# CASE: CHALLENGER ARM SAMPLING FOR EFFICIENT IN-CONTEXT REASONING

Anonymous authors

Paper under double-blind review

### ABSTRACT

The in-context learning paradigm with LLMs has been instrumental in advancing applications that require complex reasoning over natural language. An optimal selection of few-shot examples (exemplars) is essential for constructing effective prompts under a limited budget. In this paper, we frame the problem of exemplar selection for In-Context Reasoning (ICR) as a top-*m* best arms identification problem. A key challenge in this context is the exponentially large number of arms that need to be evaluated to identify the *m*-best arms. We propose CASE (Challenger Arm Sampling for Exemplar selection), a novel *selective exploration* strategy that maintains a shortlist of "challenger" arms, which are current candidates for the top-m arms. In each iteration, only the arms from this shortlist and the current top-m set are pulled, thereby reducing sample complexity and, consequently, the number of LLM evaluations. Furthermore, we model the scores of exemplar subsets (arms) using a parameterized linear scoring function, leading to a *stochastic linear bandits* setting. In this setting, CASE identifies the top-*m* arms with significantly fewer evaluations than existing state-of-the-art methods. CASE effectively works with black box LLMs and selects a static set of few-shot examples, resulting in an extremely efficient scheme for in-context reasoning. The exemplars selected with CASE show surprising performance gains of up to 15.19% compared to state-of-the-art exemplar selection methods. We release our **code and data**<sup>1</sup>.

### 028 029

031

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

### 1 INTRODUCTION

In-context learning (ICL) and Chain-of-Thought (COT) have emerged as important techniques for
 enhancing the capabilities of large language models (LLMs) across a range of tasks like question
 answering and complex reasoning. ICL allows LLMs to perform tasks by conditioning on a context
 that includes examples or instructions, without the need for additional fine-tuning, making it flexible
 and adaptable. COT facilitates complex reasoning by decomposing tasks into intermediate steps, and
 employing rationales explaining the reasoning process to the LLM. However, one key challenge in
 maximizing the effectiveness of ICL is the careful selection of few-shot examples along with their
 corresponding rationales (Lu et al., 2022; Zhao et al., 2021).

040 The inclusion of rationales alongside standard training examples has been found to be especially 041 important for complex reasoning tasks, such as multistep reasoning-based QA (Geva et al., 2021b; 042 Chen et al., 2022b; Lu et al., 2023b). We refer to the triplet of (input, rationale, output) as an exemplar 043 and define the paradigm of using rationale-augmented examples in ICL as In-Context Reasoning (ICR). The majority of existing approaches for exemplar selection rely on heuristics or trial-and-error 044 methods Fu et al. (2023); Brown et al. (2020a), with only a few attempting to address the problem in 045 a more principled way Xiong et al. (2024). Exemplars selected for ICR can be classified into two 046 categories: task-level, where a static set of exemplars representative of the task is chosen for inference, 047 and *instance-level*, where exemplars are dynamically selected for each test instance during inference. 048 Instance-level selection approaches typically use similarity (Rubin et al., 2022; Xiong et al., 2024) 049 and diversity-based measures (Ye et al., 2023b) to identify the most suitable exemplars from the training pool for each test instance, which introduces overhead during inference. Additionally, these 051 approaches do not consider the positive and negative interactions between the exemplars, as they 052 select each exemplar independently. In contrast, selecting a static set of exemplars at the task level not

<sup>&</sup>lt;sup>1</sup>https://anonymous.4open.science/r/CASE\_exemplar\_bandits-7403



Figure 1: Overview of CASE for selection of top-*m* best exemplar subsets (arms). The process begins with clustering exemplars from the training set, followed by sampling with replacement to form exemplar subsets (S). In each iteration, the top-*m* arms ( $U_t$ ) and challenger arms ( $N_t$ ) are computed from a uniformly sampled set. Each arm pull corresponds to an LLM evaluation using the exemplar subset. The process repeats until the stopping criterion is met.

066

067

068

071

only eliminates inference-time overhead but also enables prompt caching, allowing for the reuse of
key-value (KV) attention states Gim et al. (2024). In this work, we propose a principled approach for
task-level exemplar selection. The closest work that proposes a task-level exemplar approach is LENS
Li & Qiu (2023), which incurs high computational costs by valuing each exemplar in the training pool
using an informativeness measure that requires LLM call over multiple iterations. Additionally, LENS
relies on confidence measures from the LLM, making it difficult to extend to black-box models, and
does not account for interactions between exemplars. In contrast, we implicitly capture interactions
by scoring exemplar subsets as a whole rather than evaluating them independently.

080 One of the challenges in proposing a principled solution is the limited understanding of the underlying 081 mechanisms of in-context learning. Recent studies aimed at developing theoretically grounded 082 models for in-context learning (Zhang et al., 2023a) primarily focus on linear functions and linear 083 transformers. Therefore, to achieve a principled and efficient selection of exemplars, we require a surrogate model for modeling the goodness of an in-context learning procedure. In this paper, we 084 use a linear function based on sentence similarities between the in-context examples and validation 085 examples as the surrogate model for the goodness score. This approach also be interpreted as a reward model for a multi-armed bandit-based exemplar selection scheme, leading to the linear stochastic 087 bandits setting (Abbasi-Yadkori et al., 2011). 088

We formulate the selection of top-m exemplar subsets as the problem of identification of top-m arms (Réda et al., 2021) in the stochastic linear bandit setting (Réda et al., 2021). The GIFA framework 090 (Réda et al., 2021) proposed an efficient gap-index based top-m arms identification algorithm with 091 reduced sample complexity. However, the computation of *challenger arm* (current candidates for 092 top-m arms), requires computation of gap-indices between all currently estimated top-m arms and the remaining arms. This is impractical in our setting since each arm corresponds to a k-sized subset 094 of the training exemplar set, leading to an exponential number of candidate arms. Hence we need an 095 algorithm that can sample arms from the candidate sets. While some uniform sampling algorithms 096 for top-m arms identification exist, Chen et al. (2017); Kaufmann & Kalyanakrishnan (2013) for the general multi-arm bandit setting, to the best of our knowledge, there are no sampling-based 098 algorithms for identification of top-m arms in the stochastic linear bandits setting.

099 In this work, we propose CASE (Challenger Arm Sampling for Exemplar Selection) where we 100 propose a principled sampling of challenger arms to form a shortlist challenger set, pruning the space 101 of possible candidate arms (see fig. 1). Our key idea is to iteratively create a low-regret set of selected 102 challenger arms, in addition to the current top-m arms, from uniformly sampled arms. This leads to a 103 selective exploration-based algorithm. We concurrently apply the state-of-the-art gap-index-based 104 algorithm rule for selecting top-m arms out of the total exploration set. We also provide theoretical 105 arguments to justify our novel approach of combined selective exploration and gap-index based identification of top-m arms. When applied to exemplar selection, we observe improvements in task 106 performance of upto 15.19% and reduces number of LLM calls (about 10.5x) when compared to 107 state-of-the-art exemplar selection approaches. To summarize, our key contributions are:

- We propose a novel and principled gap index, MAB based approach for task level exemplar selection. Our contribution primarily lies in principled challenger arm sampling to prune the search space, which renders the process highly efficient compared to state-of-the-art MAB frameworks and other subset, exemplar selection approaches.
  - We demonstrate that CASE has significantly lower number of comparisons, gap-index computations and runtime compared to state-of-the-art gap-index based algorithms.
  - We perform extensive experiments on diverse tasks that require complex reasoning and observe that CASE is sample efficient and yields high gains compared to other state-of-the-art exemplar selection approaches.
- 118 119 2 RELATED WORK

109

110

111

112

113

114

115

116

117

120

**Exemplar Selection for ICL.** The rise of LLMs has transformed them into general-purpose answer-121 ing engines through emergent capabilities like ICL (Brown et al., 2020b; Wei et al., 2022; 2023; Wang 122 et al., 2023a; Kojima et al., 2023; Chen et al., 2022a) where a few examples are provided to LLMs 123 to demonstrate the task. To eliminate manual selection, several automated methods have emerged, 124 such as reinforcement learning (Zhang et al., 2022; Lu et al., 2023a), trained retrievers Xiong et al. 125 (2024), Determinantal Point Processes (Ye et al., 2023a) and constrained optimization (Tonglet 126 et al., 2023), which are effective for reasoning tasks. Additionally, dynamic selection methods that 127 are learning-free, such as similarity-based (Rubin et al., 2022), complexity-based (Fu et al., 2023), 128 and MMR (Ye et al., 2023b), have been explored. However, dynamic selection methods increase 129 inference-time computational costs. To address this, a pre-selected, representative set of exemplars can be chosen for ICL, akin to coreset selection methods (Guo et al., 2022), though the key difference 130 131 is that ICL does not involve parameter updates, unlike traditional deep learning training. To the best of our knowledge, there has been very little research in a principled, efficient approach for task-level 132 exemplar selection, with the closest work being LENS (Li & Qiu, 2023), which depends on LLM 133 output probabilities, is expensive in terms of number of LLM calls and is unsuitable for black-box 134 models. We propose a principled novel static exemplar selection method grounded in gap-index based 135 stochastic linear bandits that works for both black-box and open LLMs. 136

137 Identification of Top-m Arms in Stochastic Linear Bandits. The top-m arms identification problem aims to estimate a subset of m arms with the highest means. Various methods have 138 been proposed for top-m arm identification in both fixed-confidence (Kalyanakrishnan et al., 2012) 139 and fixed-budget settings (Bubeck et al., 2013). In this paper, we focus on the fixed-confidence 140 setting, where the error probability in estimating the top-m arms should be smaller than a predefined 141 parameter  $\delta \in (0, 1)$ . Adaptive sampling algorithms like UGapE (Gabillon et al., 2012) and LUCB 142 (Kalyanakrishnan et al., 2012), along with uniform sampling methods (Kaufmann & Kalyanakrishnan, 143 2013; Chen et al., 2017), have been introduced for the fixed confidence setup, but they lack efficiency 144 in terms of sample complexity. While efficient adaptive sampling methods for linear bandits, such 145 as Fiez et al. (2019), RAGE Zhang et al. (2023a), LTS Jedra & Proutiere (2020), PEPS Li et al. 146 (2023), LinGapE (Xu et al., 2017) and LinGame Degenne et al. (2020), have been proposed, they primarily address best-arm identification (m = 1). To the best of our knowledge, GIFA (Réda 147 et al., 2021) was the first unified framework for efficient top-m arm identification with low sample 148 complexity. However, algorithms implemented in GIFA framework require large number of gap-index 149 computations and comparisons, leading to high sample complexity. This is due to its challenger 150 arm sampling mechanism, which considers the complement of the current top-m estimate as the 151 challenger set, resulting in a large search space. In this work, we propose a novel and efficient 152 algorithm with applications to exemplar selection for complex reasoning in LLMs. 153

- 154
- 155
- 156 157

### 3 CASE: CHALLENGER ARM SAMPLING-BASED EXPLORATION FOR EXEMPLAR-SUBSET SELECTION

In-context learning (ICL) leverages LLMs to acquire task-specific knowledge from just a few examples, typically structured as *input-rationale-output* triplets, without updating the model's parameter.
 The rationales serve to elucidate the reasoning process behind the input-output pair, fostering reasoning capabilities within LLMs. However, due to the financial and computational costs associated with processing large contexts, providing all training examples is impractical. Therefore, a key challenge

lies in the efficient and optimal selection of exemplars, particularly in a black-box setting where
 access to the model's parameters or confidence estimates is unavailable.

To address this challenge, we propose a novel algorithm for exemplar subset selection. First, within 165 the in-context learning paradigm, the task of constructing an optimal prompt is framed as a multiple 166 exemplar subset selection problem (section 3.1). In section 3.2, we model this as a top-m arm 167 selection problem in stochastic linear bandits (Kalyanakrishnan & Stone, 2010; Réda et al., 2021), 168 where a linear function is employed to approximate the loss of each arm (prompt). Section 3.3, 169 introduces a new algorithm based on *selective exploration*—conceptually related to SETC (Lattimore 170 & Szepesvári, 2020)), which samples and stores the "challenger arms" (arms that have the potential 171 to belong to the top-m set (Réda et al., 2021)). This approach significantly reduces the computational 172 overhead by avoiding the need to evaluate all possible arms in every iteration.

173 174

175

### 3.1 MULTIPLE EXEMPLAR-SUBSET SELECTION FOR IN-CONTEXT REASONING

In-context learning (ICL) allows LLMs to be used for supervised learning tasks by leveraging a few 176 demonstration examples (called exemplars) of the task without training or parameter fine-tuning. 177 Unlike traditional supervised learning, where training examples include the input features u and output 178 labels v, exemplars sometime include chain of thought demonstrations, which are not available for test 179 examples. For simplicity, we include them in u for the training exemplars. Let  $\mathcal{X} = \{u_i, v_i\}_{i=1}^n$  be the set of all n potential training exemplars and  $(u_{test}, v_{test})$  denote a test input features and desired 181 output. A prompt P is constructed from a subset  $S \subseteq \mathcal{X}$  of k exemplars, which can be used for 182 prediction of  $u_{test}$ . Hence,  $P = [S, u_{test}] = [(u_{i_1}, v_{i_1}), ..., (u_{i_k}, v_{i_k}), u_{test}]$ . The prompt P is then passed to a response generator function f which uses a decoding mechanism  $\mathcal{G}$  to sample responses 183 184 from an LLM. Hence,  $f(P) = \mathcal{G}(\mathbb{P}_{LLM}(r|P))$ . The final step is a post-processing  $\delta$  applied to the 185 LLM-generated response f(P), in order to extract the task-specific output  $\hat{v}_{test}$ . Commonly used post-processing strategies include regular-expression matching ( $\delta_{regex}$ ) and self-consistency ( $\delta_{SC}$ ) 186 Wang et al. (2023b) based strategies. 187

188 *In-context reasoning* (ICR) uses ICL to learn prompts for tasks involving reasoning with natural 189 language, e.g. numerical reasoning based Question Answering (QA) (Chen et al., 2022b; Lu et al., 190 2023b), commonsense reasoning based QA (Geva et al., 2021b) etc. Due to the inherent difficulty of ICR, prompt generator  $\pi$  which uses multiple high-scoring subsets  $S_1, ..., S_m$  and a test input  $u_{test}$  to 191 generate a prompt P are used. Such prompts  $P = \pi(S_1, ..., S_m, u_{test})$  are expected to improve the 192 overall performance or robustness of the underlying ICR task. Examples of prompt generators are a 193 similarity-based prompt generator  $\pi_{KNN}$  and a diversity-based prompt generator  $\pi_{MMR}$  (discussed 194 in section 4). Hence, the entire output generation process can be described as: 195

$$\hat{v}_{test} = \delta(f(P(U(\mathcal{X})))), \text{ where } P(\mathcal{X}) = \pi(U(\mathcal{X}), u_{test}), \text{ where } U(\mathcal{X}) = (S_1, ..., S_m)(\mathcal{X})$$
 (1)

Here  $U(\mathcal{X}) = (S_1, ..., S_m)(\mathcal{X}))$  a set of m subsets of the training set  $\mathcal{X}$ . Let  $\mathcal{V}$  be the set of n' validation examples  $\{u'_i, v'_i\}_{i=1}^{n'}$ , and define the validation accuracy for a prompt  $P(U(\mathcal{X}))$  as:  $A(P(U(\mathcal{X})), \mathcal{V}) = \frac{1}{n'} \sum_{i=1}^{n'} \mathbb{1}(v'_i = \delta(f(\pi(U(\mathcal{X}), u'_i)))))$ . For ICR, we are interested in finding a set  $U(\mathcal{X})$  of m-subsets of  $\mathcal{X}$  such that the corresponding prompt  $P(U(\mathcal{X}))$  generated by the prompt generator  $\pi$  maximizes the total validation accuracy.

$$U^*(\mathcal{X}) = \underset{U \in \mathcal{S}(\mathcal{X})^m}{\arg \max} A(P(U), \mathcal{V}) \text{ where } \mathcal{S}(\mathcal{X}) \text{ contains all } k \text{-subsets of } \mathcal{X}$$
(2)

We call this as *multiple exemplar-subset selection* (MESS) formulation for finding the optimal prompt.

