EVALUATING THE GENERALIZATION ABILITY OF QUANTIZED LLMS: BENCHMARK, ANALYSIS, AND TOOLBOX

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) have exhibited exciting progress in multiple scenarios, while the huge computational demands hinder their deployments in lots of real-world applications. As an effective means to reduce memory footprint and inference cost, quantization also faces challenges in performance degradation at low bit-widths. Understanding the impact of quantization on LLM capabilities, especially the generalization ability, is crucial. However, the community's main focus remains on the algorithms and models of quantization, with insufficient attention given to to the impact of *data* on the generalization abilities of quantized LLMs. In this work, we fill this gap by providing a comprehensive benchmark suite for this research topic, including an evaluation system, detailed analyses, and a general toolbox. Specifically, based on the dominant pipeline in LLM quantization, we primarily explore the impact of calibration data distribution on the generalization of quantized LLMs and conduct the benchmark using more than 40 datasets within two main scenarios. Based on this benchmark, we conduct extensive experiments with well-known LLMs (LLaMA and Baichuan) and four quantization algorithms to investigate this topic in-depth, yielding several counter-intuitive and valuable findings, e.g., models quantized using a calibration set with the same distribution as the test data are not necessarily optimal. Besides, to facilitate future research, we also release a modular-designed toolbox, which decouples the overall pipeline into several separate components, e.g., base LLM module, dataset module, quantizer module, etc. and allows subsequent researchers to easily assemble their methods through a simple configuration. Our code is submitted in the supplementary materials and will be publicly available.

033 034

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

1 INTRODUCTION

036 037

In recent years, large language models (LLMs) have made groundbreaking advancements, demonstrating remarkable results and outstanding *generalization ability* across various tasks (Vaswani, 2017; Zhang et al., 2022; Achiam et al., 2023; Touvron et al., 2023). For example, given a few prompt ex-040 amples or questions, LLMs can produce insightful answers within the unseen domain (Radford et al., 041 2019; Brown et al., 2020). However, while LLMs exhibit remarkable capabilities, their substantial 042 size makes real-world implementation cost-prohibitive. To address this challenge, model quantization 043 has emerged as a prevailing technique for reducing the memory footprint of LLMs (Frantar et al., 2022; Lin et al., 2023; Dettmers et al., 2023; Chee et al., 2024; Yao et al., 2022; Xiao et al., 2023; 044 Shao et al., 2023). Specifically, quantization reduces the model size by replacing high-precision floating-point numbers with lower-precision integers (e.g., from FP16 to INT4) (Nagel et al., 2021; 046 Gholami et al., 2022; Zhu et al., 2023). Currently, to avoid the substantial retraining costs of LLMs, 047 the quantization methods for large models primarily employ post-training quantization (PTQ) (Frantar 048 et al., 2022; Lin et al., 2023; Dettmers et al., 2023; Chee et al., 2024), which leverages calibration data to optimize the error caused by the quantization. Given the prevalent view that LLM capabilities stem from their extensive parameter count (Kaplan et al., 2020), a critical question emerges:

051 052

Can the quantized LLMs still retain their strong generalization ability?

While some works have acknowledged this issue (Liu et al., 2023a; Jaiswal et al., 2023; Li et al., 2024; Huang et al., 2024a; Jin et al., 2024b), there is still a lack of systematic evaluation regarding

the generalization performance of LLMs after quantization, particularly considering the impact of *calibration data* introduced during the quantization process.

As shown in Fig. 1, the process of model quan-057 tization encompasses three distinct stages: pretraining, quantization, and inference, utilizing pre-training data, calibration data, and test data, 060 respectively. Existing quantization researches 061 typically use a standard calibration set, which is 062 usually a subset of the pre-training data (Sce-063 nario 1, S1), and evaluate on several fixed 064 datasets (Paperno et al., 2016; Clark et al., 2018; Tata & Patel, 2003; Sakaguchi et al., 2021; 065 Zellers et al., 2019). However, because using 066 task-specific data for model calibration is a more 067 reasonable choice in practical applications, the 068 relationship between the distribution of calibra-069 tion data and test data and its impact on the generalization ability of quantized models is a more 071 worthy research topic that has not been deeply 072 explored (Scenario 2, <u>S2</u>). In this work, to an-073 swer the abovementioned question and bridge 074 the gap between academic research and practi-075 cal implementation, we provide a platform to evaluate the generalization ability of quantized 076



Figure 1: We show the pipeline of model quantization and the data required at each stage (*Top*). The calibration data used in previous works generally share the same distribution with pre-training data (\underline{SI}), and the relation between calibration data and test data should be further discussed ($\underline{S2}$).

077 LLMs, covering *benchmarks*, *analyses*, and a modular-designed *toolbox*.

Benchmark evaluation. As shown in Fig. 1, we build the benchmark based on the two scenarios:

• In <u>S1</u> (Section 2), beyond the existing research, we collect the most comprehensive evaluation of test datasets to date, covering 9 categories and 26 datasets. We use C4 (Raffel et al., 2020) as the calibration dataset and quantize LLaMA2-7B (Touvron et al., 2023) model by two methods (Frantar et al., 2022; Dettmers et al., 2023) across three weight bit-widths.

• In S2 (Section 3), our benchmark covers 19 datasets with two types of distribution shifts between 084 calibration data and test data: cross-dataset and cross-subject. We consider both English and Chinese 085 domains for the cross-dataset setting. Besides, our benchmark also includes a more challenging cross-subject setting, e.g. from humanities to social science. To our knowledge, no prior work has 087 investigated the generalization of quantized models in a cross-subject setting. For all settings, our 088 benchmark builds the Independent and Identically Distribution (I.I.D) and Out-of-Distribution (OOD) 089 evaluations by adjusting the calibration data distributions. In our experiments, we quantize LLaMA family (Touvron et al., 2023) and Baichuan2-7B-Base (Du et al., 2021) for English and Chinese 091 models with four methods (Frantar et al., 2022; Dettmers et al., 2023; Lin et al., 2023; Xiao et al., 092 2023) across three weight bit-widths.

The generalization performance of quantized models is assessed using zero-shot and few-shot evaluation for all experiments, and we summarize the key features of our benchmark in Tab. 1.

Based on experimental results, we provide the following answers to the core question of the paper.

Answers. Quantized LLMs maintain strong generalization capabilities under high bit-width scenarios.
 However, their generalization performance deteriorates significantly in extremely low bit-width scenarios. Additionally, we find that using calibration data related to the test data does not significantly enhance generalization performance.

Empirical findings. Based on the experiments, we observe several counter-intuitive phenomena, *e.g.*,

Tasks vary significantly in their sensitivity to quantization; low-bit quantization can lead to improved performance for certain tasks. Scientific knowledge QA, reading comprehension, common sense reasoning, and mathematical reasoning are highly sensitive to quantization, particularly in low-bit scenarios, where quantization can lead to a significant decline in performance. Meanwhile, sentiment analysis is also very sensitive to quantization, but in low-bit situations, quantization can actually

| | | Ta | able 1: Su | mmary of t | he proposed benchmark. | | | | |
|---|------------------|---------|---------------|--------------|--|---------|--|--|--|
| ScenarioDistribution ShiftTask LanguageWeigh Precision | | | | Model | Benchmark & Dataset | | | | |
| S1 | - | English | {16, 4, 3, 2} | LLaMA2-7B | WinoGrande, WSC273, HellaSwag, SWAG, PIQA, MathQA, Mutual, Mutual_Plus, CrowS-Pairs, Toxigen, PubMedQA OpenBookQA, SciQ,ARC-Easy, ARC-Challenge, MC-TACO, RACE,QA4MRE, GLUE (6 datasets), ANLI, BLiMP | Fig. 2 | | | |
| S2 | Cross-dataset | English | {4, 3} | LLaMA family | BOSS (16 datasets) | Tab. 2 | | | |
| S2 | Cross-dataset | Chinese | {4, 3, 2} | Baichuan2-7B | C-EVAL, CMMLU | Tab. 13 | | | |
| S2 | Cross-discipline | Chinese | {4, 3, 2} | Baichuan2-7B | C-EVAL | Tab. 14 | | | |

Table 1. Summers of the proposed banchmark

lead to an improvement in performance. Tasks such as natural language inference demonstrate considerable robustness to quantization, with minimal changes in performance after quantization.

• Consistency between calibration data and test distribution does not always yield optimal performance, which is in stark contrast to the consensus during the pre-training and fine-tuning phases. The performance of quantized models using calibration data with distributions similar to or identical to those of downstream tasks is comparable to that achieved using subsets from high-quality corpora.

Toolbox. To support this work and facilitate future research, we develop a modular-designed code library. Specifically, this toolbox decouples the overall pipeline shown in Fig.1 into several separate components, e.g., LLM module, dataset module, quantizer module, etc., and provides common choices for each component and easy-to-use interface for possible extensions (see Section 4 and Fig. 4 for more details of the toolbox). This toolbox will be open-sourced along with the benchmark to facilitate future quantization applications and research.

S1: GENERALIZATION ASSESSMENT OF QUANTIZED LLMS WITH **STANDARD SETTING**



Figure 2: S1: Evaluation of quantized LLaMA2-7B on several standard datasets. Quantization methods include GPTQ and SpQR. Quantization bits include W4A16, W3A16, and W2A16, with W16A16 used as reference. The left figure shows 5-shot results, and the right figure shows 0-shot results. Different background colors represent different task types.

Experiment settings. To assess the difference in generalization ability, it is necessary to ensure that all other settings remain consistent except for the quantization process. To maintain consistency in the data encountered by the model before and after quantization, we strive to use calibration data during quantization that is as similar as possible to the data used during the pre-training phase of the LLM, namely the dataset C4 (Raffel et al., 2020) derived from pre-training data. The experimental setting is consistent with the evaluation settings used previously for quantized models (Frantar et al., 2022; Dettmers et al., 2023; Jaiswal et al., 2023; Liu et al., 2023a; Li et al., 2024). We utilize the LM Evaluation Harness (Gao et al., 2021a) with recommended parameters to conduct zero-shot and

few-shot tests on the following tasks. We provide full configurations in Appendix C, as well as code in the *supplementary materials*.

The 26 datasets we evaluated can be roughly divided into 9 categories: **O**common sense reason-165 ing, Omathematical reasoning, Omulti-turn dialogue reasoning, Obias diagnosis and mitigation, 166 **S**cientific knowledge question answering, **B**reading comprehension, **B**natural language inference, 167 **③**sentiment analysis, and **④**syntax phenomena evaluation. The common sense reasoning datasets 168 include WinoGrande (Sakaguchi et al., 2021), WSC273 (Levesque et al., 2012), GLUE-WNLI (Wang 169 et al., 2018), HellaSwag (Zellers et al., 2019), SWAG (Zellers et al., 2018), and PIQA (Tata & Patel, 170 2003). The mathematical reasoning datasets include MathQA(Amini et al., 2019). The multi-turn 171 dialogue reasoning datasets include Mutual and Mutual_Plus (Cui et al., 2020). The bias diagnosis 172 and mitigation datasets include CrowS-Pairs (Nangia et al., 2020) and Toxigen (Hartvigsen et al., 2022). The scientific knowledge question answering datasets include PubMedQA (Jin et al., 2019), 173 OpenBookQA (Mihaylov et al., 2018), SciQ (Welbl et al., 2017), ARC-Easy, ARC-Challenge (Clark 174 et al., 2018), and MC-TACO (Zhou et al., 2019). The reading comprehension datasets include 175 RACE (Lai et al., 2017) and QA4MRE (Peñas et al., 2013). The natural language inference datasets 176 include GLUE-MNLI, GLUE-MNLI-Mismatched, GLUE-RTE, GLUE-QNLI (Wang et al., 2018), 177 and ANLI (Nie et al., 2019). The sentiment analysis dataset includes GLUE-SST (Wang et al., 2018). 178 The syntax phenomena evaluation dataset includes BLiMP (Warstadt et al., 2020). 179

Results and analysis. We present the experimental results for both the 5-shot and 0-shot scenarios 180 in Fig. 2. To more clearly observe the experimental results of different downstream task types, we 181 average the decline in accuracy after quantization per task type. We present the results in Tab. 12 182 and conduct a more detailed analysis in Appendix D.2. It can be obviously observed that models 183 still retain strong generalization capabilities under low-bit quantization, but experience significant 184 performance degradation under ultra-low-bit quantization. When quantizing model weights to 185 3-4 bits, the performance degradation of all methods is not very pronounced. For some datasets, quantizing to 4 bits even leads to higher model performance compared to full precision. However, 187 when weights are quantized to 2 bits, GPTQ exhibits a significant performance drop on most tasks. 188 Compared to other methods, SPQR maintains relatively good performance at 2 bits, which may be 189 attributed to SPQR's ability to identify and isolate outlier weights. Additionally, different tasks exhibit varying sensitivities to quantization, with scientific knowledge QA, reading comprehension, common 190 sense reasoning, mathematical reasoning and sentiment analysis showing higher sensitivity while 191 tasks like natural language inference emerge lower sensitivity. For example, in extreme quantization 192 scenarios in Tab. 12, the performance of scientific knowledge QA declines by up to 40%, while the 193 performance drop for natural language inference is only 10%. Interestingly, we observe sentiment 194 analysis actually shows a significant performance improvement under extreme quantization scenarios. 195 We also demonstrate the robustness of the experiments with respect to random seeds in Appendix D.1. 196

197 198

199 200

3 S2: GENERALIZATION ASSESSMENT OF QUANTIZED LLMS WITH DOMAIN SHIFTS

This section investigates novel generalization scenarios in quantization, where different generalization 201 scenarios serve as instantiations of the framework. The distribution shift we consider primarily 202 pertains to the shift from calibration data to test data. Types of distribution shift include cross-dataset 203 distribution shift and *cross-subject* distribution shift, aimed at studying the impact of distribution 204 shift from calibration data to test data on quantized model performance. Cross-dataset distribution 205 shift refers to using different datasets as calibration set, while cross-subject distribution shift refers 206 to using different subjects from the same dataset as calibration set. Experiments will encompass 207 two main categories: English cross-dataset distribution shift experiments on the out-of-distribution 208 generalization benchmark BOSS in Sec. 3.1, and Chinese cross-dataset distribution shift experiments 209 as well as Chinese cross-subject distribution shift experiments on Chinese domain-specific tasks in 210 Sec. 3.2. We will conduct an in-depth analysis of the results from S2 in Sec. 3.3. We provide full 211 configurations in Appendix C, as well as code in the supplementary materials.

212 213

214

3.1 ENGLISH CROSS-DATASET TRANSFER TASK

Experiment settings. We evaluate *cross-dataset* distribution shift experiments on the OOD benchmark BOSS (Yuan et al., 2024) in NLP. Previous work in NLP concerning OOD mostly considers

Table 2: Cross-dataset distribution shift evaluation on BOSS. "Calib." represents the calibration 217 dataset, and "Gene." represents generalization scenario. To save space, abbreviations are used for 218 datasets. Each row presents experimental results using different datasets as calibration sets on the 219 same test dataset. Results with blue backgrounds indicate I.I.D results, while those without color 220 represent OOD results. The higher the metric, the better the performance. Bold results indicate 221 the best performance on the same test dataset. The result indicates I.I.D dataset achieve the best 222 performance. Note: "-" indicates out of memory results. 223

