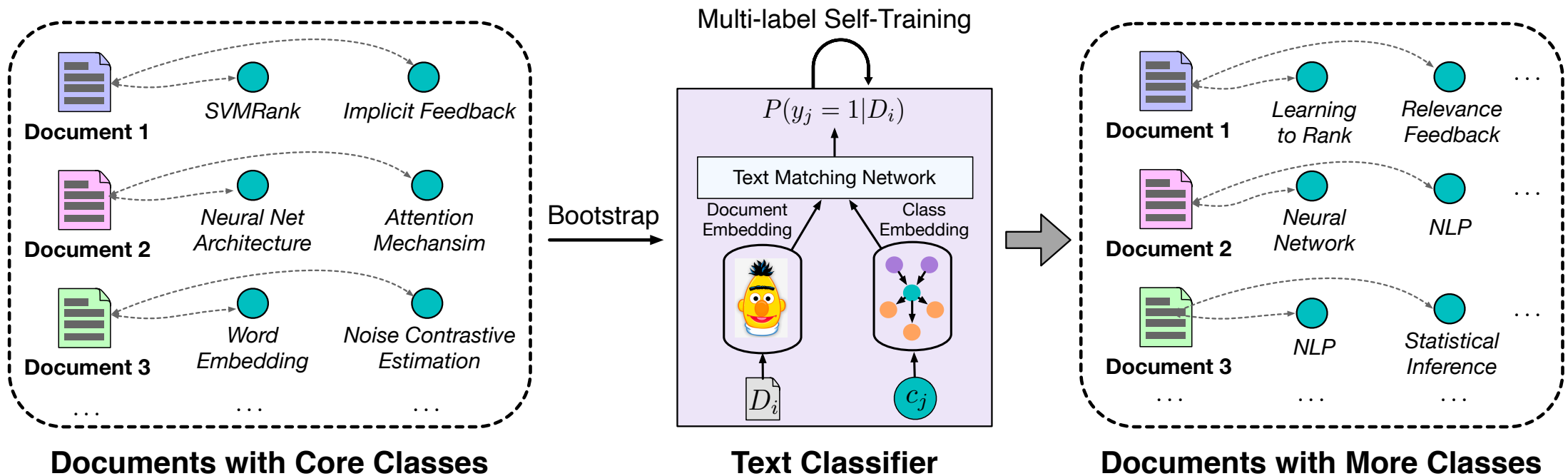
- How to use the knowledge from pre-trained LMs?
- Relevance model: BERT/RobERTa fine-tuned on the NLI task
- <https://huggingface.co/roberta-large-mnli>

After reading the **premise**, can you infer the **hypothesis**?



TaxoClass: Identify Core Classes and More Classes

- Identify document core classes in reduced label search space
- Generalize from core classes with bootstrapping and self-training



TaxoClass: Experiment Results

Weakly-supervised multi-class classification method

Semi-supervised methods using 30% of training set

Zero-shot method

Methods	Amazon		DBPedia	
	Example-F1	P@1	Example-F1	P@1
WeSHClass (Meng et al., AAAI'19)	0.246	0.577	0.305	0.536
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
Hier-0Shot-TC (Yin et al., EMNLP'19)	0.474	0.714	0.677	0.787
TaxoClass (ours)	0.593	0.812	0.816	0.894


- **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

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

Outline

- ❑ What Weakly-Supervised Text Classification Is, and Why It Matters
- ❑ Flat Text Classification
 - ❑ Embedding: WeSTClass [CIKM'18]
 - ❑ Pre-trained LM: ConWea [ACL'20], LOTClass [EMNLP'20], X-Class [NAACL'21]
- ❑ Text Classification with Taxonomy Information
 - ❑ Embedding: WeSHClass [AAAI'19]
 - ❑ Pre-trained LM: TaxoClass [NAACL'21]
- ❑ Text Classification with Metadata Information
 - ❑ Embedding: MetaCat [SIGIR'20] 
 - ❑ Pre-trained LM: MICoL [WWW'22]

MetaCat: Leveraging Metadata for Categorization

❑ Metadata is prevalent in many text sources

❑ **GitHub repositories:** User, Tag

❑ **Tweets:** User, Hashtag

❑ **Amazon reviews:** User, Product

❑ **Scientific papers:** Author, Venue

Zhang, Y., Meng, Y., Huang, J., Xu, F.F., Wang, X., & Han, J. “Minimally Supervised Categorization of Text with Metadata”, SIGIR’20.
A few labeled documents as supervision.

❑ How to leverage these heterogenous signals in the categorization process?

The screenshot shows a GitHub repository page for 'dodan'. Annotations include: 'User' pointing to the repository owner 'dodan'; 'Description (Text)' pointing to the repository description; 'Tags' pointing to the repository tags; and 'README (Text)' pointing to the README file content.

(a) GITHUB REPOSITORY

The screenshot shows a Twitter profile for 'Anna Mandelbaum' (@notdjam) and a tweet. Annotations include: 'User' pointing to the profile; 'Tweet (Text)' pointing to the tweet content; and 'Tags' pointing to the hashtags in the tweet.