203 204

196 197

### 3.2 TOP-*m* Arm Selection formulation for Prompt Generation

The MESS problem defined above is a discrete optimization problem over an exponentially large search space  $S(\mathcal{X})^m$ . Additionally, the function A is computationally expensive and nondifferentiable due to the black-box (oracle) access to the LLM response generator function f. Hence, naive or heuristic search-based algorithms for solving MESS are computationally infeasible in this setting. Further, due to the varied nature of the different prompt generators  $(\pi)$ , the accuracy of the actual generated prompt is difficult to optimize. Hence, we optimize the average accuracy of the prompts generated by each of the m-subsets. Let  $a \in \{0, 1\}^n$  such that  $||a||_1 = k$  denote the encoding for an exemplar subset of size k, which also denote a prompt  $P = \pi_r(a)$ . Here,  $\pi_r$  is a prompt generator that takes a random ordering of the exemplars in the subset denoted by a. Considering each  $a_i \in S$  as the *arms* and  $A(\pi_r(a), V)$  as the *reward* for arm a, we use multi-armed bandit (MAB)-based algorithms (Lattimore & Szepesvári, 2020) for sample-efficient exploration of the arms space. Particularly, we pose the modified MESS problem as a top-m arm identification problem (Kalyanakrishnan & Stone, 2010; Bubeck et al., 2009), also called the pure exploration problem, and design an algorithm that iteratively samples the reward for an arm a(t) at each iteration t.

222 We further assume that the observed reward  $A(\pi_r(a(t)), \mathcal{V})$  can be modeled as a linear function 223  $\rho(a(t))$  of the features of arm a(t). Intuitively, given an arm a(t), the reward scoring function  $\rho(a(t))$ 224 should be a function of the similarity between the texts of the selected exemplar,  $u_i$  such that  $a_i = 1$ , 225 and the validation exemplar  $u_i \in \mathcal{V}$ . In this work, we use a normalized BERT-based similarity score,  $\sigma_{ij} = \frac{\phi(u_i)^T \phi(u'_j)}{\|\phi(u_i)\| \|\phi(u'_j)\|}$ , where  $\phi(u)$  is the sentence encoding vector obtained from a pre-trained 226 227 transformer model, e.g. SentenceBERT. Let,  $\alpha_i, i = 1, ..., n$  denote the *i*-th training exemplar's 228 coefficient in the reward scoring function. Our linear model for the observed reward of an arm a(t) at 229 the *t*-th iteration can be written as: 230

$$\rho(\vec{\alpha}, a(t)) = \frac{1}{n'} \sum_{j=1}^{n'} \sum_{i=1}^{n} \alpha_i a(t)_i \sigma_{ij} + \eta_t = \sum_{i=1}^{n} \alpha_i a(t)_i x_i + \eta_t, \text{ where } x_i = \frac{1}{n'} \sum_{j=1}^{n'} \sigma_{i,j} \quad (3)$$

234 and  $\eta_t$  is a subgaussian noise, i.e.  $\mathbb{E}[e^{\lambda \eta_t}] \leq \exp(\lambda^2 \xi^2/2)$ , for some variance  $\xi^2$ . Under this 235 stochastic linear bandits assumption, the problem of identifying top-m arms can solved using 236 the Gap-index based algorithms (Xu et al., 2018b; Réda et al., 2021), which were unified under 237 the GIFA framework (Réda et al., 2021). The gap-index between any two arms  $i, j, B_t(i, j)$ , 238 is defined as the confidence-enhanced gap between their estimated mean rewards  $\hat{\rho}_t$ . Hence,  $B_t(i,j) = \hat{\rho}_t(i) - \hat{\rho}_t(j) + W_t(i,j)$ , where the confidence term is defined as:  $W_t(i,j) = \hat{P}_t(i) - \hat{\rho}_t(j) + \hat{P}_t(i,j)$ 239  $\left[\sqrt{2\ln\left(\frac{1}{\delta}\right) + N\ln\left(1 + \frac{(t+1)L^2}{\lambda^2 N}\right)} + \frac{\sqrt{\lambda}}{\sigma}S\right] (||x_i||_{\hat{\Sigma}_t^{\lambda}} + ||x_j||_{\hat{\Sigma}_t^{\lambda}}), \text{ where } S \text{ and } L \text{ are constants, } N \text{ is the total number of arms pulled, and } \hat{\Sigma}_t^{\lambda} = \sigma^2 (\hat{V}_t)^{-1} (\hat{V}_t \text{ is the design matrix defined below).}$ 240 241 242 243

GIFA algorithms maintain a set of estimated m-best arms  $U_t$ . In each iteration t, the most ambiguous 244 arm from  $U_t$ , say  $b_t = \arg \max_{b \in U_t} \max_{a \in U_t^c} B_t(a, b)$ , and the most ambiguous arm from  $U_t^c$ 245 (called the *challenger* arm), denoted by  $c_t = \arg \max_{c \in U_t^c} B_t(c, b_t)$ , are computed. The arm with 246 the highest variance between  $b_t$  and  $c_t$  is pulled, and the model parameters are updated. The design 247 matrix  $\hat{V}_{t+1}$  is defined as:  $\hat{V}_{t+1} = \lambda I_n + \sum_{a \in S} N_a x_a x_a^T$ , where  $N_a$  is the number of times arm a248 was pulled. The updated parameters  $\hat{\alpha}_{t+1}$  are computed as:  $\hat{\alpha}_{t+1} = (\hat{V}_{t+1})^{-1} (\sum_{l=1}^{t+1} r_l x_{a_l})$ , where 249  $r_l$  is the reward received by computing the accuracy of the prompt generated by the current arm. The 250 main problem with the GIFA framework is that the total number of arms  $|\mathcal{S}|$  is exponentially large. 251 Hence, the computation of most ambiguous arm  $b_t$  and challenger arm  $c_t$  is infeasible. Next, we 252 propose a new scheme for mitigating this problem. 253

253 254 255

231 232 233

### 3.3 CASE: CHALLENGER-ARM SAMPLING-BASED TOP-*m* ARM SELECTION ALGORITHM

256 Implementing gap-index-based schemes for top-m arm identification, for settings with exponentially 257 large number of arms is infeasible. The key problem is to identify the most ambiguous arms in 258 each iteration. We propose to mitigate this problem using: (a) identify a low-regret subset  $N_t$  of 259 m' next-best arms after  $U_t$ , and (b) use a GIFA-based algorithm to identify the top-m arms from 260 the set  $U_t \cup N_t$ . The problem of combinatorial blowup of arms in the linear bandit setting has been 261 sparsely studied, with selective exploration as one of the strategies (e.g. Algorithm 13, Chapter 23 of Lattimore & Szepesvári (2020)). Since it is impractical to explore all arms, the selective exploration 262 scheme uniformly samples arms from the unexplored set and then selects the highest-scoring arms 263 according to the current model. It then pulls the selected arms to update the model parameters. 264

Algorithm 1 describes the proposed *challenger-arm sampling* based exploration technique, called CASE. The set of top-m subsets (arms)  $U_0$  is initialized to a random set sampled from S. In lines 11 – 14, we compute the updated  $U_t$  by moving the highest scoring arm from  $N_{t-1}$  if its score is higher than that of the lowest scoring arm in  $U_{t-1}$ , which is then moved to  $N_{t-1}$ . Using the selective exploration idea, CASE uniformly samples m' arms from  $(U_t \cup N_{t-1})^c$ , to generate the set  $M_t$ , and then selects the top-m' arms from  $M_t \cup N_{t-1}$  to generate the updated  $N_t$ .  $N_t \cup U_t$  is

270 Algorithm 1: CASE: Challenger Arm Sampling for Exemplar selection 271 **1 Input:**  $\mathcal{X}$ : set of all training exemplars, k: prompt size,  $\mathcal{S}$ : all k-subsets of  $\mathcal{X}$ ,  $a \in \mathcal{S}$ : an arm or k-subset 272 2 273 **Define:**  $U_t$ : set of currently estimated top-*m* arms. 3 274  $N_t$ : set of currently estimated next best-m' arms. 4 275  $b_t$ : the most ambiguous arm from  $U_t$ 5 276  $s_t$ : the most ambiguous sampled arm from  $N_t$ 6 277 s Initialize:  $U_0 \leftarrow$  set of random m arms from S,  $t \leftarrow 1$ ,  $\vec{\alpha}_1 \leftarrow \mathcal{N}(0, 1)$ 278 while  $B_t(s_t, b_t) \leq \epsilon$  do 9 279 Construct  $U_t$  by replacing  $n_t \in U_{t-1}$  with a potentially better arm  $c_t \in N_{t-1}$ 10  $n_t = \arg\min_{a \in U_{t-1}} \hat{\rho}_t(a); c_t = \arg\max_{a \in N_{t-1}} \hat{\rho}_t(a)$ 11 281 if  $\hat{\rho}_t(c_t) \geq \hat{\rho}_t(n_t)$  then 12  $U_t, N_t \leftarrow \operatorname{swap}(n_t, c_t) \text{ from } U_{t-1}, N_{t-1}$ 13 14 end 15  $M_t \leftarrow s' \sim_{m'} (U_t \cup N_{t-1})^c$  $\triangleright$  Randomly sample from  $(U_t \cup N_{t-1})^c$ 284  $\triangleright$  Updated  $N_t$  from  $N_{t-1}$  and  $M_t$  $N_t \leftarrow \operatorname{top}_{m'}(M_t \cup N_{t-1}; \hat{\rho}_{(t)})$ 16 Compute the revised most ambiguous arms for detecting convergence 17  $b_{t+1} = \arg \max_{b \in U_t} \max_{a \in N_t} [B_t(a, b) = \hat{\rho}_t(a) - \hat{\rho}_t(b) + W_t(a, b)]$ 18 287  $s_{t+1} = \arg\max_{s \in N_t} \left[ B_t(s, b_{t+1}) = \hat{\rho}_t(s) - \hat{\rho}_t(b_t) + W_t(s, b_{t+1}) \right]$ 19 288 Pull selected arm, receive reward, and update model parameters 20 289  $a_{t+1} \leftarrow \text{selection\_rule}(U_t, N_t)$ 21 22  $r_{t+1} = A(\pi_r(a(t)), \mathcal{V})$  $\hat{V}_{t+1} = \lambda I_n + \sum_{a \in \mathcal{S}} N_a x_a x_a^T \\ \hat{\alpha}_{t+1} = (\hat{V}_{t+1})^{-1} (\sum_{l=1}^{t+1} r_l x_{a_l})$  $\triangleright \lambda$  regularized design matrix 291 23 292 24 Least-squares estimate  $t \leftarrow t + 1$ 293 25 26 end **Output:**  $U_T$ : Set of m arms from K which have the highest reward 27 295

296 297

the high-reward selected set, from which we explore using the selection rule. We use the greedy 298 selection rule proposed in (Réda et al., 2021), where we select the arm that minimizes the variance 299 between  $s_t$  and  $b_t$ :  $a^* = \arg \min_{a \in N_t \cup U_t} ||x_{b_t} - x_{s_t}||_{(\hat{V}_{t-1} + x_a x_a^T)^{-1}}$ . We call the LLM (line 22 in 300 Algorithm 1) after the selection rule, where an arm (i.e., a set of exemplars) is sampled. Specifically, 301  $A(\pi_r(a(t)), \mathcal{V})$  denotes to the computation of accuracy on the validation subset, which requires 302 LLM reasoning to generate predictions. This ensures that the LLM's reasoning process and output 303 generation are explicitly integrated into our algorithm. By doing so, we effectively leverage the 304 gap-index-based multi-arm bandit framework to optimize exemplar selection for complex reasoning 305 tasks. Steps 18 and 19 in algorithm 1 compute the revised most ambiguous arms  $b_t \in U_t$  and 306  $s_t \in N_t$ .  $s_t$  is sampled challenger arm selected from the selected set of next best arms  $N_t$  using the 307 gap index  $B_t(a,b) = \hat{\rho}_t(a) - \hat{\rho}_t(b) + W_t(a,b)$ . Finally, the revised parameters  $\hat{\alpha}_{t+1}$  are computed 308 using the revised design matrix  $\hat{V}_{t+1}$  and the least squares estimation formulae described in lines 23 and 24 of algorithm 1. Note that,  $\hat{V}_{t+1}$  can be computed from  $\hat{V}_t$  using the incremental update 310 formula  $\hat{V}_{t+1} = \hat{V}_t + x_{a_{t+1}} x_{a_{t+1}}^T$ . Hence,  $(\hat{V}_{t+1})^{-1}$  can be calculated in  $O(n^2)$  interusing the 311 Sherman-Morrison formula. We stop the updates when the convergence criteria for switching the 312 arms between  $U_t$  and  $N_t$  have been achieved, i.e.  $B_t(s_t, b_t) \leq \epsilon$ . The time complexity for each 313 iteration of our algorithm is  $O(mm' + n^2 + \text{LLM\_inference\_time})$  which is due to lines 18, 21 and 314 24. Next, we discuss some theoretical results related to our method.

315 316

317

### 3.4 SAMPLE COMPLEXITY BOUNDS FOR THE TOP-*m* SELECTION ALGORITHM

The ability of the proposed algorithm to identify the top-*m* arm and correctly estimate the model parameters  $\alpha$  rests on two arguments. Firstly, the selective exploration strategy results in a low regret set of arms  $U_t \cup N_T$ . Specifically, we assume the average regret (total regret / #iterations) of the set  $U_T \cup N_T$  to upper bounded by  $\epsilon$ . While we postpone a rigorous derivation of the regret bound for CASE to a later study, we justify our assumption by using the SETC Algorithm (Algo 13 in (Lattimore & Szepesvári, 2020)), which follows a similar uniform exploration and commitment strategy in the linear bandit setting. SETC has a regret bound of  $O(d\sqrt{n \log(n)})$  where *n* is the 324 number of time steps. Hence, the algorithm will achieve an average regret of at most  $\epsilon$  after at most 325  $\exp(W(\frac{\epsilon^2}{C^2d^2}))$  timesteps, where W is Lambert's function, and C is the constant for the regret bound. 326 Secondly, once a low regret  $U_T \cup N_T$  set is achieved, the gap-index-based algorithm correctly identifies 327 the top-m arms from the set  $U_T \cup N_T$ . Following (Réda et al., 2021), we obtain a high probability 328  $(1-\delta)$  upper bound for the sample complexity of CASE on the event  $\mathcal{E} \triangleq \bigcap_{t>0} \bigcap_{i,j \in [K]} \left( \rho_i - \rho_j \in \mathcal{E} \right)$ 329 330  $[-B_t(j,i), B_t(i,j)]$ , Let  $\mathcal{S}_m^{\star}$  be the true set of top-*m* arms. We define the true gap of an arm *i* as 331  $\Delta(i) \triangleq \rho(i) - \rho(m+1) \text{ if } i \in \mathcal{S}_m^{\star}, \rho(m) - \rho(i) \text{ otherwise } (\Delta(i) \ge 0 \text{ for any } i \in [K]).$ 332 333 **Theorem 1.** For CASE, on event  $\mathcal{E}$  on which the algorithm is  $(\varepsilon, m, \delta)$ -PAC, stopping time  $\tau_{\delta}$ satisfies  $\tau_{\delta} \leq \inf\{u \in \rho^{*+} : u > 1 + H^{\varepsilon}(\mu)C_{\delta,u}^2 + \mathcal{O}(K)\}$ , where, for algorithm with the largest 334 335 variance selection rule<sup>2</sup> :  $H^{\varepsilon}(\mu) \triangleq 4\sigma^2 \sum_{a \in [K]} \max\left(\varepsilon, \frac{\varepsilon + \Delta_a}{3}\right)^{-2}$ , 336 337 338

Above theorem essentially adopts the result in Theorem 2 from Réda et al. (2021) to the setting where we restrict ourselves to arms in  $U_T \cup N_T$ . It mentions that the  $\epsilon$ -optimal top-m arms from  $U_T \cup N_T$ are present in  $U_T$  with prob.  $1 - \delta$ , if  $T > \tau_{\delta}$ . K is the size of  $U_T \cup N_T$ .