| 004 | Method | | | | EQA | | | | | | | SA | | | | | | | NLI | | | | | | | TD | | | |
|-------|--------|----------|---------|------|-------|--------------|-------|--------------|------|--------|------|-------|-------|-------|----------------|----------|---------|------|-------|------------|-----------|-----|----------|---------|------------|----------------|-------------|-----------|-------|
| 224 | | Tect | Cone | W/A | | Ca | lib. | | Test | Cono | W/A | | Ca | dib. | | Test | Cono | W/A | | Cal | ib. | | Test | Cono | W/A | | Ca | lib. | |
| 225 | | Ita | ounc. | | SQ | AQA | NQA | SQA | nest | ocne. | | AZ | DS | SE | SST | nor | oene. | | MN | AN | WN | CN | Itsi | Ocne. | A | CC | AC | IH | TG |
| | | | 0-shot | 4/16 | 53.84 | 52.73 | 54.69 | 57.31 | | 0-shot | 4/16 | 70.81 | 17.87 | 63.18 | 72.08 | | 0-shot | 4/16 | 0.36 | 0.23 | 0.22 | - | | 0-shot | 4/16 | 23.90 | 26.96 | 52.52 | 53.32 |
| 226 | | SQ | | 3/16 | 45.31 | 48.86 | 49.49 | 50.79 | AZ | | 3/16 | 38.06 | 0.38 | 0.26 | 0.04 | MN | | 3/16 | 0.00 | 0.00 | 0.00 | - | cc | | 3/16 | 0.60 | 2.45 | 9.70 | 10.60 |
| ~~~ | | | 1-shot | 4/16 | 67.04 | 65.97 | 67.06 | 68.16 | | 3-shot | 4/16 | 83.69 | 56.66 | 80.79 | 82.55 | | 3-shot | 4/16 | 49.69 | 32.81 | 34.93 | - | | 2-shot | 4/16 | 91.80 | 87.46 | 91.71 | 91.84 |
| 227 | | | | 3/10 | 28.00 | 27.12 | 28.40 | 30.40 | | | 3/10 | /4.54 | 24.80 | 31.82 | 39.79 46 79 | | | 3/10 | 1.07 | 0.52 | 0.03 | - | | | 3/10 | 89.11 10 12 | 5.03 | 7.84 | 90.33 |
| 222 | | | 0-shot | 3/16 | 21.81 | 25 28 | 23 35 | 24 99 | | 0-shot | 3/16 | 17 59 | 1 72 | 0.01 | 0.00 | | 0-shot | 3/16 | 4 17 | 0.00 | 0.00 | | | 0-shot | 3/16 | 0.76 | 1 72 | 0.19 | 0.57 |
| 220 | anno | AQA | | 4/16 | 35.50 | 36.11 | 31.97 | 35.77 | DS | | 4/16 | 54.40 | 38.78 | 52.54 | 55,50 | AN | | 4/16 | 34.34 | 33.76 | 33.24 | - | AC | | 4/16 | 15.87 | 17.59 | 15.87 | 16.25 |
| 229 | GPIQ | | 1-shot | 3/16 | 31.39 | 29.54 | 31.60 | 32.24 | | 3-shot | 3/16 | 54.68 | 36.05 | 33.86 | 43.46 | | 3-shot | 3/16 | 30.97 | 33.69 | 33.28 | - | | 2-shot | 3/16 | 60.23 | 90.35 | 15.87 | 56.02 |
| | | | 0-shot | 4/16 | 37.94 | 38.76 | 38.63 | 38.23 | | 0-shot | 4/16 | 18.32 | 8.21 | 15.60 | 26.43 | | 0_shot | 4/16 | 0.09 | 0.04 | 0.11 | - | | 0.shot | 4/16 | 37.37 | 22.55 | 33.90 | 40.82 |
| 230 | | NOA | 0-31101 | 3/16 | 31.36 | 33.79 | 33.37 | 34.45 | SE | 0-shot | 3/16 | 4.83 | 0.09 | 0.20 | 0.01 | WN | 0-31101 | 3/16 | 0.49 | 0.00 | 0.00 | - | ш | 0-shot | 3/16 | 11.27 | 7.32 | 4.53 | 13.18 |
| 004 | | | 1-shot | 4/16 | 48.55 | 49.30 | 49.73 | 49.09 | | 3-shot | 4/16 | 42.96 | 28.55 | 42.99 | 44.75 | | 3-shot | 4/16 | 41.51 | 43.34 | 47.53 | - | | 2-shot | 4/16 | 62.36 | 63.46 | 62.00 | 62.29 |
| 231 | | | | 3/16 | 44.38 | 43.35 | 46.95 | 45.61 | | | 3/16 | 42.36 | 22.67 | 35.54 | 29.40 | | | 3/16 | 38.83 | 48.09 | 48.15 | - | | | 3/16 | 63.52 | 90.35 | 61.83 | 61.77 |
| 232 | | | 0-shot | 4/16 | 42.58 | 45.72 | 46.21 | 44.20 | | 0-shot | 2/16 | 49.15 | 20.73 | 27.12 | 44.98 | | 0-shot | 4/16 | 0.06 | 0.00 | 0.00 | - | | 0-shot | 4/16 | 48.44 | 36.72 | 44.84 | 57.97 |
| 101 | | SQA | | 1/16 | 56.04 | 61.80 | 60.92 | 62 17 | SST | | 1/16 | 60.50 | 33.25 | 15 24 | 51 11 | CN | | 1/16 | 35 23 | 36 35 | 32 44 | | TG | | J/16 | 72.03 | 75 47 | 67.81 | 68.40 |
| 233 | | | 1-shot | 3/16 | 43.46 | 42.83 | 45.17 | 48.82 | | 3-shot | 3/16 | 54.37 | 33.25 | 35.46 | 50.20 | | 3-shot | 3/16 | 29.54 | 29.03 | 33.39 | - | | 2-shot | 3/16 | 70.47 | 90.35 | 57.50 | 62.19 |
| 004 | | _ | - | | | Ca | lib. | | _ | - | | | Ca | dib. | | - | - | | -, | Cal | ib. | | _ | - | | | Ca | lib. | |
| 234 | | Test | Gene. | W/A | SQ | AQA | NQA | SQA | Test | Gene. | W/A | AZ | DS | SE | SST | Test | Gene. | W/A | MN | AN | WN | CN | Test | Gene. | W/A | CC | AC | IH | TG |
| 225 | | | 0 -1 | 4/16 | 57.03 | 49.87 | 53.00 | 54.36 | | 0 -1 | 4/16 | 63.34 | 62.46 | 72.52 | 83.14 | | 0 | 4/16 | 0.57 | 0.02 | 0.13 | - | | 0 -1 | 4/16 | 61.73 | 59.48 | 58.92 | 37.48 |
| 200 | | so | 0-snot | 3/16 | 52.37 | 45.90 | 54.55 | 58.36 | 47 | 0-snot | 3/16 | 72.38 | 55.79 | 37.28 | 27.84 | MN | 0-snot | 3/16 | 0.00 | 0.01 | 0.00 | - | 0.0 | 0-snot | 3/16 | 36.90 | 2.54 | 15.42 | 22.38 |
| 236 | | 2 | 1-shot | 4/16 | 66.45 | 66.80 | 67.41 | 67.21 | | 3-shot | 4/16 | 79.65 | 69.31 | 85.44 | 82.91 | | 3-shot | 4/16 | 36.19 | 40.45 | 41.62 | - | 00 | 2-shot | 4/16 | 90.65 | 89.27 | 91.74 | 84.69 |
| | | | | 3/16 | 65.12 | 65.55 | 68.65 | 66.95 | | | 3/16 | 83.68 | 86.30 | 72.18 | 83.50 | | | 3/16 | 32.39 | 40.31 | 38.47 | - | | | 3/16 | 87.70 | 91.76 | 86.99 | 83.56 |
| 237 | | | 0-shot | 4/16 | 30.59 | 25.11 | 27.60 | 29.50 | | 0-shot | 4/16 | 35.47 | 43.53 | 40.85 | 50.40 | | 0-shot | 4/16 | 0.86 | 0.07 | 0.28 | - | | 0-shot | 4/16 | 10.13 | 4.97 | 12.05 | 13.58 |
| 000 | | AQA | | 3/10 | 20.35 | 21.45 | 27.55 | 30.30 | DS | | 3/10 | 41.8/ | 31.17 | 15.42 | 29.10 | AN | | 3/10 | 22.17 | 22 21 | 22 70 | - | AC | | 3/10 | 2.49 | 0.70 | 15.97 | 2.87 |
| 230 | SpQR | | 1-shot | 3/16 | 34.61 | 34 75 | 37 49 | 33.10 | | 3-shot | 3/16 | 59 10 | 54.80 | 52 56 | 56.02 | | 3-shot | 3/16 | 33.66 | 31.93 | 33.14 | - | | 2-shot | 3/16 | 15.87 | 15.87 | 19 31 | 15.87 |
| 239 | | | | 4/16 | 40.30 | 38.01 | 39.40 | 38.22 | | | 4/16 | 14.62 | 23.36 | 19.85 | 33.24 | | | 4/16 | 0.28 | 0.00 | 0.00 | - | | | 4/16 | 42.21 | 41.79 | 40.12 | 31.76 |
| 200 | | | 0-shot | 3/16 | 35.79 | 33.27 | 40.80 | 38.77 | | 0-shot | 3/16 | 16.05 | 10.22 | 4.75 | 7.30 | | 0-shot | 3/16 | 0.00 | 0.06 | 0.00 | - | | 0-shot | 3/16 | 31.32 | 6.78 | 17.68 | 16.96 |
| 240 | | NQA | 1 shot | 4/16 | 49.61 | 49.12 | 49.70 | 48.47 | SE | 2 shot | 4/16 | 44.48 | 44.15 | 44.25 | 44.39 | WN | 2 shot | 4/16 | 43.28 | 43.77 | 41.79 | - | ш | 2 shot | 4/16 | 64.24 | 65.85 | 62.14 | 66.07 |
| 0.4.4 | | | 1-SHOL | 3/16 | 48.25 | 46.61 | 48.99 | 47.79 | | J-SHOL | 3/16 | 53.16 | 43.63 | 41.76 | 44.77 | | 5-shot | 3/16 | 39.09 | 47.32 | 40.77 | - | | 2-51101 | 3/16 | 62.95 | 63.14 | 63.17 | 64.37 |
| 241 | | | 0-shot | 4/16 | 46.45 | 42.62 | 44.30 | 45.10 | | 0-shot | 4/16 | 46.02 | 29.47 | 44.72 | 55.67 | | 0-shot | 4/16 | 0.00 | 0.22 | 0.45 | - | | 0-shot | 4/16 | 54.37 | 52.66 | 51.09 | 39.53 |
| 2/12 | | SOA | - | 3/16 | 36.90 | 44.57 | 42.88 | 39.31 | SST | | 3/16 | 23.08 | 14.87 | 3.65 | 6.52 | CN | | 3/16 | 0.06 | 0.00 | 0.89 | - | TG | | 3/16 | 41.88 | 9.69 | 19.38 | 37.34 |
| 272 | | - | 1-shot | 4/16 | 61.63 | 57.77 | 61.79 | 60.55 | | 3-shot | 4/16 | 55.41 | 42.37 | 58.54 | 59.32 | | 3-shot | 4/16 | 36.13 | 34.84 | 34.23 | - | | 2-shot | 4/16 | 69.84 | 76.56 | 61.41 | 77.60 |
| 243 | | | | 3/16 | 48.80 | 59.19 | 36.34 | 55.06 | | | 3/16 | 63.49 | 60.37 | 53.98 | 61.80 | | | 3/10 | 35.29 | 35.90 | 33.17 | - | | | 3/10 | /3.13 | 00.88 Cr | 08.44 | 77.03 |
| | | Test | Gene. | W/A | 50 | | NOA | SOA | Test | Gene. | W/A | AZ | | SE | SST | Test | Gene. | W/A | MN | | ID. WN | CN | Test | Gene. | W/A | CC | | шо. ПН | TG |
| 244 | | | | 4/16 | 56,73 | 55.09 | 52.09 | 50.21 | | | 4/16 | | 5.42 | 35.23 | 33.65 | | | 4/16 | 0.48 | 0.14 | 0.06 | - | | | 4/16 | 50.17 | 66.60 | 42.19 | 42.11 |
| 245 | | 60 | 0-shot | 3/16 | 48.32 | 37.95 | 44.45 | 40.30 | | 0-shot | 3/16 | - | 39.41 | 70.10 | 35.95 | 101 | 0-shot | 3/16 | 0.00 | 0.01 | 0.01 | - | 00 | 0-shot | 3/16 | 41.96 | 39.03 | 46.95 | 14.72 |
| 243 | | SQ | 1 -1 | 4/16 | 66.57 | 66.91 | 67.02 | 66.21 | AL | 2 -1 | 4/16 | - | 83.64 | 83.73 | 78.06 | MIN | 2 | 4/16 | 42.20 | 38.37 | 36.05 | - | LCC. | 2 - 1 4 | 4/16 | 91.84 | 91.63 | 90.80 | 89.31 |
| 246 | | | 1-shot | 3/16 | 59.81 | 61.81 | 61.27 | 61.38 | | J-SHOL | 3/16 | - | 88.73 | 90.16 | 88.92 | | 5-shot | 3/16 | 35.44 | 34.22 | 35.34 | - | | 2-51101 | 3/16 | 36.43 | 73.04 | 90.93 | 27.24 |
| | | | 0-shot | 4/16 | 29.73 | 29.20 | 28.34 | 27.57 | | 0-shot | 4/16 | - | 2.36 | 20.10 | 22.19 | | 0-shot | 4/16 | 0.59 | 0.07 | 0.07 | - | | 0-shot | 4/16 | 9.56 | 11.85 | 11.28 | 5.55 |
| 247 | | AQA | | 3/16 | 23.02 | 17.58 | 20.37 | 18.62 | DS | | 3/16 | - | 8.76 | 27.09 | 11.87 | AN | | 3/16 | 0.00 | 0.00 | 0.00 | - | AC | | 3/16 | 5.74 | 4.59 | 4.21 | 1.15 |
| 0.40 | AWQ | | 1-shot | 4/16 | 35.76 | 37.01 | 37.55 | 36.78 | | 3-shot | 4/16 | - | 53.91 | 55.92 | 50.95 | | 3-shot | 4/16 | 33.66 | 33.66 | 33.66 | - | | 2-shot | 4/16 | 15.87 | 15.87 | 16.06 | 16.63 |
| 240 | | | | 3/10 | 30.20 | 38.58 | 30.47 | 38 10 | | | 3/10 | - | J0.95 | 18 90 | 14.96 | | | 3/10 | 0.30 | 0.17 | 0.02 | - | <u> </u> | | 3/10 | 24.80 | 44 64 | 34.00 | 27.16 |
| 249 | | | 0-shot | 3/16 | 35.75 | 31.27 | 32.91 | 33.69 | | 0-shot | 3/16 | | 5.52 | 14.95 | 5.49 | | 0-shot | 3/16 | 0.00 | 0.00 | 0.00 | - | | 0-shot | 3/16 | 20.22 | 17.97 | 25.72 | 4.73 |
| | | NQA | | 4/16 | 43.25 | 43.18 | 43.39 | 42.56 | SE | | 4/16 | - | 45.03 | 45.44 | 43.77 | WN | | 4/16 | 40.02 | 39.40 | 38.23 | - | ш | | 4/16 | 62.36 | 62.46 | 65.03 | 64.67 |
| 250 | | | 1-shot | 3/16 | 41.02 | 40.50 | 41.27 | 41.26 | | 3-shot | 3/16 | - | 38.53 | 55.02 | 44.50 | | 3-shot | 3/16 | 37.11 | 44.38 | 37.17 | - | | 2-shot | 3/16 | 61.85 | 63.03 | 61.88 | 61.79 |
| 054 | | | 0 shot | 4/16 | 43.83 | 43.07 | 44.32 | 44.20 | | 0 shot | 4/16 | - | 2.09 | 11.47 | 19.17 | | 0 shot | 4/16 | 3.35 | 1.56 | 3.41 | - | | 0 shot | 4/16 | 49.38 | 52.5 | 40.31 | 36.56 |
| 251 | | SOA | 0-SHOE | 3/16 | 35.10 | 29.62 | 29.55 | 32.07 | SST | J-SHOT | 3/16 | - | 3.39 | 30.77 | 8.21 | CN | J-SHOT | 3/16 | 3.07 | 0.06 | 1.79 | - | TG | J-SHOL | 3/16 | 26.72 | 20.00 | 37.03 | 8.44 |
| 252 | | 2.1 | 1-shot | 4/16 | 48.12 | 48.39 | 49.37 | 47.24 | 001 | 3-shot | 4/16 | - | 58.28 | 58.80 | 51.76 | 0.1 | 3-shot | 4/16 | 33.84 | 34.51 | 33.28 | - | 1.0 | 2-shot | 4/16 | 65.31 | 65.47 | 71.25 | 75.16 |
| 232 | | | | 3/16 | 40.48 | 39.14 | 39.84 | 43.61 | | | 3/16 | - | 57.11 | 64.93 | 65.84 | | | 3/16 | 28.14 | 33.61 | 29.87 | - | | | 3/16 | 68.75 | 74.22 | 63.91 | 67.66 |
| 253 | | Test | Gene. | W/A | 60 | Ca | lib. | 604 | Test | Gene. | W/A | 17 | | dib. | COT | Test | Gene. | W/A | MN | Cal | ID. | CN | Test | Gene. | W/A | CC | Ca | lib. | TC |
| | | | | 1/8 | 35.34 | AQA 30.17 | 40.12 | 3QA 40.64 | | | 1/8 | AZ | 0.87 | 5E | 0.01 | | | //8 | 0.01 | AN 0.00 | 0.00 | CN. | | | 1/8 | 0.03 | AC | 0.00 | 0.00 |
| 254 | | | 0-shot | 3/8 | 0.01 | 0.01 | 0.00 | 0.01 | | 0-shot | 3/8 | 0.00 | 0.00 | 0.00 | 0.01 | | 0-shot | 3/8 | 0.00 | 0.00 | 0.00 | - | | 0-shot | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 |
| 055 | | SQ | | 4/8 | 28.01 | 56.13 | 55.12 | 56.59 | AZ | | 4/8 | 54.76 | 88.36 | 85.90 | 84.85 | MN | | 4/8 | 31.70 | 32.86 | 34.54 | - | cc | | 4/8 | 90.39 | 65.29 | 65.08 | 87.63 |
| 200 | | | 1-shot | 3/8 | 0.00 | 0.00 | 0.00 | 0.01 | | 3-shot | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 | | 3-shot | 3/8 | 0.00 | 0.00 | 0.00 | - | | 2-shot | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 |
| 256 | | | 0 shot | 4/8 | 14.22 | 18.10 | 18.18 | 18.17 | | 0 shot | 4/8 | 0.00 | 0.02 | 0.00 | 0.00 | | 0 shot | 4/8 | 0.03 | 0.00 | 0.00 | - | | 0 shot | 4/8 | 0.19 | 0.19 | 0.19 | 0.00 |
| 200 | | 101 | 0-snot | 3/8 | 0.01 | 0.00 | 0.00 | 0.01 | DS | 0-snot | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 | AN | 0-snot | 3/8 | 0.00 | 0.00 | 0.00 | - | AC | 0-snot | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 |
| 257 | so | -ngA | 1-shot | 4/8 | 28.01 | 28.91 | 27.96 | 29.13 | 100 | 3-shot | 4/8 | 47.01 | 50.42 | 50.00 | 34.88 | and | 3-shot | 4/8 | 32.97 | 33.93 | 33.14 | - | AC | 2-shot | 4/8 | 18.16 | 23.71 | 53.15 | 23.33 |
| 050 | | <u> </u> | | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 | | | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 | | | 3/8 | 0.00 | 0.00 | 0.00 | - | | | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 |
| 258 | | | 0-shot | 4/8 | 24.26 | 27.83 | 23.95 | 24.07 | | 0-shot | 4/8 | 0.00 | 0.00 | 0.00 | 0.01 | | 0-shot | 4/8 | 0.00 | 0.00 | 0.00 | - | | 0-shot | 4/8 | 0.09 | 0.01 | 0.00 | 0.04 |
| 250 | | NQA | | 3/8 | 0.00 | 0.00 | 0.00 | 0.01 | SE | | 3/8 | 0.00 | 24.42 | 0.00 | 0.00 | WN | | 3/8 | 0.00 | 0.00 | 0.00 | - | ІН | | 3/8 | 0.00 | 0.00 | 0.02 | 0.00 |
| 200 | | | 1-shot | 4/8 | 0.09 | 0.00 | 29.85 | 0.00 | | 3-shot | 3/8 | 47.03 | 0.00 | 45.45 | 34.27 | | 3-shot | 4/8 | +/.32 | 40.70 | 47.15 | 1 | | 2-shot | 4/8 3/8 | 0.00 | 0.00 | 39.28 | 57.42 |
| 260 | | <u> </u> | | 4/8 | 19.92 | 20.30 | 19.07 | 18.07 | | | 4/8 | 0.00 | 0.00 | 0.00 | 0.00 | \vdash | | 4/8 | 0.00 | 0.00 | 0.00 | | | | 4/8 | 0.63 | 0.00 | 0,16 | 0.00 |
| | | | 0-shot | 3/8 | 0.00 | 0.00 | 0.00 | 0.01 | 0.00 | 0-shot | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 | | 0-shot | 3/8 | 0.00 | 0.00 | 0.00 | - | | 0-shot | 3/8 | 0.00 | 0.00 | 0.31 | 0.00 |
| 261 | | SQA | 1 | 4/8 | 25.64 | 17.73 | 21.70 | 21.10 | SST | 2 | 4/8 | 26.47 | 53.06 | 55.02 | 36.90 | CN | 2 | 4/8 | 26.02 | 27.19 | 14.91 | - | TG | a | 4/8 | 58.28 | 65.31 | 59.06 | 59.38 |
| | | | 1-snot | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 | | 3-snot | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 | | 5-snot | 3/8 | 0.00 | 0.00 | 0.00 | | | ∠-snot | 3/8 | 0.00 | 0.00 | 0.00 | 0.00 |

263 264

distribution shifts from various sources, e.g. from movies to Twitter (Yu et al., 2024). GLUE-X (Yang et al., 2022) and BOSS (Yuan et al., 2024) represent pioneering efforts in benchmarking OOD generalization in NLP. BOSS, building upon GLUE-X, improves by employing SimCSE scores for 265 detection analysis and identifying dataset pairs exhibiting the lowest semantic similarity. These pairs 266 are then utilized for training and testing, constructing a benchmark consisting of five downstream 267 tasks. Each downstream task comprises and I.I.D dataset and three OOD datasets. 268