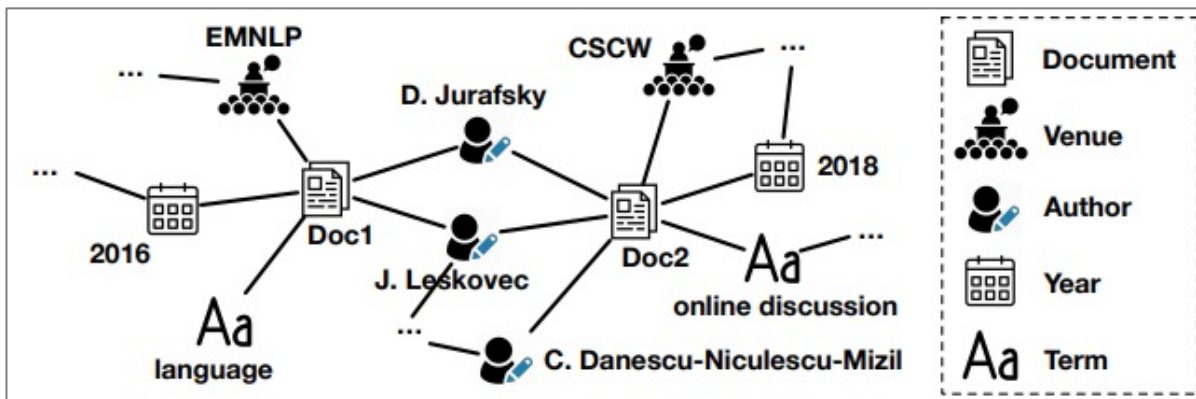
(b) TWEET

The screenshot shows an Amazon product page for the book 'Deep Learning (Adaptive Computation and Machine Learning series)'. Annotations include: 'Product' pointing to the book title; 'User' pointing to the reviewer's name; 'Title (Text)' pointing to the book title; and 'Review (Text)' pointing to the reviewer's text.

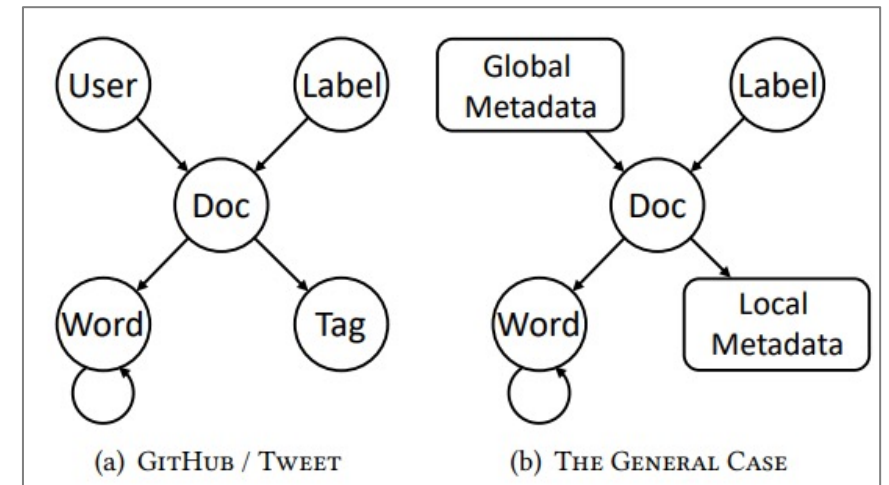
(c) AMAZON REVIEW

MetaCat: The Underlying Generative Process

- Two categories of metadata:
 - Global metadata:** user/author, product
 - “Causes” the generation of documents. (E.g., User/Author -> Document)
 - Local metadata:** tag/hashtag
 - “Describes” the documents. (E.g., Document -> Tag)
- We can also say “labels” are global, and “words” are local



A network view of corpus with metadata



A generative-process view of corpus with metadata

MetaCat: The Underlying Generative Process

□ We use GitHub/Tweet as a specific example to illustrate the process.

□ Step 1: **User** (Global Metadata) & **Label** -> **Document**

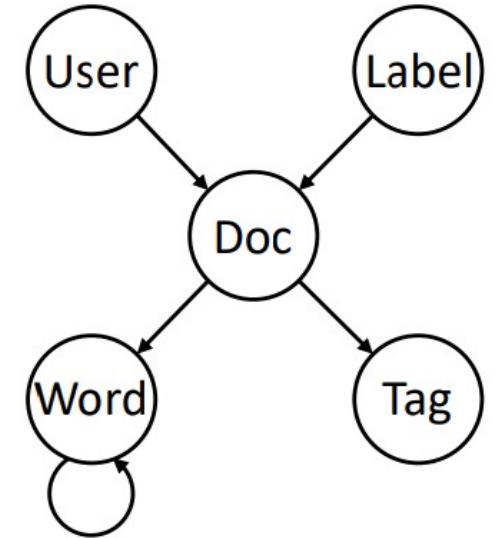
$$p(d|u, l) \propto \exp(\mathbf{e}_d^T \mathbf{e}_u) \cdot \exp(\mathbf{e}_d^T \mathbf{e}_l)$$

□ Step 2: Document -> Word

$$p(w|d) \propto \exp(\mathbf{e}_w^T \mathbf{e}_d)$$

□ Step 3: Document -> Tag (Local Metadata)

□ Step 4: Word -> Context



(a) GITHUB / TWEET

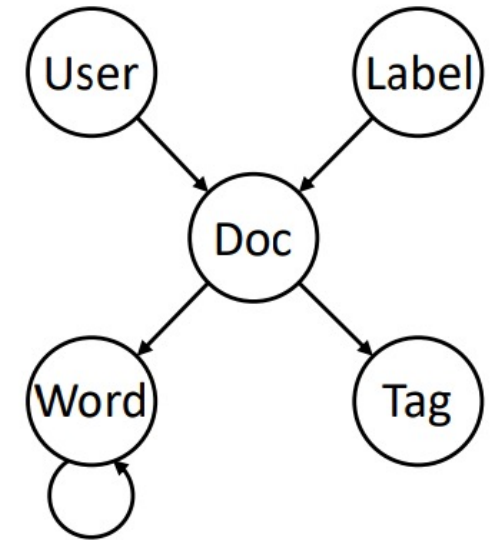
MetaCat: How to use this underlying model?

□ Embedding Learning Module

- All embedding vectors e_u, e_l, e_d, e_t, e_w are parameters of the generative process
- Learn the embedding vectors through maximizing the likelihood of observing all text and metadata

□ Training Data **Generation** Module

- e_u, e_l, e_d, e_t, e_w have been learned
- Given a label l , generate d, w and t according to the generative process



(a) GITHUB / TWEET

MetaCat: Experiment Results

□ Metadata is more helpful on smaller corpora.