**Proof Overview**: The proof builds upon the proofs for classical Top-*m* linear bandits, LinGapE Xu et al. (2018a) and LinGIFA Réda et al. (2021) while additionally accounting for the challenger arms shortlist in  $N_t$  proposed in this work. To prove it, one of the key components is the following lemma, which holds for any gap indices of the form  $B_t(i, j) \triangleq \hat{\rho}_t(i) - \hat{\rho}_t(j) + W_t(i, j)$  for  $i, j \in [K]^2$ .

**Lemma 1.** On the event  $\mathcal{E}$ , for all t > 0,  $B_t(s_t, b_t)(t) \le \min(-(\Delta(b_t) \lor \Delta(s_t)) + 2W_t(b_t, s_t), 0) + W_t(b_t, s_t)$ ; where  $a \lor b = max(a, b)$ .

The proof for Lemma 1 is provided in Appendix A.1. In summary, Theorem 1 and Lemma 1 together provide an upper bound on the expected number of arm pulls required by the algorithm, which translates to the number of LLM calls needed when applying CASE to complex reasoning tasks.

### 4 EXPERIMENTS AND RESULTS

We aim to address the following research questions:

RQ I. How computationally efficient is CASE compared to other stochastic linear bandit approaches?
 RQ II. Can CASE, a task-level exemplar selection method, achieve competitive performance across diverse tasks compared to state-of-the-art exemplar selection methods?

**RQ III.** Can exemplars selected using smaller LLMs be effectively reused for larger LLMs?

**RQ IV.** How efficient is CASE for exemplar selection compared to state-of-the-art methods?

4.1 EXPERIMENTAL SETUP

**Synthetic Experiments** For synthetic experiments, we adopt a setup similar to Réda et al. (2021) and present results on simulated data. We set  $\sigma = 0.5$ ,  $\epsilon = 0$ , and  $\delta = 0.05$  across all experiments. Each experiment is conducted over 500 simulations. We explore various numbers of arms  $K \in [7, 10, 20]$ , and set the number of top arms to be identified with the highest means as m = 3. We choose a challenger set  $N_t$  of size 3. The feature dimension is fixed at 3 for all synthetic experiments.

367 368

351

352

358

359 360

361

SETUP FOR TASK LEVEL EXEMPLAR SELECTION FOR COMPLEX REASONING

369 Datasets and Metrics: We evaluate on well-known datasets that require complex reasoning. For
 370 numerical reasoning, we use GSM8K and AquaRAT; for commonsense reasoning, we use StrategyQA;
 371 and for tabular and numerical reasoning, we use FinQA and TabMWP. Detailed descriptions of the
 372 datasets are provided in Appendix C. We report performance using the official metrics: Exact Match
 373 (EM) and Cover-EM Rosset et al. (2021) for the respective datasets.

Hyperparameters (LLMs): For LLMs, we set the temperature to 0.3 to reduce randomness in the generated outputs. To reduce repetition, we apply a presence penalty of 0.6 and a frequency penalty of 0.8. The max\_length for generation is set to 900.

<sup>&</sup>lt;sup>2</sup>or pulling both arms in  $\{b_t, c_t\}$  at time t

	Method	GSM8K	AquaRat	TabMWP	FinQA	StrategyQA
			GPT-3.5-turbo			
	Instance level					
	KNN (Rubin et al., 2022)	53.45	51.96	77.07	51.52	81.83
	KNN (S-BERT) (Rubin et al., 2022)	53.07	52.75	77.95	52.65	81.83
	MMR (Ye et al., 2023b)	54.36	51.18	77.32	49.87	82.86
	KNN+SC (Wang et al., 2023b)	80.21	62.59	83.08	54.49	83.88
	MMR+SC (Wang et al., 2023b)	78.01	59.45	81.36	50.74	83.88
	PromptPG (Lu et al., 2023a)	-	-	68.23	53.56	-
	Task level					
	Zero-Shot COT (Kojima et al., 2023)	67.02	49.60	57.10	47.51	59.75
	Manual Few-Shot COT (Wei et al., 2023)	73.46	44.88	71.22	52.22	73.06
	Auto-COT (Zhang et al., 2023b)	62.62	43.31	-	-	71.20
	Complex-COT (Fu et al., 2023)	71.04	51.18	-	-	71.63
	Random	67.79	49.80	55.89	53.70	81.02
	PS+ (Wang et al., 2023a)	59.30	46.00	-	-	-
	GraphCut (Iyer & Bilmes, 2013)	66.19	47.24	60.45	52.31	80.00
	FacilityLocation (Iyer & Bilmes, 2013)	68.61	48.43	67.66	36.79	81.63
	LENS (Li & Qiu, 2023)	69.37	48.82	77.27	54.75	79.79
	LENS+SC (Li & Qiu, 2023)	79.37	57.87	80.68	60.06	82.24
	LENS+KNN+SC (L1 & Q10, 2023)	80.06	57.87	79.79	60.68	83.46
	LENS+MMR+SC (Li & Qiu, 2023)	80.36	59.44	80.39	60.94	82.45
	Our Approach	70.01	54.50	00.40	50.50	0.4.40
	CASE	<b>79.91</b> (▲15.19%) ‡	54.72(12.09%)	83.42(17.96%)‡	<b>59.</b> 72(▲9.08%)†	84.49 (▲5.89%)†
	CASE+SC	<b>86.96</b> (▲25.36%)‡	62.20( <u>1</u> 27.41%)‡	84.91(49.89%)‡	63.64( <b>1</b> 6.24%)‡	<b>87.55</b> (▲9.73%)‡
	CASE+KNN+SC	<b>87.49</b> (▲26.12%)‡	<b>64.17</b> (▲31.44%)‡	<b>86.23</b> ( <b>11.60%</b> )‡	64.25 (A17.35%)‡	<b>85.9</b> 2( <b>▲</b> 7.68%)‡
-	CASE+MMR+SC	85.60 (123.40%)‡	02.00(▲28.23%)‡	<b>85.9I</b> ( <b>▲</b> 11.18%)‡	<b>03.4</b> /(▲15.93%)‡	<b>84.69(</b> ▲6.14%)‡
		76.10	GP1-40-mini	06.04	(0.01	02.05
	LENS (Li & Qiu, 2023)	76.19	64.56	86.34	69.31	92.85
	CASE	91.13	13.25	89.73	72.89	95.92

Table 1: Results across datasets. Percentage improvements are reported over LENS (Li & Qiu, 2023).
 † indicates statistical significance using t-test over LENS at 0.05 level and ‡ at 0.01 level.

401 402

380

403 Hyperparameters (CASE): In CASE we begin by clustering the training set into 5 clusters. We then 404 form the set of exemplar subsets (S), by sampling one exemplar from each cluster *with replacement*. 405 This approach allows us to explore various combinations of exemplars. Note that the training set 406 to form S are the same across baselines for a fair comparison. Our main results are reported for 407 |S| = 100. We set  $m = |U_t| = 10$  to identify the top scoring subsets (arms) and choose a challenger 408 set  $N_t$  of size 5. We set the number of validation examples  $\mathcal{V}$  to, 20 and  $\epsilon = 0.1$  (stopping criterion).

409 Baselines: We compare with zero-shot Kojima et al. (2023) and manual few-shot Wei et al. (2023) settings with COT rationales to highlight the importance of careful exemplar selection. Additionally, 410 we compare against instance-level exemplar selection methods, such as MMR Ye et al. (2023b) and 411 KNN Rubin et al. (2022), which use diversity and similarity-based measures to select exemplars for 412 each test example. For fair comparison, we set the number of exemplars to 5 and configure MMR 413 with  $\lambda = 0.5$  to balance similarity and diversity. We also evaluate against coreset selection methods 414 like Graphcut and Facility Location Iyer & Bilmes (2013). Finally, we compare with LENS Li & Qiu 415 (2023), a state-of-the-art task-level exemplar selection approach that is closely aligned with our work. 416

CASE Hybrid Variants: We also propose hybrid variants of CASE, where we select instance-level
exemplars subset from the top-m subsets identified by CASE using KNN or MMR, thereby reducing
the search space. While KNN and MMR typically select individual exemplars for each test instance,
our hybrid approach scores entire subsets by aggregating the similarity scores of each exemplar
within the subset to the test instance, thereby preserving the interactions between exemplars.

- 422
- 4.2 PERFORMANCE OF CASE COMPARED TO EXISTING GAP-INDEX-BASED APPROACHES

424 To address **RQ1**, we conduct synthetic experiments following the setup in Section 4.1. We evaluate 425 metrics such as average runtime, the average number of comparisons for gap index computation, for 426 multi-armed bandit (MAB) approaches including LinGapE, LinGIFA, and CASE. The results for 427 K = 20, m = 3, and N = 3 are shown in Figure 2. Other synthetic experiments are presented in 428 Appendix H. In Figure 2(a), we observe that CASE significantly reduces the number of comparisons required for gap index computations compared to state-of-the-art gap-index-based MAB approaches. 429 CASE requires significantly fewer comparisons than LinGapE. This improvement is largely due to the 430 novel and principled challenger arm sampling strategy, resulting in efficient gap index computations. 431 The efficiency gains stem from the fact that  $|N_t| < |U_t^c|$ , where existing gap-index frameworks use the



Figure 2: Top-*m* arm identification by CASE, LinGIFA and LinGapE for K=20, m=3, N=3. (a) Average number of comparisons across simulations (b) Average runtime (in seconds) (c) Gap Index  $(B_t(s_t, b_t))$  comparison and (d) Simple regret comparison for each round across simulations

449 entire  $U_t^c$  to select challenger arms, leading to a larger number of computations. In contrast, CASE, 450 prunes the space of possible challenger arms  $(N_t)$ , leading to significant efficiency improvements. 451 From Figure 2(b), we also observe that runtime of CASE is lower when compared to LinGapE (about 5.6x) and LinGIFA (about 12x) due to lower number of arm comparisons and gap-index 452 computations. In Figures 2(c) and 2(d), we analyze the gap index  $(B_t(s_t, b_t))$  and simple regret 453 across rounds (averaged across simulations) and observe that they approach 0 as the number of rounds 454 increases. This demonstrates that, CASE samples good arms with our shortlist of  $N_t$  serving as a 455 good approximation of challenger arms. CASE converges with much lower gap index computations 456 and has lower runtime compared to existing state-of-the-art gap index algorithms. We also observed 457 that all algorithms have similar average sample complexity with negligible differences. 458

### 4.3 PERFORMANCE COMPARISON ON TASK-LEVEL EXEMPLAR SELECTION

461 To address RQ II, we compare CASE and its variants against state-of-the-art task-level and instance-462 level exemplar selection methods. Table 1 shows that CASE consistently outperforms the random, 463 as well as zero-shot and manual Few-Shot COT Wei et al. (2023) methods, which rely on random or hand-picked samples without accounting for the interactions between exemplars and task-level 464 performance. Furthermore, classical coreset selection methods like Graph Cut and Facility Location, 465 which evaluate a subset's informativeness and diversity, perform worse or on par with random 466 selection. This is largely because these coreset methods were not designed for the ICL paradigm. We 467 also observe that CASE outperforms dynamic selection methods like KNN, PromptPG, and MMR. 468 Beyond the efficiency gained during inference, task-level exemplars demonstrate greater robustness 469 (see Appendix E) compared to dynamic methods, which selects exemplars for each test instance, 470 introducing more variance. The independent selection of exemplars in dynamic methods may result 471 in negative interactions or skill redundancy, where lexically different exemplars may represent similar 472 skills, limiting their utility in reasoning tasks. To further support the claim that CASE selects more 473 representative task-level exemplars, we conduct a qualitative analysis comparing exemplars chosen 474 by LENS and CASE (see Appendix G). We find that CASE consistently selects exemplars with 475 diverse skills required for solving tasks, whereas LENS selects exemplars with redundant skills.

476 477

478

445

446

447 448

459

460

### 4.4 EXEMPLAR REUSE AND SELECTION EFFICIENCY

To address RQ III, we evaluate the performance of gpt-3.5-turbo using exemplars selected by smaller
models, such as Llama2-7b (see Figure 3(a)) and Mistral-7b. In Table 1 we present the results of
exemplar reuse from Mistral-7b for CASE. We observe that exemplars selected by CASE using
smaller LLMs perform well with larger LLMs, promoting both exemplar reuse and efficiency.

To address RQ IV, we compare the number of LLM calls made by CASE and LENS (shown in Figure 3(b)). We observe that CASE reduces LLM calls significantly compared to LENS (about 6x to 10.5x) and also reduces the exemplar selection time (shown in Figure 3(c)). This difference arises because LENS iteratively eliminates each training example by computing an informativeness and

499

500 501

502

504

505

506 507

508

523

524

526

527

528

529

530 531

532

486

487



Figure 3: (a) Exemplar Reuse for LLMs, (b-c) Efficiency of CASE compared to LENS.

diversity measure based on LLM confidence estimates, resulting in an LLM call for each training instance over multiple rounds. In contrast, CASE directly scores exemplar subsets, and its novel challenger set sampling mechanism dramatically reduces the number of subsets (arms) that need to be explored and evaluated. Consequently, CASE is more efficient than even state-of-the-art gap-index-based MAB algorithms like LinGIFA and LinGapE, as demonstrated in Figure 2.

#### 4.5 ABLATION STUDIES

509 We conduct several ablation studies, including one-time sampling of all exemplar subsets and a 510 variant of CASE without exploration, to highlight the need for an exploration-based MAB approach, 511 as shown in Table 2. In the one-time sampling approach, we sample each subset once, evaluate them on the validation set, and select the subset with the lowest validation loss. The test set 512 inference results using the selected subset are presented in Table 2. We observe that one-time 513 sampling underperforms compared to CASE since it evaluates all arms only once and lacks sufficient 514 information to confidently identify the best arms. This method tends to overfit the validation set, 515 providing only a baseline understanding of the exemplar subsets (arms) and failing to lead to optimal 516 selection without incorporating exploration or exploitation mechanisms. 517

Furthermore, we conduct an ablation where we fit the linear model once on a set of S subsets (where 518 |S|=K (number of arms)=100) (as shown in Table 2) and select the subset with the highest mean for 519 test set inference. Unlike CASE, which models the top-m subsets in each round, this ablation fits the 520 linear model once for all subsets. Since, this ablation is devoid of exploration mechanism it does not 521 result in confident estimation of top-m arms resulting in suboptimal performance. 522

> Datasets GSM Tab Fin Strat Aqua One-time sampling 76.72 50.39 81.23 54.14 80.20 CASE (-exploration) (Mistral) 76.57 47.64 77.17 45.95 80.00 CASE (Llama) 77.79 56.30 83.65 57.72 82.24 79.91 CASE (Mistral) 54.72 83.42 59.72 84.49

Table 2: Ablation studies: one-time sampling, w/o exploration vs proposed exploration (CASE).

#### 5 CONCLUSION

In this work, we propose an efficient task-level exemplar subset selection method that identifies 534 highly informative exemplar subsets. Our approach significantly reduces resource consumption by minimizing the number of LLM calls, offering a clear advantage over existing state-of-the-art 536 exemplar selection methods. Additionally, we find that exemplars selected using smaller LLMs 537 can be effectively reused for larger models. CASE not only outperforms existing static exemplar selection methods but is also superior to, dynamic selection methods. In the future, we plan to derive 538 a rigorous regret bound for CASE, explore the mechanisms of ICL and the role of exemplars in tasks requiring complex reasoning.

## 540 REFERENCES

580

Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic
 bandits. *Advances in neural information processing systems*, 24, 2011.

