# DIVIDE, REWEIGHT, AND CONQUER: A LOGIT ARITHMETIC APPROACH FOR IN-CONTEXT LEARNING

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027 028 029 Paper under double-blind review

#### Abstract

In-Context Learning (ICL) emerges as a key feature for Large Language Models (LLMs), allowing them to adapt to new tasks by leveraging task-specific examples without updating model parameters. However, ICL faces challenges with increasing numbers of examples due to performance degradation and quadratic computational costs. In this paper, we propose Logit Arithmetic Reweighting Approach (LARA), a novel framework that enhances ICL by using logit-based ensembling of multiple demonstrations. Our approach divides long input demonstrations into parallelizable shorter inputs to significantly reduce memory requirements, and then effectively aggregate the information by reweighting logits of each group via a non-gradient optimization approach. We further introduce Binary LARA (B-LARA), a variant that constrains weights to binary values to simplify the search space and reduces memory usage by filtering out less informative demonstration groups. Experiments on BBH and MMLU demonstrate that LARA and B-LARA outperform all baseline methods in both accuracy and memory efficiency. We also conduct extensive analysis to show that LARA generalizes well to scenarios of varying numbers of examples from limited to many-shot demonstrations. Our codes can be found in https://anonymous.4open.science/ r/LARA-F55B.

### 1 INTRODUCTION

031 In-Context Learning (ICL) (Brown et al., 2020) is one of the emergent abilities of Large Language 032 Models (LLMs) as they are scaled to billions of parameters (Wei et al., 2022). ICL enables LLMs to 033 adapt to new tasks by utilizing task-specific examples within the input context (Dong et al., 2023), 034 and does not require any updates to or access to model parameters. While ICL has achieved impressive performance across various domains, it encounters significant challenges when dealing with an increasing number of examples. Longer context window size often leads to performance degrada-037 tion (Xiong et al., 2023). This is due to the low density of useful information within longer prompts, 038 and the reduced sensitivity to positional information, both of which diminish the capability of the model to effectively capture and utilize key content. Additionally, the quadratic growth of computational cost with the input length makes it particularly expensive for large-scale models. 040

041 Previous works primarily focus on two directions to address these challenges. The first direction is 042 input compression, which aims to shorten the input length (Jiang et al., 2023b; Pan et al., 2024; Xu 043 et al., 2023a; Wingate et al., 2022) or selectively retrieve relevant portions of demonstrations to be 044 included in the prompt (an Luo et al., 2024). However, these methods risk losing critical information, which may negatively impact model performance. The second direction involves aggregating hidden states within LLMs to simulate the effect of in-context demonstrations (Hao et al., 2022; Li et al., 046 2023b; Hendel et al., 2023). These methods, however, are not applicable to closed-source models 047 like GPT-4, as they require direct access to the model internal weights. Additionally, they contradict 048 the core advantage of in-context learning, which is the ability to operate without modifications to hidden states or model parameters. 050

In this study, we propose a novel framework, Logit Arithmetic Reweighting Approach (LARA),
 which aims to combine the strengths of both input compression and hidden state approaches. Our
 method first divides demonstrations into subgroups to allow LLMs to focus on shorter inputs and reduce computational requirements. We then design a weighted sum aggregation approach to combine



Figure 1: Illustration of the differences between few-shot in-context learning and LARA (ours) during inference. Unlike few-shot in-context learning, which concatenates all demonstrations as a prefix to the input, our method splits the in-context examples into different groups. The next token is then generated based on a weighted average of logits, with weights precomputed using the framework described in Sec. 3.3.

067 068 069

the output logits from the language model given each subgroup of examples. This ensures that the relevant information from each subgroup could potentially be captured by the language model. One 071 key innovation in LARA is that we use a non-gradient approach to optimize the weights of logits for 072 each subgroup. We employ the Covariance Matrix Adaptive Evolution Strategy (CMA-ES) (Hansen 073 & Ostermeier, 1996) to efficiently explore the weight vector space via resampling based on best-074 performing candidates. This allows us to optimize the contribution of each subgroup without any 075 gradient updates. We further develop Binary-LARA (B-LARA) by constraining the weight values 076 to  $\{0,1\}$ , which can be interpreted as a process of subgroup selection. This not only reduces the 077 computational cost but more importantly, leads to better performance due to the simplified search 078 space for the binary weight vector.

Our experiments on BBH and MMLU benchmarks show that both LARA and B-LARA consistently outperform direct in-context learning and simple retrieval-based demonstration selection across various models, with the additional benefit of lower GPU memory usage. Further analysis reveals that the method excels in both low-resource scenarios with few examples and settings with abundant demonstrations, consistently delivering superior performance. Moreover, our ablation study highlights the critical role of the reweighting steps, although even logit averaging alone outperforms standard in-context learning.

- To summarize, our main contributions are as follows:
  - To the best of our knowledge, we are the first to propose ensembling information through logit arithmetic from different ICL demonstrations. We introduce LARA, a non-gradient optimization framework that reweights the information of different demonstration groups to improve ICL performance.
  - We conduct extensive experiments on Llama3.1-8B (Dubey et al., 2024), Mistral-7B (Jiang et al., 2023a), and Gemma-7B (Mesnard et al., 2024) on BBH Srivastava et al. (2022) and MMLU Hendrycks et al. (2021), and show that LARA outperforms all baseline methods across all three models.
  - Our comprehensive analysis reveals the broad applicability and efficiency of LARA and B-LARA. We demonstrate that our methods consistently outperform baselines across a wide range of example quantities, from fewer than 5 to more than 200. We also demonstrate the applicability of our methods to black-box LLMs.
- 099 100 101 102

103

090

092

093

095

096

098

2 PRELIMARIES

**In-Context Learning.** Traditional In-Context Learning leverages N labeled examples in the input prompt, represented as  $\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N$  to provide hints for language model generation. Each pair  $(x_i, y_i)$  is converted into a semantically meaningful demonstration  $d_i = \tau(x_i, y_i)$  using a predefined template  $\tau$ . These demonstrations are then concatenated to form a comprehensive context  $\mathcal{C} = d_1 \oplus d_2 \oplus \cdots \oplus d_N$ , with appropriate separators (*e.g.*, newlines or special tokens) between



131 Figure 2: Illustration of the LARA framework. The input demonstration set  $\mathcal{D}_{train}$  is divided into subsets  $S_1, S_2, \ldots, S_k$ , which are further split into two groups: one for candidate examples and the 132 other for validation examples. For each token, logits are generated using Logit-Arithmetic Decoding, 133 which aggregates the output logits from all subsets. After generating all tokens, the cross-entropy 134 loss is computed based on the weighted-average logits and the ground truth from the validation 135 subset. The subset weights are then resampled and adjusted to minimize the loss. This process of 136 token generation, loss calculation, and weight resampling is repeated iteratively. After optimizing 137 the weights for the first group of candidate examples, the roles of the candidate and validation 138 examples are swapped. 139

143

150 151

152 153

154

156 157

each demonstration. For each test input  $x_{test}$ , the language model receives the concatenated prompt  $C \oplus x_{test}$  to generate a response.

**Logit-based Generation.** We consider decoding approaches for language generation, where the language model receives an input prompt  $C \oplus x_{test}$  and produces coherent and logical responses. The term "logit" refers to the raw, unnormalized scores output by the model before they are converted into probabilities by a softmax function. These logits are generated by passing the input sequence through the LLM. Formally, given the logit z, the probability of the next token  $x_t$  given the previous tokens  $x_{1:t-1}$  is computed using the softmax function:

$$P(x_t \mid x_{1:t-1}) = \frac{\exp(\boldsymbol{z}_{x_t})}{\sum_{x' \in V} \exp(\boldsymbol{z}_{x'})}$$
(1)

where  $z_{x_t}$  is the logit corresponding to the token  $x_t$ , and V is the vocabulary set.

#### 3 METHODOLOGY

In this section, we provide an overview of LARA. Figure 2 illustrates the overall framework of our approach. Unlike directly concatenating  $\mathcal{D}_{\text{train}}$  into a single sequence, we first divide the *N* examples into subgroups, which are used as inputs to the LLM. The output logits from these subgroups are then aggregated, and we assign weights to each subgroup using a non-gradient search algorithm. During inference, the precomputed weights are used to combine the logits from each group.

In Sec. 3.1, we explain the partition strategy to divide examples into subgroups. Then we introduce how the outputs are aggregated across different subgroups in Sec. 3.2, and the reweighting strategy for optimal combination in Sec. 3.3. Furthermore, we show in Sec. 3.4 that imposing a hard constraint for our reweighting strategy could further reduce memory usage and computational resources. Finally, we discuss in Sec. 3.5 the inference efficiency brought by our proposed approach.

3.1 PARTITION STRATEGY

168 169

167

170

Given *N*-shot in-context examples, we first split  $\mathcal{D}_{\text{train}}$  into *k* disjoint subsets each containing *L* incontext examples, such that  $\mathcal{D}_{\text{train}} = S_1 \cup S_2 \cup \ldots \cup S_k$  with  $|S_i| = L$  for all  $i \in \{1, \ldots, k\}$ . When inputting a subgroup  $S_i$  to an LLM, we concatenate all of its elements to get  $C_i = d_{(i-1)L+1} \oplus d_{(i-1)L+2} \oplus \cdots \oplus d_{iL}$ , and the complete input for the *i*-th subgroup to LLM is  $\mathcal{C}_i \oplus \boldsymbol{x}_{\text{test}}$ . We assume that *N* is divisible by *k* in our experiments, so that L = N/k. In practice, in cases where *N* is not divisible by *k*, we could truncate the last subset and only retain L(k-1) examples.

177 178

179

#### 3.2 LOGIT-ARITHMETIC DECODING

Previous studies (Li et al., 2022; Liu et al., 2024; Dekoninck et al., 2023) have utilized logit offsets to control the outputs of large language models for better generation quality or instruction following. Inspired by these work, we propose a novel method that combines information from multiple incontext demonstrations through logit-arithmetic decoding. Specifically, our approach focuses on aggregating the logits produced by the language model outputs for various contextual inputs. With the input query  $\boldsymbol{x}_{\text{test}}$  and the example subset being  $S_i$ , we can compute the logit outputs of the language model, denoted as  $f_{\theta}(S_i, \boldsymbol{x}_{\text{test}}) = \log p(y | S_i, \boldsymbol{x}_{\text{test}})$ . We then combine these logits using a weighted sum to get the generation probability over the output token:

188

189

190 191

$$p(y \mid \boldsymbol{x}_{\text{test}}, \boldsymbol{w}) = \text{softmax}\left(\sum_{i=1}^{k} w_i \cdot f_{\theta}(\mathcal{S}_i, \boldsymbol{x}_{\text{test}})\right)$$
(2)

where k is the number of example subsets, and  $w_i$  are weights that indicate the importance of the contribution of each subset, with  $\sum_{i=1}^{k} w_i = 1$ . As a baseline approach, we could set uniform weighting, where  $w_i = 1/k$ . However, this may not be optimal for all tasks, as the quality and relevance of different subgroups may vary. In the following section, we introduce a reweighting strategy to optimize these weights to enhance model performance.

3.3 Reweighting Logits by Non-Gradient Optimization

To further enhance the model performance, we employ non-gradient optimization methods to optimize the weights  $w_i$  based on the loss calculated from  $p(y \mid \boldsymbol{x}_{val})$ . Given the combined probability  $p(y \mid \boldsymbol{x}_{val})$ , our objective is to minimize a cross-entropy loss function  $\mathcal{L}(\boldsymbol{w})$  over the predicted probabilities and the ground truth. Specifically, we utilize the following cross-entropy loss function for the generation model:

205

207 208 209

197 198

199

$$\mathcal{L}(oldsymbol{w}) = -\sum_{(oldsymbol{x}_{ ext{val}},oldsymbol{y}_{ ext{val}}) \in \mathcal{D}} \sum_{t=1}^T \log p(y_t \mid oldsymbol{x}_{ ext{val}},oldsymbol{w})$$

where D represents the validation dataset, T is the length of the sequence,  $y_t$  is the true word at time step t,  $x_{val}$  is the input sequence, w denotes the weight vector, and  $p(y_t | x_{val}, w)$  represents the predicted probability of the true word  $y_t$  at time step t, given the input sequence  $x_{val}$  and the weight vector w.

To avoid introducing additional labeled data, we employ a cross-validation strategy. We partition the demonstration set S into two subsets:  $S_A = S_1 \cup S_2 \cup \ldots \cup S_{\lfloor k/2 \rfloor}$  and  $S_B = S_{\lfloor k/2 \rfloor + 1} \cup S_{\lfloor k/2 \rfloor + 2} \cup \ldots \cup S_k$ . When optimizing weights for  $S_i \in S_A$ , we use  $S_B$  as the validation set, and vice versa. We choose non-gradient optimization methods over gradient-based alternatives due to two key factors: (1) The loss function  $\mathcal{L}(w)$  is non-differentiable, since updating the weight vector w affects the logits of subsequent tokens, leading to possibly different decoding results of subsequent tokens. (2) The dimensionality of the weight vector w is relatively low, specifically equalled to the number of groups k.

In our empirical experiments, we refer to Liu et al. (2020) and employ the Covariance Matrix Adaptive Evolution Strategy (CMA-ES) (Hansen & Ostermeier, 1996). CMA-ES is a stochastic, derivative-free optimization algorithm. During each iteration, CMA-ES samples a set of candidates in the space of the weight vector w from a multivariate normal distribution, evaluates  $\mathcal{L}(w)$  for each candidate, and then updates the mean and covariance matrix of the distribution based on the best-performing candidates. This allows for an efficient exploration over the weight space.

227 228

229

3.4 BINARY CONSTRAINTS FOR LARA

We further propose a variant of LARA, named as B-LARA, by imposing a hard constraint on the weight vector w to binary values  $\{0, 1\}$ . This binary constraint offers two key advantages: first, it simplifies the search space and potentially leads to faster convergence; second, it allows for direct elimination of demonstration groups with zero weight, thereby improving inference efficiency. Intuitively, the binary optimization of w can be seen as a form of subset selection to identify the most relevant demonstrations in  $\mathcal{D}_{train}$  benefitting model performance on specific tasks.

To solve this binary optimization problem, we employ the simplest evolution strategy (1+1)-ES (Rechenberg, 1973). It involves a simple cycle: a single parent produces one offspring per generation through mutation—adding a small, random change. If this offspring performs as well or better than the parent based on a predefined fitness criterion, it becomes the new parent for the next generation. Otherwise, the original parent remains. The overall sampling procedure is shown in Algorithm 1.

The simplicity of this method in repeated mutation and selection makes it particularly suitable for our binary optimization scenario.