□ Datasets

□ GitHub-Bio: 10 categories;
876 docs

□ GitHub-AI: 14 categories;
1,596 docs

□ GitHub-Sec: 3 categories;
84,950 docs

□ Amazon: 10 categories;
100,000 docs

□ Twitter: 9 categories;
135,619 docs

Table 2: Micro F1 scores of compared algorithms on the five datasets. “-”: excessive memory requirements.

Type	Method	GitHub-Bio	GitHub-AI	GitHub-Sec	Amazon	Twitter
Text-based	CNN [12]	0.2227 ± 0.0195	0.2404 ± 0.0404	0.4909 ± 0.0489	0.4915 ± 0.0374	0.3106 ± 0.0613
	HAN [38]	0.1409 ± 0.0145	0.1900 ± 0.0299	0.4677 ± 0.0334	0.4809 ± 0.0372	0.3163 ± 0.0878
	PTE [32]	0.3170 ± 0.0516	0.3511 ± 0.0403	0.4551 ± 0.0249	0.2997 ± 0.0786	0.1945 ± 0.0250
	WeSTClass [23]	0.3680 ± 0.0138	0.5036 ± 0.0287	0.6146 ± 0.0084	0.5312 ± 0.0161	0.3568 ± 0.0178
	PCEM [36]	0.3426 ± 0.0160	0.4820 ± 0.0292	0.5912 ± 0.0341	0.4645 ± 0.0163	0.2387 ± 0.0344
	BERT [4]	0.2680 ± 0.0303	0.2451 ± 0.0273	0.5538 ± 0.0368	0.5240 ± 0.0261	0.3312 ± 0.0860
Graph-based	ESim [27]	0.2925 ± 0.0223	0.4376 ± 0.0323	0.5480 ± 0.0109	0.5320 ± 0.0246	0.3512 ± 0.0226
	Metapath2vec [5]	0.3956 ± 0.0141	0.4444 ± 0.0231	0.5772 ± 0.0594	0.5256 ± 0.0335	0.3516 ± 0.0407
	HIN2vec [6]	0.2564 ± 0.0131	0.3614 ± 0.0234	0.5218 ± 0.0466	0.4987 ± 0.0252	0.2944 ± 0.0614
	TextGCN [39]	0.4759 ± 0.0126	0.6353 ± 0.0059	-	-	0.3361 ± 0.0032
	METACAT	0.5258 ± 0.0090	0.6889 ± 0.0128	0.7243 ± 0.0336	0.6422 ± 0.0058	0.3971 ± 0.0169

Table 3: Macro F1 scores of compared algorithms on the five datasets. “-”: excessive memory requirements.

Type	Method	GitHub-Bio	GitHub-AI	GitHub-Sec	Amazon	Twitter
Text-based	CNN [12]	0.1896 ± 0.0133	0.1796 ± 0.0216	0.4268 ± 0.0584	0.5056 ± 0.0376	0.2858 ± 0.0559
	HAN [38]	0.0677 ± 0.0208	0.0961 ± 0.0254	0.4095 ± 0.0590	0.4644 ± 0.0597	0.2592 ± 0.0826
	PTE [32]	0.2630 ± 0.0371	0.3363 ± 0.0250	0.3803 ± 0.0218	0.2563 ± 0.0810	0.1739 ± 0.0190
	WeSTClass [23]	0.3414 ± 0.0129	0.4056 ± 0.0248	0.5497 ± 0.0054	0.5234 ± 0.0147	0.3085 ± 0.0398
	PCEM [36]	0.2977 ± 0.0281	0.3751 ± 0.0350	0.4033 ± 0.0336	0.4239 ± 0.0237	0.2039 ± 0.0472
	BERT [4]	0.1740 ± 0.0164	0.2083 ± 0.0415	0.4956 ± 0.0164	0.4911 ± 0.0544	0.2834 ± 0.0550
Graph-based	ESim [27]	0.2598 ± 0.0182	0.3209 ± 0.0202	0.4672 ± 0.0171	0.5336 ± 0.0220	0.3399 ± 0.0113
	Metapath2vec [5]	0.3214 ± 0.0128	0.3220 ± 0.0290	0.5140 ± 0.0637	0.5239 ± 0.0437	0.3443 ± 0.0208
	HIN2vec [6]	0.2742 ± 0.0136	0.2513 ± 0.0211	0.4000 ± 0.0115	0.4261 ± 0.0284	0.2411 ± 0.0142
	TextGCN [39]	0.4817 ± 0.0078	0.5997 ± 0.0013	-	-	0.3191 ± 0.0029
	METACAT	0.5230 ± 0.0080	0.6154 ± 0.0079	0.6323 ± 0.0235	0.6496 ± 0.0091	0.3612 ± 0.0067

HIMECat: Jointly Modeling Metadata and Hierarchy

- How to jointly leverage the label hierarchy, metadata, and text information?

Zhang, Y., Chen, X., Meng, Y., & Han, J. "Hierarchical Metadata-Aware Document Categorization under Weak Supervision", WSDM'21. A few labeled documents as supervision.

The screenshot shows the GitHub repository for 'huggingface/transformers'. Annotations include:

- Label Hierarchy:** A grid of categories like 'Computer Vision' (Semantic Segmentation, Image Classification, Object Detection, Image Generation) and 'Natural Language Processing' (Machine Translation, Language Modeling).
- User (Metadata):** The repository name 'huggingface/transformers' and the README file.
- Description (Text):** The repository's description: 'State-of-the-art Natural Language Processing for PyTorch and TensorFlow 2.0'.
- Tags (Metadata):** A list of tags including 'natural-language-processing', 'pytorch', 'tensorflow', etc.

