# Part III: Mining Document Structures: Weakly-Supervised Text Classification

**EDBT 2023 Tutorial: Mining Structures from Massive Texts by Exploring the Power of Pre-trained Language Models**

**Yu Zhang, Yunyi Zhang, Jiawei Han**

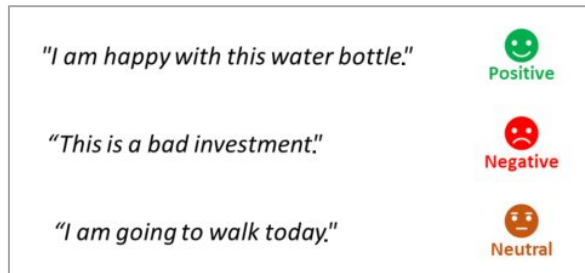**Department of Computer Science, University of Illinois at Urbana-Champaign**
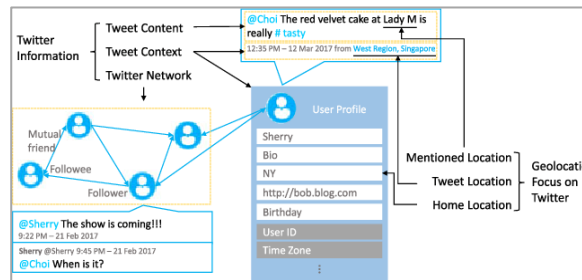
**Mar 29, 2023**

# Outline

- ❑ What Weakly-Supervised Text Classification Is, and Why It Matters

- ❑ Flat Text Classification

- ❑ Text Classification with Taxonomy Information

- ❑ Text Classification with Metadata Information

# 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

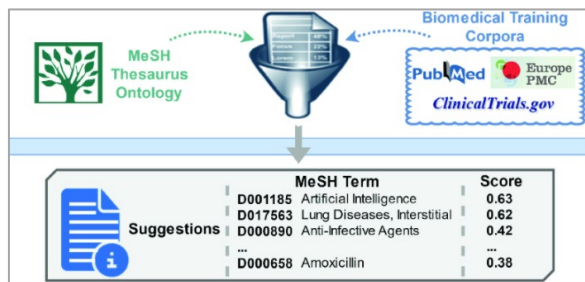❑ Text Classification has a lot of downstream applications



**Sentiment Analysis**



**Location Prediction**



**News Topic Classification**



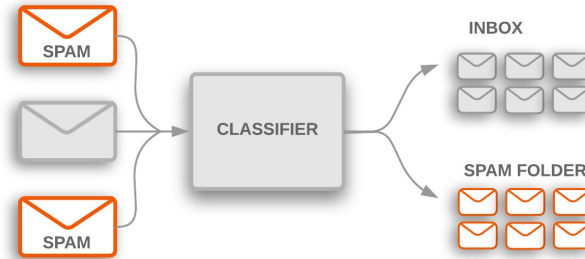**Paper Topic Classification**



**Email Intent Identification**



**Hate Speech Detection**

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

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

    ❑ E.g., Spam Detection



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

    ❑ E.g., Paper Topic Classification



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

# Different Text Classification Settings: Flat vs. Hierarchical

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

   ❑ E.g., 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/B089YFHYYS/

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

   ❑ E.g., 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 <span style="color:darkred">embedding learning</span> or <span style="color:blue">pre-trained language models</span>

❑ Pseudo training data generation

  ❑ Retrieve some unlabeled documents or synthesize some artificial documents using <span style="color:darkred">text embeddings</span> or <span style="color:blue">contextualized representations</span>

  ❑ Give them pseudo labels to train a text classifier

❑ Transfer the knowledge of <span style="color:blue">pre-trained language models</span> to classification tasks

# Outline

❑ What Weakly-Supervised Text Classification Is, and Why It Matters

❑ Flat Text Classification

    ❑ Static Embedding: WeSTClass [CIKM'18]

    ❑ Pre-trained LM: ConWea [ACL'20], LOTClass [EMNLP'20], X-Class [NAACL'21],

                      Prompt-based Classifier

❑ Text Classification with Taxonomy Information

❑ Text Classification with Metadata Information

# 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-Mishes 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: Experiment Results

**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:**

| Methods | The New York Times | | | AG's News | | | Yelp Review | | |
|---|---|---|---|---|---|---|---|---|---|
| | LABELS | KEYWORDS | DOCS | LABELS | KEYWORDS | DOCS | LABELS | KEYWORDS | DOCS |
| 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)** |

# Outline

❑ What Weakly-Supervised Text Classification Is, and Why It Matters

❑ Flat Text Classification

    ❑ Static Embedding: WeSTClass [CIKM'18]

    ❑ Pre-trained LM: ConWea [ACL'20], LOTClass [EMNLP'20], X-Class [NAACL'21],

        Prompt-based Classifier

❑ Text Classification with Taxonomy Information

❑ Text Classification with Metadata Information

# 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

# 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"

**User-Provided Seed Words**

| Class | Seed Words |
|-------|------------|
| Soccer | soccer, goal, penalty |
| Law | law, judge, court |
| … | … |

**Raw Docs**

Messi scored the penalty! …
Judge passed the order of …
The court issued a penalty …
            ……

**Extended Seed Words**

| Class | Seed Words |
|-------|------------|
| Soccer | soccer, goal$0, goal$1, penalty$0, penalty$1, |
| Law | law, judge, court$0, court$1 |
| … | … |

**Contextualized Docs**

Messi scored the **penalty$1**! …
Judge passed the order of …
The **court$1** issued a **penalty$0** …
            ……

