

# BOLT: Benchmarking Open-World Learning for Text Classification

Chuan Qin<sup>1,2,♣</sup>, Xi Chen<sup>3,1,♣</sup>, Jinpeng Li<sup>1,2,♣</sup>, Hengshu Zhu<sup>1,2,♠</sup>

<sup>1</sup> Computer Network Information Center, Chinese Academy of Sciences

<sup>2</sup> University of Chinese Academy of Sciences

<sup>3</sup> School of Computer Science and Technology, University of Science and Technology of China

{chuanqin0426, zhuhengshu}@gmail.com

chenxi0401@mail.ustc.edu.cn, lijinpeng25@mails.ucas.ac.cn

♣ Equal contribution   ♠ Corresponding author

## Abstract

Text classification has long been a cornerstone of NLP, yet most prior work and benchmarks have been limited to closed-world settings, where all classes are assumed to be known in advance. In contrast, open-world learning has recently emerged as a critical paradigm for building more robust and realistic systems. However, existing benchmarks largely focus on out-of-distribution (OOD) detection, while overlooking broader challenges such as the discovery of novel categories. To address this gap, we introduce **BOLT**, a unified **B**enchmark and evaluation toolkit supporting **O**pen-world **L**earning for **T**ext classification. BOLT encompasses two representative tasks: Open-set Text Classification (OSTC), which requires models to classify in-distribution (ID) samples while rejecting OOD inputs, and Generalized Category Discovery (GCD), which aims to identify both known and novel categories from partially labeled corpora. We carefully curate 12 publicly available datasets spanning diverse domains and benchmark 22 methods, including 15 for OSTC and 7 for GCD, under a standardized protocol that explicitly accounts for varying labeled ratios and known class ratios. Our results reveal key challenges: most current methods tend to overfit training distributions and struggle to generalize to unseen classes. Moreover, by comparing our lightweight LLM-based variants with prior open-set baselines, we demonstrate the promise of leveraging LLMs for open-world text classification. BOLT provides standardized evaluation protocols that enable fair comparison and support future research in this emerging area. All datasets, baselines, and tools are available at <https://github.com/CNIC-DSL/BOLT>.

## 1 Introduction

Text classification has long been a central task in natural language processing (NLP), supporting diverse applications such as sentiment analysis (Wankhade et al., 2022; Qin et al., 2025c),

topic labeling (Chen et al., 2020; Wang et al., 2023; Chen et al., 2024; Qin et al., 2025b), intent recognition (Zhang et al., 2024a), and document management (Shu et al., 2017). Traditional research in this area has largely relied on closed-world assumptions, where the set of target classes is predefined and assumed to be fully observable during training (Minaee et al., 2021). Such an assumption limits applicability in real-world scenarios, where label spaces are inherently open and novel categories continuously emerge (Fei and Liu, 2016).

To address these limitations, recent research has shifted toward open-world learning, a paradigm originally studied in computer vision (Li and Wechsler, 2005; Zhang et al., 2023) and later extended to NLP (Prakhya et al., 2017). Building on this foundation, open-world text classification (OTC) requires models to classify ID data while remaining robust to emerging categories. Current research on OTC primarily investigates two representative tasks. The first, OSTC, addresses the challenge of classifying ID inputs while reliably rejecting OOD samples, under the restrictive setting where only labeled ID data are available for training (Chen et al., 2026). This task requires learning robust decision boundaries that generalize beyond the observed training distribution (Shu et al., 2017; Zhang et al., 2021a). The second, GCD, extends this objective by discovering novel categories from unlabeled corpora through clustering-based pseudo-labeling, which requires not only precise ID classification but also effective structure induction to accommodate emerging classes (Zhang et al., 2021b, 2024b). By integrating these two perspectives, OTC bridges the gap between closed-world classification and the demands of dynamic real-world environments.

However, existing OSTC and GCD methods primarily rely on a few intent recognition datasets for evaluation, such as Banking and CLINC. This line of evaluation suffers from inherent limitations, as it typically focuses on a single task type and is

Task	KCR	LAR	Evaluation Metrics	Paper List
OSTC	0.250, 0.500, 0.750	1.000	F1	DOC [(Shu et al., 2017)], DeepUnk [(Lin and Xu, 2019)]
	0.250, 0.500, 0.750	1.000	ACC, F1, K-F1, N-F1	ADB [(Zhang et al., 2021a)], AB [(Lorenc et al., 2022)], KNNCon [(Zhou et al., 2022)]
	0.667	1.000	ACC, F1, Recall	SCL [(Zeng et al., 2021a)]
	0.250, 0.750	1.000	ACC, F1, K-F1, N-F1	DyEn [(Zhou et al., 2023)]
	0.250, 0.500, 0.750	0.100, 0.500, 1.000	ACC, F1, K-F1, N-F1	<b>BOLT</b>
GCD	0.750	0.100	NMI, ARI, ACC	DAL [(Zhang et al., 2021b)]
	0.250	0.100	NMI, ARI, ACC	GeoID [(Tang et al., 2024)]
	0.750	0.075	ACC	DPN [(An et al., 2023)]
	0.750	0.100	H-score, K-ACC, N-ACC	SDC [(An et al., 2025)], TAN [(An et al., 2024b)], LOOP [(An et al., 2024a)]
	0.250, 0.500	0.100	ACC, ARI, NMI	ALUP [(Liang et al., 2024)]
	0.250, 0.500, 0.750	0.100	ACC, ARI, NMI	USNID [(Zhang et al., 2023)]
	0.250, 0.500, 0.750	0.100, 0.500, 1.000	ACC, ARI, NMI, H-score, K-ACC, N-ACC	<b>BOLT</b>

Table 1: Comparison of evaluation settings for OSTC and GCD.

confined to homogeneous text domains and data characteristics; for instance, these datasets exhibit narrowly concentrated length distributions and are uniformly in English. This narrow coverage hinders systematic investigation of model generalization across diverse scenarios. In addition, existing studies often adopt inconsistent experimental protocols, such as using different proportions of known classes, varying amounts of labeled data, or divergent evaluation metrics, which further complicates fair comparison across methods as shown in Table 1. In particular, OSTC is seldom evaluated under low-resource settings, while GCD has seen only limited exploration across varying levels of supervisory signals. Consequently, **the field still lacks a comprehensive benchmark and robust evaluation toolkit** to support systematic and reliable exploration of OTC. Moreover, although LLMs have been widely adopted in NLP, their use in OTC remains limited. Current attempts mainly exploit data augmentation to improve GCD (Liang et al., 2024), while other directions, including employing LLMs as backbones for representation learning, are largely unexplored.

To this end, we present **BOLT**, a novel and comprehensive **Benchmark** and evaluation toolkit designed to support **Open-world Learning for Text** classification. Specifically, we first unify two representative OTC tasks, OSTC and GCD, within a single evaluation framework. To enhance the diversity of our benchmark, we carefully curate 12 publicly available datasets that span a wide range of text classification scenarios, including intent detection, topic categorization, and ontology classification. These datasets cover diverse text domains (e.g., news, conversational queries, and question-answer), exhibit varying text lengths (with an average document length ranging from 8 to 1,856 words), differ in label granularity (with class counts between 10 and 219), and vary in multiple languages. Along this line, we conduct a systematic

and fine-grained evaluation of 22 representative methods, including 15 designed for OSTC and 7 for GCD, under a standardized protocol that explicitly accounts for different labeled and known-class ratios. The experimental results not only reveal the limitations of existing methods but also demonstrate the promise of **leveraging LLMs for advancing different OTC tasks**.

## 2 Related Work

In this section, we review prior work on the two representative OTC tasks, i.e., OSTC and GCD.

### 2.1 Open-Set Text Classification

OSTC is a challenging task that requires models to classify ID textual inputs into known categories while detecting OOD inputs (Prakhya et al., 2017). This dual objective contrasts with conventional OOD detection, which is typically framed as a binary discrimination problem between ID and OOD samples (Liu et al., 2020). Early studies on this problem setting mainly originated in computer vision, where maximum softmax probability (MSP) (Hendrycks and Gimpel, 2016) was introduced for OOD detection. However, this approach merely rejects uncertain predictions without explicitly modeling unknown classes (Bendale and Boult, 2016). OpenMax attempted to improve MSP by fitting a Weibull distribution to classifier logits to estimate OOD probabilities (Bendale and Boult, 2016). However, its reliance on the assumption that equiprobable logits indicate unknown classes often leads to misclassification of hard ID examples. Extending to the text domain, DOC introduced a one-vs-rest sigmoid layer to mitigate open-space risk (Shu et al., 2017). Nevertheless, threshold-based approaches as a whole remain limited in capturing the intrinsic differences between known and unknown classes.

Subsequent work shifted from thresholding to enhancing representation learning. DeepUnk (Lin

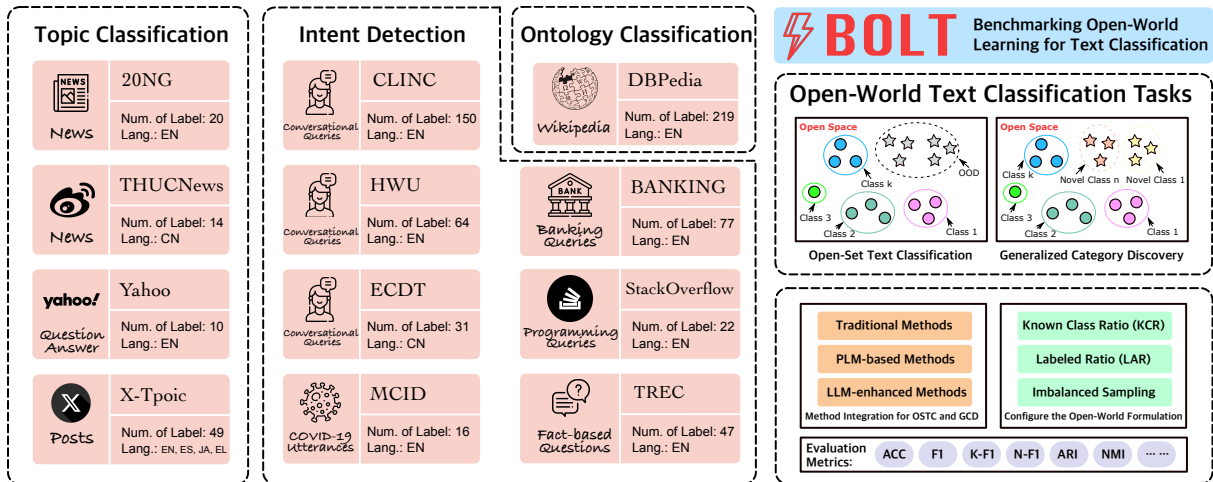


Figure 1: An overview of our BOLT benchmark for open-world text classification, covering two OTC tasks (OSTC and GCD) over 12 datasets from three scenarios (topic, intent, ontology) across eight domains. BOLT supports configurable open-world settings via KCR, LAR, and imbalanced sampling, and provides a standardized pipeline to evaluate traditional, PLM-based, and LLM-enhanced methods with consistent metrics for both tasks.

and Xu, 2019) adopted a large-margin cosine loss to better separate features, while ADB (Zhang et al., 2021a) refined decision boundaries to improve OOD discrimination. KNNCon (Zhou et al., 2022) integrated KNN-based contrastive learning with LOF-based detection, and DyEn (Zhou et al., 2023) introduced dynamic ensembling with early exits to mitigate overfitting. Recently, the progress of LLMs has significantly advanced the field of NLP, enabling notable improvements in semantic understanding and text generation (Zhao et al., 2023; Jiang et al., 2024; Tong et al., 2025; Qin et al., 2025a; Zhou et al., 2025; Huang et al., 2026; Song et al., 2026). Although a few initial studies have begun to explore the use of LLMs for OOD detection (Liu et al., 2023), their potential for OSTC has not yet been investigated. In addition, most existing OSTC evaluations focus on a narrow range of supervision settings, leaving low-resource scenarios largely underexplored.

## 2.2 Generalized Category Discovery

In contrast to OSTC, GCD aims to identify both known and novel categories from partially labeled data. Most existing approaches first learn representations from labeled known classes and then apply clustering-based pseudo-labeling to discover novel categories. Early methods such as DeepAligned (Zhang et al., 2021b) employed pairwise-similarity supervision or alignment strategies to guide clustering. However, the accuracy of pseudo-labels remains a critical bottleneck, as noise propagation can severely degrade overall

performance. Subsequent studies sought to mitigate pseudo-label noise through two main strategies. Pseudo-label calibration approaches, including GeoID (Tang et al., 2024) and SDC (An et al., 2025), alleviate the dominance of known categories during clustering. Prototype-based transfer methods, including TAN (An et al., 2024b) and DPN (An et al., 2023), facilitate knowledge transfer from known to novel classes.

Recently, LLMs have been incorporated into GCD. LOOP (An et al., 2024a) and ALUP (Liang et al., 2024) leverage LLMs to refine neighborhoods or enable exemplar-based comparison.

## 3 The BOLT Benchmark

### 3.1 Problem Definition

We consider the OTC problem, where models are required to handle both ID and OOD samples under varying supervision settings. Let the textual input space be  $\mathcal{X}$  and the label space be  $\mathcal{Y} = \mathcal{Y}^k \cup \mathcal{Y}^n$ , where  $\mathcal{Y}^k = \{1, \dots, K\}$  denotes the set of known (predefined) classes, and  $\mathcal{Y}^n = \{K + 1, \dots, K + N\}$  represents the set of novel classes that may appear during testing.

Given the dataset setup: a labeled set  $\mathcal{D}^l = \{(x_i, y_i)\}_{i=1}^{N_l}$ , where  $x_i \in \mathcal{X}$  and  $y_i \in \mathcal{Y}^k$ ; an unlabeled set  $\mathcal{D}^u = \{x_j\}_{j=1}^{N_u}$ , which may contain samples drawn from both  $\mathcal{Y}^k$  and  $\mathcal{Y}^n$ ; and a test set  $\mathcal{D}^t = \{(x_i, y_i)\}_{i=1}^{N_t}$ , where  $y_i \in \mathcal{Y}^k \cup \mathcal{Y}^n$ . Based on the availability of labeled and unlabeled data, as well as the target prediction objectives, we define two representative OTC tasks.

Scenario	Text Domain	Dataset	Lang.	# Inst.	# Lab.	Avg. Len./ Inst.	Avg. Len./ Lab.
Topic Classification	Usenet News	20NG	EN	10000	20	290.41	2.80
Topic Classification	Sina News	THUCNews	CN	9421	14	749.07	2.00
Topic Classification	Yahoo pages	Yahoo	EN	10000	10	95.19	1.80
Topic Classification	X pages	X-Topic	Multi	8587	49	26.86	1.00
Intent Detection	Online banking queries	BANKING	EN	13083	77	11.74	3.35
Intent Detection	Conversational queries	CLINC	EN	22500	150	8.23	1.94
Intent Detection	Programming queries	StackOverflow	EN	19985	20	8.36	1.05
Intent Detection	Conversational queries	HWU	EN	9677	64	6.60	2.12
Intent Detection	Fact-based questions	TREC	EN	5871	47	8.72	2.00
Intent Detection	Conversational queries	ECDT	CN	3069	31	7.87	1.00
Intent Detection	COVID-19 utterance	MCID	EN	1699	16	6.87	3.00
Ontology Classification	Wikipedia pages	DBPedia	EN	10000	219	98.31	1.00

Table 2: Dataset statistics in BOLT. We report the total number of instances (# Inst.), the number of labels (# Lab.), the average text length per instance (Avg. Len./ Inst.), and the average text length per label (Avg. Len./ Lab.).

**Open-Set Text Classification (OSTC)** In this setting, the model is trained only on the labeled data  $\mathcal{D}^l$  containing known classes  $\mathcal{Y}^k$ . During testing, it encounters OOD samples belonging to the novel label space  $\mathcal{Y}^n$ , which is not observed during training. The objective is to learn a  $K + 1$  classifier:

$$f_{\text{OOD}} : \mathcal{X} \rightarrow \mathcal{Y}^k \cup \{\text{Unknown}\}, \quad (1)$$

that correctly classifies ID samples into  $\mathcal{Y}^k$  and identifies OOD samples by assigning them to the “Unknown” class representing  $\mathcal{Y}^n$ .

**Generalized Category Discovery (GCD)** In this setting, both the labeled data  $\mathcal{D}^l$  and the unlabeled data  $\mathcal{D}^u$  are available during training. The unlabeled data contain a mixture of known and novel classes, i.e.,  $x_j \in \mathcal{X}$  with  $y_j \in \mathcal{Y}^k \cup \mathcal{Y}^n$ . The task requires the model to jointly recognize known categories and discover novel ones. Formally, the model learns from  $\mathcal{D}^l \cup \mathcal{D}^u$  to construct a model:

$$f_{\text{GCD}} : \mathcal{X} \rightarrow \mathcal{Y}^k \cup \mathcal{Y}^n, \quad (2)$$

where the number of novel categories  $|\mathcal{Y}^n|$  is assumed to be known (An et al., 2023; Shi et al., 2024). Evaluation is conducted on  $\mathcal{D}^t$ , assessing the model’s ability to distinguish and correctly assign both seen and unseen categories.

### 3.2 Dataset and Diversity Analysis

To construct a comprehensive benchmark for OTC, BOLT aims to provide a unified evaluation framework that encompasses diverse text classification scenarios and domains. In total, BOLT includes 12 datasets spanning three major scenarios: topic classification, intent detection, and ontology classification, covering a wide range of text sources

such as news articles, conversational queries, and Wikipedia entries. The datasets include 20NG (Lang, 1995), THUCNews (Sun et al., 2016), Yahoo (Zhang et al., 2015a), BANK (i.e., BANKING) (Casanueva et al., 2020), CLINC (Larson et al., 2019), S.O. (i.e., StackOverflow) (Beau and Crabbé, 2024), HWU (Liu et al., 2021), TREC (Li and Roth, 2002), ECDT (Wang et al., 2022), MCID (Arora et al., 2020a), DBPedia (Zhang et al., 2015b), and X-Topic (Antypas et al., 2024). Detailed descriptions of these datasets are provided in the Appendix B.1.

We sample 10,000 instances from each large-scale corpus, including 20 Newsgroups (20NG), THUCNews, Yahoo, and DBPedia, to ensure both fair comparison and computational efficiency. To maintain fairness and consistency, all datasets are de-duplicated and standardized into a unified format. Table 2 summarizes the statistics of the datasets provided in BOLT.

**Diverse Scenarios.** BOLT encompasses diverse text classification scenarios under open-world settings. Specifically, it includes three representative scenarios: (1) *Topic Classification*, which focuses on identifying thematic categories from longer, content-rich documents (e.g., news articles); (2) *Intent Detection*, which emphasizes fine-grained recognition of user intentions from short, query-like utterances; and (3) *Ontology Classification*, which involves assigning entities or documents to high-cardinality, knowledge-based categories (e.g., DBPedia), thereby evaluating the model’s scalability and ability for semantic abstraction.

**Diverse Text Domains.** The benchmark encompasses multiple text domains, such as news articles, question–answer datasets, and knowledge

bases, thereby ensuring coverage of both formal and conversational contexts. It consequently exhibits substantial diversity across both linguistic and structural dimensions, covering datasets with varying text lengths (ranging from short commands to long articles, with an average of 8–1856 words) and label granularities (from 10 to 219 classes), across multiple languages such as English, Chinese, Japanese and Spanish. This diversity enables comprehensive evaluation and more faithfully reflects the complexity of real-world open-world settings.

### 3.3 BOLT Framework

The BOLT framework provides a Python-based, practical, and extensible evaluation toolkit for OTC. By standardizing dataset formats, experimental configurations, and evaluation metrics, it enables reproducible comparisons of a wide range of existing methods and facilitates the integration of new datasets and algorithms, thereby advancing research toward scalable open-world applications.

**Open-World Formulation** To emulate open-world conditions, the BOLT framework introduces a unified formulation strategy grounded in diverse real-world text classification datasets. Specifically, each dataset is restructured into controlled experimental settings by adjusting two key factors: the Known Class Ratio (KCR), which partitions existing categories into known and unknown subsets to simulate varying degrees of openness, and the Labeled Ratio (LAR), which specifies the proportion of labeled samples within the known classes to reflect different levels of supervision. By specifying appropriate KCR and LAR values, users can flexibly generate customized experimental configurations, ensuring reproducibility, comparability, and scalability across studies. This design enables consistent evaluation under diverse openness and supervision conditions. Furthermore, BOLT offers a unified dataset configuration mechanism, where users only need to provide training, validation, and test sets in a standardized CSV format. By specifying the desired KCR and LAR, BOLT automatically generates the corresponding label partitions, enabling flexible experimental settings.

**Evaluation Metrics** To ensure consistent and comparable evaluation, BOLT adopts a unified evaluation protocol across both OSTC and GCD tasks, while employing task-specific metric sets. For OSTC (Shu et al., 2017; Zhang et al., 2021a), commonly used metrics include ACC, F1, K-F1,

and N-F1, whereas GCD (An et al., 2023, 2025) is evaluated with ACC, K-ACC, N-ACC, H-Score, ARI, and NMI. These metrics jointly measure model performance in classification, OOD detection, and novel category discovery. Detailed descriptions of the evaluation metrics are provided in Appendix B.4. BOLT standardizes the evaluation pipeline across datasets, ensuring fair and consistent benchmarking of models under open-set conditions.

**Method Integration** BOLT provides a unified and extensible implementation framework that supports the integration, reproduction, and evaluation of diverse algorithmic paradigms across two representative OTC tasks, OSTC and GCD. The framework prioritizes compatibility and scalability over isolated algorithmic comparisons, facilitating the seamless incorporation of both established and newly developed LLM-based methods.

For OSTC, the framework incorporates four representative families of approaches: (i) Traditional non-PLM methods, including DOC (Shu et al., 2017) and DeepUnk (Lin and Xu, 2019), which rely on shallow neural architectures or activation-based confidence estimation; (ii) PLM-based methods, including ADB (Zhang et al., 2021a), SCL (Zeng et al., 2021a), AB (Lorenz et al., 2022), Kn-Con (Zhou et al., 2022), and DyEn (Zhou et al., 2023), which enhance discriminative representation learning through PLMs and adaptive decision boundaries; (iii) our newly implemented LLM-based training framework, which fine-tunes an LLM with an MLP head as a discriminative backbone and incorporates post-hoc methods for OOD detection, including LLM-MSP (Hendrycks and Gimpel, 2016), LLM-OpenMax (Bendale and Boulton, 2016), LLM-TempScale (Guo et al., 2017), LLM-Energy (Liu et al., 2020), LLM-LogitNorm (Wei et al., 2022), LLM-Entropy (Chan et al., 2021), LLM-KLMatching (Basart et al., 2022), and LLM-MaxLogit (Basart et al., 2022), supporting open-world classification.

For GCD, BOLT integrates two major families of approaches: (i) PLM-based clustering methods, including DeepAligned (Zhang et al., 2021b), GeoID (Tang et al., 2024), SDC (An et al., 2025), and DPN (An et al., 2023), TAN (An et al., 2024b), which learn discriminative embeddings from pre-trained encoders and discover novel categories through clustering; and (ii) LLM-enhanced methods, including LOOP (An et al., 2024a) and

Metric	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.	
K-F1	DOC	60.06	42.71	30.42	35.57	45.66	68.85	58.95	47.89	0.09	42.98	22.58	6.72	38.54	
	DeepUNK	73.65	27.15	36.90	25.83	53.96	52.42	64.70	63.18	1.59	54.40	24.78	6.95	40.46	
	ADB	78.67	<b>79.94</b>	<b>55.77</b>	48.30	68.58	<b>81.72</b>	82.60	71.49	66.93	60.69	<b>73.70</b>	18.88	65.61	
	SCL	79.30	23.57	21.41	29.88	35.14	51.58	57.55	41.22	1.81	46.47	30.18	8.81	35.58	
	AB	57.59	27.23	46.17	25.49	62.57	63.57	72.93	67.19	40.31	54.72	42.39	12.29	47.70	
	KNNCon	92.45	73.72	55.27	57.19	76.05	75.78	<b>85.63</b>	<b>79.91</b>	77.08	73.90	65.60	19.49	69.34	
	DyEn	<b>94.73</b>	73.54	53.21	<b>57.72</b>	<b>76.07</b>	71.22	<u>85.41</u>	<u>76.63</u>	<b>80.06</b>	<b>75.75</b>	<u>73.12</u>	<b>19.72</b>	<b>69.77</b>	
	LLM-OpenMax	75.32	64.13	47.91	30.38	46.23	74.09	60.52	50.78	33.10	42.69	25.35	10.90	46.79	
	LLM-Entropy	72.95	64.40	53.61	47.42	62.89	67.92	72.53	66.42	53.43	50.89	50.73	13.26	56.37	
	LLM-MaxLogit	84.18	69.85	43.28	37.46	65.89	<u>80.18</u>	81.16	63.39	42.84	46.39	50.58	8.02	56.10	
	N-F1	DOC	77.93	66.77	58.53	59.45	69.29	72.25	76.41	70.33	61.79	<b>68.20</b>	64.08	<b>54.72</b>	66.64
		DeepUNK	85.29	45.11	<u>59.97</u>	51.62	58.30	72.69	62.71	56.68	43.13	53.57	43.96	32.57	55.47
ADB		67.16	<b>78.06</b>	39.91	51.49	67.75	<b>83.03</b>	<b>80.48</b>	70.20	68.43	55.25	75.39	31.90	64.09	
SCL		<u>90.64</u>	52.88	45.63	58.04	53.21	62.59	62.69	51.78	3.93	58.35	49.87	37.23	52.24	
AB		63.74	45.89	55.25	10.94	67.73	69.30	74.36	<u>75.41</u>	59.01	<u>67.10</u>	62.23	50.31	58.44	
KNNCon		87.68	54.26	28.44	53.40	59.21	67.21	78.32	59.57	65.25	56.35	61.45	28.15	58.27	
DyEn		<b>93.71</b>	61.59	21.58	51.19	64.33	53.72	80.09	57.42	66.45	67.06	64.25	21.47	58.57	
LLM-OpenMax		77.26	65.26	21.85	<u>62.52</u>	<u>71.79</u>	70.79	77.68	72.63	67.87	56.69	68.74	<u>52.20</u>	63.77	
LLM-Entropy		48.67	38.05	46.58	43.94	37.91	39.47	45.90	45.30	56.13	50.09	52.64	31.69	44.70	
LLM-MaxLogit		88.85	<b>78.37</b>	<b>64.28</b>	<b>68.92</b>	<b>76.21</b>	<b>83.44</b>	<b>84.67</b>	<b>76.34</b>	<b>72.95</b>	66.86	<b>77.47</b>	50.27	<b>74.05</b>	

Table 3: Average performance on BOLT benchmark for OSTC under  $KCR \in \{0.25, 0.50, 0.75\}$  and  $LAR = 1.0$ . Metrics include K-F1 and N-F1. The best results are highlighted in **bold**, while the second-best results are underlined.

ALUP (Liang et al., 2024), which leverage LLM-generated semantic priors or pseudo-labels to improve clustering alignment.

Overall, BOLT enables consistent evaluation and flexible extension from traditional classifiers to next-generation LLM-based open-world learners.

## 4 Experimental Results and Analysis

### 4.1 Experimental Setup

To simulate open-world conditions, we adopt a unified experimental protocol. Specifically, we vary the KCR at  $\{25\%, 50\%, 75\%\}$  to control the proportion of known classes, and adjust the LAR at  $\{10\%, 50\%, 100\%\}$ . To ensure stable and comparable evaluation when varying KCR, we uniformly partition the entire label space into five predefined subsets with minimal overlap. The final results are reported as the average performance over these subsets. More dataset details are provided in Appendix B.2.

### 4.2 Overall Performance on the OSTC Task

Table 3 presents the overall performance of various methods on the OSTC task, which is the setting adopted by most previous studies. The results show that existing PLM-based methods achieve better K-F1 than N-F1, failing to generalize to unseen classes. LLM-based approaches alleviate this issue and yield more balanced performance, though their generalization to novel classes remains limited.

**Known-class performance (K-F1).** Across all KCRs, PLM-based methods such as ADB and KNNCon exhibit remarkable robustness in known-class recognition. ADB achieves the top K-

F1 on structured topic such as DBPedia (73.70). KNNCon also performs consistently, outperforming all others on intent detection datasets (CLINC, HWU). LLM-based variants do not consistently surpass strong PLM-based baselines on OSTC, suggesting that the potential of directly leveraging LLMs for OSTC to be further explored.

**Novel-class performance (N-F1).** For novel-class detection, a different pattern emerges. LLM-based scoring methods (LLM-Energy, LLM-MaxLogit) dominate across nearly all datasets and KCRs, achieving the highest N-F1. These improvements demonstrate that our proposed LLM-based methods significantly enhance open-set robustness and uncertainty estimation, showcasing promising potential for future open-world understanding tasks. Besides, the non-PLM methods **DOC** and **DeepUnk** exhibit poor K-F1 performance on the ECDD dataset due to limited supervision and the absence of prior semantic knowledge for Chinese text. As a result, their predicted probabilities are nearly uniform, causing most samples to be classified as OOD. Consequently, the K-F1 is low, accompanied by an artificially inflated N-F1.

### 4.3 Overall Performance on the GCD Task

Table 4 presents the overall performance of various methods on the GCD task. PLM-based clustering methods exhibit strong known-class discrimination on long-text datasets, whereas LLM-enhanced methods show superior generalization and novel categories, particularly in short-text or low-resource settings. However, they still tend to overfit the known classes, leading to suboptimal performance on novel categories.

