

Understanding Survey Paper Taxonomy about Large Language Models via Graph Representation Learning

Anonymous ACL submission

Abstract

As new research on Large Language Models (LLMs) continues, it is difficult to keep up with new research and models. To help researchers synthesize the new research many have written survey papers, but even those have become numerous. In this paper, we develop a method to automatically assign survey papers to a taxonomy. We collect the metadata of 144 LLM survey papers and explore three paradigms to classify papers within the taxonomy. Our work indicates that leveraging graph structure information on co-category graphs can significantly outperform the language models in two paradigms; pre-trained language models' fine-tuning and zero-shot/few-shot classifications using LLMs. We find that our model surpasses an average human recognition level and that fine-tuning LLMs using weak labels generated by a smaller model, such as the GCN in this study, can be more effective than using ground-truth labels, revealing the potential of weak-to-strong generalization in the taxonomy classification task.

1 Introduction

Collective attention in the field of Natural Language Processing (NLP)—and the wider public—has turned to Large Language Models (LLMs). It has become so difficult to keep up with the proliferation of new models that many researchers have written survey papers to help synthesize the research progress. Survey papers are often crucial for newcomers to gain an in-depth understanding of the evolution of a research field. However, the volume of survey papers itself has become unruly for researchers—especially newcomers—to sift through. As illustrated in Figure 1, the number of survey papers has been increasing significantly. This leads to our research question, aimed at aiding the field of NLP: Is it possible to automatically reduce the barriers for newcomers in a way that can keep up with the constant influx of new information?

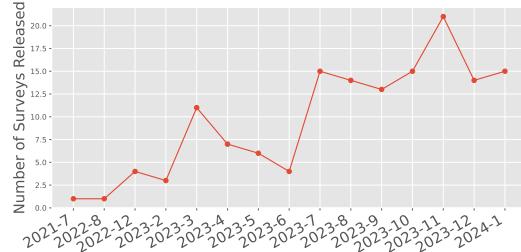


Figure 1: Trends of survey papers about large language models since 2021. Numbers reflect the year and month (e.g., 2023-3 is March 2023).

In this paper, we address the above question by developing a method that can automatically assign survey papers to a taxonomy. Such a taxonomy will help researchers see new trends in the field and focus on specific survey papers that are relevant to their research. Classifying papers into a taxonomy may seem an ordinary task, but it is actually quite challenging for the following reasons:

1. Our dataset contains 144 papers. While this number for the survey papers is uncommonly large, the number of instances in the dataset is still relatively small.
2. We propose a new taxonomy for the collected survey papers, where the distribution of each category is not uniform, which leads to a substantial class imbalance issue.
3. Authors usually use similar terminologies to describe the LLMs in the title and the abstract of these survey papers. Such textual similarity introduces significant difficulties in taxonomy classifications.

To answer our research questions, we investigate three types of attributed graphs: text graphs, co-author graphs, and co-category graphs. Extensive experiments indicate that leveraging graph structure information of co-category graphs can help better classify the survey papers to the corresponding categories in the proposed taxonomy.

Moreover, we validate that graph representation learning (GRL) can outperform language models in two paradigms, fine-tuning pre-trained language models and zero-shot/few-shot classifications using LLMs. Inspired by a recent study, which indicates that leveraging weak labels, which are generated by smaller (weaker) models, may help enhance the performance of larger (stronger) models (Burns et al., 2023), we further examine whether using weak labels, which are generated by GNNs in this study, in the fine-tuning paradigm can help the pre-trained language models. The experiments demonstrate that fine-tuning using weak labels can exceed that using ground-truth labels. For the latter paradigm, we use the results of human recognition as the baseline. The analysis demonstrates that GRL achieves higher accuracy and F1 scores, and even surpasses the average human recognition level by a substantial margin. Overall, our primary contributions can be summarized as follows:

- We collected and analyzed 144 survey papers about LLMs and their metadata.¹
- We propose a new taxonomy for categorizing the survey papers, which will be helpful for the research community, particularly newcomers and multidisciplinary research teams.
- Extensive experiments demonstrate that graph representation learning on co-category graph structure can effectively classify the papers and substantially outperform the language models and average human recognition level on a relatively small and class-imbalanced dataset with high textual similarity.
- Our results also reveal the potential of fine-tuning pre-trained language models using weak labels.

2 Related Work

Taxonomy Classification Conventional taxonomy classification is a subset of Automatic Taxonomy Generation (ATG), which aims to generate taxonomy for a corpus (Krishnapuram and Kummamuru, 2003). The main challenge in ATG is to cluster the unlabeled topics into different hierarchical structures. Thus, most existing methods in ATG are clustering-based methods. Zamir and Etzioni (1998) design a mechanism, Grouper, that dynamically groups and labels the search results. Vaithyanathan and Dom (1999) propose a model to

generate hierarchical clusters. Lawrie et al. (2001) discover the relationship among words to generate concept hierarchies. Within these methods, a subset, called co-clustering, clusters keywords and documents simultaneously (Frigui and Nasraoui, 2002; Kummamuru et al., 2003). Different from ATG, in this study, we classify survey papers into corresponding categories in the proposed taxonomy on relatively small and class-imbalanced datasets, whose text content contains similar terminologies.

Graph Representation Learning Graph representation learning (GRL) is a powerful approach for learning the representation in graph-structure data (Zhou et al., 2020), whereas most recent works achieve this goal using Graph Neural Networks (GNNs) (Veličković et al., 2018; Xu et al., 2018). Bruna et al. (2013) first introduce a significant advancement in convolution operations applied to graph data using both spatial method and spectral methods. To improve the efficiency of the eigen-decomposition of the graph Laplacian matrix, Defferrard et al. (2016) approximate spectral filters by using K-order Chebyshev polynomial. Kipf and Welling (2016) simplify graph convolutions to a first-order polynomial while yielding superior performance in semi-supervised learning. Hamilton et al. (2017) propose an inductive-learning approach that aggregates node features from corresponding fixed-size local neighbors. These GNNs have demonstrated exceptional performance in GRL, underscoring their significance in advancing this field.

3 Methodology

In this section, we first introduce the procedure of data collection and then explore the metadata. We further explain the process of constructing three types of attributed graphs and how we learn graph representation via graph neural networks.