**Contextualized & Expanded Seed Words**

| Class | Seed Words |
|-------|------------|
| Soccer | soccer, goal$0, penalty$1, … |
| Law | law, judge, court$1, penalty$0, … |
| … | … |

**Text Classifier**

**Contextualized Docs with Predictions**

Messi scored the **penalty$1**! …
Judge passed the order of …
The **court$1** issued a **penalty$0** …
            ……

**Comparative Ranking**

Law ⟷ Soccer

Cosmos ⟷ Politics

16

❑ 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-$F_1$ |
| 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** |
| 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 |
| HAN-Supervised | 0.96 | 0.92 | 0.94 | 0.82 | 0.90 | 0.88 | 0.83 | 0.83 |

Baselines

Ablations

Upper bound

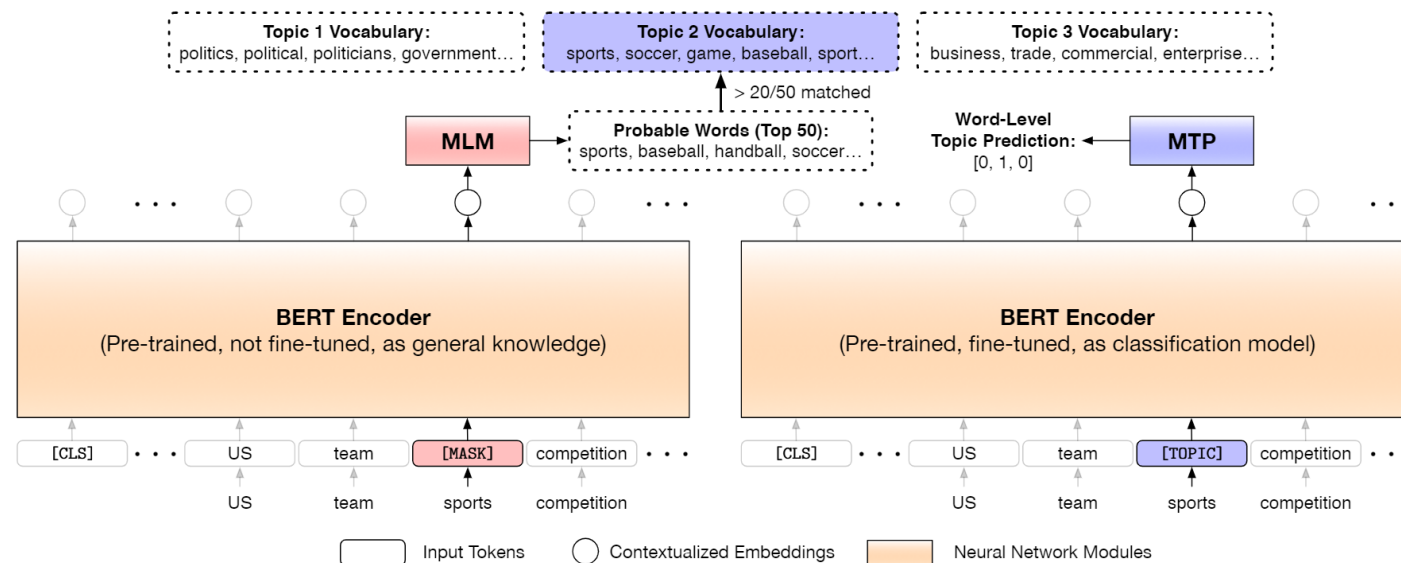# LOTClass: Find Similar Meaning Words with Label Names

❑ 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.**

❑  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
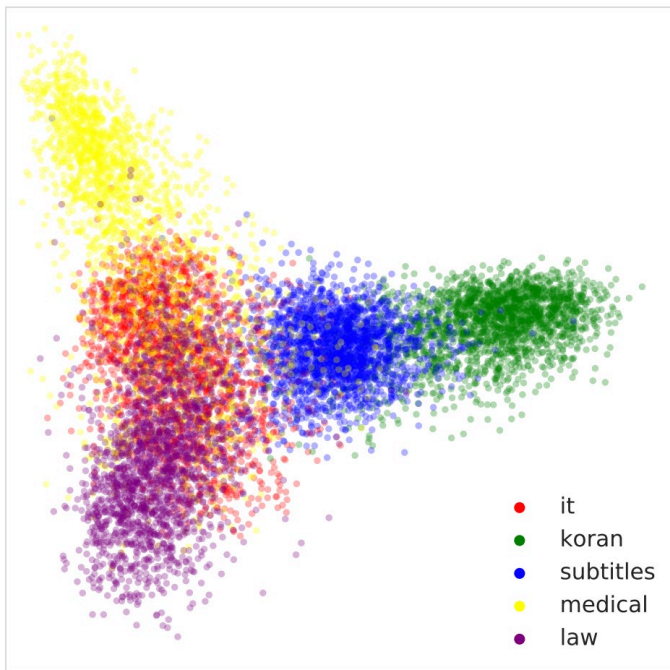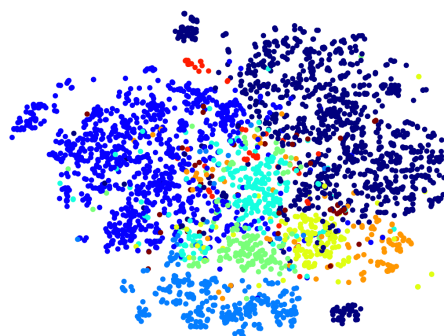


Figure 1: A 2D visualization of average-pooled BERT hidden-state sentence representations using PCA. The colors represent the domain for each sentence.