(a) **GitHub Repository.** Label Hierarchy: PaperWithCode Task Taxonomy (<https://paperswithcode.com/sota>); Text: Description and README; Metadata: User and Tag.

The screenshot shows the arXiv abstract for the paper 'Language Models are Few-Shot Learners'. Annotations include:

- Label Hierarchy:** A sidebar showing the category path: Computer Science > Artificial Intelligence > cs.AI.
- Title (Text):** The paper title 'Language Models are Few-Shot Learners'.
- Authors (Metadata):** The list of authors: Tom B. Brown, Benjamin Mann, Nick Ryder, etc.
- Abstract (Text):** The first paragraph of the abstract: 'Recent work has demonstrated substantial gains on many NLP tasks and benchmarks by pre-training on a large corpus of text followed by fine-tuning on a specific task.'

(b) **arXiv Paper.** Label Hierarchy: arXiv Category Taxonomy (https://arxiv.org/category_taxonomy); Text: Title and Abstract; Metadata: Author.

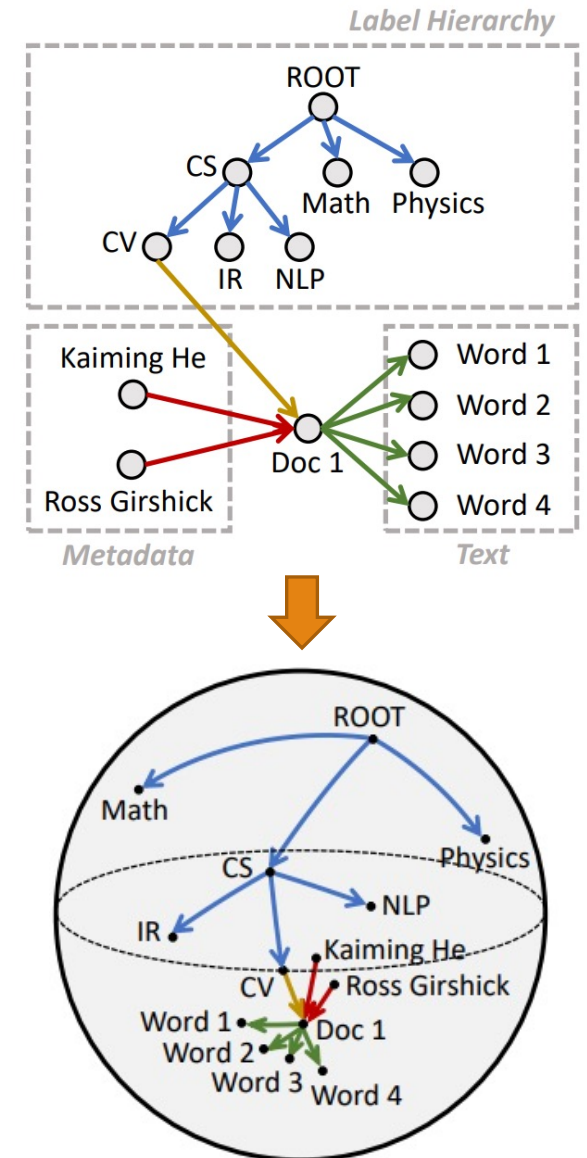
The screenshot shows an Amazon product review for 'iCrave Coffee Flavored Coffee Variety Pack'. Annotations include:

- Label Hierarchy:** A list of departments: Grocery & Gourmet Food, Coffee Beverages, Single-Serve Coffee Capsules & Pods, etc.
- User (Metadata):** The reviewer's name 'LJsquared'.
- Title (Text):** The review title 'Great flavor, strong coffee!'.
- Review (Text):** The review content: 'Deeeelicious! I was skeptical because I'm a diehard Starbucks fan and prefer my medium roast to be strong.'

(c) **Amazon Review.** Label Hierarchy: Amazon Product Catalog [24]; Text: Title and Review; Metadata: User and Product.

HIMECat: A Hierarchical Generative Process

- ❑ Step 1: Parent Label -> Child Label
- ❑ Step 2: **Leaf** label & Metadata -> Document
- ❑ Step 3: Document -> Word
- ❑ **Joint Representation Learning**
 - ❑ Embeddings are the parameters of the generative process.
 - ❑ Maximum likelihood estimation of the parameters when observing the hierarchy, metadata and text
- ❑ **Hierarchical Data Augmentation**
 - ❑ After representation learning, how to synthesize training data for each class?
 - ❑ Follow the generative process



HIMECat: Experimental Results

□ Datasets

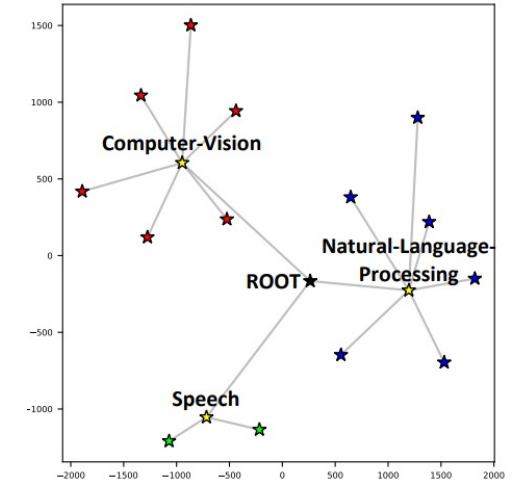
- GitHub: 3+14 categories; 1,596 docs
- ArXiv: 5+88 categories; 25,960 docs
- Amazon: 18+147 categories; 147,000 docs

□ Metrics

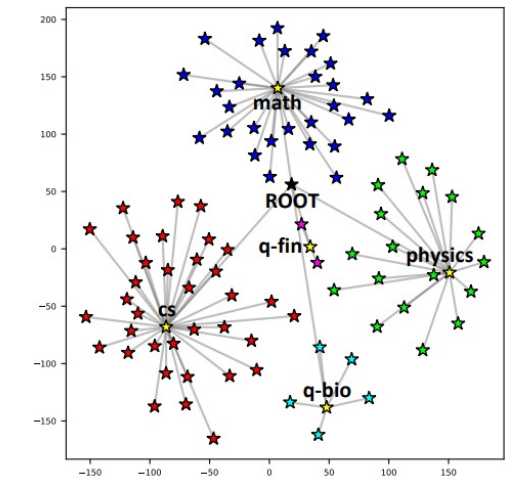
- F1 scores on leaf categories
- F1 scores on all non-root categories

Table 2: {Leaf, Overall}×{Micro, Macro} F1 scores of compared algorithms on the three datasets. *: significantly worse than HiMECAT (p-value < 0.05). **: significantly worse than HiMECAT (p-value < 0.01).