3.1 Data Collection and Exploration

We scraped the metadata of survey papers about LLMs from arXiv and further manually supplemented the metadata from Google Scholar and the ACL anthology. The papers range from July 2021 to January 2024. Given these survey papers, we designed a taxonomy and assigned each paper to a corresponding category within the taxonomy. Our motivation is that a reasonable taxonomy can provide a clear hierarchy of concepts for readers to better understand the relationship among a large num-

¹We will release all datasets and source code after this paper is accepted.

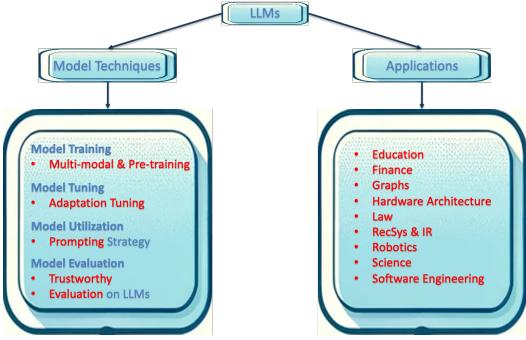


Figure 2: The mind map of survey papers about large language models. Besides "Comprehensive" and "Others" that are not included in the mind map, we highlight fourteen categories in our proposed taxonomy. The total number of categories for the 144 papers is sixteen.

ber of survey papers. Though survey papers can be taxonomized differently, we noticed two broad categories: *applications* and *model techniques*. The *applications* category further sub-divides into specific domains of focus (e.g., education or science), whereas *model techniques* further sub-divides into ways of effecting models (e.g., fine-tuning).

We visualize our proposed taxonomy and highlight fourteen classes, i.e., the leaf nodes, in Figure 2. The total classes in the labels are sixteen, including *comprehensive* and *others* (not shown in the figure). To better understand the distribution of the classes, we present the class distribution in Figure 3. The distribution indicates that the class is extremely imbalanced, introducing a challenge to the taxonomy classification task.

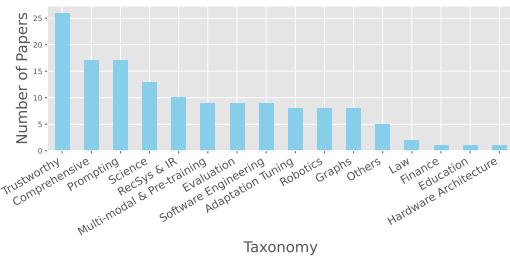


Figure 3: Distribution of classes in the taxonomy.

After visualizing the proposed taxonomy, we further explain the motivation for proposing a new taxonomy instead of using the arXiv categories. In Figure 4, we present the distribution of survey papers across different arXiv categories. Top-2 frequent categories are *cs.CL* (Computation and Language), and *cs.AI* (Artificial Intelligence), which means that most authors choose these two categories for their works. However, these choices cannot help

readers to better distinguish the survey papers. For example, papers related to model techniques are indistinguishable in arXiv categories. Thus, designing a new taxonomy is an essential step in this study.

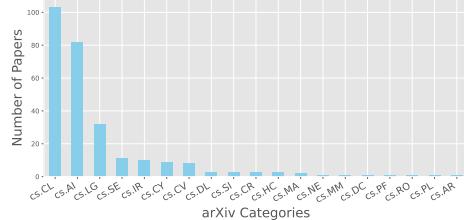


Figure 4: Distribution of survey papers that we found across different arXiv categories.

We also present the word frequency in Figure 5 to show which words have been frequently used in abstracts. These distributions suggest that the abstracts of these papers contain many similar terms, which increases the difficulty of text classification.

Attributes	Descriptions
Taxonomy	Proposed taxonomy
Title	Paper title
Authors	Lists of author's name
Release Date	First released date
Links	Links of papers
Paper ID	The arXiv paper ID
Categories	The arXiv category
Summary	Abstract of papers

Table 1: Data attributes and descriptions.

Overall, we present the data description of their attributes in Table 1. After designing the taxonomy and building the dataset, we explain how we classify documents into the taxonomy categories in the following section.

3.2 Building Attributed Graphs

The goal of building the graphs is to utilize the graph structure information to classify the taxonomy. Before building the graphs, We first define the attributed graphs as follows:

Definition 1 An attributed graph \mathcal{G} is a topological structure that represents the relationships among vertices associated with attributes. \mathcal{G} consists of a set of vertices $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$ and edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$, where N is the number of vertices in \mathcal{G} .

Given the Definition 1, we further define the matrix representation of an attributed graph as follows:

Definition 2 Given an attributed graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, the topological relationship among vertices can

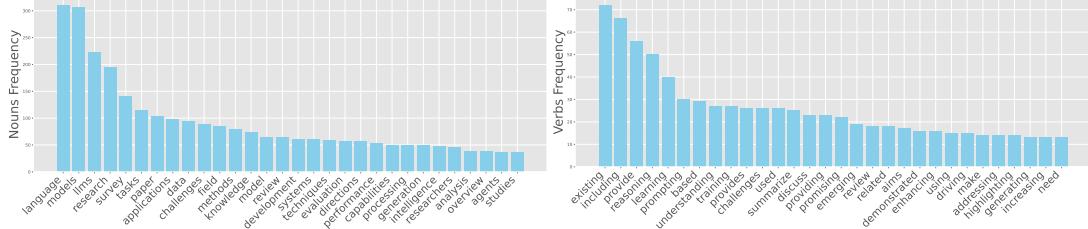


Figure 5: Top 30 keywords frequency in the summary of survey papers.

be represented by a symmetric adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$. Each vertex contains an attribute vector, a.k.a., a feature vector. All feature vectors constitute a feature matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$, where d is the number of features for each vertex. Thus, the matrix representation of an attributed graph can be formulated as $\mathcal{G}(\mathbf{A}, \mathbf{X})$.

Based on the above definitions, we build the graph by creating the term frequency-inverse document frequency (TF-IDF) feature matrices for both title and summary (i.e., abstract) columns, where the term frequency denotes the word frequency in the document, and inverse document frequency denotes the log-scaled inverse fraction of the number of documents containing the word. TF-IDF matrix is commonly used for text classification tasks because it helps capture the distinctive words that can indicate specific classes (Yao et al., 2019). After establishing the TF-IDF matrices, we apply one-hot encoding on the arXiv’s categories and then combine three matrices along the feature dimension to build the feature matrix \mathbf{X} .