Metric	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
K-ACC	DAL	92.51	83.81	57.86	49.84	70.35	75.72	87.95	72.77	63.18	74.62	80.76	21.16	69.21
	GeoID	98.34	89.72	64.02	55.57	83.01	84.18	92.53	84.17	71.06	75.57	86.79	26.56	75.96
	SDC	80.79	82.78	57.28	65.58	62.37	80.24	72.49	72.00	73.37	70.04	72.77	<b>44.66</b>	69.53
	DPN	95.86	91.72	63.32	80.80	82.16	88.24	89.46	80.71	<b>86.44</b>	81.83	87.31	35.92	80.32
	TAN	98.39	93.05	66.23	<b>82.54</b>	84.31	86.07	92.98	81.55	<u>85.04</u>	<u>83.21</u>	<b>88.83</b>	34.47	<b>81.39</b>
	LOOP	96.97	<b>94.47</b>	68.78	63.31	<b>87.41</b>	88.97	<b>95.87</b>	<u>86.35</u>	77.12	77.70	85.41	27.06	79.12
	ALUP	96.95	92.42	<b>70.64</b>	69.35	<u>86.44</u>	<u>88.59</u>	<u>94.53</u>	<u>85.38</u>	79.97	<b>84.64</b>	86.61	29.86	<u>80.45</u>
N-ACC	DAL	81.41	42.32	25.79	38.67	51.03	56.85	71.74	55.22	51.88	51.40	<b>76.22</b>	14.82	51.44
	GeoID	92.73	66.98	45.69	43.23	57.14	<u>75.55</u>	74.92	63.84	53.67	33.37	67.86	9.41	57.03
	SDC	84.69	<b>78.71</b>	<b>56.00</b>	<b>56.54</b>	58.97	<b>79.94</b>	68.53	<b>72.57</b>	<b>72.34</b>	58.49	69.26	<b>19.57</b>	<b>64.63</b>
	DPN	83.55	39.33	24.00	28.60	28.73	14.71	49.49	29.55	47.01	24.53	68.47	17.78	37.98
	TAN	<u>93.12</u>	57.83	31.09	39.92	36.38	31.31	53.39	34.57	53.03	30.31	72.20	19.05	46.02
	LOOP	87.18	66.62	51.52	<u>44.84</u>	<u>65.41</u>	78.44	<b>83.69</b>	<u>72.01</u>	62.98	58.09	74.24	17.61	<u>63.55</u>
	ALUP	<b>95.70</b>	<u>68.90</u>	<u>49.88</u>	43.53	<u>64.87</u>	75.41	<u>81.61</u>	68.65	<u>65.64</u>	53.76	<u>75.21</u>	15.61	63.23

Table 4: Average performance on BOLT benchmark for GCD across different KCRs and LARs. Metrics include K-ACC and N-ACC. The best results are highlighted in **bold**, while the second-best results are underlined.

Method	LAR=0.1		KCR=0.25 LAR=0.5		LAR=1.0		LAR=0.1		KCR=0.50 LAR=0.5		LAR=1.0		LAR=0.1		KCR=0.75 LAR=0.5		LAR=1.0		
	K-F1	N-F1	K-F1	N-F1	K-F1	N-F1	K-F1	N-F1	K-F1	N-F1	K-F1	N-F1	K-F1	N-F1	K-F1	N-F1	K-F1	N-F1	
OSTC	14.09	<b>81.13</b>	50.29	84.46	54.53	82.72	15.71	63.82	52.18	74.65	60.80	75.60	13.49	40.27	51.44	52.44	62.20	56.10	
DOC	28.33	42.02	46.95	82.71	53.99	82.48	31.19	41.05	49.28	69.99	54.49	72.56	33.28	28.91	47.49	51.63	51.60	54.13	
DeepUnk	45.11	58.99	64.15	73.99	67.20	77.92	59.26	57.05	<b>76.23</b>	71.08	<b>77.98</b>	73.16	64.92	47.18	80.65	59.85	<b>82.38</b>	61.71	
ADB	23.92	52.54	27.75	70.91	31.02	68.21	26.89	50.91	37.13	56.12	40.15	52.45	30.81	36.94	46.37	40.79	51.25	39.34	
SCL	42.73	73.33	47.83	71.57	48.31	73.60	44.43	59.63	52.18	57.94	53.44	59.63	45.44	41.19	54.50	42.07	56.16	42.00	
AB	<b>61.74</b>	61.68	59.62	58.32	60.40	61.39	<b>74.89</b>	62.77	<u>74.56</u>	62.32	<u>74.97</u>	62.40	<b>82.55</b>	<b>56.35</b>	<b>82.12</b>	53.87	53.87	<u>82.19</u>	56.01
KNNCon	18.61	65.71	50.68	69.06	58.53	71.67	25.41	65.50	62.08	70.54	70.89	70.08	26.03	43.81	64.58	52.76	74.78	56.98	
LLM-OpenMax	24.74	<u>80.43</u>	61.80	<b>89.41</b>	69.26	<b>90.42</b>	32.75	<u>70.41</u>	66.58	<u>81.50</u>	71.20	<u>83.36</u>	35.78	48.54	62.00	59.70	68.26	63.41	
LLM-Energy	26.32	79.71	<b>65.43</b>	<u>89.11</u>	<u>70.77</u>	<u>90.08</u>	37.55	<b>71.58</b>	70.69	<b>82.31</b>	74.42	<b>83.87</b>	41.52	<u>50.41</u>	67.07	<b>62.05</b>	72.50	<b>65.50</b>	
GCD	K-ACC	N-ACC	K-ACC	N-ACC	K-ACC	N-ACC	K-ACC	N-ACC	K-ACC	N-ACC	K-ACC	N-ACC	K-ACC	N-ACC	K-ACC	N-ACC	K-ACC	N-ACC	
DAL	64.18	46.62	68.01	48.10	67.19	48.16	69.24	53.13	74.51	56.10	73.64	55.00	71.57	53.98	76.09	58.70	76.82	55.76	
GeoID	73.53	<u>62.51</u>	82.81	64.40	85.84	<b>66.22</b>	75.30	<b>64.60</b>	84.70	67.12	87.10	67.15	75.24	<b>64.88</b>	85.13	66.67	88.44	<b>68.46</b>	
SDC	58.21	49.42	61.37	54.22	63.88	56.94	68.89	<u>64.21</u>	70.84	<b>68.66</b>	73.75	<b>70.24</b>	75.69	<b>71.94</b>	80.67	<b>76.43</b>	83.19	<b>77.18</b>	
DPN	74.91	34.12	81.70	33.16	82.69	32.68	<b>78.27</b>	36.34	83.49	37.36	86.20	37.97	<u>79.43</u>	41.12	86.38	40.36	88.42	39.21	
TAN	<b>79.58</b>	44.73	<b>84.80</b>	44.51	<b>87.56</b>	43.23	<b>80.15</b>	49.04	<b>87.29</b>	47.02	<u>88.32</u>	47.62	<b>81.50</b>	52.57	<b>88.21</b>	51.14	<b>89.99</b>	49.27	
LOOP	74.65	<b>63.24</b>	83.43	<u>64.86</u>	<u>86.02</u>	64.67	76.49	63.23	<b>84.55</b>	64.44	87.26	65.85	78.00	64.31	84.74	65.11	87.15	64.24	
ALUP	74.13	59.59	<u>83.46</u>	<b>66.41</b>	85.67	<b>67.79</b>	73.30	58.59	<u>86.68</u>	<u>67.49</u>	<b>88.73</b>	<u>67.74</u>	76.03	57.47	<u>87.09</u>	<u>68.00</u>	88.92	66.29	

Table 5: Average performance on BOLT benchmark for OSTC and GCD under different KCRs and LARs. Due to page limitations, the complete experimental results across all datasets and metrics are provided in Appendix E.

**Known-class performance (K-ACC).** PLM-based methods (TAN, DPN) achieve the highest K-ACC on long and information-rich corpora such as 20NG, DBPedia, and THUCNews. TAN consistently attains top performance, confirming its stability in structured topic classification tasks. In contrast, LLM-enhanced methods (LOOP, ALUP) yield superior results on short-text intent detection benchmarks (CLINC, BANKING, StackOverflow), where semantic priors help compensate for limited lexical cues.

**Novel-class performance (N-ACC).** For novel-class detection, a complementary trend emerges. LLM-enhanced models (ALUP, LOOP) outperform others indicating their stronger ability to infer unseen categories under weak supervision. Among PLM-based methods, SDC demonstrates relatively strong performance due to its distinctive calibration mechanism, suggesting that effective logit calibration is a promising strategy for GCD, particularly evident when pseudo-label noise is non-negligible.

#### 4.4 Impact of Known Class Ratio (KCR)

Table 5 summarizes the average performance under different KCRs across OSTC and GCD bench-

marks. On OSTC, increasing KCR generally improves K-F1, but poses a clear challenge to open-set recognition: most methods exhibit a noticeable drop in N-F1 as KCR increases, indicating a stronger closed-set bias when more classes are treated as known. This challenge stems from the K+1 classification formulation of OSTC, where unknown samples must be explicitly rejected from an expanding known-class set without additional supervision. On GCD, however, the negative impact of higher KCR on unknown-class performance is far less pronounced, as the availability of extra unlabeled data enables K+N pseudo-label learning that jointly refines known and novel categories.

#### 4.5 Impact of Labeled Ratio (LAR)

As shown in Table 5, increasing LAR consistently enhances performance on the known (K-F1/K-ACC) for PLM-based discriminative learners, as richer supervision sharpens decision boundaries and yields more separable known-class representations. In contrast, *open-set* robustness (N-F1/N-ACC) is less sensitive to labeled data volume, where LLM-based or LLM-enhanced methods (e.g., LLM-Energy, LLM-MaxLogit, ALUP)

Metric	LLM	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
K-F1	Qwen2.5-3B	79.68	66.35	41.37	42.25	57.42	70.15	76.15	54.31	39.28	36.50	54.72	16.49	52.89
	Qwen2.5-7B	71.37	56.77	30.42	45.51	58.04	74.25	73.68	52.30	38.79	39.89	42.11	11.67	49.57
	Qwen2.5-32B	71.90	63.37	39.32	42.34	57.14	76.05	78.26	61.57	39.50	45.96	35.45	11.03	51.82
N-F1	Qwen2.5-3B	88.53	76.46	62.45	48.87	73.54	81.01	82.61	76.47	56.80	64.46	64.84	11.11	65.60
	Qwen2.5-7B	82.07	74.76	58.79	49.88	74.87	82.30	79.77	73.76	54.01	64.64	59.96	9.80	63.72
	Qwen2.5-32B	83.80	78.40	63.84	50.14	73.69	84.10	82.82	78.12	55.87	69.18	59.37	11.07	65.86

Table 6: Performance of LLM-MaxLogit under the OSTC scenario with different Qwen2.5 model scales.

Setting	Method	Random		Similarity	
KCR, LAR	GCD	K-ACC	N-ACC	K-ACC	N-ACC
0.25, 0.1	DAL	64.18	46.62	66.67	43.24
	GeoID	73.53	62.51	73.80	<u>64.08</u>
	SDC	58.21	49.42	<b>79.95</b>	47.64
	DPN	74.91	34.12	76.75	37.89
	TAN	<b>79.58</b>	44.73	73.63	10.20
	LOOP	74.65	<u>63.24</u>	75.70	61.66
	ALUP	74.13	59.59	<u>79.94</u>	<b>68.42</b>
KCR, LAR	OSTC	K-F1	N-F1	K-F1	N-F1
0.25, 1.0	DOC	54.53	82.72	56.54	83.73
	DeepUnk	53.99	82.48	53.48	83.42
	ADB	67.20	77.92	69.59	79.98
	SCL	31.02	68.21	36.94	70.65
	AB	48.31	73.60	50.52	76.37
	KNNCon	60.40	61.39	63.17	65.86
	DyEn	63.88	65.49	68.98	71.90
	LLM-OpenMax	58.53	71.67	61.12	75.38
	LLM-Energy	<u>69.26</u>	<b>90.42</b>	<u>70.38</u>	<b>91.39</b>
LLM-MaxLogit	<b>70.77</b>	<u>90.08</u>	<b>72.73</b>	<u>91.21</u>	

Table 7: Average performance on BOLT benchmark for OTC under different label sampling strategies.

maintain superior performance even in low-label regimes. These results suggest that while PLM-based training benefits from abundant labels, LLM-driven calibration and pseudo-supervision offer a complementary pathway toward generalization.

#### 4.6 Impact of LLM Backbone

We evaluate LLM-Maxlogit with Qwen2.5 backbones of different scales (3B/7B/32B) under OSTC (Table 6). Overall, scaling does not yield monotonic improvements: 32B gives the best average OOD-aware performance (Avg. N-F1: 65.86; Avg. K-F1: 51.82), while 3B is highly competitive and even better on average K-F1 (52.89). The 7B model underperforms both in most cases, indicating that the mid-scale backbone does not offer a consistently more balanced trade-off between K-F1 and N-F1 in our OSTC setting. Dataset-wise, longer documents with clearer topical evidence (e.g., 20NG, THUCNews; Table 2) tend to benefit more from larger models on OOD detection, whereas extremely short and distribution-shifted data (X-Topic) remains challenging for all scales.

From a task perspective, larger models help more

for OOD detection than for closed-set classification on several intent benchmarks (e.g., StackOverflow, MCID), consistent with improved semantic coverage for separating in- vs. out-of-domain. However, for fine-grained label spaces with minimal label descriptions (e.g., DBPedia with 219 classes and Avg. Len./Lab.  $\approx 1$ ), scaling can even degrade both K-F1 and N-F1, implying that the bottleneck is label-space ambiguity rather than representation capacity. These results highlight that simply increasing model size is insufficient in OSTC.

#### 4.7 Impact of Imbalanced Sampling

To further analyze the impact of known label distribution, we compare two label sampling strategies: (1) **Random**, where known and unknown classes are randomly selected without considering their semantic relatedness; and (2) **Similarity**, where labels are grouped based on semantic proximity, so that known classes share closer meanings while unknown classes differ more significantly. This Similarity sampling better reflects real-world open-world scenarios, where initial known categories are often semantically coherent (e.g., within a specific domain). As shown in Table 7, semantically concentrated splits consistently outperform random ones, yielding clearer intra-class discrimination and more accurate OOD identification. These results suggest that OTC methods tend to deliver better performance in practical applications.

## 5 Conclusions

In this work, we presented **BOLT**, the first comprehensive benchmark and evaluation toolkit designed to support open-world learning for text classification, unifying the evaluation of two key tasks in OTC, namely OSTC and GCD. Compared with prior studies that focus narrowly on single tasks or limited domains, BOLT provides: (1) broader coverage, encompassing 12 diverse datasets across multiple text classification scenarios and domains, (2) a standardized evaluation protocol that explicitly accounts for varying known class ratios and

labeled ratios, and (3) systematic comparisons across traditional methods, PLM-based approaches, and LLM-enhanced models. The experimental results not only reveal the limitations of existing approaches but also demonstrate the great potential of leveraging LLMs to advance OTC.

## Limitations

While BOLT provides a comprehensive benchmark and evaluation toolkit for OTC, several promising directions for future work remain. One promising direction is to investigate the impact of different LLM backbones on performance, which could yield deeper insights into the trade-offs between efficiency and generalization. Moreover, incorporating additional multilingual datasets would further enhance the coverage and cross-lingual applicability of BOLT, making it a more versatile benchmark for global open-world learning research.

## Ethical Considerations

We present BOLT, a comprehensive benchmark and evaluation toolkit for OTC. This research was conducted in full compliance with the laws, rules, and regulations of our community, institution, workplace, and country. During the development of BOLT, we carefully considered multiple ethical aspects. All datasets used are publicly available and widely adopted in prior research; their licenses and terms of use were thoroughly reviewed to ensure compliance. No personally identifiable or sensitive information was collected, modified, or redistributed, and all preprocessing steps were performed with care to preserve the integrity and semantics of the original text. BOLT is intended solely for academic and research purposes to advance the understanding of OTC. To promote transparency and reproducibility, all benchmark configurations and experimental settings will be made publicly available, in line with the principles of open and responsible NLP research.

## Acknowledgments

This work was supported in part by the National Natural Science Foundation of China (NSFC) (Grant No. 62506352), and in part by the Strategic Priority Research Program of Chinese Academy of Sciences (Grant No. XDB1350102).

## References

- Wenbin An, Haonan Lin, Jiahao Nie, Feng Tian, Wenkai Shi, Yaqiang Wu, Qianying Wang, and Ping Chen. 2025. Unleashing the potential of model bias for generalized category discovery. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Wenbin An, Wenkai Shi, Feng Tian, Haonan Lin, Qianying Wang, Yaqiang Wu, Mingxiang Cai, Luyan Wang, Yan Chen, Haiping Zhu, and Ping Chen. 2024a. Generalized category discovery with large language models in the loop. In *Findings of the Association for Computational Linguistics: ACL 2024*.
- Wenbin An, Feng Tian, Wenkai Shi, Yan Chen, Yaqiang Wu, Qianying Wang, and Ping Chen. 2024b. Transfer and alignment network for generalized category discovery. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Wenbin An, Feng Tian, Qinghua Zheng, Wei Ding, Qianying Wang, and Ping Chen. 2023. Generalized category discovery with decoupled prototypical network. In *Proceedings of the AAAI conference on artificial intelligence*.
- Dimosthenis Antypas, Asahi Ushio, Francesco Barbieri, and Jose Camacho-Collados. 2024. Multilingual topic classification in x: Dataset and analysis. *arXiv preprint arXiv:2410.03075*.
- Abhinav Arora, Akshat Shrivastava, Mrinal Mohit, Lorena Sainz-Maza Lecanda, and Ahmed Aly. 2020a. [Cross-lingual transfer learning for intent detection of covid-19 utterances](#).
- Abhinav Arora, Akshat Shrivastava, Mrinal Mohit, Lorena Sainz-Maza Lecanda, and Ahmed Aly. 2020b. Cross-lingual transfer learning for intent detection of covid-19 utterances. In *The EMNLP 2020 Workshop NLP-COVID Submission*.
- Steven Basart, Mazeika Mantas, Mostajabi Mohamadreza, Steinhardt Jacob, and Song Dawn. 2022. Scaling out-of-distribution detection for real-world settings. In *International Conference on Machine Learning*.
- Nathanaël Beau and Benoit Crabbé. 2024. [CodeInsight: A curated dataset of practical coding solutions from Stack Overflow](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 5935–5947, Bangkok, Thailand. Association for Computational Linguistics.
- Abhijit Bendale and Terrance E. Boult. 2016. Towards open set deep networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. [Efficient intent detection with dual sentence encoders](#). In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 38–45, Online. Association for Computational Linguistics.

- Robin Chan, Matthias Rottmann, and Hanno Gottschalk. 2021. Entropy maximization and meta classification for out-of-distribution detection in semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Junyang Chen, Zhiguo Gong, and Weiwen Liu. 2020. A dirichlet process biterm-based mixture model for short text stream clustering. *Applied Intelligence*.
- Xi Chen, Chuan Qin, Chuyu Fang, Chao Wang, Chen Zhu, Fuzhen Zhuang, Hengshu Zhu, and Hui Xiong. 2024. Job-sdf: A multi-granularity dataset for job skill demand forecasting and benchmarking. *Advances in Neural Information Processing Systems*, 37:129329–129356.
- Xi Chen, Chuan Qin, Ziqi Wang, Shasha Hu, Chao Wang, Hengshu Zhu, and Hui Xiong. 2026. Beyond the known: An unknown-aware large language model for open-set text classification. In *The Fourteenth International Conference on Learning Representations*.
- Geli Fei and Bing Liu. 2016. Breaking the closed world assumption in text classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017. On calibration of modern neural networks. In *International conference on machine learning*. PMLR.
- Dan Hendrycks and Kevin Gimpel. 2016. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations*.
- Xiaohan Huang, Meng Xiao, Chuan Qin, Qingqing Long, Jinmiao Chen, Yuanchun Zhou, and Hengshu Zhu. 2026. Scihorizon-gene: Benchmarking llm for life sciences inference from gene knowledge to functional understanding. *arXiv preprint arXiv:2601.12805*.
- Feihu Jiang, Chuan Qin, Kaichun Yao, Chuyu Fang, Fuzhen Zhuang, Hengshu Zhu, and Hui Xiong. 2024. Enhancing question answering for enterprise knowledge bases using large language models. In *International Conference on Database Systems for Advanced Applications*, pages 273–290. Springer.
- Konstantin Kirchheim, Marco Filax, and Frank Ortmeier. 2022. Pytorch-ood: A library for out-of-distribution detection based on pytorch. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.
- Harold W Kuhn. 2004. The hungarian method for the assignment problem. *Naval Research Logistics (NRL)*.
- Ken Lang. 1995. Newsweeder: Learning to filter news. In *Proceedings of the Twelfth International Conference on Machine Learning*.
- Stefan Larson, Anish Mahendran, Joseph J. Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael A. Laurenzano, Lingjia Tang, and Jason Mars. 2019. An evaluation dataset for intent classification and out-of-scope prediction. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1311–1316, Hong Kong, China. Association for Computational Linguistics.
- Fayin Li and Harry Wechsler. 2005. Open set face recognition using transduction. *IEEE transactions on pattern analysis and machine intelligence*.
- Xin Li and Dan Roth. 2002. *Learning question classifiers*. In *COLING 2002: The 19th International Conference on Computational Linguistics*.
- Jingui Liang, Lizi Liao, Hao Fei, Bobo Li, and Jing Jiang. 2024. Actively learn from LLMs with uncertainty propagation for generalized category discovery. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Ting-En Lin and Hua Xu. 2019. Deep unknown intent detection with margin loss. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- Bo Liu, Liming Zhan, Zexin Lu, Yujie Feng, Lei Xue, and Xiao-Ming Wu. 2023. How good are llms at out-of-distribution detection? *arXiv preprint arXiv:2308.10261*.
- Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. 2020. Energy-based out-of-distribution detection. *Advances in neural information processing systems*.
- Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and Verena Rieser. 2021. Benchmarking natural language understanding services for building conversational agents. In *Increasing naturalness and flexibility in spoken dialogue interaction: 10th international workshop on spoken dialogue systems*, pages 165–183. Springer.
- Petr Lorenc, Tommaso Gargiani, Jan Pichl, Jakub Konrád, Petr Marek, Ondřej Kobza, and Jan Šedivý. 2022. Metric learning and adaptive boundary for out-of-domain detection. In *International Conference on Applications of Natural Language to Information Systems*. Springer.
- Shervin Minaee, Nal Kalchbrenner, Erik Cambria, Narjes Nikzad, Meysam Chenaghlu, and Jianfeng Gao. 2021. Deep learning-based text classification: a comprehensive review. *ACM computing surveys (CSUR)*.
- Sridhama Prakhya, Vinodini Venkataram, and Jugal Kalita. 2017. Open set text classification using cnns. In *Proceedings of the 14th International Conference on Natural Language Processing (ICON-2017)*.

- Chuan Qin, Xin Chen, Chengrui Wang, Pengmin Wu, Xi Chen, Yihang Cheng, Jingyi Zhao, Meng Xiao, Xi-angchao Dong, Qingqing Long, and 1 others. 2025a. Scihorizon: Benchmarking ai-for-science readiness from scientific data to large language models. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pages 5754–5765.
- Chuan Qin, Chuyu Fang, Kaichun Yao, Xi Chen, Fuzhen Zhuang, and Hengshu Zhu. 2025b. Cotr: Efficient job task recognition for occupational information systems with class-incremental learning. *ACM Transactions on Management Information Systems*, 16(2):1–30.
- Chuan Qin, Le Zhang, Yihang Cheng, Rui Zha, Dazhong Shen, Qi Zhang, Xi Chen, Ying Sun, Chen Zhu, Hengshu Zhu, and 1 others. 2025c. A comprehensive survey of artificial intelligence techniques for talent analytics. *Proceedings of the IEEE*.
- Wenkai Shi, Wenbin An, Feng Tian, Yan Chen, Yaqiang Wu, Qianying Wang, and Ping Chen. 2024. A unified knowledge transfer network for generalized category discovery. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Lei Shu, Hu Xu, and Bing Liu. 2017. Doc: Deep open classification of text documents. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*.
- Zhiheng Song, Jingshuai Zhang, Chuan Qin, Chao Wang, Chao Chen, Longfei Xu, Kaikui Liu, Xiangxiang Chu, and Hengshu Zhu. 2026. Mobilitybench: A benchmark for evaluating route-planning agents in real-world mobility scenarios. *arXiv preprint arXiv:2602.22638*.
- Maosong Sun, Jingyang Li, Zhipeng Guo, Zhao Yu, Yabing Zheng, Xiance Si, and Zhiyuan Liu. 2016. [Thuctc: An efficient chinese text classifier](#).
- Kai Tang, Junbo Zhao, Xiao Ding, Runze Wu, Lei Feng, Gang Chen, and Haobo Wang. 2024. Learning geometry-aware representations for new intent discovery. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*.
- Zhenyu Tong, Chuan Qin, Chuyu Fang, Kaichun Yao, Xi Chen, Jingshuai Zhang, Chen Zhu, and Hengshu Zhu. 2025. From missteps to mastery: Enhancing low-resource dense retrieval through adaptive query generation. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 1*, pages 1373–1384.
- Fengjun Wang, Moran Beladev, Ofri Kleinfeld, Elina Frayerman, Tal Shachar, Eran Fainman, Karen Lastmann Assaraf, Sarai Mizrachi, and Benjamin Wang. 2023. Text2topic: Multi-label text classification system for efficient topic detection in user generated content with zero-shot capabilities. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*.
- Lihua Wang, Huiting Yang, Feng Li, Wenzhong Yang, and Zhenwan Zou. 2022. Intent detection model based on dual-channel feature fusion. In *2022 IEEE 6th Information Technology and Mechatronics Engineering Conference (ITOEC)*, volume 6, pages 1862–1867. IEEE.
- Mayur Wankhade, Annavarapu Chandra Sekhara Rao, and Chaitanya Kulkarni. 2022. A survey on sentiment analysis methods, applications, and challenges. *Artificial Intelligence Review*.
- Hongxin Wei, Renchunzi Xie, Hao Cheng, Lei Feng, Bo An, and Yixuan Li. 2022. Mitigating neural network overconfidence with logit normalization. In *International conference on machine learning*. PMLR.
- Zhiyuan Zeng, Keqing He, Yuanmeng Yan, Zijun Liu, Yanan Wu, Hong Xu, Huixing Jiang, and Weiran Xu. 2021a. Modeling discriminative representations for out-of-domain detection with supervised contrastive learning. *arXiv preprint arXiv:2105.14289*.
- Zhiyuan Zeng, Keqing He, Yuanmeng Yan, Zijun Liu, Yanan Wu, Hong Xu, Huixing Jiang, and Weiran Xu. 2021b. Modeling discriminative representations for out-of-domain detection with supervised contrastive learning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*.
- Feng Zhang, Wei Chen, Fei Ding, Meng Gao, Tengjiao Wang, Jiahui Yao, and Jiabin Zheng. 2024a. From discrimination to generation: Low-resource intent detection with language model instruction tuning. In *Findings of the Association for Computational Linguistics ACL 2024*.
- Hanlei Zhang, Hua Xu, and Ting-En Lin. 2021a. Deep open intent classification with adaptive decision boundary. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Hanlei Zhang, Hua Xu, Ting-En Lin, and Rui Lyu. 2021b. Discovering new intents with deep aligned clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Hanlei Zhang, Hua Xu, Xin Wang, Fei Long, and Kai Gao. 2023. A clustering framework for unsupervised and semi-supervised new intent discovery. *IEEE Transactions on Knowledge and Data Engineering*.
- Shun Zhang, Chaoran Yan, Jian Yang, Changyu Ren, Ji-qi Bai, Tongliang Li, and Zhoujun Li. 2024b. Ronid: new intent discovery with generated-reliable labels and cluster-friendly representations. In *International Conference on Database Systems for Advanced Applications*.
- Wei-Nan Zhang, Zhigang Chen, Wanxiang Che, Guoping Hu, and Ting Liu. 2017. The first evaluation of chinese human-computer dialogue technology. *arXiv preprint arXiv:1709.10217*.

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015a. [Character-level convolutional networks for text classification](#). In *Advances in Neural Information Processing Systems*.

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015b. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.

Qi Zhou, Xi Chen, Chuyu Fang, Jianji Wang, Chuan Qin, and Fuzhen Zhuang. 2025. Enhancing dual-target cross-domain recommendation via similar user bridging. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management*, pages 4487–4497.

Yunhua Zhou, Peiju Liu, and Xipeng Qiu. 2022. Knn-contrastive learning for out-of-domain intent classification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

Yunhua Zhou, Jianqiang Yang, Pengyu Wang, and Xipeng Qiu. 2023. Two birds one stone: Dynamic ensemble for ood intent classification. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

## A LLM Usage Statement

LLMs were used exclusively to assist in refining the writing of this paper, including grammar correction, phrasing improvement, and enhancement of readability and clarity. It should be noted that LLMs were not involved in the ideation, research methodological, experimental design, or result analysis. All research ideas, findings, and conclusions originate solely from the authors. The authors affirm that the use of LLMs complies with ethical guidelines and does not contribute to plagiarism or any form of scientific misconduct.

## B Additional Details on BOLT

### B.1 Dataset Details

This section presents detailed descriptions of the 12 datasets incorporated in BOLT. These datasets span diverse text classification scenarios and domains, providing comprehensive coverage for the evaluation of open-world text classification systems.

- **20NG**<sup>1</sup>: The 20NG, also known as the 20 Newsgroups dataset, contains approximately 20,000 newsgroup documents, partitioned nearly evenly across 20 categories. We sampled 10,000 documents from it to construct our dataset.

- **THUCNews**<sup>2</sup>: The THUCNews dataset is a large-scale Chinese news classification corpus, based on historical news articles collected via Sina News RSS feeds between 2005 and 2011. We sampled 10,000 documents from it to construct our version of the dataset.

- **Yahoo**<sup>3</sup>: The Yahoo dataset, also known as the Yahoo! Answers topic classification dataset, is constructed from user-submitted questions and their best answers under the Yahoo! Webscope program, and it is categorized into 10 major topic classes. For our study, we sampled 10,000 documents from this dataset to build our version.

- **BANK (BANKING)**<sup>4</sup>: The BANK, also known as the BANKING77, consists of 13,083 English customer service queries from the BANK domain, annotated with 77 fine-grained intent classes.

- **S.O. (StackOverflow)**<sup>5</sup>: The S.O. (i.e., StackOverflow) dataset consists of 20,000 short text titles collected from the Stack Overflow question forum, each annotated with one of 20 technical topic labels.

- **CLINC**<sup>6</sup>: The CLINC dataset, also known as CLINC150, is an intent classification corpus covering 150 in-domain intent classes. These intents cover a broad spectrum of everyday tasks and informational queries, such as playing music, checking the weather, and setting alarms.

- **HWU**<sup>7</sup>: The HWU dataset, also known as HWU-64, contains English user utterances across multiple domains (e.g. news and IoT) for intent detection, covering 64 intents. It is widely adopted

<sup>1</sup>20NG, <https://people.csail.mit.edu/jrennie/20Newsgroups/>

<sup>2</sup>THUCNews, <http://thuctc.thunlp.org/>

<sup>3</sup>Yahoo, <https://www.kaggle.com/datasets/bhavigardeshna/yahoo-email-classification/data>

<sup>4</sup>BANKING, <https://huggingface.co/datasets/PolyAI/banking77>

<sup>5</sup>StackOverflow, <https://github.com/jacoxu/StackOverflow>

<sup>6</sup>CLINC, <https://archive.ics.uci.edu/dataset/570/clinc150>

<sup>7</sup>HWU, <https://service.tib.eu/ldmservice/en/dataset/hwu64>

in conversational AI research, as its utterances simulate realistic task-oriented conversational queries.

