
AutoMix: Mixing Models with Few-shot Self and Meta Verification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models (LLMs) are now available in various sizes and configura-
2 tions from cloud API providers. While this diversity offers a broad spectrum of
3 choices, effectively leveraging the options to optimize computational cost and per-
4 formance remains challenging. In this work, we present AutoMix, an approach that
5 strategically routes queries to larger LMs, based on the approximate correctness of
6 outputs from a smaller LM. Central to AutoMix is a few-shot self-verification mech-
7 anism, which estimates the reliability of its own outputs without requiring training.
8 Given that verifications can be noisy, we employ a meta verifier in AutoMix to re-
9 fine the accuracy of these assessments. Our experiments using LLAMA2-13/70B,
10 on five context-grounded reasoning datasets demonstrate that AutoMix surpasses
11 established baselines, improving the incremental benefit per cost by up to 57%.

12 1 Introduction

13 Human problem-solving inherently follows a multi-step process: generate a solution, verify its
14 validity, and refine it further based on verification outcomes. The emulation of this self-refinement and
15 reflective behavior has gained attention in the recent research (Pan et al., 2023; Madaan et al., 2023;
16 Reid and Neubig, 2022; Schick et al., 2022; Welleck et al., 2022; Shinn et al., 2023). Classic self-refine
17 paradigms consistently employ a singular model across all problem-solving stages, demonstrating
18 effectiveness in certain scenarios (Madaan et al., 2023; Welleck et al., 2022). Yet, the intrinsic
19 complexity and variability of tasks, from simplistic (e.g., binary classification on separable data)
20 to complex (e.g., code generation) and potentially unsolvable (e.g., certain forms of multi-step
21 reasoning), motivate an alternative approach. This approach iteratively queries over models of
22 disparate sizes and capabilities, verifying feedback at each step and determining whether to accept
23 the output or route to a more capable, albeit computationally intensive, model (Liu et al., 2020; Zhou
24 et al., 2020; Madaan and Yang, 2022; Geng et al., 2021; Schuster et al., 2022).

25 Past studies in model-switching strategies predominantly rely on separate models trained explicitly
26 for each step or require access to logits (Welleck et al., 2022; Reid and Neubig, 2022), which may not
27 always be feasible as modern LLMs rely on access to black box APIs. To address these challenges,
28 we propose a new method, which we call AutoMix. In contrast to existing approaches, AutoMix
29 fully leverages black-box LLM APIs, avoiding the need for separate models or access to logits
30 using few-shot learning (Brown et al., 2020) and meta-verification. Our method proposes strategies
31 for each step of problem-solving: solution generation, verification, and routing, all assuming we
32 only have access to black-box LLMs. In contrast to existing approaches, which generally delineate
33 tasks as Simple or Complex for model routing, AutoMix integrates a third category of *Unsolvable*
34 queries. These queries are likely unsolvable even by a Large Language Model (LLM) and should
35 not be routed to larger models if identified early enough. This consideration allows AutoMix to
36 judiciously allocate computational resources, preventing unwarranted computational spending on
37 these particularly challenging instances.

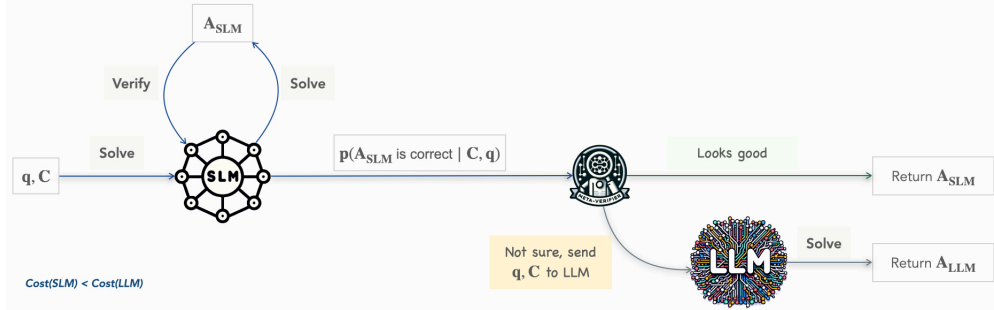


Figure 1: AutoMix: Given a context C and question q , an initial answer A_{SLM} is generated with the smaller language model (SLM). A_{SLM} is verified by the SLM, yielding a noisy verification score. Based on the meta-verifier’s decision, either A_{SLM} is returned if deemed satisfactory, or the task is rerouted to a larger language model (LLM) to enhance accuracy. The looping arrow around the SLM symbolizes self-correction before proceeding to the meta-verifier.