To leverage the topological information among vertices, we proceed to construct the graph structures to connect the attribute vectors. In this study, we are interested in three types of graphs: text graph, co-author graph, and co-category graph. We explain each type as follows.

Text Graph We follow the same settings as TextGCN (Yao et al., 2019) to build the text graph. Specifically, the edges of the text graph are built based on word occurrence (paper-word edges) in the paper’s text data, including both title and summary, and word co-occurrence (word-word edges) in the whole text corpus. To obtain the global word co-occurrence information, we slide a fixed-size window on all papers’ text data. Moreover, we calculate the edge weight between a paper vertex and a word vertex using the TF-IDF value of the word in the paper and calculate the edge weight between two-word vertices using point-wise mutual

information (PMI), a popular metric to measure the associations between two words. Note that in the text graph, we don’t use the above feature matrix because only paper vertices contain attribute vectors. To retain consistency, we set all values in the feature matrix as one. For the same reason, only the paper vertices are assigned labels, whereas all word vertices are labeled as a new class, which is not touched during the training or testing phase.

Co-author Graph In the co-author graph, we introduce an edge connecting two vertices (papers) if they share at least one common author.

Co-category Graph In the co-category graph, an edge is added between two vertices with at least one common arXiv category. In the co-authorship and co-category graphs, each vertex is assigned one class (taxonomy) as the label. Note that in this study all edges are undirected.

3.3 Taxonomy Classification via Graph Representation Learning

Given the well-built attributed graphs $\mathcal{G}(\mathbf{A}, \mathbf{X})$, we aim to investigate whether graph representation learning (GRL) using graph neural networks (GNNs) can help classify survey papers into the taxonomy. Before feeding the matrix representation, \mathbf{A} and \mathbf{X} , of the attributed graphs \mathcal{G} into GNNs, we first preprocess the adjacency matrix \mathbf{A} as follows:

$$\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}, \quad (1)$$

where $\tilde{\mathbf{A}} = \mathbf{A} + I_N$, $\tilde{\mathbf{D}} = \mathbf{D} + I_N$. I_N is an identity matrix. $\mathbf{D}_{i,i} = \sum_j \mathbf{A}_{i,j}$ is a diagonal degree matrix.

After preprocessing, we utilize GNNs to learn graph representation. The layer-wise message-passing mechanism of GNNs can be generally formulated as follows:

$$f_{\mathbf{W}^{(l)}}(\hat{\mathbf{A}}, \mathbf{H}^{(l)}) = \sigma(\hat{\mathbf{A}} \mathbf{H}^{(l)} \mathbf{W}^{(l)}), \quad (2)$$

where $\mathbf{H}^{(l)}$ is a node hidden representation in the l -th layer. The dimension of $\mathbf{H}^{(l)}$ in the input layer,

222
223
224
225
226
227
228

229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249

250
251
252
253
254
255
256
257
258
259
260
261
262

263
264
265
266
267
268
269
270
271

272
273
274
275
276
277
278
279
280

281
282
283
284
285
286
287
288
289

middle layer, and output layer is the number of features d , hidden units h , and classes K , respectively. $\mathbf{H}^{(0)} = \mathbf{X}$. $\mathbf{W}^{(l)}$ is the weight matrix in the l -th layer. σ denotes a non-linear activation function, such as ReLU.

In general node classification tasks, GNNs are trained with ground-truth labels $\mathbf{Y} \in \mathbb{R}^{N \times 1}$. In this study, we build the ground-truth labels based on our proposed taxonomy. To simplify the problem, each paper is assigned one primary category as the label, even if the paper sometimes may belong to more than one category. During training, we optimize GNNs with cross-entropy.

In brief, we address the Taxonomy Classification problem via GRL approaches in this study and formally state the problem as follows:

Problem 1 After building an attributed graph $\mathcal{G}(\hat{\mathbf{A}}, \mathbf{X})$ and the ground-truth labels \mathbf{Y} based on the survey metadata, we train a graph neural network (GNN) on the train data and evaluate the taxonomy classification performance on the test data. Our goal is to design a method to better understand (classify) the taxonomy of the survey papers.

4 Experiment

In this section, we evaluate the graph representation learning (GRL)'s effectiveness compared with two paradigms using language models.

Subsets	Graphs	$ \mathcal{V} $	$ \mathcal{E} $	$ F $	$ C $
$Data_{Nov23}$	Text	737	94,943	737	16
	Co-author	112	204	3,065	15
	Co-category	112	4,908	3,065	15
$Data_{Jan24}$	Text	951	137,709	951	17
	Co-author	144	332	3,542	16
	Co-category	144	8,140	3,542	16
$Data_{subset}$	Text	905	128,575	905	12
	Co-author	134	302	3,394	11
	Co-category	134	6,964	3,394	11

Table 2: Statistics of graph datasets. $|\mathcal{V}|$, $|\mathcal{E}|$, $|F|$, and $|C|$ denote the number of nodes, edges, features, and classes, respectively.

Experimental Settings To examine the generalization of our method on various graph structures, we investigate three types of attributed graphs: text graphs, co-author graphs, and co-category graphs, and compare the classification performance of GRL with that of fine-tuning pre-trained language models across three subsets of our data. Both $Data_{Nov23}$ and $Data_{Jan24}$ con-

tain survey papers collected at the end of corresponding months (November 2023 and January 2024). $Data_{Jan24}$ includes a new category, *Hardware Architecture*. We further construct the third subset $Data_{subset}$ by removing some proposed categories with fewer instances in $Data_{Jan24}$; these categories are *Law*, *Finance*, *Education*, *Hardware Architecture*, and *Others*. The motivation for constructing three subsets is to validate the generalization of our method across different subsets since the classification performance may significantly change on small datasets. Also, new categories may emerge at any period because research on LLMS is developing rapidly, and so are related survey papers. For example, a new category, *Hardware Architecture*, emerges in $Data_{Jan24}$. The change of categories may affect the performance as well. Therefore, we investigate our method on three subsets that contain different categories. The statistics of our dataset and corresponding attributed graphs are presented in Table 2. Recall that the text graph consists of paper vertices and word vertices, and thus contains one additional class because all word vertices are labeled as a new class, which is not touched during the training or testing phase.