- 245 3.5 COMPUTATIONAL COMPLEXITY
- 246

244

5.5 COMI O TATIONAL COMI LEAT

246

We analyze the computational complexity of LARA and B-LARA compared to standard ICL. During inference, the self-attention mechanism in Transformer models is the primary bottleneck for GPU memory requirement, with the memory complexity being  $O(n^2)$ , where *n* is the input sequence length. This quadratic scaling is due to the pairwise interactions between tokens in the attention matrix.

By splitting the input sequence into k groups, each of length around  $\frac{n}{k}$ , LARA and B-LARA can leverage parallel computing resources more effectively. The complexity for LARA becomes  $O(\frac{n^2}{k} * k) = O(\frac{n^2}{k})$ . B-LARA further reduces computational complexity by selecting only a subset of groups. If m out of k subgroups are assigned non-zero weights, then the complexity of B-LARA becomes  $O(\frac{mn^2}{k^2})$ . We show the empirical GPU memory usage in Sec. 5.5.

4 Experime

259 260 261

258

### EXPERIMENTS

In this section, we provide details of our main experiments. We first give an overview of the experimental setup and implementation details in Sec. 4.1, and then present our findings along with the results in Sec. 4.2.

4.1 EXPERIMENTAL SETUP

- 265 266
- 267 268

**Datasets and Evaluation.** We evaluate our methods using two well-established benchmarks: Big-Bench Hard (BBH) (Srivastava et al., 2022) and Massive Multitask Language Understanding

(MMLU) (Hendrycks et al., 2021). BBH tests models on challenging reasoning tasks across domains including arithmetic reasoning, commonsense reasoning, and linguistics. MMLU measures
generalization across 57 diverse subjects, covering both humanities and STEM fields, offering a
comprehensive evaluation of knowledge and problem-solving abilities of LLMs. For both benchmarks, we use exact match (EM) as our evaluation criterion, which requires model predictions to
perfectly match the correct answers. We report the accuracy scores in our experiment results. The
details about dataset analysis and prompts can be found in Appendix A.

277 Models. Our proposed LARA for in-context learning is applicable to any LLM. To demonstrate its 278 generality, we evaluate it on three open-source, decoder-only models: Llama3.1-8B (Dubey et al., 279 2024), Mistral-7B (Jiang et al., 2023a), and Gemma-7B (Mesnard et al., 2024). Llama-3.1-8B is 280 known for strong performance across various NLP tasks, Mistral-7B is optimized for efficiency and 281 is balanced between computational cost and accuracy. Gemma-7B focuses on advanced reasoning 282 and language comprehension. These models represent diverse architectures and training strategies, 283 allowing us to test the adaptability of our methods. By using open-source models in evaluation, 284 we ensure the reproducibility of our proposed method and validate its broad applicability across 285 state-of-the-art model architectures. 286

**Hyperparameter Setting.** In our main experiment, we use  $\mathcal{D}_{\text{train}}$  consisting of N = 32 in-context examples for our methods. For each task,  $\mathcal{D}_{\text{train}}$  is split into subsets of size  $L \in \{2, 4, 8\}$ , and for each L we perform up to J = 20 iterations for weight optimization. We compare the minimum validation loss across different settings of L to determine the optimal configuration, including L and corresponding w, for the final inference phase. The baseline methods also use the same  $\mathcal{D}_{\text{train}}$  as input. For our method and all baselines, we set the temperature to 0 to enforce greedy decoding. Our experiments are conducted on a single A100 80GB GPU.

**Compared Methods.** We introduce several primary baseline methods: Direct In-Context Learn-295 ing (ICL), KNN-Augmented In-ConText Example Selection Liu et al. (2022) (KATE), Rationale-296 Augmented Ensembles (RAE) (Wang et al., 2022) and In-context Vector (ICV) (Liu et al., 2023) and 297 StructICL (Hao et al., 2022) as the representative of parameter access methods. We use the same 32 298 in-context examples as inputs to all baseline methods as our proposed method. For Direct ICL, all 299 32 examples are concatenated with the prompt. For KATE, we apply the Top-K selection from Liu et al. (2022) that uses a smaller model<sup>1</sup> to retrieve the most similar input-output pairs from  $\mathcal{D}_{train}$ 300 as in-context demonstrations. We evaluate KATE with 2, 4, and 8 demonstrations as baselines. For 301 RAE, we divide the examples into different groups and use each group as in-context examples to 302 generate separate results. The final output is determined by applying majority voting across these 303 individual group-based results. For StructICL, we also present the results with varying numbers of 304 groups: 2, 4, and 8. In ICV, we follow the original paper to set  $\lambda = 0.1$  and average the ICV given 305 by all 32 examples. We report results with group sizes of 2, 4, and 8 to ensure the same memory 306 usage as our method. 307

- 4.2 MAIN RESULTS
- 308 309 310
- 311

Results from Table 1 demonstrate the effectiveness of our proposed methods, LARA, and B-LARA, 312 across BBH and MMLU benchmarks. B-LARA consistently outperforms most of baseline methods 313 across three model architectures. Notably, B-LARA achieves the highest accuracy and improves 314 over direct ICL by 2.05, 5.67, and 2.12 points on BBH dataset across three models respectively. 315 Moreover, our methods and can consistently outperform retrieval or simple ensemble baselines like 316 KATE and RAE, indicating that our method is more effective in combining information from mul-317 tiple demonstration subgroups. Compared to the ICV and StructICL baseline, which has the advan-318 tage of access to model parameters, our methods still achieve better performance without access to the hidden state, which further demonstrates the efficacy of our methods in aggregating information 319 without direct access to model internal parameters. 320

An interesting finding is that B-LARA performs better than LARA despite a more constrained search space for the weight vector. We believe this is because we only use 20 iterations for weight opti-

<sup>&</sup>lt;sup>1</sup>https://huggingface.co/sentence-transformers/all-distilroberta-v1

324 Table 1: Accuracy of all methods on BBH and MMLU. The results shown are the average perfor-325 mance across datasets within each benchmark. Please refer to appendix D.2 for breakdown results 326 of each dataset. The subscript of KATE indicates the number of selected ICL demonstrations as input to LLMs. 327

		$BBH_{average}$		Μ	$MLU_{average}$		
	Llama3.1-8B Gemma-7B			Llama3.1-8B	Gemma-7B	Mistral-7B	
Black-Box Meth	nod:						
ICL	45.64	37.08	42.91	65.63	61.44	62.84	
$KATE_2$	43.60	37.07	43.16	66.62	56.28	63.99	
$KATE_4$	44.03	38.83	43.16	66.75	55.78	63.48	
KATE <sub>8</sub>	44.47	37.03	42.96	67.19	54.13	63.93	
$RAE_2$	44.59	40.24	43.95	66.88	65.18	62.99	
$RAE_4$	45.23	40.44	44.49	66.40	65.01	62.99	
RAE <sub>8</sub>	44.06	39.85	44.07	67.09	64.80	63.61	
LARA (ours)	47.46	41.77	44.77	66.54	64.36	63.93	
B-LARA (ours)	47.69	42.75	45.03	67.80	65.56	64.12	
White-Box Meth	nod:						
ICV	45.93	42.16	44.50	66.97	64.99	64.02	
StructICL <sub>2</sub>	46.64	39.54	44.68	66.78	64.34	63.52	
StructICL <sub>4</sub>	46.98	40.53	44.89	66.97	64.46	63.99	
StructICL <sub>8</sub>	46.57	41.46	43.99	66.56	65.16	63.46	

mization, and the binary constraint brings more benefits by introducing a simplified optimization landscape and providing a regularization effect to prevent overfitting.

#### 5 ANALYSIS

In this section, we present a deep analysis of our proposed method LARA under various conditions.

#### 5.1 HOW DOES LARA EXTEND TO **GENERATION TASKS?**

In previous experiments, we mainly focus on the 356 classification or single token generation tasks. Here 357 we extended our experiments to generation tasks 358 like GSM8K (Cobbe et al., 2021) for math rea-359 soning. We follow the experiment setting used 360 in FocusICL (Yuan et al., 2024) and evaluate our 361 method against Llama-3-8B-Instruct, LongChat-7B-362 V1.5-32K (Li et al., 2023a) and Vicuna-7B-V1.5-

Table 2:	Accuracy	across	different r	nodels on
Gsm8k.				

	Llama-3	LongChat	Vicuna
ICL	66.64	9.93	16.30
EarlyStop	71.21	11.14	17.44
StructICL	69.43	11.25	17.12
FocusICL	71.89	12.28	17.74
B-LARA	73.86	12.23	18.12

16K (Dong et al., 2023). Following FocusICL, we randomly select 80 examples from the training 364 set and split them into 10 groups for our methods.

The results in Table 2 show B-LARA outperforms FocusICL, which is the previous state-of-the-art 366 method, in 2 out of 3 models. Notably, this is achieved without relying on hidden states, highlighting 367 the simplicity and efficiency of our method on generation task. 368

369 370

345

346

347 348 349

350 351

352 353

354

355

#### 5.2 CAN LARA PERFORM WELL WITH MORE EXAMPLES?

371 We investigate the performance of LARA with an increased number of demonstrations, leveraging 372 the LongICLBench (Li et al., 2024), a benchmark tailored for addressing challenges in long in-373 context learning. For our experiments, we select two datasets: GoEmotion and TacRED. Following 374 the LongICLBench setup, we employ multiple rounds of examples, where each round includes sev-375 eral examples, each labeled with a distinct class. To align with the input limit constraints of ICL, we sampled 8 rounds (224 examples) of examples for GoEmotions and 4 rounds (164 examples) 376 for TacRED. For LARA and B-LARA, we choose 4, 8, and 16 as the potential candidate number of 377 groups. We report the accuracy of different methods on these datasets in Table 3.

		GoEmotion			TacRED		
	Llama3.1-8B Gemma-7B M		Mistral-7B	Llama3.1-8B	Gemma-7B	Mistral-7B	
Black-Box Meth	od:						
ICL	18.60	15.60	17.80	38.20	43.80	55.40	
$RAE_2$	22.20	22.20	21.60	43.80	45.40	55.40	
$RAE_4$	21.00	22.40	21.40	45.60	45.00	52.40	
RAE <sub>8</sub>	21.20	19.00	20.40	36.20	39.00	49.20	
LARA (ours)	21.00	20.80	19.20	48.60	47.40	54.20	
B-LARA (ours)	24.00	22.80	23.80	48.60	49.00	59.00	
White-Box Meth	od:						
ICV	18.80	20.80	18.40	44.40	46.80	54.40	
StructICL <sub>2</sub>	19.00	21.00	18.60	46.60	47.60	55.80	
StructICL <sub>4</sub>	19.80	21.40	19.00	47.60	48.00	56.40	
StructICL <sub>8</sub>	20.60	22.00	19.80	44.80	48.60	56.20	

378 Table 3: Accuracy of methods on GoEmotion and TacRED. The subscript of RAE means the number 379 of groups in RAE.

The experimental results clearly highlight the advantages of LARA, which demonstrates consistent improvements over baseline methods across both GoEmotion and TacRED datasets, showcasing its effectiveness in diverse tasks. Notably, the B-LARA variant further amplifies this performance, outperforming all competing approaches on both datasets and across various models. This suggests that B-LARA can work well in many shot settings.

400 401 402

403

394

396

397

398

399

#### 5.3 CAN LARA PERFORM WELL WITH LIMITED IN-CONTEXT EXAMPLES?

404 In previous experiments, we primarily explore the many-shot in-context learning (ICL) setting. In 405 this subsection, we focus on a more constrained scenario, where only a limited number of in-context 406 examples are available. This analysis aims to understand the relationship between the number of 407 demonstrations and the performance of LARA compared to baseline methods with limited examples.

408 We set the number of examples N within 409  $\{2, 4, 8, 16\}$  and compare our proposed method 410 with ICL on the BBH dataset with Mistral-411 7B. Figure 3 demonstrates that both LARA 412 and B-LARA consistently outperform the baseline ICL, and the performance gap increases 413 with the number of examples used. Note that 414 we do not plot the performance of LARA and 415 B-LARA under N = 2. This is because 416 LARA and B-LARA are simplified to our non-417 reweighting ablation when the size of each sub-418 group becomes 1 and no reweighting is re-419 quired. We also show the performance of per-420 formance without reweighing here. We set the 421 number of group k as 2 in this experiment. 422 While there is a significant gap between the 423 non-reweight version and B-LARA, the nonreweight version still demonstrates effective-424 ness compared to ICL. 425



Figure 3: Accuracy of LARA on BBH using different numbers of examples. B-LARA uses different settings due to differences in example usage during training and inference. We use two lines to highlight this difference. The accuracy means the average accuracy on BBH dataset.

426 Since B-LARA has a weight constraint of  $\{0, 1\}$ , subgroups with zero-weights are pruned during 427 inference for efficiency. As shown in Figure 3, the real number of examples used by B-LARA 428 in inference is substantially lower than other methods. In the 32-shot setting, only about 45% of 429 subgroups of B-LARA are assigned non-zero weights, reducing more than half of the computational load without compromising performance. Additionally, as the total number of examples increases, 430 the proportion of examples used in inference decreases, indicating that B-LARA is particularly 431 suitable for resource-constrained environments.

# 432 5.4 IS LARA APPLICABLE TO BLACK-BOX LLMS?

One advantage of our method is that it could
also be applied to LLM APIs, since it only uses
output logits for example reweighting or selection. In these scenarios, techniques such as incontext vector or task vector, which often rely
on internal state visibility, cannot be applied.

Table 4: Average performance of various method	ds
of GPT-40-mini on the BBH benchmark.	

ICL	LARA	B-LARA
53.17	56.06	57.41

We evaluate our method with GPT-4o-mini<sup>2</sup> on BBH dataset. The results in Table 4 demonstrate that LARA and B-LARA outperform ICL. We note that the OpenAI API only provides top 20 logits for each output token, while our methods are still able to achieve competitive results. This indicates that our method generalizes well to black-box LLMs, and can be applied to situations where internal weights of models are restricted and only output logits are available.

445 446

447

5.5 How Does LARA ENHANCE MEMORY EFFICIENCY?

We empirically evaluate the computational
efficiency of LARA by measuring GPU
memory usage with different input sequence lengths and subgroup configurations. We set the number of groups k with
1,2,4,8. Specifically, when k is set as 1,
LARA will degrade to ICL.