- **ECDT**<sup>8</sup>: The ECDT dataset, also known as SMP-ECDT, is a Chinese intent classification corpus introduced in (Zhang et al., 2017).
- **TREC**<sup>9</sup>: The TREC is a question classification corpus of open-domain, fact-based questions.
- **MCID**<sup>10</sup>: The MCID dataset contains 6,871 utterances across multiple languages (English, Spanish, French, German, and Spanglish) with 16 COVID-19 related intents (Arora et al., 2020b). In our work, we select only the English utterances to form our version of the dataset.
- **DBpedia**<sup>11</sup>: DBpedia is a project that extracts structured content from Wikipedia. This dataset is an ontology classification corpus constructed by selecting 219 non-overlapping classes. For our study, we sampled 10,000 documents from this dataset to build our version.
- **X-Topic**<sup>12</sup>: X-Topic is a topic classification dataset based on X (formerly Twitter), featuring 19 topic labels. The classification task is multi-label, with tweets available in four languages: English, Japanese, Spanish, and Greek.

Moreover, to justify dataset selection, we conducted a systematic survey of 50 representative OTC-related papers and collected dataset usage statistics. The aggregated distribution is shown in Table S1. The results indicate that all high-frequency datasets adopted in prior OTC studies are included in BOLT, confirming that our benchmark covers the core datasets used by the community.

## B.2 Dataset Split

To ensure consistent and reproducible evaluation across all tasks, we adopt a unified five-fold dataset partitioning strategy. Unlike prior studies that rely on random class splits, which often introduce instability and overlap between folds, our benchmark

<sup>8</sup>ECDT, <https://github.com/ygwpz/SMP2019-ECDT-NLU>

<sup>9</sup>TREC, <https://huggingface.co/datasets/CogComp/trec>

<sup>10</sup>MCID, <https://github.com/fanolabs/IntentBert>

<sup>11</sup>DBpedia, [https://huggingface.co/datasets/DeveloperOats/DBpedia\\_Classes](https://huggingface.co/datasets/DeveloperOats/DBpedia_Classes)

<sup>12</sup>X-Topic, [https://huggingface.co/datasets/cardiffnlp/tweet\\_topic\\_multilingual](https://huggingface.co/datasets/cardiffnlp/tweet_topic_multilingual)

Dataset	# OTC Papers Using It	Percentage (%)
CLINC	44	88
BANKING	41	82
StackOverflow	26	52
HWU	12	24
MCID	3	6
20NG	1	2

Table S1: Dataset usage statistics from 50 OTC-related papers.

employs a **uniform label partitioning** scheme that minimizes class overlap while maintaining semantic balance. Specifically, the entire label space of each dataset is evenly divided into five disjoint subsets. Each fold serves as the known-class set once, while the remaining labels are treated as unknown classes for testing. This design ensures comparability across experiments and reduces randomness caused by arbitrary label sampling.

Within each known-class subset, samples are further divided according to the **Labeled Data Ratio (LAR)** and **Known Class Ratio (KCR)**:

- KCR (0.25, 0.5, 0.75) controls the proportion of known classes relative to the total label set, simulating different levels of open-world exposure.
- LAR (0.1, 0.5, 1.0) determines the proportion of labeled samples among the known-class instances, reflecting different supervision levels.

For each configuration, we maintain a fixed **train/validation/test** split within the known-class set to ensure balanced learning and reliable evaluation. The unlabeled pool includes both known and unknown instances, allowing open-set and generalized category discovery tasks to share the same structural foundation.

As summarized in Table S2, all datasets are processed under identical sampling rules, ensuring that model performance reflects genuine open-world generalization rather than artifacts of random data selection. This systematic partitioning provides a stable and comparable basis for evaluating BOLT and all baseline methods across OSTC and GCD settings.

## B.3 Method Details

This section provides detailed descriptions of the seven methods for GCD and the fifteen methods for OSTC implemented and evaluated in the BOLT benchmark. All methods are reproduced

KCR		0.25			0.50			0.75			ALL			
LAR	Dataset	#label	!Train	!Evall	#label	!Train	!Evall	#label	!Train	!Evall	#label	!Train	!Evall	!Test!
0.1	20NG	5	183	25	10	351	48	15	518	71	20	701	96	2,000
	THUCNews	4	200	28	7	350	49	10	500	70	14	700	98	2,000
	Yahoo	2	140	20	5	350	50	8	560	80	10	700	100	2,000
	TREC	12	324	35	24	435	51	35	468	62	47	487	71	490
	BANK	19	229	28	38	472	55	58	704	81	77	900	105	3,080
	S.O.	5	300	50	10	600	100	15	900	150	20	1,200	200	5,991
	CLINC	38	457	76	75	901	150	112	1,345	224	150	1,801	300	2,250
	HWU	16	187	26	32	361	51	48	557	79	64	770	109	1,032
	ECDT	8	97	12	16	137	20	23	188	28	31	209	33	770
	MCID	4	31	4	8	62	8	12	94	12	16	122	16	331
	DBPedia	55	252	55	110	480	110	164	647	164	219	736	218	2,000
X-Topic	12	531	75	24	582	87	37	598	100	49	610	112	1,767	
0.5	20NG	5	911	130	10	1,750	249	15	2,581	366	20	3,498	496	2,000
	THUCNews	4	1,000	144	7	1,750	252	10	2,500	360	14	3,500	504	2,000
	Yahoo	2	700	100	5	1,750	250	8	2,800	400	10	3,500	500	2,000
	TREC	12	1,610	175	24	2,176	239	35	2,343	261	47	2,426	273	490
	BANK	19	1,150	130	38	2,368	265	58	3,522	391	77	4,498	499	3,080
	S.O.	5	1,500	250	10	3,000	500	15	4,500	750	20	6,000	1,000	5,991
	CLINC	38	2,280	304	75	4,500	600	112	6,722	896	150	9,002	1,200	2,250
	HWU	16	943	110	32	1,808	214	48	2,784	328	64	3,853	451	1,032
	ECDT	8	487	48	16	684	68	23	938	92	31	1,046	103	770
	MCID	4	148	23	8	298	45	12	450	67	16	595	90	331
	DBPedia	55	1,204	177	110	2,315	341	164	3,111	463	219	3,517	537	2,000
X-Topic	12	2,657	378	24	2,917	415	37	2,970	432	49	2,982	444	1,767	
1.0	20NG	5	1,825	261	10	3,505	501	15	5,166	737	20	7,000	1,000	2,000
	THUCNews	4	2,000	287	7	3,500	502	10	5,000	716	14	7,000	1,000	2,000
	Yahoo	2	1,400	200	5	3,500	500	8	5,600	800	10	7,000	1,000	2,000
	TREC	12	3,224	350	24	4,350	475	35	4,683	513	47	4,849	532	490
	BANK	19	2,303	258	38	4,740	527	58	7,048	783	77	9,003	1,000	3,080
	S.O.	5	2,999	499	10	5,997	998	15	8,997	1,498	20	11,996	1,998	5,991
	CLINC	38	4,560	570	75	9,000	1,125	112	13,440	1,680	150	18,000	2,250	2,250
	HWU	16	1,886	227	32	3,616	436	48	5,575	674	64	7,712	933	1,032
	ECDT	8	976	92	16	1,372	130	23	1,879	180	31	2,099	200	770
	MCID	4	298	43	8	601	85	12	905	127	16	1,198	170	331
	DBPedia	55	2,387	345	110	4,603	661	164	6,195	886	219	7,000	1,000	2,000
X-Topic	12	5,318	756	24	5,836	825	37	5,937	849	49	5,954	866	1,767	

Table S2: **Detailed dataset split statistics** in the BOLT benchmark. Each dataset is partitioned under three Known Class Ratios (KCR = 0.25/0.5/0.75) and three Labeled Data Ratios (LAR = 0.1/0.5/1.0). For each configuration, we report the number of labels, training, evaluation, and test samples. This unified partitioning scheme ensures minimal class overlap across folds and consistent evaluation across both OSTC and GCD tasks.

within a unified framework to ensure fair comparison, using consistent training–evaluation protocols and hyperparameter settings.

#### OSTC Methods. *Non-PLM-based Methods:*

- **DOC** (Shu et al., 2017): A classical open-intent detection model that replaces the softmax layer with one-vs-rest sigmoid classifiers to reduce open-space risk, enabling explicit rejection of unknown intents.
- **DeepUnk** (Lin and Xu, 2019): Introduces a large-margin cosine loss to enhance separation between known and unknown samples, yielding more discriminative intent representations.

#### OSTC Methods. *PLM-based Methods.*

- **ADB** (Zhang et al., 2021a): Employs an adaptive spherical margin loss to dynamically adjust class decision boundaries in the pre-trained language model (PLM) feature space.
- **SCL** (Zeng et al., 2021a): Utilizes supervised contrastive learning to optimize intra-class compactness and inter-class separability.
- **AB** (Lorenc et al., 2022): Enhances metric-based classification through angular boundary scaling, improving generalization to unseen categories.
- **KNNCon** (Zhou et al., 2022): Combines KNN-based neighborhood contrastive learning with outlier detection using the Local Outlier Factor (LOF).

- **DyEn** (Zhou et al., 2023): Mitigates PLM “overthinking” through dynamic ensembling and confidence-based early exits for efficient OOD rejection.

#### **OSTC Methods.** *Proposed LLM-based Methods.*

Since no LLM-based methods have been explored for OSTC to date, we further propose a series of LLM fine-tuning frameworks, each incorporating a different post-hoc OOD detection approach, to investigate the potential of modern LLMs in addressing OSTC. Specifically, we employ eight different OOD detection methods, described as follows:

- **MSP** (Hendrycks and Gimpel, 2016): Uses the maximum softmax probability as a simple yet effective post-hoc confidence score for OOD detection.
- **LLM-OpenMax** (Bendale and Boult, 2016): Fits Weibull distributions to logits to estimate OOD probabilities, refining uncertainty estimation beyond MSP.
- **TempScale** (Guo et al., 2017): Applies temperature scaling to logits for post-hoc calibration and more reliable confidence estimation.
- **EnergyBased** (Liu et al., 2020): Reformulates confidence estimation as an energy-based scoring function, offering a theoretically grounded OOD detection approach.
- **LogitNorm** (Wei et al., 2022): Normalizes logits during training to mitigate overconfidence and improve model calibration.
- **Entropy** (Chan et al., 2021): Employs predictive LLM-Entropy as a task-agnostic uncertainty measure for identifying unknown samples.
- **KLMatching** (Basart et al., 2022): Compares class-wise probability distributions using KL divergence for fine-grained OOD differentiation.
- **MaxLogit** (Basart et al., 2022): Scores samples based on their maximum logit values, achieving strong empirical performance with minimal computational overhead.

#### **GCD Methods.** *PLM-based Methods:*

- **DAL** (Zhang et al., 2021b): A strong GCD baseline employing alignment-based pseudo-labeling and feature consistency constraints for novel class discovery.

- **GeoID** (Tang et al., 2024): Enhances GCD via geometry-aware clustering, calibrating pseudo-labels through similarity-based distribution alignment.

- **SDC** (An et al., 2025): Introduces self-debiasing calibration to suppress dominant known-category influence and strengthen novel-class separation.

- **DPN** (An et al., 2023): Utilizes dual-path representation learning to decouple shared and category-specific knowledge for improved clustering quality.

- **TAN** (An et al., 2024b): Implements transfer and alignment mechanisms to bridge known and novel feature spaces through prototype-based knowledge transfer.

#### **GCD Methods.** *LLM-enhanced Methods:*

- **LOOP** (An et al., 2024a): Leverages LLM querying and neighborhood contrastive learning to align cluster semantics and enhance pseudo-label reliability.

- **ALUP** (Liang et al., 2024): Adopts comparison-based prompting with LLMs for active pseudo-label refinement, improving clustering stability.

All implementations in BOLT follow official codebases or public releases when available. Specifically, for these LLM-based OOD detection methods, we first fine-tuned the LLM on the training set and classified samples with an OOD score exceeding a predefined threshold as OOD. Following prior work (Zeng et al., 2021b), we set this threshold to 0.5. The baselines were implemented using the pytorch-ood (Kirchheim et al., 2022)<sup>13</sup>. Pytorch-ood consists of various OOD Detection with Deep Neural Networks based on PyTorch. This library is open-source and publicly available and distributed under the Apache-2.0 license.

#### **B.4 Evaluation Metrics**

We employ distinct yet complementary evaluation metrics designed for the characteristics of the OSTC and GCD tasks.

**OSTC.** Following prior studies (Zhang et al., 2021a; Zhou et al., 2022, 2023), we adopt metrics that jointly assess classification performance on ID classes and detection capability on OOD samples.

<sup>13</sup><https://github.com/kkirchheim/pytorch-ood>

Task	Method	Backbone	Training Time
OSTC	DOC	BERT	5m 20s
	DeepUnk	BERT	5m 20s
	ADB	BERT	6m 40s
	SCL	BERT	6m 40s
	AB	BERT	5m 10s
	KnnCon	BERT	9m 10s
	DyEn	BERT	18m 50s
	LLM+OOD detection	LLaMA-3.1-8B	19m 2s
GCD	DAL	BERT	14m 0s
	Geoid	BERT	47m 10s
	SDC	BERT	11m 20s
	DPN	BERT	28m 50s
	TAN	BERT	27m 50s
	LOOP	BERT + LLM API	1h 0m 20s
	ALUP	BERT + LLM API	1h 22m 20s

Table S3: Training time on the BANKING dataset (KCR=0.25, LAR=0.1).

- **ACC:** ACC measures the proportion of correctly predicted samples in the entire test set:  $ACC = \frac{1}{N_t} \sum_{i=1}^{N_t} \mathbb{I}(\hat{y}_i = y_i)$ , where  $\mathbb{I}(\cdot)$  equals 1 if the condition holds and 0 otherwise. This metric reflects global classification performance but may be biased when ID and OOD sample distributions are imbalanced.
- **Macro-F1 (F1):** To account for class imbalance, we adopt the macro-averaged F1 score computed over the full label space  $\mathcal{Y}^k \cup \{\text{unknown}\}$ , where all novel classes  $\mathcal{Y}^n$  are collapsed into a single “unknown” label. For each label  $c$ , we compute:

$$\text{Precision}_c = \frac{TP_c}{TP_c + FP_c}, \quad (3)$$

$$\text{Recall}_c = \frac{TP_c}{TP_c + FN_c}, \quad (4)$$

$$\text{F1}_c = \frac{2 \times \text{Precision}_c \times \text{Recall}_c}{\text{Precision}_c + \text{Recall}_c}. \quad (5)$$

Here,  $TP_c$ ,  $FP_c$ , and  $FN_c$  denote the numbers of true positives, false positives, and false negatives for class  $c$ , respectively. The overall macro-F1 is defined as:  $\text{F1} = \frac{1}{|\mathcal{Y}^k|+1} \sum_{c \in \mathcal{Y}^k \cup \{\text{unknown}\}} \text{F1}_c$ .

- **Known-Class F1 (K-F1):** K-F1 isolates model performance on in-distribution (known) classes:  $\text{K-F1} = \frac{1}{|\mathcal{Y}^k|} \sum_{c \in \mathcal{Y}^k} \text{F1}_c$ . It quantifies how well the model preserves classification accuracy on ID samples while performing OOD detection.
- **Novel-Class F1 (N-F1):** To evaluate OOD detection capability, all novel classes  $\mathcal{Y}^n$  are collapsed into a single “unknown” class, and F1 is computed over this merged category:  $\text{N-F1} =$

$\text{F1}_{\text{unknown}}$ . Higher N-F1 indicates stronger ability to reject or correctly identify unseen classes as “unknown”.

**GCD.** For GCD, we follow the standard evaluation protocol introduced by Zhang et al. (2021b); An et al. (2024b, 2025) and adopted in subsequent studies. This protocol simultaneously evaluates the recognition of known categories and the discovery of novel categories:

- **K/N-ACC:** Clustering accuracy is computed separately for known and novel categories to assess performance on both in-distribution and unseen data. Given the Hungarian matching between clusters and labels (Kuhn, 2004), let  $\mathcal{Y}^k$  and  $\mathcal{Y}^n$  denote the sets of known and novel labels, respectively. The accuracy on each subset is given by:  $ACC_S = \frac{1}{|\mathcal{Y}^S|} \sum_{y \in \mathcal{Y}^S} \frac{|\hat{C}_y \cap C_y|}{|C_y|}$ ,  $S \in \{k, n\}$ , where  $C_y$  and  $\hat{C}_y$  represent the ground-truth and predicted sample sets for label  $y$ . We report  $\text{K-ACC} = \text{ACC}_k$  and  $\text{N-ACC} = \text{ACC}_n$ .
- **H-Score:** The harmonic mean of K-ACC and N-ACC:  $\text{H-Score} = \frac{2 \cdot \text{K-ACC} \cdot \text{N-ACC}}{\text{K-ACC} + \text{N-ACC}}$ , which rewards balanced performance across known and novel classes.
- **ACC:** Clustering accuracy over all samples after Hungarian matching. While intuitive, this measure can be dominated by known-class performance under class imbalance.
- **Normalized Mutual Information (NMI):** Quantifies the dependency between predicted clusters  $U$  and ground-truth labels  $V$ :  $\text{NMI}(U, V) = \frac{2I(U;V)}{H(U)+H(V)}$ , where  $I(U;V)$  is mutual information and  $H(\cdot)$  denotes LLM-Entropy. Values range from 0 (no agreement) to 1 (perfect match).
- **Adjusted Rand Index (ARI):** Measures pairwise consistency between predicted and true partitions, adjusted for chance:  $\text{ARI} = \frac{\sum_{ij} \binom{n_{ij}}{2} - \frac{\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}}{\binom{n}{2}}}{\frac{1}{2} [\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}] - \frac{\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}}{\binom{n}{2}}}$ , where  $n_{ij}$  is the number of samples shared by cluster  $i$  and class  $j$ , and  $a_i, b_j$  are marginal sums. ARI values range from  $-1$  (mismatch) to 1 (perfect alignment).

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
K-F1	0.25	DOC	27.49	34.16	17.18	30.75	0.06	57.45	0.00	0.00	0.36	0.00	0.00	1.23	14.06
		DeepUNK	33.54	14.28	26.59	26.95	37.49	25.44	46.57	48.02	3.45	35.92	28.10	5.92	27.69
		ADB	43.99	<b>61.09</b>	39.23	18.65	34.73	<b>70.76</b>	69.33	49.04	35.09	36.75	40.45	14.22	42.78
		SCL	<b>96.57</b>	15.05	22.49	8.58	18.75	33.15	23.17	26.83	1.79	20.11	9.71	4.27	23.37
		AB	52.25	23.71	<b>42.56</b>	27.43	57.24	55.20	67.50	60.87	25.92	41.22	25.47	4.96	40.36
		KNNCon	88.67	59.43	41.14	50.69	63.07	59.39	<b>76.88</b>	<b>71.59</b>	62.07	<b>58.44</b>	61.90	16.34	59.13
		DyEn	<u>93.77</u>	<b>63.08</b>	39.54	<b>62.77</b>	<b>63.17</b>	50.77	<u>76.55</u>	<u>65.97</u>	<b>77.94</b>	<b>65.99</b>	<b>66.66</b>	<b>19.57</b>	<b>62.15</b>
		LLM-MaxSoftmax	21.94	38.44	25.52	5.83	35.17	43.32	47.47	33.22	15.90	10.49	6.19	3.81	23.94
		LLM-OpenMax	20.16	33.97	24.10	4.19	7.24	52.62	13.67	7.37	11.17	13.62	0.53	3.88	16.04
		LLM-TempScale	21.65	38.32	25.25	5.64	35.00	42.28	47.09	32.29	14.85	9.18	6.01	4.24	23.48
	LLM-Energy	23.71	37.24	23.04	5.00	31.97	61.17	44.05	28.42	8.21	7.46	3.25	2.82	23.03	
	LLM-LogitNorm	22.46	35.37	25.13	5.50	39.19	38.93	58.01	34.34	15.09	10.61	4.63	4.16	24.45	
	LLM-Entropy	22.47	40.81	26.49	5.49	35.49	46.37	47.41	33.75	16.76	8.12	6.24	4.25	24.47	
	LLM-KLMatching	21.73	37.74	24.84	6.40	32.86	40.10	43.07	31.85	17.25	12.54	6.27	5.37	23.34	
	LLM-MaxLogit	24.44	38.56	23.42	5.20	37.16	<u>63.51</u>	50.56	31.81	10.30	7.88	4.40	2.84	25.01	
	0.5	DOC	19.19	37.60	4.61	0.98	0.08	68.29	43.03	0.00	0.42	0.00	0.00	0.25	14.54
		DeepUNK	30.36	12.00	11.46	27.72	48.99	26.42	59.97	58.03	2.04	50.95	27.79	4.21	30.00
		ADB	70.45	<b>81.23</b>	<u>56.07</u>	23.66	57.86	<b>82.56</b>	72.86	59.14	47.55	41.33	58.15	15.91	55.56
		SCL	<b>98.47</b>	14.51	22.88	12.96	22.57	47.49	32.10	27.25	1.39	29.21	8.13	5.09	26.84
		AB	53.74	25.71	43.59	18.48	61.04	63.20	70.10	64.81	27.48	43.75	23.96	12.28	42.34
KNNCon		92.65	<u>77.87</u>	<b>58.07</b>	<b>54.95</b>	<b>77.67</b>	<u>80.34</u>	<b>86.93</b>	<b>81.54</b>	<b>80.99</b>	<b>79.12</b>	<u>66.14</u>	<b>21.02</b>	<b>71.44</b>	
DyEn		<u>93.76</u>	75.91	55.16	<u>51.40</u>	<u>77.63</u>	74.38	<u>85.94</u>	<u>77.92</u>	<b>81.57</b>	<u>76.21</u>	<b>72.12</b>	<u>17.59</u>	<u>69.97</u>	
LLM-MaxSoftmax		78.65	54.14	43.72	21.59	57.66	71.72	75.29	51.24	24.11	8.53	10.36	6.75	41.98	
LLM-OpenMax		68.11	35.45	34.10	8.88	7.60	71.74	22.21	10.43	6.16	10.60	0.58	2.96	23.24	
LLM-TempScale		78.63	54.16	43.73	21.48	57.45	71.47	74.95	51.24	24.62	8.26	10.38	7.31	41.97	
LLM-Energy	80.84	44.30	31.41	11.16	41.02	74.87	62.65	25.88	10.16	5.59	3.52	3.54	32.91		
LLM-LogitNorm	79.64	51.65	42.04	19.57	54.98	73.09	73.70	43.03	21.73	10.56	6.34	4.64	40.10		
LLM-Entropy	79.55	55.18	44.74	20.64	57.51	72.61	74.67	52.58	23.80	6.98	10.49	6.83	42.13		
LLM-KLMatching	76.05	53.82	43.60	21.11	56.12	68.69	71.89	51.28	23.40	12.18	10.46	8.41	41.42		
LLM-MaxLogit	85.08	47.26	33.33	13.07	48.95	78.67	70.35	32.30	13.25	5.95	5.18	4.24	36.47		
0.75	DOC	31.92	27.21	5.70	23.15	0.05	72.87	0.00	0.00	0.00	0.00	0.00	0.02	13.41	
	DeepUNK	52.06	7.10	14.14	18.27	49.98	32.75	63.62	60.03	1.28	62.72	22.22	5.32	32.46	
	ADB	88.49	<u>84.66</u>	62.84	24.86	63.66	84.19	81.02	65.87	51.55	48.75	65.95	13.89	61.31	
	SCL	<b>98.39</b>	<u>14.65</u>	16.55	15.18	23.70	63.65	55.21	41.50	1.12	27.29	9.08	6.00	31.03	
	AB	57.20	24.99	48.33	15.46	62.71	66.60	71.79	64.71	28.40	43.16	22.42	11.64	43.12	
	KNNCon	<u>97.71</u>	<b>85.35</b>	<b>66.76</b>	<b>72.68</b>	<b>86.60</b>	<b>87.34</b>	<b>92.85</b>	<b>86.50</b>	<b>88.22</b>	<b>84.72</b>	<u>73.74</u>	<b>20.65</b>	<b>78.59</b>	
	DyEn	95.58	84.04	<u>65.75</u>	<u>62.85</u>	<b>86.96</b>	85.80	<u>92.39</u>	<u>82.94</u>	<b>90.28</b>	<b>84.78</b>	<b>79.43</b>	<u>17.71</u>	<u>77.38</u>	
	LLM-MaxSoftmax	91.67	72.72	55.87	28.49	65.93	84.11	86.33	61.50	32.59	36.74	12.52	8.95	53.12	
	LLM-OpenMax	51.59	49.35	36.73	8.13	11.23	73.33	30.90	9.65	4.43	8.86	1.18	2.34	23.98	
	LLM-TempScale	91.37	72.72	55.87	28.34	65.94	84.10	86.31	61.46	32.98	36.83	12.72	9.20	53.15	
LLM-Energy	81.54	53.74	29.44	13.55	34.68	76.84	68.00	32.54	9.10	17.46	3.87	3.34	35.34		
LLM-LogitNorm	93.36	72.59	54.45	21.13	49.56	<u>86.76</u>	78.47	46.16	26.44	30.22	6.07	6.89	47.68		
LLM-Entropy	92.04	73.44	56.28	29.14	65.92	83.92	85.83	62.79	31.94	36.37	12.97	9.15	53.32		
LLM-KLMatching	90.62	72.73	55.98	28.71	65.96	82.38	85.65	61.74	34.20	38.70	13.94	10.32	53.41		
LLM-MaxLogit	87.51	58.59	32.79	16.42	43.63	80.39	75.51	40.04	11.67	21.35	6.19	4.54	39.89		

Table S4: K-F1 performance of OSTC on the BOLT benchmark under different KCRs (LAR = 0.1)

## B.5 Implementation Details

All experiments were conducted on a workstation equipped with eight NVIDIA RTX 5090 GPUs (each with 32 GB of memory) and 512 GB of system RAM. For each KCR setting, the known-class label space was uniformly partitioned into five splits with minimal overlap to ensure stable and comparable evaluation. Experiments were conducted on all five splits, and the final results were reported as the average performance across these splits.

We adopted **LLaMA-3.1-8B**<sup>14</sup> as the default LLM backbone for all LLM-based methods. For PLM-based discriminative baselines, we used the **bert-base-uncased**<sup>15</sup> model for English datasets and **bert-base-chinese**<sup>16</sup> for Chinese datasets to ensure linguistic consistency. For LLM-enhanced variants, we further incorporated **DeepSeek-V3.1**

<sup>14</sup><https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

<sup>15</sup><https://huggingface.co/google-bert/bert-base-uncased>

<sup>16</sup><https://huggingface.co/google-bert/bert-base-chinese>

as a generative verification and semantic calibration module, which supports pseudo-label validation and open-set consistency checking.

All implementations in BOLT followed official codebases or public releases when available. Each configuration was pre-trained for 100 epochs when applicable and fine-tuned for 50 epochs, with early stopping applied. Specifically, in our implementation of LLM-based OOD detection, the models were fine-tuned on the training set, and OOD detection was performed using the pytorch-ood library<sup>17</sup>, an open-source framework.