38 We use context-grounded few-shot entailment to quantify the uncertainty in an answer’s correct-
 39 ness (Poliak, 2020; Dagan et al., 2022). However, recognizing that verifications can sometimes be
 40 inconsistent or noisy, we introduce a *meta-verifier* to evaluate the reliability of the initial verification.
 41 The meta-verifier acts as a secondary check, providing an additional layer of confidence assessment
 42 to ensure that the decision to route a task to a larger or smaller model is well-founded.

43 In summary, our contributions are:

- 44 • We introduce AutoMix, a method that strategically leverages black-box LLM APIs for generating
 45 a solution, verifying the solution, and switching to a larger language model, everything without
 46 access to model weights, gradients, or logits.
- 47 • We also show that context-grounded entailment is a reasonable albeit noisy proxy for self-
 48 verification. To deal with this noise, we propose a POMDP-based meta-verification mechanism
 49 that helps improve the reliability of the final decision.
- 50 • We propose and introduce the *Incremental Benefit Per Cost* (IBC) metric, a novel measure that
 51 quantifies the efficiency of integrating smaller and larger language models.
- 52 • We present empirical evidence from experiments on five context-grounded reasoning datasets
 53 using the language models LLAMA2-13B and LLAMA2-70B as the SLM and LLM. Our results
 54 demonstrate that AutoMix surpasses established baselines, enhancing the incremental benefit per
 55 cost by up to 57%.

56 2 AutoMix: Few-shot Self-Verification and Meta-Verification

```

Context: {context}

Question: {question}

AI Generated Answer: {generated_answer}

Instruction: Your task is to evaluate if the AI Generated Answer is correct, based
↪ on the provided context and question. Provide the judgement and reasoning for
↪ each case. Choose between Correct or Incorrect.

Evaluation: "

```

Listing 1: **Verification Prompt.** The verification process is framed as a natural language entailment task, where the model determines the validity of the model-generated answer with respect to the context and question. We use a generic few-shot prompt for all tasks (prompt in appendix C.1).

57 **Task and setup** We tackle the problem of context-grounded question answering, where given
58 a context \mathcal{C} (e.g., stories, newswire, or research article) and a question q , the model is tasked to
59 generate an accurate and coherent answer, consistent with the provided context. We deploy two
60 distinct models: a smaller, cost-efficient model, denoted as SLM, and a larger, more accurate yet
61 costly model, LLM. Our objective is to optimize performance while staying economical. We use a
62 verifier, \mathcal{V} , to ascertain the validity of SLM’s outputs and decide if a query should be redirected to
63 LLM. We start by generating an initial answer, \mathcal{A}_s , using the smaller SLM. Next, we need to assess
64 the trustworthiness of \mathcal{A}_s . To this end, we use a few-shot verifier. Our choice of tasks is motivated
65 by two key concerns. First, longer queries are more computationally demanding, underscoring the
66 need for an approach like AutoMix to navigate the cost-accuracy trade-off. Second, Context allows
67 the verifier to cross check the preliminary answers with available information, aiding in identifying
68 inconsistencies as ungrounded is challenging (Pan et al., 2023; Huang et al., 2023).

69 Verification is framed as an entailment task (Poliak, 2020; Dagan et al., 2022). The objective is to
70 determine if the answer generated by SLM aligns with the provided context. Specifically, the verifier
71 gauges $v = p(\text{correct} = 1 \mid \mathcal{A}_s, \mathcal{C}, q)$, with $\text{correct} = 1$ indicating that \mathcal{A}_s is correct. The verification
72 prompt is outlined in Figure 1. We use the same verification prompt for all tasks.