Results in Figure 4 demonstrate that 455 LARA is more memory-efficient com-456 pared to standard ICL, especially when 457 handling long sequences. Standard ICL 458 results in Out-of-Memory (OOM) errors 459 when the input length exceeds 10k tokens 460 on a Mistral-7B model with a batch size of 461 4 on an A100 80GB GPU. In contrast, our 462 method handles input lengths over 25k to-463 kens with 4 and 8 subgroups, demonstrating that LARA efficiently utilizes larger 464 amounts of training data. 465

466 467

468

469

485

#### 5.6 How does the Reweighting Step Affect Model Performance?

We conduct an ablation study to assess the effectiveness of the reweighting step, denoted as "w/o reweight" which simply averages over the output logits of the LLM across different demonstration groups.

In our ablation study, removing the reweighting
step used in LARA also demonstrated its value
by outperforming traditional baseline methods.
For instance, it achieved a notable 67.58 with
Llama3.1-8B in the MMLU benchmark, which
is better than directly ICL (65.63). This perfor-



Figure 4: GPU Memory usage of LARA in gigabytes on a single A100 80GB GPU with different input sequence lengths and number of subgroups. Note that when the number of subgroups equals to 1, the setting is the same as ICL. The sequence length is denoted in thousands of tokens. We set the batch size equal to 4. Data points indicating Out-Of-Memory (OOM) are omitted.

Table 5: Average performance of Llama3.1-8B our methods without reweighting. For the ablation "w/o reweight", the subscript means the size L of each group of demonstrations. The results for other models are shown in Appendix D.1

	1	1
Method	BBH	MMLU
LARA	47.46	66.54
B-LARA	47.69	67.80
w/o reweight <sub>2</sub>	43.33	67.23
w/o reweight <sub>4</sub>	44.50	67.58
w/o reweight <sub>8</sub>	43.02	67.43

mance highlights that logit-arithmetic can successfully combine the information in different groups of demonstrations.

The results further emphasize the importance of the reweighting step in LARA. LARA outperforms the non-reweight version in most settings. This underscores the reweighting process as critical for

<sup>&</sup>lt;sup>2</sup>gpt-4o-mini-2024-07-18

enhancing model accuracy. The worse performance of non-reweight offers clear evidence of how significant reweighting is to optimizing the model's contextual handling.

489

491 492

493

6 RELATED WORK

## 6.1 LONG IN-CONTEXT LEARNING

Recent studies on long-context learning problems in LLMs can be categorized into two main strate-494 gies: enhancing the impact of in-context examples and compressing input sequences. Structured 495 prompting leverages rescaled attention mechanisms to effectively integrate grouped examples (Hao 496 et al., 2022). Methods such as task vectors (Hendel et al., 2023) and function vectors (Todd et al., 497 2023) further refine this strategy by generating vectors that assess the contribution of each example 498 based on the offset of hidden state, which improves model adaptability. Liu et al. (2023) generate 499 task-specific vectors that steer model behavior in latent space based on the in-context examples. 500 Regarding input compression, methods like prompt pruning (Jiang et al., 2023b; Pan et al., 2024) 501 and additional summarization models (Xu et al., 2023a; Gilbert et al., 2023) directly shorten in-502 puts while maintaining essential content. Soft prompt-based compression (Wingate et al., 2022; Mu 503 et al., 2023) intends to generate a soft-prompt that includes most of the information. For many-shot in-context learning problems, previous studies have proposed group-based methods, such as Struc-504 tICL (Hao et al., 2022) and FocusICL (Yuan et al., 2024), which refine attention maps by utilizing 505 subgroup structures within the demonstrations. 506

507 508

### 6.2 LOGIT ARITHMETIC

509 Several works have employed logit arithmetic across various domains and downstream tasks. Con-510 trastive decoding (Li et al., 2022) improves performance by utilizing the difference in logits from 511 models of different sizes. Proxy tuning (Liu et al., 2024) enhances a larger model's capabilities by 512 adding the logit differences of a smaller model, recorded before and after training, to simulate train-513 ing effects. In model arithmetic (Dekoninck et al., 2023), logits adjusted with various prompts steer 514 the generation processes of large language models. Huang et al. (2024) propose using logit subtrac-515 tion to facilitate the selective forgetting of knowledge in LLMs. Additionally, logit arithmetic has 516 been leveraged to enhance the safety of generated outputs (Xu et al., 2024).

517 518 519

### 6.3 NON-GRADIENT OPTIMIZATION OF LLMS

520 Due to the high memory requirements associated with gradient-based optimization methods, recent 521 research has shifted towards non-gradient techniques for neural network optimization. Zhang et al. 522 (2024); Malladi et al. (2023) propose training large language models (LLMs) using non-gradient 523 methods to mitigate these memory constraints. These approaches have also been applied in federated learning, exploring their effectiveness in distributed settings (Xu et al., 2023b). Additionally, a 524 gradient-free method has been used to optimize manifold neural networks (Zhang et al., 2022). 525 Similarly, LoraHub (Huang et al., 2023) utilizes non-gradient techniques to dynamically reweight 526 different LoRA modules, enhancing adaptation to new downstream tasks. Guo et al. (2023) also 527 introduces non-gradient methods to prompt engineering to search for better prompts. 528

529 530

## 7 CONCLUSION

531

532 We proposed LARA, a novel framework that enhances in-context learning by ensembling logits 533 from multiple demonstrations, improving performance without requiring parameter updates. Our 534 method reduces computational complexity while achieving better accuracy. Additionally, Binary 535 LARA further optimizes efficiency by selectively removing less informative demonstrations. Ex-536 periments on BBH and MMLU benchmarks show that both LARA and B-LARA outperform tra-537 ditional ICL methods in terms of efficiency and performance. Future research directions include extending our study to combine logits from different sources beyond just in-context learning (ICL) 538 examples—such as different models or varying instructions—and building a distributed inference system based on LARA.

# 540 REPRODUCIBILITY STATEMENT

We provide detailed information on the dataset we used in Appendix A. The codes for our main experiment can be found in https://anonymous.4open.science/r/LARA-F55B.

#### References

542

545

546

550

561

569

570

571

572

580

581

582

583

588

- an Luo, Xin Xu, Yue Liu, Panupong Pasupat, and Mehran Kazemi. In-context learning with retrieved demonstrations for language models: A survey. *ArXiv preprint*, abs/2401.11624, 2024. URL https://arxiv.org/abs/2401.11624.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-551 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-552 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, 553 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, 554 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCan-555 dlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual 558 Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 559 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Jasper Dekoninck, Marc Fischer, Luca Beurer-Kellner, and Martin T. Vechev. Controlled text generation via language model arithmetic. ArXiv preprint, abs/2311.14479, 2023. URL https: //arxiv.org/abs/2311.14479.
  - Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. A survey for in-context learning. *ArXiv*, abs/2301.00234, 2023. URL https://api.semanticscholar.org/CorpusID:263886074.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, and etc. The Ilama 3 herd of models. ArXiv, abs/2407.21783, 2024. URL https://api.semanticscholar.org/ CorpusID:271571434.
  - Henry Gilbert, Michael Sandborn, Douglas C. Schmidt, Jesse Spencer-Smith, and Jules White. Semantic compression with large language models. 2023 Tenth International Conference on Social Networks Analysis, Management and Security (SNAMS), pp. 1–8, 2023.
- Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian,
   Yujiu Yang, Tsinghua University, and Microsoft Research. Connecting large language models
   with evolutionary algorithms yields powerful prompt optimizers. *ArXiv preprint*, abs/2309.08532,
   2023. URL https://arxiv.org/abs/2309.08532.
  - Nikolaus Hansen and Andreas Ostermeier. Adapting arbitrary normal mutation distributions in evolution strategies: the covariance matrix adaptation. *Proceedings of IEEE International Conference on Evolutionary Computation*, 1996.
- Yaru Hao, Yutao Sun, Li Dong, Zhixiong Han, Yuxian Gu, and Furu Wei. Structured prompting:
   Scaling in-context learning to 1, 000 examples. ArXiv, abs/2212.06713, 2022. URL https: //api.semanticscholar.org/CorpusID:254591686.

- Roee Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. ArXiv preprint, abs/2310.15916, 2023. URL https://arxiv.org/abs/2310.15916.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
   Steinhardt. Measuring massive multitask language understanding. In *Proc. of ICLR*. OpenReview.net, 2021. URL https://openreview.net/forum?id=d7KBjmI3GmQ.
- Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub:
   Efficient cross-task generalization via dynamic lora composition. *ArXiv preprint*, abs/2307.13269,
   2023. URL https://arxiv.org/abs/2307.13269.
- James Y. Huang, Wenxuan Zhou, Fei Wang, Fred Morstatter, Sheng Zhang, Hoifung Poon, and Muhao Chen. Offset unlearning for large language models. ArXiv preprint, abs/2404.11045, 2024. URL https://arxiv.org/abs/2404.11045.
- Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L'elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. ArXiv preprint, abs/2310.06825, 2023a. URL https://arxiv.org/abs/2310.06825.
- Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Llmlingua: Compressing
   prompts for accelerated inference of large language models. In *Conference on Empirical Methods in Natural Language Processing*, 2023b.
- Dacheng Li, Rulin Shao, Anze Xie, Ying Sheng, Lianmin Zheng, Joseph Gonzalez, Ion Stoica, Xuezhe Ma, and Hao Zhang. How long can context length of open-source LLMs truly promise? In NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following, 2023a. URL https://openreview.net/forum?id=LywifFNXV5.
- Mukai Li, Shansan Gong, Jiangtao Feng, Yiheng Xu, Jinchao Zhang, Zhiyong Wu, and Lingpeng Kong. In-context learning with many demonstration examples. *ArXiv preprint*, abs/2302.04931, 2023b. URL https://arxiv.org/abs/2302.04931.
- Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. Long-context llms struggle with
   long in-context learning. *ArXiv preprint*, abs/2404.02060, 2024. URL https://arxiv.org/
   abs/2404.02060.
- Kiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke
   Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. In Annual Meeting of the Association for Computational Linguistics, 2022.
- Alisa Liu, Xiaochuang Han, Yizhong Wang, Yulia Tsvetkov, Yejin Choi, and Noah A. Smith. Tuning language models by proxy. ArXiv preprint, abs/2401.08565, 2024. URL https://arxiv. org/abs/2401.08565.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for GPT-3? In *Proceedings of Deep Learning Inside Out (Dee-LIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pp. 100–114, Dublin, Ireland and Online, 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.deelio-1.10. URL https://aclanthology.org/2022. deelio-1.10.
- Jialin Liu, A. Moreau, Mike Preuss, Baptiste Rozière, Jérémy Rapin, Fabien Teytaud, and Olivier
   Teytaud. Versatile black-box optimization. *Proceedings of the 2020 Genetic and Evolutionary Computation Conference*, 2020.
- Sheng Liu, Haotian Ye, Lei Xing, and James Y. Zou. In-context vectors: Making in context learning more effective and controllable through latent space steering. *ArXiv preprint*, abs/2311.06668, 2023. URL https://arxiv.org/abs/2311.06668.
- Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alexandru Damian, Jason D. Lee, Danqi Chen,
   and Sanjeev Arora. Fine-tuning language models with just forward passes. ArXiv preprint,
   abs/2305.17333, 2023. URL https://arxiv.org/abs/2305.17333.

658

659

660

661 662

663

664

673

682

688

689

690

691

692

- Gemma Team Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, L. Sifre, Morgane Riviere, Mihir Kale, J Christopher Love, Pouya Dehghani Tafti, L'eonard Hussenot, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Am'elie H'eliou, and et al. Gemma: Open models based on gemini research and technology. *ArXiv preprint*, abs/2403.08295, 2024. URL https://arxiv.org/ abs/2403.08295.
- Jesse Mu, Xiang Lisa Li, and Noah D. Goodman. Learning to compress prompts with gist tokens.
   *ArXiv preprint*, abs/2304.08467, 2023. URL https://arxiv.org/abs/2304.08467.
  - Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor Rühle, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and et al. Llmlingua-2: Data distillation for efficient and faithful task-agnostic prompt compression. *ArXiv preprint*, abs/2403.12968, 2024. URL https://arxiv.org/abs/2403.12968.
  - Ingo Rechenberg. Evolutionsstrategie : Optimierung technischer systeme nach prinzipien der biologischen evolution. 1973. URL https://api.semanticscholar.org/CorpusID: 60975248.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, and et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. ArXiv preprint, abs/2206.04615, 2022. URL https://arxiv.org/abs/2206.04615.
- Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau.
   Function vectors in large language models. *ArXiv preprint*, abs/2310.15213, 2023. URL https:
   //arxiv.org/abs/2310.15213.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou.
   Rationale-augmented ensembles in language models. *ArXiv preprint*, abs/2207.00747, 2022. URL https://arxiv.org/abs/2207.00747.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed Huai hsin Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. ArXiv preprint, abs/2206.07682, 2022. URL https://arxiv.org/abs/2206.
  07682.
- David Wingate, Mohammad Shoeybi, and Taylor Sorensen. Prompt compression and contrastive conditioning for controllability and toxicity reduction in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 5621–5634, Abu Dhabi, United Arab Emirates, 2022. Association for Computational Linguistics. URL https://aclanthology.org/2022.findings-emnlp.412.
  - Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang, Prajjwal Bhargava, Rui Hou, Louis Martin, Rashi Rungta, Karthik Abinav Sankararaman, Barlas Oğuz, Madian Khabsa, Han Fang, Yashar Mehdad, Sharan Narang, Kshitiz Malik, Angela Fan, Shruti Bhosale, Sergey Edunov, Mike Lewis, Sinong Wang, and Hao Ma. Effective long-context scaling of foundation models. ArXiv preprint, abs/2309.16039, 2023. URL https://arxiv.org/abs/2309.16039.
- Fangyuan Xu, Weijia Shi, and Eunsol Choi. Recomp: Improving retrieval-augmented lms with compression and selective augmentation. ArXiv preprint, abs/2310.04408, 2023a. URL https://arxiv.org/abs/2310.04408.
- Mengwei Xu, Dongqi Cai, Yaozong Wu, Xiang Li, and Shangguang Wang. Fwdllm: Efficient fedllm
   using forward gradient. *arXiv preprint arXiv:2308.13894*, 2023b.
- Zhangchen Xu, Fengqing Jiang, Luyao Niu, Jinyuan Jia, Bill Yuchen Lin, and Radha Poovendran.
   Safedecoding: Defending against jailbreak attacks via safety-aware decoding. ArXiv preprint, abs/2402.08983, 2024. URL https://arxiv.org/abs/2402.08983.