544 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-546 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, 547 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz 548 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In 549 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neu-550 ral Information Processing Systems, volume 33, pp. 1877-1901. Curran Associates, Inc., 551 2020a. URL https://proceedings.neurips.cc/paper\_files/paper/2020/ 552 file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf. 553

554 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-555 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec 558 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In 559 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 1877–1901. Curran Associates, Inc., 561 2020b. URL https://proceedings.neurips.cc/paper\_files/paper/2020/ file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf. 563

- Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in multi-armed bandits problems.
   In Algorithmic Learning Theory: 20th International Conference, ALT 2009, Porto, Portugal, October 3-5, 2009. Proceedings 20, pp. 23–37. Springer, 2009.
- 567
   568
   569
   569
   569
   569
   570
   570
   58–265, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL https://proceedings.
   571
   mlr.press/v28/bubeck13.html.
- Lijie Chen, Jian Li, and Mingda Qiao. Nearly Instance Optimal Sample Complexity Bounds for Top-k Arm Selection. In Aarti Singh and Jerry Zhu (eds.), *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 101–110. PMLR, 20–22 Apr 2017. URL https://proceedings.mlr. press/v54/chen17a.html.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting:
   Disentangling computation from reasoning for numerical reasoning tasks, 2022a.
- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. Finqa: A dataset of numerical reasoning over financial data, 2022b.
- 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, 2021.
- Rémy Degenne, Pierre Menard, Xuedong Shang, and Michal Valko. Gamification of pure exploration for linear bandits. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 2432–2442. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/ degenne20a.html.
- Tanner Fiez, Lalit Jain, Kevin Jamieson, and Lillian Ratliff. Sequential experimental design for transductive linear bandits, 2019. URL https://arxiv.org/abs/1906.08399.

594 595 596	Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting for multi-step reasoning. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=yflicZHC-19.
597 598 599 600 601 602	Victor Gabillon, Mohammad Ghavamzadeh, and Alessandro Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger (eds.), <i>Advances in Neural Information Processing Systems</i> , volume 25. Curran Associates, Inc., 2012. URL https://proceedings.neurips.cc/paper_files/paper/2012/file/8b0d268963dd0cfb808aac48a549829f-Paper.pdf.
603 604 605	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. <i>Transactions of the Association for Computational Linguistics</i> , 9:346–361, 2021a.
606 607	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies, 2021b.
609 610 611	In Gim, Guojun Chen, Seung seob Lee, Nikhil Sarda, Anurag Khandelwal, and Lin Zhong. Prompt cache: Modular attention reuse for low-latency inference, 2024. URL https://arxiv.org/abs/2311.04934.
612 613	Chengcheng Guo, Bo Zhao, and Yanbing Bai. Deepcore: A comprehensive library for coreset selection in deep learning, 2022.
614 615 616	Rishabh K Iyer and Jeff A Bilmes. Submodular optimization with submodular cover and submodular knapsack constraints. <i>Advances in neural information processing systems</i> , 26, 2013.
617 618	Yassir Jedra and Alexandre Proutiere. Optimal best-arm identification in linear bandits. Advances in Neural Information Processing Systems, 33:10007–10017, 2020.
619 620 621	Shivaram Kalyanakrishnan and Peter Stone. Efficient selection of multiple bandit arms: Theory and practice. In <i>ICML</i> , volume 10, pp. 511–518, 2010.
622 623 624 625	Shivaram Kalyanakrishnan, Ambuj Tewari, Peter Auer, and Peter Stone. Pac subset selection in stochastic multi-armed bandits. In <i>Proceedings of the 29th International Coference on International Conference on Machine Learning</i> , ICML'12, pp. 227–234, Madison, WI, USA, 2012. Omnipress. ISBN 9781450312851.
626 627 628 629 630	Emilie Kaufmann and Shivaram Kalyanakrishnan. Information complexity in bandit subset selection. In Shai Shalev-Shwartz and Ingo Steinwart (eds.), <i>Proceedings of the 26th Annual Conference</i> <i>on Learning Theory</i> , volume 30 of <i>Proceedings of Machine Learning Research</i> , pp. 228–251, Princeton, NJ, USA, 12–14 Jun 2013. PMLR. URL https://proceedings.mlr.press/ v30/Kaufmann13.html.
631 632 633	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners, 2023.
634 635	Tor Lattimore and Csaba Szepesvári. Bandit algorithms. Cambridge University Press, 2020.
636 637	Xiaonan Li and Xipeng Qiu. Finding support examples for in-context learning. In <i>The 2023</i> Conference on Empirical Methods in Natural Language Processing, 2023.
638 639	Zhaoqi Li, Kevin Jamieson, and Lalit Jain. Optimal exploration is no harder than thompson sampling, 2023. URL https://arxiv.org/abs/2310.06069.
640 641 642 643 644 645	Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 158–167, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1015. URL https://aclanthology.org/P17-1015.
646 647	Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning, 2023a.

648	Pan Lu, Liang Oiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Raipurohit, Peter
649	Clark, and Ashwin Kalvan. Dynamic prompt learning via policy gradient for semi-structured
650	mathematical reasoning, 2023b.
651	6,

- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered 652 prompts and where to find them: Overcoming few-shot prompt order sensitivity. In Smaranda 653 Muresan, Preslav Nakov, and Aline Villavicencio (eds.), Proceedings of the 60th Annual Meeting 654 of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 8086-8098, Dublin, 655 Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long. 656 556. URL https://aclanthology.org/2022.acl-long.556. 657
- 658 Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. Measuring 659 and narrowing the compositionality gap in language models, 2023. 660
- 661 Clémence Réda, Emilie Kaufmann, and Andrée Delahaye-Duriez. Top-m identification for linear 662 bandits. In International Conference on Artificial Intelligence and Statistics, pp. 1108–1116. PMLR, 2021. 663
- 664 Corby Rosset, Chenyan Xiong, Minh Phan, Xia Song, Paul Bennett, and Saurabh Tiwary. Knowledge-665 aware language model pretraining, 2021. 666
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context 668 learning. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 2655–2671, Seattle, United 670 States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main. 191. URL https://aclanthology.org/2022.naacl-main.191.
- 672 Clémence Réda, Emilie Kaufmann, and Andrée Delahaye-Duriez. Top-m identification for linear 673 bandits, 2021. 674
- 675 Jonathan Tonglet, Manon Reusens, Philipp Borchert, and Bart Baesens. Seer : A knapsack approach 676 to exemplar selection for in-context hybridga, 2023. 677
- 678 Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 679 Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), Proceedings of the 680 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 681 pp. 2609–2634, Toronto, Canada, July 2023a. Association for Computational Linguistics. doi: 682 10.18653/v1/2023.acl-long.147. URL https://aclanthology.org/2023.acl-long. 683 147. 684
- 685 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha 686 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In The Eleventh International Conference on Learning Representations, 2023b. 688
- 689 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, 690 Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models, 2022. 691
- 692 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, 693 and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. 694
- Jing Xiong, Zixuan Li, Chuanyang Zheng, Zhijiang Guo, Yichun Yin, Enze Xie, Zhicheng YANG, 696 Qingxing Cao, Haiming Wang, Xiongwei Han, Jing Tang, Chengming Li, and Xiaodan Liang. DQ-697 lore: Dual queries with low rank approximation re-ranking for in-context learning. In The Twelfth 698 International Conference on Learning Representations, 2024. URL https://openreview. net/forum?id=qAoxvePSlq. 699
- 700

667

669

Liyuan Xu, Junya Honda, and Masashi Sugiyama. Fully adaptive algorithm for pure exploration in 701 linear bandits, 2017. URL https://arxiv.org/abs/1710.05552.

- Liyuan Xu, Junya Honda, and Masashi Sugiyama. A fully adaptive algorithm for pure exploration in linear bandits. In Amos Storkey and Fernando Perez-Cruz (eds.), *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pp. 843–851. PMLR, 09–11 Apr 2018a. URL https://proceedings.mlr.press/v84/xu18d.html.
- Liyuan Xu, Junya Honda, and Masashi Sugiyama. A fully adaptive algorithm for pure exploration in linear bandits. In *International Conference on Artificial Intelligence and Statistics*, pp. 843–851.
   PMLR, 2018b.
- Greg Yang, Edward J. Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder,
   Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural networks
   via zero-shot hyperparameter transfer, 2022.
- Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. Compositional exemplars for in-context learning, 2023a.
- Xi Ye, Srinivasan Iyer, Asli Celikyilmaz, Veselin Stoyanov, Greg Durrett, and Ramakanth Pasunuru. Complementary explanations for effective in-context learning. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 4469–4484, Toronto, Canada, July 2023b.
   Association for Computational Linguistics. URL https://aclanthology.org/2023. findings-acl.273.
- Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models in-context,
   2023a. URL https://arxiv.org/abs/2306.09927.
- Yiming Zhang, Shi Feng, and Chenhao Tan. Active example selection for in-context learning. In
   Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 9134–9148, Abu Dhabi, United Arab
   Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
   emnlp-main.622. URL https://aclanthology.org/2022.emnlp-main.622.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. In *The Eleventh International Conference on Learning Representations*, 2023b. URL https://openreview.net/forum?id=5NTt8GFjUHkr.
- Tony Z. Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models, 2021.

PROOFS

A.1 PROOF OF LEMMA 1

756

757 758

759

760

806

А

#### *Proof.* To prove the above case, we first introduce the following property. We primarily follow the proof structure of GIFA framework (Réda et al., 2021) with some modifications required for CASE 761 due to the shortlist $N_t$ and our swapping rule to compute $U_t$ . 762 Let $\mathcal{S}_m^{\star}$ be the true set of top-*m* arms and $(S_m^{\star})^c$ denote the true set remaining worst arms. 763 764 **Property 1:** For $b_t \in U_t$ and $s_t \in N_t$ it holds that $\hat{\rho}_t(b_t) \geq \hat{\rho}_t(s_t)$ . Hence, it follows that 765 $B_t(s_t, b_t) = \hat{\Delta}_t(s_t, b_t) + W_t(b_t, s_t) \leq W_t(b_t, s_t)$ as $\hat{\Delta}_t(s_t, b_t) < 0$ From property 1, we can 766 establish that $B_t(s_t, b_t) \leq W_t(b_t, s_t)$ . Hence, to show that 767 $B_t(s_t, b_t) \le -(\Delta(b_t) \lor \Delta(s_t)) + 3W_t(b_t, s_t)$ 768 769 we consider the following scenarios: 770 (i) $b_t \in \mathcal{S}_m^{\star}$ and $s_t \notin \mathcal{S}_m^{\star}$ : In that case, 771 772 $\Delta(b_t) = \rho(b_t) - \rho(m+1); \Delta(s_t) = \rho(m) - \rho(s_t)$ 773 774 As event $\mathcal{E}$ holds, 775 $B_t(s_t, b_t) = -B_t(b_t, s_t) + 2W_t(b_t, s_t) < \Delta(s_t, b_t) + 2W_t(b_t, s_t)$ 776 777 As $s_t \notin \mathcal{S}_m^{\star}$ , 778 $\rho(s_t) < \rho(m+1)$ 779 $\Delta(s_t, b_t) < \rho(m+1) - \rho(b_t) = -\Delta(b_t)$ 780 781 782 But as $b_t \in \mathcal{S}_m^{\star}$ , it also holds that $\rho(b_t) \ge \rho(m)$ , and $\Delta(s_t, b_t) \le \rho(s_t) - \rho(m) = -\Delta(s_t)$ . 783 Hence $B_t(s_t, b_t) \leq -(\Delta(b_t) \vee \Delta(s_t)) + 2W_t(b_t, c_t) \leq -(\Delta(b_t) \vee \Delta(s_t)) + 3W_t(b_t, c_t)$ . 784 785 (ii) $b_t \notin \mathcal{S}_m^{\star}$ and $s_t \in \mathcal{S}_m^{\star}$ : 786 787 $\Delta(s_t) = \rho(s_t) - \rho(m+1); \Delta(b_t) = \rho(m) - \rho(b_t)$ 788 789 By Property 1, 790 $B_t(s_t, b_t) < W_t(b_t, s_t) < \hat{\Delta}_t(b_t, s_t) + W_t(b_t, s_t) = B_t(b_t, s_t)$ 791 as $\hat{\rho}_t(b_t) \geq \hat{\rho}_t(s_t)$ . Further, as $\mathcal{E}$ holds. 792 793 $B_t(b_t, s_t) = -B_t(s_t, b_t) + 2W_t(b_t, s_t) \le \Delta(b_t, s_t) + 2W_t(b_t, s_t)$ 794 795 As $b_t \notin \mathcal{S}_m^{\star}$ , $\rho(b_t) \leq \rho(m+1)$ and hence $\Delta(b_t, s_t) \leq \rho(m+1) - \rho(s_t) = -\Delta(s_t)$ As $s_t \in \mathcal{S}_m^{\star}$ , 796 $\rho(s_t) \ge \rho(m)$ and hence $\Delta(b_t, s_t) \le \rho(b_t) - \rho(m) = -\Delta(b_t)$ . 797 Hence $B_t(s_t, b_t) \leq -(\Delta(b_t) \vee \Delta(s_t)) + 2W_t(b_t, c_t) \leq -(\Delta(b_t) \vee \Delta(s_t)) + 3W_t(b_t, c_t).$ 798 799 (iii) $b_t \notin S_m^*$ and $s_t \notin S_m^*$ : We state that there exists a $b \in S_m^*$ that belongs to $N_t$ . At any time t, 800 $M_t \leftarrow s' \sim_{m'} (U_t \cup N_{t-1})^c$ 801 $N_t \leftarrow \operatorname{top}_{m'}(M_t \cup N_{t-1}; \hat{\rho}_{(t-1)})$ 802 803 Due to the above sampling approach adopted for $N_t$ which captures the next m' arms with the highest 804 means, we posit that $N_t$ captures at least one arm in $\mathcal{S}_m^{\star}$ . Assuming the event $\mathcal{E}$ holds and $b \in \mathcal{S}_m^{\star}$ , 805 $W_t(b_t, s_t) > B_t(s_t, b_t) > B_t(b, b_t)$

As by the definition of  $s_t$ , which is one of the most ambiguous arms with largest gap to  $b_t$ 808  $B_t(s_t, b_t) \ge B_t(b, b_t)$ . Hence,  $B_t(s_t, b_t) \ge B_t(b, b_t)$ . From this and event  $\mathcal{E}$  it follows that 809

$$B_t(s_t, b_t) \ge B_t(b, b_t) \ge \rho(b) - \rho(b_t) \ge \rho(m) - \rho(b_t)$$

)