To evaluate the generalization ability of quantized models in cross-dataset distribution shift experi-269 ments, we randomly sample 300 samples from the test set of each OOD dataset within the BOSS

270 benchmark as its corresponding training set, serving as the calibration set for the quantization process. 271 For each downstream task, we utilize the training set from different datasets as the calibration set for 272 the quantization process and test on the corresponding I.I.D and OOD test sets. In our experiments, 273 we employ LLaMA2-7B, LLaMA2-13B and LLaMA3-8B (Touvron et al., 2023) as the target for 274 quantization and selected four PTQ methods: GPTQ (Frantar et al., 2022), AWQ (Lin et al., 2023), SpQR (Dettmers et al., 2023), and SmoothQuant (Xiao et al., 2023). Given that there is not much 275 difference in performance between excessively high bits and full precision, and too low a bit has 276 already lost basic performance in these tasks, we only present the results of LLaMA2-7B weights quantizing to 3-4 bits with SmoothQuant quantizing the activations to 8 bits. Results for more models 278 are in Appendix D.4, full precision results are in Appendix D.5 and 2-bit results are in Appendix D.6. 279 We test two forms: 0-shot and few-shot. 280

Results and Analysis. We present the results in Tab. 2. We evaluate four downstream tasks in BOSS: 281 EQA, SA, NLI, and TD. Each downstream task consists of four datasets, with each dataset tested 282 using four datasets as calibration set. The following conclusions can be observed: For the same 283 test dataset, it's not necessarily the case that using I.I.D dataset as calibration set yield superior 284 performance. We can clearly observe that the overlap frequency between the I.I.D results with 285 background and the bolded best performance results in Tab. 2 is not high, and may even be quite 286 low. This suggests that the I.I.D dataset, when used as a calibration set, does not reliably enhance the 287 capabilities of the quantized model. This is a counter-intuitive conclusion, as it is a well-established 288 fact in fields such as pre-training and instruction fine-tuning that I.I.D datasets used as training sets 289 can improve performance on corresponding downstream tasks. For instance, when testing SQ dataset 290 in EQA task with GPTQ, the performance when using the SQ dataset as the calibration set is not 291 ideal; on the contrary, using the SQA dataset as the calibration set yields better performance. We also demonstrate the robustness of the experiments with respect to random seeds in Appendix D.1. 292

293 294

295

3.2 CHINESE CROSS-DATASET AND CROSS-DISCIPLINE TASKS

296 **Experiment settings**. We evaluate *cross-dataset* distribution shift experiments and *cross-subject* 297 distribution shift experiments on the Chinese domain-specific datasets C-EVAL (Huang et al., 2024b) 298 and CMMLU (Li et al., 2023). C-EVAL serves as a comprehensive benchmark for evaluating Chinese 299 LLM. It consists of 13,948 multiple-choice questions covering 52 different subjects categorized into Humanities, Social Sciences, STEM, and Other. CMMLU is another Chinese evaluation dataset 300 designed specifically to assess the advanced knowledge and reasoning abilities of LLM in the context 301 of the Chinese language and culture. It encompasses 67 different subjects categorized into Humanities, 302 Social Sciences, STEM, and Chinese specific and others. 303

304 Both C-EVAL and CMMLU, two Chinese-specific domain datasets, include Humanities, Social Sciences, and STEM three subject categories. We design cross-dataset distribution shift experiments 305 based on the same subject categories. For each subject test, we respectively utilize the corresponding 306 subjects from C-EVAL and CMMLU as calibration set to assess the impact of different datasets 307 as calibration set on the test results. Additionally, we conducted cross-subject distribution shift 308 experiments on the C-EVAL dataset. For each subject test, we use Humanities, Social Sciences, and 309 STEM as calibration set to evaluate the influence of different subject subsets as calibration set on 310 the test results. Since both C-EVAL and CMMLU lack training datasets, we used the validation 311 dataset of C-EVAL as the training dataset and randomly sampled 300 samples from the test dataset of 312 CMMLU as the training dataset. We utilize the Chinese LLM Baichuan2-7B-Base (Yang et al., 2023) 313 as the quantization target and selecte four PTQ methods: GPTQ (Frantar et al., 2022), AWQ (Lin 314 et al., 2023), SpQR (Dettmers et al., 2023), and SmoothQuant (Xiao et al., 2023). We quantize the 315 weights to 2-4 bits, with SmoothQuant quantizing the activations to 8 bits, and test both 0-shot and 5-shot forms. 316

Results and Analysis. The results of cross-dataset distribution shift experiments on C-EVAL and
 CMMLU are presented in Tab. 13. We observed that *the performance of the I.I.D dataset is slightly higher than that of the OOD dataset, but it still tends to be random*. In most cases, the performance of
 the I.I.D dataset as a calibration set is better at times, while at other times the OOD dataset performs
 better, or both perform equally well. For example, when testing C-EVAL with GPTQ, the I.I.D
 dataset C-EVAL significantly outperforms CMMLU as the calibration set; however, when testing
 CMMLU with GPTQ, the performance of the I.I.D and OOD datasets tends to be random. There are

C-EVAL dataset with GPTQ, the I.I.D dataset C-EVAL serves as a better calibration set, whereas
 when testing the C-EVAL dataset with AWQ, the OOD dataset CMMLU appears to yield better
 performance for the quantized model.

The results of cross-subject distribution shift experiments on C-EVAL are presented in Tab. 14. The 328 results tend to be more random, and no conclusion can be drawn that using an I.I.D dataset as 329 calibration set results in higher test accuracy. An intuitive idea is that the subjects of humanities 330 and social sciences are closely related and differ significantly from STEM, suggesting that using 331 similar datasets would outperform those that are more disparate. However, our experimental results 332 indicate that when testing on humanities and social sciences, there are instances where using STEM 333 as a calibration set can also lead to higher performance. Similarly, when testing on STEM, there are 334 occasions when using humanities and social sciences as calibration sets yields better performance than STEM itself. Thus, the conclusion that cross-disciplinary approaches result in a significant 335 performance drop does not hold true. 336

337 338

339

345 346 347

348

3.3 S2: COMPARATIVE ANALYSIS OF ALL RESULTS

Overall Findings. Consistency between calibration data and test distribution does not always yield
 optimal performance. We integrate Tab. 2,Tab. 13and Tab. 14 and calculate the wining rates for two
 strategies: using I.I.D dataset and OOD dataset as the test set. To compare the model's performance
 when using I.I.D and OOD datasets as calibration sets, the calculation method for the IID wining rate
 is cauculated as:

$$\frac{Num(win)_{I.I.D}}{Num(all)} \tag{1}$$

 $\begin{array}{ll} \begin{array}{l} 349\\ 350\\ 350\\ 351\\ 351\\ 351\\ 352\\ 353\\ 353\\ 353\\ \end{array} \begin{array}{l} Num(win)_{I.I.D} \text{ refers to the number of samples where the performance of the I.I.D calibration set exceeds that of the OOD calibration set, while <math>Num(all)$ denotes the total number of data samples. The performance of the OOD calibration set is calculated as the average performance across all OOD calibration sets to eliminate discrepancies in the number of OOD calibration sets. The results are displayed in Tab. 3. \end{array}

354 We are surprised to find that the calibration set using I.I.D datasets do not achieve better results in 355 more than half of the settings that the winning rates are all near to 0.5, which is a rather counter-356 intuitive finding. Based on previous research on OOD generalization, using I.I.D data for pre-training 357 or fine-tuning typically yields performance that far exceeds that of OOD data (Yuan et al., 2024). 358 Additionally, experimental results (Wang et al., 2024; Albalak et al., 2024) indicate that fine-tuning 359 large models using domain-relevant data yields better performance for specific downstream tasks. 360 However, during the quantization phase, using datasets with the same or similar distribution as the 361 test dataset did not significantly improve the performance of the quantized model.

362 **Guess.** LLMs may not require highly relevant data related to downstream tasks to recover the 363 performance loss due to quantization. For the quantization task, the current design of quantization 364 algorithms aims to restore the performance of full-precision models, rather than adapt to downstream 365 tasks. Therefore, the model does not need data that closely resembles the downstream task, but rather 366 a small amount of data to restore its capabilities, and it is not very sensitive to this subset of data. Thus, 367 when the distribution difference between the calibration and test data is not significant, I.I.D data cannot yield better results. This also explains why the continuous increase of the calibration dataset 368 leads to diminishing returns, rather than following the scaling laws like pre-training data (Williams & 369 Aletras, 2023). 370

Table 3: The winning rate of I.I.D calibration data against OOD calibration data in three groups of experiments in S2.

| 374Cross-dataset English Task0.457375Cross-dataset Chinese Task0.61T | esults | ing Rate | Distribution Shift | 373 |
|--|--------|----------|----------------------------|-----|
| 375 Cross-dataset Chinese Task 0.61 T | ab. 2 | 0.45 | Cross-dataset English Task | 374 |
| | b. 13 | 0.61 | Cross-dataset Chinese Task | 375 |
| 376 Cross-subject Chinese Task 0.43 Ta | b. 14 | 0.43 | Cross-subject Chinese Task | 376 |

378 Clue. We conduct a comparative ex-379 periment using C4 (Raffel et al., 2020) 380 as the calibration set on BOSS and 381 present the result in Tab 4, which is a 382 standard setting usually used for quantization. We observed that the perfor-383 mance of C4 is similarly close and 384 random to that of I.I.D/OOD datasets, 385 without a unified conclusion indicat-386 ing which one consistently performs 387 better. This result further validates our 388 speculation that the choice of data dur-389 ing the quantization phase is more ro-390 bust compared to other stages of data 391 selection. 392



Figure 3: Normalized accuracy on EQA task using GPTQ and AWQ method. The left figure displays GPTQ results, while the right figure displays AWQ results.

Table 4: Results of C4 compared to I.I.D and OOD dataset as calibration set. We use GPTQ and test on the EQA task on BOSS. C4' is selected using a different random seed. The two best performances are denoted in descending order with red and orange respectively.

| | Tost | Cono | W/A | | Calib. | | | | Tost | Gana | W/A | Calib. | | | | | | |
|---|------|---------|------|-------|--------|-------|-------|-------|-------|------|---------|--------|-------|-------|-------|-------|-------|-------|
| | itsi | Gene. | | C4 | C4' | SQ | AQA | NQA | SQA | lest | Gene. | , WA | C4 | C4' | SQ | AQA | NQA | SQA |
| | | 0 shot | 4/16 | 54.50 | 51.71 | 53.84 | 52.73 | 54.69 | 57.31 | | 0 shot | 4/16 | 39.79 | 38.98 | 37.94 | 38.76 | 38.63 | 38.23 |
| | 50 | 0-51101 | 3/16 | 54.29 | 54.73 | 45.31 | 48.86 | 49.49 | 50.79 | NOA | 0-51101 | 3/16 | 35.79 | 36.68 | 31.36 | 33.79 | 33.37 | 34.45 |
| | SQ. | 1_shot | 4/16 | 67.73 | 67.42 | 67.04 | 65.97 | 67.06 | 68.16 | nQA | 1 shot | 4/16 | 48.80 | 49.00 | 48.55 | 49.30 | 49.73 | 49.09 |
| | | 1-51101 | 3/16 | 63.72 | 64.64 | 60.76 | 58.84 | 63.34 | 63.01 | | 1-51101 | 3/16 | 45.66 | 45.03 | 44.38 | 43.35 | 46.95 | 45.61 |
| - | | 0 shot | 4/16 | 30.80 | 27.61 | 28.00 | 27.12 | 28.40 | 30.40 | | 0 shot | 4/16 | 46.29 | 46.03 | 42.58 | 45.72 | 46.21 | 44.20 |
| | AQA | 0-51101 | 3/16 | 25.04 | 27.91 | 21.81 | 25.28 | 23.35 | 24.99 | 501 | 0-51101 | 3/16 | 31.13 | 33.38 | 30.19 | 26.99 | 28.49 | 33.73 |
| | | 1 shot | 4/16 | 36.43 | 36.23 | 35.50 | 36.11 | 31.97 | 35.77 | SQA | 1 shot | 4/16 | 62.48 | 61.44 | 56.04 | 61.89 | 60.92 | 62.17 |
| | | 1- shot | 3/16 | 33.45 | 34.01 | 31.39 | 29.54 | 31.60 | 32.24 | | 1-snot | 3/16 | 53.25 | 53.00 | 43.46 | 42.83 | 45.17 | 48.82 |

Interesting Findings. For a specific algorithm, there may exist one or more datasets that enhance the performance of the quantized model, independent of whether the dataset is I.I.D or OOD. We present the results of normalizing the performance of the GPTQ and AWQ methods to the range [0, 1] on the EQA task in Fig. 3. The detailed normalization process is provided in Appendix C.5. It is evident that for GPTQ, the SQA dataset consistently exhibits good performance, while the SQ dataset consistently shows poor performance. However, this conclusion does not hold for AWQ, as all datasets demonstrate more random performance under AWQ. This may be related to the differing utilization of calibration data by algorithms.

412 413 414

415 416

406

407

408

409

410

411

394

395

4 MI-OPTIMIZE: LLM QUANTIZATION TOOLBOX

Overview. MI-optimize is a versatile tool designed for the quantization and evaluation of LLMs.

417 The library's seamless integration of var-418 ious quantization methods and evaluation techniques empowers users to customize 419 their approaches according to specific re-420 quirements and constraints, providing a 421 high level of flexibility. Although LLMs 422 excel in various NLP tasks, their compu-423 tational and memory demands may limit 424 their deployment in real-time applications 425 and on resource-constrained devices. MI-426 optimize addresses this challenge by em-427 ploying quantization techniques to com-428 press these models, ensuring they maintain 429 performance while remaining adaptable to a wide range of scenarios. Fig. 4 illustrates 430 the framework of MI-optimize, which com-431 prises five main modules: the Configura-



Figure 4: Overview of the Quantization and Evaluation Framework.

Table 5: Perplexity (PPL) of the LLaMA-2-7B model using SmoothQuant and a combination of
SmoothQuant for activations and GPTQ for weight quantization on the WikiText-2 (Wiki2), Penn
Treebank (PTB), and C4 datasets.

| Method | W/A | Wiki2 | C4 | РТВ |
|------------------|-------|----------|----------|----------|
| Baseline | 16/16 | 5.47 | 37.92 | 7.22 |
| Smoothquant | 8/8 | 19.70 | 3026.75 | 11.27 |
| Smoothquant+GPTQ | 8/8 | 21.18 | 3110.05 | 11.27 |
| Smoothquant | 4/8 | 34.87 | 5133.82 | 20.82 |
| Smoothquant+GPTQ | 4/8 | 22.95 | 1359.59 | 13.39 |
| Smoothquant | 3/8 | 24041.06 | 42625.86 | 29585.39 |
| Smoothquant+GPTQ | 3/8 | 290.77 | - | 231.02 |

Table 6: Comparison with other quantization toolboxs.

| ToolBox | Number of Methods | Number of Benchmark&Datasets | Quantization Backend | Quantization Method Combination |
|-----------------------------|----------------------|---------------------------------|-------------------------|------------------------------------|
| Hugging Face Quanto library | 1 | - | \checkmark | × |
| qllm-eval (Li et al., 2024) | 3 | 30+ | × | × |
| TensorRT-LLM | 5 | - | \checkmark | × |
| VLLM (Kwon et al., 2023) | 4 | - | \checkmark | × |
| llama.cpp | 1 | - | \checkmark | × |
| Our Toolbox | 10+ | 40+ | \checkmark | \checkmark |

tion, Quantization, Evaluation, Inference,

and Execution modules. Tab. 6 presents the differences between our toolbox and other toolboxes. We provide a more comprehensive explanation of our toolbox in Appendix A.

Experimental Setup and Results. To validate the framework's capability of combining mixed 457 quantization methods, we conduct experiments using the LLaMA-2-7B model (Touvron et al., 2023). 458 We test the model using SmoothQuant and a combination of SmoothQuant for activations and GPTQ 459 for weight quantization on WikiText-2 (Wiki2) (Merity et al., 2016), Penn Treebank (PTB) (Marcus 460 et al., 1994), and C4 (Raffel et al., 2020) datasets, and measure the perplexity (PPL) of the quantized 461 models. Quantization is implemented using PyTorch. All quantization experiments are exclusively 462 conducted on the LLaMA-2-7B model, utilizing a single NVIDIA V100 GPU. For calibration, we 463 utilize a dataset consisting of 128 random segments, each containing 512 tokens, extract from the 464 C4 dataset. These segments represent generic text data, sourced from randomly crawled websites, 465 ensuring that the quantization process does not rely on task-specific information. Our quantization 466 setup employ SmoothQuant with default activation quantization of 8 bits. We utilize groupwise quantization with a group size of 128. 467

The results presented in Tab. 5 indicate several key findings. Comparing SmoothQuant with
SmoothQuant + GPTQ configurations, it is evident that the latter consistently outperforms the
former across all bit-width settings. This suggests that the combined use of SmoothQuant and
GPTQ leads to a notable improvement in model performance. Particularly, at bit-widths of 4 and 3,
the SmoothQuant + GPTQ method demonstrates a significant reduction in perplexity compared to
SmoothQuant alone, indicating the pronounced effectiveness of GPTQ in reducing perplexity.