702 703 704 705	Peiwen Yuan, Shaoxiong Feng, Yiwei Li, Xinglin Wang, Yueqi Zhang, Chuyi Tan, Boyuan Pan, Heda Wang, Yao Hu, and Kan Li. Focused large language models are stable many-shot learners. ArXiv, abs/2408.13987, 2024. URL https://api.semanticscholar.org/ CorpusID:271957090.
706	Liong Thong Dingsong Li Viron Koshy Thakymperempil Sewsong Oh, and Nico Ho. Drzerow
707	Private fine-tuning of language models without backpronagation. In <i>Forty-first International Con</i>
708	ference on Machine Learning, 2024.
709	<u> </u>
710	Rui Zhang, Ziheng Jiao, Hongyuan Zhang, and Xuelong Li. Manifold neural network with non-
/11	gradient optimization. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , PP:1–1, 2022
/12	2022.
713	
714	
710	
710	
710	
710	
720	
721	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

## A DATASET DETAILS

#### A.1 PROMPTS FOR INFERENCE Table 6: Prompt examples for each dataset in One-shot learning. Dataset Prompt BBH Question: {question} Answer: {answer} Question: {*question*}

Answer: MMLU The following are multiple choice questions (with answers) about {*subject*}. Question: {*question*} Answer: {*answer*} Question: {question} Answer: GoEmotion Given a comment, please predict the emotion category of this comment. The predict answer must come from the demonstration examples with the exact format. The examples are as follows: comment: {*question*} emotion category: {*answer*} comment: {*question*} emotion category: TacRED Given a sentence and a pair of subject and object entities within the sentence, please predict the relation between the given entities. You can only select from the following words: {*potential relation*} sentence: {question} the relation between the two entities is: {*answer*} sentence: {question} the relation between the two entities is: A.2 DATASET STATISTICS Table 7. Dataset Statistics 

Dataset	#Tokens/Shot	Description
BBH	55	A collection of challenging tasks from the BIG-Bench Hard benchmark.
MMLU	65	Multiple-choice questions across various subjects.
GoEmotion	28	Annotated Reddit comments for emotion classification.
TacRED	80	A dataset for relation extraction tasks.

Model	Llama	3-8B-Instruct	LongCh	at-7B-V1.5-32K	Vicuna-7B-V1.5-16K			
	ARC	GSM8K	ARC	GSM8K	ARC	GSM8K		
ICL	90.00	66.64	62.43	9.93	77.11	16.30		
EarlyStop	90.47	71.21	62.43	11.14	78.14	17.44		
StructICL	90.70	69.43	64.05	11.25	78.05	17.12		
FocusICL	91.02	71.89	64.55	12.28	78.51	17.74		
<b>B-LARA</b>	90.89	73.86	64.27	12.23	78.79	18.12		

820

810

#### Algorithm 1 B-LARA Optimization Algorithm with Updated Index

821 Input:  $\mathcal{D}_{\text{train}}$ : In-context examples  $\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N$ . 822 Parameter: k: Number of subgroups. J: Number of iterations. 823 **Output:**  $w^*$ : Optimized binary weight vector. 824 Split  $\mathcal{D}_{\text{train}}$  into k groups:  $\{S_1, S_2, \dots, S_k\}$   $S_A \leftarrow \{S_1, \dots, S_{\lfloor k/2 \rfloor}\}$   $S_B \leftarrow \{S_{\lfloor k/2 \rfloor+1}, \dots, S_k\}$ 825 for  $r \in \{A, B\}$  do 826 Initialize  $w^{(0)}$  as a random binary vector of length  $|S_r|$ 827 for j = 1 to J do 828 for m = 1 to  $dim(\boldsymbol{w}^{(j-1)})$  do 829  $u_m \leftarrow \text{Uniform}(0,1)$ 830  $w'_m \leftarrow w_m^{(j-1)} \oplus \mathbb{I}(u_m < 1/\dim(\boldsymbol{w}^{(j-1)}))$ 831 end 832 Compute  $\mathcal{L}(\boldsymbol{w}')$  using  $\mathcal{S}_{r'}$ , where  $r' \neq r$ 833 if  $\mathcal{L}(\boldsymbol{w}') \leq \mathcal{L}(\boldsymbol{w}^{(j-1)})$  then

834

#### 835 836 837

838 839 840

841 842

#### 843 844

845 846

847

848

#### **COMPARISION TO PREVIOUS METHODS** В

С ALGORITHM

#### 849 FULL RESULTS D 850

851 D.1 FULL ABLATION STUDY 852

 $w^{(j)} \leftarrow w'$ 

 $\pmb{w}^{(j)} \leftarrow \pmb{w}^{(j-1)}$ 

else

end

 $oldsymbol{w}_r^* \leftarrow oldsymbol{w}^{(J)}$ 

 $oldsymbol{w}^* \leftarrow [oldsymbol{w}_A^*, oldsymbol{w}_B^*]$ 

end

return  $w^*$ 

end

853 D.2 FULL MAIN RESULTS 854

Here we will show the full results of our three models in BBH and MMLU benchmark. The methods 855 include LARA, B-LARA, KATE, ICL, LAG(logit-average-generation which is the ablation study in 856 our paper, together with RAE and ICV. 857

858

- 859
- 860 861
- 862

Table 9: Ablation Study Results

		$BBH_{average}$		Μ	$MLU_{average}$			
	Llama3.1-8B	Gemma-7B	Mistral-7B	Llama3.1-8B	Gemma-7B	Mistral-7		
ICL	45.64	37.08	42.91	65.63	61.44	62.84		
$KATE_2$	43.60	37.07	43.16	66.62	56.28	63.99		
KATE <sub>4</sub>	44.03	38.83	43.16	66.75	55.78	63.48		
KATE <sub>8</sub> 44.47         37.03           RAE <sub>2</sub> 44.59         40.24		42.96	67.19	54.13	63.93			
		40.24	43.95	66.88	65.18	62.99		
$RAE_4$	45.23	40.44	44.49	66.40	65.01	62.99		
RAE <sub>8</sub>	44.06	39.85	44.07	67.09	64.80	63.61		
ICV	45.93	42.16	44.50	66.97	64.99	64.02		
LARA	47.46	41.77	44.77	66.54	64.36	63.93		
B-LARA	47.69	42.75	45.03	67.80	65.56	64.12		
w/o reweight <sub>2</sub>	43.33	43.56	42.83	67.23	65.61	62.95		
w/o reweight <sub>4</sub>	44.50	41.98	44.78	67.58	65.87	63.32		
w/o reweight <sub>8</sub>	43.02	39.35	44.84	67.43	65.04	63.55		

920 921																				
922																			I	I
923	C	54	55	8	.45	23	.35	53	34	23	.98	.45	4	53	57	34	45	.65	50	
924	Ι	24	65	58	39	38	52	60	99	Ś	14	52	63	35	29	86	32	31	4	
925	~	~	<del>. +</del>	0	<del></del>	~		~	~		5	~	0	_	~	_	_	$\sim$	~	
926	<b>AE</b>	5.23	5.12	6.42	3.92	1.2	7.8(	4.7	0.9	0.55	6.0	3.6	9.72	7.6	2.5	4.6	5.1	7.98	4.0	
927	$\mathbf{R}_{I}$	2	Ö	Ñ	ŝ	4	ŝ	Ŵ	Ñ	Ξ	÷,	ŝ	9	è	ŝ	ò	ŝ	6	4	
928	4	6	0	2	_	ŝ	6	6	6		9	2	0		6	4	6	Ś	6	
929	AE	3.3	3.3	6.4	7.6	1.2	4.5	9.5	5.0	1.4	6.0	6.4	7.8	0.3	11	6.6	3.4	1.6	4.	
930	R	2	9	ŝ	ŝ	4	ŝ	ŝ	ŝ	Τ	-	ŝ	ŝ	4	ŝ	×	ŝ	$\mathfrak{c}$	4	
931	5	-	×	9	6	4	-	Ō	2	6		×	Ŀ.	4	0	3	×	Ģ	S	
932	AE	26.6	2.6	55.9	88.9	36.8	33.2	58.9	54.1	2.3	6.5	51.3	51.4	36.2	5.7.5	32.0	35.7	4.4	13.9	
933	R	(1	U	4,1	61	61	4,	4,1	43	_	_	4,	U	6.1	(1	00	61	61		
934	8	13	33	6	5	5	33	3	80	4	6	33	2	2	90	井	9	33	35	_
935	'AG	4.7	7	56.9	2.5	35.3	57.5	1.1	52.5	13.4	15.5	57.5	38.	6.0	27.5	85.4	32.2	29.(	4.	<u> </u>
936	Γ						.,	• •				• •	Ū	4		•••				al-,
937	7	57	58	53	86	59	92	28	50	53	98	30	4	25	4	03	18	26	78	istr
938	JA(	25.	72.	53.	40.	36.	48.	62.	62.	10.	13.	54.	63.	39.	27.	87.	31.	32.	4	Ð
939	Ι																			H
940	$5^2$	Ξ	68	23	25	15	16	3	52	75	52	46	4	17	26	54	3	26	84	BB
941	LA	30.	59.	53.	39.	34.	45.	57.	56.	10.	14	49.	59.	38.	32.	80.	36.	32.	42.	in.
942																				sks
943	CL	.74	.28	.60	.48	.15	.05	.49	.83	22	.52	.61	.52	.71	.81	.65	.33	.49	.91	s ta
944	Ĩ	17	68	58	35	34	65	46	4	10	14	51	64	38	25	88	33	28	4	SO
945	x	x	~	×	_		2	×	S	6	ŝ	0	4	_		S	x		6	acı
946	TE	8.2	8.8	5.3	4.	5.3	9.4	0.8	4.3	1.2	5.0:	0.0	<u>'</u> .'	8.7	5.5	5.9	Ξ	3.8	2.9	res
947	KA	-	9	ŝ	ŝ	$\mathfrak{C}$	4	Ś	Ś	Ξ	-	ŝ	9	ŝ	0	×	ξ	ŝ	4	sco
948																				ce
949	$\mathbf{E}_4$	.89	4	.30	.26	.02	.91	.39	.35	.29	.28	8.0	2	.63	.88	.78	.26	.96	.16	nan
950	<b>EAT</b>	19	3	5	32	39	55	54	54	Ξ	$\frac{18}{8}$	50	5	37	26	8	32	51	4	
951	Ł																			erf
952	$\mathbb{E}_2$	19	83	15	63	59	53	89	72	75	20	62	43	17	34	95	<del>8</del>	<del>4</del> 9	16	Ч
953	AT	24.	61.	52.	37.	36.	53.	57.	52.	10.	17.	4	70.	38.	26.	85.	35.	28.	43.	E E
954	K																			ble
955	V	<del></del>	0	4		0	6	e	e	0	2	œ	З		0	Э	2	2	3	Ë
956	AR.	5.8	2.72	4.8	8.7	7.8	2.6	0.5	5.4	8.6	5.0	1.0	4.7	8.7	2.8	7.0	0.6	3.8	5.0	
957	- <b>L</b> /	(1	Û	<b>W</b> )	<b>G</b> )	<b>G</b> )	<b>W</b> )	Ð	<b>v</b> )		—	<b>W</b> )	(-	<b>G</b> )	<b>G</b> 1	œ	<b>G</b> 1	<i>a</i> )	4	
958	В																			
929	Y	32	5	90	20	15	39	57	39	66	4	38	80	33	33	76	20	\$	5	
960	AR	40.	1.1	58.(	32.1	34.	<del>7</del> 8	66.(	67.	6	13.4	51.(	59.0	33.	29.0	76.	36.5	35.4	4	
967	Γ																			
963																	et			
964							es		s.				n				Ð			
965	ıme	þ	Ø	q	bj3	ble	lap	1	me	bj7	bj5	ŝ	ato	S		ec	ßEı	10f		
966	$\mathbf{N}_{\mathbf{a}}$	pSe	ĝ	Ūn	КÕ	Tal	nSI	ķ	Ŋa	КÕ	КÕ	)ed	erb	)ed	)ed	ieR	ran	Ŝ	age	
967	ask	em	isa	ate	rac	eng	feor	nar	uin	rac	rac	ogl	[yp	lgo	lgo	lov	alT	eas	ver	
968	H	Ē	ά	Ω	Ξ	đ	0	S	R	Η	Η	Г	Ħ	Γ	L	Z	Ñ	R	A	
969																				