	GitHub				ArXiv				Amazon			
	Leaf		Overall		Leaf		Overall		Leaf		Overall	
	Micro	Macro	Micro	Macro	Micro	Macro	Micro	Macro	Micro	Macro	Micro	Macro
HierSVM [7]	0.1861**	0.1388**	0.4862**	0.2457**	0.0538**	0.0460**	0.4066**	0.0750**	0.0248**	0.0217**	0.2218**	0.0494**
WeSHClass [29]	0.1727**	0.1559**	0.3332**	0.1924**	0.0604**	0.0602**	0.3077**	0.0797**	0.0483**	0.0500**	0.1234**	0.0640**
PCEM [48]	0.2519**	0.1234**	0.5299*	0.1786**	0.1090**	0.0717**	0.4440	0.0963**	0.0675**	0.0439**	0.2189**	0.0659**
HiGitClass [53]	0.3984	0.3902*	0.5073**	0.4084**	0.1738**	0.1656**	0.3928**	0.1880**	0.0903**	0.0876**	0.1677**	0.1040**
MetaCat [51]	0.3762**	0.3403**	0.5411*	0.3863**	0.0790**	0.0768**	0.3071**	0.0935**	0.1008**	0.0994**	0.1703**	0.1083**
Metapath2vec [6]	0.2814**	0.2805**	0.4592**	0.3212**	0.1360**	0.1344**	0.3419**	0.1534**	0.0669**	0.0666**	0.1334**	0.0800**
Poincaré [32]	0.2750**	0.1980**	0.4302**	0.2218**	0.1336**	0.1296**	0.2995**	0.1454**	0.0645**	0.0607**	0.1202**	0.0739**
BERT [5]	0.2889**	0.2561**	0.4675**	0.3007**	0.1316**	0.0812**	0.4203**	0.1100**	0.0891**	0.0520**	0.2361**	0.0771**
HiMECAT	0.4254	0.4209	0.5820	0.4535	0.2038	0.1938	0.4509	0.2191	0.1552	0.1553	0.2748	0.1770




(a) GitHub



(b) ArXiv

Outline

- ❑ What Weakly-Supervised Text Classification Is, and Why It Matters
- ❑ Flat Text Classification
 - ❑ Embedding: WeSTClass [CIKM'18]
 - ❑ Pre-trained LM: ConWea [ACL'20], LOTClass [EMNLP'20], X-Class [NAACL'21]
- ❑ Text Classification with Taxonomy Information
 - ❑ Embedding: WeSHClass [AAAI'19]
 - ❑ Pre-trained LM: TaxoClass [NAACL'21]
- ❑ Text Classification with Metadata Information
 - ❑ Embedding: MetaCat [SIGIR'20]
 - ❑ Pre-trained LM: MICoL [WWW'22] 

MICoL: Metadata-Induced Contrastive Learning for Zero-Shot Multi-Label Text Classification

Input

- A set of labels. Each label has its name and description.
- A large set of unlabeled documents associated with metadata (e.g., authors, venue, references) that can connect the documents together.

Output

- A multi-label text classifier. Given some new documents, the classifier can predict relevant labels for each document.

Webgraph Label Name

105 Publications 64,901 Citations*

Definition

The webgraph describes the directed links between pages of the World Wide Web. A graph, in general, consists of several vertices, some pairs connected by edges. In a directed graph, edges are directed lines or arcs. The webgraph is a directed graph, whose vertices correspond to the pages of the WWW, and a directed edge connects page X to page Y if there exists a hyperlink on page X, referring to page Y.

Label Description

(a) Label “Webgraph” from Microsoft Academic (<https://academic.microsoft.com/topic/2777569578/>).

Betacoronavirus MeSH Descriptor Data 2021

Label Name MeSH Tree Structures Concepts

MeSH Heading Betacoronavirus

Tree Number(s) B04.820.578.500.540.150.113

Unique ID D000073640

RDF Unique Identifier <http://id.nlm.nih.gov/mesh/D000073640>

Annotation infection: coordinate with CORONAVIRUS INFECTIONS

Scope Note A genus of the family CORONAVIRIDAE which causes respiratory or gastrointestinal disease in a variety of mostly mammals. Human betacoronaviruses include HUMAN ENTERIC CORONAVIRUS; HUMAN CORONAVIRUS OC43; MERS VIRUS; and SARS VIRUS. Members have either core transcription regulatory sequences of 5'-CUAAAC-3' or 5'-CUAAAC-3' and mostly have no ORF downstream to the N protein gene.

Entry Term(s) HCoV-HKU1
Human coronavirus HKU1
Pipistrellus bat coronavirus HKU5
Rousettus bat coronavirus HKU9
Tylonycteris bat coronavirus HKU4

Label Description

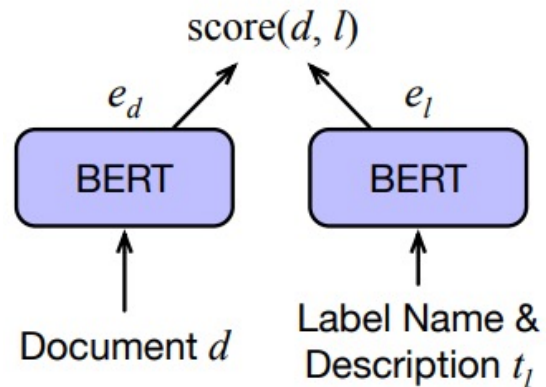
Synonyms (also viewed as Label Names)

(b) Label “Betacoronavirus” from PubMed (<https://meshb.nlm.nih.gov/record/ui?ui=D000073640>).