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

Aharoni, R., & Goldberg, Y. "Unsupervised domain clusters in pretrained language models." ACL'20.

# 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.**

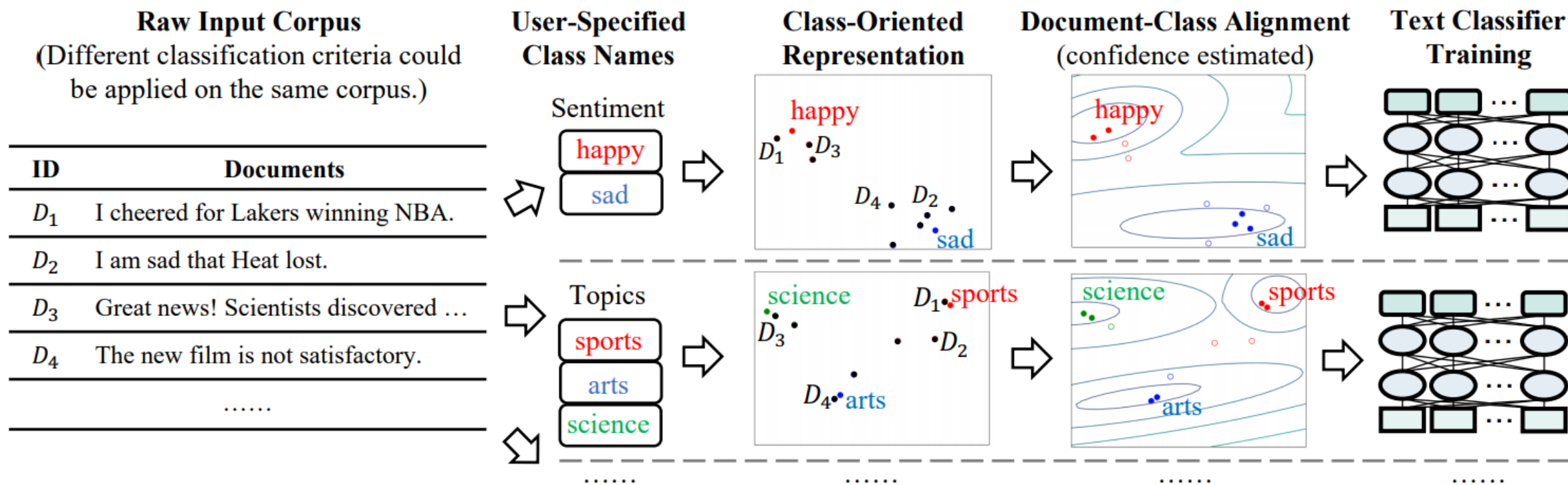# X-Class: Class-Oriented BERT Representations

❑ Clustering for classification based on class-oriented representations
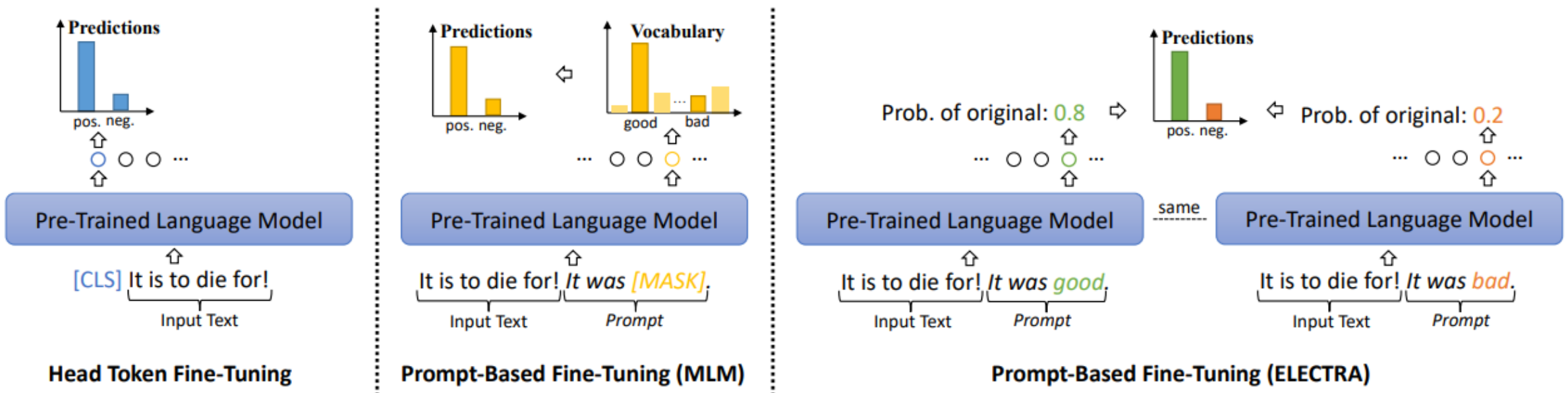
# X-Class: Experiment Results

❑ 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 |

# Prompt-based Fine-tuning for Text Classification

❑ **Head token fine-tuning** randomly initializes a linear classification head and directly predicts class distribution using the [CLS] token, which needs a substantial amount of training data.

❑ **Prompt-based fine-tuning for MLM-based PLM** converts the document into the masked token prediction problem by reusing the pre-trained MLM head.

❑ **Prompt-based fine-tuning for ELECTRA-style PLM** converts documents into the replaced token detection problem by reusing the pre-trained discriminative head.

# Integrating Head Token & Prompt-based Fine-tuning

❑ Why do we need prompts to get pseudo training data?

❑ Simple keyword matching may induce errors.

❑ E.g., "*die*" is a negative word, but a food review "It is to *die* for!" implies a strong positive sentiment.