972 973																					
974	1	I																			I I
975	CV	21	20	.52	8.	.50	.50	8.	.33	.78	.30	8.	.20	.48	8.	.43	.50	.10	.56	00.	
976	Ĩ	30	55	63	50	5	F	50	48	27	63	36	50	58	46	60	43	37	30	74	
977	x	-	9	~	4	0	Ś	9	9	2	2	×	0	3	0	6	0		9	0	
978	AE	6.1	0.5	8.1	4.	1.0	3.7	5.5	5.5	2.2	4	1.5	5.0	6.7	6.0	0.4	9.5	3.8	0.5	2.0	
979	R	$ $ $\infty$	9	9	4			ŝ	ŝ	2	9	ŝ		ŝ	4	9	ŝ	ŝ	ŝ		
960	4	9	5	Ξ	4	Q	2	ø	4	Q	2	4	Q	6	Q	2	Q	3	3	Q	
201	AF	30.5	5.13	5.5	4.4	70.0	13.7	52.7	4.4	25.0	2.7	36.8	75.0	37.8	36.0	58.0	£0.0	0.63	33.3	57.0	
302	R	01	0	0	7		( -	4,	7	( I	0	0.1	( -	41	0.1	47	7	( I	0.1	0	
984	57	33	75	32	67	8	25	78	4	78	53	51	53	76	8	3	8	4	57	8	
985	<b>[A</b> ]	33.	57.	66	41.	69.(	71.5	52,	4	27.	.43	34.	72.7	54.0	40.0	58.0	41.0	27.4	41.0	74.0	
986	I																				
987	ß	56	75	2	67	50	8	33	8	78	30	Ξ	22	4	8	02	50	45	33	8	
988	ΓV	30.	57.	63.	41.	70.	75.	58.	50.	27.	63.	45.	72.	56.	4	58.	39.	27.	33.	76.	art
989																					- D
990	$G_4$	Ξ	.15	.91	.67	8.	00.	.56	0.0	22	22	.95	4	.73	8.	.26	00.	54.	.33	00.	- <u>7</u> B
991	ΓA	36	59	65	41	5	75	55	50	22	4	28	69	56	4	59	41	27	33	75	ral-
992						_					_				_	_	_			_	list
993	$\mathbf{G}_2$	7.78	1.75	5.1	l.67	.50	3.75	2.78	.22	22	3.30	l.58	52	1.97	3.00	.49	.50	.42	4.	00. <del>1</del>	5
994	LA	5	ŝ	6	4	Ч	5	ŝ	4	8	6	3	F	ŝ	3	90	¥	2	4	Ľ	F
995	. 1	9	5	4	~	0	5	0	_	_	4	ŝ	~	6	0	ŝ	0	0	_	0	Įξ
996	ICI	0.5	1.1	3.6	7.2	0.0	3.7	0.0	6.1	6.1	5.1	8.9	2.2	7.8	0.0	1.7	0.0	7.10	6.1	4.5	
997		ς β	Ś	9	4			Ś	ξ	ξ	9	0	2	ŝ	4	9	4	$\mathfrak{C}$	ŝ	2	cs i
998	Е <sub>8</sub>	56	34	91	8	50	75	11	67	Π	30	32	4	73	8	96	50	<del>4</del> 8	89	50	task
999	Π	30.	56.	65.	50.	73.	78.	61.	4	36.	63.	26.	.69	56.	40.	62.	40.	35.	38.	76.	ss 1
1000	K																				cro
1001	4	e	9	5	Q	0	Q	ŝ	0	c	0	-	0	9	Q	2	0	ò	ŝ	Ō	es a
1002	TE	33.3	50.5	2.7	50.0	71.5	75.0	58.3	17.2	33.3	53.3	2.2	72.2	55.5	4.0	58.0	±1.5	35.4	33.3	75.0	COL
1003	KA	01	0	0	41		( -	4,	7	01	0	0.1	( -	41	7	47	7	01	0.1	( -	e se
1004			~		_	_			~	~	~		~		_	~	_		~	_	anc
1005	$\mathbf{IE}_2$	2.23	1.9	26.7	0.0	8.5(	3.75	5.56	2.78	8.8	2.39	5.82	6.67	5.56	8.00	1.73	8.00	2.26	7.78	7.00	L m
1005	KA	6	9	ìO	ñ	ö	Ē	ŝ	ŝ	õ	6	õ	õ	ŝ	ŝ	9	õ	é	'n	ŕ	l fi
1007																					Pe
1008	RA	.56	.56	.05	8	.50	.50	Ξ.	.89	8.	.39	.47	8	.56	8.	.20	.50	.26	.89	.50	1 ::
1010	[Y]	30	9	67	50	5	5	61	38	52	62	39	75	55	46	4	41	32	38	74	le
1011	<b>B-</b> ]																				Tal
1012		~						~	~	0	+	~	<del>. +</del>	<b>~</b>		5		5	_		
1013	RA	7.78	9.1.	7.05	0.0	9.0(	6.25	8.3	1.67	2.23	5.12	9.4	4.6	7.89	0.0	2.96	2.00	2.26	6.11	5.00	
1014	LA	6	ŝ	Ó	Ñ	9	7	ŝ	4	6	Ö	ŝ	9	ŝ	Ñ	9	4	ŝ	õ	1	
1015									c).												
1016									Suc							-	tics				
1017						c)			scie	ics				s		ing	ma			Ň	
1018		ra				gb		try	ter	nat	ne		rity	vsic		leel	uthe			30lc	
1019		jeb)			hics	wl	Q0	mis	ndi	her	lici	sics	ecu	hd	S	ngi	т	ల		bi	
1020	me	alg	>	ny	etl	knc	biol	chei	COL	nat	nec	yhc	SL-S	ual.	etri	ul_ei	ary	igo	acts	lool	
1021	Na	act	<b>m</b> o	IOU	less	Cal	ge_]	e.	ge_	ge_	ge_1	ge_]	oute	ept	Õ	ricé	ent	al	J.L	scł	
1022	ask	bstı	natu	stro	usir	inic	olle	olle	olle	alle	olle	olle	luic	onc	COD	ect	em	m	lob	igh.	
1023	Ĥ	a	a	ä	ã	5	ప	ర	ర	ర	ర	ప	ప	ప	ð	e	e	f	50	μ	
1024																					
1025																					

1026 1027																					
1028	>		ŝ	ŝ	_	Ś	0	0	_	×	0	2	6	3	4	Ś	<del>_</del>	ŝ	6	2	
1029	IC1	8.2	8.3	2.7	7.3	6.0	5.5(	9.5(	4.4	3.6	0.5(	3.9,	4 2	6.13	4.	0.1:	7.2	2.7	8.7	7.9	
1030		4	ŝ	-	×	×	9	ε	9	З	×	ŝ	×	×	9	-			~	4	
1031	$\mathbf{F}_8$	8	Ξ	5	58	72	8	50	52	33	8	95	71	03	30	13	19	27	78	83	
1032	[Y]	46.	61.	78.	83.	83.	62.	38.	65.	33.	80.	53.	80.	78.	67.	73.	77.	77.	H.	45.	
1033	Ι																				
1034	$\mathbf{E}_4$	00	33	51	84	72	50	8	67	48	50	32	86	03	55	12	4	27	84	67	
1035	RA	41.	58.	79.	82.	83.	61.	39.	66.	34.	80.	51.	82.	78.	68.	76.	75.	77.	83.	41.	
1036																					
1037	$\mathbf{E}_2$	96.	.56	:23	.34	.50	00.	.50	.37	.48	.50	.97	.43	.61	.30	.61	4	00.	.81	.67	
1038	RA	53	55	F	81	84	59	43	64	34	80	51	81	78	67	5	75	75	80	41	
1039																					
1040	$G_8$	96.	Ξ	52	.34	.95	.50	.50	.67	.18	.50	.32	8.	.61	.81	.12	.19	.27	.81	.92	
1041	LA	53	61	78	81	8	62	37	66	32	80	51	80	78	69	76	F	F	80	4	t 2
1042		_		- )		- )	_	_	- `		_						~~	_	_		Par
1043	<b>\G</b> ₄	5.40	8.33	23	3.58	212	0.0	0.1	7.82	2.18	0.1	2.63	2.14	S.03	3.55	t.63	3.68	0.0	9.80	5.83	B B
1044	LA	5.5	ŝ	22	ŝ	õ	90	4	6	б	8	ŝ	õ	22	õ	2	5	Ľ	5	4	1-7
1045	2	~	m		6	~			6	$\sim$		_	<del>. +</del>	ŝ		~	ŝ		6	~	stra
1046	₹Ğ	4.6	8.3	0.2	5.0	3.7	0.0	1.0	6.0	6.7	0.5(	4.6	5.1	8.0	7.3(	6.1.	1.9	7.2	8.7	1.6	Mis Mis
1047	$\mathbf{L}_{i}$	ý	ŝ	õ	òò	òò	Ø	4	õ	õ	õ	Ń	òò	1	Q	7	L-	í-	1	4	
1048	L	0	9	ŝ	×	2	0	0		×	0	6	0	S		2	0	Э	×	3	E
1049	ICI	1.8	5.5	7.2	3.5	3.7	2.5	1.5	2.0	4.4	0.5	3.2	0.0	0.3	6.6	6.1	0.7	2.7	L.L	5.8	\$
1050		ŝ	ŝ	-	œ	œ	0	e)	0	e	œ	ŝ	œ	œ	0	-	œ	-	-	ব	II.
1051	$\mathbf{F}_8$	40	89	21	57	50	8	50	37	03	50	92	29	46	67	4	20	55	78	92	sks
1052	AT	55.	63.	62	86.	2.	63.	35.	2.	31.	81.	55.	79.	E	66.	71.	80.	79.	E.	47.	tas
1053	К																				oss
1054	4	5		<del></del>	4	2	0	0	4	4	0	<del></del>	0	2	0	e	S	0	×	2	acr
1055	ΤE	2.5	9.9	9.2	2.8	6.8	2.5	3.0	4.9	8.7	1.5	4.6	0.0	5.7	7.3	4.6	8.9	5.0	L.L	3.7	es
1055	KA	<b>u</b> )	Û	(~	00	œ	U	<b>G</b> 1	U	(1	00	<b>W</b> )	œ	(-	Ð	(~	(~	(~	(~	4	COI
1057		~~	_	- )			_	_	_		_							~			s s
1050	$\mathbb{IE}_2$	t.68	<u>8</u> .	23	3.58	5.05	3.50	3.50	33	5.78	0.0	3.55	55	S.03	3.55	7.61	<u>.</u> 93	7.27	.81	3.33	anc
1060	(V)	5	6	22	õ	×	6	ŝ	õ	ĕ	<u>∞</u>	ŝ	ò	22	õ	ŀ	1	ŀ	<u></u>	$\tilde{\mathcal{D}}$	
1061	1																				erfc
1062	¥٨	12	89	33	57	05	8	50	82	78	50	89	4	03	92	15	19	55	20	92	- B
1063	ΨY	56.	63.	5	86.	86.	62.	29.	67.	36.	80.	57.	82.	78.	67.	70.	5	79.	69.	47.	12
1064	B-I																				ble
1065							_	_													Tal
1066	RA	.83	8.	:33	.82	.27	8	8	.39	.48	.50	.92	.14	.46	.92	-64	.19	.55	17.	.08	
1067	<b>TA</b>	56	63	F	85	85	62	32	68	<del>2</del>	80	55	8	F	67	7	F	79	50	52	
1068																					
1069			nce	ory			cs		S												
1070			cie	isto		<u>ц</u>	Ē	8	m					Ŋ.							
1071		Ŷ	SL'S	n_h	hy	nen	0UQ	atio	ono		ğ		<u>v</u>	isto							
1072		uist	Jute	pea	rap	rnn	)90.	lem	0ec	ics	lol	štic	sto	dh						<b>"</b> ••	
1073		nen	Įmc	lol	1905	DVel	laci	lath	licr	hysi	sycl	atis	s-hi	orl		lity	law		8	ning	
1074		l_cl	Ŭ_I	l_eı	ğ	ğ	l m	l m	l m	I_p	l_p	Lst	l_u	l_w	gu	ual	al	nce	acie	arr	
1075	1ME	100	100	100	100	100	100	100	100	100	100	100	100	100	agi	sex	tion	abt	fall	e_le	
1076	Ň	sc	SC	SC	Sc	SC	Sc	Sc	SC	Sc	SC	Sc	SC	Sc	an	an	na	pru	al	hiņ	
1077	ask	igh	iigh	ligh	igh	ligh	igh	igh	igh	igh	igh	igh	ligh	ligh	um	mm	ntei	uris	Dic	nac	
1078	Γ	q	Ч	Ч	Ч	Ч	Ч	Ч	Ч	Ч	Ч	Ч	Ч	Ч	Ч	Ч	. <b>I</b>	·-,	ľ	u	
1079																					

080																					
1081																					
1082																					I I
1083	СV	.18	5	.11	.50	8	.50	.50	8.	.50	8	.50	8.	8.	60.	69.	.32	.89	.96	.11	.03
1084	I	87	88	86	80	74	4	76	74	68	53	46	69	69	76	69	88	88	51	84	64
1085	~	~		4	0	0	0	0	0	0	0	0	0	0	_		_	6	~	~	
1007	ΑĘ	7.13	6.4	4.6	5.0	3.5	5.5	3.5(	4.0	0.0	0.0	2.5	6.0	0.0	3.9	1.2	0.5	8.8	3.9	6.9	3.6
1087	R	8	×	9	$\infty$	~	4	5	-	5	Ś	ŝ	9	5		L-	6	×	Ś	×	6
1000	4	8	6	2	0	0	0	0	0	0	0	0	0	0	6	Э	4	~	2	×	6
1009	AE	7.1	5.2	2.2	1.0	3.0	9.5	3.0	2.0	9.5	1.5	0.5	5.0	9.5	6.0	4.0	1.2	1.6	3.9	5.9	2.9
1090	R	8	×	-	×	-	ŝ			9	ŝ	ŝ	9	9	~	-	6	6	Ś	×	9
1091	2	8	×		Q	0	0	Q	0	0	Q	Q	Q	0	6	ò	<del></del>	6	4	2	6
1092	AE	37.1	5.8	9.9	1.0	2.5	5.6	6.0	0.5	8.5	0.0	1.0	6.0	5.6	6.0	3.4	0.5	8.8	2.9	5.0	2.9
1093	R	8	œ	Û	00	(~	<b>G</b> 1	(~	(-	U	<b>W</b> )	<b>u</b> )	Ð	Q	(~	(~	5	œ	<b>W</b> )	00	
1094	8	52	L L	2	2	õ	0	2	2	2	0	0	2	0	1	33	17	6	2	-	5
1095	'AG	34.6	36.4	2.2	31.0	2.5	t3.5	73.0	73.0	57.0	18.5	51.5	58.0	20.5	73.5	74.0	90.5	88.8	53.5	¥.1	53.5
1090	Τ	~	~	( ·	~	( ·	7			C	7	47	U				0,	~	47	~	
1002	44	8	÷	4	00	2	2	2	20	20	00	2	20	2	6	33	51	68	2	3	22
1090	'YC	87.	86.4	5 <del>.</del> 7	80.4	74.(	47.(	73.(	70.4	58.5	51.5	<del>1</del> 9.(	67.5	71.(	76.(	72.5	90.4	88.2	53.9	85.(	63.3
1100	Ι			-		`	4	`	`	-		4	-	`	`	`	•				
1101	25	18	4	67	20	8	8	8	8	50	8	8	50	8	60	59	51	89	32	11	96
1101	)VC	87.	86.	<u>66.</u>	80.	74.0	45.(	74.(	70.	67.	49.0	49.0	65.	69.	76.0	47	90.	88.	53.	84.	62.
1102	Ι																				
1104	T	18	4	67	8	8	8	8	8	50	8	50	8	8	91	03	32	Ξ	98	05	84
1105	I	87.	88.	66.	81.	73.	42.	74.	69.	68.	49.	51.	68.	70.	73.	74.	88.	86.	50.	85.	62.
1106			_		_	_	_	_	_	_	_	_	_	_	_						
1107	$\mathbb{TE}_8$	2.05	4.	53	3.00	3.50	0.0	00. <del>1</del>	t.50	3.50	0.0	1.50	7.50	0.0	<b>60</b> .0	.72	9.05	5.11	5.86	<u>†</u> .1	3.93
1108	<b>ZA</b> J	8	×	F	õ	2	4	Ľ	2	30	S.	S	6	F	2	Й	×	8	Š	8	63
1109	1																				
1110	$\mathbf{E}_4$	62	7	53	8	8	50	50	50	8	50	8	50	50	57	38	78	Ξ	88	05	48
1111	AT	84.	<u>8</u> .	72.	79.	72.	39.	71.	73.	67.	50.	53.	66.	72.	69.	72.	89.	86.	55.	85.	63.
1112	К																				
1113	2	6	_	4	0	0	0	0	0	0	0	0	0	0	<del></del>	Э	S	~	2	-	6
1114	ΤE	9.4	7.7	4.6	2.5	2.5	0.0	4.0	1.5	57.5	0.5	2.0	9.5	5.6	3.9	2.9	9.0	1.6	3.9	4.1	3.9
1115	KA	(-		U	æ	(~	7	(-	(-	U	4,	4,1	U	U	(~	(~	00	0,	4,1	8	
1116		_			_	_	_	_	_	_	_	_	_	_							
1117	RA	.49	5	2	8.	1.50	.50	50	8.0	0.0	.50	8	0.0	8	.43	27	.51	6	.98	.98	112
1118	ΓV	75	8	6	8	77	4	5	5	6	4	5	2	5	80	7	90	91	50	85	64
1119	B																				
1120		(	~	~	_	_	_	_	_	_	_	_	_	_	~	_	~	~	~	~	~
1121	RA	9.49	5.8	53	3.0(	5.5(	3.0(	2.5(	3.5(	6.5(	9.5(	1.5(	8.00	2.5(	0.43	9.6	9.78	1.67	0.9	5.98	3.95
1122	LA	E-	òò	È,	òò	1	4	1	1	õ	4	ŝ	Ø	1	õ	9	õ	6	Ñ	ŝ	6
1123														~							
1124											ling		e	0g							
1125											III		icir	hol							
1126				ics		ş	SO				000	ME	led	syc	SL	SS		icy		S	
1127		nt		net	SUC	ute	ari				ul_a	l_l	l_n	l_p	tiol	ibi		pol		ion	
1128	Ime	me	gu	-ge	nec	lisp	cen	u	yhc	iry	0 <b>n</b> £	0n£	0 <b>n</b> £	0 <b>n</b> £	ela.	sti	<b>N</b>	E C	~	elig	
1129	$\mathbf{N}_{\mathbf{B}}$	age	keti	ical	ella	al_d	al s	itio	loso	isto	essi	essi	essi	essi	ic_r	rity	gol	irei	<u></u>	d_h	age
1130	ask	an	lar	ledi	uisc	lor	lor	utri	hild	reh	rof	rof	rof	rof	ldu	) Scur	čio	s_fo	lo	orl	ver
1131	E	B	Ξ	H	ũ	H	ũ	n	d	d	d	d	d	d	d	š	S	ñ	2	M	a,
1132																					-