. Hence  $W_t(b_t, s_t) \geq B_t(s_t, b_t) \geq \Delta(b_t)$ . Using event  $\mathcal{E}$ ,  $B_t(s_t, b_t) \le \Delta(s_t, b_t) + 2W_t(b_t, s_t) = (\rho(s_t) - \rho(m)) + (\rho(m) - \rho(b_t)) + 2W_t(b_t, s_t)$ (4)From 4 and since  $B_t(s_t, b_t) \ge \Delta(b_t)$ ,  $B_t(s_t, b_t) < -\Delta(s_t) + \Delta(b_t) + 2W_t(b_t, s_t) < -\Delta(s_t) + 3W_t(b_t, s_t)$ Also from Property 1 and  $W_t(b_t, s_t) \ge \Delta(b_t)$ , it holds that  $B_t(s_t, b_t) < W_t(b_t, s_t) = -W_t(b_t, s_t) + 2W_t(b_t, s_t) < -\Delta(b_t) + 2W_t(b_t, s_t) < -\Delta(b_t) + 3W_t(b_t, s_t)$ Hence  $B_t(s_t, b_t) \leq -(\Delta(b_t) \vee \Delta(s_t)) + 3W_t(b_t, c_t).$ (iv)  $b_t \in S_m^*$  and  $s_t \in S_m^*$ : Then there exists a  $s \notin S_m^*$  and  $s \in U_t$  In that case,  $\Delta(b_t) = \rho(b_t) - \rho(m+1); \Delta(s_t) = \rho(s_t) - \rho(m+1)$ Also by definition of  $b_t$  and  $s_t$ , it holds that  $B_t(s_t, b_t) = \max_{i \in U_t} \max_{i \in N_t} [B_t(j, i)]$  Since there exists  $s \in U_t$  and  $s_t \in N_t$ ,  $B_t(s_t, b_t) = \max_{i \in U_t} \max_{j \in N_t} [B_t(j, i)] \ge \max_{j \in N_t} B_t(j, s) \ge B_t(s_t, s) \ge \rho(s_t) - \rho(s) \ge \rho(s_t) - \rho(m+1)$ As  $\rho(s_t) - \rho(m+1) = \Delta(s_t)$ ,  $B_t(s_t, b_t) \ge \Delta(s_t)$  By property 1,  $B_t(s_t, b_t) \le W_t(b_t, s_t)$ . Hence,  $\Delta(s_t) \le B_t(s_t, b_t) \le W_t(b_t, s_t)$ On event  $\mathcal{E}$  it follows that  $B_t(s_t, b_t) \leq \rho(s_t) - \rho(b_t) + 2W_t(b_t, s_t)$  as  $(B(s_t, b_t) \leq W_t(b_t, s_t)$  Then  $\rho(s_t) - \rho(b_t)$  can be expressed as  $\rho(s_t) - \rho(m+1) + \rho(m+1) - \rho(b_t)$ . hence,  $B_t(s_t, b_t) \le \rho(s_t) - \rho(m+1) + \rho(m+1) - \rho(b_t + 2W_t(b_t, s_t) \le \Delta(s_t) - \Delta(b_t) + 2W_t(b_t, s_t)$ We already know that  $B_t(s_t, b_t) \ge \Delta(s_t)$  resulting in,  $(a)B_t(s_t, b_t) \le -\Delta(b_t) + 3W_t(b_t, s_t)$ Now to prove  $B_t(s_t, b_t) \leq -\Delta(s_t) + 3W_t(b_t, s_t)$ , we rely on property 1,  $B(s_t, b_t) \le W_t(b_t, s_t) \le -W_t(b_t, s_t) + 2W_t(b_t, s_t)$ As  $W_t(b_t, s_t) \geq \Delta(s_t), -W_t(b_t, s_t) \leq -\Delta(s_t)$ . Therefore,  $(b)B(s_t, b_t) \le W_t(b_t, s_t) \le -W_t(b_t, s_t) + 2W_t(b_t, s_t) \le -\Delta(s_t) + W_t(b_t, s_t) \le -\Delta(s_t) + 3W_t(b_t, s_t) + 3W_t(b_t, s_t) \le -\Delta(s_t) + 3W_t(b_t, s_t) + 3W_t(b_t,$ From (a) and (b)  $B_t(s_t, b_t) \leq -(\Delta(b_t) \vee \Delta(s_t)) + 3W_t(b_t, c_t)$ . В **PROOF STRUCTURE FOR THEOREM 1** *Proof.* Combining Lemma 4 with stopping rule  $B_t(s_t, b_t) \leq \epsilon$  following Lemma 8 in Réda et al. (2021) directly yields

$$N_{a_t}(t) \leq 4\sigma^2 C_{\delta,t}^2 \max\left(\varepsilon, \frac{\varepsilon + \Delta_{a_t}}{3}\right)^{-2}$$

where  $N_{a_t}(t)$  is the number of times arms a is sampled. This is equivalent to the sample complexity term  $\mathrm{H}^{\varepsilon}(\mu)$  in Theorem 1. Hence, maximum number of samplings on event  $\mathcal{E}$  is upper-bound by inf $_{u \in \mathbb{R}^{*+}} \{ u > 1 + \mathrm{H}^{\varepsilon}(\mu) C_{\delta,u}^2 \}$ , where  $\mathrm{H}^{\varepsilon}(\mu) \triangleq 4\sigma^2 \sum_{a \in [K]} \max\left(\varepsilon, \frac{\varepsilon + \Delta_a}{3}\right)^{-2}$ .

Dataset	#Train	#Test	Example Question	Description
GSM8K Cobbe et al. (2021)	7473	1319	The red car is 40% cheaper than the blue car. The price of the blue car is \$100. How much do both cars cost?	multi-step arithmetic word problems
AquaRat Ling et al. (2017)	97467	254	John found that the average of 15 numbers is 40. If 10 is added to each number then the mean of number is?	multi-step arithmetic word problems
TabMWP Lu et al. (2023b)	23059	7686	A newspaper researched how many grocery stores there are in each town. What is the mean of the numbers?	Table based numerical reasoning
FinQA Chen et al. (2022b)	6251	1147	what is the percentage change in the the gross liability for unrecognized tax benefits during 2008 compare to 2007?	Table and Text based numerical reasoning
StrategyQA Geva et al. (2021a)	) 1800	490	Does the United States Secretary of State answer the phones for the White House?	multi-step reasoning

Table 3: Overview of the Complex QA datasets used in this study.

#### DATASETS DESCRIPTION С

879 An overview of the dataset statistics and examples are shown in Table 3.

880 **FinOA**: Comprises financial questions over financial reports that require numerical reasoning with structured and unstructured evidence. Here, 23.42% of the questions only require the information 882 in the text to answer; 62.43% of the questions only require the information in the table to answer; 883 and 14.15% need both the text and table to answer. Meanwhile, 46.30% of the examples have one 884 sentence or one table row as the fact; 42.63% has two pieces of facts; and 11.07% has more than two 885 pieces of facts. This dataset has 1147 questions in the evaluation set.

886 AquaRat: It comprises 100,000 algebraic word problems in the train set with dev and test set each 887 comprising 254 problems. The problems are provided along with answers and rationales providing the step-by-step solution to the problem. An examples problem is shown in Table 3. 889

**TabMWP**: It is a tabular-based math word problem-solving dataset with 38,431 questions. TabMWP 890 is rich in diversity, where 74.7% of the questions in TabMWP belong to free-text questions, while 891 25.3% are multi-choice. We evaluate on the test set with 7686 problems. 892

893 **GSM8K**: This dataset consists of linguistically diverse math problems that require multi-step reasoning. The dataset consists of 8.5K problems and we evaluate on the test set of 1319 questions. 894

895 **StrategyQA**: To prove the generality of our approach for reasoning tasks, we evaluate on StrategyQA 896 Geva et al. (2021b), a dataset with implicit and commonsense reasoning questions. Since there is no 897 public test set with ground truth answers, we perform stratified sampling done on 2290 full train set 898 to split into 1800 train and 490 test.

899 Metrics: For TabMWP and StrategyQA we employ cover-EM Rosset et al. (2021); Press et al. (2023), 900 a relaxation of Exact Match metric which checks whether the ground truth answer is contained in the 901 generated answer. This helps handle scenarios where LLM generates "24 kilograms" and the ground 902 truth is "24". For other numerical reasoning datasets, we employ Exact match.

903 904 905

874

875 876 877

878

#### D **RESULTS USING ALTERNATE OPEN SOURCE LLMS**

906 We also report the performance of exemplars from CASE on open-source models like Mistral-7b 907 and LLama2-7b. The results are shown in Table 4. We observe that the absolute performance across 908 baselines and CASE is lower for smaller LLMs like Llama2 and Mistral-7b when compared to 909 gpt-3.5-turbo or gpt-4o-mini. We observe that this is due to the scale of the Language models as 910 Mistral and LIAMA2 models have 7 billion parameters while gpt-3.5-turbo is of much larger scale 911 and the emergent capabilities like ICL, reasoning capabilities are more pronounced in large scale 912 models Wei et al. (2022).

913 However, we still observe that CASE leads to reasonable performance gains over other static exemplar 914 selection methods across the smaller open-source LLMs. We also observe that CASE is competitive 915 with instance-level/dynamic exemplar selection methods. 916

Our main experiments are carried out in an exemplar reuse setup where exemplar selection is done 917 using small open source LLMs and transferred to larger LLMs. This is done to reduce the LLM

inference cost during exemplar selection. This setup also leverages the reasoning and emergent capabilities of large scale LLMs. This philosophy is inspired from the work  $\mu$ P Yang et al. (2022) where the language model hyperparameters are tuned using smaller LM and transferred to a larger LM for the task under consideration.

Method	GSM8K	AquaRat	TabMWP	FinQA	StrategyQA	
		Mistral-7B				
Instance-level						
KNN Rubin et al. (2022)	28.00	23.16	45.3	9.06	78.27	
MMR Ye et al. (2023b)	28.97	18.11	47.61	10.11	79.95	
Task-level						
Zero-shot-COT Kojima et al. (2023)	7.42	18.89	38.96	1.74	35.37	
Manual Few-shot COT	22.36	14.90	41.93	3.22	62.55	
LENS Li & Qiu (2023)	26.08	14.17	41.82	5.14	76.12	
Our Approach						
CASE	32.6	21.2	45.55	11.24	77.75	
		Llama2-7B				ĺ
Instance-level						
KNN Rubin et al. (2022)	22.51	23.62	43.02	10.37	76.35	
MMR Ye et al. (2023b)	21.60	21.65	41.66	12.20	76.32	
Task-level						
Zero-shot-COT Kojima et al. (2023)	6.14	6.29	12.64	1.67	53.27	
Manual Few-shot COT	19.26	20.47	23.62	2.87	64.29	
LENS Li & Qiu (2023)	17.06	19.29	33.20	6.62	73.06	
Our Approach						
CASE	21.91	24.02	44.69	9.59	77.55	

Table 4: Results across datasets on MISTRAL-7B and LLAMA-2-7B (5-shot exemplars).

#### Ε ROBUSTNESS OF EXEMPLARS SELECTED BY CASE

We compare the robustness of CASE to other exemplar selection methods. We measure standard deviation of performance across different subsets of the evaluation set through 10-fold cross validation, as shown in Table 5. We observe that in 3 out of 4 datasets, exemplars chosen by CASE has less variance in task performance when compared to other exemplar selection methods. Exemplars selected through instance-level approaches are not optimized for the task but rather on a per-testexample basis. Consequently, this leads to greater variance in final task performance. Hence, CASE helps select exemplars for the task which are more robust than other static methods or instance-level selection methods.

#### F PROMPTS

We also demonstrate the instructions issued to the LLM for different tasks discussed in this work, along with some exemplars selected using CASE. An example of prompt construction for FinQA is shown in Figure 8. We also showcase example prompts for AquaRat (Figure 7), GSM8K (Figure 6), TabMWP (Figure 9) and StrategyQA (Figure 10).

#### G EXEMPLAR QUALITATIVE ANALYSIS

We provide a qualitative analysis of exemplars and compare the exemplars selected using CASE with exemplars selected using LENS Li & Qiu (2023), the recent state-of-the-art approach. The final set of exemplars chosen by LENS vs CASE for the AquaRat dataset is shown in Table 6. We observe that Question 4 and Question 5 in the set of exemplars chosen by LENS are redundant in that they are very similar problems that require similar reasoning steps and are also similar thematically. Both the

Datasets	GSM8K	AquaRat	TabMWP	FinQA	StrategyQA
Zero-Shot COT	±5.18	±7.08	±1.84	±4.50	±4.19
Few-Shot COT	$\pm 4.48$	$\pm 12.03$	$\pm 1.66$	$\pm 4.76$	$\pm 5.67$
KNN	$\pm 3.76$	$\pm 5.49$	$\pm 1.27$	$\pm 4.17$	$\pm 4.85$
MMR	$\pm 4.00$	$\pm 10.53$	$\pm 1.68$	$\pm 6.10$	$\pm 5.70$
Graph Cut	$\pm 6.38$	$\pm 8.18$	$\pm 2.03$	$\pm 5.29$	$\pm 7.62$
Facility Location	$\pm 4.23$	$\pm 6.71$	$\pm 1.74$	$\pm 4.94$	$\pm 5.93$
LENS	$\pm 5.04$	$\pm 6.67$	$\pm 1.59$	$\pm 5.81$	$\pm 3.98$
CASE	$\pm$ 3.47	$\pm 6.86$	$\pm 0.88$	$\pm 3.72$	$\pm 2.91$

Table 5: Comparison of robustness of CASE to other approaches. We report standard deviation (lower is better) with scores from different splits of the evaluation set.



Figure 4: Top-*m* arm identification by CASE, LinGIFA and LinGapE for K=7, m=3, N=3. (a) Average number of comparisons across simulations (b) Average runtime (in seconds) (c) Gap Index  $(B_t(s_t, b_t))$  comparison and (d) Simple regret comparison for each round across simulations 

questions are centered on the theme of work and time and are phrased in a similar manner. Hence, they do not add any additional information to solve diverse problems the LLM may encounter during inference. However, we observe that the exemplars chosen by CASE are problems that require diverse reasoning capabilities and are also different thematically. 

We also compare the exemplars chosen by CASE with LENS for the FinQA dataset. We observe that the exemplars chosen by CASE comprises diverse set of problems. We also observe that CASE also contains exemplars that require composite numerical operations with multi-step reasoning rationales to arrive at the solutions, whereas LENS mostly has exemplars with single-step solutions. 

The exemplars chosen by LENS compared to CASE for TabMWP are shown in Table 9. We observe that exemplar 1 and exemplar 3 chosen by LENS are redundant, as they represent the same reasoning concept of computing median for a list of numbers. However, we observe that CASE selects diverse exemplars, with each exemplar representing a different reasoning concept. We also demonstrate the exemplars for GSM8K and StrategyQA in Table 8 and Table 10 respectively. 

SYNTHETIC EXPERIMENTS - EFFICIENCY AND CONVERGENCE ANALYSIS Η

We further present the results for synthetic experiments for scenarios where K = 7, m = 3, N = 3and K = 10, m = 3, N = 3 in Figures 4 and 5. We observe that in both the cases, CASE drastically reduces the number of gap-index computations and comparison based operations (comparing arms). For instance, for K = 10, m = 3 scenario, on average across 500 simulations CASE only requires 20366.84 comparisons, whereas LinGapE requires 948206.10 comparisons and LinGIFA requires 2180251.73 comparisons. This is due to the shortlist of challenger arms  $N_t$  maintained by the proposed approach CASE. We also observe that this results in significant reduction in time (approx. **6x** lower compared to LinGIFA and **2.5x** compared to LinGapE for K = 10 case) due to low number of comparisons. We observe that the gap index and simple regret approaches 0 in a similar trend for





Figure 7: Prompt for AquaRat

**Test Input** : Question: Options: Explanation: [INS] Answer: [INS]

#### 1080 1082 **FinQA** Prompt 1084 Instruction: You are a helpful, respectful and honest assistant helping 1085 solve math word tasks requiring reasoning, using the information 1086 from given table and text. 1087 Exemplars : 1088 *Read the following table, and then answer the question:* 1089 [Table]: beginning balance as of december 1 2007 | 201808 gross increases in unrecognized tax benefits 2013 prior year tax positions | 14009 1090 gross increases in unrecognized tax benefits 2013 current year tax positions | 11350 1091 ending balance as of november 28 2008 | 139549 1092 [Question]: what is the percentage change in the the gross liability for unrecognized tax benefits during 1093 2008 compare to 2007? 1094 [Explanation]: x0 = 139549 - 201808, ans = x0/2018081095 [Answer]: -30.9% 1096 Test Input: Read the table and answer the question: Table: Question: Explanation: [INS] Answer: 1098 [INS] 1099 1100 1101 Figure 8: Prompt for FinQA 1102 1103 1104 1105 1106 1107 1108 TabMWP Prompt 1109 Instruction: You are a helpful, respectful and honest assistant helping 1110 to solve math word problems or tasks requiring reasoning or math, 1111 using the information from the given table. Solve the given 1112 problem step by step providing an explanation for your answer. 1113 Exemplars : 1114 [Table]: Town | Number of stores 1115 Mayfield | 9 1116 Springfield | 9 Riverside | 6 1117 Chesterton | 5 1118 Watertown | 2 1119 [Question]: A newspaper researched how many grocery stores there are in each town. What is the range 1120 of the numbers? 1121 [Explanation]: Read the numbers from the table. 9, 9, 6, 5, 2First, find the greatest number. The greatest number is 9. 1122 Next, find the least number. The least number is 2. 1123 Subtract the least number from the greatest number: 9 - 2 = 71124 [Answer]: The range is 7 1125 . . . 1126 1127 **Test Input** : Table: Question: 1128 Explanation: [INS] Answer: [INS] 1129 1130 1131