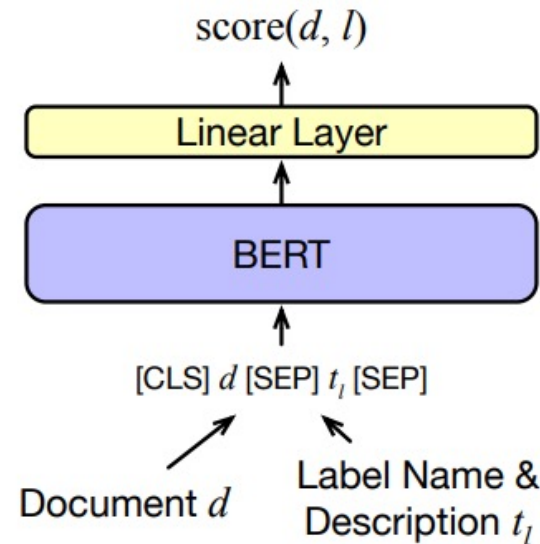
Zhang, Y., Shen, Z., Wu, C., Xie, B., Wang, Y., Wang, K. & Han, J. "Metadata-Induced Contrastive Learning for Zero-Shot Multi-Label Text Classification", To appear in WWW'22. **Category names and descriptions as supervision.**

Pre-trained Language Models for Multi-Label Text Classification

- If we could have some labeled documents, ...
 - We can use relevant (document, label) pairs to fine-tune the pre-trained LM.
 - Both Bi-Encoder and Cross-Encoder are applicable.



(a) Bi-Encoder

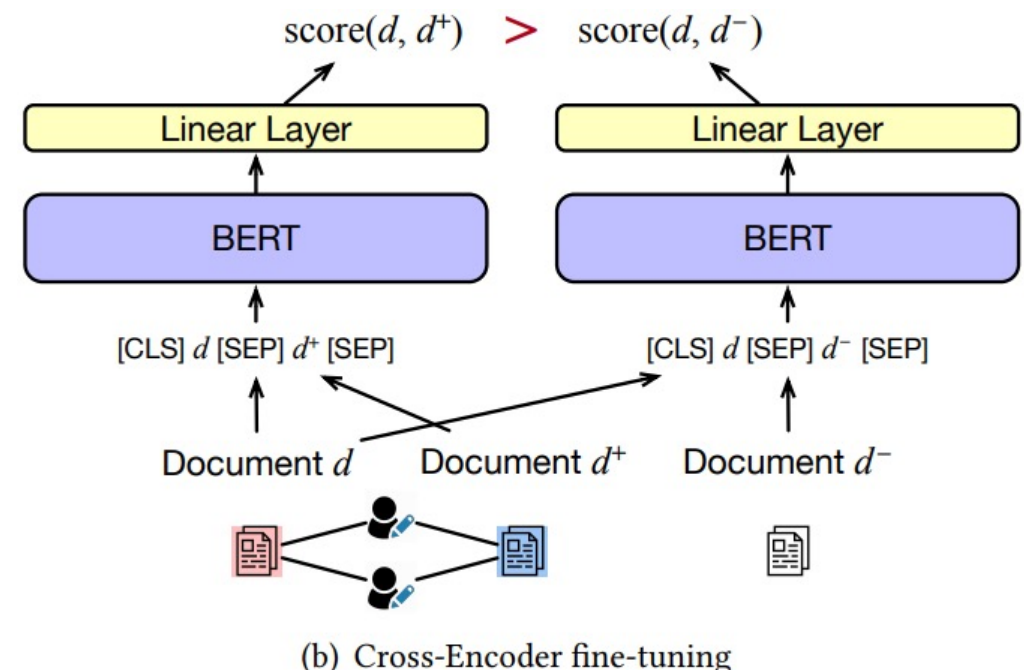
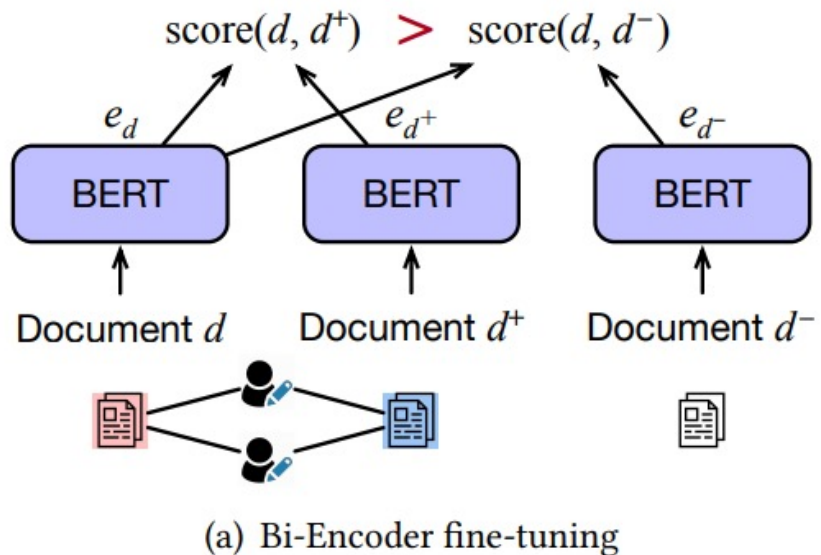
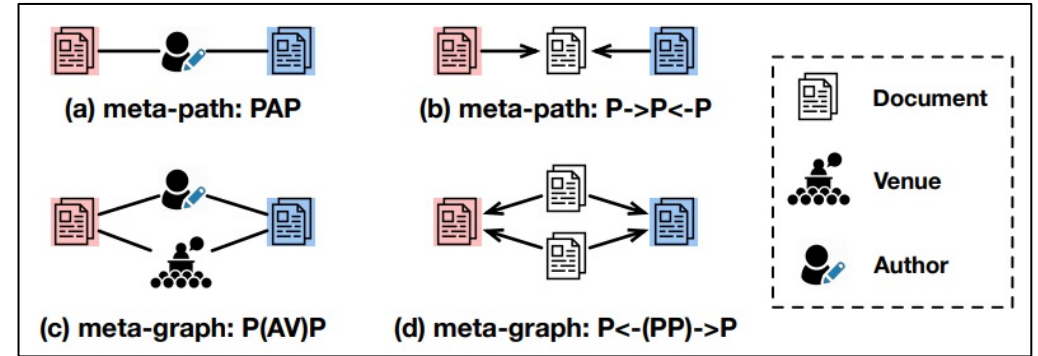