Table 13: Performance scores across tasks in MMLU (Mistral-7B) Part 2

Task Name	LARA	<b>B-LARA</b>	$KATE_2$	$\mathbf{KATE}_4$	<b>KATE</b> <sub>8</sub>	ICL	$\mathbf{LAG}_2$	$\mathbf{LAG}_4$	$\mathbf{LAG}_{8}$	$\mathbf{RAE}_2$	$\mathbf{RAE}_4$	$\mathbf{RAE}_8$	ICV
TempSeq	25.27	19.89	18.82	25.27	16.13	17.20	19.89	17.74	16.67	20.18	18.81	15.14	14.52
DisambQA	65.59	67.74	55.38	61.83	65.05	73.66	67.74	68.28	65.59	61.47	65.60	68.35	67.74
DateUnd	48.92	54.30	38.17	43.55	40.32	40.32	54.30	52.69	48.39	49.54	47.71	48.62	54.84
TrackObj3	32.26	32.80	39.78	38.71	29.03	37.10	38.71	34.95	34.41	32.57	30.28	37.61	34.95
PengTable	42.68	45.12	40.24	50.00	53.66	43.90	42.68	46.34	41.46	38.60	40.35	49.12	43.90
GeomShapes	54.30	47.31	51.08	56.99	48.92	63.98	46.77	48.92	46.77	44.04	53.67	50.00	64.52
Snarks	55.26	58.77	58.77	53.51	48.25	50.88	57.89	57.02	52.63	56.85	60.27	58.90	59.65
RuinNames	36.96	31.52	25.00	27.17	27.17	28.26	38.59	31.52	30.43	31.02	29.17	25.93	36.96
TrackObj7	12.90	11.83	15.05	14.52	13.98	12.90	12.37	15.59	10.75	5.96	11.93	13.76	11.29
TrackObj5	17.20	15.05	17.20	13.44	14.52	20.97	16.67	18.82	20.43	17.43	15.14	18.81	17.74
LogDed3	39.25	46.77	37.10	40.86	40.32	43.55	50.54	49.46	43.55	49.08	45.87	43.12	52.69
Hyperbaton	68.82	66.13	59.14	56.99	55.38	46.77	65.05	65.05	58.06	61.93	55.50	61.93	60.75
LogDed5	32.80	34.95	30.65	27.96	20.97	23.12	40.86	33.87	26.88	33.03	33.49	28.44	41.40
LogDed7	41.94	46.24	19.89	21.51	17.74	17.20	43.01	30.11	24.73	39.45	33.94	22.48	32.26
MovieRec	72.43	75.68	54.59	61.62	69.73	76.22	71.35	70.81	73.51	68.20	66.82	64.06	87.03
SalTransErrDet	28.49	30.65	29.03	29.03	30.11	0.00	32.26	34.41	32.26	30.73	38.99	28.90	0.00
ReasColObj	34.95	41.94	40.32	37.10	38.17	34.41	41.94	38.17	42.47	44.04	39.91	42.20	36.56
Average	41.77	42.75	37.07	38.83	37.03	37.08	43.57	41.99	39.35	40.24	40.44	39.85	42.16

(Gemma-7B)
ks in BBH
s across tas
ance score
4: Perform
Table 1₄

1188 1189 1190																					
1191	^	0	1	2	12	õ	5	0	90	68	6	Ξ	68	21	2	7	2	6	0	0	
1192	IC	25.(	67.6	<u>6</u> 6	<u>56.</u> (	73.5	81.2	50.0	55.5	38.8	62.3	5.1	63.8	57.3	56.(	70.3	43.0	48.5	25.0	81.0	
1193			-	-	-	`					-	1	-			Ì		4			
1194	$\mathbf{E}_{8}$	.56	.15	.86	.78	.50	00.	.33	.33	.89	.30	.84	8	.23	8	.73	.50	:39	.56	00.	
1195	RA	30	59	5	52	88	75	58	33	38	63	36	75	90	4	61	4	48	30	79	
1196				~	_	_	_			~	~	~	_	_	_	•	_	_	~	_	
1197	<b>AE</b>	0.50	7.7	5.13	0.0	0.0	7.5(	1.1	4.	3.33	5.5	1.58	5.00	3.72	5.0	8.02	6.00	8.39	8.8	8.5(	
1198	$\mathbf{R}_{1}$	3	ì	È,	Ñ	7	i-	9	4	č	Ó	ŝ	Ë	Ö	4	ŝ	4	4	õ	F	
1199	5	9	ŝ	0	×	0	0	9	×		6	8	0	6	0	9	0	6		0	
1200	AE	30.5	27.7	75.0	52.7	9.0	12.5	5.5	52.7	H.6	52.3	31.5	75.0	51.9	12.0	9.2	15.0	18.3	36.1	0.6	
1201	R	а,	<b>U</b> )	(~	<b>u</b> )	Q	(~	<b>v</b> )	<b>W</b> )	7	Φ	<b>a</b> )	(~	Ð	7	<b>W</b> )	7	7	<b>a</b> )	(~	
1202	~8	57	4	2	20	2	75	26	78	26	30	35	2	35	2	20	2	2	Ξ	20	
1203	'YC	41.(	56.	75.(	55.5	66.(	73.7	55.5	52.7	30.4	63.	28.5	12.2	59.6	4	5.5	45.(	50.0	36.	80.4	t
1204	Ι	-													-		-				Pai
1206	$\mathbf{G}_4$	33	34	8	56	50	25	56	78	89	22	95	8	.16	8	26	8	Ľ	4	8	B B
1207	LA	33.	56.	75.	55.	7.	76.	55.	52.	38.	64.	28	75.	63.	46.	59.	46.	46.	4	79.	a-7
1208																					un l
1209	$G_2$	00.9	.34	.73	.78	8	.25	.33	0.0	6.11	.30	.21	52	66	8.	.56	i.50	0.0	.67	.50	Gei
1210	LA	25	56	6	52	7	76	58	50	36	69	34	5	61	4	55	4	50	4	78	Η
1211	. 1	~	0	~	ŝ	0	0	ŝ	_	~	S	_	~	ŝ	0	5	0	4	9	0	BB
1212	ICI	7.7	0.7	9.3	8.3	6.5	0.0	8.3	6.1	2.2	0.5	2.1	7.7	6.7	0.0	<u>.</u>	8.0	$1.9^{\circ}$	0.5	6.5	<u>.</u>
1213		7	ŝ	9	ŝ	9		Ś	ŝ	0	9	4		Ś	ŝ	ŝ	ŝ	4	ξ	2	sks
1214	8	8	Ξ	F	33	50	8	8	89	33	89	51	8	56	8	79	50	39	33	50	tas
1215	AT	25.	52.	4.	58.	69.	75.	50.	38.	33.	67.	34.	75.	55.	4	56.	46.	48.	33.	76.	sso
1216	K																				acr
1217	4	33	=	×	90	0	5	33	2	5	2	Ξ	5	22	2	6	0	9	4	õ	res
1218	ATE	33.3	52.1	58.1	55.5	58.5	73.7	58.3	47.2	41.6	64.2	5.1	90.C	59.6	46.(	56.7	45.5	45.1	4.	77.0	sco
1219	$\mathbf{K}_{\ell}$			-		-	Ì		4	4	-	4	-		7			4	4		ce
1220	5	2	9	6	<u> </u>	0	0	×	0	_		<u> </u>	0	0	0	ŝ	0	S	6	0	nan
1221	ΤE	2.2	0.5	1.5	1.1	6.5	0.0	2.7	0.0	6.1	4.	4.2	5.0	1.4	5.0	1.7	7.5	3.5	8.8	7.5	orn
1223	KA	0	9	( <b>-</b>	9	9	00	Ś	ŝ	e	9	e	-	9	ŝ	9	ব	J	e o	-	Perf
1224		_		_		_	_		_	~	_	~~	- )	_	_		_			_	5: H
1225	RA	5.00	6.6	1.59	5.56	1.0	5.00	5.56	0.0	1.67	2.39	1.58	2.23	26.1	8.00	5.43	3.5(	3.23	1.67	8.5(	e 1
1226	-LA	6	Ň		ŝ		Ë,	ŝ	ñ	4	<del>ن</del>	ξ	È,	9	4	6	4	ŝ	4	Ë	abl
1227	ġ																				
1228	¥	8	č	90	8	0	0	33	6	4	6	Ξ	4	2	2	9	0	0	-	õ	
1229	AR	27.7	5.5	73.8	52.7	53.5	75.0	58.3	38.8	4.4	52.3	27.7	59.4	50.8	<del>1</del> 8.C	59.2	4	50.0	36.1	75.0	
1230	Γ		• •		• •	Ū	Ì	• •		7	Ũ		Ū	Ũ	7		4			Ì	
1231									e								S				
1232									ien							<u>5</u> 0	atic				
1233						ş		>	SC.	tics			k	S		erin	lem			g	
1234		bra			S	led	Ŋ	str	uter	Sma	ine	S	uri	iysi		ine	lath			iolc	
1235	e	lgel		1	thi	NOI	olog	emi	npi	athe	<b>idic</b>	ysid	sec	l_pl	·ics	eng	y_n	ji	s	d_lc	
1236	am	ct_a	λ	Ś	ss_e	L kn	_bid	Ę	<u>[00</u>	т	m	hq-	ter_	tua	netı	cal_	tar.	log	fac	ihot	
1237	k N	trac	ton	JUO.	ine	ical	ege	ege	ege	ege	ege	ege	ndt	cep	non	Ĕ	nen	nal	Dal_	h_SC	
1230	Tas	abs	ana	astı	bus	clin	coll	coll	coll	coll	coll	coll	con	con	eco)	elec	elen	for	glol	higl	
1240	-				_	-	-	-	-	-	-	-	-	-	-	-	-			_	I
19/1																					