## C Runtime Comparison Across Methods

Table S3 reports the end-to-end training time on BANKING with KCR=0.25 and LAR=0.1. For OSTC, most BERT-based baselines are relatively lightweight, finishing within about 5–9 minutes (e.g., AB: 5m10s; DOC/DeepUnk: 5m20s; ADB/SCL: 6m40s; KnnCon: 9m10s). In contrast, methods involving heavier ensembling or addi-

<sup>17</sup><https://github.com/kkirchheim/pytorch-ood>

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.		
K-F1	0.25	DOC	74.47	43.80	38.71	29.59	60.96	58.47	77.20	65.11	0.00	56.39	38.09	9.41	46.02		
		DeepUNK	82.13	16.60	43.12	24.10	63.86	55.91	70.89	64.56	1.41	51.86	30.01	9.07	42.79		
		ADB	69.34	73.67	47.47	45.74	67.06	<b>76.68</b>	80.15	71.86	<u>60.74</u>	50.48	<b>74.20</b>	<u>19.41</u>	<b>61.40</b>		
		SCL	68.20	20.17	17.71	28.90	20.43	40.75	39.42	29.90	2.35	38.61	22.47	9.00	28.16		
		AB	53.40	25.90	42.02	26.00	58.20	59.32	70.82	67.40	40.12	58.26	45.38	8.38	46.27		
		KNNCon	85.52	54.62	42.36	38.75	62.27	60.40	76.40	72.35	60.35	<u>60.02</u>	56.78	16.37	57.18		
		DyEn	<b>92.41</b>	63.88	39.20	42.45	64.46	55.67	76.54	68.79	<b>62.67</b>	<b>64.35</b>	<u>63.42</u>	<b>23.21</b>	59.75		
		LLM-MaxSoftmax	53.72	48.62	40.49	<u>46.97</u>	49.36	47.32	56.74	60.88	46.57	47.34	50.15	12.47	46.72		
		LLM-OpenMax	79.85	58.68	40.31	33.57	56.44	57.02	70.50	61.48	37.19	42.25	18.17	11.47	47.24		
		LLM-TempScale	53.72	48.62	40.53	46.39	49.07	47.78	56.74	60.88	45.57	46.86	50.14	12.60	46.58		
		LLM-Energy	<u>88.07</u>	<u>74.88</u>	<u>49.66</u>	39.71	69.68	74.20	<u>85.36</u>	<u>72.83</u>	49.89	55.28	52.78	5.63	59.83		
		LLM-LogitNorm	53.46	47.70	41.28	46.52	65.27	50.06	82.68	69.16	45.97	46.57	59.31	11.10	51.59		
		LLM-Entropy	55.50	49.86	43.22	<b>48.51</b>	48.65	48.83	55.64	61.45	48.22	50.82	49.84	12.41	47.75		
		LLM-KLMatching	50.79	46.76	40.06	42.26	46.24	44.52	53.30	59.09	43.59	46.14	46.30	13.51	44.38		
		LLM-MaxLogit	87.58	<b>76.17</b>	<b>50.37</b>	40.68	<b>71.77</b>	<u>76.55</u>	<b>86.28</b>	<b>74.48</b>	52.93	57.46	56.47	7.32	<b>61.50</b>		
		K-F1	0.5	DOC	71.82	48.20	39.74	41.15	64.91	72.75	79.48	69.43	0.00	62.97	16.32	8.43	47.93
				DeepUNK	86.98	38.56	44.72	24.50	58.77	61.07	67.23	65.28	1.22	51.31	22.24	7.56	44.12
				ADB	86.02	<b>84.80</b>	56.99	62.24	77.57	<u>85.26</u>	86.42	77.66	<u>81.33</u>	72.21	<b>84.26</b>	<u>19.43</u>	<b>72.85</b>
SCL	79.46			26.56	15.72	35.58	34.75	50.43	74.35	38.73	1.71	55.16	36.39	9.82	38.22		
AB	59.80			29.00	46.80	26.72	64.22	65.12	74.80	68.94	47.20	60.56	47.70	15.88	50.56		
KNNCon	92.69			75.17	57.59	55.10	<u>79.40</u>	81.16	86.60	<b>81.44</b>	<b>81.65</b>	<u>77.14</u>	64.31	<b>20.38</b>	<u>71.05</u>		
DyEn	<u>95.30</u>			76.83	56.18	56.19	77.64	74.84	86.63	78.14	77.85	<b>77.27</b>	74.40	19.34	70.88		
LLM-MaxSoftmax	79.12			72.10	<u>61.22</u>	<u>62.47</u>	71.07	71.56	78.57	76.17	69.99	67.18	74.17	14.52	66.51		
LLM-OpenMax	92.34			<u>81.48</u>	58.97	37.10	55.02	85.18	73.90	62.76	43.16	52.12	28.87	13.36	57.02		
LLM-TempScale	79.12			72.22	61.04	62.19	70.98	71.56	78.73	76.04	69.65	66.91	74.17	14.92	66.46		
LLM-Energy	94.70			76.88	45.83	39.89	74.82	84.96	86.67	73.42	46.61	63.27	62.64	6.73	63.03		
LLM-LogitNorm	83.76			72.80	58.38	<b>63.58</b>	<b>80.99</b>	76.73	<b>91.61</b>	<u>79.73</u>	72.64	67.63	<u>77.13</u>	10.50	69.62		
LLM-Entropy	79.93			72.73	<b>62.05</b>	62.45	70.10	71.50	77.29	76.54	70.05	67.98	73.25	14.93	66.57		
LLM-KLMatching	76.49			69.91	58.05	59.88	68.99	70.56	75.25	74.55	67.88	65.91	71.19	18.87	64.79		
LLM-MaxLogit	<b>95.40</b>			79.10	49.51	47.45	77.51	<b>86.21</b>	<u>89.63</u>	77.10	54.15	66.35	71.67	9.23	66.94		
K-F1	0.75			DOC	77.65	43.92	36.91	41.66	66.98	78.98	80.49	71.88	0.00	62.12	0.00	8.03	47.39
				DeepUNK	92.53	38.12	46.17	21.07	48.39	64.67	63.01	60.84	0.17	50.93	19.44	5.11	42.54
				ADB	94.41	<b>87.63</b>	65.47	66.35	82.81	86.33	90.40	82.14	<u>84.86</u>	77.09	<b>88.22</b>	19.70	<u>72.12</u>
		SCL	72.52	28.76	26.02	40.09	56.76	69.12	80.32	47.91	1.31	65.67	49.66	10.44	45.72		
		AB	63.74	29.72	51.04	27.36	67.82	68.68	77.80	71.06	50.04	60.88	49.16	16.00	52.78		
		KNNCon	<b>97.45</b>	85.00	66.16	69.33	<b>86.99</b>	86.65	<b>92.79</b>	<b>86.30</b>	<b>87.65</b>	<b>84.95</b>	<u>73.32</u>	<b>20.89</b>	<b>78.12</b>		
		DyEn	95.37	84.93	64.57	64.01	86.21	84.69	92.05	84.18	84.26	84.59	80.95	17.98	76.98		
		LLM-MaxSoftmax	92.40	82.28	<u>70.09</u>	70.45	84.07	85.68	90.17	<u>84.59</u>	80.91	77.16	84.66	17.05	76.63		
		LLM-OpenMax	93.57	<u>85.37</u>	68.85	38.94	60.49	<b>89.22</b>	77.66	67.69	31.28	54.42	35.42	13.25	59.68		
		LLM-TempScale	92.40	82.37	<b>70.17</b>	<u>70.84</u>	83.94	85.90	90.34	84.50	80.84	76.93	<u>84.80</u>	17.42	76.70		
		LLM-Energy	90.25	81.04	33.87	40.28	74.35	82.53	82.62	73.25	42.66	47.43	67.14	7.39	60.23		
		LLM-LogitNorm	<u>96.48</u>	84.56	67.75	62.48	81.37	<u>88.99</u>	91.98	81.91	79.34	72.07	75.58	10.41	74.41		
		LLM-Entropy	92.70	82.46	69.11	<b>70.88</b>	83.27	85.70	89.64	84.36	81.41	77.64	84.72	17.62	76.63		
		LLM-KLMatching	91.71	81.49	69.37	70.52	82.96	83.59	89.22	83.86	79.75	76.48	83.14	<u>20.31</u>	76.03		
		LLM-MaxLogit	92.73	83.03	44.15	44.64	77.96	85.02	87.90	77.20	49.36	53.32	74.93	9.92	65.01		

Table S5: K-F1 performance of OSTC on the BOLT benchmark under different KCRs (LAR = 0.5)

tional components incur noticeably higher costs: DyEn requires 18m50s.

For GCD, training is generally more time-consuming than OSTC. This is partly because GCD typically leverages additional unlabeled data and performs pseudo-label learning, which introduces extra computation beyond supervised training. Under this setting, SDC and DAL finish in 11m20s and 14m00s, while DPN and TAN take around 28 minutes and Geoid reaches 47m10s. Finally, LLM-enhanced methods are the most expensive, as they additionally depend on external LLMs during training: LOOP and ALUP require 1h00m20s and 1h22m20s, respectively, underscoring the substantial overhead of LLM-augmented training.

## D Unified Experiment Runner and Configuration Example

This work provides a unified experiment entry point that supports multiple task types (e.g., GCD and OSTC) and multiple datasets within a single framework. It enables grid-style experiments over the known class ratio (kcr) and labeled ratio (lar),

while allowing different methods.

### D.1 Command-line Entry

```
bolt-grid \
  --config ./config.yaml \
  --output-dir ./exp \
  --model-dir ./pretrained_models
```

where:

- `-config`: the YAML configuration file specifying tasks, methods, datasets, and grid settings.
- `-output-dir`: the output directory. Logs, aggregated results, and model artifacts are written to this directory for reproducibility and archiving.
- `-model-dir`: the local directory of pre-trained models. Users should download and place required models in advance (e.g., bert-base-chinese, bert-base-uncased, Meta-Llama-3.1-8B-Instruct). The configuration file only references model names/identifiers.

## D.2 Configuration Example (config.yaml)

**Minimal example.** The following configuration illustrates the main fields with a small subset of methods and datasets for readability.

```
maps:
  gcd: [alup, deepaligned, dpn, geoid, loop,
        ↪ plm_gcd, sdc, tan]
  openset: [ab, adb, clap, deepunk, doc, dyen,
            ↪ knncon, plm_ood, scl]

methods: [alup, deepaligned]
datasets: [banking, clinic]
result_file: summary

grid:
  known_cls_ratio: [0.25, 0.5, 0.75]
  labeled_ratio:   [0.1, 0.5, 1.0]
  fold_idx: [0, 1, 2, 3, 4]
  fold_num: 5
  seed: 2025
  cluster_num_factor: [1.0]

run:
  gpus: [0, 1]
  max_workers: 2
  num_pretrain_epochs: 100
  num_train_epochs: 50

paths:
  results_dir: results
  logs_dir: logs
```

## D.3 Field Descriptions

- **maps:** assigns each method to a task type (e.g., GCD vs. OpenSet), enabling a unified runner across tasks.
- **methods:** the methods to execute in the current run. Experiments are instantiated by combining methods, datasets, and grid.
- **datasets:** the datasets to evaluate. The runner is dataset-agnostic and supports multiple datasets in one configuration.
- **grid:** grid variables for controlled evaluation. `known_cls_ratio` corresponds to `kcr` and `labeled_ratio` corresponds to `lar`. `fold_*` and `seed` control reproducibility, and `cluster_num_factor` is used by some clustering/class-discovery methods.

## E Additional Details of Experimental Results

### E.1 Additional Results on OSTC

Table S4-S15 provide all experimental results across all datasets for the OSTC task, covering four evaluation metrics: K-F1, N-F1, ACC, and F1, under different settings of  $LAR = \{0.1, 0.5, 1.0\}$  and  $KCR = \{0.25, 0.5, 0.75\}$ .

### E.2 Additional Results on GCD

Table S16-S33 provide all experimental results across all datasets for the GCD task, including six evaluation metrics ACC, ARI, NMI, H-Score, K-ACC, and N-ACC evaluated under the same configurations of  $LAR = \{0.1, 0.5, 1.0\}$  and  $KCR = \{0.25, 0.5, 0.75\}$ .

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
0.25		DOC	76.90	48.18	39.98	42.81	66.74	57.17	80.21	70.46	0.02	62.61	52.03	12.53	50.80
		DeepUNK	<u>92.69</u>	34.13	43.76	35.37	70.48	62.65	75.07	73.46	1.70	58.44	33.07	11.21	49.34
		ADB	73.70	73.76	48.93	<u>54.51</u>	66.44	76.90	82.33	73.27	<u>71.86</u>	62.47	<u>78.94</u>	<u>21.76</u>	65.41
		SCL	69.02	25.11	17.04	34.56	27.80	41.11	50.33	39.31	2.74	43.48	38.18	10.30	33.25
		AB	53.52	26.07	42.27	28.18	58.56	59.34	70.90	67.30	42.04	60.05	51.71	10.32	47.52
		KNNCon	86.81	62.57	40.13	44.55	64.16	60.07	78.33	71.88	61.22	58.93	58.21	17.30	58.68
		DyEn	<b>96.23</b>	63.06	39.40	52.56	64.42	56.51	80.46	66.86	<b>74.14</b>	64.30	62.86	<b>23.21</b>	62.00
		LLM-MaxSoftmax	59.20	51.25	39.62	47.16	50.92	46.55	58.12	61.96	50.23	51.86	59.18	13.73	49.15
		LLM-OpenMax	81.27	64.07	39.02	42.39	68.78	63.70	81.59	<b>78.36</b>	51.03	53.90	41.91	14.56	56.71
		LLM-TempScale	58.81	51.37	39.70	46.97	50.89	46.63	58.13	61.81	48.49	51.23	59.43	13.57	48.92
		LLM-Energy	88.48	<u>77.16</u>	<u>53.97</u>	48.59	<u>72.77</u>	<u>82.38</u>	<u>88.76</u>	75.94	56.89	<u>64.40</u>	71.05	9.23	<u>65.80</u>
		LLM-LogitNorm	59.32	50.01	40.63	<b>55.93</b>	64.35	45.17	82.62	71.73	46.43	49.05	<b>79.38</b>	11.20	54.65
		LLM-Entropy	60.66	52.35	42.78	48.05	49.58	47.22	56.85	61.79	52.80	54.47	58.55	13.74	49.90
		LLM-KLMatching	56.65	49.79	38.43	44.43	47.51	44.97	54.14	60.02	44.25	49.22	54.40	14.86	46.56
		LLM-MaxLogit	88.64	<b>77.74</b>	<b>54.32</b>	53.44	<b>74.92</b>	<b>82.92</b>	<b>89.70</b>	<u>77.61</u>	60.38	<b>66.60</b>	75.89	10.57	<b>67.73</b>
K-F1	0.5	DOC	78.55	51.70	46.65	53.68	74.69	73.21	84.12	75.64	0.00	70.02	49.03	10.01	55.61
		DeepUNK	95.61	41.68	49.91	30.33	62.73	69.18	71.20	71.46	1.81	63.45	22.51	8.41	49.02
		ADB	87.12	<u>84.64</u>	58.13	<u>65.76</u>	81.32	<b>85.92</b>	88.60	80.46	<u>82.35</u>	<u>77.11</u>	<b>84.63</b>	<b>22.60</b>	<b>74.89</b>
		SCL	55.15	36.41	19.70	36.44	48.01	55.43	78.09	48.95	2.10	64.48	41.61	10.70	41.42
		AB	60.44	29.60	47.36	28.70	64.64	65.42	74.82	68.54	49.18	61.88	56.60	14.00	51.77
		KNNCon	92.92	77.29	58.84	58.78	77.68	79.83	86.79	80.95	<b>83.85</b>	76.38	64.48	<u>20.80</u>	71.55
		DyEn	95.45	76.71	56.44	57.70	78.42	74.07	86.01	79.95	81.57	<b>79.24</b>	75.29	19.64	71.71
		LLM-MaxSoftmax	79.39	70.11	<u>64.27</u>	64.60	72.16	70.65	77.32	78.91	70.87	71.05	75.95	18.77	67.84
		LLM-OpenMax	93.72	82.98	58.97	47.83	73.14	83.66	85.35	78.72	59.48	73.85	48.03	18.15	66.99
		LLM-TempScale	79.39	70.11	64.05	64.60	72.16	70.76	77.32	78.91	70.34	70.97	76.08	19.05	67.81
		LLM-Energy	<b>98.44</b>	84.45	50.02	55.17	79.66	83.71	89.64	79.05	58.23	66.34	71.00	9.65	68.78
		LLM-LogitNorm	85.19	73.12	61.57	<b>71.66</b>	<b>83.75</b>	75.08	<b>92.40</b>	<b>85.19</b>	80.28	73.66	<u>83.31</u>	10.57	<u>72.98</u>
		LLM-Entropy	79.97	70.24	<b>66.11</b>	63.14	70.85	70.29	76.22	78.35	71.16	72.32	74.83	19.57	67.75
		LLM-KLMatching	78.21	69.17	60.00	63.58	70.88	69.86	75.19	78.18	66.96	69.89	72.42	20.60	66.24
		LLM-MaxLogit	<u>97.04</u>	<b>85.06</b>	56.29	59.56	<u>81.84</u>	<u>84.62</u>	<u>91.40</u>	<u>81.49</u>	62.95	69.30	78.66	11.80	71.67
0.75		DOC	82.52	49.60	44.31	56.33	76.51	80.46	86.04	78.52	0.00	72.68	47.71	10.61	57.11
		DeepUNK	96.92	41.87	52.22	24.18	44.94	73.66	64.78	66.97	1.19	64.05	17.64	5.77	46.18
		ADB	94.48	<b>88.02</b>	66.79	72.90	<u>85.79</u>	86.91	92.33	83.99	87.08	80.06	<b>88.54</b>	<u>22.96</u>	<b>79.15</b>
		SCL	75.88	30.89	34.55	56.60	63.45	63.11	84.96	70.62	1.74	74.25	56.36	13.67	52.17
		AB	64.18	30.36	51.58	31.06	68.66	69.26	77.86	71.12	52.40	62.74	59.08	17.14	54.62
		KNNCon	97.60	<u>86.14</u>	66.39	69.91	<b>86.59</b>	86.82	<u>93.09</u>	<b>86.66</b>	<u>87.68</u>	<b>85.44</b>	71.49	21.68	78.29
		DyEn	94.69	84.32	62.68	69.56	85.71	84.29	92.12	84.94	<b>90.28</b>	<u>85.00</u>	82.96	19.26	77.98
		LLM-MaxSoftmax	93.14	82.32	<b>71.79</b>	<b>78.67</b>	85.12	85.14	89.90	86.35	84.53	83.26	86.35	20.87	78.95
		LLM-OpenMax	97.30	85.78	70.14	52.43	76.10	<b>90.37</b>	88.93	80.60	54.00	74.61	53.49	18.17	70.16
		LLM-TempScale	93.14	82.39	<u>71.73</u>	<u>78.67</u>	85.13	85.14	90.04	<u>86.36</u>	84.45	83.18	<u>86.38</u>	21.06	<u>78.97</u>
		LLM-Energy	<b>99.45</b>	81.18	37.06	51.53	75.55	81.81	86.38	74.51	63.30	65.91	74.66	9.03	66.70
		LLM-LogitNorm	96.90	85.84	70.20	74.18	83.99	<u>88.40</u>	<b>93.55</b>	84.08	86.85	82.26	84.66	10.79	78.48
		LLM-Entropy	93.69	82.53	71.72	78.52	84.61	84.87	89.22	86.20	84.74	83.32	85.65	20.87	78.83
		LLM-KLMatching	92.25	81.51	70.35	77.64	84.10	84.37	88.79	85.84	83.42	82.37	84.77	<b>23.44</b>	78.24
		LLM-MaxLogit	<u>99.19</u>	83.15	45.37	56.65	79.30	83.73	89.12	78.44	70.57	69.28	81.83	11.76	70.70

Table S6: K-F1 performance of OSTC on the BOLT benchmark under different KCRs (LAR = 1.0)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.		
0.25		DOC	87.39	<b>84.38</b>	76.33	67.42	81.92	84.59	85.50	<u>85.83</u>	69.09	<b>86.02</b>	80.16	<b>82.25</b>	<b>80.91</b>		
		DeepUNK	88.01	62.66	<b>83.35</b>	49.85	4.16	86.06	9.96	0.15	0.00	24.35	12.02	7.48	35.67		
		ADB	46.79	74.82	43.06	34.81	44.46	85.89	85.41	72.11	63.51	47.99	70.66	50.46	60.00		
		SCL	<b>99.03</b>	73.80	47.05	51.66	34.18	86.67	57.84	42.71	2.14	35.08	42.19	41.08	51.12		
		AB	84.72	52.99	<u>80.08</u>	0.00	84.70	83.81	88.91	<b>91.61</b>	<b>82.62</b>	<u>83.60</u>	83.58	<u>78.14</u>	74.56		
		KNNCon	90.36	55.74	34.78	54.45	67.07	61.43	84.02	71.83	53.06	44.92	69.14	28.85	59.64		
		DyEn	<u>96.06</u>	65.10	24.23	34.54	70.80	47.93	86.57	74.58	<u>72.99</u>	74.46	69.25	29.36	62.16		
		LLM-MaxSoftmax	46.36	42.27	39.16	56.75	61.68	50.71	69.30	61.26	64.26	52.54	72.98	49.67	55.58		
		LLM-OpenMax	42.37	66.18	17.64	60.94	86.15	67.90	86.22	85.12	<u>80.22</u>	47.10	<b>85.31</b>	75.82	66.75		
		LLM-TempScale	50.51	41.65	36.12	54.29	61.77	48.71	69.13	59.09	63.68	59.29	74.06	54.34	56.05		
		LLM-Energy	76.40	<u>81.60</u>	73.32	<b>75.72</b>	<b>87.65</b>	<b>87.92</b>	<u>89.14</u>	84.56	80.21	73.49	<u>84.61</u>	69.72	<u>80.36</u>		
		LLM-LogitNorm	31.32	27.06	47.32	63.18	76.67	21.36	88.05	70.23	38.64	43.12	<u>81.60</u>	46.89	52.95		
		LLM-Entropy	54.31	53.52	60.38	62.42	61.64	61.42	68.64	64.28	72.80	64.10	73.02	51.71	62.35		
		LLM-KLMatching	28.37	37.04	29.24	13.26	39.75	30.83	42.46	35.62	14.74	22.50	41.17	27.03	30.17		
		LLM-MaxLogit	74.74	80.59	74.00	<u>72.60</u>	<u>87.54</u>	<u>87.49</u>	<b>89.60</b>	83.26	78.26	72.74	83.30	66.33	79.20		
		N-F1	0.5	DOC	69.73	69.71	65.41	33.83	64.68	78.17	73.14	68.69	54.33	67.21	53.07	52.85	62.57
				DeepUNK	71.12	61.03	66.89	32.62	31.21	70.13	29.93	32.60	0.38	19.40	10.11	28.31	37.81
				ADB	56.37	<b>79.39</b>	42.02	26.73	60.38	<b>85.94</b>	71.25	67.82	44.19	43.98	<b>68.71</b>	29.96	56.39
SCL	<b>98.62</b>			63.88	48.79	48.22	44.74	75.48	45.37	45.47	2.98	28.27	44.55	39.06	48.79		
AB	65.25			47.69	56.91	0.00	72.48	72.21	76.81	<b>80.66</b>	63.28	68.60	64.08	<b>55.92</b>	60.32		
KNNCon	85.82			57.97	29.83	54.24	58.30	73.55	78.45	62.84	<b>73.14</b>	<b>69.79</b>	63.04	37.24	62.02		
DyEn	<u>89.13</u>			64.49	23.46	<b>73.75</b>	63.22	58.45	78.35	70.60	<u>69.66</u>	<u>69.69</u>	65.14	16.55	61.87		
LLM-MaxSoftmax	48.51			43.57	38.10	56.61	56.40	49.55	63.44	63.48	63.35	53.38	60.50	39.43	53.03		
LLM-OpenMax	73.30			64.97	59.43	<u>67.33</u>	67.97	71.44	69.32	67.80	67.13	55.78	65.25	52.89	65.22		
LLM-TempScale	48.44			43.58	38.09	51.80	55.47	48.64	62.38	60.75	61.07	55.01	60.05	40.41	52.14		
LLM-Energy	87.08			72.65	<b>67.47</b>	66.38	<u>72.89</u>	84.11	79.18	70.39	66.63	53.58	<u>65.83</u>	<u>55.46</u>	70.14		
LLM-LogitNorm	55.72			30.08	19.46	63.07	72.40	57.62	<b>82.14</b>	71.28	63.65	37.53	65.11	48.53	55.55		
LLM-Entropy	52.18			48.13	43.47	57.66	55.22	52.88	60.51	64.64	65.85	57.62	59.99	39.38	54.79		
LLM-KLMatching	36.09			40.37	34.44	31.68	45.42	33.33	47.61	46.93	12.80	14.96	35.35	16.78	32.98		
LLM-MaxLogit	88.91			<u>72.95</u>	<u>67.21</u>	65.88	<b>73.87</b>	<u>85.44</u>	<u>81.54</u>	<u>71.44</u>	66.74	53.90	65.62	53.56	<b>70.59</b>		
0.75				DOC	48.08	48.12	34.10	29.62	34.79	61.14	40.43	42.60	42.27	40.74	23.47	31.21	39.71
				DeepUNK	53.29	44.76	34.92	21.72	36.25	45.42	37.15	8.28	12.25	1.42	11.85	19.49	27.23
				ADB	57.34	<b>69.60</b>	33.76	22.20	49.22	<b>71.32</b>	62.86	55.49	38.73	26.80	<u>59.96</u>	<b>40.97</b>	49.02
		SCL	<b>95.89</b>	43.06	30.27	35.52	28.85	55.99	37.20	30.93	3.54	27.31	28.86	23.56	36.75		
		AB	45.86	34.17	34.56	0.00	51.08	53.12	58.06	<b>60.10</b>	43.60	44.85	36.59	30.46	41.04		
		KNNCon	87.78	54.71	30.78	<b>52.48</b>	<u>51.64</u>	67.16	<u>71.44</u>	53.67	<b>69.90</b>	<u>56.74</u>	55.24	18.21	<b>55.81</b>		
		DyEn	<u>94.81</u>	50.18	20.38	<u>47.97</u>	<b>60.39</b>	56.02	<b>75.52</b>	24.08	<u>55.61</u>	<b>65.44</b>	<b>60.76</b>	12.59	<u>51.98</u>		
		LLM-MaxSoftmax	50.18	38.81	30.88	37.76	43.94	45.22	56.29	47.77	45.84	36.25	38.92	25.93	41.48		
		LLM-OpenMax	50.67	49.54	35.60	36.99	40.87	61.04	44.49	40.49	43.85	39.15	39.01	31.23	42.74		
		LLM-TempScale	47.79	38.81	30.79	37.87	43.19	45.20	56.24	47.75	44.83	35.92	38.70	26.98	41.17		
		LLM-Energy	69.91	55.46	<u>37.56</u>	39.08	45.28	65.58	59.60	46.14	44.86	41.71	40.09	<u>31.62</u>	48.07		
		LLM-LogitNorm	67.64	44.16	27.55	43.03	48.34	66.88	66.20	48.79	47.77	38.15	40.31	29.97	47.40		
		LLM-Entropy	52.50	43.03	32.73	39.55	42.48	43.38	52.96	50.67	47.13	38.76	38.44	24.31	42.16		
		LLM-KLMatching	39.71	35.38	28.97	31.15	37.50	28.14	50.03	42.00	21.41	29.49	30.25	14.94	32.41		
		LLM-MaxLogit	75.11	<u>57.06</u>	<b>38.04</b>	40.39	47.12	<u>68.55</u>	64.23	48.01	45.25	41.69	40.38	31.34	49.76		

Table S7: N-F1 performance of OSTC on the BOLT benchmark under different KCRs (LAR = 0.1)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
0.25		DOC	94.02	81.51	70.17	84.29	86.66	78.54	92.07	80.71	86.43	80.20	89.10	<b>80.77</b>	83.71
		DeepUNK	<u>95.68</u>	67.52	76.38	<u>64.38</u>	89.36	89.29	92.83	91.20	86.47	<b>89.37</b>	66.42	36.75	78.80
		ADB	73.34	82.65	55.79	59.46	82.93	89.54	90.90	79.11	76.61	53.84	84.76	40.78	72.48
		SCL	89.34	66.06	61.80	83.43	82.00	84.57	82.11	81.38	8.68	<u>86.05</u>	69.84	31.53	68.90
		AB	78.72	54.78	73.12	14.72	81.22	83.46	87.58	88.62	73.06	83.92	79.34	68.34	72.24
		KNNCon	86.76	35.37	31.58	50.99	62.03	63.39	83.92	60.64	55.41	46.14	65.99	26.05	55.69
		DyEn	<b>96.15</b>	68.63	22.70	62.72	71.49	55.36	86.43	60.69	69.05	70.44	71.29	36.17	64.26
		LLM-MaxSoftmax	38.96	28.00	19.02	45.31	39.03	33.96	49.96	43.33	56.54	37.17	62.65	41.66	41.30
		LLM-OpenMax	83.31	59.20	25.23	73.42	89.29	53.46	92.85	87.96	67.50	60.64	86.97	68.01	70.65
		LLM-TempScale	38.94	28.00	19.36	43.13	38.02	35.45	49.96	43.30	53.30	34.95	62.61	44.03	40.92
		LLM-Energy	94.50	<b>90.32</b>	<b>82.24</b>	<b>84.65</b>	<b>90.82</b>	<u>93.55</u>	<b>95.09</b>	<b>92.59</b>	<b>90.36</b>	81.51	<b>92.51</b>	<u>71.59</u>	<b>88.31</b>
		LLM-LogitNorm	33.46	24.96	30.62	64.73	77.54	36.05	90.95	75.43	56.90	36.05	90.96	53.65	55.94
		LLM-Entropy	43.77	33.87	32.36	48.39	35.67	38.40	46.11	43.12	61.42	49.15	59.32	40.12	44.31
		LLM-KLMatching	24.88	17.55	15.48	29.05	25.16	19.23	37.97	32.80	40.96	29.80	47.83	28.35	29.09
		LLM-MaxLogit	93.85	<u>90.19</u>	<u>81.72</u>	83.44	<u>90.37</u>	<b>93.83</b>	<u>94.69</u>	<u>92.51</u>	<u>90.21</u>	<u>80.95</u>	<u>92.38</u>	<u>68.92</u>	<u>87.75</u>
N-F1	0.5	DOC	83.28	71.70	68.52	74.14	76.75	74.55	83.36	76.94	66.33	74.57	69.60	<u>55.06</u>	72.90
		DeepUNK	90.17	39.30	67.66	67.12	77.85	77.38	82.74	80.36	51.21	74.71	57.32	51.00	68.07
		ADB	76.23	83.88	28.08	66.42	76.40	87.77	83.97	74.95	<b>82.28</b>	68.07	80.55	20.40	69.08
		SCL	88.83	50.58	54.98	73.20	57.09	56.15	73.54	57.40	2.97	74.90	54.48	42.53	57.22
		AB	63.58	48.64	55.76	0.00	68.74	71.26	75.96	77.42	56.44	68.52	65.04	49.76	58.43
		KNNCon	87.17	48.60	26.39	51.05	65.09	74.97	77.61	62.90	74.70	63.75	58.28	30.98	60.12
		DyEn	<b>93.87</b>	65.15	21.31	58.22	63.87	56.10	80.26	58.36	65.74	66.54	63.00	18.55	59.25
		LLM-MaxSoftmax	44.00	36.96	43.04	48.17	37.94	34.12	48.99	49.01	51.03	46.76	57.48	28.37	43.82
		LLM-OpenMax	89.94	75.50	19.85	75.53	75.46	82.10	82.78	79.24	76.74	62.79	70.76	<b>56.10</b>	70.57
		LLM-TempScale	44.00	37.47	41.57	46.93	37.50	34.12	49.62	47.67	50.41	44.91	57.49	26.60	43.19
		LLM-Energy	93.86	83.49	72.79	78.13	81.61	<u>88.40</u>	89.23	<u>82.62</u>	78.43	<u>77.46</u>	83.54	54.37	<u>80.33</u>
		LLM-LogitNorm	<u>65.67</u>	42.44	24.66	76.59	80.41	58.49	<u>90.10</u>	80.56	72.93	59.55	<b>87.73</b>	43.53	65.22
		LLM-Entropy	48.54	39.44	48.35	47.86	30.94	32.71	42.25	47.77	50.25	53.04	<b>52.89</b>	23.65	43.14
		LLM-KLMatching	30.36	24.18	8.35	33.12	24.06	27.86	30.99	32.84	43.72	37.34	43.77	10.04	28.89
		LLM-MaxLogit	93.53	<b>84.00</b>	<b>73.24</b>	<b>78.53</b>	<b>82.23</b>	<b>88.93</b>	<b>90.65</b>	<b>83.77</b>	<u>79.72</u>	<b>78.45</b>	<u>86.21</u>	53.27	<b>81.04</b>
0.75		DOC	67.11	49.59	38.25	47.66	54.49	62.71	65.32	57.74	42.27	53.04	39.29	30.01	50.62
		DeepUNK	83.46	31.06	38.49	36.45	56.60	56.02	64.94	59.36	42.28	50.33	47.86	19.15	48.83
		ADB	70.77	<b>72.24</b>	31.05	51.24	<u>65.16</u>	<u>73.65</u>	73.82	62.70	<b>68.72</b>	56.32	<b>72.68</b>	30.12	60.71
		SCL	85.13	34.04	26.71	40.89	44.73	41.11	58.92	41.89	3.37	54.14	40.01	31.21	41.85
		AB	46.38	36.10	32.72	0.00	51.28	52.88	59.46	59.82	41.38	50.52	45.70	27.08	41.94
		KNNCon	85.92	50.00	9.10	50.99	54.32	63.74	71.90	50.84	63.35	<u>58.53</u>	55.48	23.17	53.11
		DyEn	<b>89.89</b>	55.23	21.41	45.75	56.65	49.96	72.09	45.93	60.56	<b>59.80</b>	58.40	12.81	52.37
		LLM-MaxSoftmax	47.42	33.50	44.90	38.39	38.61	43.57	43.10	42.55	47.92	45.65	53.36	23.35	41.86
		LLM-OpenMax	84.93	67.44	0.20	47.88	52.91	<b>75.49</b>	64.20	58.81	52.39	50.05	47.06	<b>32.72</b>	52.84
		LLM-TempScale	47.42	34.08	<u>45.22</u>	38.40	37.72	45.23	44.88	41.91	47.67	42.83	55.04	20.59	41.75
		LLM-Energy	87.04	69.23	39.21	51.44	61.48	69.95	70.63	61.68	54.67	48.86	66.69	31.31	59.35
		LLM-LogitNorm	79.96	54.59	40.19	<b>59.82</b>	<b>66.16</b>	68.46	<b>79.92</b>	<b>66.01</b>	<u>65.07</u>	56.11	70.47	28.87	<u>61.30</u>
		LLM-Entropy	50.24	35.70	<b>49.11</b>	35.04	30.11	43.12	36.43	39.20	49.47	50.22	50.69	21.89	40.94
		LLM-KLMatching	39.40	25.25	7.91	30.50	26.31	20.80	31.87	30.71	38.66	37.93	39.65	6.09	27.92
		LLM-MaxLogit	<u>87.85</u>	<u>71.01</u>	41.76	<u>53.25</u>	63.90	72.28	<u>75.63</u>	<u>64.28</u>	57.01	50.26	<u>71.70</u>	<u>31.32</u>	<b>61.69</b>