To evaluate our model, we split the train, validation, and test data as 60%, 20%, and 20%. Due to the potential for random splits to result in an easier task for our model, we ran the experiments five times using random seed IDs from 0 to 4 and reported the mean values with corresponding standard deviations, mean (std). We evaluate the classification performance by accuracy and weighted f1 score. Accuracy is a common metric on classification tasks, whereas the weighted f1 score provides a balanced measure of the class-imbalanced dataset.

4.1 Leveraging Graph Structure Information for Taxonomy Classification

We investigate whether leveraging the graph structure information can help better classify the papers to their corresponding categories in the proposed taxonomy. In this experiment, we construct the attributed graphs based on the text data (including the title and summary) and the relationship of the co-authorship and co-category. To examine the generalization of GRL, we employ GCN (Kipf and Welling, 2016) as a backbone GNN on various graph structures across three subsets. According to Table 3, GNNs fail to learn graph representation on both the text graph and the co-author graph. For

	<i>Data_{Nov23}</i>		<i>Data_{Jan24}</i>		<i>Data_{subset}</i>	
	Accuracy	Weighted-F1	Accuracy	Weighted-F1	Accuracy	Weighted-F1
Text	20.91 (5.45)	14.20 (4.41)	17.86 (7.14)	16.31 (4.49)	23.08 (4.87)	18.82 (1.50)
Co-author	33.04 (8.06)	33.06 (8.69)	20.00 (8.56)	19.24 (8.79)	29.63 (7.03)	29.24 (6.02)
Co-category (All)	63.48 (18.36)	62.82 (16.96)	75.17 (5.52)	74.60 (4.81)	79.26 (6.87)	77.88 (7.52)
Co-category (Rm cs.CL)	70.43 (9.28)	68.46 (9.63)	67.59 (15.36)	65.81 (17.03)	76.30 (12.31)	73.83 (14.53)
Co-category (Rm cs.AI)	73.91 (18.03)	72.41 (18.28)	75.86 (8.99)	75.79 (9.62)	77.04 (3.63)	74.15 (4.11)
Co-category (Rm cs.CL, cs.AI)	26.09 (10.65)	20.19 (10.56)	37.93 (8.45)	35.97 (7.92)	49.63 (7.63)	47.32 (7.59)
Co-category (Rm cs.IR)	63.48 (18.36)	62.82 (16.96)	75.17 (5.52)	74.60 (4.81)	79.26 (6.87)	77.88 (7.52)
Co-category (Rm cs.R0)	63.48 (18.36)	62.82 (16.96)	75.17 (5.52)	74.60 (4.81)	79.26 (6.87)	77.88 (7.52)
Co-category (Rm cs.SE)	65.22 (11.00)	63.21 (10.37)	74.48 (8.33)	75.10 (6.77)	82.96 (6.02)	82.93 (6.22)
Co-category (Rm cs.IR, cs.R0)	63.48 (18.36)	62.82 (16.96)	75.17 (5.52)	74.60 (4.81)	79.26 (6.87)	77.88 (7.52)
Co-category (Rm cs.IR, cs.SE)	65.22 (11.00)	63.21 (10.37)	74.48 (8.33)	75.10 (6.77)	82.96 (6.02)	82.93 (6.22)
Co-category (Rm cs.R0, cs.SE)	65.22 (11.00)	63.21 (10.37)	74.48 (8.33)	75.10 (6.77)	82.96 (6.02)	82.93 (6.22)
Co-category (Rm cs.IR, cs.R0, cs.SE)	65.22 (11.00)	63.21 (10.37)	74.48 (8.33)	75.10 (6.77)	82.96 (6.02)	82.93 (6.22)

Table 3: Evaluation of graph representation learning on three types of attributed graphs across three subsets of our data. We also conducted ablation studies on the graph structure of co-category graphs by removing (denoted by Rm) some arXiv categories. We ran the experiments five times and presented the mean (std).

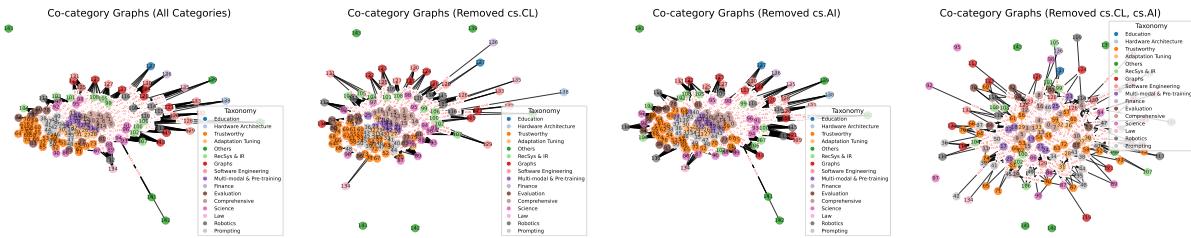


Figure 6: Visualization of four graph structures on co-category graphs by removing the categories.

the text graph, we argue that the degradation may be caused by excessively similar words in the summary of survey papers. When constructing the text graph, these word vertices connect with many paper vertices, resulting in the paper vertices being less distinguishable. For the co-author graph, we conjecture that it is challenging to categorize papers solely based on the sparse co-authorship in this dataset. Furthermore, we observe that some co-authorships come from a common mentor in the same lab whereas two first authors work on the survey papers in two distinct categories. These reasons weaken the effectiveness of using graph structure information. GNNs, in contrast, are very reliable (evaluated by both accuracy and weighted F1 score) in most co-category graphs.

Ablation Analysis We further examine the graph structures of co-category graphs by conducting ablation studies. First, according to Figure 4, most papers are assigned as cs.CL and cs.AI in the arXiv categories. Thus, we study how the categories, cs.CL and cs.AI, affect the performance by muting these two categories in a combinatorial manner. In Table 3, we observe that GNNs can maintain a comparable performance after removing either cs.CL or cs.AI. However, the performance

dramatically drops after removing both categories. This is possible since most node connections are significantly sparsified after these two categories are removed. Even though both cs.CL and cs.AI do not directly map to the existing classes, either one can connect the nodes and further strengthen the message-passing in GNNs, allowing GNNs to learn better node representations.

We visualize co-category graphs in *Data_{Jan24}* in Figure 6. The visualization indicates that most nodes are clustered well even if we remove the category either cs.CL or cs.AI. However, after removing these two categories simultaneously, we observe that node classifications gradually become disordered and several nodes are then isolated. This visualization illustrates the effectiveness of GRL.