**Two fine-tuning strategies for pre-trained language model**

**Head Token Fine-Tuning**

Positive sentiment
⇧
○ ○ ○ ...
⇧
Pre-Trained Language Model
⇧
[CLS] It is to die for!
Input Text

**Prompt-Based Fine-Tuning**

0.8 (original)
⇧
... ○ ○ ○ ...
⇧
Pre-Trained Language Model
⇧
It is to die for! *It was good.*
Input Text  *Prompt*



Initial Pseudo Labels $P^0$

Zero-Shot Prompting → Head Token Fine-Tuning $P_0^i$ → Sampling → Prompt-Based Fine-Tuning $P_1^i$ → Intersection → Updated Pseudo Labels $P^i$

... Prompt-Based Fine-Tuning $P_r^i$

Unlabeled Corpus

Use updated pseudo labels to repeat the process

**(1) Zero-Shot Prompting for Pseudo Label Acquisition**

**(2) Iterative Classifier Training and Pseudo Label Expansion**

# Experimental Results

❑ Integrating head token and prompt-based fine-tuning for weakly supervised text classification with category names only.
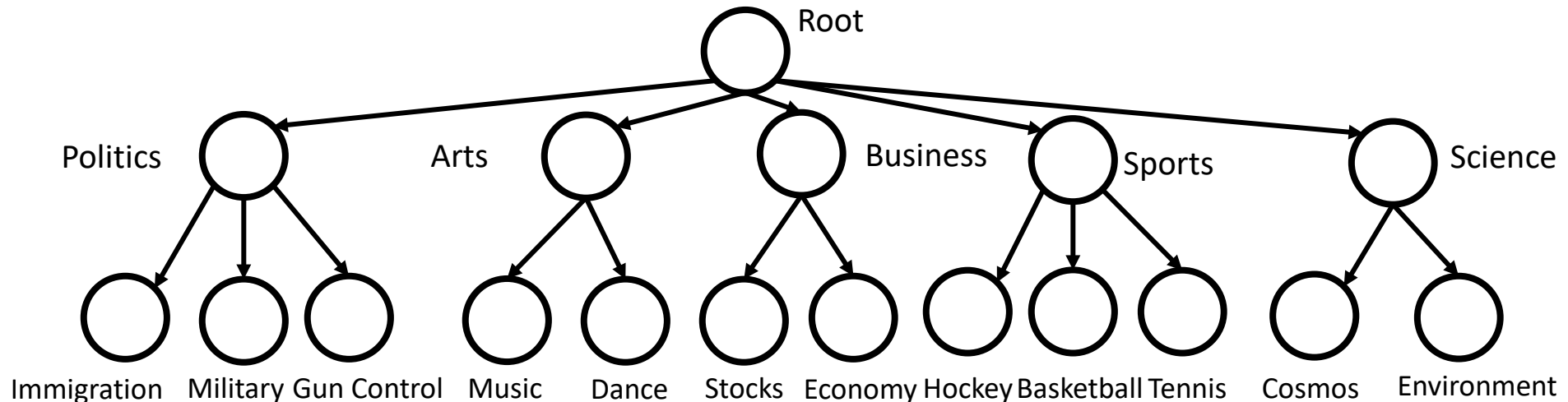
| Methods | AGNews | | 20News | | Yelp | | IMDB | |
|---|---|---|---|---|---|---|---|---|
| | Micro-F1 | Macro-F1 | Micro-F1 | Macro-F1 | Micro-F1 | Macro-F1 | Micro-F1 | Macro-F1 |
| WeSTClass | 0.823 | 0.821 | 0.713 | 0.699 | 0.816 | 0.816 | 0.774 | - |
| ConWea | 0.746 | 0.742 | 0.757 | 0.733 | 0.714 | 0.712 | - | - |
| LOTClass | 0.869 | 0.868 | 0.738 | 0.725 | 0.878 | 0.877 | 0.865 | - |
| XClass | 0.857 | 0.857 | 0.786 | 0.778 | 0.900 | 0.900 | - | - |
| ClassKG[†] | 0.881 | 0.881 | <u>0.811</u> | **0.820** | 0.918 | 0.918 | 0.888 | 0.888 |
| RoBERTa (0-shot) | 0.581 | 0.529 | $0.507^{\ddagger}$ | $0.445^{\ddagger}$ | 0.812 | 0.808 | 0.784 | 0.780 |
| ELECTRA (0-shot) | 0.810 | 0.806 | 0.558 | 0.529 | 0.820 | 0.820 | 0.803 | 0.802 |
| PromptClass | | | | | | | | |
|   ELECTRA+BERT | <u>0.884</u> | <u>0.884</u> | 0.789 | 0.791 | 0.919 | 0.919 | 0.905 | 0.905 |
|   RoBERTa+RoBERTa | **0.895** | **0.895** | $0.755^{\ddagger}$ | $0.760^{\ddagger}$ | <u>0.920</u> | <u>0.920</u> | <u>0.906</u> | <u>0.906</u> |
|   ELECTRA+ELECTRA | <u>0.884</u> | <u>0.884</u> | **0.816** | <u>0.817</u> | **0.957** | **0.957** | **0.931** | **0.931** |
| Fully Supervised | 0.940 | 0.940 | 0.965 | 0.964 | 0.957 | 0.957 | 0.945 | - |

# Outline

❑ What Weakly-Supervised Text Classification Is, and Why It Matters

❑ Flat Text Classification

❑ Text Classification with Taxonomy Information

   ❑ Static Embedding: WeSHClass [AAAI'19]

   ❑ Pre-trained LM: TaxoClass [NAACL'21]