474 475

432

5 RELATED WORK

476

477 Quantization of LLMs. Quantization techniques for LLMs mainly include Post-Training Quanti-478 zation (PTQ) and Quantization-Aware Training (QAT). PTQ does not require retraining the model 479 and is typically suitable for situations with limited computational resources (Frantar et al., 2022; 480 Dettmers et al., 2023; Lin et al., 2023; Xiao et al., 2023; Chee et al., 2024; Yao et al., 2022; Shao et al., 481 2023). QAT simulates the effects of quantization throughout the entire training process, enabling 482 the model to adapt to low-precision representations during training, which typically leads to higher performance (Liu et al., 2023b; Dettmers et al., 2024). It's worth noting that in this paper, we consider 483 applying quantization directly on the pretrained LLMs instead of performing quantization-aware 484 finetuning for the quantized LLMs (such as variants of QLoRA (Dettmers et al., 2024; Yi et al., 2024; 485 Xu et al., 2023)) because the latter typically needs the former for initialization.

| | Table 7: Comparison with related works. | | | | | | | |
|---|---|------------|------------------------------|--|--|--|--|--|
| _ | Work | Scenario | Number of Benchmark&Datasets | | | | | |
| | Jaiswal et al. (2023) | S 1 | 5 | | | | | |
| , | Williams & Aletras (2023) | S 1 | 10 | | | | | |
| | Liu et al. (2023a) | S 1 | 4 | | | | | |
| | Jin et al. (2024a) | S 1 | 10 | | | | | |
| | Li et al. (2024) | S 1 | 19 | | | | | |
| | Huang et al. (2024a) | S 1 | 9 | | | | | |
| | Our Method | S1&S2 | 40+ | | | | | |

Evaluation of guantized LLMs. Numerous studies have undertaken evaluations of the performance 497 of quantized LLMs (Frantar et al., 2022; Dettmers et al., 2023; Lin et al., 2023; Xiao et al., 2023; Chee 498 et al., 2024; Jaiswal et al., 2023; Williams & Aletras, 2023; Li et al., 2024; Liu et al., 2023a; Jin et al., 499 2024b; Huang et al., 2024a). The majority of assessments employ fixed calibration set, primarily 500 focusing on language modeling tasks (Raffel et al., 2020; Marcus et al., 1994; Merity et al., 2016) and 501 standard NLP tasks (Zellers et al., 2019; Paperno et al., 2016; Tata & Patel, 2003; Clark et al., 2018; 502 Sakaguchi et al., 2021; Mihaylov et al., 2018; Mostafazadeh et al., 2016). Certain investigations have deviated from the practice of using fixed calibration set, extending them to encompass a broader 504 spectrum of crawled web text and pre-training data, while also conducting multiple random samplings 505 for calibration set selection (Williams & Aletras, 2023). Additionally, certain studies have conducted assessments encompassing a broader array of downstream task types and datasets, approaching the 506 evaluation from various angles (Liu et al., 2023a; Jaiswal et al., 2023; Li et al., 2024). 507

508 Differences between our work and related work. We have provided a detailed presentation of 509 the differences between our work and related work in Tab. 7. The related work primarily addresses 510 scenarios similar to S1 in our experiments and does not involve experiments related to the S2 scenario. 511 Additionally, while some studies have considered distribution shifts (such as zero-shot and in-context learning), the scope of these shifts is limited, and the calibration datasets are fixed. This limitation 512 results in a lack of systematic analysis regarding generalization capability and distribution shifts. 513 These evaluations did not account for high-level generalization scenario classifications or assess 514 variations in generalization ability across different settings. (Liu et al., 2023a) investigated the impact 515 of quantization on model emergent abilities, evaluating OOD generalization tasks including zero-shot 516 and in-context learning (e.g., ICL, CoT, instruction following). The evaluation types are similar to 517 those in our S1 scenario but with limited scope. (Jaiswal et al., 2023) argued that perplexity (PPL) 518 is not a good evaluation metric and thus evaluated numerous popular zero-shot tasks, similar to our 519 S1 experiments but with limited scope and extent of distribution shifts. Li et al. (2024) evaluated 520 a broader range of tasks and capabilities compared to previous work, but still focused on our S1 521 scenario and did not include tests for various distribution shifts.

522 523

524

486

6 FUTURE WORK AND CONCLUSION

Despite comprehensive evaluation on over 50 datasets, our study acknowledges the need for a more thorough assessment of models and quantization algorithms. Future work could involve a more extensive evaluation framework. Additionally, the developed toolbox does not yet support all quantization algorithms and large models. Further development is warranted to expand its capabilities.

We investigat the generalization ability of quantized LLMs, proposing two evaluation scenarios and 530 testing them on our own implemented platform. S1 demonstrates that the quantized LLM maintains 531 its generalization capability under all situations except those involving extreme bit quantization. 532 Building upon this foundation, we introduce a novel generalization shift evaluation framework in S2, 533 which investigates methods for enhancing generalization ability from a data perspective. Drawing 534 from our evaluation results, we find that quantized models do not benefit from the alignment between calibration and test distributions. Further investigation revealed that this may be attributed to the 536 fact that quantized models do not require a substantial amount of data relevant to downstream tasks 537 to recover performance. Our work unveils the relationship between calibration data and test data, prompting the development of novel methods for optimizing calibration data collection, which is 538 overlooked in the current field of model quantization. Lastly, we provided a modular and scalable toolbox to this topic to facilitate future research.

540 REFERENCES

552

553

554

555

578

579

580

581

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

- Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi Wang, Niklas Muennighoff, Bairu Hou, Liangming Pan, Haewon Jeong, et al. A survey on data selection for language models. *arXiv preprint arXiv:2402.16827*, 2024.
- Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi.
 Mathqa: Towards interpretable math word problem solving with operation-based formalisms.
 arXiv preprint arXiv:1905.13319, 2019.
 - Max Bartolo, Alastair Roberts, Johannes Welbl, Sebastian Riedel, and Pontus Stenetorp. Beat the ai: Investigating adversarial human annotation for reading comprehension. *Transactions of the Association for Computational Linguistics*, 8:662–678, 2020.
- Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced metrics
 for measuring unintended bias with real data for text classification. In *Companion proceedings of the 2019 world wide web conference*, pp. 491–500, 2019.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. Quip: 2-bit quantization of
 large language models with guarantees. *Advances in Neural Information Processing Systems*, 36,
 2024.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- Leyang Cui, Yu Wu, Shujie Liu, Yue Zhang, and Ming Zhou. Mutual: A dataset for multi-turn dialogue reasoning. *arXiv preprint arXiv:2004.04494*, 2020.
- Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashkboos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. Spqr: A sparse-quantized representation for near-lossless llm weight compression. *arXiv preprint arXiv:2306.03078*, 2023.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning
 of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm: General language model pretraining with autoregressive blank infilling. *arXiv preprint arXiv:2103.10360*, 2021.
- Matthew Dunn, Levent Sagun, Mike Higgins, V Ugur Guney, Volkan Cirik, and Kyunghyun Cho.
 Searchqa: A new q&a dataset augmented with context from a search engine. *arXiv preprint arXiv:1704.05179*, 2017.
- 585
 586 Mai ElSherief, Caleb Ziems, David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, and Diyi Yang. Latent hatred: A benchmark for understanding implicit hate speech. arXiv preprint arXiv:2109.05322, 2021.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence
 Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, et al. A framework for few-shot
 language model evaluation. *Version v0. 0.1. Sept*, pp. 8, 2021a.

594 Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence 595 Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, et al. A framework for few-shot 596 language model evaluation. Version v0. 0.1. Sept, pp. 8, 2021b. 597 Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A 598 survey of quantization methods for efficient neural network inference. In Low-Power Computer Vision, pp. 291–326. Chapman and Hall/CRC, 2022. 600 601 Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. 602 Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. 603 arXiv preprint arXiv:2203.09509, 2022. 604 605 Wei Huang, Xudong Ma, Haotong Oin, Xingyu Zheng, Chengtao Ly, Hong Chen, Jie Luo, Xiaojuan Qi, Xianglong Liu, and Michele Magno. How good are low-bit quantized llama3 models? an 606 empirical study. arXiv preprint arXiv:2404.14047, 2024a. 607 608 Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, 609 Chuancheng Lv, Yikai Zhang, Yao Fu, et al. C-eval: A multi-level multi-discipline chinese 610 evaluation suite for foundation models. Advances in Neural Information Processing Systems, 36, 611 2024b. 612 613 Ajay Jaiswal, Zhe Gan, Xianzhi Du, Bowen Zhang, Zhangyang Wang, and Yinfei Yang. Compressing llms: The truth is rarely pure and never simple. arXiv preprint arXiv:2310.01382, 2023. 614 615 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. Pubmedqa: A 616 dataset for biomedical research question answering. arXiv preprint arXiv:1909.06146, 2019. 617 618 Renren Jin, Jiangcun Du, Wuwei Huang, Wei Liu, Jian Luan, Bin Wang, and Deyi Xiong. A 619 comprehensive evaluation of quantization strategies for large language models. arXiv preprint 620 arXiv:2402.16775, 2024a. 621 Renren Jin, Jiangcun Du, Wuwei Huang, Wei Liu, Jian Luan, Bin Wang, and Deyi Xiong. A 622 comprehensive evaluation of quantization strategies for large language models. arXiv preprint 623 arXiv:2402.16775, 2024b. 624 625 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott 626 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. 627 arXiv preprint arXiv:2001.08361, 2020. 628 Yuta Koreeda and Christopher D Manning. Contractnli: A dataset for document-level natural language 629 inference for contracts. arXiv preprint arXiv:2110.01799, 2021. 630 631 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph 632 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model 633 serving with pagedattention. In Proceedings of the 29th Symposium on Operating Systems 634 Principles, pp. 611-626, 2023. 635 636 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading 637 comprehension dataset from examinations. arXiv preprint arXiv:1704.04683, 2017. 638 Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In 639 Thirteenth international conference on the principles of knowledge representation and reasoning, 640 2012. 641 642 Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy 643 Baldwin. Cmmlu: Measuring massive multitask language understanding in chinese. arXiv preprint 644 arXiv:2306.09212, 2023. 645 Shiyao Li, Xuefei Ning, Luning Wang, Tengxuan Liu, Xiangsheng Shi, Shengen Yan, Guohao Dai, 646 Huazhong Yang, and Yu Wang. Evaluating quantized large language models. arXiv preprint 647

arXiv:2402.18158, 2024.

| 648 649 650 | Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. Awq: Activation- aware weight quantization for llm compression and acceleration. <i>arXiv preprint arXiv:2306.00978</i> , 2023. |
|---------------------------------|---|
| 651 652 653 | Alisa Liu, Swabha Swayamdipta, Noah A Smith, and Yejin Choi. Wanli: Worker and ai collaboration for natural language inference dataset creation. <i>arXiv preprint arXiv:2201.05955</i> , 2022. |
| 654 655 656 | Peiyu Liu, Zikang Liu, Ze-Feng Gao, Dawei Gao, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen. Do emergent abilities exist in quantized large language models: An empirical study. <i>arXiv preprint arXiv:2307.08072</i> , 2023a. |
| 657 658 659 660 | Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. Llm-qat: Data-free quantization aware training for large language models. <i>arXiv preprint arXiv:2305.17888</i> , 2023b. |
| 661 662 663 664 | Mitch Marcus, Grace Kim, Mary Ann Marcinkiewicz, Robert MacIntyre, Ann Bies, Mark Ferguson, Karen Katz, and Britta Schasberger. The penn treebank: Annotating predicate argument structure. In <i>Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March</i> 8-11, 1994, 1994. |
| 665 666 667 | Julian McAuley and Jure Leskovec. Hidden factors and hidden topics: understanding rating dimensions with review text. In <i>Proceedings of the 7th ACM conference on Recommender systems</i> , pp. 165–172, 2013. |
| 668 669 670 | Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. <i>arXiv preprint arXiv:1609.07843</i> , 2016. |
| 671 672 673 | Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. <i>arXiv preprint arXiv:1809.02789</i> , 2018. |
| 674 675 676 677 | Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Van- derwende, Pushmeet Kohli, and James Allen. A corpus and evaluation framework for deeper understanding of commonsense stories. <i>arXiv preprint arXiv:1604.01696</i> , 2016. |
| 678 679 680 | Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart Van Baalen, and Tij- men Blankevoort. A white paper on neural network quantization. <i>arXiv preprint arXiv:2106.08295</i> , 2021. |
| 681 682 | Preslav Nakov, Alan Ritter, Sara Rosenthal, Fabrizio Sebastiani, and Veselin Stoyanov. Semeval-2016 task 4: Sentiment analysis in twitter. <i>arXiv preprint arXiv:1912.01973</i> , 2019. |
| 684 685 686 | Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R Bowman. Crows-pairs: A challenge dataset for measuring social biases in masked language models. <i>arXiv preprint arXiv:2010.00133</i> , 2020. |
| 687 688 | Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial nli: A new benchmark for natural language understanding. <i>arXiv preprint arXiv:1910.14599</i> , 2019. |
| 690 691 692 | Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset: Word prediction requiring a broad discourse context. <i>arXiv preprint arXiv:1606.06031</i> , 2016. |
| 693 694 695 696 697 | Anselmo Peñas, Eduard Hovy, Pamela Forner, Álvaro Rodrigo, Richard Sutcliffe, and Roser Morante. Qa4mre 2011-2013: Overview of question answering for machine reading evaluation. In <i>Infor-</i> <i>mation Access Evaluation. Multilinguality, Multimodality, and Visualization: 4th International</i> <i>Conference of the CLEF Initiative, CLEF 2013, Valencia, Spain, September 23-26, 2013. Proceed-</i> <i>ings 4</i> , pp. 303–320. Springer, 2013. |
| 698 699 700 | Christopher Potts, Zhengxuan Wu, Atticus Geiger, and Douwe Kiela. Dynasent: A dynamic benchmark for sentiment analysis. <i>arXiv preprint arXiv:2012.15349</i> , 2020. |
| | ALL DESCRIPTION WE DESCRIPTION STREET DESCRIPTION STREET |

701 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

702 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi 703 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text 704 transformer. Journal of machine learning research, 21(140):1-67, 2020. 705 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for 706 machine comprehension of text. arXiv preprint arXiv:1606.05250, 2016. 708 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An 709 adversarial winograd schema challenge at scale. Communications of the ACM, 64(9):99-106, 710 2021. 711 Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, 712 Peng Gao, Yu Qiao, and Ping Luo. Omniquant: Omnidirectionally calibrated quantization for large 713 language models. arXiv preprint arXiv:2308.13137, 2023. 714 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and 715 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. 716 In Proceedings of the 2013 conference on empirical methods in natural language processing, pp. 717 1631–1642, 2013. 718 719 Sandeep Tata and Jignesh M Patel. Piqa: An algebra for querying protein data sets. In 15th 720 International Conference on Scientific and Statistical Database Management, 2003., pp. 141–150. 721 IEEE, 2003. 722 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 723 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and 724 efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023. 725 Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordoni, Philip Bachman, and 726 Kaheer Suleman. Newsqa: A machine comprehension dataset. arXiv preprint arXiv:1611.09830, 727 2016. 728 729 A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017. 730 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: 731 A multi-task benchmark and analysis platform for natural language understanding. arXiv preprint 732 arXiv:1804.07461, 2018. 733 734 Jiahao Wang, Bolin Zhang, Qianlong Du, Jiajun Zhang, and Dianhui Chu. A survey on data selection for llm instruction tuning. arXiv preprint arXiv:2402.05123, 2024. 735 736 Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and 737 Samuel R Bowman. Blimp: The benchmark of linguistic minimal pairs for english. Transactions 738 of the Association for Computational Linguistics, 8:377–392, 2020. 739 Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. 740 arXiv preprint arXiv:1707.06209, 2017. 741 742 Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus for 743 sentence understanding through inference. arXiv preprint arXiv:1704.05426, 2017. 744 Miles Williams and Nikolaos Aletras. How does calibration data affect the post-training pruning and 745 quantization of large language models? arXiv preprint arXiv:2311.09755, 2023. 746 747 Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: 748 Accurate and efficient post-training quantization for large language models. In International 749 Conference on Machine Learning, pp. 38087–38099. PMLR, 2023. 750 Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhensu Chen, 751 Xiaopeng Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language 752 models. arXiv preprint arXiv:2309.14717, 2023. 753 Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, 754 Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. arXiv preprint 755 arXiv:2309.10305, 2023.

| 756 757 758 | Linyi Yang, Shuibai Zhang, Libo Qin, Yafu Li, Yidong Wang, Hanmeng Liu, Jindong Wang, Xing Xie, and Yue Zhang. Glue-x: Evaluating natural language understanding models from an out-of-distribution generalization perspective. <i>arXiv preprint arXiv:2211.08073</i> , 2022. |
|--------------------------|--|
| 759 760 761 762 | Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong He. Zeroquant: Efficient and affordable post-training quantization for large-scale transformers. <i>Advances in Neural Information Processing Systems</i> , 35:27168–27183, 2022. |
| 763 764 765 | Ke Yi, Yuhui Xu, Heng Chang, Chen Tang, Yuan Meng, Tong Zhang, and Jia Li. One quantllm for all: Fine-tuning quantized llms once for efficient deployments. <i>arXiv preprint arXiv:2405.20202</i> , 2024. |
| 766 767 768 | Han Yu, Jiashuo Liu, Xingxuan Zhang, Jiayun Wu, and Peng Cui. A survey on evaluation of out-of-distribution generalization. <i>arXiv preprint arXiv:2403.01874</i> , 2024. |
| 769 770 771 | Lifan Yuan, Yangyi Chen, Ganqu Cui, Hongcheng Gao, Fangyuan Zou, Xingyi Cheng, Heng Ji, Zhiyuan Liu, and Maosong Sun. Revisiting out-of-distribution robustness in nlp: Benchmarks, analysis, and Ilms evaluations. <i>Advances in Neural Information Processing Systems</i> , 36, 2024. |
| 772 773 774 | Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. Swag: A large-scale adversarial dataset for grounded commonsense inference. <i>arXiv preprint arXiv:1808.05326</i> , 2018. |
| 775 776 | Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.07830</i> , 2019. |
| 777 778 779 | Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> , 2022. |
| 780 781 782 783 | Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. " going on a vacation" takes longer than" going for a walk": A study of temporal commonsense understanding. <i>arXiv preprint arXiv:1909.03065</i> , 2019. |
| 784 785 786 | Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large language models. <i>arXiv preprint arXiv:2308.07633</i> , 2023. |
| 787 | |
| 788 | |
| 789 | |
| 790 | |
| 791 | |
| 792 | |
| 793 | |
| 794 | |
| 795 | |
| 796 | |
| 797 | |
| 798 | |
| 799 | |
| 800 | |
| 801 | |
| 802 | |
| 803 | |
| 804 | |
| 805 | |
| 806 | |
| 807 | |
| 800 | |
| 003 | |

A MORE DETAILS OF OUR TOOLBOX

Fig. 4 illustrates the framework of MI-optimize, which comprises five main modules: the Configuration, Quant, Evaluation, Inference, and Execution modules. Combining these modules forms a cohesive pipeline that provides researchers with a reliable experimental environment, with each module responsible for a specific step in the pipeline. The subsequent sections will provide a detailed description of the implementation of each module.

- **Configuration Module**: Manages all parameters involved in the framework, including default settings, quantization configurations, and evaluation configurations.
 - Model Module: Contains various pre-trained models such as LLaMA (Touvron et al., 2023), Baichuan (Yang et al., 2023), ChatGLM (Du et al., 2021), and custom user models.
- **Dataset Module**: Handles different datasets, including Chinese domain-specific datasets (e.g., C-EVAL (Huang et al., 2024b) and CMMLU (Li et al., 2023)), the BOSS benchmark (Yuan et al., 2024), general datasets (e.g., Amazon reviews, Dynasent), and LM-EVAL datasets (e.g., Winogrande, WSC273).
 - **Quant Module**: Responsible for loading pre-trained models, applying various quantization methods (e.g., GPTQ (Frantar et al., 2022), AWQ (Lin et al., 2023), SPQR (Dettmers et al., 2023)), and performing the actual model quantization.
 - Inference & Eval Module: Exports the quantized model, runs inference using engines such as VLLM (Kwon et al., 2023) and TensorRT, and evaluates benchmark performance.
 - **Execution Module**: Oversees the primary tasks of model quantization, benchmarking, and the combined process of quantization and evaluation.

Key Features Supported by MI-optimize.