Figure 9: Prompt for TabMWP

1134	
1135	
1136	
1137	
1138	
1130	
1133	
1140	
1141	
1142	
1143	
1144	
1145	
1146	
1147	
1148	
1149	
1150	
1151	StrategyOA Prompt
1152	Sume5, QITTOMPC
1153	Instruction: You are a helpful, respectful and honest assistant
1154	helping to solve commonsense problems requiring reasoning. Follow
1155	the given examples that use the facts to answer a question by
1156	decomposing into sub-questions first and then predicting the final
1157	answer as "Yes" or "No" only.
1158	Exemplars :
1159	[Facts]: The role of United States Secretary of State carries out the President's foreign policy. The
	White House has multiple phone lines managed by multiple people.
1160	
1160 1161	[Question]: Does the United States Secretary of State answer the phones for the White House?
1160 1161 1162	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1?
1160 1161 1162 1163	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1?
1160 1161 1162 1163 1164	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No
1160 1161 1162 1163 1164 1165	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No
1160 1161 1162 1163 1164 1165 1166	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No 
1160 1161 1162 1163 1164 1165 1166 1167	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question:
1160 1161 1162 1163 1164 1165 1166 1167 1168	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  <b>Test Input</b> : Facts: Question: Sub-question: [INS] Answer: [INS]
1160 1161 1162 1163 1164 1165 1166 1167 1168	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No   <b>Test Input</b> : Facts: Question: Sub-question: [INS] Answer: [INS]
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No   <b>Test Input</b> : Facts: Question: Sub-question: [INS] Answer: [INS]
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No   <b>Test Input</b> : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No   <b>Test Input</b> : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No   <b>Test Input</b> : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No   <b>Test Input</b> : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No   <b>Test Input</b> : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No   <b>Test Input</b> : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180	[Question]: Does the United States Secretary of State answer the phones for the White House?         [Sub-question 1]: What are the duties of the US Secretary of State?         [Sub-question 2]: Are answering phones part of #1?         [Answer]: No            Test Input       : Facts: Question:         Sub-question: [INS] Answer: [INS]    Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181 1182	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181 1182 1183	[Question]: Does the United States Secretary of State answer the phones for the White House?         [Sub-question 1]: What are the duties of the US Secretary of State?         [Sub-question 2]: Are answering phones part of #1?         [Answer]: No            Test Input       : Facts: Question:         Sub-question: [INS] Answer: [INS]    Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1177 1178 1179 1180 1181 1182 1183 1184	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181 1182 1183 1184 1185	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA
1160 1161 1162 1163 1164 1165 1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177 1178 1179 1180 1181 1182 1183 1184 1185 1186	[Question]: Does the United States Secretary of State answer the phones for the White House? [Sub-question 1]: What are the duties of the US Secretary of State? [Sub-question 2]: Are answering phones part of #1? [Answer]: No  Test Input : Facts: Question: Sub-question: [INS] Answer: [INS] Figure 10: Prompt for StrategyQA

Method	Exemplars
LENS	Question: A cat chases a rat 6 hours after the rat runs. cat takes 4 hours to reach the rat. If
	average speed of the cat is 90 kmph, what s the average speed of the rat?
	<b>Options</b> : ['A)32kmph', 'B)26kmph', 'C)35kmph', 'D)36kmph', 'E)32kmph']
	Rationale: Cat take 10 hours and rat take 4 hoursthen Distance is 90*4.so speed of rat is (90*4)
	= 36kmph <b>Answer: D</b>
	Question: A business executive and his client are charging their dinner tab on the executive
	expense account. The?
	<b>Options</b> : ['A)69.55\$', 'B)50.63\$', 'C)60.95\$', 'D)52.15\$', 'E)53.15']
	Rationale: let x is the cost of the food 1.07x is the gross bill after including sales tax 1.
	1.07x=75 <b>Answer</b> : C
	Question: John and David were each given X dollars in advance for each day they were expecte
	perform at a community festival. John eventually,?
	<b>Options</b> : 'A)11Y', 'B)15Y', 'C)13Y', 'D)10Y', 'E)5Y' <b>Rationale: Answer:</b> A
	Question: A contractor undertakes to do a piece of work in 40 days. He engages 100 men at
	beginning and 100 more after 35 days and completes the work in stipulated time. If he had
	engaged the additional men, how many days behind schedule would it be finished??
	<b>Options</b> : 'A)2', 'B)5', 'C)6', 'D)8', 'E)9' <b>Rationale:</b> $[(100 \times 35)+(200 \times 5)]$ men can finish
	work in 1 day therefore 4500 men can finish the work in 1 day. 100 men can finish it in $\frac{4500}{100}$ =
	days. This is 5 days behind Schedule Answer: A
	Question: A can do a job in 9 days and B can do it in 27 days. A and B working together
	finish twice the amount of work in ——- days?
	<b>Options</b> : 'A)22 days', 'B)18 days', 'C)22 6/2 days', 'D)27 days', 'E)9 days' <b>Rationale:</b> 1
	$1/27 = 3/27 = 1/9 \ 9/1 = 9 \approx 2 = 18 \ day $ <b>Answer:</b> B
CASE	Question: In a 1000 m race, A beats B by 50 m and B beats C by 100 m. In the same race, by
	many meters does A beat C?
	<b>Options</b> : 'A)156 m', 'B)140 m', 'C)145 m', 'D)169 m', 'E)172 m' <b>Rationale:</b> By the time
	covers 1000 m, B covers $(1000 - 50) = 950$ m. By the time B covers 1000 m, C covers $(1000 - 100)$
	= 900 m. So, the ratio of speeds of A and C = $1000/950 * 1000/900 = 1000/855$ . So, by the tim
	covers 1000 m, C covers $855$ m. So in 1000 m race A beats C by $1000 - 855 = 145$ m. Answer
	<b>Question:</b> Count the numbers between 10 - 99 which yield a remainder of 3 when divided by 9
	also yield a remainder of 2 when divided by 5? <b>Options</b> : 'A)Two', 'B)Five', 'C)Six', 'D)Fo
	'E)One'
	<b>Rationale:</b> Numbers between $10 - 99$ giving remainder 3 when divided by $9 = 12, 21, 30, 39$
	57, 66, 75, 84, 93. The Numbers giving remainder 2 when divided by $5 = 12$ , $57 = 2$ Answer:
	Question: A train running at the speed of 60 km/hr crosses a pole in 3 seconds. Find the lengt
	the train. <b>Options:</b> (A)60°, (B)50°, (C)/5°, (D)100°, (E)120° <b>Rationale:</b> Speed = $60^{\circ}(5/18)$ m
	= 50/3 m/sec. Length of Irain (Distance) = Speed * Time $(50/3)$ * 3 = 50 meter. Answer: B
	Question: If n is an integer greater than 7, which of the following must be divisible by $3?$
	<b>Options:</b> A)1. $n(n+1)(n-4)^{2}$ , B)2. $n(n+2)(n-1)^{2}$ , C)3. $n(n+3)(n-5)^{2}$ , D)4. $n(n+4)(n-2)^{2}$ , C)5. $n(n+3)(n-5)^{2}$ , D)4. $n(n+4)(n-2)^{2}$ , D)5. $n(n+3)(n-5)^{2}$ , D)4. $n(n+4)(n-2)^{2}$ , D)5. $n(n+3)(n-5)^{2}$ , D)5. $n(n+3)(n-5)(n-5)(n-5)(n-5)(n-5)(n-5)(n-5)(n-5$
	n(n+3)(n-6) <b>Kalionale:</b> we need to find out the number which is divisible by three, in eve
	consecutive integers, there must contain 1 multiple of 3. So $n+4$ and $n+1$ are same if we nee
	nnd out the 3 s multiple. replace all the numbers which are more than or equal to three Ans
	D Overtions A merchant seise enlager in a herrein a contain sum. In a control herrein herrein
	<b>Question:</b> A merchant gains of loses, in a bargain, a certain sum. In a second bargain, he gains dellars, and in a third loses 20. In the end he finds he has pointed 120 dellars, he the three tasks
	uonars, and, in a united loses 20. In the end he finds he has gained 120 dollars, by the three toge
	now much did ne gain or lose by the first ( <b>Options</b> : A)80, B)-140, C)140, D)120, E)NO
	<b>Kauonaie:</b> In this sum, as the profit in morked 1, the loss must be I at y = the sum a minute of the profit is marked 1, the loss must be I at y = the sum a minute
	by contrary signs. If the profit is marked +, the loss must be Let $x =$ the sum required. I
	according to the statement $x + 280 - 20 = 120$ . And $x = -140$ . Answer: B

Table 6: Qualitative analysis of exemplars for AquaRat dataset selected by LENS vs CASE.
 Rationale is not completely shown for some questions to conserve space. However, in our experiments

all exemplars include rationales.