❑ Text Classification with Metadata Information

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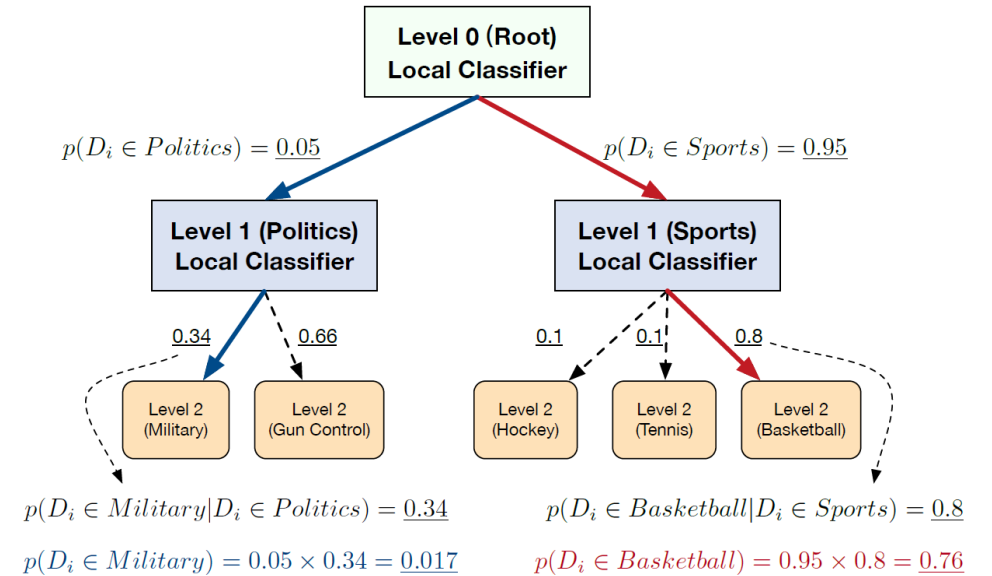**Applicable to both keyword-level and document-level supervision.**

# Hierarchical Classification Model

- Local Classifier Per Node

  - Essentially a flat classification task

  - Follow WeSTClass

- 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 $k$



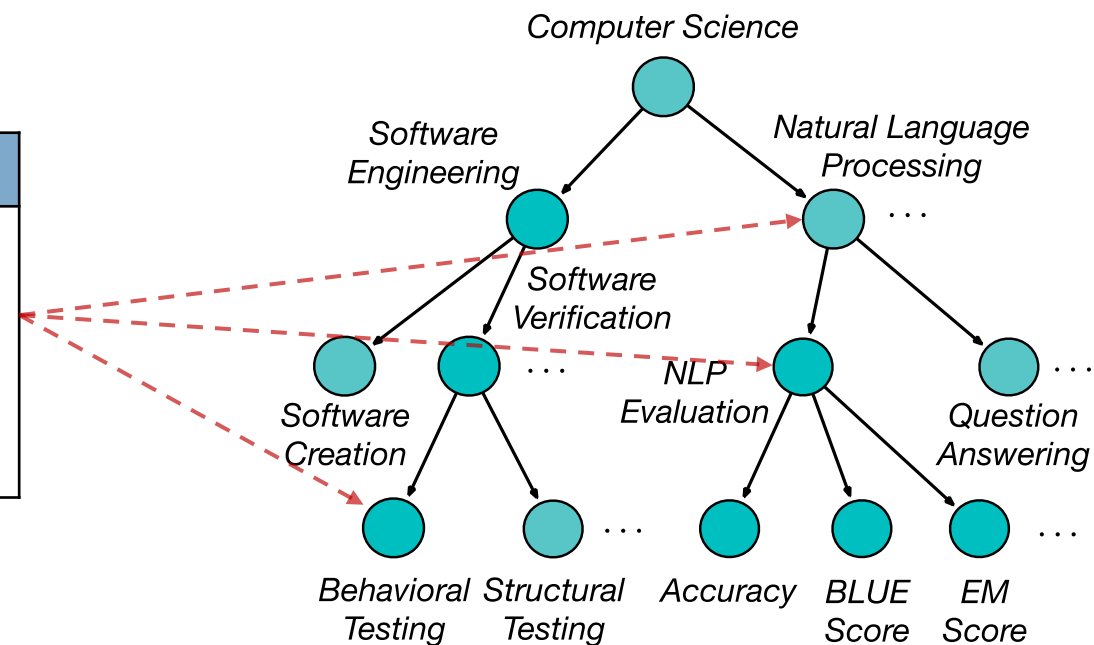| Methods | NYT | | | | arXiv | | | | Yelp Review | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | KEYWORDS | | DOCS | | KEYWORDS | | DOCS | | KEYWORDS | | DOCS | |
| | Macro | Micro | Macro Avg. (Std.) | Micro Avg. (Std.) | Macro | Micro | Macro Avg. (Std.) | Micro Avg. (Std.) | Macro | 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

❑ Text Classification with Taxonomy Information

  ❑ Static Embedding: WeSHClass [AAAI'19]

  ❑ Pre-trained LM: TaxoClass [NAACL'21]