We further visualize GCNs' hidden representation on the above co-category graphs in *Data_{Jan24}* in Figure 7. The figures show that the nodes are well-classified in the hidden space even if either the category cs.CL or cs.AI is removed. However, the distribution of nodes tends to become chaotic when both of these two categories are removed simultaneously, shown in Table 3.

For completeness, we conducted another ablation study to examine how the categories, cs.IR,

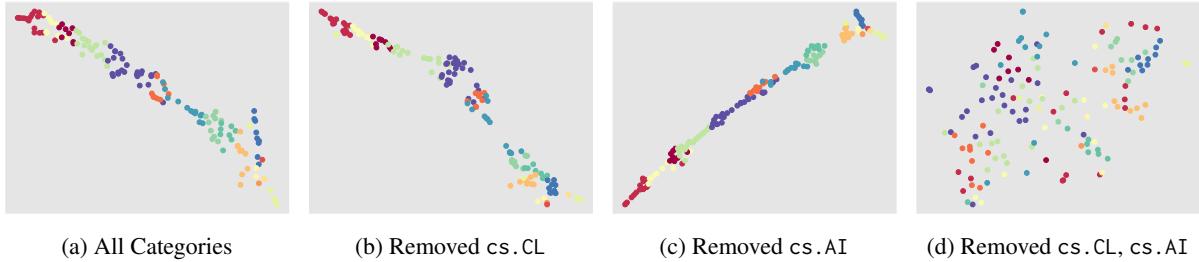


Figure 7: Visualization of GCNs’ hidden representation on co-category graphs in $Data_{Jan24}$ via t-SNE. Each dot represents one node and is labeled with one color.

	$Data_{Nov23}$		$Data_{Jan24}$		$Data_{subset}$	
	Accuracy	Weighted-F1	Accuracy	Weighted-F1	Accuracy	Weighted-F1
BERT (Kenton and Toutanova, 2019)	30.43 (18.45)	25.70 (19.91)	43.45 (18.84)	41.50 (22.31)	58.74 (6.87)	57.51 (6.94)
RoBERTa (Liu et al., 2019)	41.74 (20.32)	39.17 (22.84)	35.86 (17.53)	27.23 (23.67)	25.93 (12.17)	17.02 (15.16)
DistilBERT (Sanh et al., 2019)	57.39 (8.87)	55.59 (10.66)	53.10 (2.76)	52.07 (4.47)	59.26 (7.41)	58.15 (8.82)
XLNet (Yang et al., 2019)	25.22 (16.82)	21.59 (20.54)	27.59 (14.30)	21.59 (15.39)	28.52 (9.37)	21.51 (9.65)
Electra (Clark et al., 2019)	23.04 (4.76)	19.06 (4.06)	44.83 (7.23)	42.01 (8.39)	20.01 (6.87)	12.03 (8.45)
Albert (Lan et al., 2019)	11.30 (8.06)	4.85 (7.21)	15.17 (4.68)	5.14 (2.87)	20.74 (9.83)	11.41 (11.15)
BART (Lewis et al., 2020)	51.30 (17.48)	50.30 (17.62)	51.72 (3.08)	50.62 (2.79)	58.25 (8.11)	57.68 (8.90)
DeBERTa (He et al., 2020)	24.78 (8.06)	19.61 (10.36)	26.21 (11.03)	20.92 (14.25)	25.93 (10.73)	24.30 (10.52)
Llama2 (Touvron et al., 2023)	14.48 (8.72)	4.77 (4.35)	19.22 (5.90)	6.03 (4.21)	23.45 (8.72)	12.59 (7.23)

Table 4: Evaluation of fine-tuning pre-trained language models on the text data across three subsets.