- Quantization of LLMs to reduce computational and memory requirements: MI-optimize focuses on reducing the computational and memory footprint of large language models through advanced quantization techniques, making them more suitable for deployment in resource-limited environments.
- Support for various quantization algorithms: The framework supports a wide range of quantization algorithms, including RTN, GPTQ Frantar et al. (2022), AWQ Lin et al. (2023), SpQR Dettmers et al. (2023), ZeroQuant Yao et al. (2022), SmoothQuant Xiao et al. (2023), QuIP Chee et al. (2024), and FP8. This flexibility allows users to choose the most appropriate method for their specific use case, optimizing performance and resource usage.
- Evaluation on OOD tasks using benchmarks: MI-optimize includes tools for evaluating quantized models on out-of-distribution (OOD) tasks using established benchmarks such as BOSS. This ensures that the models maintain their performance even when encountering data that differs from their training set.
 - Support for multiple datasets: The framework supports multiple datasets for both calibration and testing purposes. Users can also incorporate custom datasets to better align the model's performance with their specific requirements.
 - Command-line interface for easy integration and automation: MI-optimize provides a command-line interface that facilitates easy integration into existing workflows and automation of the quantization and evaluation processes, streamlining the deployment pipeline.
 - Support for combination of quantization methods: The framework allows for the combination of different quantization methods within the same model. Different layers can apply different quantization algorithms, and even multiple quantization algorithms can be applied to the same layer. This granular control helps optimize model performance and efficiency.
 - Ease of adding new quantization algorithms: Researchers can easily add new quantization algorithms to the MI-optimize repository. This modularity ensures that the framework remains up-to-date with the latest advancements in quantization techniques.
- Customer tools for model quantization and evaluation: Customers can install the tools provided by MI-optimize to quantize and evaluate their own models. This empowers users to tailor the framework to their specific needs, ensuring optimal model performance in their applications.

B DATASETS

864

865 866

867

868

869 870 871

872

In this section, we present all the datasets utilized in the experiments, encompassing their evaluated tasks and abilities, assessment metrics, and dataset sizes. Tab. 8 and 9 provide a comprehensive summary of all the datasets.

B.1 DATASETS IN S1

873 **Common sense reasoning.** WinoGrande (Sakaguchi et al., 2021) is a large-scale coreference reso-874 lution task dataset derived from extensive internet text, aimed at addressing ambiguous and complex 875 coreference relationships. WSC273 (Levesque et al., 2012) comprises 273 coreference resolution 876 problems derived from the classic Winograd Schema Challenge, primarily assessing the common-877 sense reasoning capabilities of natural language understanding systems. GLUE-WNLI (Wang 878 et al., 2018) is designed to test coreference resolution capability, which involves determining which noun a pronoun in a sentence refers to. It is sourced from the Winograd Schema Challenge. Hel-879 laSwag (Zellers et al., 2019) is generated from web videos and Wikipedia articles and is used to infer 880 the most suitable continuation for text segments in multiple-choice tasks. SWAG (Zellers et al., 2018) 881 is generated based on video descriptions, aiming to predict plausible subsequent scenarios for video 882 events. PIQA (Tata & Patel, 2003) is a dataset for reasoning about physical common sense, derived 883 from physics problems and solutions, designed to evaluate algorithms' reasoning abilities in physical 884 environments. 885

Mathematical reasoning. MathQA (Amini et al., 2019) is collected from the MathQA website, consisting of 37,200 mathematical questions, with the task being to automatically answer mathematical questions.

Multi-turn dialogue reasoning. MuTual (Cui et al., 2020) and Mutual_plus (Cui et al., 2020) is a
 retrieval-based dataset for multi-turn dialogue reasoning, which is modified from Chinese high school
 English listening comprehension test data.

Bias diagnosis and mitigation. CrowS-Pairs (Nangia et al., 2020) is derived from a wide range of internet text and is designed to evaluate social biases in language models. Toxigen (Hartvigsen et al., 2022) is for implicit hate speech detection.

Scientific knowledge question answering. PubMedQA (Jin et al., 2019) is a biomedical question 896 answering dataset sourced from PubMed articles, aimed at evaluating systems' understanding and 897 answering capabilities of biomedical texts. OpenBookQA (Mihaylov et al., 2018) is a new kind 898 of question-answering dataset modeled after open book exams for assessing human understanding 899 of a subject. It originates from open science education resources. SciQ (Welbl et al., 2017) is 900 a high-quality, science-themed multiple-choice dataset constructed manually. ARC-Easy (Clark 901 et al., 2018) originates from science exams administered in American elementary through high 902 schools, assessing fundamental scientific knowledge. ARC-Challenge (Clark et al., 2018) presents 903 challenging scientific questions aimed at testing higher-level scientific comprehension and reasoning abilities. MC-TACO (Zhou et al., 2019) consists of temporal common-sense questions sourced from 904 a wide range of internet texts, designed for temporal common-sense reasoning tasks. 905

Reading comprehension. RACE (Lai et al., 2017) is a large-scale reading comprehension dataset sourced from English exams for Chinese middle school and high school students, aimed at testing reading comprehension abilities. QA4MRE (Peñas et al., 2013) is created for the CLEF 2011/2012/2013 shared tasks, aimed at testing cross-domain reading comprehension abilities.

910 Natural language inference. GLUE-MNLI (Wang et al., 2018) is a natural language inference 911 dataset comprising pairs of sentences sourced from various text genres such as novels, telephone 912 conversations, and news articles. GLUE-MNLI-Mismatched (Wang et al., 2018) is utilized to 913 evaluate the generalization capability of models on unseen text genres, with sentence pairs sourced 914 from the same origins as GLUE-MNLI. GLUE-RTE (Wang et al., 2018) is sourced from news reports 915 and Wikipedia. GLUE-QNLI (Wang et al., 2018) originates from the Stanford University's SQuAD dataset. ANLI (Nie et al., 2019) is a large-scale adversarial natural language inference dataset divided 916 into three difficulty levels. It is constructed by employing adversarial search techniques to generate 917 challenging questions based on human annotations.

Sentiment analysis. GLUE-SST (Wang et al., 2018) is sourced from movie reviews, and its task involves sentiment classification, which entails determining the emotional inclination of a sentence.

Syntax phenomena evaluation. BLiMP (Warstadt et al., 2020) is a challenge set for evaluating what language models know about major grammatical phenomena in English. BLiMP consists of 67 sub-datasets, each containing 1000 minimal pairs isolating specific contrasts in syntax, morphology, or semantics. The data is automatically generated according to expert-crafted grammars.

- 925
- 926

B.2 DATASETS IN S2

927 928

929 Extractive question answering in BOSS. SQuAD (Rajpurkar et al., 2016) is a collection of question-930 answer pairs derived from Wikipedia articles. AdversarialQA (Bartolo et al., 2020) formulates 931 adversarial questions within the SQuAD context, utilizing a collaborative process involving both 932 human annotators and models. NewsQA (Trischler et al., 2016) crafts questions based on CNN news articles, each demanding reasoning for answers, rather than relying solely on lexical overlap and 933 textual entailment. SearchQA (Dunn et al., 2017) employs a reverse construction approach, utilizing 934 the Google search engine to fetch pertinent contexts for each question-answer pair from the J!Archive 935 website. 936

Sentiment analysis in BOSS. Amazon (McAuley & Leskovec, 2013) is a dataset comprising reviews across 29 distinct product categories from the Amazon website. DynaSent (Potts et al., 2020) constructs a dataset by identifying challenging sentences from existing collections and generating adversarial counterparts through human-and-model collaborative annotation. SemEval (Nakov et al., 2019) offers a three-class sentiment analysis dataset centered on Twitter content. SST (Socher et al., 2013) features sentence-level movie reviews sourced from the Rotten Tomatoes website.

943 Natural language inference in BOSS. MNLI (Williams et al., 2017) offers sentence pairs across
944 ten diverse categories of written and verbal communication, showcasing various styles, topics, and
945 formalities. ANLI (Nie et al., 2019) is an adversarial dataset created using a human-and-model-in946 the-loop method, featuring premises primarily sourced from Wikipedia and hypotheses crafted by
947 human adversaries. ContractNLI (Koreeda & Manning, 2021) treats individual contracts as premises
948 and applies a consistent set of hypotheses across the dataset. WANLI (Liu et al., 2022) is generated
949 by GPT-3, containing examples that include challenging patterns initially identified in MNLI.

Toxic detection in BOSS. Civil Comments (Borkan et al., 2019) features public comments from the Civil Comments platform, encompassing a diverse user base and various subtypes of toxic text.
AdvCivil introduces a new toxic dataset, derived from Civil Comments through textual adversarial attacks within an automated model-in-the-loop adversarial pipeline. Implicit Hate (ElSherief et al., 2021) includes toxic tweets that are both explicit and implicit, with the latter capable of evading keyword-based toxic detection systems. ToxiGen (Hartvigsen et al., 2022) is generated by GPT-3 and contains subtly and implicitly toxic texts targeting 13 minority groups.

Chinese domain-specific. C-Eval (Huang et al., 2024b) is a comprehensive Chinese evaluation suite
 for foundation models. It consists of 13948 multi-choice questions spanning 52 diverse disciplines
 and four difficulty levels, primarily encompassing humanities, social sciences, STEM, and other 4
 categories. CMMLU (Li et al., 2023) is a comprehensive Chinese evaluation benchmark designed
 specifically to assess language models' knowledge and reasoning abilities within Chinese contexts.
 CMMLU covers 67 topics ranging from fundamental subjects to advanced professional levels. It
 encompasses topics such as STEM requiring calculation and reasoning, humanities and social sciences
 necessitating knowledge, and everyday knowledge such as Chinese driving rules.

- 964
- 965 966

C EXPERIMENT DETAILS

- 967 968
- 969

In this section, we will present all the details of our experiment, including hardware resources,
 experimental setup, hyperparameter selection, and data selection. Besides, Our benchmark suite is available in the supplementary materials.

| Table 8: Summary of the datasets in S | 1 |
|---------------------------------------|---|
|---------------------------------------|---|

| | Scenario | Task&Ability | Dataset | Gene. | Metric | Size |
|---|----------|---|---|-------|----------------|-------|
| | S1 | Common sense reasoning | WinoGrande Sakaguchi et al. (2021) | 0/5 | Acc | 1267 |
| | S1 | Common sense reasoning | WSC273 Levesque et al. (2012) | 0/5 | Acc | 273 |
| | S1 | Common sense reasoning | GLUE-WNLI Wang et al. (2018) | 0/5 | Acc | 71 |
| | S1 | Common sense reasoning | HellaSwag Zellers et al. (2019) | 0/5 | Acc | 10042 |
| | S1 | Common sense reasoning | SWAG Zellers et al. (2018) | 0/5 | Acc | 20006 |
| | S1 | Common sense reasoning | PIQA Tata & Patel (2003) | 0/5 | Acc | 1838 |
| | S1 | Mathematical reasoning | MathQA Amini et al. (2019) | 0/5 | Acc | 2985 |
| | S1 | Multi-turn dialogue reasoning | Mutual Cui et al. (2020) | 0/5 | R2 | 886 |
| | S1 | Multi-turn dialogue reasoning | Mutual_Plus Cui et al. (2020) | 0/5 | R2 | 886 |
| | S1 | Bias diagnosis and mitigation | CrowS-Pairs Nangia et al. (2020) | 0 | Pct_stereotype | 6708 |
| | S1 | Bias diagnosis and mitigation | Toxigen Hartvigsen et al. (2022) | 0/5 | Acc | 940 |
| | S1 | Scientific knowledge question answering | PubMedQA Jin et al. (2019) | 0/5 | Acc | 1000 |
| | S1 | Scientific knowledge question answering | OpenBookQA Mihaylov et al. (2018) | 0/5 | Acc | 500 |
| | S1 | Scientific knowledge question answering | SciQ Welbl et al. (2017) | 0/5 | Acc | 1000 |
| | S1 | Scientific knowledge question answering | ARC-Easy Clark et al. (2018) | 0/5 | Acc | 2376 |
| | S1 | Scientific knowledge question answering | ARC-Challenge Clark et al. (2018) | 0/5 | Acc | 1172 |
| | S1 | Scientific knowledge question answering | MC-TACO Zhou et al. (2019) | 0/5 | F1 | 9442 |
| | S1 | Reading comprehension | RACE Lai et al. (2017) | 0/5 | Acc | 1045 |
| | S1 | Reading comprehension | QA4MRE Peñas et al. (2013) | 0/5 | Acc | 564 |
| | S1 | Natural language inference | GLUE-MNLI Wang et al. (2018) | 0/5 | Acc | 9815 |
| | S1 | Natural language inference | GLUE-MNLI-Mismatched Wang et al. (2018) | 0/5 | Acc | 9832 |
| | S1 | Natural language inference | GLUE-RTE Wang et al. (2018) | 0/5 | Acc | 277 |
| | S1 | Natural language inference | GLUE-QNLI Wang et al. (2018) | 0/5 | Acc | 5463 |
| | S1 | Natural language inference | ANLI Nie et al. (2019) | 0/5 | Acc | 3200 |
| _ | S1 | Sentiment analysis | GLUE-SST Wang et al. (2018) | 0/5 | Acc | 872 |
| _ | S1 | Syntax phenomena evaluation | BLiMP Warstadt et al. (2020) | 5 | Acc | 67000 |
| | | | | | | |

Table 9: Summary of the datasets in S2.

| Scenario | Task&Ability | Dataset | Gene. | Metric | Size |
|----------|-------------------------------|---------------------------------------|-------|--------|-------|
| S2 | Extractive question answering | SQuAD Rajpurkar et al. (2016) | 0/1 | F1 | 10570 |
| S2 | Extractive question answering | AdversarialQA Bartolo et al. (2020) | 0/1 | F1 | 2694 |
| S2 | Extractive question answering | NewsQA Trischler et al. (2016) | 0/1 | F1 | 3912 |
| S2 | Extractive question answering | SearchQA Dunn et al. (2017) | 0/1 | F1 | 16680 |
| S2 | Sentiment analysis | Amazon McAuley & Leskovec (2013) | 0/3 | Acc | 38905 |
| S2 | Sentiment analysis | DynaSent Potts et al. (2020) | 0/3 | Acc | 4020 |
| S2 | Sentiment analysis | SemEval Nakov et al. (2019) | 0/3 | Acc | 20322 |
| S2 | Sentiment analysis | SST Socher et al. (2013) | 0/3 | Acc | 767 |
| S2 | Natural language inferenc | MNLI Williams et al. (2017) | 0/3 | Acc | 9815 |
| S2 | Natural language inferenc | ANLI Nie et al. (2019) | 0/3 | Acc | 2900 |
| S2 | Natural language inferenc | ContractNLI Koreeda & Manning (2021) | 0/3 | Acc | 1791 |
| S2 | Natural language inferenc | WANLI Liu et al. (2022) | 0/3 | Acc | 4700 |
| S2 | Toxic detection | Civil Comments Borkan et al. (2019) | 0/2 | Acc | 97320 |
| S2 | Toxic detection | AdvCivil | 0/2 | Acc | 523 |
| S2 | Toxic detection | Implicit Hate ElSherief et al. (2021) | 0/2 | Acc | 21180 |
| S2 | Toxic detection | ToxiGen Hartvigsen et al. (2022) | 0/2 | Acc | 641 |
| S2 | Chinese domainspecific | CEVAL Huang et al. (2024b) | 0/5 | Acc | 13948 |
| S2 | Chinese domainspecific | CMMLU Li et al. (2023) | 0/5 | Acc | 11917 |

1026 C.1 HARDWARE RESOURCES

1027

1031

In our experiments, we utilize one computer with 8 AMD Aldebaran GPUs and two computers with
2 NVIDIA Tesla V100 GPUs each. Specifically, each AMD Aldebaran GPU has 64GB of memory,
totaling 512GB. Each NVIDIA Tesla V100 GPU has 32GB of memory, totaling 128GB.

1032 C.2 EXPERIMENT DETAILS IN S1

Experimental Setup. We quantize LLaMA2-7B (Touvron et al., 2023) using the GPTQ (Frantar et al., 2022), SpQR (Dettmers et al., 2023) methods. We quantize the weights to 2-4 bits and test 16 bits as reference. The quantization is implemented using our custom toolbox, maintaining consistency with the original method in all experimental details.

Hyperparameter Selection. For the GPTQ (Frantar et al., 2022) method, we set the group-size parameter to 128 and apply block-sequential as well as layer-sequential quantization. For the SpQR (Dettmers et al., 2023) method, we set the group-size parameter to 128 and apply block-sequential quantization. Throughout the quantization process, we use 128 calibration examples. In the few-shot setting, the number of selected examples corresponds to LM Evaluation Harness (Gao et al., 2021b), remaining at 5-shot.

Data Selection. We follow GPTO (Frantar et al., 2022) and randomly sample 128 samples from 1044 C4-en-val (Raffel et al., 2020) as the calibration set with a random seed of 42. For the selection 1045 of test data, we use the test splits of ANLI (Nie et al., 2019), ARC (Clark et al., 2018), CrowS-1046 Pairs (Nangia et al., 2020), GLUE-MNLI-Mismatched (Wang et al., 2018), MathQA (Amini et al., 1047 2019), MCTACO (Zhou et al., 2019), OpenBookQA (Mihaylov et al., 2018), RACE (Lai et al., 1048 2017), SciQ (Welbl et al., 2017), Toxigen (Hartvigsen et al., 2022), and WSC273 (Levesque et al., 1049 2012) as the test set. We use the validation splits of GLUE-SST, GLUE-MNLI, GLUE-QNLI, 1050 GLUE-WNLI, GLUE-RTE (Wang et al., 2018), HellaSwag (Zellers et al., 2019), Mutual (Cui et al., 1051 2020), PIQA (Tata & Patel, 2003), SWAG (Zellers et al., 2018), WinoGrande (Sakaguchi et al., 2021) 1052 as the test set. Additionally, we use the train splits of BLiMP (Warstadt et al., 2020), PubMedQA (Jin 1053 et al., 2019), and QA4MRE (Peñas et al., 2013) as the test set. For the selection of examples in the 1054 few-shot setting, we use the default setting.

1055

1056 C.3 EXPERIMENT DETAILS IN S2

1058 C.3.1 BOSS

Experimental Setup. We quantize LLaMA2-7B (Touvron et al., 2023) using the GPTQ (Frantar et al., 2022), SpQR (Dettmers et al., 2023), AWQ (Lin et al., 2023), and SmoothQuant (Xiao et al., 2023) methods. We quantize the weights to 3-4 bits, and for smoothquant, we further quantize the activations to 8 bits. The quantization is implemented using our custom toolbox, maintaining consistency with the original method in all experimental details.

Hyperparameter Selection. For the GPTQ (Frantar et al., 2022) method, we set the group-size parameter to 128 and apply block-sequential as well as layer-sequential quantization. For the SpQR (Dettmers et al., 2023) method, we set the group-size parameter to 128 and apply block-sequential quantization. For the AWQ (Lin et al., 2023) method, we set the group-size parameter to 128. Throughout the quantization process, we use 128 calibration examples. In the few-shot setting, the number of selected examples corresponds to those in BOSS. Specifically, EQA is 1-shot, SA and NLI are 3-shot, and TD is 2-shot. The prompt template is presented in Tab. 10.