1242 1243																					
1244				_			_	_			_										I
1245	C	.55	.78	8.	.81	.15	8.	8.	.26	.38	.50	.53	.14	171	:55	-64	:21	.45	5.7	.92	
1246	Ī	57	22	0	88	8	69	4	71	4	8	8	87	86	8	71	2	2	76	4	
1247	æ	2	9	_	×		0	0	2	×	0	6	6	ŝ	ŝ	6	Ś	ŝ	×	0	
1248	AE	7.5	5.5	9.2	3.5	8.3	6.0	8.0	5.5	6.7	4.0	7.8	4.2	5.5	2.3	2.6	8.9	9.5	7.7	0.0	
1240	R	5	Ś		×	×	9	ε	9	ŝ	×	ŝ	×	×		9				Ś	
1250	4	9	З	0	2	2	0	0	6	Э	0	2	9	6	9	S	0	×	Э	0	
1251	AE	3.9	8.3	0.2	6.5	6.8	3.5	6.0	8.3	0.2	3.0	5.9	2.8	4.3	2.9	0.1	0.7	8.1	2.7	0.0	
1252	R	ŝ	ŝ	×	×	×	9	ŝ	9	4	×	ŝ	×	×			×	9		ŝ	
1253	2	9	ŝ	×	5	0	Q	Q	6	<sub>∞</sub>	Q	Ξ	Q	5	9	2	0	Q	4	2	
1254	AF	3.9	8.3	79.1	36.5	37.6	52.0	36.0	<b>8</b> .3	36.7	33.0	2.7	35.0	2.5	2.9	0.1	S0.7	15.0	13.7	52.1	
1255	R	<b>u</b> )	<b>U</b> )	(~	00	00	U	<b>G</b> 1	U	<b>G</b> )	00	<b>W</b> )	œ	œ	(~	(~	00	(-	(~	<b>U</b> )	
1256	- <u>«</u>	90	ŝ	0	6	0	0	0	5	ŝ	0	8	6	5	Ξ	8	0	×.	4	33	
1257	'AG	53.5	58.3	30.2	32.0	37.6	2.5	39.5	56.6	t0.2	32.0	56.5	2.5	2.5	74.2	2.1	30.7	58.1	73.7	58.3	
1258	Γ	47	47	~	~	~	C		C	7	~	47	~	~		C	~	C		4,	1 7 T
1259	77	24	78	17	58	05	8	8	76	63	50	53	59	39	28	15	20	45	75	4	Pa
1260	'YC	53.	52.	83.	83.	86.0	63.(	39.0	68.	35.0	82.	60.:	84.	84.	73.	70.	80.	70.	4.	00.	B B
1261	Ι																				la-7
1262	$\mathbf{\tilde{G}}$	12	56	19	82	05	50	50	Ξ	63	50	87	8	97	33	16	21	8	74	92	un l
1263	Γ¥	56.	55.	81.	85.	86.	61.	4	70.	35.	83.	59.	85.	<u></u>	5	67.	<u></u>	75.	73.	4	Ge
1264																					Η
1265	CL	.92	.56	:22	.84	.27	00.	8.	.52	.63	8.	.03	.86	.39	.81	.18	.95	.64	.72	8	3B]
1266	Ī	48	55	78	82	85	59	26	65	35	81	48	82	84	69	4	78	63	71	50	in I
1267	æ	0	×	0	2	2	0	0	-	Э	0	2	9	9	4	9	6	S	×	2	ks
1268	TE	8.2	2.7	0.3	5.8	6.8	2.0	3.5	0.1	3.3	5.5	8.5	2.8	2.6	8.	7.1	7.1	0.4	7.7	6.7.	tas
1269	KA	Ā	ŝ	-	œ	œ	0	e)	-	e	œ	ŝ	œ	œ	-	0	~	-	-	ব	ssc
1270							_	_					_	_							acre
1271	$\mathbb{E}_4$	6.83	5.56	.23	82	5.82	00.9	2.50	339	.08	0.1	.24	0.0	.50	.33	5.67	3.95	3.18	.78	3.75	es
1272	(A)	56	5	F	8	8	69	3	68	ő	8	ŝ	8	8	F	65	22	30	F	4	COL
1273	H																				se s
1274	$\mathbf{E}_2$	96	Ξ	53	33	92	8	50	39	03	8	61	86	39	92	15	95	73	E	08	anc
1275	AT	53.	61.	78	<b>8</b> 4	89.	60.	35.	68	31.	<b>8</b> 4.	54.	82	<u></u>	67.	70.	78.	5	76.	52.	
1276	K																				erfc
1277	¥	0	9	0	-	2	Q	Q	-	<sub>∞</sub>	0	6	Q	2	3	S	S	Q	4	Q	P
1278	<b>AR</b>	51.8	5.5	<u>80</u> .2	37.3	36.8	2.0	4.0	0.1	36.7	33.5	57.8	35.0	35.5	2.3	0.1	78.9	15.0	13.7	0.0	16
1279	-T	47	47	~	~	~	C	7			~	47	~	~						4,	ble
1280	H																				Tal
1281	¥	83	78	19	33	82	8	50	82	63	8	61	43	99	33	64	46	45	74	17	
1282	AF.	56.	52.	81.	84.	86.	62.	4	67.	35.	84.	54.	86.	82.	72.	71.	82.	70.	73.	54.	
1283	Τ																				
1284			e	Ŋ			S		s												
1200			ien	sto			mi		nic					N.							
1997		y	r_SC	id	Ŋ	ent	ouo	itic	IOU		S		v	stor							
1288		str	ute	ean	aph	um	)ec	Sma	éco	S	olo	ics	tor	hi							
1289		em	du	rop	1go	ver	JCL	ath	cro	ysid	ych	utist	his	nrld		ţ	aw		c n	ing	
1290		ch	<u> </u>	-eu	e e	8	, m	-m	m	hq-	-ps	St	us.	-wc	ğ	ıali	al_l	JCe	cie	arn	
1291	me	loo	lool	00	lool	lool	lool	lool	lool	lool	lool	lool	lool	lool	iĝ	<b>Sex</b>	ion	den	alla	le	
1292	Na	sch	sch	sch	sch	sch	sch	sch	sch	sch	sch	sch	sch	sch	3m.	5 UE	nati	pru	al f	uine	
1293	ask	gh	gh	da Ha	dg,	ď	ď	ď	dg,	ď	ď	dg,	gh	g	imi	Im	ter	uris	gic	act	
1294	Ë	hi	Ę	hi	ġ	'n	ġ	'n	ġ	ġ	'n	ġ	h	ġ	þ	ц	Е.	.Ē	lo	ш	
1295																					

	12	<u>~</u>	+	•	<u> </u>	<u> </u>	0	0	0	0	0	0	0	0	_	~	0	,0	+	0		
00	ICV	87.18	88.2	88.8	82.00	78.5(	41.5(	78.00	75.5(	69.5(	48.00	46.0(	69.5(	72.5(	73.9	72.9	88.32	80.50	52.9	86.92	64.99	
01	~	-		_	_	_	_	_	_	_	_	_	_	_	~	~	~	~				
J2	<b>≜</b> E <sub>s</sub>	9.72	7.06	5.00	4.5(	5.0	9.0	3.00	2.5(	4.5(	9.00	6.00	9.5(	9.0	9.57	5.69	4.67	1.67	1.96	0.65	4.8(	
о л	R	∞ ∞	ò	Ē	ò	È,	τ,	È.	È,	Ļ,	4	4	ŝ	6	6	Ë,	ò	6	S	6	é	
	4	4	0	0	0	0	0	0	0	0	0	0	0	0	6	S	4		8	×	1	
	AE	9.7	0.0	2.2	7.0	2.0	H.0	6.0	2.5	3.5	0.5	1.0	2.0	0.0	57.3	7.3	3.9	1.6	5.8	5.9	5.0	
	R		5	(~	œ	(~	J	(~	(~	(-	<b>W</b> )	<b>v</b> )	Q	(-	C	(-	œ	5	<b>v</b> )	æ	ę	
	5	31	8	53	20	20	20	20	8	8	50	8	50	8	39	45	13	89	6	85	18	
	[Y]	89.	87.	72.	<b>2</b>	72.	35.	78.	73.	73.	50.	51.	65.	1.	67.	Ë	86.	88.	54.	87.	65.	
	-																					
	<sup>8</sup>	62	8	22	8	8	.50	8	.50	8.	.50	.50	8.	8.	.74	.03	4.	.89	.92	.85	.04	
	LA	8	8	72	8	73	4	F	71	73	46	48	2	69	71	74	85	88	53	87	65	
		Ι.			_	_	_	_	_	_	_	_	_	_	_	_		_	_			
	$G_4$	.74	.41	22	<u> </u>	3.00	0.0	7.50	2.50	3.50	2.50	1.50	5.50	0.0	.57	5.69	5.86	8.89	96. <del>1</del>	7.85	5.87	
	$\mathbf{L}^{\prime}$	∞	Š.	Ê	õ	F	4	ĥ	F	F	ŝ	Ň	9	Ř	6	Ľ	×	õ	Ň	ò	6;	
	5		4	4	0	0	0	0	0	0	0	0	0	0		0	3		2	6	1	
	AG	2.3	8.2	4.6	4.0	3.0	0.5	5.0	4.0	4.5	4.0	4.5	8.5	0.6	9.5	7.9	6.1.	1.6	3.9	8.7	5.6	
	Γ	6	×	9	×		4				ŝ	ŝ	9	9	9		×	6	Ś	×	9	
	Г	2	-	Q	Q	0	Q	Q	Q	0	0	Q	Q	0	6	6	ò	3	21	×	4	
	IC	34.6	2.7	75.0	32.0	55.5	39.0	70.0	9.6	70.5	17.5	2.0	57.0	57.5	57.3	75.6	32.4	33.3	53.9	35.9	51.4	
			~		~	U			Ŭ		7		4,	C	U		~	~	47	~		
	$\mathbf{E}_8$	<i>T</i> 9	4.	.78	8	.50	8.	.50	8	.50	8	8	.50	8.	8.	.34	.96	.78	.16	.45	.13	
	IAT	71	86	77	81	20	32	33	2	26	31	2	20	35	50	19	37	27	4	36	54	
	¥																					
	$\mathbf{E}_4$	64	4	22	50	8	8	8	50	8	50	8	50	50	65	34	31	89	12	60	78	
	AT	79.	88.	72.	84.	69.	25.	37.	20.	33.	35.	22.	20.	40.	45.	19.	34.	38.	4	48.	55.	
	K																					
	5	17	3	4	0	8	8	00	0	8	0	00	0	2	8	76	66	Ξ	ß	73	58	
	II	92.3	87.6	e9. <sup>2</sup>	83.1	69.(	25.(	39.5	27	35.(	30.5	22.7	20.1	37.(	47.8	23.7	38.0	36.	37.2	46.7	56.2	
	K																					
	-	-	ŝ	4	0	0	0	0	0	0	0	0	0	0	2	S	6	4	2	2	6	
	AR.	2.3	37.6	9.4	3.5	3.5	5.5	5.0	3.5	3.0	2.5	3.0	6.0	9.0	5.2	7.3	7.5	4.4	3.9	6.9	5.5	
	· <b>I</b>	5	œ	Ð	œ	(~	<b>G</b> )	(~	(~	(-	<b>W</b> )	<b>v</b> )	Q	Û	C	(-	œ	5	<b>v</b> )	æ	ę	
	B																					
	A	2	4	5	50	8	20	8	8	50	50	8	50	8	57	80	4	89	8	98	36	
	AR.	89.	86.	72.	82.	71.	30.	4	5.2	71.	52.	53.	68.	66.	69.	76.	85.	88.	40.	85.	64.	
	Π																					
											gu			<u>y</u> s								
											ntiı		ine	olo								
				S			s				cou	A	edic	ych	<b>د</b>			Ś				
		t l		hetik	ns	Ites	urio				ac	Llav	ñ	_ps	ion	die		olii		ons		
	me	nen	g	gen	neo	ispu	ent	_	hy	ŗ	na	na	na	na	elat	stu	4	<u>n p</u>		iligi		
	Naj	Iger	etin	cal	lla	I_di	l_sc	tior	sop	stol	ssic	ssic	ssic	ssic	C_r	ity_	log l	reig	ğ	l_re	ge	
	ìsk	ana	ark	edi	isce	ora	ora	Itri	lilo	·ehi	ofe.	.ofe	ofe.	ofe.	ildi	cur	cio	fo	rolc	orlc	'era	
	Ë	E	H	Ш	H	ш	Ξ	III	þ	Id	Id	Id	Id	Id	īd	Se	SO	ŝ	<u>v</u>	M	aı	

Table 17: Performance scores across tasks in MMLU (Mistral-7B) Part 3

25

Under review	as a conference	paper at ICLR 2025
		T T T T T

1352																				
1353																				
1354	~		~	6	Ś		6	ŝ		_	~		~	0	~	6		~	m	
1355	5	5.5	1.0	5.6	7.9	7.8	4 0	0.5	8.8	1.5	3.1	1.4	3.1	9.2	4.6	6.4	%	8.9	5.9	
1356		9	Ś	Ś	0	ŝ	4	9	0	0	0	4		ŝ	4	×	ŝ	4	4	
1357	8	2	5	00	8	6	6	6	2	4	80	96	8	57	2	22	8	8	9	
1358	<b>IA</b>	19.7	59.1	55.5	30.2	46.4	38.9	59.5	48.	12.8	14.6	55.9	72.4	40.5	43.]	56.8	4.0	43.5	14	
1359	H				` '	4	` '	•••	1				`	4	4	-	1	4		
1360	Н 4	12	26	50	32	72	37	01	69	30	51	80	5	5	50	81	20	8	53	
1361	<b>V</b>	19.	58.	55.	35.	37.	40.	63.	47.	13.	16.	57.	4.	39.	4	12.	4	50.	45.	
1362	H																			
1363	$\mathbf{E}_2$	72	60	92	32	4	Ξ	59	76	30	43	88	60	20	12	89	01	29	59	
1364	RA	19.	60.	50.	35.	39.	32.	59.	46.	13.	17.	56.	78.	<u>4</u>	43.	71.	38.	52.	4.	
1365																				
1366	ß	.12	.32	.56	.33	00.	.17	.81	29	.68	.81	.55	21	.74	.58	.80	.46	.88	0.02	B
1367	LA	×	39	52	33	50	40	4	51	13	15	58	75	4	46	58	38	4	43	1-8
1368																				a3.
1369	$G_4$	24	9.	.26	.47	0.0	.72	.35	.98	.25	.38	.26	.64	.31	.01	.52	.31	.71	.50	ama
1370	LA	16	35	57	35	4	48	3	4	13	$\frac{18}{18}$	57	55	4	4	99	4	51	4	Ē
1371											_				_					H
1372	₽G2	7.52	1.45	.85	3.76	5.38	3.21	5.56	3.28	283	20.7	).26	7.78	2.74	5.30	5.52	.32	3.85	3.33	BB
1373	LA	12	4	5(	3	4	8	ŝ	4	2	5	90	F	4	4	90	č	ŝ	4	Е.
1374	. 1	+	6	~	<del>. +</del>		~	<del>. +</del>	ŝ		<del>. +</del>	0	<del>. +</del>	$\sim$	Ś	5		5	+	sks
1375	CI	0.6	2.3	6.42	3.92	0.0	2.5	$1.6^{2}$	4.6	4.1	5.12	7.2	2.8	3.58	7.10	9.2(	8.0	7.16	5.64	s tag
1376		6	<u>ن</u>	Ñ	ŝ	Ñ	Ó	9	Ŵ	Ξ	÷	4	9	4	ŝ	Ē	ŝ	ŝ	4	SOS
1377	8	4	33	4	57	4	2	33	60	93	4	2	4	66	16	<del>1</del> 3	20	80	4	acı
1378	Ĩ	15.	61.6	56.4	32.5	56.	84	62.	47.0	Ξ	12.5	50.0	12.6	38.	37.	71.4	36.	43.1	4	res
1379	K						-		-									-		SCO
1201	4	2	4		ŝ	Ś	×	_		0	×	ŝ	_	ŝ	_	2	4		6	ce
1292	ΤE	4.2	2.8	4.1	3.0	1.7	9.0	3.0	1.6	3.3	4.6	2.7	8.8	8.5	7.6	9.1	6.2	0.3	4.0	nan
1382	KA	-	9	9	ŝ	ŝ	4	9	4	Ξ	-	ŝ	9	ŝ	ŝ	9	ω	ব	4	orn
1384																				erf
1385	$\mathbb{E}_2$	.64	90.	.83	.65	.75	.71	0.	.30	.47	.47	.50	.48	.45	.37	.13	6.	.91	.60	Ч.
1386	E	8	65	51	31	51	4	63	46	Π	11	55	5	39	4	3	82	39	4	18
1387	¥																			ble
1388	Y	76	3	80	36	8	91	38	20	01	90	93	3	95	4	58	95	83	69	Ta
1389	AR	30.	72.0	57.	29.	50.0	39.0		53.	11.0	16.	61.9	72.0	4	45.	75.	4	40.	47.0	
1390	B-L																			
1391																				
1392	R	10	.05	45	.61	.12	5	81	63	68	51	63	:23	.58	41	43	.95	88	.46	
1393	[Y]	21	55	56	26	49	4	67	54	4	16	59	75	43	45	71	4	56	47	
1394	Ι																			
1395																	)et			
1396	e		_		m		oes		S	~	5		n				Έ	j.		
1397	am	ba	ğ	pı	)j	ble	hal		ame	jį	įq	13	oatc	15	4	čec	nsE	Q	e	
1398	N	ŊS	ımt	eUr	ckC	gla	Sm	rks	Ž	ckC	ckC	Dec	hert	Dec	Dec	vieł	<b>Ura</b> ı	Š	rag	
1399	Lasl	Tem	)is£	Dat	l'ra	Pen	3e0	<b>Jna</b>	Zui.	Tra	l'ra	g	<b>Hyp</b>	g	g	Mo	Sal	Rea	<b>Ive</b>	
1400			-	Ι		-	-		-	<u> </u>	_	-	-	Γ	Γ	2			4	
1401																				