438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
10010
10011
10012
10013
10014
10015
10016
10017
10018
10019
10020
10021
10022
10023
10024
10025
10026
10027
10028
10029
10030
10031
10032
10033
10034
10035
10036
10037
10038
10039
10040
10041
10042
10043
10044
10045
10046
10047
10048
10049
10050
10051
10052
10053
10054
10055
10056
10057
10058
10059
10060
10061
10062
10063
10064
10065
10066
10067
10068
10069
10070
10071
10072
10073
10074
10075
10076
10077
10078
10079
10080
10081
10082
10083
10084
10085
10086
10087
10088
10089
10090
10091
10092
10093
10094
10095
10096
10097
10098
10099
100100
100101
100102
100103
100104
100105
100106
100107
100108
100109
100110
100111
100112
100113
100114
100115
100116
100117
100118
100119
100120
100121
100122
100123
100124
100125
100126
100127
100128
100129
100130
100131
100132
100133
100134
100135
100136
100137
100138
100139
100140
100141
100142
100143
100144
100145
100146
100147
100148
100149
100150
100151
100152
100153
100154
100155
100156
100157
100158
100159
100160
100161
100162
100163
100164
100165
100166
100167
100168
100169
100170
100171
100172
100173
100174
100175
100176
100177
100178
100179
100180
100181
100182
100183
100184
100185
100186
100187
100188
100189
100190
100191
100192
100193
100194
100195
100196
100197
100198
100199
100200
100201
100202
100203
100204
100205
100206
100207
100208
100209
100210
100211
100212
100213
100214
100215
100216
100217
100218
100219
100220
100221
100222
100223
100224
100225
100226
100227
100228
100229
100230
100231
100232
100233
100234
100235
100236
100237
100238
100239
100240
100241
100242
100243
100244
100245
100246
100247
100248
100249
100250
100251
100252
100253
100254
100255
100256
100257
100258
100259
100260
100261
100262
100263
100264
100265
100266
100267
100268
100269
100270
100271
100272
100273
100274
100275
100276
100277
100278
100279
100280
100281
100282
100283
100284
100285
100286
100287
100288
100289
100290
100291
100292
100293
100294
100295
100296
100297
100298
100299
100300
100301
100302
100303
100304
100305
100306
100307
100308
100309
100310
100311
100312
100313
100314
100315
100316
100317
100318
100319
100320
100321
100322
100323
100324
100325
100326
100327
100328
100329
100330
100331
100332
100333
100334
100335
100336
100337
100338
100339
100340
100341
100342
100343
100344
100345
100346
100347
100348
100349
100350
100351
100352
100353
100354
100355
100356
100357
100358
100359
100360
100361
100362
100363
100364
100365
100366
100367
100368
100369
100370
100371
100372
100373
100374
100375
100376
100377
100378
100379
100380
100381
100382
100383
100384
100385
100386
100387
100388
100389
100390
100391
100392
100393
100394
100395
100396
100397
100398
100399
100400
100401
100402
100403
100404
100405
100406
100407
100408
100409
100410
100411
100412
100413
100414
100415
100416
100417
100418
100419
100420
100421
100422
100423
100424
100425
100426
100427
100428
100429
100430
100431
100432
100433
100434
100435
100436
100437
100438
100439
100440
100441
100442
100443
100444
100445
100446
100447
100448
100449
100450
100451
100452
100453
100454
100455
100456
100457
100458
100459
100460
100461
100462
100463
100464
100465
100466
100467
100468
100469
100470
100471
100472
100473
100474
100475
100476
100477
100478
100479
100480
100481
100482
100483
100484
100485
100486
100487
100488
100489
100490
100491
100492
100493
100494
100495
100496
100497
100498
100499
100500
100501
100502
100503
100504
100505
100506
100507
100508
100509
100510
100511
100512
100513
100514
100515
100516
100517
100518
100519
100520
100521
100522
100523
100524
100525
100526
100527
100528
100529
100530
100531
100532
100533
100534
100535
100536
100537
100538
100539
100540
100541
100542
100543
100544
100545
100546
100547
100548
100549
100550
100551
100552
100553
100554
100555
100556
100557
100558
100559
100560
100561
100562
100563
100564
100565
100566
100567
100568
100569
100570
100571
100572
100573
100574
100575
100576
100577
100578
100579
100580
100581
100582
100583
100584
100585
100586
100587
100588
100589
100590
100591
100592
100593
100594
100595
100596
100597
100598
100599
100600
100601
100602
100603
100604
100605
100606
100607
100608
100609
100610
100611
100612
100613
100614
100615
100616
100617
100618
100619
100620
100621
100622
100623
100624
100625
100626
100627
100628
100629
100630
100631
100632
100633
100634
100635
100636
100637
100638
100639
100640
100641
100642
100643
100644
100645
100646
100647
100648
100649
100650
100651
100652
100653
100654
100655
100656
100657
100658
100659
100660
100661
100662
100663
100664
100665
100666
100667
100668
100669
100670
100671
100672
100673
100674
100675
100676
100677
100678
100679
100680
100681
100682
100683
100684
100685
100686
100687
100688
100689
100690
100691
100692
100693
100694
100695
100696
100697
100698
100699
100700
100701
100702
100703
100704
100705
100706
100707
100708
100709
100710
100711
100712
100713
100714
100715
100716
100717
100718
100719
100720
100721
100722
100723
100724
100725
100726
100727
100728
100729
100730
100731
100732
100733
100734
100735
100736
100737
100738
100739
100740
100741
100742
100743
100744
100745
100746
100747
100748
100749
100750
100751
100752
100753
100754
100755
100756
100757
100758
100759
100760
100761
100762
100763
100764
100765
100766
100767
100768
100769
100770
100771
100772
100773
100774
100775
100776
100777
100778
100779
100780
100781
100782
100783
100784
100785
100786
100787
100788
100789
100790
100791
100792
100793
100794
100795
100796
100797
100798
100799
100800
100801
100802
100803
100804
100805
100806
100807
100808
100809
100810
100811
100812
100813
100814
100815
100816
100817
100818
100819
100820
100821
100822

the taxonomy classification task. Recently, Burns et al. (2023) verified that training stronger models with pseudo labels, a.k.a. weak labels, generated by weaker models can enable the stronger models to achieve comparable performance as closely as those trained with ground-truth labels. In this experiment, we first generate weak labels by GCN on co-category graphs and then fine-tune pre-trained language models with weak labels. We present the results on *DataJan24* in Figure 8 as an example. The results indicate that performance achieved through training with weak labels can surpass that of training with ground-truth labels. One possible reason is that training the model using noisy labels with a low noise ratio can be equivalent to a kind of regularization, improving the classification results (Zhuang and Al Hasan, 2022).²

This experiment demonstrates that leveraging weak labels generated by smaller models may effectively enhance the performance of larger models. This is one of the applications related to "weak-to-strong generalization" (Burns et al., 2023).

4.3 LLM Zero-shot/Few-shot Classification and Human Evaluation

	Accuracy	Weighted-F1
Human	58.73 (19.16)	59.50 (19.13)
Claude w.o. hints	11.61 (1.27)	12.66 (0.14)
Claude w. hints	10.27 (3.15)	12.81 (2.10)
GPT 3.5 w.o. hints	47.32 (3.25)	43.21 (4.33)
GPT 3.5 w. hints	53.57 (2.81)	53.16 (3.13)
GPT 4 w.o. hints	29.76 (7.22)	26.91 (9.66)
GPT 4 w. hints	33.04 (5.57)	27.78 (7.76)

Table 5: Evaluation of zero-shot and few-shot classification capabilities of large language models, Claude, GPT 3.5, and GPT 4. We compare these results with those from human recognition.

Major	CS	DS	Security	Math	Chemistry
#Students	9	4	2	2	1

Table 6: The number of students in each major. CS, DS, and Security denote the major in Computer Science, Data Science, and Cyber-security, respectively.

In this experiment, we evaluate zero-shot/few-shot classification capabilities of LLMs Claude (Bai et al., 2022), GPT 3.5 (Brown et al., 2020), and GPT 4 (Achiam et al., 2023), on the text data, which contains both title and summary, on

²Weak labels can be regarded as a kind of noisy label.

DataNov23 as an example. We also compare the results with human participants. We recruited 18 students across five different majors in a graduate-level course. The number of students in each major is shown in Table 6. Each participant was given the titles and abstracts of survey papers and was asked to assign a category to each paper from our taxonomy. We present the mean value with the corresponding standard deviation in Table 5. For the LLMs, we ran the experiments five times. The standard deviation in human recognition is relatively large as some students do not have a strong technical background so they perform worse in this test. Among the LLMs, GPT 3.5 outperforms the other two models given that all models have not seen the data before (zero-shot). We further provide some hints to the models before classification (few-shot). For example, we release the keywords of the class "Trustworthy" to the models before classification. In this setting, both GPT 3.5 and GPT 4 can achieve higher accuracy and a weighted F1 score after obtaining some hints. In brief, GRL can outperform all three LLMs and human recognition, whereas these LLMs couldn't surmount human recognition, which reveals that LLMs still have much room to improve in taxonomy classification.