1071 Data Selection. For the calibration set, we use 128 calibration examples. For SQuAD (Rajpurkar 1072 et al., 2016) dataset in EQA, Amazon (McAuley & Leskovec, 2013) dataset in SA, MNLI (Williams 1073 et al., 2017) dataset in NLI, and Civil Comments (Borkan et al., 2019) dataset in TD, as the original 1074 datasets include train and test splits, we directly select the first 128 instances from the train split as 1075 the calibration set. For the remaining datasets, given that the original datasets exclusively contain a test split, we randomly sample 300 instances from the test split to form a train split, subsequently removing the sampled data from the test split. We use the first 128 instances from the sampled 1077 train split as the calibration set. The random seed is set to 42. The code for processing the original 1078 BOSS benchmark will be placed in our GitHub repository. Concerning the selection of examples 1079 in the few-shot setting, we maintain consistency with BOSS. For datasets lacking examples, we

appropriately select suitable samples from the portion of the train split not chosen as part of the calibration set. For the test data, we use the test split of each dataset as the testing dataset.

1083 C.3.2 CHINESE DOMAIN-SPECIFIC

Experimental Setup. We quantize Baichuan2-7B-Base (Yang et al., 2023) using the GPTQ (Frantar et al., 2022), SpQR (Dettmers et al., 2023), AWQ (Lin et al., 2023), and Smoothquant (Xiao et al., 2023) methods. We quantize the weights to 3-4 bits, and for smoothquant, we further quantize the activations to 8 bits. The quantization is implemented using our custom toolbox, maintaining consistency with the original method in all experimental details. Since the test split of C-EVAL was not publicly available, we upload the test answers to the official platform to obtain the results.

Hyperparameter Selection. For the GPTQ (Frantar et al., 2022) method, we set the group-size parameter to 128 and apply block-sequential as well as layer-sequential quantization. For the SpQR (Dettmers et al., 2023) method, we set the group-size parameter to 128 and apply block-sequential quantization. For the AWQ (Lin et al., 2023) method, we set the group-size parameter to 128. Throughout the quantization process, we use 128 calibration examples. In the few-shot setting, the number of selected examples corresponds to those in C-EVAL (Huang et al., 2024b) and CMMLU (Li et al., 2023), remaining at 5-shot.

Data Selection. For the calibration set, we use 128 calibration examples. For C-EVAL (Huang et al., 2024b), we utilize its validation split as the calibration set. For CMMLU (Li et al., 2023), we randomly select 300 instances from its test split for the train split, subsequently removing the sampled data from the test split. We use the first 128 instances from the sampled train split as the calibration set. The random seed is set to 42. As for the selection of examples in the few-shot setting, we remain consistent with the official standards of C-EVAL and CMMLU. The prompt template is presented in Tab. 10.

- 1104
- 1105 C.4 EXPERIMENT DETAILS OF COMPARISION BETWEEN C4 AND I.I.D/OOD CALIBRATION 1106 SET 1107

Experimental Setup. We quantize LLaMA2-7B (Touvron et al., 2023) using the GPTQ (Frantar et al., 2022) method and test on the EQA task in BOSS. We quantize the weights to 3-4 bits. The quantization is implemented using our custom toolbox, maintaining consistency with the original method in all experimental details.

Hyperparameter Selection. For the GPTQ (Frantar et al., 2022) method, we set the group-size parameter to 128 and apply block-sequential as well as layer-sequential quantization. Throughout the quantization process, we use 128 calibration examples. In the few-shot setting, EQA task is 1-shot. The prompt template is presented in Tab. 10.

Data Selection. Regarding the data selection from C4 (Raffel et al., 2020) as the calibration set, We follow GPTQ (Frantar et al., 2022) and randomly sample 128 samples from C4-en-val as the calibration set with the random seeds of 42 and 567. Regarding the data selection from BOSS Yuan et al. (2024) as the calibration set, we remain consistent with the previous experiments in sec. C.3.1 and set the random seed to 42. For the test data, we use the test split of each dataset as the testing dataset.

1122

1123 C.5 EXPERIMENT DETAILS OF NORMALIZATION

Experimental Setup. To eliminate performance differences between different downstream tasks and situations and better evaluate the quality of different datasets as calibration datasets, we normalize the results based on the same method, generation scenario, bits and different test sets, specifically the 1/4 row in Tab 2. We used Min-max normalization, calculated as follows:

1129 1130

1131

$$x' = \frac{x - X_{min}}{X_{max} - X_{min}} \tag{2}$$

1132 The normalized results range from [0, 1]. This allows for better comparison of the performance of 1133 different calibration sets and enables visualization of the performance differences between calibration sets, rather than just their rankings.

| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. #QA ### Format ### Passage: {{Passage}} / Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{Inpu_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Inpu_1} // Prediction: {{Prediction}} ### Format ### Text: {{Inpu}} // Prediction: {{Prediction}} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Inpu_1} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Inpu_1} // Hypothesis: {{Inpu_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Inpu_1} // Hypothesis: {{Inpu_2} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Inpu_1} // Prediction: {{Prediction}} ### Format ### Premise: {{Inpu_1} // Prediction: {{Prediction}} | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {[input_1]} // Question: {{input_2}} // Answer: ### Format ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. star Text: {[Text]} // Prediction: {{Prediction}} ### Format ### Text: {[Text]} // Prediction: {{Prediction}} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Instruction ### Solve the NLI task. Options for toxicity: benign, toxic. ### Format ### Premise: {{Premise} } // Hypothesis: {{Input_2} // Prediction: ### Format ### Text: {{Input} // Prediction: {{Prediction} } // Hypothesis} // Prediction: ### Format ### Text: {{Input} // Prediction: {{Prediction} } // H## Text: | | |
|--|--|---------------------|---|
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{Input_1}} // Question: {{input_2}} // Answer: {{Answer}}. ### Format ### Passage: {{Input_1}} // Question: {{input_2}} // Answer: {{Answer}}. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Input}} // Prediction: {{Prediction}} ### Format ### Text: {{Input}} // Prediction: {{Prediction: {### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1}} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Prediction: {{Prediction}} ### Noto ### Premise: { | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Instruction ### Passage: {{input_1}} // Prediction: {{Input_2}} // Answer: ### Format ### Passage: {{input_1} // Prediction: {{Prediction}} // ## Format ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction} ### Format ### Solve the NL1 task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Instruction ### Solve the NL1 task. Options for toxicity: benign, toxic. ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2} // Prediction: ### Format ### Premise: {{Premise}} // Prediction: {{Prediction: ### | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Instruction ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### format ### Passage: {{Input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Instruction ### Solve the schument analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Iext}} // Prediction: {{Prediction}} ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Instruction ### Solve the NLI task. Options for toxicity: benign, toxic. ### Instruction ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Prediction: | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {[Passage]} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {[input_1] // Question: {{input_2}} // Answer: {{Answer}}. ### Input ### Passage: {[input_1] // Question: {{input_2}} // Answer: ### Format ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{input} // Prediction: {{Prediction}} ### Input ### Text: {{input} // Prediction: {{Prediction} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Premise: {{Input_1} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Prediction: {{Prediction} ### Format ### Premise: {{Input_1} // Prediction: {{Prediction} ### Format ### Premise: {{Input_2 | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. ### Format ## Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ## Passage: {{Input_1}} // Question: {{Question}} // Answer: ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Text: {{Input}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction: ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction: {Premise: {Premise}} NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with 1 official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Ipassage}} // Question: {{Question}} // Answer: {{Answer}}. ### Instruction ### Passage: {{Ipassage}} // Question: {{input_2}} // Answer: Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Iprut]} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Format ### Premise: {{Iprut_1}} // Hypothesis: {{iprut_2} // Prediction: ### Format ### Premise: {{Iprut_1} // Hypothesis: {{iprut_2} // Prediction: ### Format ### Premise: {{Iprut_1} // Prediction: {{Prediction: ### Format ### Premise: {{Iprut_1} // Prediction: {{Prediction: ### Format ### | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with for official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. #QA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ## Passage: {{Inpu_1}} // Question: {{input_2} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text} // Prediction: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text} // Prediction: ### Format ### Text: {{Text} // Prediction: ### Format ### Premise: {{Premise} // Hypothesis: {{Input_2} // Prediction: ### Format ### Premise: {{Premise} // Hypothesis: {{Input_2} // Prediction: ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: ### Format ### <t< td=""><td>Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with i official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ## Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{input_2}} // Answer: ### Input ### Passage: {{Input_1}} // Prediction: {{Prediction}} ### Input ### Fext: {{Input_1}} // Prediction: ### Input ### Text: {{Text}} // Prediction: ### Input ### Text: {{Input}} // Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: ### Tormat ### Premise: {{Premise}} // Hypothesis: {{input_2}} // Prediction: ### Format ### Text: {{Input}} // Prediction: {{Prediction}} ### Format ### Premise: {{Premise} // Hypothesis: {{input_2} // Pre</td><td></td><td></td></t<> | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with i official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ## Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{input_2}} // Answer: ### Input ### Passage: {{Input_1}} // Prediction: {{Prediction}} ### Input ### Fext: {{Input_1}} // Prediction: ### Input ### Text: {{Text}} // Prediction: ### Input ### Text: {{Input}} // Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: ### Tormat ### Premise: {{Premise}} // Hypothesis: {{input_2}} // Prediction: ### Format ### Text: {{Input}} // Prediction: {{Prediction}} ### Format ### Premise: {{Premise} // Hypothesis: {{input_2} // Pre | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with tofficial template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Text: {{Text}} // Prediction: ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{Premise}} // Hypothesis: {{Input_2} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: ### Format ### Text: {{Tex | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Instruction ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: {{Answer}}. ### Instruction ### Passage: {{Input_1}} // Question: {{Prediction}} Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Instruction ### NLI ### Format ### Premise: {{Premise} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Premise} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Premise} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Premise} // Hypothesis: {{Input_2} // Pre | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. ### Format ### Passage: {Passage} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {[input_1]} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {[input_1]} // Question: {{Input_2}} // Answer: ### Format ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Input ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{[Premise}] // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Format ### Text: {{Text}} // Prediction: {{Predic | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with 1 official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ## Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text} }/ / Prediction: {{Prediction}} ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text} }/ Prediction: {{Prediction}} ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text} }/ Prediction: {{Prediction}} ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text} }/ Prediction: ### Format ### Premise: {{Prediction: ### Instruction ### Solve the toxic detection task. O | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. ### Format ### Passage: {{Input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Instruction ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Instruction ### Solve the estiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Instruction ### Solve the Schipt // Prediction: {{Prediction}} ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Instruction ### Solve the value ### Solve the NL1 task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Format ### Text: {{Text}} // Prediction task. Options for toxicity: benign, toxic. ### Instruction ### Solve the toxic detection task. Options for | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with a official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ## Passage: {{Ipust_1}} // Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Iput_1}} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Premise} // Hypothesis: {{Iput_2}} // Prediction: {{Prediction}} ### Format ### Premise: {{Premise} // Hypothesis: {{Iput_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Premise} // Hypothesis: {{Iput_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Premise} // Hypothes | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b).Task Prompt### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be.### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}.### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer:### Input ### Passage: {{Passage}} // Question: {{Question}} // Answer:### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer:### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral.### Format ### Text: {{Text}} // Prediction: {{Prediction}}### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral.### Format ### Text: {{Text}} // Prediction: {{Prediction}}### Format ### Text: {{Input}} // Prediction:### Format ### Text: {{Input}} // Prediction:### Format ### Text: {{Input}} // Prediction:### Format ### Text: {{Input}} // Hypothesis: {{Input_2}} // Prediction:### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradic-tion.MLI### Format ### Premise: {{Input_1}} // Hypothesis: {{Input_2}} // Prediction:### Format # | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with 1 official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Input_1} // Prediction: {{Prediction}} ### Input ### Text: {{Input_1} // Prediction: {{Prediction}} # | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b).TaskPrompt### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be.EQA### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}.### Input ### Passage: {{Passage}} // Question: {{Question}} // Answer:### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer:### Format ### Text: {{Text}} // Prediction: {{Prediction}}### Format ### Text: {{Input}} // Prediction: {{Prediction}}### Format ### Text: {{Input}} // Prediction:NLI### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction:MI### Format ### Text: {{Input_1}} // Hypothesis: {{Input_2}} // Prediction:TDText: {{Input}} // Prediction: {{Prediction}### Format ### Text: {{Input_1}} // Prediction: {{Prediction:### Format ### Text: {{Input_1}} // Prediction: {{Prediction:### Format ### Text: {{Input_1}} // Prediction: {{Prediction:### Format ### Text: {{Input_1}} // Prediction: {{Prediction}}### Format ### Text: {{Input_1}} // Prediction: {{Prediction}}### Format ### Text: {{Input_1}} // Prediction: {{Prediction}}### Format ### Text: {{Input_1} // Prediction: {{Prediction}}### Format ### Text: {{Input_1} // Prediction: {{Prediction}}### Format ### Text: {{Input_1} // Prediction: <td>Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Input}} // Prediction: {{Prediction}} ### Instruction ### Solve the sentiment analysis task. Options for sentiment, neutral, contradiction. NLI ### Format ### Text: {{Iext}} // Prediction: ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction: ### Format ### Solve the toxic detection task. Options for toxicity: benign, toxic. ###</td> <td></td> <td></td> | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Input}} // Prediction: {{Prediction}} ### Instruction ### Solve the sentiment analysis task. Options for sentiment, neutral, contradiction. NLI ### Format ### Text: {{Iext}} // Prediction: ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction: ### Format ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. ### Format ### Passage] // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Text} // Prediction: {{Prediction}} ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Prediction: {{Prediction}} ### Format ### Format ### | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with 1 official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. #QA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2} // Answer: Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Format ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Format ### Text: {{Text}} // Prediction: {{Prediction}} | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{Input_2}} // Answer: Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### NLI ### Format ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Prediction: {{Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Prediction: {{Prediction} ### Format ### <td>Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Imput} // Prediction: {{Prediction}} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Input_1}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1}} // Prediction: {{Prediction: ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: ### Format ### Premise: {{Input_</td> <td></td> <td></td> | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Imput} // Prediction: {{Prediction}} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Input_1}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1}} // Prediction: {{Prediction: ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: ### Format ### Premise: {{Input_ | | |
| Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with to official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1} // Question: {{Input_2}} // Answer: ### Format ### Passage: {{Input_1} // Prediction: {{Prediction}} // Answer: ### Format ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} // ### Input ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Input ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### < | Table 10: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with 1 official template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. #QA ### Format ### Passage: {Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: {{Prediction}} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Input_1}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1}} // Prediction: {{Prediction: ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: ### Format ### Premise: {{Input_1} // Prediction: {{Prediction}} ### Format ### <th></th> <th></th> | | |
| Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. FQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Question}} // Answer: ### Input ### Passage: {{Input_2}} // Answer: solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input}} // Prediction: {{Prediction}} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Input ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Input ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### <td< th=""><th>Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. #QA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{Input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Input ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Input_1} // Prediction: {{Prediction}} ### Input ###</th><th>Table 1 official</th><th>0: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with t template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b).</th></td<> | Task Prompt ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. #QA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{Input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Input ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Input_1} // Prediction: {{Prediction}} ### Input ### | Table 1 official | 0: Prompts for BOSS and Chinese domain-specific tasks. We maintain consistency with t template provided by BOSS Yuan et al. (2024) and C-EVAL Huang et al. (2024b). |
| Image ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1} // Question: {{Input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction} ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Text: {{Text}} // Prediction: {{Prediction}} | Image ### Instruction ### Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Format ### Passage: {{input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{Input_2}} // Answer: ### Input ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: ### Input ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Solve the toxic detection task. Options for toxicity: benign, | Task | Promnt |
| Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Input ### Passage: {{input_1}} // Prediction: {{input_2}} // Answer: ### Input ### Postage: {{input_1}} // Prediction: {{input_2}} // Answer: ### Input ### Text: {{Text}} // Prediction: {{Prediction} ### Format ### Text: {{Text}} // Prediction: ### Format ### Text: {{Input}} // Prediction: ### Format ### Premise: {{Premise}} // Hypothesis: {{Input_2}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: ### Format ### Premise: {{Input_1} // Prediction: {{Prediction} <td< td=""><td>Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Hypothesis: {{input_2}} // Prediction: ### Input ### Premise: {{Input_1} // Hypothesis: {{input_2} // Prediction: ### Input ### Premise: {{Input_1} // Hypothesis: {{input_2} // Prediction: ### Input ### Premise: {{Input_1} // Prediction: {{Prediction} ### Input ### Text: {{Text} // Prediction: {{Predicti</td><td>LUDI</td><td>### Instruction ###</td></td<> | Solve the extractive question answering task. Refering to the passage below and extract answer for the question. The answer should be the shortest phrase as it can be. EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Hypothesis: {{input_2}} // Prediction: ### Input ### Premise: {{Input_1} // Hypothesis: {{input_2} // Prediction: ### Input ### Premise: {{Input_1} // Hypothesis: {{input_2} // Prediction: ### Input ### Premise: {{Input_1} // Prediction: {{Prediction} ### Input ### Text: {{Text} // Prediction: {{Predicti | LUDI | ### Instruction ### |
| answer for the question. The answer should be the shortest phrase as it can be. ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Input ### Passage: {{Input_1}} // Question: {{Input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ## Text: {{Input}} // Prediction: ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: ### Instruction ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: ### Instruction ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: ### Instruction ### Premise: {{Input_1}} // Prediction: {{Prediction: ### Input ### Text: {{Input_1} // Prediction: {{Prediction}} ### Format ### | answer for the question. The answer should be the shortest phrase as it can be. ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{Iput_1}} // Question: {{Iput_2}} // Answer: ### Input ### Passage: {{Iput_1}} // Question: {{Iput_2}} // Answer: ### Input ### Passage: {{Iput_1}} // Question: {{Iput_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Iput_2}} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Iput_1}} // Hypothesis: {{Input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Instruction ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}}< | | Solve the extractive question answering task. Refering to the passage below and extract |
| EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Input_1}} // Prediction: {{Prediction: ### Instruction ### Premise: {{Input_1}} // Prediction: {{Prediction: ### Instruction ### Premise: {{Input_1}} // Prediction: {{Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: With Prediction: Solve the toxic detection task. Options for toxici | EQA ### Format ### Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input}} // Prediction: CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | | answer for the question. The answer should be the shortest phrase as it can be. |
| Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Instruction ### Text: {{Text}} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Premise: {{Premise} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text} // Prediction: Wif Linput ### <t< td=""><td>Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction} CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。</td><td>EQA</td><td>### Format ###</td></t<> | Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction} CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | EQA | ### Format ### |
| ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ## Text: {{Text}} // Prediction: ### Input ## Text: {{Input}} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1}} // Hypothesis: {{Input_2}} // Prediction: ### Format ### Premise: {{Input_1}} // Hypothesis: {{Input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Text} // Prediction: Up Teal ### | ### Input ### Passage: {{input_1}} // Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Instruction ### Text: {{Text}} // Prediction: ### Instruction ### Text: {{input}} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: ### Format ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: UND UND ### Format ### Text: {{Input_1} // Prediction: Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### <t< td=""><td></td><td>Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}.</td></t<> | | Passage: {{Passage}} // Question: {{Question}} // Answer: {{Answer}}. |
| Passage: {{input_1}}// Question: {{input_2}} // Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Instruction ### Premise: {{Premise} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Premise: {{Input_1} // Prediction: {{Prediction: {{Prediction: {{Prediction: {{Prediction: {{### Format ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Text: {{Input} // Prediction: Kolve the toxic detection: {{Prediction}} ### Format ### Text: {{Text} // Prediction: Kolve the toxic detection: {{Prediction}} ### Format ### Tex | Passage: {{input_1}}// Question: {{input_2}}// Answer: ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: {{Prediction}} ### Format ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: UND UND Detemate ### | | ### Input ### |
| SA ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Input ### Text: {{Text}} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Hypothesis: {{input_2} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: \$UTD UT ### Format ### Text: {{Input} // Prediction: CDS UT | ### Instruction ### Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Input ### Text: {{Input} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction}} ### Input ### Premise: {{Premise} // Hypothesis: {{Input_2} // Prediction: ### Instruction ### Premise: {{Input_1} // Hypothesis: {{input_2} // Prediction: ### Format ### Premise: {{input_1} // Prediction: {{Prediction: ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Text: {{Input} // Prediction: CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | | Passage: {{input_1}} // Question: {{input_2}} // Answer: |
| SA Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Input ### Text: {{input}} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1}} // Hypothesis: {{Input_2} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input}} // Prediction: CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | SA Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1} // Hypothesis: {{input_2} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Format ### Text: {{Text} // Prediction: {{Prediction}} ### Text: {{Text} // Prediction: {{Prediction}} ### Text: {{Text} // Prediction: {{Prediction}} ### Input ### Text: {{Input}} // Prediction: CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | | ### Instruction ### |
| SA ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{Input_1} // Hypothesis: {{Input_2} // Prediction: ### Format ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: {{Prediction}} ### Input ### | SA ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Format ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input}} // Prediction: | | Solve the sentiment analysis task. Options for sentiment: negative, positive, neutral. |
| Text: {{ [iptt] } // Prediction: {{ Prediction} } } ### Input ### Text: {{ [input] } // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise} // Hypothesis: {{Hypothesis} // Prediction: {{Prediction} } ### Input ### Premise: {{Input_1} // Hypothesis: {{ input_2} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{ Text} } // Prediction: {{ Prediction} } ### Input ### Text: {{ Text} } // Prediction: {{ Prediction} } ### Format ### Text: {{ Text} } // Prediction: {{ Prediction} } ### Input ### Text: {{ Input } // Prediction: UND UND | Text: {{ [IExt} } // Prediction: {{ Prediction} } } ### Input ### Text: {{ input} } // Prediction: ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise} // Hypothesis: {{ Hypothesis } // Prediction: {{Prediction} } ### Input ### Premise: {{ input_1 } // Hypothesis: {{ input_2 } // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{ Text } // Prediction: {{ Prediction} } ### Input ### Text: {{ Input } // Prediction: {{ Prediction} } ### Input ### Text: {{ Input } // Prediction: {{ Prediction} } ### Input ### Text: {{ Input } // Prediction: CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | SA | ### Formal ### |
| Image: Image | Image: Image | | ### Input ### |
| ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | ### Instruction ### Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | | Text: {{input}} // Prediction: |
| NLI Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | NLI Solve the NLI task. Options for entailment relationship: entailment, neutral, contradiction. NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | | ### Instruction ### |
| tion. ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input}} // Prediction: {{Prediction}} UNTLE 以下是中国考试的单项选择题,请选出其中的正确答案。 | tion. ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input}} // Prediction: CDS 以下是中国考试的单项选择题, 请选出其中的正确答案。 | | Solve the NLI task. Options for entailment relationship: entailment, neutral, contradic- |
| NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: TD ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | NLI ### Format ### Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | | tion. |
| Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Input} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | NLI | ### Format ### |
| ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | ### Input ### Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | | Premise: {{Premise}} // Hypothesis: {{Hypothesis}} // Prediction: {{Prediction}} |
| Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | | ### Input ### |
| ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | ### Instruction ### Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | | Premise: {{input_1}} // Hypothesis: {{input_2}} // Prediction: |
| TD Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | TD Solve the toxic detection task. Options for toxicity: benign, toxic. ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | | ### Instruction ### |
| TD ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | TD ### Format ### Text: {{Text}} // Prediction: {{Prediction}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | | Solve the toxic detection task. Options for toxicity: benign, toxic. |
| rext. {{Text}} // Frediction: {{Text: {{input}} // Prediction: Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | Itext: {{Itext}} // Prediction: {{Itextended}} ### Input ### Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | TD | ### FOIIIal ### Text: [[Text]] // Prediction: [[Prediction]] |
| Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | Text: {{input}} // Prediction: CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | | |
| CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | CDS 以下是中国考试的单项选择题,请选出其中的正确答案。 | | Text: {{input}} // Prediction: |
| CD5 以于是于国行战的千次起汗感,引起田兴于的正端日来。 | CD3 以于是于国内战的千次起汗感,谓远田兴于的正确日来。 | CDS | 以下是中国老试的单项选择题。 |
| | | CDS | 以下足下国行风的平须起汗越,谓远田兴于时正朔百米。 |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