1404 1405																					
1400	I																				I
1407	C	00.	2	.18	Ξ.	.50	.75	8.	0.0	.56	57	5.7	4.	.48	8	-07	.50	.55	.33	.50	
1400	Ι	25	5	68	61	55	25	SC SC	50	30	2	4	66	58	52	74	4	4	б.	81	
1409	×	Ļ	×	0		0	0	0	×	2	0	4	9	0	0	~	0	~	6	0	
1410	AE	6.1	33.3	5.0	<u>6</u> .6	7.5	7.5	1.2	2.7	1.2	2.7	7.7	0.5	1.4	4.0	0.3	ŀ6.0	ŀ6.7	8.8	31.0	
1/10	R	01	0	( -	0		( -	7	47	7	C	7	$\sim$	0	4,	( -	7	7	0.1	$\sim$	
1413	4	52	38	73	11	8	52	8	2	2	30	4	78	22	8	14	8	32	20	20	
1414	<b>EAI</b>	47.	63.	2	61.	79.(	76.	50.0	4.	4.	63	39.4	Ľ.	60.5	48.(	.69	47.(	40	30.1	81.1	
1415	H																				
1416	$\mathbf{E}_2$	89	38	8	56	50	50	8	8	89	30	84	56	65	8	14	50	8	11	50	
1417	RA	38.	63.	75.	55.	5	77.	50.	50.	38.	63.	36.	80.	59.	46.	69.	46.	50.	36.	81.	
1418																					
1419	$\mathbf{G}_{\mathbf{s}}$	.11	20	.86	.56	.50	8.	:22	4	:22	.47	.7	8	.89	8	.90	.50	.84	Ξ.	.50	<b>1</b>
1420	ΓA	36	99	73	55	78	75	47	4	47	61	4	75	57	52	67	45	54	36	80	Par
1421		_		_	_	_	_			_	_			_	_		_		_	_	B
1422	$G_4$	3.89	3.38	00.0	3.89	<u>8</u> .00	7.50	.22	7.22	0.0	3.30	.47	7.78	66.1	<u>5.00</u>	9.14	9.00	3.55	8.89	).50	1-8
1423	LA	38	69	5	69	22	F	4	4	50	69	ŝ	5	6	Š	66	4	4	38	80	a3.
1424	~	2	$\sim$	<del>. +</del>	5				<del>. +</del>	<del></del>	~	<del>. +</del>	5	5		<del>. +</del>		~	~		am
1425	₽Ğ	1.67	3.38	6.12	5.56	8.5(	8.75	0.0	4.	4.	4.	4.7	0.56	9.06	2.00	9.12	7.0(	6.7.	8.8	1.00	E
1426	$\Gamma$	4	6	7	ŝ	ř	F	Ñ	4	4	9	4	õ	ŝ	ŝ	6	4	4	õ	×	Ę
1427	L	5	0	2	×	0	0	0	Э	×	0	0	0	×	0	ŝ	0	0	9	0	ΙĮ
1428	IC	2.2	6.2	9.3	L.L.	15.5	7.5	0.0	8.3	L.L.	33.3	0.0	2.2	8.4	4.0	5.4	3.5	7.1	0.5	8.0	Ξ
1429		(1	Q	Q	(1	(~	(-	<b>W</b> )	<b>W</b> )	(1	Ο	<b>v</b> )	(~	<b>W</b> )	<b>W</b> )	U	7	<b>G</b> 1	<b>G</b> )	(~	s in
1430	$\mathbf{E}_8$	11	67.	59	4	8	25	4	22	8	30	58	22	14	8	67	8	61	4	00	ask
1431	AT	36.	4	71.	69	76.	76.	4	4	50.	63.	31.	72	56.	<del>4</del> 8	66.	43.	51.	4	79.	ss ti
1432	K																				CT OS
1433	4	0	5	13	Ξ	00	15	4	2	5	33	4	2	8	2	00	2	22	5	õ	s ac
1434	<b>T</b> I	25.(	51.5	72.7	61.1	78.5	78.7	4.	50.0	41.6	59.6	4	75.(	58.4	50.0	71.6	4	40.5	41.6	80.0	ore
1435	$\mathbf{K}$		-		-	-	-						-			-					sc
1436	0	8	6	Ś	6	0	5	x	×	~	0	4	0	9	0	4	0	0	ŝ	0	nce
1437	ΤE	T.T	4	0.4	3.8	7.5	8.7	2.7	2.7	7.2	3.3	4.7	5.0	9.0	6.0	9.1	3.0	0.0	3.3	0.0	ma
1430	KA	0	9		9			ŝ	ŝ	4	9	4		ŝ	4	9	4	ŝ	ŝ	×	for
1440			_		_	_	_		_	_	_		_		_		_			_	Pei
1441	RA	5.11	5.20	2.73	3.89	00.7	7.50	.22	00.0	00.0	2.39	.47	00.0	).23	5.00	t.07	9.00	t.84	5.11	).50	6
1442	ΓA	36	Ğ	F	6	F	ŀ	4	S(	S.	6	č	Ļ	6	ñ	Ļ	4	Š	õ	8	le l
1443	ġ																				Lab
1444	-	3	8	6	9	0	S	0	×	2	8	<del></del>	0	Э	0	9	0	_	9	0	
1445	<b>AR</b>	3.3	33.3	1.5	5.5	15.5	6.2	0.0	2.7	H.6	6.8	5.1	5.0	0.2	-8.0	2.9	4.0	51.6	0.5	8.5	
1446	$\Gamma$	с,	Q	(-	<b>W</b> )	(~	(~	<b>W</b> )	<b>W</b> )	7	<b>v</b> )	7	(~	Q	7	U	7	<b>W</b> )	<b>G</b> )	(~	
1447									e								s				
1448									enc							50	atic				
1449						ē			Sci	tics			v	S		rin	em			<u>y</u> g	
1450		ra			s	edg	~	stry	lter	mai	ne	s	<b>irit</b>	isi		nee	ath			iolo	
1451		geb			hic	0wl	log	imi	ndu	the	dici	vsic	Seci	-hd	ics	igu	<b>m</b> _7	ં	s	l_b	
1452	ame	t_al	v	my	s_et	Ŕ	bio	che	COL	ma	me	hd	er	ual	etr	al_e	ary	log	act	h00	
1453	Ž	rac	Om	ono	nes	cal	je je	ğ	<u>s</u> e	- Se-	ş.	ŝ	put	čept	nom	tric	lent	lal_	ali	SC	
1454	lash	bst	nat	str	isn	lini	olle	olle	olle	olle	olle	olle	Om	OUC	COL	leci	len	orn	lob	ligh	
1455			c3	ವ	2	J	J	3	J	J	J	J	J	J	e	e	e	÷	<b>CI</b> )	4	
1450																					
1407																					

e LARA B-LARA KATE <sub>2</sub> KATE <sub>4</sub> KATE <sub>8</sub> ICL LAG <sub>2</sub> LAG <sub>4</sub> LAG <sub>8</sub> RAE <sub>2</sub> RAE <sub>4</sub> RA 01-chemistry 57.55 61.87 60.43 57.55 61.87 58.27 56.83 59.71 56.83 56.12 0.00 C	e LARA B-LARA KATE2 KATE4 KATE8 ICL LAG2 LAG4 LAG8 RAE2 RAE4 RAF	<b>J.chemistry</b> 57.55 61.87 60.43 57.55 61.87 58.27 56.83 59.71 56.83 56.12 0.00 0.00 <b>J.computer_science</b> 63.89 61.11 63.89 72.22 58.33 58.33 61.11 58.33 63.89 61.11 63.89 61.11	ol_european history 78.22 77.23 75.25 74.26 75.25 77.23 77.23 77.23 76.24 77.23 78.22 77.23	ol_geography 85.82 88.81 83.58 88.06 86.57 87.31 84.33 85.82 85.07 81.34 85.82 86.57	ol_government 88.37 85.27 90.70 89.15 86.05 86.05 86.82 86.82 88.37 86.05 86.05 86.05	<b>J_macroeconomics</b> 69.00 69.50 63.50 67.50 70.50 70.00 70.50 68.50 68.00 68.00 68.50	<b>J_mathematics</b> 36.50 43.50 42.00 39.00 42.50 39.50 44.50 45.00 41.00 44.50 42.00 41.50	ol_microeconomics 71.26 75.29 70.11 71.26 72.99 72.41 71.26 72.99 74.71 72.41 74.14 75.86	ol_physics 48.28 52.87 49.43 45.98 49.43 43.68 50.57 50.57 45.98 48.28 49.43 48.28	Jl.psychology 86.50 84.50 86.50 87.00 87.00 84.50 86.50 86.50 87.00 87.00 84.50 87.00	ol_statistics 52.63 54.61 55.92 57.89 54.61 55.92 52.63 53.29 53.29 54.61 52.63 53.95	J.us.history 85.71 85.71 84.29 87.86 88.57 84.29 85.71 87.14 85.71 85.71 86.43 87.86	Juworld history 85.55 83.82 84.39 87.28 84.39 86.13 86.71 85.55 85.55 85.55 86.13 83.24	ing 71.70 71.70 68.55 69.18 71.70 70.44 69.18 69.8 69.81 69.81 67.30 70.44 72.33	vuality 76.12 76.12 76.12 76.12 79.10 70.15 77.61 76.12 79.10 76.12 74.63	nal.law 78.95 82.46 82.46 85.96 85.96 84.21 80.70 84.21 85.96 78.95 85.96 80.70	ence 63.64 65.91 72.73 68.18 68.18 72.73 70.45 65.91 68.18 70.45 68.18 68.18	lactes 77.78 78.79 79.80 78.79 78.79 78.79 78.79 79.80 75.76 78.79 77.78 79.80	earning 43.75 39.58 37.50 45.83 43.75 47.92 37.50 37.50 43.75 43.75 45.83
high solved commuter solons	high_school_chemistry	Inglise roundings induced selection	high school european histor	high_school_geography	high_school_government	high_school_macroeconomic	high_school_mathematics	high_school_microeconomics	high_school_physics	high_school_psychology	high_school_statistics	high_school_us_history	high_school_world_history	human_aging	human_sexuality	international_law	jurisprudence	logical_fallacies	machine_learning

) Part 2	
I-8B	
(Llama3.1	
MMLU	
-8	
tasks	
across	
scores	
Performance	
20:	
able	

1556 1557 1558 1559 1560 1561 1562 1563 1564 1565	1553 1554 1555	1549 1550 1551 1552	1545 1546 1547 1548	1542 1543 1544	1538 1539 1540	1535 1536 1537	1532 1533 1534	1528 1529 1530 1531	1525 1526 1527	1522 1523 1524	1519 1520 1521	1516 1517 1518	1513 1514 1515	1512
Task Name	LARA	<b>B-LARA</b>	$\mathbf{KATE}_2$	$\mathbf{KATE}_4$	<b>KATE</b> <sub>8</sub>	ICL	$\mathbf{LAG}_2$	$\mathbf{LAG}_4$	$\mathbf{LAG}_{8}$	$\mathbf{RAE}_2$	$\mathbf{RAE}_4$	$\mathbf{RAE}_8$	ICV	
management	89.74	87.18	79.49	82.05	82.05	89.74	84.62	87.18	89.74	84.62	87.18	89.74	89.74	
marketing	92.94	90.59	91.18	90.59	91.18	88.24	89.41	90.00	92.94	90.00	90.59	90.59	90.59	
medical genetics	88.89	88.89	88.89	88.89	86.11	86.11	88.89	88.89	88.89	88.89	88.89	88.89	88.89	
miscellaneous	81.50	82.00	80.50	80.50	84.50	83.50	82.00	83.00	82.50	82.00	82.00	82.50	83.50	
moral_disputes	74.50	76.00	75.00	77.00	77.50	74.00	74.50	75.00	75.50	74.00	75.50	76.00	76.50	
moral_scenarios	46.50	46.00	34.00	38.00	40.50	40.50	37.50	41.50	49.00	40.00	46.00	45.50	41.50	
nutrition	79.00	80.50	77.00	78.00	77.50	76.50	79.00	80.50	79.50	79.00	77.00	79.50	77.50	
philosophy	72.50	71.00	74.00	72.50	73.00	74.00	72.50	72.00	72.50	74.00	75.00	75.00	73.00	
prehistory	70.00	68.50	69.00	70.50	72.50	68.50	69.50	70.00	70.00	68.00	67.00	71.00	69.50	
professional_accounting	50.00	48.50	48.00	48.50	51.00	53.00	50.00	50.50	47.50	50.50	50.50	48.50	47.50	
professional law	53.00	55.50	53.50	54.50	52.50	46.50	52.50	53.50	56.00	54.50	53.00	55.50	52.00	
professional_medicine	65.00	65.00	71.00	67.00	67.50	71.00	68.50	67.00	66.50	66.50	67.00	65.00	68.50	
professional_psychology	70.50	72.50	72.00	73.50	76.00	00.69	72.50	74.00	74.00	73.00	73.50	73.50	71.00	
public_relations	73.91	76.09	80.43	78.26	78.26	76.09	73.91	76.09	76.09	73.91	76.09	78.26	78.26	
security_studies	74.59	76.24	74.59	75.14	76.80	75.69	77.90	76.80	75.14	76.24	77.35	75.69	74.59	
sociology	91.97	91.24	89.05	87.59	88.32	88.32	91.24	91.97	91.97	91.24	91.97	91.97	88.32	
us_foreign_policy	88.89	88.89	77.78	77.78	83.33	88.89	88.89	86.11	88.89	86.11	88.89	88.89	86.11	
virology	53.92	51.96	51.96	52.94	52.94	51.96	52.94	53.92	51.96	52.94	52.94	52.94	51.96	
world_religions	84.11	84.11	83.18	85.05	84.11	84.11	84.11	84.11	84.11	84.11	83.18	85.05	85.05	
average	66.54	67.80	66.62	66.75	67.19	65.64	66.71	66.88	69.99	66.22	66.91	67.24	67.10	
		Table 21:	Performan	ice scores a	cross task:	IMM ui s	LU (Ilama)	3.1-8B) P	art 3					

(llama3.1-8B) Part 3	
MMLU	
tasks in	
across	
scores	
Performance	
le 21:	