Table S8: N-F1 performance of OSTC on the BOLT benchmark under different KCRs (LAR = 0.5)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
N-F1	0.25	DOC	94.62	74.28	67.86	80.49	85.58	74.72	92.33	82.63	86.80	78.04	91.23	<b>78.43</b>	82.25
		DeepUNK	<u>97.97</u>	39.64	72.71	86.41	89.85	89.51	93.51	<b>91.98</b>	86.85	<b>89.29</b>	74.53	60.68	81.08
		ADB	73.88	82.42	61.55	74.29	83.77	90.36	91.91	82.39	86.49	67.72	87.20	49.48	77.62
		SCL	91.19	60.92	70.96	84.06	79.51	69.14	78.95	68.43	5.82	85.82	71.75	43.47	67.50
		AB	79.68	55.50	75.30	48.78	81.83	83.69	87.66	87.46	71.63	84.51	79.25	67.72	75.25
		KNNCon	89.65	59.71	31.13	59.80	67.29	64.66	85.76	64.15	55.56	46.22	72.56	33.56	60.84
		DyEn	<b>98.27</b>	68.01	19.53	48.99	72.27	56.98	89.28	74.71	79.13	68.28	71.87	36.17	65.29
		LLM-MaxSoftmax	35.04	21.80	23.40	39.80	37.50	28.33	48.57	38.29	52.27	36.37	56.22	43.69	38.44
		LLM-OpenMax	85.88	67.29	25.70	77.46	90.90	71.30	94.74	88.45	82.84	60.98	90.71	66.78	75.25
		LLM-TempScale	33.99	22.20	23.87	39.32	37.33	28.64	48.60	37.71	48.06	34.36	56.75	45.13	38.00
		LLM-Energy	93.16	<b>89.45</b>	<b>84.95</b>	<b>88.18</b>	<b>91.54</b>	<b>94.73</b>	<b>96.16</b>	<u>91.98</u>	<u>92.56</u>	<u>86.68</u>	<b>94.63</b>	<u>71.34</u>	<b>89.61</b>
		LLM-LogitNorm	37.02	14.36	43.05	66.98	73.93	21.51	90.22	69.97	45.00	31.33	92.26	57.43	53.59
		LLM-Entropy	39.34	27.30	38.74	42.30	30.70	30.99	43.64	36.81	56.24	45.66	53.13	41.74	40.55
		LLM-KLMatching	21.22	12.38	16.77	29.81	21.35	20.79	32.93	28.38	29.62	24.75	40.09	14.81	24.41
		LLM-MaxLogit	92.68	<u>89.02</u>	<u>84.89</u>	<u>87.62</u>	<u>91.21</u>	<u>94.64</u>	<u>96.07</u>	91.60	<b>92.69</b>	86.25	<u>94.59</u>	68.51	<u>89.15</u>
N-F1	0.5	DOC	86.06	71.17	67.26	70.04	78.82	72.34	85.25	76.74	66.33	77.22	76.89	<u>53.24</u>	73.45
		DeepUNK	96.09	31.88	62.84	68.70	78.49	79.72	84.83	82.31	66.35	78.10	67.07	46.29	70.22
		ADB	78.76	83.76	33.12	73.47	79.72	<b>88.26</b>	86.42	74.93	<b>82.44</b>	73.53	81.15	18.23	71.15
		SCL	81.16	49.80	44.12	65.43	61.80	53.43	70.24	53.78	2.32	75.32	52.83	<b>57.47</b>	55.64
		AB	63.34	47.90	55.80	29.88	67.64	70.70	75.80	75.02	57.34	68.76	62.14	48.70	60.25
		KNNCon	88.16	65.60	35.32	55.04	54.20	71.50	77.67	57.00	76.29	59.16	63.42	32.61	61.33
		DyEn	<u>97.92</u>	66.27	18.21	53.70	64.22	55.15	79.54	55.02	69.66	70.57	63.83	18.16	59.35
		LLM-MaxSoftmax	32.88	29.28	53.29	39.70	34.44	27.32	40.32	35.37	53.28	40.84	50.85	31.19	39.06
		LLM-OpenMax	92.10	79.56	5.87	74.24	81.18	78.59	88.04	80.19	80.89	74.44	77.60	48.13	71.74
		LLM-TempScale	32.88	29.28	52.00	39.70	34.44	27.97	40.34	35.37	51.87	40.21	51.48	29.56	38.76
		LLM-Energy	<b>98.06</b>	<u>86.82</u>	<u>74.53</u>	80.09	83.96	87.56	90.95	85.44	81.12	78.59	86.77	50.55	<u>82.04</u>
		LLM-LogitNorm	66.10	45.27	29.71	76.87	80.74	51.27	90.62	80.69	79.74	63.66	<b>88.55</b>	46.88	66.68
		LLM-Entropy	35.86	29.13	63.38	34.60	25.02	24.84	34.17	30.05	54.00	46.44	44.73	24.60	37.24
		LLM-KLMatching	24.40	23.48	11.45	33.75	25.60	22.29	27.84	26.80	36.77	33.03	32.42	7.66	25.46
		LLM-MaxLogit	96.00	<b>86.88</b>	<b>75.47</b>	<b>80.87</b>	<b>84.52</b>	<u>87.85</u>	<b>91.72</b>	<b>86.25</b>	<u>82.32</u>	<b>79.50</b>	<u>88.26</u>	47.39	<b>82.25</b>
N-F1	0.75	DOC	71.05	50.47	38.87	47.52	59.92	63.45	70.25	61.10	42.27	56.72	53.92	28.62	53.68
		DeepUNK	91.81	28.15	36.45	37.29	60.91	60.71	68.49	63.89	42.35	55.18	48.50	24.00	51.48
		ADB	70.98	<b>73.78</b>	30.79	54.80	67.71	<u>74.56</u>	77.76	62.32	<b>72.91</b>	58.99	72.87	6.71	60.35
		SCL	86.53	33.79	26.03	39.99	45.97	40.77	60.02	44.01	3.51	58.25	44.34	25.19	42.37
		AB	46.14	35.28	32.96	5.12	50.62	52.54	58.98	57.96	41.78	50.62	44.32	26.70	41.92
		KNNCon	87.47	60.68	27.07	51.54	52.96	64.47	74.13	52.22	65.82	<b>61.92</b>	49.87	22.64	55.90
		DyEn	87.27	51.22	23.03	35.09	56.05	47.53	72.76	52.82	55.61	58.31	54.69	12.85	50.60
		LLM-MaxSoftmax	57.75	30.79	<u>45.66</u>	31.06	35.89	31.19	36.72	34.61	48.11	44.17	49.10	18.51	38.63
		LLM-OpenMax	92.82	57.69	7.17	48.92	61.38	<b>75.77</b>	76.45	65.60	59.30	59.30	56.03	<b>38.09</b>	58.21
		LLM-TempScale	57.75	31.51	45.02	31.06	35.95	31.21	38.49	34.63	47.64	43.42	49.53	17.75	38.66
		LLM-Energy	<b>98.03</b>	72.16	39.77	55.30	62.49	70.24	74.89	63.37	61.31	56.29	70.46	31.90	63.02
		LLM-LogitNorm	81.65	59.99	38.70	<b>63.66</b>	<b>69.23</b>	63.40	<b>82.64</b>	<b>67.14</b>	<u>71.18</u>	<u>60.75</u>	<b>75.62</b>	27.94	<b>63.49</b>
		LLM-Entropy	61.27	32.37	<b>50.71</b>	27.62	29.39	27.53	28.40	31.13	48.04	45.81	41.52	17.80	36.80
		LLM-KLMatching	50.32	20.32	9.82	23.98	21.65	21.35	22.66	21.57	37.74	32.15	34.55	9.29	25.45
		LLM-MaxLogit	<u>96.98</u>	<u>73.59</u>	42.20	<u>57.73</u>	65.17	71.98	<u>77.90</u>	<u>65.96</u>	64.38	57.99	<u>74.79</u>	31.79	<b>65.04</b>

Table S9: N-F1 performance of OSTC on the BOLT benchmark under different KCRs (LAR = 1.0)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
F1	0.25	DOC	37.47	44.20	36.90	34.15	4.15	61.98	2.19	5.05	8.00	17.20	1.43	7.46	21.68
		DeepUNK	42.62	23.96	<u>45.51</u>	28.93	35.82	35.54	45.63	45.20	3.06	33.61	27.81	6.04	31.14
		ADB	44.46	<b>63.83</b>	40.50	20.04	35.22	<b>73.28</b>	69.74	50.39	38.25	39.00	40.99	17.07	44.40
		SCL	<b>96.98</b>	26.80	30.68	12.77	19.52	42.07	24.06	27.76	1.83	23.10	10.29	7.11	26.91
		AB	57.65	29.56	<b>55.04</b>	27.43	58.62	59.96	68.07	62.64	32.23	49.72	26.50	11.50	44.91
		KNNCon	88.95	58.70	39.02	<u>50.91</u>	<u>63.27</u>	59.73	<b>77.07</b>	<b>71.60</b>	<u>61.06</u>	<u>55.74</u>	<u>62.03</u>	<u>17.37</u>	<u>58.79</u>
		DyEn	<u>94.15</u>	63.48	34.44	<b>60.40</b>	<b>63.55</b>	50.29	76.81	66.48	<b>77.39</b>	<b>67.69</b>	<b>66.70</b>	<b>20.32</b>	<b>61.81</b>
		LLM-MaxSoftmax	26.01	39.21	30.07	10.19	36.50	44.55	48.03	34.87	21.27	18.90	7.38	7.36	27.03
		LLM-OpenMax	23.87	40.41	21.95	8.77	11.18	55.16	15.53	11.95	18.84	20.32	2.05	9.44	19.96
		LLM-TempScale	26.46	38.99	28.87	9.78	36.34	43.36	47.66	33.86	20.28	19.20	7.23	8.12	26.68
		LLM-Energy	32.49	46.11	39.80	11.51	34.75	65.63	45.20	31.72	16.21	20.67	4.70	7.98	29.73
		LLM-LogitNorm	23.94	33.71	32.53	10.47	41.07	36.01	58.78	36.45	17.71	17.11	6.01	7.46	26.77
	LLM-Entropy	27.78	43.36	37.79	10.37	36.80	48.87	47.95	35.54	22.98	19.32	7.44	7.92	28.84	
	LLM-KLMatching	22.83	37.60	26.31	6.96	33.20	38.56	43.05	32.08	16.97	14.53	6.89	7.08	23.84	
	LLM-MaxLogit	32.82	46.97	40.28	11.48	39.68	<u>67.51</u>	51.56	34.84	17.85	20.85	5.81	7.74	31.45	
	0.5	DOC	23.78	41.61	14.74	2.42	1.74	69.19	43.43	2.08	3.59	7.47	0.48	2.35	17.74
		DeepUNK	34.07	18.13	20.70	27.93	48.54	30.40	59.58	57.26	1.95	47.44	27.63	5.18	31.57
		ADB	69.17	<b>81.00</b>	<b>53.73</b>	23.79	57.92	<b>82.87</b>	72.83	59.40	47.35	41.63	58.24	16.47	55.37
		SCL	<b>98.49</b>	20.68	27.20	14.67	23.14	50.04	32.27	27.81	1.48	29.11	8.45	6.45	28.32
		AB	54.82	28.46	45.80	18.48	61.32	64.05	70.20	65.31	29.58	46.54	24.35	16.62	43.79
		KNNCon	92.03	<u>75.38</u>	<u>53.36</u>	<b>54.92</b>	<u>77.18</u>	<u>79.73</u>	<b>86.82</b>	<b>80.97</b>	<u>80.53</u>	<b>78.08</b>	<u>66.11</u>	<b>21.40</b>	<b>70.54</b>
		DyEn	<u>93.34</u>	74.48	49.88	<u>52.34</u>	<b>77.26</b>	72.93	<u>85.84</u>	<u>77.70</u>	<b>80.87</b>	<u>75.48</u>	<b>72.05</b>	<u>17.54</u>	<u>69.14</u>
		LLM-MaxSoftmax	75.91	52.82	42.78	23.16	57.63	69.71	75.14	51.61	26.42	13.51	10.81	8.06	42.30
		LLM-OpenMax	68.58	39.14	38.33	11.71	9.15	71.71	22.83	12.16	9.75	15.62	1.16	4.96	25.43
		LLM-TempScale	75.89	52.84	42.79	22.85	57.40	69.40	74.79	51.52	26.77	13.46	10.83	8.63	42.26
		LLM-Energy	81.41	47.84	37.42	13.80	41.84	75.71	62.87	27.22	13.48	10.92	4.08	5.61	35.18
		LLM-LogitNorm	77.46	48.95	38.28	21.56	55.43	71.68	73.81	43.88	24.19	13.56	7.07	6.40	40.19
	LLM-Entropy	77.06	54.30	44.53	22.30	57.45	70.82	74.49	52.94	26.27	12.60	10.93	8.13	42.65	
	LLM-KLMatching	72.41	52.14	42.07	21.59	55.85	65.47	71.57	51.15	22.78	12.48	10.69	8.75	40.58	
	LLM-MaxLogit	85.43	50.47	38.98	15.57	49.59	79.29	70.50	33.48	16.40	11.28	5.72	6.22	38.58	
	0.75	DOC	32.93	29.11	8.86	23.36	0.64	72.14	0.36	0.87	1.76	3.13	0.14	0.85	14.51
		DeepUNK	52.14	10.52	16.45	18.38	49.74	33.54	63.39	58.98	1.74	58.00	22.15	5.69	32.56
		ADB	86.54	<b>83.29</b>	59.61	24.78	63.42	83.38	80.86	65.66	51.02	47.07	65.91	14.10	60.47
		SCL	<b>98.24</b>	17.23	18.08	15.86	23.79	63.17	55.05	41.28	1.22	27.29	9.20	6.46	31.41
		AB	56.52	25.80	46.80	15.46	62.51	65.75	71.69	64.65	29.04	43.29	22.52	14.10	43.18
		KNNCon	<u>97.09</u>	<u>82.56</u>	<b>62.76</b>	<b>72.03</b>	<u>86.01</u>	<b>86.08</b>	<b>92.66</b>	<b>85.83</b>	<u>87.45</u>	<u>82.57</u>	<u>73.62</u>	<b>20.59</b>	<b>77.44</b>
DyEn		95.53	80.96	<u>60.71</u>	<u>62.38</u>	<b>86.51</b>	83.94	<u>92.24</u>	<u>81.74</u>	<b>88.83</b>	<b>83.29</b>	<b>79.31</b>	<u>17.57</u>	<u>76.08</u>	
LLM-MaxSoftmax		89.07	69.64	53.10	28.78	65.56	81.68	86.06	61.22	33.14	36.70	12.68	9.40	52.25	
LLM-OpenMax		51.53	49.37	36.60	9.10	11.74	72.56	31.02	10.27	6.08	11.19	1.41	3.10	24.50	
LLM-TempScale		88.65	69.64	53.08	28.63	65.55	81.67	86.05	61.18	33.48	36.76	12.88	9.66	52.27	
LLM-Energy		80.81	53.90	30.34	14.42	34.86	76.13	67.93	32.82	10.59	19.32	4.09	4.08	35.77	
LLM-LogitNorm		91.75	70.00	51.46	21.84	49.53	<u>85.51</u>	78.36	46.21	27.33	30.83	6.28	7.50	47.22	
LLM-Entropy	89.57	70.67	53.66	29.48	65.52	81.38	85.54	62.54	32.57	36.55	13.12	9.55	52.51		
LLM-KLMatching	87.44	69.33	52.98	28.78	65.47	78.99	85.34	61.34	33.67	37.99	14.04	10.44	52.15		
LLM-MaxLogit	86.74	58.45	33.37	17.23	43.69	79.65	75.41	40.20	13.07	22.92	6.40	5.25	40.20		

Table S10: F1 performance of OSTC on the BOLT benchmark under different KCRs (LAR = 0.1)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
F1	0.25	DOC	77.73	51.34	49.20	33.15	62.24	61.82	77.58	66.03	9.60	61.16	39.00	14.90	50.31
		DeepUNK	84.39	26.79	54.21	27.60	65.13	61.47	71.45	66.13	10.86	59.37	30.66	11.20	47.44
		ADB	70.01	75.47	50.24	46.79	67.86	<u>78.83</u>	80.43	72.29	<u>62.50</u>	51.15	<b>74.39</b>	<u>21.05</u>	62.58
		SCL	71.72	29.35	32.41	33.97	23.51	48.05	40.52	32.93	3.05	48.10	23.32	10.78	33.14
		AB	57.66	31.70	52.38	26.78	59.34	63.32	71.26	68.62	43.76	<u>63.40</u>	45.98	13.88	49.84
		KNNCon	85.73	50.77	38.76	39.71	62.26	60.90	76.59	71.66	59.80	<u>57.24</u>	56.94	17.20	56.46
		DyEn	<b>94.32</b>	64.65	33.70	44.19	63.87	55.62	76.79	68.31	<b>65.43</b>	<b>65.57</b>	<u>63.56</u>	<b>24.21</b>	60.02
		LLM-MaxSoftmax	51.26	44.50	33.34	46.69	48.84	45.09	56.56	59.84	47.68	45.31	50.37	14.80	45.36
		LLM-OpenMax	80.43	58.79	35.28	36.82	58.08	56.42	71.07	63.04	40.56	45.93	19.40	15.95	48.48
		LLM-TempScale	51.26	44.50	33.47	45.97	48.52	45.72	56.56	59.84	46.43	44.48	50.36	15.11	45.19
		LLM-Energy	<u>89.14</u>	<u>77.97</u>	<u>60.52</u>	43.66	70.74	77.43	<u>85.61</u>	<u>73.99</u>	54.38	60.52	53.49	10.72	<u>63.18</u>
		LLM-LogitNorm	50.13	43.15	<u>37.73</u>	<u>47.94</u>	65.89	47.73	82.90	69.53	47.18	44.46	59.87	14.45	50.91
	LLM-Entropy	53.55	46.66	39.60	<b>48.33</b>	48.00	47.09	55.40	60.37	49.69	50.49	50.01	14.64	46.99	
	LLM-KLMatching	46.47	40.92	31.87	41.14	45.18	40.30	52.91	57.54	43.29	42.87	46.33	14.76	41.96	
	LLM-MaxLogit	88.62	<b>78.97</b>	<b>60.82</b>	44.45	<b>72.70</b>	<b>79.43</b>	<b>86.49</b>	<b>75.54</b>	57.07	62.16	57.11	12.09	<b>64.62</b>	
	0.5	DOC	72.86	51.14	44.53	41.34	65.21	72.92	79.53	69.66	3.90	64.26	16.80	10.29	49.37
		DeepUNK	87.27	38.65	48.54	26.46	59.26	62.55	67.44	65.73	4.16	53.91	22.56	9.30	45.49
		ADB	85.13	<b>84.69</b>	52.17	<u>62.38</u>	77.54	<u>85.49</u>	86.39	77.58	<b>81.39</b>	71.75	<b>84.23</b>	19.46	<b>72.35</b>
		SCL	80.31	29.56	22.26	37.38	35.32	50.95	74.34	39.30	1.78	57.35	36.56	11.13	39.69
		AB	60.12	31.46	48.32	26.72	64.32	65.68	74.78	69.18	47.74	61.44	47.86	19.74	51.45
		KNNCon	92.19	71.85	52.39	54.90	<u>79.04</u>	80.60	86.48	<b>80.88</b>	<u>81.24</u>	<u>75.65</u>	64.26	<b>20.80</b>	<u>70.02</u>
		DyEn	<u>95.17</u>	75.65	50.37	56.31	77.28	72.38	86.54	77.54	77.69	<b>76.07</b>	74.29	19.31	69.88
		LLM-MaxSoftmax	75.93	67.71	<u>58.19</u>	61.78	70.22	68.16	78.18	75.35	68.88	64.91	74.02	15.07	64.87
		LLM-OpenMax	92.12	<u>80.73</u>	52.45	38.93	55.54	84.90	74.02	63.26	45.13	53.30	29.25	15.07	57.06
		LLM-TempScale	75.93	67.88	57.79	61.46	70.12	68.16	78.35	75.18	68.51	64.47	74.02	15.39	64.77
		LLM-Energy	94.63	77.70	50.33	41.75	74.99	85.27	86.71	73.70	48.48	64.85	62.83	8.63	64.16
		LLM-LogitNorm	82.11	69.00	52.76	<b>64.18</b>	<b>80.97</b>	75.08	<b>91.59</b>	<u>79.76</u>	72.66	66.73	<u>77.22</u>	11.82	68.66
	LLM-Entropy	77.08	68.57	<b>59.77</b>	61.75	69.10	67.98	76.83	75.67	68.89	66.32	73.07	15.28	65.03	
	LLM-KLMatching	72.29	64.20	49.76	58.67	67.84	66.68	74.67	73.28	66.46	62.73	70.95	18.52	62.17	
	LLM-MaxLogit	<b>95.23</b>	79.72	53.47	48.96	77.64	<b>86.46</b>	<u>89.65</u>	77.30	55.66	67.70	71.80	10.99	67.88	
	0.75	DOC	76.99	44.43	37.06	41.62	66.77	77.97	80.36	71.59	1.76	61.42	0.24	8.61	47.40
		DeepUNK	91.97	37.48	45.32	21.56	48.53	64.13	63.03	60.81	1.93	50.89	19.62	5.48	42.56
		ADB	92.94	<b>86.23</b>	61.65	65.88	82.51	85.53	90.25	81.74	84.19	75.49	<b>88.13</b>	19.46	<u>76.17</u>
		SCL	73.31	29.24	26.10	40.13	56.55	67.37	80.13	47.79	1.39	64.79	49.60	10.99	45.62
		AB	62.62	30.28	49.02	27.36	67.52	67.70	77.66	70.82	49.68	60.10	49.14	18.06	52.50
		KNNCon	<b>96.73</b>	81.82	59.82	68.72	<b>86.44</b>	85.22	<b>92.60</b>	<b>85.57</b>	<b>86.64</b>	<b>82.92</b>	73.21	<b>20.78</b>	<b>76.71</b>
DyEn		95.02	82.14	59.77	63.43	85.71	82.52	<u>91.87</u>	83.40	<u>84.29</u>	<u>82.68</u>	80.82	17.85	75.79	
LLM-MaxSoftmax		89.59	77.84	<u>67.29</u>	69.44	83.30	83.05	89.76	<u>83.73</u>	79.53	74.73	84.48	17.21	75.00	
LLM-OpenMax		93.03	<u>83.74</u>	61.22	39.25	60.36	<b>88.36</b>	77.54	67.51	32.16	54.08	35.49	13.76	58.87	
LLM-TempScale		89.59	77.98	<b>67.40</b>	<b>69.82</b>	83.16	83.35	89.94	83.63	79.45	74.31	<u>84.62</u>	17.50	75.06	
LLM-Energy		90.05	79.97	34.46	40.67	74.13	81.74	82.51	73.01	43.16	47.54	67.13	8.02	60.20	
LLM-LogitNorm		<u>95.45</u>	81.84	64.69	62.41	81.12	<u>87.70</u>	91.87	81.58	78.75	70.84	75.55	10.90	73.56	
LLM-Entropy	90.05	78.21	66.89	<u>69.75</u>	82.37	83.04	89.17	83.43	80.08	75.53	84.51	17.73	75.06		
LLM-KLMatching	88.44	76.38	62.54	69.26	82.00	79.67	88.71	82.77	78.04	73.51	82.88	<u>19.93</u>	73.68		
LLM-MaxLogit	92.42	81.93	43.88	44.94	77.73	84.23	87.79	76.94	49.68	53.08	74.91	10.49	64.83		

Table S11: F1 performance of OSTC on the BOLT benchmark under different KCRs (LAR = 0.5)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
F1	0.25	DOC	79.85	53.40	49.28	43.67	67.68	60.09	80.52	71.18	9.66	65.70	52.73	17.60	54.28
		DeepUNK	<u>93.57</u>	35.23	53.41	39.86	71.45	67.13	75.54	74.55	11.17	64.61	33.81	15.01	52.95
		ADB	73.73	75.49	53.13	56.04	67.30	79.14	82.58	73.81	<u>73.49</u>	63.52	<u>79.09</u>	<u>23.89</u>	66.77
		SCL	72.71	32.27	35.01	39.13	30.38	45.79	51.07	41.02	3.08	51.94	38.77	12.93	37.84
		AB	57.86	31.95	53.28	30.21	59.71	63.41	71.34	68.51	45.31	64.95	52.21	15.60	51.20
		KNNCon	87.29	62.00	37.13	45.74	64.32	60.83	78.52	71.43	60.59	56.39	58.46	18.64	58.44
		DyEn	<b>96.57</b>	64.05	32.78	52.21	64.81	56.59	80.69	67.32	<b>74.70</b>	65.10	63.03	<b>24.21</b>	61.84
		LLM-MaxSoftmax	55.17	45.36	34.21	46.55	50.25	43.51	57.87	60.56	50.46	48.76	59.13	16.12	47.33
		LLM-OpenMax	82.04	64.71	34.58	45.30	69.88	64.97	81.93	<b>78.96</b>	54.56	55.32	42.78	18.69	57.81
		LLM-TempScale	54.67	45.53	34.42	46.34	50.21	43.63	57.89	60.39	48.44	47.86	59.38	16.08	47.07
		LLM-Energy	89.26	<u>79.62</u>	<u>64.30</u>	52.32	73.71	<u>84.44</u>	<u>88.95</u>	76.89	60.85	<u>68.86</u>	71.47	14.03	<u>68.73</u>
		LLM-LogitNorm	55.60	42.88	41.44	<b>56.80</b>	64.83	41.22	82.82	71.63	46.27	45.51	<b>79.61</b>	14.80	53.62
	LLM-Entropy	57.10	47.34	41.43	47.55	48.63	44.51	56.51	60.32	53.18	52.71	58.46	15.98	48.64	
	LLM-KLMatching	50.75	42.31	31.21	43.25	46.21	40.94	53.60	58.16	42.62	44.33	54.15	14.98	43.54	
	LLM-MaxLogit	89.31	<b>79.99</b>	<b>64.51</b>	<u>56.64</u>	<b>75.73</b>	<b>84.87</b>	<b>89.87</b>	<u>78.44</u>	63.97	<b>70.53</b>	76.22	15.06	<b>70.43</b>	
	0.5	DOC	79.23	54.14	50.08	51.31	74.80	73.13	84.13	75.68	3.90	70.82	49.28	11.74	56.52
		DeepUNK	95.65	40.45	52.07	32.14	63.13	70.14	71.38	71.79	5.60	65.07	22.92	9.92	50.02
		ADB	86.36	84.53	53.96	<u>66.04</u>	81.28	<b>86.13</b>	88.57	80.29	<u>82.36</u>	<u>76.71</u>	<b>84.59</b>	<b>22.42</b>	<b>74.44</b>
		SCL	57.51	38.08	23.77	37.87	48.37	55.24	77.98	49.09	2.12	65.68	41.71	12.57	42.50
		AB	60.70	31.88	48.76	29.26	64.72	65.92	74.84	68.72	49.64	62.62	56.64	16.82	52.54
		KNNCon	92.48	75.83	54.92	58.55	77.08	79.07	86.68	80.22	<b>83.41</b>	74.47	64.47	<u>21.28</u>	70.70
		DyEn	95.67	75.40	50.07	57.48	78.06	72.35	85.92	79.20	80.87	<b>78.28</b>	75.19	19.58	70.67
		LLM-MaxSoftmax	75.16	65.01	<u>62.44</u>	63.45	71.20	66.71	76.83	77.59	69.84	67.69	75.73	19.26	65.91
		LLM-OpenMax	93.57	82.56	50.12	49.08	73.35	83.20	85.39	78.76	60.74	73.91	48.30	19.35	66.53
		LLM-TempScale	75.16	65.01	62.04	63.45	71.20	66.87	76.84	77.59	69.26	67.55	75.86	19.47	65.86
		LLM-Energy	<b>98.41</b>	<u>84.75</u>	54.10	56.37	79.77	84.06	89.66	79.24	59.57	67.70	71.14	11.28	69.67
		LLM-LogitNorm	83.45	69.64	56.26	<b>71.87</b>	<b>83.67</b>	72.92	<b>92.37</b>	<b>85.06</b>	80.25	72.55	<u>83.36</u>	12.02	71.95
	LLM-Entropy	75.96	65.10	<b>65.66</b>	61.87	69.67	66.16	75.67	76.89	70.15	69.44	74.56	19.77	65.91	
	LLM-KLMatching	73.32	63.46	51.90	62.26	69.72	65.54	74.57	76.62	65.18	65.79	72.06	20.08	63.37	
	LLM-MaxLogit	<u>96.95</u>	<b>85.29</b>	59.49	60.59	<u>81.91</u>	<u>84.92</u>	<u>91.40</u>	<u>81.63</u>	64.09	70.43	78.75	13.23	<u>72.39</u>	
	0.75	DOC	81.80	49.68	43.70	54.64	76.22	79.40	85.90	78.16	1.76	71.45	47.75	11.09	56.80
		DeepUNK	96.60	40.62	50.47	24.61	45.21	72.85	64.81	66.91	2.90	63.37	17.83	6.25	46.04
		ADB	93.01	<b>86.73</b>	62.79	72.33	<u>85.48</u>	86.14	92.20	83.55	86.49	78.44	<b>88.45</b>	<u>22.53</u>	<b>78.18</b>
		SCL	76.54	31.15	33.60	56.08	63.15	61.71	84.74	70.08	1.82	73.02	56.29	13.98	51.85
		AB	63.00	30.80	49.52	31.00	68.34	68.20	77.66	70.86	51.96	61.80	59.02	19.20	54.28
		KNNCon	96.97	<u>83.82</u>	62.03	69.33	<b>86.01</b>	85.42	<u>92.92</u>	<b>85.96</b>	<u>86.78</u>	<b>83.63</b>	71.36	21.62	77.15
DyEn		94.22	81.31	58.27	68.48	85.21	81.99	91.95	84.28	<b>88.83</b>	<u>82.95</u>	82.79	19.10	76.61	
LLM-MaxSoftmax		90.93	77.63	<u>68.89</u>	<b>77.14</b>	84.29	81.76	89.42	85.30	83.01	80.25	86.12	20.80	77.13	
LLM-OpenMax		97.02	83.23	63.14	52.32	75.85	<b>89.46</b>	88.82	80.29	54.22	73.43	53.51	18.69	69.17	
LLM-TempScale		90.93	77.77	68.76	<u>77.14</u>	84.29	81.77	89.58	<u>85.31</u>	82.92	80.12	<u>86.16</u>	20.97	77.14	
LLM-Energy		<b>99.36</b>	80.36	37.37	51.68	75.33	81.09	86.28	74.28	63.22	65.17	74.63	9.63	66.53	
LLM-LogitNorm		95.94	83.49	66.70	73.83	83.74	<u>86.84</u>	<b>93.46</b>	83.73	86.20	80.61	84.61	11.24	<u>77.53</u>	
LLM-Entropy	91.66	<u>77.97</u>	<b>69.39</b>	76.89	83.67	81.29	88.68	85.07	83.21	80.44	85.39	20.79	<u>77.04</u>		
LLM-KLMatching	89.63	75.94	63.63	75.92	83.04	80.43	88.20	84.53	81.52	78.51	84.47	<b>23.07</b>	75.74		
LLM-MaxLogit	<u>99.05</u>	82.28	45.02	56.71	79.06	82.99	89.02	78.19	70.32	68.41	81.79	12.28	70.43		