73 2.1 Meta-verifier

74 Given the potential inconsistency or noise in verifier outcomes, a secondary evaluation mechanism,
75 which we term the *meta-verifier*, is crucial to vet the verifier’s conclusions. In particular, the verifier
76 is tasked with determining whether the SLM’s answer is entailed by the context, and this decision
77 is made without considering the inherent difficulty of the problem. Notably, routing unsolvable
78 queries for the LLM is resource-inefficient without enhancing performance. While ascertaining the
79 ground truth of query difficulty is non-trivial, verification probability and historical data can provide
80 insightful guidance. Formally, we define the meta-verifier’s outputs as $m(v, \mathcal{A}_s, \mathcal{C}, q) \rightarrow \{0, 1\}$,
81 where $m = 1$ implies the verifier’s output can be trusted.

82 Addressing the notable challenges of self-correction in large language models (Madaan et al., 2023;
83 Huang et al., 2023), our method employs a non-LLM setup for meta-verification to avoid escalating
84 issues like hallucination and reasoning errors (Dziri et al., 2023). The versatile meta-verifier can
85 adopt various advanced learning strategies, from supervised to reinforcement learning, explored
86 further in upcoming sections. Subsequent sections provide a deeper exploration into two particular
87 implementations of this strategy.

88 **Thresholding** In this simplistic meta-verifier approach, the decision is made based on probability
89 of verifier being correct with a threshold t , defined as $H(t) = 0$ for $t < 0$ and $H(t) = 1$ for $t \geq 0$.
90 For black-box language models, the probability of correctness can be derived by sampling $k > 1$
91 samples at a higher sampling temperature.

92 **Using a POMDP** In the context of meta-verifier, we observe that all the queries in this two language
93 model setup could be categorized in three different categories: *Simple*, *Complex*, and *Unsolvable*.
94 The simple queries are addressable by SLM itself, the complex queries are addressable by LLM but
95 not by SLM and Unsolvable queries are so complex that they can’t be addressed by either LLM or
96 SLM. Hence, a ground truth oracle should route only the complex queries but not unsolvable queries.
97 Since the ground truth state is not known and unobserved, we formulate this decision problem as
98 a Partially Observable Markov Decision Process (POMDP) (Monahan, 1982). POMDP presents a
99 robust framework, offering a structured way to manage and navigate through the decision spaces
100 where the system’s state is not fully observable. A POMDP is defined by a tuple (S, A, T, R, Ω, O) ,
101 where S is a set of states, A is a set of actions, T represents the state transition probabilities, R is the
102 reward function, Ω is a set of observations, and O is the observation function. See Appendix A.1 for
103 more details.

104 Another advantage of the POMDP-based meta-verifier is its interpretability and customizability via
105 reward assignment. For instance, in a *Complex* state, assigning a reward of +50 for invoking the
106 LLM indicates a preference for accurate solutions over computational cost. Although the POMDP
107 framework inherently handles sequences of decisions, we confine our approach to a single-decision
108 scenario (horizon or episode length 1) for simplicity, with potential for extension to streaming settings
109 for optimizing across multiple queries or a fixed time duration.

```

procedure ANSWERQUERY( $\mathcal{C}, q$ )
  ▷ Context, Question
   $\mathcal{A}_s \leftarrow \text{solve}(\text{SLM}, \mathcal{C}, q)$ 
   $v \leftarrow \text{self-verify}(\mathcal{A}_s, \mathcal{C}, q)$ 
  if  $m(v, \mathcal{A}_s, \mathcal{C}, q)$  then
    ▷ Check meta-verifier decision
    if  $v \geq 0.5$  then
      return  $\mathcal{A}_s$ 
    else
      return  $\text{solve}(\text{LLM}, \mathcal{C}, q)$ 
    end if
  else
    return  $\mathcal{A}_s$ 
  end if
end procedure

```

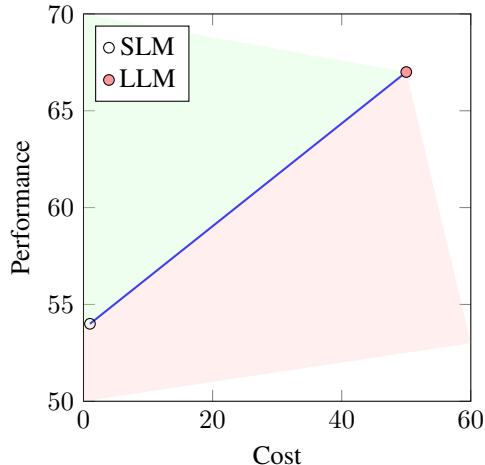


Figure 2: **Left:** AutoMix algorithm. **Right:** Performance vs. Cost curve. The slope between SLM and LLM provides a way to the Incremental Benefit per Cost (IBC) for methods that mix models. Methods with a steeper slope than this reference when plotted against SLM have a positive IBC (green region), whereas those below the reference have a negative IBC (red region), falling into the red region.