5 Conclusion

In this work, we aim to develop a method to automatically assign survey papers about Large Language Models (LLMs) to a taxonomy. To achieve this goal, we first collected the metadata of 144 LLM survey papers and proposed a new taxonomy for these papers. We further explored three paradigms to classify survey papers into the categories in the proposed taxonomy. After investigating three types of attributed graphs, we observed that leveraging graph structure information on co-category graphs can significantly help the taxonomy classification. Furthermore, our analysis validates that graph representation learning outperforms pre-trained language models' fine-tuning, zero-shot/few-shot classifications using LLMs, and even surpasses an average human recognition level. Last but not least, our experiments indicate that fine-tuning pre-trained language models using weak labels, which are generated by a weaker model, such as GCN, can be more effective than using ground-truth labels, revealing the potential for weak-to-strong generalization in the taxonomy classification task.

Limitations & Future Work	Constructing a graph structure may encounter certain constraints. For instance, we build co-category graphs based on the arXiv categories. When papers come from distinct fields, such as biology, physics, and computer science, the graph structure may be very sparse, weakening the effectiveness of GRL.	615 616 617 618 619 620 621
In the future, our primary motivation extended from this study is to tailor GPT-based applications to assist readers in understanding survey papers more effectively. We also plan on further exploring the weak-to-strong generalization which could potentially have many important applications.		622 623 624 625
References		626 627 628 629
Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	Hichem Frigui and Olfa Nasraoui. 2002. Simultaneous categorization of text documents and identification of cluster-dependent keywords. In <i>2002 IEEE World Congress on Computational Intelligence. 2002 IEEE International Conference on Fuzzy Systems. FUZZ-IEEE'02. Proceedings (Cat. No. 02CH37291)</i> , volume 2, pages 1108–1113. IEEE.	630 631 632 633
Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022. Constitutional ai: Harmlessness from ai feedback. <i>arXiv preprint arXiv:2212.08073</i> .	Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In <i>Advances in neural information processing systems</i> , pages 1024–1034.	634 635 636
Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. <i>Advances in neural information processing systems</i> , 33:1877–1901.	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. <i>arXiv preprint arXiv:2006.03654</i> .	637 638 639 640
Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. 2013. Spectral networks and locally connected networks on graphs. <i>arXiv preprint arXiv:1312.6203</i> .	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of naacL-HLT</i> , volume 1, page 2.	641 642 643 644 645
Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, et al. 2023. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. <i>arXiv preprint arXiv:2312.09390</i> .	Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. <i>arXiv preprint arXiv:1609.02907</i> .	646 647 648 649 650
Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2019. Electra: Pre-training text encoders as discriminators rather than generators. In <i>International Conference on Learning Representations</i> .	Raghu Krishnapuram and Krishna Kummaru. 2003. Automatic taxonomy generation: Issues and possibilities. In <i>International Fuzzy Systems Association World Congress</i> , pages 52–63. Springer.	651 652 653 654 655
Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In <i>Advances in neural information processing systems</i> , pages 3844–3852.	Krishna Kummaru, Ajay Dhawale, and Raghu Krishnapuram. 2003. Fuzzy co-clustering of documents and keywords. In <i>The 12th IEEE International Conference on Fuzzy Systems, 2003. FUZZ'03.</i> , volume 2, pages 772–777. IEEE.	656 657 658 659 660 661 662 663
Yinhan Liu, Myle Ott, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880.	Mike Lewis, Yinhan Liu, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	664 665 666 667 668

669 Victor Sanh, Lysandre Debut, Julien Chaumond, and
670 Thomas Wolf. 2019. Distilbert, a distilled version
671 of bert: smaller, faster, cheaper and lighter. *arXiv*
672 preprint arXiv:1910.01108.

673 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-
674 bert, Amjad Almahairi, Yasmine Babaie, Nikolay
675 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti
676 Bhosale, et al. 2023. Llama 2: Open founda-
677 tion and fine-tuned chat models. *arXiv preprint*
678 arXiv:2307.09288.

679 Shivakumar Vaithyanathan and Byron Dom. 1999.
680 Model selection in unsupervised learning with ap-
681 plications to document clustering. In *ICML*, pages
682 433–443.

683 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
684 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
685 Kaiser, and Illia Polosukhin. 2017. Attention is all
686 you need. *Advances in neural information processing*
687 systems, 30.

688 Petar Veličković, Guillem Cucurull, Arantxa Casanova,
689 Adriana Romero, Pietro Liò, and Yoshua Bengio.
690 2018. Graph attention networks. In *International*
691 *Conference on Learning Representations*.

692 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie
693 Jegelka. 2018. How powerful are graph neural net-
694 works? *arXiv preprint arXiv:1810.00826*.

695 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-
696 bonell, Russ R Salakhutdinov, and Quoc V Le. 2019.
697 Xlnet: Generalized autoregressive pretraining for lan-
698 guage understanding. *Advances in Neural Infor-
699 mation Processing Systems*, 32.

700 Liang Yao, Chengsheng Mao, and Yuan Luo. 2019.
701 Graph convolutional networks for text classification.
702 In *Proceedings of the AAAI conference on artificial*
703 *intelligence*, volume 33, pages 7370–7377.

704 Oren Zamir and Oren Etzioni. 1998. Web document
705 clustering: A feasibility demonstration. In *Pro-
706 ceedings of the 21st annual international ACM SIGIR*
707 *conference on Research and development in informa-*
708 *tion retrieval*, pages 46–54.

709 Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan
710 Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang,
711 Changcheng Li, and Maosong Sun. 2020. Graph
712 neural networks: A review of methods and applica-
713 tions. *AI open*, 1:57–81.

714 Jun Zhuang and Mohammad Al Hasan. 2022. Robust
715 node classification on graphs: Jointly from bayesian
716 label transition and topology-based label propagation.
717 In *Proceedings of the 31st ACM International Con-
718 ference on Information & Knowledge Management*,
719 pages 2795–2805.

A APPENDIX

In the appendix, we present the GNN and pre-trained language models’ hyper-parameters and the hardware and software. We also include the additional comparison results about fine-tuning using weak labels and additional visualization of co-category graphs.