Table S12: F1 performance of OSTC on the BOLT benchmark under different KCRs (LAR = 1.0)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
ACC	0.25	DOC	78.63	<b>74.88</b>	62.74	<b>70.41</b>	68.13	76.00	74.67	<u>75.19</u>	63.35	<b>75.46</b>	67.00	<b>78.63</b>	<u>72.09</u>
		DeepUNK	79.80	49.02	<b>73.43</b>	47.70	19.39	76.51	26.10	17.74	7.45	27.42	21.34	6.55	37.70
		ADB	43.36	69.19	39.75	31.84	37.25	80.12	79.16	61.48	55.37	40.99	61.60	52.18	54.36
		SCL	<b>98.51</b>	65.50	46.15	48.13	38.69	78.13	57.54	45.80	7.10	38.47	42.46	42.00	50.71
		AB	75.92	40.39	<u>69.85</u>	63.72	76.37	74.94	<b>83.02</b>	<b>84.95</b>	<u>72.92</u>	<u>73.99</u>	<u>73.83</u>	<u>77.26</u>	<b>72.26</b>
		KNNCon	88.50	55.96	36.31	61.18	63.23	57.43	79.82	67.60	55.08	49.35	66.95	36.72	59.84
		DyEn	94.97	61.85	29.57	<u>68.53</u>	64.94	46.85	81.99	68.17	<b>76.71</b>	69.20	66.19	31.60	63.38
		LLM-MaxSoftmax	36.24	39.77	32.65	<u>44.22</u>	51.77	46.51	61.47	49.19	50.52	35.50	55.15	38.32	45.11
		LLM-OpenMax	31.27	54.33	21.21	51.16	75.88	63.01	76.20	74.27	67.71	31.87	<b>74.69</b>	75.71	58.11
		LLM-TempScale	39.94	39.39	30.98	41.05	52.47	45.53	61.70	47.74	51.01	41.09	57.00	44.60	46.04
		LLM-Energy	62.89	<u>71.30</u>	60.90	62.76	<b>78.92</b>	<b>81.20</b>	81.81	74.22	67.43	57.25	73.34	62.16	69.51
		LLM-LogitNorm	26.55	31.29	41.12	47.82	65.15	30.94	80.90	57.09	29.51	29.25	67.96	35.57	45.26
	LLM-Entropy	43.08	47.04	47.96	49.42	52.33	54.74	61.07	52.15	60.13	45.87	55.12	40.45	50.78	
	LLM-KLMatching	24.71	36.88	27.59	14.42	35.09	34.55	40.43	30.77	18.49	17.62	24.82	21.73	27.26	
	LLM-MaxLogit	60.91	70.32	61.51	59.25	<u>78.92</u>	<u>81.06</u>	<u>82.66</u>	72.48	64.86	56.19	70.87	57.63	68.05	
	0.5	DOC	56.58	59.66	48.95	23.01	45.40	73.79	64.25	52.33	45.77	50.62	36.43	45.27	50.17
		DeepUNK	59.38	45.57	51.78	43.07	41.46	57.59	46.84	43.83	12.57	39.11	29.36	23.78	41.20
		ADB	63.65	<b>79.73</b>	48.84	37.51	58.06	<b>84.25</b>	71.12	62.71	51.26	43.14	68.56	42.83	59.30
		SCL	<b>98.59</b>	51.65	44.55	50.61	46.09	67.27	47.15	46.45	11.87	37.97	40.42	39.65	48.52
		AB	59.24	36.37	49.56	52.56	66.62	67.33	73.56	<b>74.15</b>	52.68	58.79	54.71	<b>54.86</b>	58.37
		KNNCon	88.92	70.18	48.23	<b>67.02</b>	<u>68.97</u>	76.41	<b>82.08</b>	70.71	<u>77.28</u>	<b>74.20</b>	<b>71.20</b>	43.88	<b>69.92</b>
		DyEn	<u>90.25</u>	<u>70.78</u>	43.33	<u>64.99</u>	<b>69.41</b>	66.84	<u>81.08</u>	<u>73.74</u>	<b>81.72</b>	<u>72.32</u>	<u>69.68</u>	33.58	<u>68.14</u>
		LLM-MaxSoftmax	65.44	50.08	41.20	61.18	56.43	62.67	69.13	57.40	53.14	33.60	42.59	38.17	52.59
		LLM-OpenMax	71.50	54.36	48.98	63.92	52.42	71.24	56.19	52.85	53.53	34.80	49.02	48.12	54.74
		LLM-TempScale	65.41	50.09	41.20	57.06	55.87	62.18	68.52	55.62	52.60	35.29	42.03	41.11	52.25
		LLM-Energy	85.19	63.29	<u>55.44</u>	62.49	63.03	80.86	74.04	58.45	52.57	33.84	50.24	47.79	60.60
		LLM-LogitNorm	69.04	44.55	35.05	62.78	65.71	66.19	79.24	62.54	51.95	23.44	49.11	41.39	54.25
	LLM-Entropy	67.13	52.34	43.86	62.24	55.76	64.37	67.61	58.41	55.84	37.04	41.74	38.96	53.77	
	LLM-KLMatching	60.57	48.80	39.82	43.55	51.05	56.10	61.59	48.33	28.21	14.08	22.45	23.85	41.53	
	LLM-MaxLogit	87.70	64.21	<b>55.52</b>	62.94	65.55	<u>82.88</u>	<u>77.88</u>	60.64	53.64	33.90	49.77	45.99	61.72	
	0.75	DOC	43.82	41.04	22.54	55.08	16.05	68.98	25.33	27.09	28.10	25.58	13.57	24.12	32.61
		DeepUNK	54.28	30.71	26.24	34.07	47.17	40.27	58.31	48.55	13.69	51.83	31.67	19.42	38.02
		ADB	81.69	<b>80.01</b>	56.03	48.16	60.33	80.06	76.09	63.93	57.76	43.08	74.47	43.73	63.78
		SCL	<b>97.76</b>	32.28	26.08	46.39	32.99	62.41	53.34	44.28	13.79	35.24	26.61	30.73	41.83
		AB	53.22	28.83	44.13	40.96	59.28	61.90	67.70	63.71	42.54	43.69	38.20	38.06	48.52
		KNNCon	95.06	<u>78.13</u>	<b>60.34</b>	<b>81.16</b>	79.54	<b>82.29</b>	<u>87.56</u>	<b>77.80</b>	<u>85.93</u>	<u>77.80</u>	<b>79.16</b>	<b>48.49</b>	<b>77.77</b>
DyEn		<u>97.15</u>	76.54	<u>59.03</u>	<u>75.41</u>	<b>80.83</b>	79.00	<b>87.96</b>	<u>70.26</u>	<b>88.46</b>	<b>79.56</b>	78.03	45.05	76.44	
LLM-MaxSoftmax		83.16	65.08	50.96	63.80	60.50	76.17	79.35	59.17	51.19	38.07	29.20	37.87	57.88	
LLM-OpenMax		52.49	50.51	36.83	44.94	29.12	69.07	39.01	28.70	33.61	27.13	25.05	29.51	38.83	
LLM-TempScale		82.61	65.08	50.94	63.27	60.38	76.17	79.33	59.09	51.17	38.43	29.00	40.36	57.99	
LLM-Energy		78.65	55.57	34.62	47.47	41.27	73.34	66.18	42.38	37.69	31.66	27.45	30.84	47.26	
LLM-LogitNorm		87.19	65.64	49.27	60.24	50.37	<u>81.57</u>	75.32	50.35	47.66	34.56	28.34	31.76	55.19	
LLM-Entropy	83.75	66.47	51.59	64.08	60.46	75.81	78.56	60.78	51.87	38.91	28.75	38.62	58.30		
LLM-KLMatching	81.00	64.64	50.91	62.12	59.84	72.87	77.96	58.08	41.58	36.98	23.62	33.23	55.24		
LLM-MaxLogit	84.25	58.76	36.57	51.71	46.48	76.70	72.57	46.78	39.35	32.93	28.65	32.27	50.58		

Table S13: ACC performance of OSTC on the BOLT benchmark under different KCRs (LAR = 0.1)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.		
0.25		DOC	90.43	71.63	59.21	78.98	79.74	70.78	88.10	73.04	77.09	72.63	82.73	<b>78.12</b>	76.87		
		DeepUNK	<u>93.14</u>	52.91	66.20	56.24	83.31	82.77	88.22	85.97	77.17	<b>82.66</b>	54.34	30.82	71.15		
		ADB	68.65	78.15	49.49	56.86	76.45	84.92	87.06	72.69	69.75	49.13	79.52	39.63	67.69		
		SCL	85.20	59.27	57.08	79.51	71.55	76.57	74.36	71.70	12.94	<u>77.72</u>	62.53	36.61	63.75		
		AB	69.72	42.76	62.36	59.28	72.72	75.54	81.82	81.50	62.66	76.18	70.70	62.64	68.16		
		KNNCon	84.02	45.05	34.85	51.70	60.18	59.60	79.56	59.58	56.39	49.48	62.57	37.92	56.74		
		DyEn	<b>95.13</b>	64.93	28.84	62.83	65.63	52.53	81.83	57.41	65.17	65.06	66.64	38.92	62.08		
		LLM-MaxSoftmax	43.36	39.07	26.99	49.42	41.74	39.48	49.97	43.93	52.42	39.88	58.28	37.22	43.48		
		LLM-OpenMax	82.20	58.63	30.04	74.83	82.77	56.04	88.59	81.47	61.19	53.47	78.33	<u>69.33</u>	68.07		
		LLM-TempScale	43.35	39.07	27.15	48.54	41.12	40.40	49.97	43.91	49.56	38.73	58.25	41.04	43.42		
		LLM-Energy	92.14	<u>86.09</u>	<b>73.21</b>	<b>81.84</b>	<b>85.78</b>	<u>89.85</u>	<b>92.60</b>	<b>88.35</b>	<b>85.06</b>	73.60	<b>88.63</b>	65.61	<b>83.56</b>		
		LLM-LogitNorm	41.76	37.65	33.81	60.58	71.16	43.18	87.40	69.34	52.88	39.82	86.91	45.97	55.87		
		LLM-Entropy	46.50	42.03	34.02	51.12	39.90	42.28	47.49	44.01	57.14	47.37	55.66	36.84	45.36		
		LLM-KLMatching	35.83	34.29	25.33	43.10	34.40	32.00	42.64	37.89	41.53	36.01	47.60	30.29	36.74		
		LLM-MaxLogit	91.33	<b>86.11</b>	<u>72.66</u>	<u>80.03</u>	<u>85.39</u>	<b>90.36</b>	<u>92.16</u>	<u>88.31</u>	<u>84.96</u>	73.23	<u>88.60</u>	62.61	<u>82.98</u>		
		ACC	0.5	DOC	79.75	63.46	57.99	77.14	72.14	73.15	81.81	73.62	50.34	69.43	58.16	<u>54.01</u>	67.58
				DeepUNK	89.04	37.22	58.59	62.37	70.32	71.61	76.76	75.56	38.55	67.13	45.42	50.04	61.88
				ADB	80.64	<b>84.00</b>	45.20	73.06	76.13	86.50	84.58	74.52	<b>82.47</b>	69.00	81.87	35.09	72.75
SCL	86.60			48.23	43.54	74.88	50.42	59.31	73.98	52.33	13.30	68.28	56.66	46.58	56.18		
AB	60.76			38.42	50.30	47.54	65.56	67.56	74.62	72.40	51.58	63.68	61.64	47.44	58.46		
KNNCon	89.28			65.93	46.80	66.96	72.39	77.36	81.51	69.93	<u>78.23</u>	70.60	69.22	44.41	69.39		
DyEn	<b>94.37</b>			71.42	43.88	71.52	70.11	65.89	82.48	66.50	70.67	71.55	71.01	33.83	67.77		
LLM-MaxSoftmax	64.45			58.66	53.61	64.57	57.81	57.91	65.67	61.61	63.04	58.61	68.07	36.64	59.22		
LLM-OpenMax	90.66			78.61	45.66	80.08	68.56	83.43	79.59	74.81	71.71	58.31	61.72	<b>58.51</b>	70.97		
LLM-TempScale	64.45			58.87	53.07	63.76	57.64	57.91	65.98	61.09	62.68	57.82	68.08	35.94	58.94		
LLM-Energy	94.02			81.06	64.40	81.22	78.79	<u>87.04</u>	88.33	79.77	73.87	<u>72.15</u>	81.91	50.23	<u>77.73</u>		
LLM-LogitNorm	74.63			60.51	45.96	81.18	<b>80.05</b>	68.51	<b>90.30</b>	79.59	73.87	63.87	<b>87.60</b>	43.37	70.79		
LLM-Entropy	66.27			59.83	56.64	64.00	55.36	57.52	62.89	61.26	62.81	61.15	65.75	33.47	58.91		
LLM-KLMatching	59.25			54.29	41.87	57.18	53.05	55.63	58.68	55.39	59.71	54.98	61.91	30.97	53.58		
LLM-MaxLogit	93.97			<u>82.04</u>	<b>65.58</b>	<b>81.84</b>	<u>80.05</u>	<b>87.77</b>	<u>90.18</u>	<b>81.47</b>	75.92	<b>73.78</b>	<u>85.55</u>	51.79	<b>79.16</b>		
0.75				DOC	75.05	48.73	38.03	69.92	62.75	74.08	75.90	68.31	28.10	59.27	24.95	35.74	55.07
				DeepUNK	89.94	35.53	43.26	39.35	53.77	61.57	65.71	63.70	28.13	51.54	38.81	24.79	49.67
				ADB	89.03	<b>82.93</b>	58.09	77.79	78.04	82.33	85.92	76.41	<b>82.55</b>	70.88	<b>85.49</b>	41.48	<u>75.91</u>
		SCL	77.81	38.57	31.07	66.54	55.33	64.31	74.73	48.32	11.84	62.30	61.48	44.72	53.08		
		AB	58.14	32.34	45.98	43.48	63.04	63.46	72.38	67.48	48.74	56.98	56.94	37.78	53.90		
		KNNCon	<u>94.40</u>	77.41	58.60	<b>81.70</b>	<b>80.14</b>	80.94	<u>87.67</u>	<u>77.32</u>	<u>81.68</u>	<u>78.46</u>	78.94	<u>47.01</u>	<b>77.02</b>		
		DyEn	<b>95.23</b>	77.68	58.29	76.45	<u>79.45</u>	77.10	86.83	75.22	77.51	<b>78.69</b>	80.10	45.01	75.63		
		LLM-MaxSoftmax	83.13	72.64	<u>64.07</u>	77.51	<u>75.27</u>	77.52	80.62	75.06	74.65	69.85	79.91	46.75	73.08		
		LLM-OpenMax	91.25	<u>80.07</u>	60.39	72.86	57.92	<b>85.31</b>	73.06	67.42	54.99	53.35	48.61	<b>49.42</b>	66.22		
		LLM-TempScale	83.13	72.81	<b>64.22</b>	77.59	75.08	77.94	80.99	74.92	74.65	69.37	80.19	46.08	73.08		
		LLM-Energy	90.32	77.30	38.77	72.29	70.23	79.10	79.35	70.99	58.57	49.85	77.50	35.18	66.62		
		LLM-LogitNorm	92.01	77.35	61.16	<u>81.51</u>	76.84	<u>84.00</u>	<b>88.61</b>	<b>77.73</b>	78.52	67.49	81.60	34.15	75.08		
		LLM-Entropy	83.75	73.26	64.05	76.29	73.84	77.61	79.40	74.53	75.30	71.24	79.40	46.23	72.91		
		LLM-KLMatching	81.43	71.01	60.99	76.29	73.22	73.31	78.51	72.87	72.65	68.70	76.72	44.99	70.89		
		LLM-MaxLogit	91.92	79.23	45.36	74.86	73.55	81.42	84.50	74.22	62.99	53.60	<u>82.47</u>	39.63	70.31		

Table S14: ACC performance of OSTC on the BOLT benchmark under different KCRs (LAR = 0.5)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
ACC	0.25	DOC	91.41	67.67	57.33	75.55	79.18	67.42	88.68	75.93	77.32	70.94	86.25	<b>75.12</b>	76.07
		DeepUNK	<u>96.85</u>	34.31	62.46	79.63	84.52	83.54	89.29	<u>87.56</u>	77.43	<b>82.99</b>	62.12	58.21	74.91
		ADB	69.48	77.96	55.02	71.41	77.34	85.91	88.46	76.51	81.78	61.75	82.89	45.01	72.79
		SCL	87.29	57.54	63.47	79.82	69.79	65.68	72.61	61.88	10.81	78.01	66.06	47.52	63.37
		AB	70.54	43.18	64.79	54.54	73.43	75.91	81.80	80.30	61.10	77.04	70.76	61.42	67.90
		KNNCon	87.22	59.92	33.78	57.85	63.69	60.26	81.63	62.33	55.55	50.53	67.90	40.00	60.06
		DyEn	<b>97.53</b>	64.26	26.85	66.88	66.36	53.80	85.38	68.15	78.02	63.21	67.16	38.92	64.71
		LLM-MaxSoftmax	41.22	36.60	29.49	42.88	41.13	36.88	49.16	42.06	50.27	40.60	53.94	40.71	42.08
		LLM-OpenMax	83.93	64.94	30.89	75.47	85.86	68.29	91.86	83.58	75.96	56.25	85.07	<u>71.54</u>	72.80
		LLM-TempScale	40.44	36.82	29.70	42.49	41.03	37.07	49.18	41.67	46.88	39.40	54.32	43.25	41.85
		LLM-Energy	91.00	<b>85.38</b>	<b>77.23</b>	<b>85.27</b>	<b>87.00</b>	<b>91.95</b>	<b>94.22</b>	<b>87.82</b>	<u>88.39</u>	<u>80.36</u>	<u>91.88</u>	65.00	<b>85.46</b>
		LLM-LogitNorm	43.18	33.57	41.35	62.90	68.24	33.59	86.60	65.12	43.74	37.79	88.86	48.61	54.46
LLM-Entropy	44.30	39.25	38.40	44.24	37.45	38.41	46.05	41.41	53.94	46.22	51.73	39.14	43.38		
LLM-KLMatching	33.75	32.57	26.06	36.94	32.84	33.16	39.99	36.59	35.45	34.35	43.32	21.74	33.90		
LLM-MaxLogit	90.50	<u>85.02</u>	<u>77.16</u>	<u>84.76</u>	<u>86.74</u>	<u>91.88</u>	<u>94.16</u>	87.43	<b>88.66</b>	80.09	<b>91.93</b>	61.97	<u>85.02</u>		
ACC	0.5	DOC	83.65	63.67	58.56	73.88	76.63	72.28	84.58	75.21	50.34	74.14	73.06	53.40	69.95
		DeepUNK	95.93	35.82	56.46	65.35	71.85	75.71	79.34	78.31	50.49	72.69	50.92	47.48	65.03
		ADB	82.34	83.90	47.55	78.82	79.68	<b>87.07</b>	86.93	75.29	<b>83.10</b>	74.44	82.58	32.73	74.54
		SCL	69.66	50.77	40.29	70.60	57.85	61.12	73.85	54.28	9.64	71.36	54.76	<b>58.21</b>	56.03
		AB	60.96	38.38	50.62	46.02	65.10	67.30	74.50	70.54	52.64	64.58	61.56	47.20	58.28
		KNNCon	89.75	72.20	49.26	69.51	67.84	75.35	81.59	67.40	80.66	68.85	71.43	45.48	69.94
		DyEn	<u>97.78</u>	71.86	42.84	70.23	70.33	65.49	82.04	63.15	<u>81.72</u>	<u>74.66</u>	72.17	34.38	68.89
		LLM-MaxSoftmax	59.97	56.13	58.79	60.16	57.71	55.73	62.23	57.75	65.97	59.52	65.59	41.41	58.41
		LLM-OpenMax	92.63	81.24	42.48	82.12	78.04	80.48	86.95	79.01	78.03	73.84	73.62	52.84	75.11
		LLM-TempScale	59.97	56.13	58.18	60.16	57.71	55.97	62.24	57.75	65.06	59.27	65.93	40.36	58.23
		LLM-Energy	<b>98.08</b>	<u>85.78</u>	<u>66.63</u>	<u>84.08</u>	<u>82.03</u>	86.34	90.45	<u>83.45</u>	77.71	74.26	86.12	45.93	<u>80.07</u>
		LLM-LogitNorm	74.74	62.49	50.23	81.80	81.38	64.99	<u>90.93</u>	80.99	80.62	69.00	<b>88.57</b>	40.80	72.21
LLM-Entropy	61.12	56.28	64.31	58.45	54.46	54.95	59.89	55.99	66.42	62.05	62.92	38.20	57.92		
LLM-KLMatching	56.97	54.24	43.88	58.20	54.66	54.10	57.69	54.86	58.47	56.68	58.32	32.97	53.42		
LLM-MaxLogit	96.13	<b>86.04</b>	<b>68.83</b>	<b>84.94</b>	<b>83.10</b>	<u>86.79</u>	<b>91.48</b>	<b>84.55</b>	79.53	<b>75.71</b>	<u>88.07</u>	44.62	<b>80.82</b>		
ACC	0.75	DOC	79.47	51.25	42.55	71.84	70.99	75.56	81.40	73.55	28.10	67.67	61.33	36.85	61.71
		DeepUNK	95.55	37.38	47.11	45.39	53.99	68.99	68.12	68.84	28.31	61.03	36.75	33.98	53.79
		ADB	89.09	<b>83.59</b>	59.09	80.86	<b>81.01</b>	<u>83.05</u>	<u>88.36</u>	77.21	<u>84.61</u>	73.78	<b>85.69</b>	42.47	77.40
		SCL	80.20	39.54	37.82	72.27	60.47	61.65	78.63	64.37	13.30	69.39	64.69	47.10	57.45
		AB	58.40	32.54	46.48	43.06	63.50	63.72	72.28	66.94	49.52	58.60	60.26	37.50	54.40
		KNNCon	94.86	79.70	59.82	81.65	79.67	81.31	88.26	<u>78.00</u>	83.78	<b>79.42</b>	77.89	49.07	<b>77.79</b>
		DyEn	94.06	76.59	57.03	<u>81.88</u>	78.91	76.31	87.11	76.75	<b>88.46</b>	<u>78.79</u>	82.44	44.65	76.91
		LLM-MaxSoftmax	85.28	72.29	<u>65.99</u>	<u>77.27</u>	75.95	75.30	79.59	74.55	77.12	75.29	79.52	46.23	73.70
		LLM-OpenMax	96.16	79.55	62.62	77.31	71.37	<b>86.31</b>	85.24	77.13	68.21	69.49	65.80	<b>57.14</b>	74.69
		LLM-TempScale	85.28	72.43	65.86	77.27	75.96	75.31	79.95	74.57	76.99	75.11	79.56	46.29	73.72
		LLM-Energy	<b>99.02</b>	78.75	40.92	76.69	71.47	78.84	83.50	72.40	69.87	63.26	81.04	37.23	71.08
		LLM-LogitNorm	92.83	79.44	64.23	<b>84.57</b>	<u>79.84</u>	82.41	<b>90.64</b>	<b>79.21</b>	84.00	76.37	<b>86.25</b>	29.13	<u>72.41</u>
LLM-Entropy	86.27	72.78	<b>66.76</b>	76.57	74.94	74.78	78.13	74.19	77.53	75.77	78.05	46.53	73.53		
LLM-KLMatching	83.47	70.54	62.40	76.12	73.72	73.74	77.14	72.64	74.99	73.17	76.36	48.31	71.88		
LLM-MaxLogit	<u>98.53</u>	<u>80.48</u>	46.57	78.86	74.89	80.64	86.26	75.64	74.08	65.92	85.30	41.23	74.03		

Table S15: ACC performance of OSTC on the BOLT benchmark under different KCRs (LAR = 1.0)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.		
K-ACC	0.25	DAL	87.09	81.43	52.09	42.75	60.46	70.18	83.80	60.54	63.29	57.48	79.57	21.04	63.31		
		GeoID	<b>98.82</b>	73.77	62.98	42.90	71.84	71.87	88.86	70.20	52.37	60.04	81.48	21.10	66.35		
		SDC	73.88	81.72	54.25	57.80	53.92	74.42	75.20	64.94	69.60	40.46	73.44	<b>33.33</b>	62.75		
		DPN	94.11	90.39	61.05	<u>66.57</u>	72.18	<b>86.64</b>	84.29	69.59	<u>75.82</u>	<u>70.82</u>	<b>82.85</b>	22.88	73.10		
		TAN	95.13	<u>91.67</u>	<u>67.87</u>	<b>78.70</b>	71.54	82.29	89.17	69.39	69.16	70.45	81.91	<u>26.00</u>	<u>74.44</u>		
		LOOP	95.01	<b>92.30</b>	60.38	49.94	<b>78.35</b>	83.09	<b>93.53</b>	<u>79.31</u>	71.93	56.06	77.93	24.32	71.85		
		ALUP	<u>98.10</u>	90.22	<b>71.52</b>	52.37	<u>73.82</u>	82.40	88.60	<u>75.69</u>	<b>76.55</b>	<b>74.84</b>	<u>82.42</u>	22.99	<u>74.13</u>		
		K-ACC	0.5	DAL	93.52	85.31	50.87	52.27	65.31	81.18	85.67	69.09	52.32	70.02	79.51	18.23	66.94
				GeoID	<u>97.87</u>	87.76	54.80	41.02	77.11	81.24	91.73	79.67	60.35	49.69	81.34	25.72	69.02
				SDC	85.96	81.08	54.22	61.43	57.67	78.08	70.36	70.53	67.22	67.35	70.41	<b>45.30</b>	67.47
				DPN	95.61	90.29	58.03	<u>78.13</u>	75.80	82.63	85.12	74.97	<u>78.81</u>	<b>77.03</b>	<u>83.51</u>	<u>36.15</u>	76.34
				TAN	<b>99.38</b>	<b>93.87</b>	59.48	<b>78.24</b>	78.57	80.91	89.11	74.25	<b>79.24</b>	76.81	<b>87.52</b>	35.79	<b>77.76</b>
LOOP	96.57			93.59	<u>64.53</u>	56.34	<b>80.29</b>	<b>87.25</b>	<b>93.88</b>	<b>80.90</b>	66.27	51.31	80.39	21.89	72.77		
ALUP	87.48			90.24	<b>68.02</b>	53.16	<u>79.69</u>	<u>83.44</u>	<u>92.86</u>	78.23	64.81	<u>76.87</u>	82.13	22.70	73.30		
K-ACC	0.75			DAL	95.48	87.27	64.35	45.91	66.67	76.28	87.14	71.10	59.90	75.42	80.13	18.17	68.99
				GeoID	97.00	85.43	62.85	53.44	79.83	81.32	90.77	79.52	66.67	52.20	83.31	22.52	71.24
				SDC	91.67	85.86	63.34	67.11	72.25	83.41	80.09	77.83	76.30	74.79	78.06	<b>50.98</b>	75.14
				DPN	<b>99.96</b>	88.79	62.95	76.24	76.85	81.64	90.06	75.56	<b>84.58</b>	76.90	84.65	<u>37.93</u>	78.01
				TAN	<u>99.94</u>	84.93	63.51	<b>78.22</b>	80.08	82.81	90.40	76.75	<u>81.59</u>	<u>79.01</u>	<b>87.47</b>	36.36	<b>78.42</b>
		LOOP	99.52	<b>92.39</b>	<u>64.83</u>	49.26	<b>83.13</b>	<b>87.83</b>	<u>93.58</u>	<u>80.42</u>	66.68	71.15	81.73	21.45	74.33		
		ALUP	94.57	<u>91.27</u>	<b>68.06</b>	57.65	<u>80.30</u>	<u>85.00</u>	<u>91.61</u>	<u>80.32</u>	73.04	<b>79.45</b>	84.05	27.05	76.03		