LENS	Table:  increase (decrease)   average yield   2.75% (2.75%)   volume   0.0 to 0.25   energy services   2013   fuel recovery fees   0.25   recycling processing and commodity sales   0.25 to 0.5   acquisitions / divestitures net   1.0   total change   4.25 to 4.75% (4.75%)   Question: what is the ratio of the acquisitions / divestitures net to the fuel recovery fees as part of the expected 2019 revenue to increase Rationale: ans=(1.0/0.25) Answer: The answer is 4Table:(in millions)   2009   2008   2007   sales and transfers of oil and gas produced net of production andadministrative costs   -4876 (4876)   -6863 (6863)   -4613 (4613)   Question: were total revisions of estimates greater than accretion of discounts? Rationale: Answer: The answer is yes Table:Table:12007   2008   change   capital gain distributions received   22.115.61-16.5 (16.5)   other than temporary impairments recognized   -3 (.3)   -91.3 (.91.3)   -91.0 (.91.0)   net gains (.losses) realized onfund dispositions   5.51 + 4.5 (.4.5)   -10.0 (.10.0)   net gain (.loss) Question: what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31, 2009? Rationale: (.1.4 / 2.2) Answer: The answer is 64%Table:in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   Question: north american printing papers net sales where what percent of total printing paper sales in 2009? Rationale: x0=(2.8 * 1000), ans=(.x0 * 5680) Answer: The answer is 49%Table:in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total drs   2351   2583   nonperforming   1119   1370   Question: what was the change in specific reserves in all between december 31, 2015 and decemb
CASE	services   2013   fuel recovery fees   0.25   recycling processing and commodity sales   0.25 to 0.5   acquisitions / divestitures net   1.0   total change   4.25 to 4.75% ( 4.75% )   <b>Question</b> : what is the ratio of the acquisitions / divestitures net to the fuel recovery fees as part of the expected 2019 revenue to increase <b>Rationale:</b> ans=(1.0 / 0.25 ) <b>Answer:</b> The answer is 4 <b>Table:</b> (in millions)   2009   2008   2007   sales and transfers of oil and gas produced net of production andadministrative costs   -4876 ( 4876 )   -6863 ( 6863 )   -4613 ( 4613 )   <b>Question</b> : were total revisions of estimates greater than accretion of discounts? <b>Rationale: Answer:</b> The answer is yes <b>Table:</b> [ 2007   2008   change   capital gain distributions received   22.1   5.61 - 16.5 ( 16.5 )   other than temporary impairments recognized  3 ( .3 )   -91.3 ( 91.3 )   -91.0 ( 91.0 )   net gains ( losses ) realized onfund dispositions   5.5   -4.5 ( 4.5 )   -10.0 ( 10.0 )   net gain ( loss ) <b>Question</b> : what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31 , 2009? <b>Rationale:</b> ( 1.4 / 2.2 ) <b>Answer:</b> The answer is 64% <b>Table:</b> in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   <b>Question</b> : north american printing papers net sales where what percent of total printing paper sales in 2009? <b>Rationale:</b> x0=( 2.8 * 1000 ), ans=( x0 * 5680 ) <b>Answer:</b> The answer is 49% <b>Table:</b> in millions   december 312015   december 312014   total consumer lending   1917   2041   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1191   1370   <b>Question</b> : what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? <b>Rationale:</b> (.34 ) <b>Answer:</b> The answer is -0.1 <b>Table:</b> in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 200
CASE	acquisitions / divestitures net   1.0   total change   4.25 to 4.75% ( 4.75% )   Question: what is the ratio of the acquisitions / divestitures net to the fuel recovery fees as part of the expected 2019 revenue to increase Rationale: ans=(1.0 / 0.25 ) Answer: The answer is 4 Table: ( in millions )   2009   2008   2007   sales and transfers of oil and gas produced net of production andadministrative costs   -4876 ( 4876 )   -6863 ( 6863 )   -4613 ( 4613 )   Question: were total revisions of estimates greater than accretion of discounts? Rationale: Answer: The answer is yes Table:  2007   2008   change   capital gain distributions received   22.1   5.6   -16.5 ( 16.5 )   other than temporary impairments recognized   -3 ( .3 )   -91.3 ( 91.3 )   -91.0 ( 91.0 )   net gain ( losse ) ) realized onfund dispositions   5.5   -4.5 ( 4.5 )   -10.0 ( 10.0 )   net gain ( loss ) Question: what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31 , 2009? Rationale: ( 1.4 / 2.2 ) Answer: The answer is 64% Table: in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   Question: north american printing papers net sales where what percent of total printing paper sales in 2009? Rationale: x0=( 2.8 * 1000 ), ans=( x0 * 5680 ) Answer: The answer is 49% Table: in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   Question: what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? Rationale: ( .34 ) Answer: The answer is -0.1 Table: in millions   total   balance december 31 2000   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008 411 payments   -38 ( 38 )   balance december 31 2009   \$ 3   Question:
- - CASE	Question: what is the ratio of the acquisitions / divestitures net to the fuel recovery fees as part of the expected 2019 revenue to increase <b>Rationale:</b> ans= $(1.0/0.25)$ <b>Answer:</b> The answer is 4 <b>Table:</b> (in millions   2009   2008   2007   sales and transfers of oil and gas produced net of production andadministrative costs   $-4876$ ( $4876$ )   $-6863$ ( $6863$ )   $-4613$ ( $4613$ )   <b>Question:</b> were total revisions of estimates greater than accretion of discounts? <b>Rationale:Answer:</b> The answer is yes <b>Table:</b>  2007   2008   change   capital gain distributions received   22.1   $5.6$   $-16.5$ ( $16.5$ )   other than temporary impairments recognized   $-3$ ( $.3$ )   $-91.3$ ( $91.3$ )   $-91.0$ ( $91.0$ )   net gains (losses ) realized onfund dispositions   $5.5$   $-4.5$ ( $4.5$ )   $-10.0$ ( $10.0$ )   net gain (loss ) <b>Question:</b> what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31 , 2009? <b>Rationale:</b> ( $1.4/2.2$ ) <b>Answer:</b> The answer is $64\%$ <b>Table:</b> in millions   2009   2008   2007   sales   $5680   6810   6530   operating profit   1091   474   839    Question: north american printing papers net sales where what percent of total printing paper sales in 2009? Rationale: x0=(2.8 * 1000), ans=(x0 * 5680 ) Answer: The answer is 49\%Table: in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370  Question: what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2007   46   additional provision   82   payments   -87 (78 )   balance december 31 2006   \$  124   payments   -78 (78 )   balance december 31 2009   \$  3  Question: in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest Rationale: ans=( 200 / 400 ) Answer$
CASE	the expected 2019 revenue to increase <b>Rationale:</b> ans= $(1.0/0.25)$ <b>Answer:</b> The answer is 4 <b>Table:</b> (in millions)   2009   2008   2007   sales and transfers of oil and gas produced net of production andadministrative costs   4876 ( 4876 )   -6863 ( 6863 )   -4613 ( 4613 )   <b>Question:</b> were total revisions of estimates greater than accretion of discounts? <b>Rationale: Answer:</b> The answer is yes <b>Table:</b>   2007   2008   change   capital gain distributions received   22.1   5.6   -16.5 ( 16.5 )   other than temporary impairments recognized  3 ( .3 )   -91.3 ( 91.3 )   -91.0 ( 91.0 )   net gains ( losses ) realized onfund dispositions   5.5   -4.5 ( 4.5 )   -10.0 ( 10.0 )   net gain ( loss ) <b>Question:</b> what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31 , 2009? <b>Rationale:</b> ( 1.4 / 2.2 ) <b>Answer:</b> The answer is 64% <b>Table:</b> in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   <b>Question:</b> north american printing papers net sales where what percent of total printing paper sales in 2009? <b>Rationale:</b> $x0=(2.8 \times 1000)$ , ans= $(x0 \times 5680)$ <b>Answer:</b> The answer is 49% <b>Table:</b> in millions   december 31 2015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   <b>Question:</b> what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? <b>Rationale:</b> ( .34 ) <b>Answer:</b> The answer is -0.1 <b>Table:</b> in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   411   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities inter
CASE	<b>Rationale:</b> ans= $(1.0/0.25)$ <b>Answer:</b> The answer is 4 <b>Table:</b> (in millions)   2009   2008   2007   sales and transfers of oil and gas produced net of production andadministrative costs   -4876 ( 4876 )   -6863 ( 6863 )   -4613 ( 4613 )   <b>Question:</b> were total revisions of estimates greater than accretion of discounts? <b>Rationale: Answer:</b> The answer is yes <b>Table:</b>   2007   2008   change   capital gain distributions received   22.1   5.6   -16.5 ( 16.5 )   other than temporary impairments recognized  3 ( .3 )   -91.3 ( 91.3 )   -91.0 ( 91.0 )   net gains ( losses ) realized onfund dispositions   5.5   -4.5 ( 4.5 )   -10.0 ( 10.0 )   net gain ( loss ) <b>Question:</b> what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31 , 2009? <b>Rationale:</b> ( 1.4 / 2.2 ) <b>Answer:</b> The answer is 64% <b>Table:</b> in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   <b>Question:</b> north american printing papers net sales where what percent of total printing paper sales in 2009? <b>Rationale:</b> $x0=(2.8 * 1000$ ), ans=( $x0 * 5680$ ) <b>Answer:</b> The answer is 49% <b>Table:</b> in millions   december 31 2015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   <b>Question:</b> what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? <b>Rationale:</b> ( .34 ) <b>Answer:</b> The answer is -0.1 <b>Table:</b> in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans=( 200 / 400 ) <b></b>
CASE	Table: (in millions )   2009   2008   2007   sales and transfers of oil and gas produced net of production andadministrative costs   -4876 (4876)   -6863 (6863)   -4613 (4613)   Question: were total revisions of estimates greater than accretion of discounts?Rationale: Answer: The answer is yesTable:   2007   2008   change   capital gain distributions received   22.1   5.6   -16.5 (16.5 )   other than temporary impairments recognized  3 (.3)   -91.3 (91.3)   -91.0 (91.0)   net gains (losses) realized onfund dispositions   5.5   -4.5 (4.5)   -10.0 (10.0)   net gain (loss)Question: what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31, 2009? Rationale: (1.4 / 2.2) Answer: The answer is 64%Table: in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   Question: north american printing papers net sales where what percent of total printing paper sales in 2009? Rationale: $x0=(2.8 * 1000)$ , $ans=(x0 * 5680)$ Answer: The answer is 49%Table: in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   Question: what was the change in specific reserves in all between december 31 , 2015 and december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   Question: in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest Rationale: ans=( 200 / 400 ) Answer: 0.5
CASE	production andadministrative costs $ -48/6 (48/6) -6863 (6863) -4613 (4613) $ <b>Question</b> : were total revisions of estimates greater than accretion of discounts? <b>Rationale</b> : <b>Answer</b> : The answer is yes <b>Table</b> : $ 2007 2008 $ change   capital gain distributions received $ 22.1 5.6 -16.5 (16.5) $ other than temporary impairments recognized $ -3 (.3) -91.3 (91.3) -91.0 (91.0) $ net gains (losses) is realized onfund dispositions $ 5.5 -4.5 (4.5) -10.0 (10.0) $ net gain (loss) <b>Question</b> : what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31, 2009? <b>Rationale</b> : $(1.4/2.2)$ <b>Answer</b> : The answer is $64\%$ <b>Table</b> : in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   <b>Question</b> : north american printing papers net sales where what percent of total printing paper sales in 2009? <b>Rationale</b> : $x0=(2.8 * 1000)$ , ans= $(x0 * 5680)$ <b>Answer</b> : The answer is $49\%$ <b>Table</b> : in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   <b>Question</b> : what was the change in specific reserves in all between december 31, 2015 and december 31, 2014 in billions? <b>Rationale</b> : $(.34)$ <b>Answer</b> : The answer is -0.1 <b>Table</b> : in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question</b> : in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale</b> : ans= $(200/400)$ <b>Answer</b> : 0.5
CASE	<b>Question:</b> were total revisions of estimates greater than accretion of discounts? <b>Rationale: Answer:</b> The answer is yes <b>Table:</b>   2007   2008   change   capital gain distributions received   22.1   5.6   -16.5 (16.5 )   other than temporary impairments recognized  3 (.3)   -91.3 (91.3)   -91.0 (91.0)   net gains (losses) realized onfund dispositions   5.5   -4.5 (4.5)   -10.0 (10.0)   net gain (loss) <b>Question:</b> what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31, 2009? <b>Rationale:</b> (1.4 / 2.2) <b>Answer:</b> The answer is 64% <b>Table:</b> in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   <b>Question:</b> north american printing papers net sales where what percent of total printing paper sales in 2009? <b>Rationale:</b> x0=(2.8 * 1000), ans=(x0 * 5680) <b>Answer:</b> The answer is 49% <b>Table:</b> in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   <b>Question:</b> what was the change in specific reserves in all between december 31, 2015 and december 31, 2014 in billions? <b>Rationale:</b> (.34) <b>Answer:</b> The answer is -0.1 <b>Table:</b> in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007 146   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans=( 200 / 400 ) <b>Answer:</b> 0.5 <b>Bable:</b>   2018   2017   2016   allewance for the ratio funde word during acartmetion   \$ 244 \$ 104 \$ 15   5
CASE	<b>Table:</b> 1 <b>Answer:</b> The answer is yes <b>Table:</b> $ 2007 2008 $ change   capital gain distributions received $ 22.1 5.6 -16.5(16.5) $ other than temporary impairments recognized $ -3(.3) -91.3(91.3) -91.0(91.0) $ net gains (losses) realized onfund dispositions $ 5.5 -4.5(4.5) -10.0(10.0) $ net gain (loss) <b>Question:</b> what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31, 2009? <b>Rationale:</b> ( $1.4/2.2$ ) <b>Answer:</b> The answer is $64\%$ <b>Table:</b> in millions $ 2009 2008 2007 $ sales $ 5680 6810 6530 $ operating profit $ 1091 474 $ 839  <b>Question:</b> north american printing papers net sales where what percent of total printing paper sales in 2009? <b>Rationale:</b> $x0=(2.8*1000)$ , ans= $(x0*5680)$ <b>Answer:</b> The answer is $49\%$ <b>Table:</b> in millions $ $ december $312015 $ december $312014 $ total consumer lending $ $ 1917 2041  total commercial lending $ $ 434 542  total tdrs $ $ 2351 2583  nonperforming $ $ 1119 1370  <b>Question:</b> what was the change in specific reserves in alll between december 31, 2015 and december 31, 2014 in billions? <b>Rationale:</b> ( $.34$ ) <b>Answer:</b> The answer is -0.1 <b>Table:</b> in millions   total   balance december 31 2006   \$ 124   payments   -78(78)   balance december 31 2007   46  additional provision   82   payments   -87(87)   balance december 31 2008   41  payments   -38(38)   balance december 31 2009   \$ 3   <b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans= $(200/400)$ <b>Answer:</b> 0.5 <b>Bable:</b> $ 2018 2017 2016 $ allownerge for other funds used during construction   \$ 2445 1045 15
CASE	table: $12007 + 2008 + \text{change} + \text{capital gain distributions received} + 22.1 + 3.0 + 16.3 (+16.3) + otherthan temporary impairments recognized  3 (.3)   -91.3 (91.3)   -91.0 (91.0)   net gains (losses)realized onfund dispositions   5.5   -4.5 (4.5)   -10.0 (10.0)   net gain (loss)Question: what percentage of tangible book value is made up of cash and cash equivalents andmutual fund investment holdings at december 31, 2009? Rationale: (1.4 / 2.2) Answer: Theanswer is 64%Table: in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474  839  Question: north american printing papers net sales where what percent of total printing paper salesin 2009? Rationale: x0=(2.8 * 1000), ans=(x0 * 5680) Answer: The answer is 49%Table: in millions   december 312015   december 312014   total consumer lending   1917   2041  total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370  Question: what was the change in specific reserves in all between december 31, 2015 anddecember 31, 2014 in billions? Rationale: (.34) Answer: The answer is -0.1Table: in millions   total   balance december 31 2006   $ 124   payments   -78 ( 78 )   balancedecember 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008  41   payments   -38 ( 38 )   balance december 31 2009   $ 3  Question: in 2006 what was the ratio of the class a shares and promissory notes international papercontributed in the acquisition of borrower entities interest Rationale: ans=( 200 / 400 ) Answer: 0.5$
CASE	that temporary impartments tecoginized 15 ( $5$ ) 1-91.5 ( $.91.5$ ) ( $.91.$
CASE	<b>Question:</b> what percentage of tangible book value is made up of cash and cash equivalents and mutual fund investment holdings at december 31, 2009? <b>Rationale:</b> $(1.4 / 2.2)$ <b>Answer:</b> The answer is 64% <b>Table:</b> in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   <b>Question:</b> north american printing papers net sales where what percent of total printing paper sales in 2009? <b>Rationale:</b> $x0=(2.8 * 1000)$ , ans= $(x0 * 5680)$ <b>Answer:</b> The answer is 49% <b>Table:</b> in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   <b>Question:</b> what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? <b>Rationale:</b> $(.34)$ <b>Answer:</b> The answer is -0.1 <b>Table:</b> in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans= $(200 / 400)$ <b>Answer:</b> 0.5
CASE	Table: in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   Question: north american printing papers net sales where what percent of total printing paper sales in 2009? Rationale: $x0=(2.8 \times 1000)$ , ans= $(x0 \times 5680)$ Answer: The answer is 49% Table: in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   Question: what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? Rationale: $(.34)$ Answer: The answer is -0.1 Table: in millions   total   balance december 31 2006   $\$$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   $\$$ 3   Question: in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest Rationale: ans= $(200 / 400)$ Answer: 0.5 Table: $  2017   2016   allaurenees for other funds used during construction   \$ 24   \$ 10   \$ 15  $
CASE	Table: in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474   839   Question: north american printing papers net sales where what percent of total printing paper sales in 2009? Rationale: $x0=(2.8 \times 1000)$ , ans= $(x0 \times 5680)$ Answer: The answer is 49% Table: in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   Question: what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? Rationale: $(.34)$ Answer: The answer is -0.1 Table: in millions   total   balance december 31 2006   $$$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   $$$ 3   Question: in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest Rationale: ans= $(200 / 400)$ Answer: 0.5 Table: $  -2018   -2016   $
CASE	Table: in millions   2009   2008   2007   sales   5680   6810   6530   operating profit   1091   474  839  Question: north american printing papers net sales where what percent of total printing paper salesin 2009? Rationale: $x0=(2.8 * 1000)$ , ans= $(x0 * 5680)$ Answer: The answer is 49%Table: in millions   december 312015   december 312014   total consumer lending   1917   2041  total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370  Question: what was the change in specific reserves in all between december 31 , 2015 anddecember 31 , 2014 in billions? Rationale: (.34 ) Answer: The answer is -0.1Table: in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balancedecember 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3  Question: in 2006 what was the ratio of the class a shares and promissory notes international papercontributed in the acquisition of borrower entities interest Rationale: ans=( 200 / 400 ) Answer:0.5Table: $  -2018   -2016   -allaurence for other function and during construction   $ 24   $ 10   $ 15  $
CASE	<b>Question:</b> north american printing papers net sales where what percent of total printing paper sales in 2009? <b>Rationale:</b> $x0=(2.8 \times 1000)$ , ans= $(x0 \times 5680)$ <b>Answer:</b> The answer is 49% <b>Table:</b> in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   <b>Question:</b> what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? <b>Rationale:</b> (.34) <b>Answer:</b> The answer is -0.1 <b>Table:</b> in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans= $(200 / 400)$ <b>Answer:</b> 0.5
CASE	Question: north american printing papers net sales where what percent of total printing paper sales in 2009? Rationale: $x0=(2.8 \times 1000)$ , ans= $(x0 \times 5680)$ Answer: The answer is 49% Table: in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   Question: what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? Rationale: $(.34)$ Answer: The answer is -0.1 Table: in millions   total   balance december 31 2006   $$ 124  $ payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   $$ 3  $ Question: in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest Rationale: ans= $( 200 / 400 )$ Answer: 0.5
CASE	in 2009? <b>Rationale:</b> x0=(2.8 * 1000), ans=(x0 * 5680) <b>Answer:</b> The answer is 49% <b>Table:</b> in millions   december 312015   december 312014   total consumer lending   1917   2041   total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   <b>Question:</b> what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? <b>Rationale:</b> (.34) <b>Answer:</b> The answer is -0.1 <b>Table:</b> in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans=( 200 / 400 ) <b>Answer:</b> 0.5 <b>Table:</b> $  2018   2017   2016   alleurones for other funds used during construction   $ 24   $ 10   $ 15   Answer:$
CASE	Table: in millions   december 312015   december 312014   total consumer lending   1917   2041  total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370  Question: what was the change in specific reserves in all between december 31 , 2015 anddecember 31 , 2014 in billions? Rationale: (.34 ) Answer: The answer is -0.1Table: in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balancedecember 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3  Question: in 2006 what was the ratio of the class a shares and promissory notes international papercontributed in the acquisition of borrower entities interest Rationale: ans=( 200 / 400 ) Answer:0.5Table: 12018   2017   2016   alleurones for other funds used during construction   \$ 24 \$ 10 \$ 15 \$ 15 \$ 15 \$ 15 \$ 15 \$ 15 \$ 15
CASE	total commercial lending   434   542   total tdrs   2351   2583   nonperforming   1119   1370   Question: what was the change in specific reserves in all between december 31 , 2015 and december 31 , 2014 in billions? Rationale: (.34 ) Answer: The answer is -0.1 Table: in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   Question: in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest Rationale: ans=( 200 / 400 ) Answer: 0.5 Table: $  2018   2017   2016   alleurones for other funds used during construction   $ 24   $ 10   $ 15   }$
CASE	Question: what was the change in specific reserves in all between december 31, 2015 and december 31, 2014 in billions? Rationale: (.34) Answer: The answer is -0.1 Table: in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   Question: in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest Rationale: ans=( 200 / 400 ) Answer: 0.5 Table:   2018   2017   2016   allowanes for other funds used during construction   \$ 24   \$ 10   \$ 15
CASE	december 31, 2014 in billions? <b>Rationale:</b> (.34) <b>Answer:</b> The answer is -0.1 <b>Table:</b> in millions   total   balance december 31 2006   \$ 124   payments   -78 ( 78 )   balance december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans=( 200 / 400 ) <b>Answer:</b> 0.5 Table:   2018   2017   2016   allowanes for other funds used during construction   \$ 24   \$ 10   \$ 15
CASE	<b>Table</b> : in millions   total   balance december 31 2006   $$ 124   payments   -78 (78)   balance december 31 2007   46   additional provision   82   payments   -87 (87)   balance december 31 2008   41   payments   -38 (38)   balance december 31 2009   $ 3  Question: in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest Rationale: ans=(200 / 400) Answer: 0.5$
	december 31 2007   46   additional provision   82   payments   -87 ( 87 )   balance december 31 2008   41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question</b> : in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans=( 200 / 400 ) <b>Answer:</b> 0.5
	41   payments   -38 ( 38 )   balance december 31 2009   \$ 3   <b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans=( 200 / 400 ) <b>Answer:</b> 0.5 Tables   2018   2018   2016   allowanes for other funds used during construction   \$ 24   \$ 10   \$ 15
	<b>Question:</b> in 2006 what was the ratio of the class a shares and promissory notes international paper contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans= $(200 / 400)$ <b>Answer:</b> 0.5
	contributed in the acquisition of borrower entities interest <b>Rationale:</b> ans= $(200/400)$ <b>Answer:</b> 0.5
	0.5 Tables 1 2018 1 2017 1 2016 1 allower as for other funds used during construction $   ^{6} 24   ^{6} 10   ^{6} 15   _{10}$
	' oblog 1' MIV 1' MIV 1' MIA Lollow on a ton other tunde wood dyming construction 1V' MIV 1V 1V 15 1
	<b>Table:</b> $12018 + 2017 + 2010$ + anowance for other funds used during construction + $524 + 519 + 515$
	allowance for borrowed funds used during construction   15   8   6
	Question: by now much did anowance for other funds used during construction increase from $2016$ to $20182$ <b>Deticipale</b> : $x_0 = (24, 15)$ and $(24, 15)$ another and $(24, 15)$ and $(24$
	<b>Table:</b> (dollars in millions) $2001(1)$ 2000 $1999(2)$   change 00-01   adjusted change 00-01 (
	3) servicing fees $ $ \$ 1624   \$ 1425   \$ 1170   14% (14%)   14% (14%)   management fees
	511 + 581 + 600 + -12 + (12) + -5 + (5) + foreign exchange trading + 368 + 387 + 306 + -5 + (5) +
	processing fees and other $  329   272   236   21   21   total fee revenue   2832   2665   \$ 2312   6  $
	8   Question: what is the growth rate in total fee revenue in 2001? Rationale: x0=(2832 - 2665
	),ans=( x0 / 2665 ) <b>Answer:</b> 6.30%
	Table:   increase (decrease)   average yield   2.75% (2.75%)   volume   0.0 to 0.25   energy services
	2013   fuel recovery fees   0.25   recycling processing and commodity sales   0.25 to 0.5   acquisitions
	/ divestitures net   1.0   total change   4.25 to 4.75% ( 4.75% )   Question: what is ratio of insurance
	recovery to incremental cost related to our closed bridgeton landfill <b>Rationale:</b> ans=(40.0/12.0)
	Answer: 3.33
	<b>Table:</b> \$ in millions   as of december 2018   as of december 2017   fair value of retained interests
	\$ 3151   \$ 2071   weighted average life ( years )   7.2   6.0   constant prepayment rate   11.9% ( 11.9
	%)   9.4% ( 9.4 % )   impact of 10% ( 10 % ) adverse change   \$ -27 ( 27 )   \$ -19 ( 19 )   impact of
	20% ( $20%$ ) adverse change   \$ -53 ( 53 )   \$ -35 ( 35 )   discount rate   $4.7%$ ( $4.7%$ )   $4.2%$ (
	4.2%   1mpact of 10% (10%) adverse change   \$ -75 (75)   \$ -35 (35)   impact of 20% (20)
	%) adverse change $ \$-147(147) \$-70(70) $ Question: what was the change in fair value of
	%) adverse change $ \$-147(147) \$-70(70) $ Question: what was the change in fair value of retained interests in billions as of december 2018 and december 2017? Rationale: ans= $(3.28 - 212)$ Answer 1.15