1188 D MORE EXPERIMENTS

1190 D.1 THE ROBUSTNESS OF DATA SELECTION WITH RESPECT TO RANDOM SEED

In the experiments conducted in the main text, we employ a random seed for the selection of train split and calibration set. In this section, we will alter the random seed to observe the sensitivity of the experiments to the random seed.







Figure 6: S1: evaluation of quantized LLaMA2-7B retested on several standard datasets. Quantization methods include GPTQ. Quantization bits include W4A16, W3A16, and W2A16, with W16A16 used as reference. The left figure shows 5-shot results, while the right figure shows 0-shot results. Different background colors represent different task types. The random seed is 567.

In S1, we randomly sampled 128 samples from c4-en-val as the calibration set and set the random seed to 42. We then modify the random seed to 567 and retest the GPTQ (Frantar et al., 2022) method. The results are presented in Fig. 5 and 6. We observe that the vast majority of datasets exhibited strong robustness to the selection of the calibration set, with performance trends remaining nearly identical across different random seeds. In the comparisons presented in Fig. 5 and Fig. 6, we can still observe that different task types exhibit varying sensitivities to quantization, and this conclusion remains consistent across different random seeds. We see that task types such as scientific knowledge QA,

reading comprehension, common sense reasoning, mathematical reasoning and sentiment analysis
 are highly sensitive to quantization, while task types like natural language inference demonstrate
 robustness against quantization.

1245 In Cross-dataset distribution shift evaluation on BOSS in S2, we randomly sample some examples 1246 from the test split as the train split and use them as the calibration set, setting the random seed to 42. 1247 We modify the random seed to 567 and retest the SA and NLI experiments using GPTQ (Frantar et al., 1248 2022) method. We present the average results with random seeds 42 and 567 in Tab. 11. The results 1249 indicate a certain robustness of the distribution shift experiment on BOSS towards the selection of the 1250 calibration set. For SA task, performance remains consistently better when using Amazon (McAuley 1251 & Leskovec, 2013) as the calibration set across different random seeds, and using SemEval (Nakov 1252 et al., 2019) as the calibration set performs better in most cases. However, the performance has consistently been poor when using DynaSent (Potts et al., 2020) as the calibration set. For NLI task, 1253 performance remains consistently better when using MNLI (Williams et al., 2017) as the calibration 1254 set across different random seeds. 1255

1256

1270 1271 1272

1274

1276

Table 11: Cross-dataset distribution shift evaluation retested on Boss. The result represents the average values obtained with random seeds 42 and 567. "Calib." represents the calibration dataset, and "Gene." represents generalization scenario. To save space, abbreviations are used for datasets. Each row presents experimental results using different datasets as calibration sets on the same test dataset. The higher the metric, the better the performance. The two best performances are denoted in descending order with red and orange respectively. Note: Some datasets could not be used as calibration sets due to insufficient memory resources.

| Meth | d | | | SA | | | | | | | NLI | | | |
|------|------|---------|------|-------|--------|-------|-------|-------|---------|------|-------|-------|--------|----|
| | Test | Gene | W/A | | Cal | ib. | | Test | Cono | W/A | | Ca | ib. | |
| | Test | Gene. | W/A | AZ | DS | SE | SST | 1051 | Gene. | WA | MN | AN | WN | CN |
| | | 0-shot | 4/16 | 65.84 | 46.90 | 66.49 | 53.61 | | 0-shot | 4/16 | 0.25 | 0.31 | 0.25 | - |
| | 47 | 0-31101 | 3/16 | 19.14 | 0.50 | 21.41 | 0.03 | MN | 0-31101 | 3/16 | 0.03 | 0.00 | 0.00 | - |
| | 112 | 3-shot | 4/16 | 78.47 | 70.35 | 81.43 | 80.32 | 14114 | 3-shot | 4/16 | 43.28 | 34.18 | 41.46 | - |
| | | 5 5000 | 3/16 | 80.73 | 41.28 | 70.79 | 70.23 | | 5 51100 | 3/16 | 32.95 | 33.02 | 32.01 | - |
| | | 0-shot | 4/16 | 41.85 | 30.55 | 40.89 | 25.15 | | 0-shot | 4/16 | 0.74 | 0.57 | 0.74 | - |
| | DS | 0 shot | 3/16 | 8.80 | 1.17 | 10.57 | 0.00 | AN | 0 shot | 3/16 | 2.26 | 0.00 | 0.00 | |
| GPT |) | 3-shot | 4/16 | 53.88 | 45.50 | 54.15 | 52.38 | | 3-shot | 4/16 | 34.1 | 33.52 | 33.76 | - |
| 011 | د | 0 shot | 3/16 | 53.86 | 40.25 | 44.26 | 48.91 | | e shot | 3/16 | 32.25 | 33.33 | 34.19 | - |
| | | 0-shot | 4/16 | 19.97 | 14.07 | 22.08 | 14.27 | | 0-shot | 4/16 | 0.09 | 0.08 | 0.10 | - |
| | SE | o shot | 3/16 | 2.48 | 0.10 | 8.25 | 0.02 | WN | o shot | 3/16 | 0.27 | 0.00 | 0.00 | - |
| | ~ | 3-shot | 4/16 | 41.09 | 36.48 | 43.41 | 44.05 | | 3-shot | 4/16 | 42.16 | 42.15 | 39.925 | - |
| | | 5-51101 | 3/16 | 42.69 | 27.98 | 38.57 | 36.48 | | | 3/16 | 43.16 | 43.36 | 46.97 | - |
| | | 0-shot | 4/16 | 44.13 | 33.505 | 37.16 | 25.56 | | 0-shot | 4/16 | 0.03 | 0.50 | 0.00 | - |
| | SST | | 3/16 | 3.93 | 0.52 | 5.09 | 0.00 | CN | | 3/16 | 0.03 | 0.56 | 0.73 | - |
| | | 3-shot | 4/16 | 54.83 | 44.01 | 52.61 | 48.11 | | 3-shot | 4/16 | 35.93 | 36.67 | 32.27 | - |
| | | | 3/16 | 57.17 | 44.33 | 46.68 | 52.29 | | | 3/16 | 28.28 | 20.41 | 26.13 | - |

1278 1279 1280

1281

D.2 SUPPLEMENTARY RESULTS OF S1

1282 In this subsection, we present supplementary results in S1. Fig. 2 primarily illustrates the performance 1283 of each dataset under different quantization situations. Based on this, we categorized these datasets 1284 into 9 types of downstream tasks and calculated the average accuracy decline for each downstream 1285 task type under various quantization situations in Tab. 12. We can clearly observe that different 1286 downstream task types exhibit varying sensitivities to quantization. Scientific knowledge QA, 1287 reading comprehension, common sense reasoning, and mathematical reasoning are highly sensitive to quantization, with low-bit quantization leading to significant performance drops. For example, in 1288 the case of scientific knowledge QA and reading comprehension, the performance decline can reach 1289 up to 40%, while for common sense reasoning and mathematical reasoning, it may drop by around 1290 25%. Notably, although sentiment analysis task also demonstrate high sensitivity to quantization, 1291 low-bit quantization can lead to substantial performance improvements. In contrast, tasks like natural language inference and multi-turn dialogue readoning are less sensitive to quantization, showing 1293 minimal performance variation across different methods and bit-widths, typically not exceeding 10%. 1294

1295 The varying sensitivity to quantization across different task types may be attributed to differences in full-precision performance. For tasks such as Scientific Knowledge QA, which can achieve a

Table 12: The specific percentages of performance degradation after quantization for each task type in S1. Performance degradation that are significantly high or low for all task types are marked in red red and blue blue, respectively

| Gene. | Method& Bits | Scientific knowledge QA | Reading comprehension | Natural language inference | Sentiment analysis | Bias diagnosis mitigation | Syntax phenomena evaluation | Common sense reasoning | Mathematical reasoning | Multi-turn dialogue reasoning |
|--------|--------------|----------------------------|--------------------------|----------------------------------|-----------------------|---------------------------------|-----------------------------------|------------------------------|------------------------|-------------------------------------|
| | GPTQ-4bit | 0.68 | 3.09 | 2.36 | 0.76 | 1.18 | - | 1.93 | 1.44 | -0.16 |
| | GPTQ-3bit | 4.958 | 3.09 | 8.64 | -10.71 | 3.60 | - | -0.01 | 6.98 | 1.35 |
| | GPTQ-2bit | 40.60 | 34.81 | 9.33 | 13.19 | 8.23 | - | 27.57 | 24.22 | 15.71 |
| 0-shot | SPQR-4bit | 0.83 | 0.55 | 2.36 | 4.02 | -0.02 | - | 0.35 | -0.12 | 0.10 |
| | SPQR-3bit | 1.22 | 3.11 | 2.36 | -8.99 | 1.89 | - | -1.26 | 7.11 | 0.51 |
| | SPQR-2bit | 5.25 | 4.18 | 5.56 | 4.21 | 1.05 | - | -0.39 | 4.33 | 1.37 |
| | GPTQ-4bit | 0.87 | -1.98 | 4.54 | -2.21 | 1.15 | 0.00 | -0.67 | 2.81 | -0.17 |
| | GPTQ-3bit | 6.37 | 2.87 | 4.72 | -1.55 | -0.38 | 2.09 | 1.06 | 9.85 | 1.77 |
| | GPTQ-2bit | 39.11 | 36.81 | 12.42 | 2.88 | -0.38 | 35.51 | 33.36 | 25.09 | 15.09 |
| 5-shot | SPQR-4bit | 1.65 | -0.55 | 1.23 | -0.89 | 0.38 | 0.47 | -0.72 | 1.17 | -0.40 |
| | SPQR-3bit | 2.75 | 1.87 | 4.47 | 1.11 | 0.77 | 1.05 | 3.33 | 1.06 | -0.13 |
| | SPQR-2bit | 5.72 | 4.40 | 3.87 | -22.79 | 0.58 | 1.02 | 1.90 | 6.45 | 1.56 |

- maximum accuracy of up to 90%, extreme quantization results in a significant drop in performance, leading to a marked decrease in relative performance. In contrast, for natural language inference tasks with generally lower full-precision performance of around 30% to 40%, the models may not even meet the threshold for effectively solving natural language inference tasks, leaving little room for performance degradation under extreme quantization.

D.3 FULL RESULTS OF CHINESE CROSS-DATASET AND CROSS-SUBJECT TRANSFER TASKS

In this subsection, we present all the results for the Chinese domain-specific tasks in S2 in sec 3.2. The experimental setup is described in detail in C.

D.4 RESULTS FOR MORE MODELS IN S2

In this subsection, we expand the range of models and further validated our conclusions. Tab. 15 presents some experimental results obtained using LLaMA3-8B and LLaMA2-13B (Touvron et al., 2023) on BOSS in S2. We employed the GPTQ method (Frantar et al., 2022) for quantization, keeping the experimental settings consistent with those in Sec. 3.1. We can still observe that the I,I,D results highlighted with background colors have a low overlap with the bolded optimal performance results. This indicates that the same conclusion can be drawn across different series and scales of LLaMA models: the similarity in distribution between calibration data and test datasets does not significantly improve performance.

- D.5 RESULTS FOR FULL PRECISION ON BOSS IN S2

In this subsection, we present the full-precision results for S2 in Tab. 16. As existing quantization methods can achieve performance comparable to full precision in 4-bit setting, we include all full-precision results in the appendix to save space. These full-precision results serve as a baseline for evaluating the impact of quantization on model performance and provide a reference for future research. Overall, although advancements in quantization techniques enable 4-bit models to approach full-precision performance on certain tasks, performance degradation remains a challenge in 3-bit quantization settings.