Table S16: K-ACC performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.1)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
K-ACC	0.25	DAL	87.96	81.19	52.44	50.61	65.93	74.21	86.28	69.55	60.41	67.73	79.36	22.56	66.52
		GeoID	<b>100.00</b>	<u>94.01</u>	51.50	48.82	80.79	86.40	85.44	87.11	67.69	82.72	<u>87.20</u>	24.69	74.70
		SDC	54.69	<u>72.65</u>	51.92	51.42	45.04	72.92	61.87	57.30	64.37	62.30	68.84	<u>30.58</u>	57.83
		DPN	88.58	91.04	65.83	<b>79.56</b>	79.17	<b>91.03</b>	84.84	80.43	<b>87.67</b>	81.63	85.98	27.59	78.61
		TAN	95.86	<u>93.22</u>	70.40	<u>79.21</u>	84.82	86.07	91.58	82.68	<u>84.28</u>	82.96	<b>87.97</b>	29.23	80.69
		LOOP	92.59	<b>96.62</b>	<u>72.75</u>	<u>67.75</u>	<b>87.78</b>	<u>89.65</u>	<b>96.53</b>	<b>88.80</b>	79.20	<u>84.30</u>	<u>86.84</u>	29.15	<b>81.00</b>
	ALUP	<u>97.84</u>	89.51	<b>74.40</b>	74.35	<u>86.62</u>	89.16	<u>95.06</u>	<u>87.70</u>	73.66	<b>86.30</b>	84.82	<b>30.67</b>	<u>80.84</u>	
	0.5	DAL	<b>98.77</b>	83.84	60.00	54.53	74.93	81.47	88.76	76.93	65.11	82.53	81.74	21.66	72.52
		GeoID	92.32	93.62	<u>68.39</u>	46.60	86.64	81.10	94.49	86.44	74.75	85.17	<u>89.30</u>	28.23	77.25
		SDC	81.32	79.41	52.84	68.64	61.98	80.35	65.80	72.33	72.11	75.11	69.67	<b>46.94</b>	68.88
		DPN	95.34	90.44	59.18	<u>84.76</u>	83.88	89.52	93.22	83.32	<u>86.70</u>	<u>87.12</u>	88.82	<u>39.99</u>	81.86
		TAN	97.40	<b>95.23</b>	65.86	<b>87.84</b>	86.16	87.58	94.95	84.76	<b>91.45</b>	<u>84.66</u>	<b>90.40</b>	35.82	<b>83.51</b>
		LOOP	<u>98.76</u>	<u>94.78</u>	67.66	64.72	<u>89.22</u>	<u>90.22</u>	<b>96.64</b>	<u>87.72</u>	76.82	86.23	88.62	26.48	80.66
	ALUP	97.49	94.48	<b>72.28</b>	76.59	<b>90.13</b>	<b>91.02</b>	<u>95.91</u>	<b>89.48</b>	83.15	<b>89.20</b>	88.17	29.58	<u>83.12</u>	
	0.75	DAL	97.10	86.95	61.38	48.21	76.99	80.51	92.01	77.89	71.61	84.66	83.53	20.85	73.47
		GeoID	99.06	<u>94.44</u>	68.12	53.66	89.40	88.96	94.17	<b>88.60</b>	78.12	88.21	90.04	26.48	79.94
		SDC	93.92	86.95	65.24	78.28	78.69	86.68	81.26	82.47	80.76	82.38	79.33	<b>51.49</b>	78.95
		DPN	<b>100.00</b>	90.99	67.46	<b>87.17</b>	88.59	88.51	93.88	85.69	<b>92.02</b>	85.72	<u>90.85</u>	<u>43.07</u>	<u>84.50</u>
		TAN	99.55	93.00	67.92	<u>86.39</u>	89.02	87.68	95.07	86.61	<u>89.79</u>	86.30	<b>91.81</b>	41.25	<b>84.53</b>
		LOOP	95.99	<b>94.45</b>	<u>69.04</u>	61.31	<b>91.23</b>	<u>89.61</u>	<b>96.69</b>	88.04	80.91	<u>88.31</u>	87.86	28.38	80.98
	ALUP	<u>99.98</u>	93.44	<b>70.06</b>	74.23	<u>91.07</u>	<b>91.23</b>	<u>95.88</u>	<u>88.26</u>	87.70	<b>89.64</b>	89.74	31.52	83.56	

Table S17: K-ACC performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.5)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
K-ACC	0.25	DAL	86.64	76.05	57.55	49.65	71.08	61.07	84.18	73.45	62.51	63.35	79.65	22.51	65.64
		GeoID	<b>100.00</b>	95.74	<u>72.63</u>	71.30	79.47	85.67	94.39	87.68	74.37	83.94	87.00	29.24	80.12
		SDC	68.47	79.69	48.30	52.74	43.13	74.25	64.60	62.26	69.74	65.61	65.97	30.94	60.47
		DPN	95.30	<u>96.07</u>	64.55	76.48	83.74	<b>93.26</b>	85.09	81.54	<b>89.04</b>	81.01	87.45	<u>35.48</u>	80.75
		TAN	<u>98.24</u>	95.12	69.39	<b>82.35</b>	89.37	88.57	95.33	84.01	<u>86.65</u>	<b>89.56</b>	<u>87.95</u>	29.16	<u>82.98</u>
		LOOP	98.00	<b>96.22</b>	<b>76.50</b>	75.20	<b>91.16</b>	<u>91.20</u>	<b>97.75</b>	<b>91.65</b>	77.39	85.88	87.24	27.47	82.97
	ALUP	98.16	92.38	68.60	<u>79.56</u>	<u>90.66</u>	90.80	<u>95.79</u>	<u>89.58</u>	79.80	<u>86.90</u>	<b>88.20</b>	<b>36.56</b>	<b>83.08</b>	
	0.5	DAL	92.18	82.99	55.42	54.96	72.89	78.09	91.02	77.26	63.01	83.00	81.35	23.05	71.27
		GeoID	<b>100.00</b>	<b>95.92</b>	70.06	71.60	90.66	89.99	<u>97.65</u>	88.83	79.01	87.12	<u>90.47</u>	30.36	82.64
		SDC	82.66	84.51	56.92	70.59	67.21	83.41	70.33	75.51	77.16	77.11	69.19	<b>51.54</b>	72.18
		DPN	96.22	93.67	62.08	<b>88.08</b>	87.74	90.68	93.67	86.67	89.04	87.71	89.85	<u>36.97</u>	83.53
		TAN	<b>100.00</b>	95.73	63.11	<u>84.68</u>	87.75	89.10	95.15	87.07	<b>90.30</b>	<b>89.85</b>	<b>91.61</b>	36.48	84.24
		LOOP	98.46	95.24	<b>73.28</b>	74.31	<b>92.58</b>	90.01	97.32	<b>90.13</b>	87.41	86.72	89.50	30.48	83.79
	ALUP	<u>99.92</u>	94.95	<u>70.58</u>	80.72	<u>92.25</u>	<b>92.20</b>	<b>97.87</b>	<u>88.85</u>	<u>89.72</u>	<u>87.90</u>	89.21	34.67	<b>84.90</b>	
	0.75	DAL	93.84	89.22	66.68	49.70	78.92	78.45	92.70	79.13	70.44	87.42	82.04	22.34	74.24
		GeoID	<b>100.00</b>	86.81	64.88	70.75	91.35	91.05	95.24	89.51	86.21	<b>91.06</b>	90.98	30.73	82.38
		SDC	94.54	93.11	68.52	82.22	81.42	88.64	82.93	84.86	83.06	85.27	80.01	<b>60.81</b>	82.12
		DPN	97.60	93.82	68.74	<b>90.23</b>	91.47	90.28	95.01	88.66	<b>94.31</b>	88.54	<u>91.84</u>	<u>43.26</u>	<b>86.15</b>
		TAN	<b>100.00</b>	94.70	68.55	<u>87.24</u>	91.49	89.64	96.07	88.41	<u>92.89</u>	89.28	<b>92.86</b>	40.12	<u>85.94</u>
		LOOP	97.87	94.63	<u>70.05</u>	70.99	<u>92.91</u>	<u>91.88</u>	<u>96.94</u>	<u>90.14</u>	87.51	89.37	88.57	33.96	83.74
	ALUP	<u>99.04</u>	<b>95.28</b>	<b>72.22</b>	75.51	<b>93.40</b>	<b>92.08</b>	<b>97.21</b>	<b>90.31</b>	91.34	<u>90.65</u>	90.76	32.99	85.07	

Table S18: K-ACC performance of GCD on the BOLT benchmark under different KCRs (LAR = 1.0)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
N-ACC	0.25	DAL	63.87	35.52	19.23	29.98	46.34	48.19	67.60	46.34	42.33	43.79	74.83	13.52	44.29
		GeoID	78.30	<u>72.81</u>	40.25	33.86	47.59	73.94	71.55	59.59	44.88	34.40	68.87	20.65	53.89
		SDC	67.34	56.45	32.25	<b>37.43</b>	37.77	62.42	58.03	57.55	48.64	35.53	59.23	<b>24.96</b>	48.13
		DPN	75.32	46.53	17.27	29.88	23.57	8.52	41.72	33.13	44.67	26.89	67.57	16.19	35.94
		TAN	83.56	54.85	20.14	31.82	28.47	13.83	47.05	37.72	51.31	29.00	66.79	<u>22.52</u>	40.59
		LOOP	<u>86.29</u>	65.64	<u>55.75</u>	<u>36.94</u>	<u>63.72</u>	<u>75.61</u>	<b>82.88</b>	<b>69.47</b>	<b>63.64</b>	40.77	70.26	19.16	<u>60.84</u>
	ALUP	79.63	<b>74.68</b>	<u>51.00</u>	34.91	<u>60.09</u>	72.66	<u>79.17</u>	<u>64.35</u>	<u>63.41</u>	<u>42.80</u>	<u>72.35</u>	20.01	<u>59.59</u>	
	0.5	DAL	89.38	40.91	25.44	36.50	50.81	62.51	72.30	52.68	51.80	50.08	<u>75.36</u>	9.06	51.40
		GeoID	92.35	67.51	51.50	39.91	55.13	73.87	72.71	63.67	51.13	39.29	70.25	9.46	57.23
		SDC	91.25	<b>80.49</b>	<u>56.01</u>	<b>56.48</b>	57.80	75.50	70.67	67.92	<b>69.08</b>	<u>55.64</u>	68.00	<u>15.38</u>	<u>63.69</u>
		DPN	81.74	33.21	18.42	24.08	33.22	13.69	54.72	29.79	46.85	21.77	68.46	12.34	36.52
		TAN	<b>97.73</b>	58.74	31.37	45.58	37.13	27.00	54.40	36.58	56.35	31.85	72.39	<b>17.10</b>	47.18
		LOOP	85.07	68.37	54.01	38.64	<u>64.98</u>	<b>78.99</b>	<u>82.25</u>	<b>73.10</b>	62.54	42.86	74.12	12.42	61.45
	ALUP	<u>97.62</u>	<u>73.16</u>	42.10	41.11	52.24	64.49	<u>81.33</u>	55.05	<u>63.89</u>	<u>54.76</u>	70.54	6.76	58.59	
	0.75	DAL	86.18	49.69	20.61	42.34	51.78	59.32	72.17	59.64	49.21	53.11	74.11	16.65	52.90
		GeoID	97.36	67.08	<u>55.25</u>	49.61	58.42	79.24	75.61	58.25	56.12	22.35	61.50	3.85	57.05
		SDC	97.87	<b>86.28</b>	<b>60.78</b>	<b>72.30</b>	<b>69.79</b>	<b>83.30</b>	72.77	<b>81.17</b>	<b>77.28</b>	60.90	<b>76.29</b>	18.15	<b>71.41</b>
		DPN	93.14	36.54	26.33	30.16	35.83	28.74	57.90	37.10	45.46	25.56	68.01	17.30	41.84
		TAN	<b>99.73</b>	70.14	42.50	<u>70.22</u>	38.96	31.89	55.93	34.83	59.29	25.56	70.55	18.20	51.48
		LOOP	82.55	66.49	53.25	53.54	66.27	77.14	<b>83.96</b>	69.80	62.91	<u>62.01</u>	72.66	<b>20.69</b>	64.27
	ALUP	<u>98.57</u>	45.70	48.85	52.37	<u>66.45</u>	68.29	72.98	<u>70.62</u>	<u>65.82</u>	27.06	69.03	3.85	57.47	

Table S19: N-ACC performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.1)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
N-ACC	0.25	DAL	70.21	37.92	19.37	30.88	48.09	48.21	71.89	47.72	46.98	50.43	76.32	17.99	47.17
		GeoID	89.64	69.22	45.56	36.89	63.32	75.75	76.78	61.78	50.98	40.80	73.42	19.94	58.67
		SDC	69.80	62.22	41.87	46.80	41.85	72.90	58.74	60.09	59.67	41.04	63.36	26.64	53.75
		DPN	75.27	43.02	19.35	28.00	24.63	4.92	42.78	24.08	45.52	23.58	66.60	25.96	35.31
		TAN	86.13	57.64	24.99	33.09	35.03	19.82	53.39	29.27	55.69	29.66	73.98	23.47	43.51
		LOOP	87.77	67.01	52.95	38.15	65.86	77.51	83.71	70.64	61.28	59.28	74.73	19.98	63.24
		ALUP	96.25	71.42	53.21	37.65	65.31	77.01	81.24	68.94	66.97	59.20	77.79	20.25	64.60
	0.5	DAL	92.73	46.22	29.74	38.83	53.55	60.42	73.73	59.42	53.18	56.36	77.16	10.18	54.29
		GeoID	95.79	69.83	52.10	42.67	56.60	74.92	71.91	64.22	54.61	43.45	69.07	4.05	58.27
		SDC	88.55	83.26	60.66	61.26	60.51	84.51	67.93	73.78	76.99	66.86	71.45	15.75	67.63
		DPN	84.93	44.20	23.58	30.20	28.00	12.22	44.83	27.91	51.64	18.70	70.38	13.14	37.48
		TAN	94.67	60.88	28.56	33.72	37.26	32.29	53.80	33.13	51.46	29.56	73.35	15.10	45.31
		LOOP	83.12	57.30	55.26	42.69	66.11	77.10	85.57	72.55	65.42	64.25	75.78	14.58	63.31
		ALUP	96.76	68.87	52.13	46.11	65.83	80.14	81.68	71.68	67.06	61.93	76.83	13.45	65.21
	0.75	DAL	94.74	52.14	37.85	45.39	52.89	59.27	74.74	62.07	61.09	56.43	77.19	17.99	57.65
		GeoID	98.60	49.63	23.50	52.63	57.37	68.56	72.46	59.65	56.96	11.76	65.93	0.00	51.42
		SDC	94.53	91.18	70.40	64.57	76.60	88.55	79.33	85.26	85.40	75.53	77.78	17.74	75.57
		DPN	94.40	37.65	31.45	32.85	32.87	25.65	55.51	33.11	44.06	24.48	67.95	16.52	41.37
		TAN	95.20	55.93	36.78	42.63	39.08	54.91	55.86	38.80	48.86	36.83	72.82	17.80	49.63
		LOOP	83.12	73.20	47.40	55.59	63.48	78.82	82.74	72.97	62.78	64.43	75.58	18.03	65.63
		ALUP	94.42	79.15	53.05	51.55	67.76	81.16	85.16	70.31	64.74	52.83	77.23	19.91	66.44

Table S20: N-ACC performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.5)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
N-ACC	0.25	DAL	64.23	34.86	18.85	33.60	49.62	52.47	70.66	50.61	52.24	48.31	77.08	16.86	47.45
		GeoID	83.44	69.71	47.56	38.16	55.43	75.75	75.42	68.15	59.02	36.00	69.38	21.34	58.28
		SDC	67.75	72.13	48.14	44.55	42.91	76.40	61.03	60.64	67.87	48.13	61.53	25.20	56.36
		DPN	72.77	36.97	16.81	27.15	23.27	4.48	42.63	20.02	44.67	24.29	67.69	20.60	33.45
		TAN	89.41	56.95	28.53	31.36	32.38	15.83	50.82	29.61	50.47	27.84	73.03	22.12	42.36
		LOOP	83.89	62.37	52.45	38.57	64.83	79.16	84.23	70.55	62.28	65.46	74.67	19.42	63.16
		ALUP	98.45	70.39	53.84	38.37	67.54	79.56	83.56	70.18	69.48	60.78	78.54	18.06	65.73
	0.5	DAL	88.41	38.98	26.64	41.26	54.78	65.10	72.41	60.16	52.50	50.92	78.04	9.89	53.26
		GeoID	99.11	64.27	45.00	44.28	56.03	80.77	74.84	71.17	54.89	41.67	67.30	5.41	58.73
		SDC	90.04	86.75	62.14	67.88	64.04	86.57	67.95	77.71	79.75	63.22	68.04	16.07	69.18
		DPN	84.60	43.13	31.62	29.07	24.90	11.51	51.00	29.72	54.13	25.45	68.68	16.48	39.19
		TAN	96.42	55.69	32.48	32.13	39.63	33.10	55.09	33.52	53.79	30.28	74.56	16.26	46.08
		LOOP	87.93	70.42	50.88	45.69	68.53	83.04	82.88	75.16	62.37	59.84	75.63	13.50	64.66
		ALUP	99.79	70.12	51.58	41.92	69.41	76.63	84.68	73.71	62.72	65.12	77.71	16.66	65.84
	0.75	DAL	82.91	44.67	34.35	49.24	51.39	56.18	70.14	58.31	57.63	53.16	75.88	21.20	54.59
		GeoID	100.00	72.78	50.50	51.09	64.34	77.18	82.98	68.07	54.43	30.59	65.04	0.00	59.75
		SDC	95.05	89.62	71.75	57.57	79.42	89.32	80.32	88.98	86.41	79.54	77.70	16.24	75.99
		DPN	89.76	32.76	31.15	25.97	32.32	22.68	54.35	31.11	46.06	30.01	70.87	21.46	40.71
		TAN	95.24	49.69	34.44	38.74	39.46	53.11	54.18	37.69	50.01	32.17	72.29	18.90	47.99
		LOOP	95.54	68.82	41.75	53.76	64.95	78.63	84.98	73.85	63.56	63.90	74.74	20.72	65.43
		ALUP	99.78	66.61	43.20	47.79	69.18	78.71	84.70	73.05	66.64	59.38	76.89	21.54	65.62

Table S21: N-ACC performance of GCD on the BOLT benchmark under different KCRs (LAR = 1.0)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
H-Score	0.25	DAL	73.44	48.88	27.74	33.33	52.37	55.88	74.76	52.38	50.00	49.19	77.09	14.83	50.82
		GeoID	87.83	74.47	51.04	42.04	57.25	73.29	80.20	64.46	48.33	43.86	74.63	18.00	59.62
		SDC	68.73	64.10	33.36	40.66	38.87	66.59	63.38	59.21	55.95	36.77	63.90	27.38	51.57
		DPN	83.25	60.57	25.55	40.71	35.27	15.09	55.66	44.19	55.47	37.12	74.12	18.07	45.42
		TAN	87.76	66.33	30.40	45.12	40.18	23.32	60.87	47.99	58.41	40.13	73.40	18.04	49.33
		LOOP	90.16	76.45	57.24	40.51	70.22	79.04	87.86	73.91	66.85	45.94	73.76	18.61	65.05
		ALUP	88.66	83.75	59.96	43.55	66.25	77.22	83.62	69.56	68.62	55.83	77.36	18.45	66.07
	0.5	DAL	91.34	55.27	33.46	40.47	57.01	70.57	78.38	59.73	51.47	57.19	77.33	10.74	56.91
		GeoID	95.03	77.52	53.10	44.20	64.29	77.38	82.52	70.78	58.80	43.88	77.10	13.98	63.22
		SDC	88.06	79.93	52.72	54.84	53.15	76.34	67.85	68.65	67.25	57.69	66.57	15.38	62.37
		DPN	87.90	47.29	25.48	36.65	46.03	23.31	66.31	42.22	58.02	33.73	75.18	16.31	46.54
		TAN	98.50	72.11	40.22	57.12	50.40	40.44	67.55	48.98	65.80	44.92	79.21	20.69	57.16
		LOOP	89.98	78.80	58.49	42.82	71.75	82.80	87.66	76.72	64.25	45.96	77.06	13.85	65.84
		ALUP	93.21	80.81	53.19	48.00	63.60	72.75	87.00	64.62	67.27	64.73	75.90	11.00	65.17
	0.75	DAL	89.83	61.91	30.58	42.72	58.12	66.47	78.91	64.76	53.72	60.71	76.75	14.04	58.21
		GeoID	97.10	78.30	59.42	50.95	67.80	80.15	82.50	67.24	61.21	31.56	71.21	6.57	62.83
		SDC	94.56	85.87	59.75	67.49	68.75	83.25	70.72	78.73	76.42	61.85	74.72	23.37	70.46
		DPN	95.99	48.05	32.32	43.05	48.57	40.25	70.44	49.57	57.94	36.95	75.04	22.61	51.73
		TAN	99.84	76.77	50.50	72.74	52.35	45.60	69.04	47.81	68.28	38.22	77.99	22.58	60.14
		LOOP	88.59	76.90	57.35	48.11	73.67	81.80	88.45	74.54	64.26	63.49	76.77	17.75	67.64
		ALUP	96.53	61.47	56.88	56.55	72.72	75.73	81.24	75.99	69.25	40.95	76.10	6.73	64.18

Table S22: H-Score performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.1)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
H-Score	0.25	DAL	77.93	51.19	27.99	37.08	55.47	58.15	78.42	56.29	52.28	57.05	77.76	16.00	53.80
		GeoID	<u>94.54</u>	<b>80.38</b>	48.35	46.27	71.00	80.77	81.93	73.46	58.15	54.65	79.70	19.62	65.73
		SDC	57.75	66.29	41.37	45.12	38.44	72.39	58.36	57.28	61.09	43.15	63.74	<u>22.50</u>	52.29
		DPN	81.14	57.54	29.08	40.52	37.39	9.31	56.28	36.32	59.31	34.15	74.91	<b>25.62</b>	45.13
		TAN	90.62	70.71	35.40	46.30	49.50	31.73	67.37	43.11	66.61	43.30	<u>80.35</u>	20.12	53.76
		LOOP	89.70	78.86	<u>60.75</u>	<u>46.80</u>	<b>75.23</b>	<b>83.05</b>	<b>89.66</b>	<b>78.66</b>	<u>68.77</u>	<u>69.22</u>	<u>80.27</u>	21.11	<u>70.17</u>
	ALUP	<b>96.93</b>	<u>79.29</u>	<b>61.65</b>	<b>49.48</b>	<u>74.41</u>	<u>82.62</u>	<u>87.52</u>	<u>76.97</u>	<b>69.37</b>	<b>69.64</b>	<b>81.06</b>	21.87	<b>70.90</b>	
	0.5	DAL	95.57	59.24	39.73	43.85	62.42	69.06	80.53	67.03	58.46	65.80	79.30	12.36	61.11
		GeoID	94.02	<u>79.54</u>	<b>61.39</b>	50.81	68.47	81.64	81.67	75.15	64.09	57.98	78.56	7.09	66.70
		SDC	84.02	<b>80.40</b>	53.89	<b>58.44</b>	55.46	82.29	61.93	71.43	<u>72.75</u>	68.80	68.19	16.81	64.53
		DPN	89.53	58.95	32.97	44.07	41.76	20.96	60.32	41.70	64.37	30.22	78.48	<u>17.45</u>	48.40
		TAN	<u>95.95</u>	73.41	38.66	48.36	51.78	46.17	68.54	47.46	65.60	43.73	<b>80.98</b>	<b>18.17</b>	56.57
		LOOP	90.04	70.19	<u>60.55</u>	48.26	<u>75.92</u>	<u>82.96</u>	<b>90.76</b>	<u>79.29</u>	70.36	<b>73.22</b>	<u>81.67</u>	16.43	<u>69.97</u>
	ALUP	<b>97.09</b>	79.34	59.88	<u>55.52</u>	<b>76.02</b>	<b>85.19</b>	<b>88.13</b>	<b>79.47</b>	<b>74.12</b>	<u>72.40</u>	<b>82.04</b>	16.83	<b>72.17</b>	
	0.75	DAL	95.90	63.78	45.75	44.45	62.20	67.77	82.37	68.73	65.20	65.95	80.03	15.40	63.13
		GeoID	<b>98.83</b>	65.28	34.95	54.26	69.89	77.44	81.90	71.39	65.88	20.76	76.24	0.00	59.74
		SDC	93.95	<b>88.97</b>	<b>67.35</b>	<b>62.73</b>	74.35	<b>87.56</b>	75.40	<b>83.07</b>	<b>82.12</b>	<b>74.76</b>	75.17	21.43	<b>73.90</b>
		DPN	96.75	51.95	41.80	46.69	47.51	39.71	69.67	47.54	58.59	37.26	77.56	21.14	53.01
TAN		<u>97.05</u>	68.20	47.01	55.87	54.12	67.11	70.19	53.32	62.71	50.44	81.05	<b>21.47</b>	60.71	
LOOP		94.19	82.30	55.70	56.03	<u>74.59</u>	<u>83.55</u>	<u>89.15</u>	<u>79.61</u>	70.26	<u>72.33</u>	<u>81.14</u>	19.21	71.51	
ALUP	96.76	<u>85.60</u>	<u>59.92</u>	<u>59.62</u>	<b>77.43</b>	<u>85.86</u>	<b>90.16</b>	78.04	<u>74.37</u>	64.27	<b>82.96</b>	21.35	<u>73.03</u>		

Table S23: H-Score performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.5)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
H-Score	0.25	DAL	73.50	47.42	28.30	39.08	58.38	55.02	76.80	59.74	56.63	54.27	78.29	16.00	53.62
		GeoID	90.97	<b>80.48</b>	<u>59.93</u>	<u>49.51</u>	65.31	80.40	83.84	76.63	65.82	50.83	79.21	22.14	67.09
		SDC	67.02	75.44	43.54	44.47	39.49	74.66	61.49	60.04	67.19	51.54	61.42	22.22	55.71
		DPN	82.49	51.92	25.11	38.59	35.84	8.52	56.40	31.66	58.09	37.02	76.11	<b>23.70</b>	43.79
		TAN	<u>93.50</u>	70.45	39.10	45.20	47.49	26.83	66.16	43.67	63.59	42.34	79.74	18.74	53.07
		LOOP	90.33	75.59	<b>61.83</b>	48.67	<u>75.73</u>	<u>84.58</u>	<b>90.48</b>	<b>79.71</b>	68.39	<b>73.98</b>	<u>80.37</u>	19.60	<u>70.77</u>
	ALUP	<b>98.30</b>	<u>79.56</u>	59.92	<b>51.56</b>	<b>77.31</b>	<b>84.78</b>	<u>89.21</u>	<u>78.59</u>	<b>74.18</b>	<u>71.11</u>	<b>83.04</b>	<u>23.38</u>	<b>72.58</b>	
	0.5	DAL	90.18	52.89	35.68	45.40	62.42	70.91	80.64	67.51	57.02	61.83	79.61	12.22	59.69
		GeoID	<u>99.55</u>	77.36	56.08	53.91	69.25	<u>85.09</u>	84.75	78.99	65.94	56.37	77.76	9.18	67.85
		SDC	85.44	<b>85.06</b>	56.09	<b>65.24</b>	61.20	84.84	64.35	75.46	<b>76.89</b>	64.75	64.52	17.38	66.77
		DPN	89.70	56.60	41.57	43.30	38.67	19.59	65.78	43.59	67.08	38.94	77.82	<u>20.76</u>	50.28
		TAN	98.10	70.01	41.76	45.92	54.56	46.79	69.45	48.15	67.30	45.12	82.17	20.30	57.47
		LOOP	92.30	<u>80.78</u>	<b>59.90</b>	53.55	<u>78.74</u>	<b>86.26</b>	<u>89.48</u>	<b>81.94</b>	72.74	69.86	81.95	17.04	<u>72.05</u>
	ALUP	<b>99.85</b>	80.55	<u>59.11</u>	<u>53.97</u>	<b>79.20</b>	83.47	<b>90.79</b>	<u>80.54</u>	<u>73.75</u>	<b>74.45</b>	<b>83.02</b>	<b>21.43</b>	<b>73.34</b>	
	0.75	DAL	87.57	59.12	44.69	48.04	61.71	64.90	79.72	66.91	63.09	64.90	78.63	19.21	61.54
		GeoID	<b>100.00</b>	<u>81.83</u>	<u>56.79</u>	<u>57.93</u>	75.59	83.42	88.69	78.23	66.15	45.79	75.86	0.00	67.52
		SDC	94.50	<b>90.73</b>	<b>69.79</b>	53.18	<u>77.94</u>	<b>88.94</b>	77.96	<b>86.58</b>	<b>83.95</b>	<b>79.51</b>	75.46	<u>23.57</u>	<b>75.18</b>
		DPN	93.22	43.99	40.78	39.88	47.55	33.90	69.07	45.78	60.85	44.31	79.86	<b>25.21</b>	52.03
TAN		97.23	64.25	44.99	52.26	54.84	66.32	69.09	52.52	64.27	46.89	<u>81.22</u>	22.40	59.69	
LOOP		96.68	78.35	51.30	<b>59.46</b>	76.30	84.40	<b>90.54</b>	<u>81.00</u>	73.47	<u>73.10</u>	80.96	22.72	<u>72.36</u>	
ALUP	<u>99.40</u>	78.08	53.07	56.98	<b>79.41</b>	<u>84.72</u>	<u>90.51</u>	80.63	<u>76.93</u>	69.88	<b>83.16</b>	22.28	<u>72.92</u>		

Table S24: H-Score performance of GCD on the BOLT benchmark under different KCRs (LAR = 1.0)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
ARI	0.25	DAL	63.10	33.83	8.76	18.46	37.74	38.46	62.57	36.89	36.36	<u>33.77</u>	<u>73.93</u>	3.52	37.28
		GeoID	77.00	60.78	26.76	27.40	41.24	51.27	69.57	50.80	38.46	26.88	68.56	7.00	45.48
		SDC	58.80	44.50	12.32	<u>36.35</u>	29.07	41.49	48.76	45.28	41.73	17.17	51.16	<b>17.47</b>	37.01
		DPN	76.82	44.23	8.82	<u>29.83</u>	23.84	12.74	39.76	26.63	40.70	19.43	64.50	3.46	32.56
		TAN	83.39	51.48	11.35	<b>61.21</b>	27.45	16.54	46.06	29.46	45.93	22.05	66.10	<u>15.77</u>	39.73
		LOOP	<u>86.91</u>	<u>61.33</u>	<u>34.04</u>	23.63	<b>56.95</b>	<b>63.15</b>	<b>79.66</b>	<b>61.23</b>	<u>51.48</u>	26.55	68.08	7.09	<u>51.68</u>
	ALUP	82.20	<b>70.67</b>	<u>32.99</u>	26.99	<b>53.30</b>	<b>58.86</b>	<u>75.39</u>	<u>55.79</u>	<b>57.22</b>	<b>34.32</b>	<u>72.66</u>	7.74	<b>52.34</b>	
	0.5	DAL	90.60	48.89	18.64	24.53	47.12	56.41	71.46	48.38	41.76	<u>47.40</u>	76.10	4.20	47.96
		GeoID	92.38	<u>69.13</u>	33.61	32.05	55.16	60.82	77.10	59.72	45.65	25.20	74.93	9.47	52.94
		SDC	84.26	65.86	30.18	<u>48.42</u>	45.30	59.37	59.85	55.90	<b>59.41</b>	43.76	60.47	<b>18.80</b>	52.63
		DPN	87.81	47.52	18.04	39.65	42.07	26.88	60.90	36.27	55.98	31.45	70.93	12.80	44.19
		TAN	<b>98.23</b>	64.96	23.52	<b>70.14</b>	44.06	37.46	62.34	37.14	58.93	35.64	78.45	<u>15.64</u>	52.21
		LOOP	90.90	68.82	<u>36.49</u>	25.91	<u>62.21</u>	<b>70.10</b>	<b>82.91</b>	<b>66.85</b>	49.92	27.65	74.20	6.75	<u>55.23</u>
	ALUP	<u>93.71</u>	<b>74.79</b>	<u>34.35</u>	39.24	<u>55.20</u>	57.84	<u>81.28</u>	55.32	54.16	<b>49.07</b>	<u>76.67</u>	7.82	<b>56.62</b>	
	0.75	DAL	93.69	65.48	32.66	31.81	52.71	59.26	77.05	56.31	47.37	<u>58.65</u>	77.39	6.81	54.93
		GeoID	97.26	72.94	39.00	39.49	65.55	64.44	82.10	61.77	54.20	29.75	78.48	10.68	57.97
		SDC	92.05	74.56	37.89	<u>65.28</u>	59.96	69.02	70.46	66.31	<u>69.40</u>	55.72	71.48	<b>29.52</b>	63.47
		DPN	98.34	65.05	33.04	60.66	54.47	43.04	74.85	51.59	<b>69.62</b>	46.30	76.69	18.29	57.66
TAN		<b>99.82</b>	72.39	35.58	<b>72.22</b>	58.45	54.31	74.36	51.39	69.34	47.96	<b>81.52</b>	<u>20.67</u>	61.50	
LOOP		95.76	<u>74.76</u>	<u>40.15</u>	31.58	<b>69.42</b>	<b>74.25</b>	<b>86.93</b>	<u>68.53</u>	50.79	53.06	76.65	8.51	60.87	
ALUP	94.39	71.03	<b>40.92</b>	48.14	67.00	67.09	80.88	<b>69.33</b>	58.99	<u>56.40</u>	80.03	10.92	<u>62.09</u>		