Table 7: Qualitative analysis of exemplars for FinQA dataset selected by LENS vs CASE. Rationale
 is not completely shown for some questions to conserve space. However, in our experiments all exemplars include rationales.

Method	Exemplars
LENS	Question: Michael wants to dig a hole 400 feet less deep than twice the depth of the hole that his
	father dug. The father dug a hole at a rate of 4 feet per hour. If the father took 400 hours to dig his
	hole?
	<b>Rationale:</b> Since the father dug a hole with a rate of 4 feet per hour, if the father took 400 hours
	digging the hole, he dug a hole $4*400 = 1600$ feet deep Michael will have to work for $2800/4 = 7001$
-	700 hours. Answer: 700
	<b>Question:</b> when Effick went to the market to sell his fruits, he realized that the price of lemons had risen by 4 for each lemon. The price of grapes had also increased by helf the price that $-2$
	Itself by 4 for each lemon. The price of grapes had also increasing by 4 is $8 \pm 4 - 12$ For the 80 lemons
	Example 140 * 9 = 1260 From the sale of all of his fruits. Erick received $1260 + 960 =$
	2220. Answer: 2220
-	Question: James decides to build a tin house by collecting 500 tins in a week. On the first day, he
	collects 50 tins. On the second day, he manages to collect 3 times that number?
	Rationale: On the second day, he collected 3 times the number of tins he collected on the first day
	which is $3 * 50 = 150$ tins he'll need to collect $200/4 = 50$ tins per day to reach his goal.
	Answer: 50
-	Question: Darrel is an experienced tracker. He can tell about an animal by the footprints it leaves
	behind. Based on the impressions, he could tell the animal was traveling east at 15 miles/hour
	<b>Rationale:</b> If we let x be the amount of time, in hours, it will take for Darrel to catch up to the
-	coyote, If we subtract 1 x from each side, we get x=1, the amount of time in hours. Answer:
	<b>Question:</b> Martina needs to paint all four walls in ner 12 foot by 10 foot kitchen, which has 10 foot kitchen, which has 10 foot bigh callings If Martha can paint 40 square fast per hour, how many hours will it take her to
	paint kitchen? <b>Pationale:</b> There are two walls that are 12' by 10' and two walls that are 16' by 10
	by how many hours she needs to finish: $1680 \text{ sq ft} / 40 \text{ sq ft} / hour = 42 \text{ hours Answer: } 42$
G 4 G F	
LASE	<b>Question:</b> Each class uses 200 sneets of paper per day. The school uses a total of 9000 sneets of paper avery weak. If there are 5 days of school days, how many classes are there in the school?
	<b>Bationale:</b> Each class uses $200 \times 5 = 1000$ sheets of paper in a week. Thus, there are $9000/1000 =$
	9 classes in the school <b>Answer</b> 9
	<b>Ouestion:</b> If Jenna has twice as much money in her bank account as Phil does, and Phil has
	one-third the amount of money that Bob has in his account, and Bob has \$60 in his account, how
	much less money does Jenna have in her account than Bob? Rationale: If Phil has one-third o
	the amount that Bob does, so he has \$60/3= \$20 in his account. Since Jenna has twice as much
	money as Phil, so she has \$20*2= 40 in her account. Since Bob has \$60 in his account, so he has
	\$60-\$40=\$20 more than Jenna. <b>Answer</b> : 20
	Question: Carlos bought a box of 50 chocolates. 3 of them were caramels and twice as many were
	nougais. The number of truffles was equal to the number of caramels plus 6 If Carlos picks a
	find the number of nougats by doubling the number of caramals 2 caramals * 2 nougats/caramal
	find the number of nougats by doubling the number of carameters. 5 catallets $\cdot$ 2 flougats/catallet = 6 nougats. Then find the number of truffles by adding 7 to the number of carameter 3 carameter = 1
	= 9 Answer: 64
	<b>Ouestion:</b> Janet has 60 less than four times as many siblings as Masud Carlos has 3/4 times a
	many siblings as Masud. If Masud has 60 siblings, how many more siblings does Janet have mor
	than Carlos? <b>Rationale:</b> If Masud has 60 siblings, and Carlos has 3/4 times as many siblings a
	Masud, Carlos has 3/4*60=45 siblings. Four times as many siblings as Masud has is 4*60=240
	Janet has 60 less than four times as many siblings as Masud, a total of 240-60=180 sibling
	180-45=135 Answer: 135
	Question: Gavin has had 4 dreams every day for a year now. If he had twice as many dreams last
	year as he had this year, calculate the total number of dreams he's had in the two years. Rational
	If Gavin has been having 4 dreams every day for a year now, he has had $4*365 = 1460$ dreams this
	year. Gavin had twice as many dreams last as he had this year, meaning he had $2*1460 = 2920$
	dreams last year. The total number of dreams he has had in the two years is 2920+1460=4380

Table 8: Qualitative analysis of exemplars for **GSM8K** dataset selected by LENS vs CASE. Rationale

is not completely shown for some questions to conserve space. However, in our experiments all exemplars include rationales.

1352	Method	Exemplars
1353 1354 1355	LENS	<b>Table</b> :   Name   Age (years)   Jessica   2   Dalton   7   Kelsey   5   Lamar   8   Alexis   2 <b>Question</b> : A girl compared the ages of her cousins. What is the median of the numbers? <b>Rationale</b> : Read the numbers from the table: 2, 7, 5, 8, 2. First, arrange the numbers from least to greatest: 2, 2, 5, 7, 8.
1356		Now find the number in the middle. The number in the middle is 5. The median is 5. Answer: 5
1357		Table:         City   Number of houses sold   Melville   878   New Hamburg   871   Charles Falls   881
1358		Pennytown   81 / Question: A real estate agent looked into how many houses were sold in different cities. Where were the fewest houses sold? <b>Pationale</b> : Find the least number in table. The least
1359		number is 817 Now find the corresponding city. Pennytown corresponds to 817 Answer: 817
1360		<b>Table:</b>   Day   Number of new customers   Saturday   2   Sunday   2   Monday   9   Tuesday   4
1301		Wednesday   10   Thursday   3   Friday   6 Question: A cable company analyst paid attention to how
1302		many new customers it had each day. What is the median of the numbers? Rationale: Find the
1267		number in the middle. The number in the middle is 4. The median is 4. Answer: 4
1365		Table:         Day   Number of cups   Friday   8   Saturday   4   Sunday   10   Monday   6   Tuesday   6           West of the second secon
1366		wednesday   1   1 hursday   0 Question: Nancy wrote down how many cups of lemonade she sold in the past 7 days. What is the range of the numbers? <b>Patianale:</b> Pased the numbers from the tables
1367		8 4 10 6 6 1 0 Subtract the least number from the greatest number: 10-0-10 The range is
1368		10 <b>Answer:</b> 10
1369		Table:         Price         Quantity demanded         Quantity supplied         \$700         9,800         22,600         \$740         8,000         22,800
1370		\$780   6,200   23,000  \$820   4,400   23,200  \$860   2,600   23,400 Question: At a price of \$860, is
1371		there a shortage or a surplus? Rationale: At price of \$860, quantity demanded is less than quantity
1372		supplied So, there is a surplus. Answer: surplus
1373	CASE	Table: Number of siblings   Frequency 0   19 1   12 2   13 3   9
1374		Question: The students in Mr. Robertson's class recorded the no. of siblings that each has. How
1375		many students have fewer than 2 siblings? <b>Rationale:</b> Find the rows for 0 and 1 sibling. Add the
1376		trequencies for these rows. $19 + 12 = 31$ , $31$ students have fewer than 2 siblings. Answer: $31$
1377		blind taste test on some of her friends in order to determine if organic fruits tasted different than
1378		non-organic fruits. Each friend ate one type of fruit. What is the probability that a randomly
1379		selected friend preferred organic and tasted peaches? Rationale: Let A be the event "the friend
1380		preferred organic" and B be the event "the friend tasted peaches" Answer: Jul-19
1381		Table:       dance performance ticket   \$29.00 play ticket   \$32.00 figure skating ticket   \$41.00 ballet
1382		ticket   \$37.00 opera ticket   \$76.00 orchestra ticket   \$58.00 Question: How much money does
1203		Hannan need to buy a ballet ticket and / orcnestra tickets? <b>Kationale:</b> Find the cost of / orcnestra tickets $\$58,00 \$7 = \$406,00$ . Now find the total cost $\$37,00 + \$406,00 = \$443,00$ . Hannah needs
1385		\$443.00 <b>Answer:</b> 443
1386		Table: Price   Quantity demanded   Quantity supplied \$665   15,500   16,200 \$855   13,700   17,300
1387		\$1,045   11,900   18,400 \$1,235   10,100   19,500 \$1,425   8,300   20,600 <b>Question</b> : Look at the
1388		table. Then answer the question. At a price of \$1,045, is there a shortage or a surplus? Rationale:
1389		At the price of \$1,045, the quantity demanded is less than the quantity supplied. There is too much
1390		of the good or service for sale at that price. So, there is a surplus. Answer: surplus
1391		<b>Iable:</b> Comfy Pillows Resort   4:15 A.M   2:30 P.M   10:00 P.M Skyscraper City   4:45 A.M   3:00 P.M   10:20 P.M   10:20 P.M Pilloregaster Land
1392		15:45 A M   4:00 PM   11:30 PM Floral Gardens   6:45 A M   5:00 PM   12:30 A M Chickenville
1393		7:15 A.M   5:30 P.M   1:00 A.M Happy Cow Farm   7:45 A.M   6:00 P.M   1:30 A.M
1394		Question: Look at the following schedule. Marshall got on the train at Rollercoaster Land at
1395		5.45 A.M. What time will he get to Floral Gardens? Rationale: Find 5:45 A.M. in the row for
1396		Rollercoaster Land. That column shows the schedule for the train that Marshall is on. Look down
1397		the column until you find the row for Floral Gardens. Marshall will get to Floral Gardens at 6:45
1398		A.M. Answer: 6:45 A.M.
1399		

Table 9: Qualitative analysis of exemplars for TabMWP dataset selected by LENS vs CASE.
 Rationale is not completely shown for some questions to conserve space. However, in our experiments

1402 all exemplars include rationales.

Method	Exemplars
LENS	Facts: Penguins are native to the deep, very cold parts of the southern hemisphere. Miami is locat
	in the northern hemisphere and has a very warm climate.
	Question: Would it be common to find a penguin in Miami?
	Rationale: Where is a typical penguin's natural habitat? What conditions make #1 suitable f
	penguins? Are all of #2 present in Miami? Answer: No
	<b>Facts:</b> Shirley Bassey recorded the song Diamonds are Forever in 1971. Over time, diamon
	degrade and turn into graphite. Graphite is the same chemical composition found in pencils.
	<b>Question:</b> Is the title of Shirley Bassey's 19/1 diamond song a true statement? <b>Rationale</b> : what
	the fifthe to Shirley Bassey's 1971 diamond song? Do diamonds fast for the period in #1? Answe
	NO Factor The first six numbers in the Fibonacci sequence are 1.1.2.3.5.8. Since 1 is doubled there is
	only five different single digit numbers <b>Ouestion</b> : Are there five different single digit Fibona
	numbers?
	<b>Rationale:</b> What are the single-digit numbers in the Fibonacci sequence? How many uniq
	numbers are in #1? Does #2 equal 5? Answer: Yes
	<b>Facts:</b> Katy Perry's gospel album sold about 200 copies. Katy Perry's most recent pop albums so
	over 800,000 copies. Question: Do most fans follow Katy Perry for gospel music? Rationa
	What type of music is Katy Perry known for? Is Gospel music the same as #1? Answer: No
	Facts: The Italian Renaissance was a period of history from the 13th century to 1600. A theocra
	is a type of rule in which religious leaders have power. Friar Girolamo Savonarola was
	ruler of Florence, after driving out the Medici family, from November 1494 †23 May 14
	Question: Was Florence a Theocracy during Italian Renaissance? Rationale: When was the Ital
	Renaissance?When did Friar Girolamo Savonarola rule Florence? Is #2 within the span of #1? I
	Friar Girolamo Savonarola belong to a religious order during #3? Answer: Yes
CASE	Facts: U2 is an Irish rock band that formed in 1976. The Polo Grounds was a sports stadium the
	was demolished in 1964. Question: Did U2 play a concert at the Polo Grounds? Rationale: Wh
	was U2 (Irish rock band) formed? When was the Polo Grounds demolished? Is #1 before #
	Answer: No
	<b>Facts:</b> The capacity of Tropicana Field is 36,973. The population of Auburn, NY is 27,6
	<b>Question:</b> Can you fill every resident of Auburn, New York, in Tropicana Field? <b>Rationale</b> : while the population of Auburn NV2 is #1 greater than t
	Answer: Ves
	<b>Facts:</b> Door to door advertising involves someone going to several homes in a residential area
	make sales and leave informational packets <b>Ouestion</b> : During the pandemic is door to do
	advertising considered inconsiderate? <b>Rationale:</b> What does door to door advertising involv
	person to do? During the COVID-19 pandemic, what does the CDC advise people to do in terms
	traveling? Does doing #1 go against #2 and #3? <b>Answer:</b> Yes
	<b>Facts:</b> Mosquitoes cannot survive in the climate of Antarctica. Zika virus is primarily spre
	through mosquito bites. Question: Do you need to worry about Zika virus in Antarctica? Rationa
	What animal spreads the Zika Virus? What is the climate of Antarctica? Can #1 survive in #
	Answer: No
	Facts: Bob Marley had 9 children. Kublai Khan had 23 children. Many of Bob Marley's children
	became singers, and followed his themes of peace and love. The children of Kublai Khan follow
	in his footsteps and were fierce warlords. Question: Could Bob Marley's children hypothetica
	win tug of war against Kublai Khan's children? Rationale: How many children did Bob Mar
	have? How many children did Kublai Khan have? Is #1 greater than #2? <b>Answer:</b> No