❑ Text Classification with Metadata Information

# 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

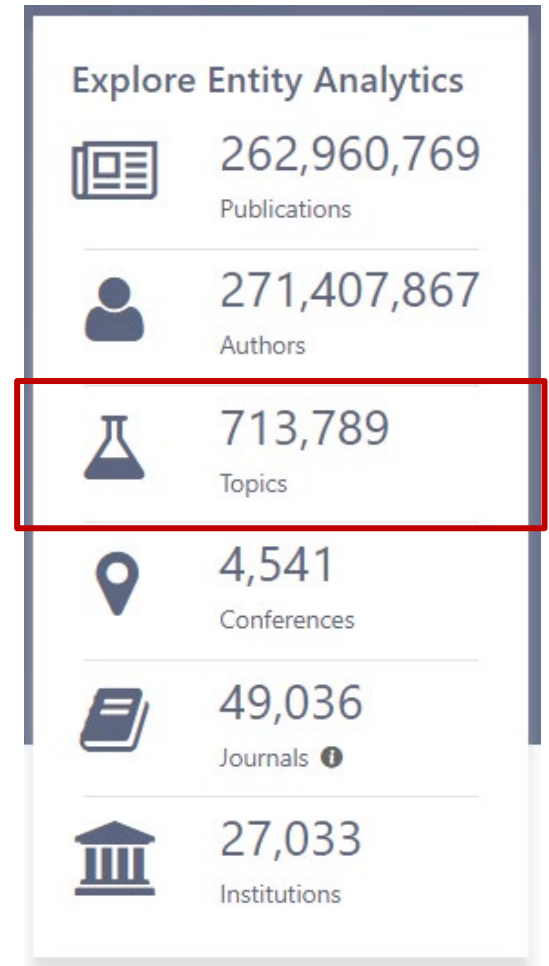❑ Category names can be phrases and may not appear in the corpus



**Document**

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

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

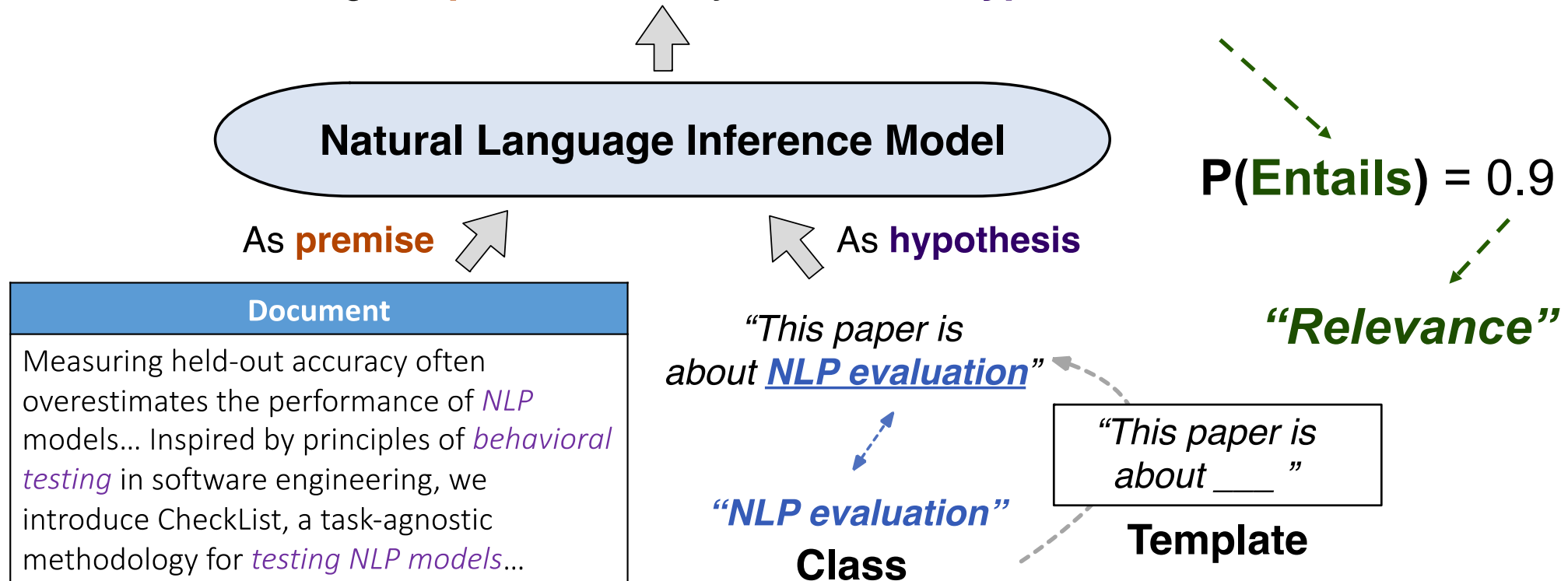❑ Impossible for users to provide even a small set of (e.g., 3) keywords/labeled documents for each category

**Explore Entity Analytics**

262,960,769 Publications

271,407,867 Authors

713,789 Topics

4,541 Conferences

49,036 Journals ❶

27,033 Institutions

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**?

**Natural Language Inference Model**

As **premise**　　　　As **hypothesis**

P(**Entails**) = 0.9

*"Relevance"*

**Document**

Measuring held-out accuracy often overestimates the performance of *NLP models*… Inspired by principles of *behavioral testing* in software engineering, we introduce CheckList, a task-agnostic methodology for *testing NLP models*…