D.6 RESULTS FOR 2-BIT ON BOSS IN S2

In this subsection, we present part of the results for the 2-bit quantization in Tab. 17 on BOSS in S2. We observe that at 3 bits, the zero-shot performance experiences a significant decline, while

1351Table 13: Cross-dataset distribution shift in Chinese domain specific task. To save space, abbreviations1352are used for datasets. Each row presents the 0-shot and 5-shot experimental results using different1353datasets as calibration sets on the same test dataset. Results with colored backgrounds indicate1354I.I.D results, while those without color represent OOD results. The higher the metric, the better the1355performance. Bold results indicate the best performance on the same test dataset.

| 1300 | Mathad | | | 0. | that | | hot | | | 0.0 | hot | 5. | hot |
|------|--------|--------|-------------|--------------|---------------|---------------|---------------|-------------|--------------|---------------|--------------|---------------|---------------|
| 1356 | Method | | | 0-2 | Ca | | shot | | | 0-5 | Ca | | |
| 1257 | | Test | W/A | СЕ-НМ | СМ-НМ | СЕ-НМ | СМ-НМ | Test | W/A | СЕ-НМ | СМ-НМ | СЕ-НМ | СМ-НМ |
| 1557 | | CE-HM | 4/16 | 39.4 | 37.9 | 53.2 | 52.1 | CM-HM | 4/16 | 50.0 | 50.7 30.6 | 59.1 | 59.1 |
| 1358 | | CL-IIM | 2/16 | 25.1 | 24.4 | 23.9 | 23.4 | CM-IIM | 2/16 | 25.3 | 23.7 | 25.9 | 24.4 |
| 1359 | | Test | W/A | CESS | Ca | lib. | CM SS | Test | W/A | CE SS | Ca | lib. | CM SS |
| 1360 | GPTQ | | 4/16 | 36.9 | 35.4 | 58.8 | 57.5 | | 4/16 | 53.9 | 54.0 | 63.1 | 63.8 |
| 1361 | | CE-SS | 3/16 | 34.6 | 30.3 | 51.9 | 47.5 | CM-SS | 3/16 | 32.8 | 34.3 | 55.4 | 54.6 |
| 1000 | | | 2/16 | 25.1 | 23.9 Ca | 25.9 lib. | 24.7 | | 2/16 | 25.7 | 26.2 Ca | 25.6 lib. | 25.3 |
| 1362 | | Test | W/A | CE-ST | CM-ST | CE-ST | CM-ST | Test | W/A | CE-ST | CM-ST | CE-ST | CM-ST |
| 1363 | | CE-ST | 4/16 | 30.4 28.1 | 26.0 25.7 | 41.8 33.9 | 39.2 35.5 | CM-ST | 4/16 3/16 | 39.3 29.9 | 35.2 25.7 | 43.1 38.6 | 43.8 37.7 |
| 1364 | | | 2/16 | 24.6 | 25.4 | 24.5 | 25.0 | | 2/16 | 26.2 | 25.7 | 24.5 | 25.2 |
| 1365 | | Test | W/A | CE IIM | | lib. | CM IIM | Test | W/A | CE IIM | | lib. | CM HM |
| 1000 | | | 4/16 | 38.5 | 36.3 | 53.8 | 52.5 | | 4/16 | 52.9 | 49.3 | 59.0 | 59.5 |
| 1300 | | CE-HM | 3/16 | 36.0 | 34.6 | 47.9 | 46.6 | CM-HM | 3/16 | 49.5 | 38.1 | 57.1 | 56.9 |
| 1367 | | | 2/16 | 30.1 | 30.9 Ca | 37.4 lib. | 34.5 | 75 4 | 2/16 | 39.3 | 26.0 Ca | 47.5 lib. | 46.3 |
| 1368 | | lest | W/A | CE-SS | CM-SS | CE-SS | CM-SS | Test | W/A | CE-SS | CM-SS | CE-SS | CM-SS |
| 1369 | SpQR | CE-SS | 4/16 | 38.2 39.8 | 38.9 34.7 | 60.0 56.1 | 57.7 53.3 | CM-SS | 4/16 3/16 | 54.8 52.8 | 54.3 51.1 | 63.8 59.4 | 64.7 60.2 |
| 1270 | | | 2/16 | 30.1 | 32.1 | 39.5 | 37.3 | | 2/16 | 38.8 | 39.7 | 44.2 | 47.1 |
| 1370 | | Test | W/A | CE-ST | Ca CM-ST | lib. CE-ST | CM-ST | Test | W/A | CE-ST | Ca CM-ST | lib. | CM-ST |
| 1371 | | | 4/16 | 32.2 | 30.3 | 41.5 | 41.1 | | 4/16 | 40.4 | 39.5 | 43.7 | 43.3 |
| 1372 | | CE-ST | 3/16 | 31.1 | 28.4 | 37.5 | 37.8 | CM-ST | 3/16 | 37.4 | 37.8 | 40.8 | 41.4 |
| 1373 | | Test | 2/10 W/A | 2/10 | Ca | 1ib. | 50.0 | Trat | 2/10 | 51.0 | Ca | lib. | 55.0 |
| 1374 | | Test | W/A | СЕ-НМ | CM-HM | CE-HM | CM-HM | Test | WA | CE-HM | CM-HM | CE-HM | СМ-НМ |
| 1074 | | CE-HM | 3/16 | 36.5 26.7 | 35.6 29.7 | 47.7 41.1 | 49.0 40.8 | CM-HM | 4/16 3/16 | 47.8 | 53.2 50.5 | 58.5 48.0 | 58.2 49.5 |
| 1375 | | | 2/16 | 24.2 | 24.3 | 24.0 | 23.3 | | 2/16 | 25.9 | 42.4 | 25.8 | 23.4 |
| 1376 | | Test | W/A | CE-SS | CM-SS | lib. CE-SS | CM-SS | Test | W/A | CE-SS | CM-SS | lib. CE-SS | CM-SS |
| 1377 | AWQ | | 4/16 | 32.2 | 34.9 | 57.5 | 56.7 | | 4/16 | 51.3 | 52.4 | 62.2 | 61.4 |
| 1378 | | CE-SS | 3/16 | 32.6 | 31.5 | 42.7 | 40.5 | CM-SS | 3/16 | 40.1 | 42.1 | 50.5 24.8 | 50.8 24.7 |
| 1070 | | Test | W/A | 24.0 | Ca | 24.) lib. | 23.7 | Test | W/A | 24.0 | Ca | lib. | 24.7 |
| 1379 | | | 4/16 | CE-ST | CM-ST 29.4 | CE-ST 20.1 | CM-ST | 100 | 4/16 | CE-ST | CM-ST | CE-ST | CM-ST |
| 1380 | | CE-ST | 3/16 | 26.2 | 27.1 | 31.9 | 34.0 | CM-ST | 3/16 | 31.7 | 31.7 | 36.3 | 35.5 |
| 1381 | | | 2/16 | 25.1 | 24.9 | 25.7 | 25.2 | | 2/16 | 24.6 | 24.6 | 24.1 | 24.5 |
| 1382 | | Test | W/A | СЕ-НМ | Ca CM-HM | hb. CE-HM | СМ-НМ | Test | W/A | СЕ-НМ | Ca CM-HM | nb. CE-HM | СМ-НМ |
| 1383 | | | 4/8 | 27.2 | 27.2 | 24.7 | 24.5 | | 4/8 | 31.6 | 29.8 | 29.4 | 27.1 |
| 1004 | | CE-HM | 3/8 | 25.5 27.1 | 25.5 24.2 | 24.9 25.5 | 23.9 24.2 | CM-HM | 3/8 2/8 | 24.7 | 24.8 25.5 | 25.3 24.8 | 23.9 25.3 |
| 1384 | | Test | W/A | | Ca | lib. | | Test | W/A | | Ca | lib. | |
| 1385 | 50 | | 4/8 | CE-SS | CM-SS 26.7 | CE-SS 24.4 | CM-SS 24.5 | | 4/8 | CE-SS 33.1 | CM-SS | CE-SS 28.7 | CM-SS 25.8 |
| 1386 | ~~ | CE-SS | 3/8 | 26.1 | 25.0 | 26.2 | 24.4 | CM-SS | 3/8 | 25.0 | 25.1 | 24.7 | 24.6 |
| 1387 | | | 2/8 | 26.6 | 25.1 | 25.3 | 23.3 | | 2/8 | 24.3 | 25.3 | 25.2 | 25.3 |
| 1007 | | Test | W/A | CE-ST | CM-ST | CE-ST | CM-ST | Test | W/A | CE-ST | CM-ST | CE-ST | CM-ST |
| 1300 | | CE CT | 4/8 | 32.2 | 26.2 | 25.5 | 23.9 | CMET | 4/8 | 28.2 | 27.7 | 26.9 | 43.3 |
| 1389 | | CE-ST | 2/8 | 31.1 27.8 | 27.4 26.8 | 24.8 24.9 | 25.6 26.8 | CM-ST | 3/8 2/8 | 25.4 24.8 | 24.2 24.9 | 24.4 24.6 | 41.4 35.6 |
| 1390 | | | | | | | | | | | | | |

few-shot learning can substantially improve the performance of the quantized model. However, at 2
bits, both zero-shot and few-shot performance face a marked deterioration, with few-shot learning no
longer able to significantly enhance model performance. We believe it is challenging to derive useful
performance insights from results that are nearly zero; therefore, we do not include cases with 2-bit
or lower quantization in the distribution shift experiments.

1405

1406 Table 14: Cross-subject distribution shift in Chinese domain-specific task. To save space, abbrevi-1407 ations are used for datasets. Each row presents the experimental results using different datasets as calibration sets on the same test dataset. Results with colored backgrounds indicate I.I.D results, 1408 while those without color represent OOD results. The higher the metric, the better the performance. 1409 Bold results indicate the best performance on the same test set. 1410

| Moth | Test | Cono | W/A | | Calib | | Test | Gene | W/A | | Calib. | | Test | Cono | W/A | | Calib | • |
|----------|------|---------|------|-------------|-------|--------------|------|---------|------|--------------|--------------|------|------|---------|------|------|-------|------|
| wieun. | Test | Gene. | WA | HM | SS | ST | lest | Gene. | WA | HM | SS | ST | Test | Gene. | WA | HM | SS | ST |
| | | | 4/16 | 39.4 | 36.4 | 37.6 | | | 4/16 | 38.8 | 36.9 | 38.9 | | | 4/16 | 30.4 | 28.4 | 30.4 |
| | | 0-shot | 3/16 | 30.0 | 30.5 | 29.2 | | 0-shot | 3/16 | 29.6 | 34.6 | 30.4 | | 0-shot | 3/16 | 25.9 | 28.3 | 28.1 |
| GPTO | нм | | 2/16 | 25.1 | 24.1 | 26.2 | SS | | 2/16 | 27.3 | 25.1 | 25.2 | ST | | 2/16 | 24.9 | 24.8 | 24.6 |
| c | | | 4/16 | 53.2 | 52.9 | 52.2 | | | 4/16 | 58.9 | 58.8 | 60.1 | ~ - | | 4/16 | 40.9 | 40.4 | 41.8 |
| | | 5-shot | 3/16 | 38.1 | 43.5 | 39.9 | | 5-shot | 3/16 | 42.5 | 51.9 | 48.2 | | 5-shot | 3/16 | 29.7 | 34.1 | 33.9 |
| | | | 2/16 | 23.9 | 26.2 | 23.7 | | | 2/16 | 24.3 | 25.9 | 24.6 | | | 2/16 | 27.3 | 25.1 | 24.5 |
| | | | 4/16 | 38.5 | 38.0 | 40.9 | | | 4/16 | 39.3 | 38.2 | 41.3 | | | 4/16 | 30.3 | 29.9 | 32.2 |
| | | 0-shot | 3/16 | 36.0 | 39.0 | 38.9 | | 0-shot | 3/16 | 34.8 | 39.8 | 39.0 | | 0-shot | 3/16 | 30.5 | 29.1 | 31.1 |
| SpQR | HM | | 2/16 | 30.1 | 29.9 | 29.2 | SS | | 2/16 | 28.7 | 30.1 | 30.6 | ST | | 2/16 | 26.1 | 26.6 | 27.8 |
| | | 5 shot | 4/10 | 55.8 | 51.0 | 52.0 46.5 | | 5 shot | 2/16 | 59.5 | 00.0 56 1 | 52.0 | | 5 shot | 4/10 | 41.4 | 41.0 | 41.5 |
| | | 5-51101 | 2/16 | 37 / | 45.0 | 40.5 37 7 | | 5-51101 | 2/16 | 10.6 | 30.1 | 45.0 | | 5-51101 | 2/16 | 28.3 | 28.0 | 37.5 |
| | | | 4/10 | 265 | 24.2 | 22.4 | | | 4/10 | 25.0 | 22.2 | 21.4 | | | 4/10 | 20.5 | 20.0 | 26.6 |
| | | 0 shot | 4/10 | 30.5 | 34.2 | 33.4 27.5 | | 0 shot | 2/16 | 35.2 28 2 | 32.2 | 31.4 | | 0 shot | 4/10 | 28.5 | 20.5 | 20.0 |
| | | 0-51101 | 2/16 | 20.7 | 24.2 | 21.5 | | 0-51101 | 2/16 | 20.5 | 32.0 24.8 | 20.2 | | 0-51101 | 2/16 | 21.1 | 20.9 | 20.2 |
| AWQ | HM | | 4/16 | 47.7 | 49 7 | 51.2 | SS | | 4/16 | 53.4 | 24.0 57.5 | 56.6 | ST | | 4/16 | 37.7 | 38.5 | 39.1 |
| | | 5-shot | 3/16 | 41.1 | 38.4 | 37.4 | | 5-shot | 3/16 | 44.0 | 42.7 | 38.7 | | 5-shot | 3/16 | 31.9 | 31.0 | 31.9 |
| | | | 2/16 | 24.0 | 24.6 | 23.8 | | | 2/16 | 23.9 | 24.9 | 25.1 | | | 2/16 | 25.2 | 25.3 | 25.7 |
| | | | 4/8 | 27.2 | 28.9 | 27.4 | | | 4/8 | 28.3 | 27.4 | 28.2 | | | 4/8 | 26.8 | 28.0 | 25.4 |
| | | 0-shot | 3/8 | 25.5 | 23.9 | 26.4 | | 0-shot | 3/8 | 26.4 | 26.1 | 25.5 | | 0-shot | 3/8 | 26.6 | 25.2 | 26.7 |
| 60 | | | 2/8 | 27.1 | 25.2 | 24.8 | | | 2/8 | 26.2 | 26.6 | 26.4 | CT | | 2/8 | 26.4 | 26.4 | 25.7 |
| SQ | HM | | 4/8 | 24.7 | 24.2 | 24.9 | 55 | | 4/8 | 26.0 | 24.4 | 24.3 | 51 | | 4/8 | 24.8 | 24.3 | 25.5 |
| | | 5-shot | 3/8 | 24.9 | 26.4 | 26.2 | | 5-shot | 3/8 | 24.7 | 26.2 | 25.9 | | 5-shot | 3/8 | 26.8 | 25.3 | 24.8 |
| | | | 2/8 | 25.5 | 26.4 | 24.2 | | | 2/8 | 26.5 | 25.3 | 24.9 | | | 2/8 | 26.6 | 26.8 | 24.9 |

1435

1437 Table 15: Cross-dataset distribution shift evaluation on BOSS in S2. The models used for these 1438 experiments are LLaMA3-8B and LLaMA2-13B. "Calib." represents the calibration dataset, and 1439 "Gene." represents generalization scenario. To save space, abbreviations are used for datasets. Each 1440 row presents experimental results using different datasets as calibration sets on the same test dataset. 1441 Results with colored backgrounds indicate I.I.D results, while those without color represent OOD 1442 results. The higher the metric, the better the performance. Bold results indicate the best performance on the same test dataset. 1443

| 1444 | | | | | | | | | | | | | | | | |
|-------|-------------|------|---------|-------|-------|-------|-------|-------|---------------|--------|---------|-------|-------|-------|-------|-------|
| 1///5 | Model | Tost | Gene | W/A | | Ca | lib. | | Model | Test | Gene | W/A | | Ca | lib. | |
| 1440 | wiodei | itst | oche. | •••/A | AZ | DS | SE | SST | Model | nət | Gene. | m/A | AZ | DS | SE | SST |
| 1446 | | | 0 shot | 4/16 | 85.15 | 88.06 | 81.24 | 86.77 | | | 0 shot | 4/16 | 88.74 | 79.73 | 87.59 | 88.69 |
| 1447 | | 47 | 0-51101 | 3/16 | 0.00 | 0.00 | 0.00 | 0.00 | | 47. | 0-shot | 3/16 | 0.01 | 81.49 | 81.64 | 85.76 |
| 1448 | | | 3-shot | 4/16 | 85.67 | 85.73 | 85.74 | 85.81 | | | 3-shot | 4/16 | 81.64 | 81.05 | 84.44 | 82.05 |
| 1440 | | | | 3/16 | 0.00 | 0.00 | 0.00 | 0.00 | | | | 3/16 | 44.74 | 86.57 | 86.07 | 85.71 |
| 1449 | | | 0-shot | 4/16 | 56.94 | 58.76 | 54.88 | 59.03 | | | 0-shot | 4/16 | 62.79 | 40.52 | 57.99 | 62.76 |
| 1450 | | DS | 0 51100 | 3/16 | 0.00 | 0.00 | 1.57 | 0.00 | | DS | 0 31100 | 3/16 | 0.00 | 61.04 | 48.03 | 53.73 |
| 1451 | | 0.5 | 3-shot | 4/16 | 54.78 | 56.84 | 56.09 | 60.55 | | 05 | 3-shot | 4/16 | 47.14 | 24.35 | 53.46 | 48.51 |
| 1450 | LL 9MA 3-8R | | 5 5000 | 3/16 | 0.00 | 0.00 | 0.00 | 1.34 | LLaMA2-13R | | 5 51100 | 3/16 | 49.73 | 64.38 | 63.43 | 65.25 |
| 1432 | ELawing-ob | | 0-shot | 4/16 | 43.35 | 45.80 | 37.08 | 45.50 | ELuiville 15D | | 0-shot | 4/16 | 42.29 | 21.30 | 37.27 | 44.55 |
| 1453 | | SF | 0-31101 | 3/16 | 0.00 | 0.00 | 0.32 | 0.00 | | SF | 0-31101 | 3/16 | 0.00 | 38.03 | 36.01 | 44.81 |
| 1454 | | 512 | 3-shot | 4/16 | 51.15 | 49.74 | 45.19 | 51.22 | | 5L | 3-shot | 4/16 | 54.27 | 42.60 | 50.56 | 55.44 |
| 1/55 | | | 5-3100 | 3/16 | 0.00 | 0.00 | 0.00 | 0.00 | | 5-8100 | 3/16 | 34.20 | 52.04 | 56.30 | 48.11 | |
| 1433 | | | 0 shot | 4/16 | 61.02 | 62.84 | 55.28 | 61.41 | | | 0 shot | 4/16 | 62.32 | 18.64 | 43.55 | 64.41 |
| 1456 | | SST | 0-51101 | 3/16 | 0.00 | 0.00 | 1.56 | 0.00 | | SST | 0-51101 | 3/16 | 0.00 | 64.02 | 50.33 | 63.62 |
| 1457 | | 551 | 3_shot | 4/16 | 61.54 | 60.37 | 61.80 | 67.67 | | 551 | 3-shot | 4/16 | 32.07 | 17.47 | 33.51 | 31.03 |
| | | | 5-500 | 3/16 | 0.00 | 0.00 | 0.00 | 0.26 | | | 5-51101 | 3/16 | 26.86 | 64.15 | 59.19 | 53.85 |

 1468
 Table 16: Full precision results on BOSS in S2. "Gene." represents generalization scenario. To save space, abbreviations are used for datasets.

 1470
 FOA

| | | EQA | | | | | | SA | | | |
|----------|-------|-------|-------|-------|-------|----------|-------|-------|-------|-------|-------|
| Gene. | W/A | SQ | AQA | NQA | SQA | Gene. | W/A | AZ | DS | SE | SST |
| 0-shot | 16/16 | 54.00 | 28.68 | 39.54 | 46.34 | 0-shot | 16/16 | 74.75 | 50.40 | 27.64 | 45.76 |
| Few-shot | 16/16 | 67.93 | 37.21 | 49.40 | 62.41 | Few-shot | 16/16 | 79.90 | 53.13 | 43.72 | 55.54 |
| | | NLI | | | | | | TD | | | |
| Gene. | W/A | MN | AN | MN | CN | Gene. | W/A | CC | AC | IH | TG |
| 0-shot | 16/16 | 0.47 | 1.55 | 0.14 | 0.06 | 0-shot | 16/16 | 62.25 | 17.60 | 48.94 | 59.69 |
| Few-shot | 16/16 | 44.81 | 33.72 | 43.32 | 34.62 | Few-shot | 16/16 | 91.25 | 16.44 | 63.76 | 73.59 |

| 1 | 489 | |
|---|-----|--|
| 1 | 490 | |
| 1 | 491 | |

Table 17: 2-bit results on BOSS in S2. We present some results of Sentiment Analysis. "Gene." represents generalization scenario. To save space, abbreviations are used for datasets.

| Calib. | Gene. | SQUAD | AdvQA | NeswQA | SearchQA |
|----------|--------|-------|-------|--------|----------|
| AdvOA | 0-shot | 0.46 | 0.22 | 0.39 | 0.07 |
| AuvQA | 1-shot | 0.52 | 0.22 | 0.20 | 0.06 |
| SaarahOA | 0-shot | 1.35 | 0.76 | 0.70 | 0.90 |
| SearchQA | 1-shot | 1.86 | 1.02 | 1.18 | 1.20 |