(b) Cross-Encoder

- However, we do not have any labeled documents!!!

Metadata-Induced Contrastive Learning

- Contrastive learning: Instead of training the model to know “what is what” (e.g., relevant (document, label) pairs), train it to know “what is similar with what” (e.g., similar (document, document) pairs).
- Using metadata to define similar (document, document) pairs.



MICoL: Experimental Results

- MICoL significantly outperforms text-based contrastive learning baselines.
- MICoL is competitive with the supervised SOTA trained on 10K–50K labeled documents.

	Algorithm	MAG-CS [49]					PubMed [24]				
		P@1	P@3	P@5	NDCG@3	NDCG@5	P@1	P@3	P@5	NDCG@3	NDCG@5
Zero-shot	Doc2Vec [31]	0.5697**	0.4613**	0.3814**	0.5043**	0.4719**	0.3888**	0.3283**	0.2859**	0.3463**	0.3252**
	SciBERT [2]	0.6440**	0.5030**	0.4011**	0.5545**	0.5061**	0.4427**	0.3572**	0.3031**	0.3809**	0.3510**
	ZeroShot-Entail [61]	0.6649**	0.5003**	0.3959**	0.5570**	0.5057**	0.5275**	0.4021	0.3299	0.4352	0.3913
	SPECTER [8]	0.7107**	0.5381**	0.4184**	0.5979**	0.5365**	0.5286**	0.3923**	0.3181**	0.4273**	0.3815**
	EDA [53]	0.6442**	0.4939**	0.3948**	0.5471**	0.5000**	0.4919	0.3754*	0.3101*	0.4058*	0.3667*
	UDA [57]	0.6291**	0.4848**	0.3897**	0.5362**	0.4918**	0.4795**	0.3696**	0.3067**	0.3986**	0.3614**
	MICoL (Bi-Encoder, $P \rightarrow P \leftarrow P$)	0.7062*	0.5369*	0.4184*	0.5960*	0.5355*	0.5124**	0.3869*	0.3172*	0.4196*	0.3774*
	MICoL (Bi-Encoder, $P \leftarrow (PP) \rightarrow P$)	0.7050*	0.5344*	0.4161*	0.5937*	0.5331*	0.5198**	0.3876*	0.3172*	0.4215*	0.3786*
	MICoL (Cross-Encoder, $P \rightarrow P \leftarrow P$)	0.7177	0.5444	0.4219	0.6048	0.5415	0.5412	0.4036	0.3257	0.4391	0.3906
	MICoL (Cross-Encoder, $P \leftarrow (PP) \rightarrow P$)	0.7061	0.5376	0.4187	0.5964	0.5357	0.5218	0.3911	0.3172*	0.4249	0.3794
Supervised	MATCH [68] (10K Training)	0.4423**	0.2851**	0.2152**	0.3375**	0.3003**	0.6915	0.3869*	0.2785**	0.4649	0.3896
	MATCH [68] (50K Training)	0.6215**	0.4280**	0.3269**	0.4987**	0.4489**	0.7701	0.4716	0.3585	0.5497	0.4750
	MATCH [68] (100K Training)	0.8321	0.6520	0.5142	0.7342	0.6761	0.8286	0.5680	0.4410	0.6405	0.5626
	MATCH [68] (Full, 560K+ Training)	0.9114	0.7634	0.6312	0.8486	0.8076	0.9151	0.7425	0.6104	0.8001	0.7310

MICoL: Experimental Results

- Most meta-paths and meta-graphs used in MICoL can improve the classification performance upon un-fine-tuned SciBERT.