*"This paper is about NLP evaluation"*

*"NLP evaluation"*
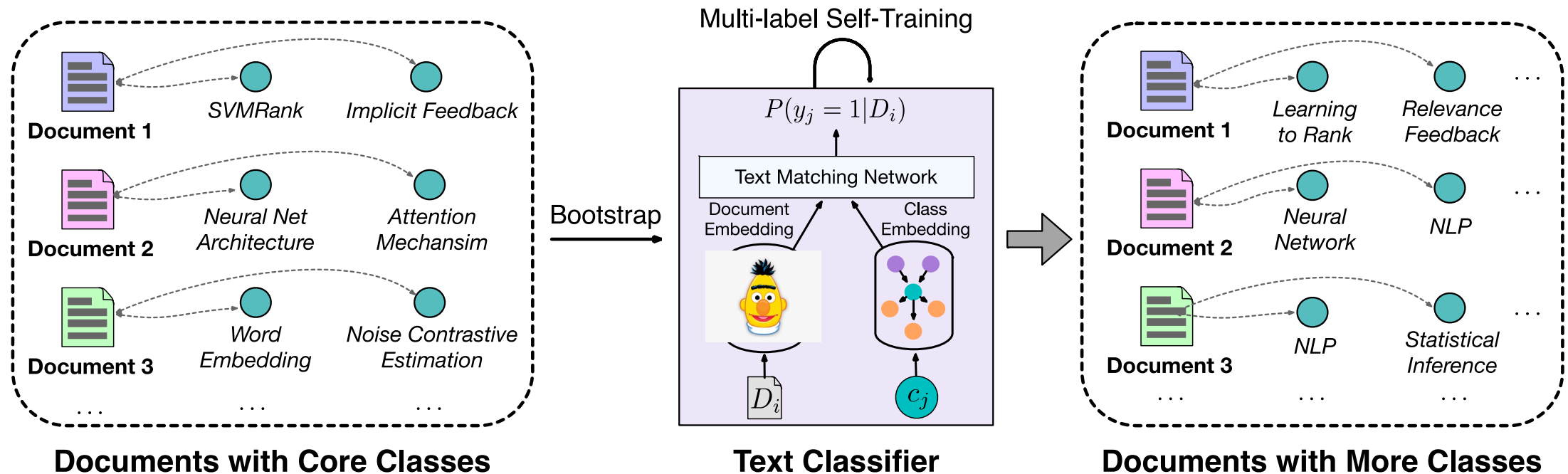**Class**

*"This paper is about ___ "*
**Template**

# TaxoClass: Top-Down Exploration

- How to use the taxonomy?
- Shrink the label search space with top-down exploration
  - Use a relevance model to filter out completely irrelevant classes



**Document-class Relevance**

$$rel(D_i, c_j)$$

**Relevance Model**
(e.g., BM25, doc2vec, BERT-NLI)

**Document** $D_i$    **Candidate Class** $c_j$

**Document**

rel=0.75

Computer Science

Information retreival

Theory

Data Mining

Learning to Rank

Query Expansion

Text Mining

Graph Mining

**Reduced Label Search Space**

35

# 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



**Documents with Core Classes**     **Text Classifier**     **Documents with More Classes**

# 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

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

# Outline

❑ What Weakly-Supervised Text Classification Is, and Why It Matters

❑ Flat Text Classification

❑ Text Classification with Taxonomy Information

❑ Text Classification with Metadata Information

  ❑ Static Embedding: MetaCat [SIGIR'20]

  ❑ Pre-trained LM: MICoL [WWW'22]

# MetaCat: Leveraging Metadata for Classification

❑ Metadata is prevalent in many text sources

    ❑ **GitHub repositories**: User, Tag     ❑ **Amazon reviews:** User, Product

    ❑ **Tweets**: User, Hashtag     ❑ **Scientific papers:** Author, Venue

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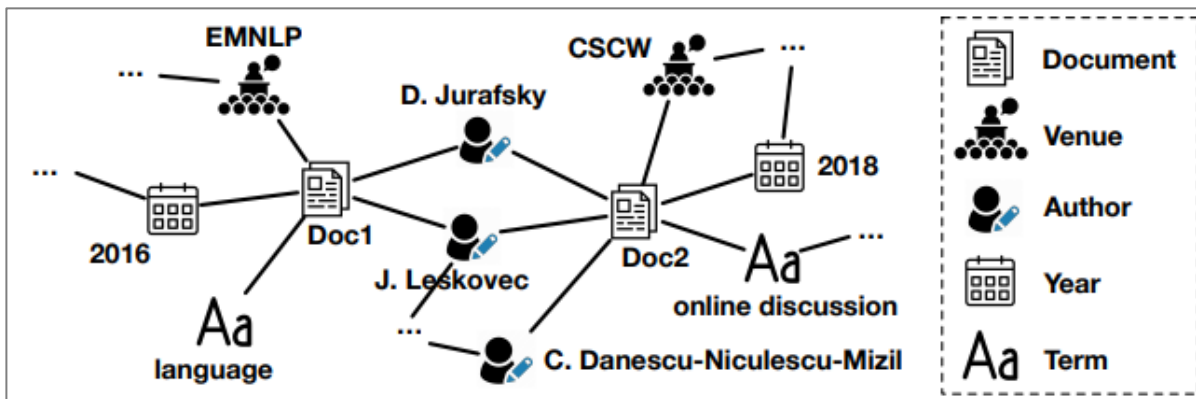