Hyper-parameters and Settings We employ a two-layer GCN (Kipf and Welling, 2016) with 200 hidden units and a ReLU activation function as the backbone GNN to examine the effectiveness of GRL. The GNN is trained by the Adam optimizer with a learning rate, 1×10^{-2} for both co-author graphs and co-category graphs and 2×10^{-2} for text graphs, and converged within 500 training epochs on all subsets. The dropout rate is 0.5.

Language Models	Model Size
BERT (Kenton and Toutanova, 2019)	109.49M
RoBERTa (Liu et al., 2019)	124.66M
DistILBERT (Sanh et al., 2019)	66.97M
XLNet (Yang et al., 2019)	117.32M
Electra (Clark et al., 2019)	109.49M
Albert (Lan et al., 2019)	11.70M
BART (Lewis et al., 2020)	140.02M
DeBERTa (He et al., 2020)	139.20M
Llama2 (Touvron et al., 2023)	6.61B

Table 7: Model size of the pre-trained language models. For example, BERT has 109.49 million parameters.

We fine-tune the pre-trained language models using the Adam optimizer with a 1×10^{-4} learning rate. We chose the batch size of 8 for the Llama2 and fixed the batch size of 16 for the rest of the models. We implement the pre-trained language models using HuggingFace packages (we choose the base version for all models) and report the model size in Table 7. All models are tuned with 30 epochs.

Hardware and Software The experiment is conducted on a server with the following settings:

- Operating System: Ubuntu 22.04.3 LTS
- CPU: Intel Xeon w5-3433 @ 4.20 GHz
- GPU: NVIDIA RTX A6000 48GB
- Software: Python 3.11, PyTorch 2.1, HuggingFace 4.31, dgl 1.1.2+cu118.

Computational Budgets Based on the above computing infrastructure and settings, computational budgets in our experiments are described as follows. The experiment presented in Table 3

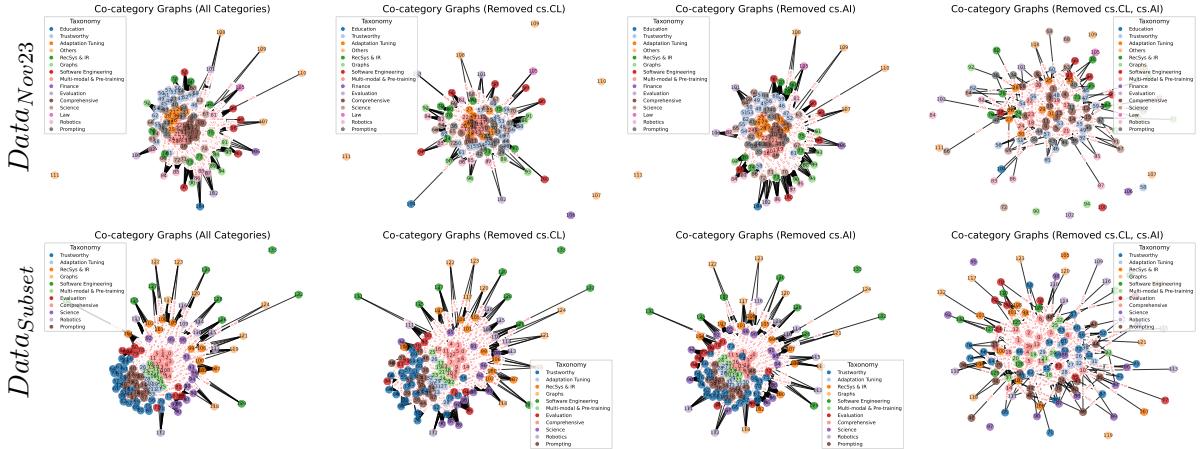


Figure 9: Additional visualization of co-category graphs (extended from Figure 6).

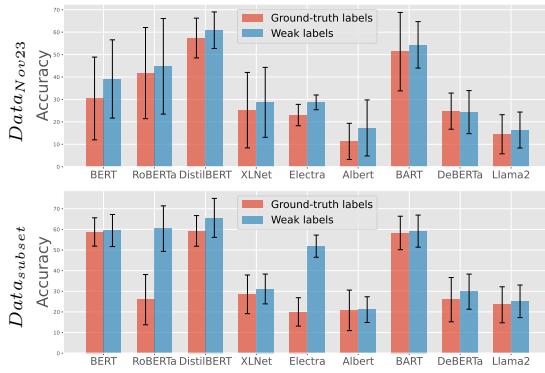


Figure 10: Additional comparison in the fine-tuning paradigm using weak labels (extended from Figure 8).

can be reproduced within one hour. The experiment displayed in Table 4 may take 93 hours to complete. Due to limited GPU memory, we implemented Llama2 using the CPU. This consumes around 90 hours in total. The experiment shown in Table 5 (excluding human recognition) can be finished in one hour.

Additional Visualization of Co-category Graphs
Besides visualizing four graph structures in *DataJan24* in Figure 6, we additionally present the visualization of four corresponding co-category graphs in both *DataNov23* and *Datasubset* in Figure 9. The visualization verifies the generalization of GRL across three subsets.

Additional Comparison of Fine-tuning Using Weak Labels
Besides the results in Figure 8, we supplement the comparisons on both *DataNov23* and *Datasubset* in Figure 10. The comparisons across nice pre-trained language models further validate the effectiveness of fine-tuning using weak labels.

Ethical and Broader Impacts We confirm that we fulfill the author’s responsibilities and address the potential ethical issues. In this work, we aim to help researchers quickly and better understand a new research field. Many researchers in academia or industry may potentially benefit from our work.

Statement of Data Privacy Our dataset contains the authors’ names in each paper. This information is publicly available so the collection process doesn’t infringe on personal privacy.

Disclaimer Regarding Human Subjects Results
In Table 5, we include partial results with human subjects. We already obtained approval from the Institutional Review Board (IRB). The protocol number is IRB24-056. We recruited volunteers from a graduate-level course. Before the assessment, we have disclaimed the potential risk (our assessment has no potential risk) and got consent from participants.

755									
756									
757									
758									
759									
760									
761									
762									
763									
764									
765									
766									
767									
768									
769									
770									
771									
772									
773									
774									
775									
776									
777									
778									
779									
780									
781									
782									
783									
784									
785									