Algorithm	MAG-CS [49]					PubMed [24]				
	P@1	P@3	P@5	NDCG@3	NDCG@5	P@1	P@3	P@5	NDCG@3	NDCG@5
Un-fine-tuned SciBERT	0.6599**	0.5117**	0.4056**	0.5651**	0.5136**	0.4371**	0.3544**	0.3014**	0.3775**	0.3485**
MICoL (Bi-Encoder, PAP)	0.6877**	0.5285**	0.4143**	0.5852**	0.5280**	0.4974**	0.3818**	0.3154*	0.4122**	0.3727**
MICoL (Bi-Encoder, PVP)	0.6589**	0.5123**	0.4063**	0.5656**	0.5145**	0.4440**	0.3507**	0.2966**	0.3761**	0.3458**
MICoL (Bi-Encoder, $P \rightarrow P$)	0.7094	0.5391	0.4190	0.5982	0.5367	0.5200*	0.3903*	0.3195	0.4240*	0.3808*
MICoL (Bi-Encoder, $P \leftarrow P$)	0.7095*	0.5374*	0.4178*	0.5970*	0.5356*	0.5195**	0.3905*	0.3192	0.4240*	0.3806*
MICoL (Bi-Encoder, $P \rightarrow P \leftarrow P$)	0.7062*	0.5369*	0.4184*	0.5960*	0.5355*	0.5124**	0.3869*	0.3172*	0.4196*	0.3774*
MICoL (Bi-Encoder, $P \leftarrow P \rightarrow P$)	0.7039*	0.5379*	0.4187*	0.5963*	0.5356*	0.5174**	0.3886*	0.3187*	0.4220*	0.3795*
MICoL (Bi-Encoder, $P(AA)P$)	0.6873**	0.5272**	0.4130**	0.5840**	0.5269**	0.4963**	0.3794**	0.3139**	0.4101**	0.3711**
MICoL (Bi-Encoder, $P(AV)P$)	0.6832**	0.5263**	0.4135**	0.5823**	0.5263**	0.4894**	0.3743**	0.3099**	0.4045**	0.3664**
MICoL (Bi-Encoder, $P \rightarrow (PP) \leftarrow P$)	0.7015**	0.5334**	0.4160**	0.5920**	0.5322**	0.5163**	0.3879*	0.3172*	0.4211*	0.3781*
MICoL (Bi-Encoder, $P \leftarrow (PP) \rightarrow P$)	0.7050*	0.5344*	0.4161*	0.5937*	0.5331*	0.5198**	0.3876*	0.3172*	0.4215*	0.3786*
MICoL (Cross-Encoder, PAP)	0.7034*	0.5355	0.4168	0.5943	0.5337	0.5212**	0.3921*	0.3207	0.4255*	0.3818*
MICoL (Cross-Encoder, PVP)	0.6720*	0.5203*	0.4103*	0.5750*	0.5210*	0.4668**	0.3633**	0.3051**	0.3908**	0.3574**
MICoL (Cross-Encoder, $P \rightarrow P$)	0.7033*	0.5391	0.4201	0.5971*	0.5365*	0.5266	0.3946	0.3207	0.4286	0.3830
MICoL (Cross-Encoder, $P \leftarrow P$)	0.7169	0.5430	0.4214	0.6033	0.5406	0.5265	0.3924	0.3186	0.4268	0.3811
MICoL (Cross-Encoder, $P \rightarrow P \leftarrow P$)	0.7177	0.5444	0.4219	0.6048	0.5415	0.5412	0.4036	0.3257	0.4391	0.3906
MICoL (Cross-Encoder, $P \leftarrow P \rightarrow P$)	0.7045	0.5356*	0.4168*	0.5944*	0.5336*	0.5243*	0.3932*	0.3190*	0.4271*	0.3814*
MICoL (Cross-Encoder, $P(AA)P$)	0.7028	0.5351	0.4171	0.5939	0.5338	0.5290*	0.3937	0.3201	0.4285*	0.3830
MICoL (Cross-Encoder, $P(AV)P$)	0.7024*	0.5354*	0.4177	0.5940*	0.5343*	0.5164**	0.3897*	0.3195*	0.4225*	0.3797*
MICoL (Cross-Encoder, $P \rightarrow (PP) \leftarrow P$)	0.7076*	0.5379*	0.4188	0.5971*	0.5363*	0.5186	0.3924*	0.3184*	0.4254*	0.3800*
MICoL (Cross-Encoder, $P \leftarrow (PP) \rightarrow P$)	0.7061	0.5376	0.4187	0.5964	0.5357	0.5218	0.3911	0.3172*	0.4249	0.3794

Summary

Method	Flat vs. Hierarchical	Single-label vs. Multi-label	Supervision Format	Embedding vs. Pretrained LM
WeSTClass	Flat	Single-label	Both types	Embedding
ConWea	Flat	Single-label	Category Names	Pretrained LM
LOTClass	Flat	Single-label	Category Names	Pretrained LM
X-Class	Flat & Hierarchical	Single-label & Path	Category Names	Pretrained LM
WeSHClass	Hierarchical	Path	Both types	Embedding
TaxoClass	Hierarchical	Multi-label	Category Names	Pretrained LM
MetaCat	Flat	Single-label	A Few Labeled Docs	Embedding
HIMECat	Hierarchical	Path	A Few Labeled Docs	Embedding
MICoL	Flat	Multi-label	Category Names	Pretrained LM

References

- ❑ Meng, Y., Shen, J., Zhang, C., & Han, J. “Weakly-supervised neural text classification”, CIKM’18
- ❑ Mekala, D. & Shang, J. “Contextualized Weak Supervision for Text Classification”, ACL’20
- ❑ Meng, Y., Zhang, Y., Huang, J., Xiong, C., Ji, H., Zhang, C., & Han, J. “Text Classification Using Label Names Only: A Language Model Self-Training Approach”, EMNLP’20
- ❑ Wang, Z., Mekala, D., & Shang, J. “X-Class: Text Classification with Extremely Weak Supervision”, NAACL’21
- ❑ Meng, Y., Shen, J., Zhang, C., & Han, J. “Weakly-Supervised Hierarchical Text Classification”, AACL’19
- ❑ Shen, J., Qiu, W., Meng, Y., Shang, J., Ren, X., & Han, J., “TaxoClass: Hierarchical Multi-Label Text Classification Using Only Class Names”, NAACL’21
- ❑ Zhang, Y., Meng, Y., Huang, J., Xu, F.F., Wang, X., & Han, J. “Minimally Supervised Categorization of Text with Metadata”, SIGIR’20
- ❑ Zhang, Y., Chen, X., Meng, Y., & Han, J. “Hierarchical Metadata-Aware Document Categorization under Weak Supervision”, WSDM’21
- ❑ Zhang, Y., Shen, Z., Wu, C., Xie, B., Wang, Y., Wang, K. & Han, J. "Metadata-Induced Contrastive Learning for Zero-Shot Multi-Label Text Classification", To appear in WWW’22



Q&A