(a) GitHub Repository      (b) Tweet      (c) Amazon Review

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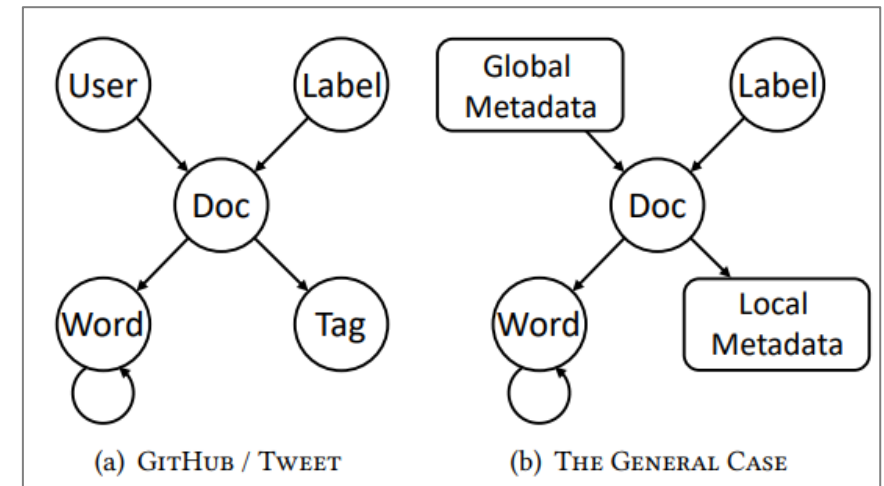**A few labeled documents as supervision.**

# 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**



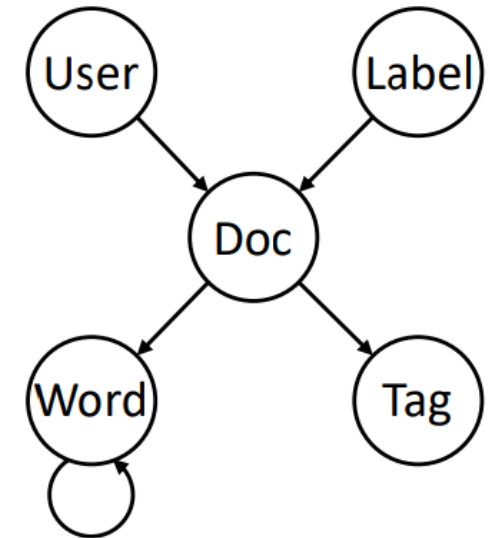**A generative-process view of corpus with metadata**

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

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

# Outline

❑    What Weakly-Supervised Text Classification Is, and Why It Matters

❑     Flat Text Classification

❑    Text Classification with Taxonomy Information

❑    Text Classification with Metadata Information

   ❑    Static 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.



(a) Label "Webgraph" from Microsoft Academic (https://academic.microsoft.com/topic/2777569578/).

(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", 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.

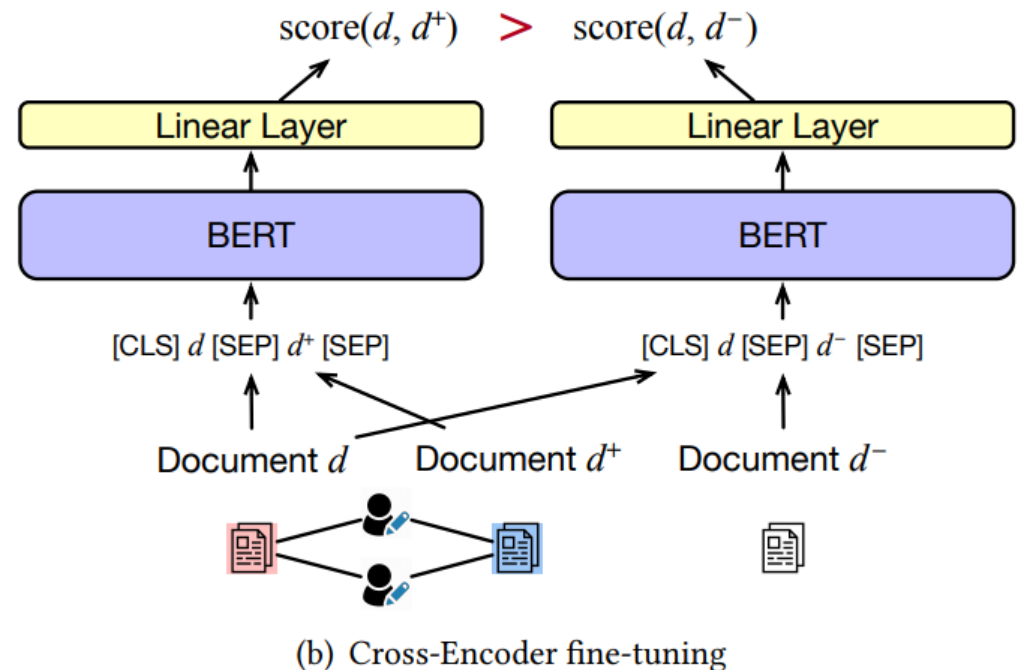    ❑ Both Bi-Encoder and Cross-Encoder are applicable.



(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.



(a) meta-path: PAP
(b) meta-path: P->P<-P
(c) meta-graph: P(AV)P
(d) meta-graph: P<-(PP)->P

Document
Venue
Author

$score(d, d^+) > score(d, d^-)$



$e_d$   $e_{d^+}$   $e_{d^-}$

BERT   BERT   BERT

Document $d$   Document $d^+$   Document $d^-$

(a) Bi-Encoder fine-tuning

$score(d, d^+) > score(d, d^-)$

Linear Layer   Linear Layer

BERT   BERT

[CLS] $d$ [SEP] $d^+$ [SEP]   [CLS] $d$ [SEP] $d^-$ [SEP]

Document $d$   Document $d^+$   Document $d^-$

(b) Cross-Encoder fine-tuning

Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. A simple framework for contrastive learning of visual representations. ICML'20.

# 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 |

# 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 |
| MICoL | Flat | Multi-label | Category Names | Pretrained LM |

48

# 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", AAAI'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., Shen, Z., Wu, C., Xie, B., Wang, Y., Wang, K. & Han, J. "Metadata-Induced Contrastive Learning for Zero-Shot Multi-Label Text Classification", WWW'22.

# Q&A