Table S25: ARI performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.1)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
ARI	0.25	DAL	70.64	34.50	9.06	20.73	41.24	41.04	67.23	39.55	38.78	41.09	75.18	4.80	40.32
		GeoID	88.77	64.78	26.86	25.75	57.09	62.76	73.54	57.23	48.63	34.52	73.84	7.96	51.81
		SDC	54.78	45.82	18.10	36.89	31.42	52.56	46.77	45.83	47.35	26.66	53.32	18.37	39.82
		DPN	77.10	42.21	11.13	30.98	27.31	12.78	40.65	19.70	41.98	19.78	62.97	9.03	32.97
		TAN	87.04	55.20	15.05	30.15	36.21	22.43	52.48	27.95	50.08	25.71	72.75	17.87	41.08
		LOOP	87.88	62.72	35.05	25.81	61.56	69.68	81.65	65.41	49.54	49.93	74.58	7.86	55.97
	ALUP	96.32	66.15	33.66	32.25	59.72	65.20	78.53	62.35	59.54	49.02	77.73	8.20	57.39	
	0.5	DAL	95.32	50.05	23.39	29.66	53.23	57.09	74.66	56.42	44.86	59.81	78.42	5.50	52.37
		GeoID	92.40	71.10	39.23	37.77	60.97	70.53	77.63	65.73	55.30	49.26	78.88	11.24	59.17
		SDC	82.49	67.35	33.51	54.02	51.74	67.51	58.38	60.08	65.09	53.84	61.72	18.84	56.21
		DPN	90.58	55.96	22.14	45.85	44.39	37.78	59.37	39.95	60.92	38.27	74.14	13.02	48.53
		TAN	95.77	66.85	26.00	48.07	51.44	44.57	66.22	44.88	61.09	40.74	77.52	15.70	53.24
		LOOP	90.87	63.74	37.53	32.41	68.80	73.78	86.26	71.99	59.87	62.22	78.80	8.57	61.24
	ALUP	97.26	71.63	37.04	46.10	68.23	72.58	83.88	69.68	63.88	61.46	80.66	8.72	63.43	
	0.75	DAL	96.72	66.86	33.44	32.74	61.30	63.17	82.46	63.59	53.27	67.52	80.83	7.42	59.11
		GeoID	99.14	75.14	39.41	39.49	72.90	73.96	85.44	70.19	65.40	57.53	83.62	12.30	64.54
		SDC	94.00	78.53	42.68	75.53	68.77	75.62	74.56	72.29	74.46	65.93	73.51	31.45	68.94
		DPN	98.86	67.41	37.62	69.90	65.32	59.84	78.32	60.46	72.26	56.64	80.28	24.39	64.28
TAN		98.43	73.18	38.66	65.18	66.71	66.11	79.64	62.65	70.72	59.60	83.21	22.34	65.54	
LOOP		95.14	78.56	42.47	42.79	76.48	77.61	89.76	76.33	65.99	70.22	81.93	10.30	67.30	
ALUP	98.52	80.70	43.73	64.98	77.12	77.82	89.82	75.26	74.68	68.35	83.65	12.69	70.61		

Table S26: ARI performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.5)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
ARI	0.25	DAL	65.78	28.97	9.78	21.78	42.86	40.18	65.67	43.78	42.08	38.45	75.80	5.82	40.08
		GeoID	85.53	64.28	31.55	30.11	50.77	63.73	72.75	62.13	51.10	32.51	73.56	8.98	52.25
		SDC	58.37	59.58	21.13	37.23	31.91	55.58	49.22	46.77	56.53	33.85	49.51	17.16	43.07
		DPN	76.21	39.03	9.70	26.80	27.59	13.25	41.19	19.58	42.50	20.70	64.23	6.88	32.30
		TAN	89.88	54.99	17.20	31.70	35.80	20.92	51.17	28.70	46.23	26.44	71.57	17.68	41.02
		LOOP	86.41	58.44	34.56	26.46	62.10	71.96	82.48	66.44	52.06	55.23	74.11	8.58	56.57
	ALUP	98.27	63.56	32.71	32.89	62.14	68.20	80.97	64.17	64.00	50.90	78.78	9.56	58.85	
	0.5	DAL	90.81	46.04	20.92	29.84	53.90	56.65	74.99	56.63	44.64	55.87	78.04	5.58	51.16
		GeoID	99.08	70.63	35.66	44.84	64.39	72.97	80.15	70.50	55.23	48.84	78.89	12.02	61.10
		SDC	84.70	74.14	35.55	58.56	56.30	71.50	60.96	64.90	69.52	51.73	60.20	20.03	59.01
		DPN	90.57	56.41	24.86	46.58	44.35	37.42	64.50	42.45	63.74	41.37	72.69	15.41	50.03
		TAN	98.26	62.06	26.87	38.72	53.15	45.84	67.49	46.33	60.00	44.81	79.38	14.99	53.16
		LOOP	93.32	71.31	38.93	39.95	71.40	76.61	85.76	74.10	66.31	62.27	79.47	9.34	64.06
	ALUP	99.70	69.47	36.13	48.42	71.66	72.88	86.79	71.98	68.69	62.29	81.26	13.15	65.20	
	0.75	DAL	90.72	64.86	38.10	39.31	62.10	59.88	82.22	63.90	53.05	69.63	79.74	9.73	59.44
		GeoID	100.00	75.92	41.95	56.27	76.65	78.16	88.60	75.60	68.54	61.79	84.54	17.19	68.77
		SDC	94.55	84.66	45.09	81.30	73.28	78.05	76.32	76.48	77.55	70.48	74.83	31.06	71.97
		DPN	95.66	69.71	39.71	69.60	67.77	60.24	79.30	61.33	73.33	60.68	82.68	27.67	65.64
TAN		98.99	72.16	40.01	69.86	68.89	62.08	79.96	64.00	72.50	61.36	84.16	21.05	66.67	
LOOP		97.42	77.28	42.50	54.53	78.88	79.44	90.95	78.55	75.23	72.06	82.46	15.89	70.43	
ALUP	99.25	78.31	43.08	61.35	79.64	78.92	90.76	77.92	77.32	70.65	84.92	15.03	71.43		

Table S27: ARI performance of GCD on the BOLT benchmark under different KCRs (LAR = 1.0)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
NMI	0.25	DAL	80.41	49.71	14.10	56.76	70.47	51.67	88.24	71.80	67.14	58.07	92.33	26.64	60.61
		GeoID	85.74	72.83	34.46	63.61	72.13	68.22	91.24	78.25	67.77	51.59	90.75	33.54	67.51
		SDC	79.04	63.37	27.16	60.78	67.44	61.87	84.86	76.85	68.51	42.40	86.88	36.12	62.94
		DPN	90.22	61.47	15.62	59.58	61.79	27.63	79.68	64.10	65.82	43.42	91.02	24.00	57.03
		TAN	91.14	65.69	17.76	69.36	62.79	27.84	81.10	65.26	66.07	46.61	90.78	27.47	59.32
		LOOP	94.17	73.53	42.32	65.33	81.89	74.77	94.11	83.90	80.19	50.91	90.71	34.62	72.20
	ALUP	90.59	79.30	41.74	64.41	79.79	71.73	92.87	81.59	79.60	57.73	92.16	35.40	72.24	
	0.5	DAL	96.00	63.23	25.37	63.10	76.23	67.20	91.50	77.55	72.46	69.28	92.97	27.31	68.52
		GeoID	95.61	78.27	40.44	69.33	80.24	73.80	93.32	82.98	76.35	49.95	92.56	36.16	72.42
		SDC	92.78	76.10	38.40	70.16	76.07	69.17	88.98	81.20	77.02	65.74	89.25	34.99	71.66
		DPN	95.78	65.83	27.72	65.56	73.93	48.94	89.20	71.87	73.24	53.90	92.23	28.64	65.57
		TAN	99.10	75.64	30.94	77.34	73.79	48.72	88.04	71.48	76.54	58.20	93.54	30.30	68.64
		LOOP	96.71	78.09	43.90	67.84	84.40	77.56	95.05	86.11	79.55	51.93	92.33	34.13	73.97
	ALUP	97.40	81.61	42.25	71.15	80.83	72.31	94.39	81.63	79.79	67.35	92.96	35.38	74.75	
	0.75	DAL	97.91	75.77	39.46	68.02	79.38	72.44	93.24	81.27	74.38	77.45	93.16	34.21	73.89
		GeoID	98.83	80.56	45.16	72.80	85.06	75.72	94.93	83.99	78.72	52.00	93.47	38.94	75.02
		SDC	96.76	81.09	43.69	77.09	82.55	74.43	91.89	85.65	80.89	72.89	91.88	41.31	76.68
		DPN	99.23	77.11	39.53	74.74	79.66	66.69	92.88	79.64	79.05	66.66	93.46	35.92	73.71
TAN		99.88	80.82	41.39	80.87	81.14	62.65	92.21	78.61	80.30	66.30	94.20	36.23	74.55	
LOOP		98.54	81.39	45.68	72.30	87.10	79.20	96.18	86.98	80.14	71.42	93.13	37.86	77.49	
ALUP	97.93	80.04	46.21	76.90	86.23	76.16	94.88	86.79	82.53	74.90	94.08	37.97	77.89		

Table S28: NMI performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.1)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
NMI	0.25	DAL	85.70	50.92	15.31	60.13	72.95	54.91	89.93	72.80	70.57	64.78	92.78	29.92	63.39
		GeoID	93.19	<u>76.31</u>	35.53	64.02	80.96	73.90	92.58	81.70	76.34	56.98	92.45	35.91	71.66
		SDC	77.36	64.61	32.59	63.63	70.50	67.63	84.97	77.73	72.92	53.42	88.21	36.60	65.85
		DPN	90.49	60.21	17.00	62.78	66.55	27.62	82.56	61.17	68.60	46.71	91.53	28.53	58.65
		TAN	94.30	68.20	22.56	63.57	68.80	33.30	84.23	63.97	74.20	50.08	92.47	30.39	62.17
		LOOP	<u>95.09</u>	74.45	<b>42.99</b>	<u>67.82</u>	<b>84.01</b>	<b>77.18</b>	<b>94.69</b>	<b>85.61</b>	<u>80.08</u>	<b>69.08</b>	92.56	36.70	<u>75.02</u>
	ALUP	<b>98.09</b>	<b>77.73</b>	<u>41.88</u>	<b>67.97</b>	<u>83.02</u>	<u>74.61</u>	<u>93.77</u>	<u>84.14</u>	<b>80.98</b>	<u>68.36</u>	<b>93.35</b>	<b>37.13</b>	<b>75.09</b>	
	0.5	DAL	98.25	64.49	30.48	68.89	79.45	67.89	92.56	81.48	76.10	<b>79.05</b>	93.46	31.44	71.96
		GeoID	96.85	<u>79.04</u>	<u>44.31</u>	72.29	83.07	77.50	93.70	85.48	79.50	68.75	<u>93.77</u>	<b>38.73</b>	76.08
		SDC	92.78	76.86	40.54	<b>75.66</b>	80.19	75.14	88.63	83.63	81.42	71.57	89.74	36.36	74.38
		DPN	96.81	71.28	30.93	70.61	77.59	49.65	88.68	74.26	77.27	61.39	93.42	32.03	68.66
		TAN	<u>98.25</u>	76.46	32.47	72.23	78.00	54.56	89.76	74.92	78.20	61.67	93.58	32.63	70.23
		LOOP	96.48	75.05	<b>44.41</b>	72.37	<b>87.25</b>	<b>79.64</b>	<b>95.93</b>	<b>88.35</b>	<u>83.03</u>	<u>77.00</u>	93.75	<u>37.73</u>	<u>77.58</u>
	ALUP	<b>99.02</b>	<b>79.98</b>	<u>43.88</u>	<u>74.12</u>	<u>86.70</u>	<u>77.96</u>	<u>95.39</u>	<u>87.10</u>	<b>83.37</b>	76.83	<b>94.15</b>	37.48	<b>78.00</b>	
	0.75	DAL	98.82	76.68	40.97	70.91	83.53	74.46	94.93	84.62	80.34	<b>82.45</b>	94.19	35.64	76.46
		GeoID	<b>99.71</b>	83.52	45.50	76.85	88.49	80.41	96.12	87.38	84.08	74.93	94.89	<u>41.31</u>	79.43
		SDC	98.10	<u>83.92</u>	<u>47.75</u>	<b>83.02</b>	87.40	79.65	93.19	88.63	85.97	79.30	92.37	<b>42.22</b>	80.13
		DPN	<u>99.58</u>	78.79	44.36	80.49	84.99	67.77	94.25	83.49	83.64	74.15	94.74	40.25	77.21
TAN		99.38	81.30	43.99	79.10	85.23	72.79	94.18	84.04	83.36	75.15	<b>95.07</b>	38.92	77.71	
LOOP		98.41	83.62	<u>47.85</u>	<u>77.24</u>	<u>90.13</u>	<b>81.67</b>	<u>97.02</u>	<b>89.90</b>	85.08	<u>82.30</u>	94.61	38.47	<u>80.52</u>	
ALUP	99.18	<b>85.14</b>	<b>48.33</b>	<u>80.56</u>	<b>90.15</b>	<u>81.46</u>	<b>97.08</b>	<u>89.49</u>	<b>86.54</b>	81.02	<u>94.97</u>	39.71	<b>81.14</b>		

Table S29: NMI performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.5)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
NMI	0.25	DAL	83.99	45.43	16.39	59.96	73.79	54.03	89.52	75.13	72.38	62.58	92.78	32.01	63.17
		GeoID	92.33	<u>75.34</u>	39.24	67.61	78.09	72.46	92.07	84.07	84.07	56.17	92.42	35.98	71.97
		SDC	80.09	72.31	35.78	64.22	70.82	70.48	85.69	78.58	78.29	58.32	87.38	34.21	68.01
		DPN	90.50	58.02	16.63	61.49	66.78	27.53	82.42	59.97	66.73	46.72	90.83	29.43	58.09
		TAN	<u>95.24</u>	<u>67.57</u>	24.33	64.43	68.74	32.31	83.70	64.54	71.17	49.26	92.35	29.96	61.97
		LOOP	94.26	72.02	<b>43.13</b>	<b>68.41</b>	<b>84.53</b>	<b>78.48</b>	<b>94.98</b>	<b>86.31</b>	<u>80.31</u>	<b>72.31</b>	92.61	<b>37.54</b>	<u>75.41</u>
	ALUP	<b>98.96</b>	<b>75.84</b>	<u>41.35</u>	<u>68.34</u>	<u>84.18</u>	<u>76.19</u>	<u>94.48</u>	<u>84.91</u>	<b>82.68</b>	<u>69.28</u>	<b>93.63</b>	<u>37.50</u>	<b>75.61</b>	
	0.5	DAL	96.78	60.51	29.84	68.64	80.34	69.87	92.73	81.71	75.90	76.19	93.50	32.08	71.51
		GeoID	99.25	<u>79.73</u>	42.40	73.12	84.81	<u>78.58</u>	94.34	87.50	80.95	70.08	93.86	<b>40.09</b>	77.06
		SDC	93.96	<b>80.83</b>	43.06	<b>77.66</b>	82.86	77.65	89.32	85.90	<b>84.71</b>	71.67	89.05	37.28	76.16
		DPN	96.67	72.49	35.37	71.35	78.53	51.49	90.53	74.82	78.44	62.78	93.35	33.10	69.91
		TAN	<u>99.31</u>	72.49	33.96	70.35	79.02	55.84	90.30	75.63	78.06	64.25	94.12	33.14	70.54
		LOOP	97.30	79.67	<b>45.87</b>	74.42	<b>88.13</b>	<b>80.91</b>	<u>95.93</u>	<b>89.31</b>	<u>84.39</u>	<b>78.28</b>	93.96	37.82	<b>78.83</b>
	ALUP	<b>99.76</b>	78.78	<u>43.89</u>	<u>75.56</u>	<u>88.12</u>	78.44	<b>96.17</b>	<u>88.12</u>	83.46	<u>76.96</u>	<b>94.22</b>	<u>40.02</u>	<u>78.62</u>	
	0.75	DAL	96.46	75.64	43.98	72.64	84.12	73.56	95.10	84.74	80.53	<b>84.29</b>	94.02	38.01	76.92
		GeoID	<b>100.00</b>	83.58	46.99	79.04	90.18	81.40	97.02	89.87	85.49	77.47	95.28	42.86	80.77
		SDC	98.20	<b>87.32</b>	<b>50.01</b>	<b>85.96</b>	89.35	81.38	93.68	90.47	<b>88.87</b>	81.95	92.81	<b>43.52</b>	<b>81.96</b>
		DPN	98.69	80.70	45.87	<u>81.72</u>	86.41	68.56	94.89	84.80	84.87	76.47	95.22	41.51	78.31
TAN		99.55	80.59	45.23	81.35	86.47	72.93	94.51	84.66	84.63	75.81	95.34	40.59	78.47	
LOOP		99.11	83.17	47.86	80.72	<b>91.10</b>	<b>82.85</b>	<b>97.48</b>	<b>90.99</b>	87.32	<u>83.68</u>	94.65	42.98	81.83	
ALUP	<u>99.72</u>	<u>84.32</u>	<u>48.26</u>	81.22	<u>91.00</u>	<u>82.42</u>	<u>97.27</u>	<u>90.65</u>	<u>87.68</u>	82.58	<b>95.39</b>	42.92	<u>81.95</u>		

Table S30: NMI performance of GCD on the BOLT benchmark under different KCRs (LAR = 1.0)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
ACC	0.25	DAL	69.62	48.59	25.80	30.60	49.83	53.68	71.70	49.91	46.42	<u>47.22</u>	<u>76.20</u>	12.32	48.49
		GeoID	83.52	<u>74.77</u>	46.15	38.96	53.57	74.04	76.53	62.21	48.37	40.79	72.37	17.00	57.36
		SDC	69.11	63.65	36.65	47.50	41.76	65.42	62.38	59.54	55.24	36.64	63.42	<b>32.07</b>	52.78
		DPN	80.03	59.00	26.02	41.21	35.56	28.05	52.50	42.27	52.47	37.71	72.05	15.83	45.23
		TAN	86.40	65.43	29.69	<b>60.71</b>	39.09	30.96	57.72	45.62	55.62	39.16	71.68	<u>32.02</u>	51.17
		LOOP	<u>88.53</u>	<u>73.22</u>	<u>56.68</u>	37.24	<u>67.33</u>	<u>77.48</u>	<b>85.58</b>	<b>71.92</b>	<u>64.64</u>	44.47	72.52	17.19	<u>63.07</u>
	ALUP	84.53	<b>80.64</b>	<u>55.35</u>	39.96	<u>63.47</u>	<u>75.10</u>	<u>81.56</u>	<u>67.15</u>	<b>66.63</b>	<b>51.96</b>	<u>76.05</u>	18.10	<b>63.38</b>	
	0.5	DAL	91.47	63.04	38.16	38.60	57.96	71.84	78.99	60.82	52.85	60.03	77.88	13.10	58.73
		GeoID	94.96	79.14	53.15	41.02	65.97	77.55	84.04	71.22	58.83	44.41	77.67	19.57	63.96
		SDC	88.72	80.77	55.12	<u>59.28</u>	57.74	76.79	70.52	69.11	<u>69.28</u>	<u>61.30</u>	70.36	<b>34.16</b>	66.10
		DPN	88.60	61.66	38.23	<u>50.27</u>	54.23	48.16	69.92	52.26	64.14	49.45	76.49	26.87	56.69
		TAN	<b>98.56</b>	76.28	45.42	<b>68.89</b>	57.58	53.95	71.75	54.56	68.14	54.06	<b>81.96</b>	<u>29.23</u>	63.37
		LOOP	90.98	<u>80.95</u>	<u>59.27</u>	40.45	<u>72.54</u>	<b>83.12</b>	<b>88.07</b>	<b>76.98</b>	64.11	47.06	77.76	16.77	<u>66.51</u>
	ALUP	93.95	<b>81.70</b>	<u>57.15</u>	51.18	<u>66.56</u>	73.96	<u>87.42</u>	65.99	66.58	<b>66.77</b>	<u>78.20</u>	18.46	<b>67.33</b>	
	0.75	DAL	93.35	76.47	55.61	41.75	62.99	72.04	83.35	68.38	56.88	69.76	79.26	15.71	64.63
		GeoID	97.10	82.41	61.69	50.73	75.26	80.80	86.93	73.64	63.98	45.62	81.19	21.55	68.41
		SDC	93.22	<u>85.97</u>	<u>62.82</u>	<u>70.16</u>	71.64	83.38	78.23	78.49	<b>77.40</b>	<u>71.23</u>	79.21	<b>45.88</b>	74.80
		TAN	98.41	73.81	55.62	64.99	66.73	68.41	81.91	66.19	74.95	63.98	81.24	<u>32.84</u>	69.09
LOOP		<b>99.90</b>	80.67	59.30	<b>74.91</b>	69.94	70.08	81.66	65.42	<u>75.10</u>	65.33	<b>85.22</b>	31.41	71.58	
ALUP		95.67	84.95	62.51	45.10	<b>78.97</b>	<u>85.16</u>	<b>91.14</b>	77.91	64.71	68.78	79.95	19.55	71.20	
ALUP	95.70	80.00	<b>65.15</b>	57.76	<u>76.88</u>	80.82	86.89	<b>79.46</b>	72.30	69.49	<u>82.97</u>	23.08	<u>72.54</u>		

Table S31: ACC performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.1)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
ACC	0.25	DAL	74.67	50.22	25.98	32.06	52.50	54.71	75.54	53.13	50.27	54.73	77.39	14.23	51.29
		GeoID	<u>92.13</u>	<b>76.89</b>	46.75	39.23	67.63	78.46	80.40	68.90	58.78	51.06	77.62	18.88	63.06
		SDC	66.23	65.22	43.88	<b>50.07</b>	42.63	72.91	59.54	59.54	61.85	46.09	65.43	<b>35.26</b>	55.72
		DPN	78.64	56.69	28.64	<u>44.59</u>	38.08	26.45	53.44	38.18	56.05	37.97	72.42	22.06	46.10
		TAN	88.63	67.75	34.07	43.25	47.31	36.38	63.07	42.77	62.15	42.96	77.78	<u>34.44</u>	53.38
		LOOP	89.07	75.44	<u>56.91</u>	40.79	<b>71.27</b>	<b>80.54</b>	<b>86.96</b>	<b>75.23</b>	<u>64.55</u>	<u>65.46</u>	<u>78.28</u>	18.42	<u>66.91</u>
	ALUP	<b>96.66</b>	<u>76.58</u>	<b>57.45</b>	44.37	<u>70.57</u>	<u>80.05</u>	<u>84.74</u>	<u>73.63</u>	<b>68.76</b>	<b>65.89</b>	<b>80.11</b>	19.34	<b>68.18</b>	
	0.5	DAL	95.76	64.98	44.87	42.00	64.10	70.94	81.24	68.21	58.83	69.51	79.97	15.03	62.95
		GeoID	94.15	<b>81.68</b>	<b>63.40</b>	49.39	71.43	82.17	83.20	76.65	66.29	64.95	82.33	22.42	69.84
		SDC	85.06	81.31	56.75	<b>65.18</b>	61.24	82.43	66.86	73.00	<b>75.94</b>	70.82	71.31	<b>35.62</b>	68.79
		DPN	90.11	67.26	41.38	58.69	55.58	50.87	69.03	55.62	69.80	53.04	80.07	26.68	59.84
		TAN	<u>96.07</u>	78.03	47.21	<u>58.86</u>	61.39	59.93	74.37	58.89	71.47	57.16	82.13	<u>27.92</u>	64.45
		LOOP	91.03	75.94	61.46	46.69	<u>77.52</u>	<u>83.66</u>	<b>91.11</b>	<u>80.26</u>	71.65	<u>75.27</u>	<u>82.50</u>	19.97	<u>71.42</u>
	ALUP	<b>97.11</b>	<u>81.64</u>	<u>62.21</u>	55.71	<b>77.82</b>	<b>85.58</b>	<u>88.80</u>	<b>80.46</b>	<u>75.58</u>	<b>75.57</b>	<b>83.32</b>	19.73	<b>73.63</b>	
	0.75	DAL	96.52	76.96	56.67	43.14	71.04	75.20	87.64	74.13	67.57	77.57	82.51	17.62	68.88
		GeoID	<b>98.95</b>	82.19	59.20	54.08	81.49	83.86	88.67	80.81	75.63	68.58	86.50	23.80	73.65
		SDC	94.13	88.15	<u>66.27</u>	<b>80.16</b>	78.18	<u>87.15</u>	80.77	82.97	<b>83.07</b>	<u>80.54</u>	80.71	<b>48.42</b>	<b>79.21</b>
		DPN	98.74	75.72	60.26	<u>75.92</u>	74.84	72.79	84.16	72.75	78.93	70.33	85.57	<u>36.73</u>	73.90
TAN		98.57	82.38	61.69	71.86	76.70	79.49	85.14	74.81	78.76	73.84	<b>87.28</b>	33.77	75.36	
LOOP		95.14	<u>88.35</u>	64.71	55.59	<u>84.38</u>	86.91	<u>93.15</u>	<b>84.38</b>	74.95	<b>82.24</b>	85.21	22.60	76.47	
ALUP	98.72	<b>89.34</b>	<b>66.66</b>	66.65	<b>85.32</b>	<b>88.71</b>	<b>93.16</b>	<u>83.93</u>	<u>81.04</u>	80.26	<u>86.90</u>	25.53	<u>78.85</u>		

Table S32: ACC performance of GCD on the BOLT benchmark under different KCRs (LAR = 0.5)

Metric	KCR	Method	20NG	THUCNews	Yahoo	TREC	BANK	S.O.	CLINC	HWU	ECDT	MCID	DBPedia	X-Topic	Avg.
ACC	0.25	DAL	69.80	46.59	26.59	34.49	54.92	54.62	74.08	56.36	54.28	52.06	78.03	14.51	51.36
		GeoID	87.42	<b>77.16</b>	54.25	42.53	61.36	78.23	80.22	73.08	63.25	48.34	77.25	21.49	63.72
		SDC	67.95	74.29	48.17	<b>49.88</b>	42.96	75.86	61.93	61.10	<u>69.28</u>	52.38	63.30	<b>34.70</b>	58.48
		DPN	78.43	53.81	26.36	41.55	38.19	26.68	53.39	35.56	56.07	38.43	73.12	19.98	45.13
		TAN	<u>91.60</u>	67.79	36.70	44.42	46.44	34.02	62.10	43.42	59.39	43.20	77.46	<u>33.59</u>	53.34
		LOOP	87.42	72.01	<b>57.26</b>	41.51	<u>71.33</u>	<u>82.17</u>	<b>87.66</b>	<b>75.91</b>	65.75	<b>70.51</b>	<u>78.15</u>	20.81	<u>67.54</u>
	ALUP	<b>98.39</b>	<u>76.64</u>	<u>56.79</u>	<u>46.37</u>	<b>73.25</b>	<b>82.37</b>	<u>86.66</u>	<u>75.04</u>	<b>72.53</b>	<u>67.25</u>	<b>81.32</b>	22.48	<b>69.92</b>	
	0.5	DAL	90.33	60.92	41.03	43.30	63.72	71.59	81.71	68.64	57.74	66.99	80.10	15.00	61.76
		GeoID	<u>99.53</u>	80.71	59.70	55.02	73.12	<u>85.38</u>	86.27	79.98	67.20	64.05	82.51	24.01	71.46
		SDC	86.48	<b>85.61</b>	59.53	<b>68.45</b>	65.60	<u>84.99</u>	69.14	76.61	<b>79.27</b>	69.97	69.84	<b>37.46</b>	71.08
		DPN	90.35	68.35	46.85	<u>58.74</u>	55.91	51.09	72.34	58.04	72.12	56.62	79.61	<u>28.36</u>	61.53
		TAN	98.24	75.64	47.79	52.95	63.37	61.10	75.12	60.21	71.91	60.12	<u>83.67</u>	28.20	64.86
		LOOP	93.38	<u>82.80</u>	<b>62.08</b>	53.55	<u>80.40</u>	<b>86.52</b>	<u>90.10</u>	<b>82.68</b>	74.93	<u>73.29</u>	<u>82.90</u>	20.40	<u>73.59</u>
	ALUP	<b>99.85</b>	82.52	<u>61.08</u>	57.47	<b>80.68</b>	84.41	<b>91.27</b>	<u>81.30</u>	<u>76.36</u>	<b>76.56</b>	<b>83.90</b>	24.84	<b>75.02</b>	
	0.75	DAL	91.16	76.46	60.21	46.53	72.13	72.89	86.99	74.03	66.16	78.77	81.04	19.67	68.84
		GeoID	<b>100.00</b>	84.10	62.00	62.53	84.73	87.83	92.13	84.65	76.18	75.53	87.22	28.72	77.14
		SDC	94.74	<b>92.10</b>	<b>69.17</b>	<b>83.84</b>	80.93	<b>88.81</b>	82.27	85.76	<b>84.66</b>	<b>83.75</b>	81.23	<b>48.86</b>	<b>81.34</b>
		DPN	95.71	76.35	61.22	<u>76.74</u>	76.88	73.38	84.71	74.53	80.34	73.84	87.04	<u>40.17</u>	75.08
TAN		98.92	81.80	61.73	75.10	78.65	80.51	85.46	75.72	80.94	74.93	<b>88.17</b>	31.94	76.16	
LOOP		97.27	<u>87.21</u>	64.39	63.76	<u>86.01</u>	88.56	<u>93.91</u>	<b>86.24</b>	80.76	<u>82.90</u>	85.45	28.01	78.71	
ALUP	<u>99.22</u>	87.08	<u>66.42</u>	66.94	<b>87.42</b>	<u>88.73</u>	<b>94.04</b>	<u>86.11</u>	<u>83.88</u>	82.78	<u>87.70</u>	26.53	<u>79.74</u>		

Table S33: ACC performance of GCD on the BOLT benchmark under different KCRs (LAR = 1.0)