110 3 Cost-Performance Efficiency Analysis

111 In our approach to leveraging model performance, it is essential to consider not only the raw accuracy
 112 of predictions but also the associated computational or monetary costs. To that end, we introduce a
 113 metric to understand the efficiency of the models in terms of cost. We use C_M and P_M to denote the
 114 cost and performance of a method M . We also use C_{SLM} and C_{LLM} , and P_{SLM} and P_{LLM} , to denote the
 115 cost and performance of using the SLM and LLM, respectively.

116 **Incremental Benefit Per Cost (IBC)** We introduce methods, denoted by M , to optimally in-
 117 tegrate SLM and LLM. For each method M , we associate a cost C_M and performance P_M . To
 118 quantify the utility of M over SLM, we define the metric *Incremental Benefit Per Cost* (IBC) as
 119 IBC_M (Equation (1)).

$$\text{IBC}_M = \frac{P_M - P_{\text{SLM}}}{C_M - C_{\text{SLM}}}, \quad \text{IBC}_{\text{BASE}} = \frac{P_{\text{LLM}} - P_{\text{SLM}}}{C_{\text{LLM}} - C_{\text{SLM}}}, \quad \Delta_{\text{IBC}}(M) = \frac{\text{IBC}_M - \text{IBC}_{\text{BASE}}}{\text{IBC}_{\text{BASE}}} \times 100 \quad (1)$$

120 The IBC metric captures the efficiency of performance enhancement relative to the additional cost.
 121 For comparative evaluation, we set a baseline IBC, IBC_{BASE} , representing the benefit of *always* using
 122 LLM over SLM. Finally, we compare methods using Δ_{IBC} , which compares the IBC of a specific
 123 method with IBC_{BASE} . A positive IBC lift suggests that M achieves performance increments more
 124 cost-effectively than a standalone LLM, whereas a negative lift indicates reduced efficiency (Figure 2)

125 **Cost Calculation** The total cost, C_M , for a method M utilizing the Small Language Model (SLM)
 126 for initial answer generation and verification, and the Large Language Model (LLM) as needed, is
 127 computed as: $C_M = 2 \cdot C_{\text{SLM}} + w_{\text{LLM}} \cdot C_{\text{LLM}}$. Here, $w_{\text{LLM}} \in [0, 1]$ represents the proportion of LLM
 128 usage, with values indicating exclusive ($w_{\text{LLM}} = 1$) or no usage ($w_{\text{LLM}} = 0$) of LLM. While we
 129 utilize SLM for verification, it is worth noting that a different verifier model could also be employed,
 130 which would alter the cost calculations accordingly.

131 4 Experiments

132 **Setup** We experiment with open-source pair LLAMA2-13B and LLAMA2-70B (Touvron et al.,
 133 2023). We assume a cost of 1 unit for the SLM, and 50 units for the LLM, following the price disparity
 134 between the small and large models offered by LLM API providers like OpenAI and Together¹.

¹<https://openai.com/pricing>, <https://together.ai/>

Method	CNLI			Quality			QASPER			NarrativeQA			COQA		
	C	P	Δ_{IBC}	C	P	Δ_{IBC}	C	P	Δ_{IBC}	C	P	Δ_{IBC}	C	P	Δ_{IBC}
SLM	1	40.1	-	1	29.8	-	1	14.0	-	1	20.3	-	1	48.1	-
FrugalGPT	37.4	59.2	66.1	49.7	42.0	-2.1	49.3	27.7	-1.1	45.9	26.0	2.5	30.3	57.1	13.1
SC	43.6	51.2	-17.3	11.8	32.3	-9.8	46.6	27.5	2.6	23.4	23.1	1.2	16.9	54.6	49.6
AutoMix +T	51.9	55.4	-4.5	24.8	36.1	3.5	38.4	24.9	-12.0	13.2	21.9	2.4	8.3	51.0	31.7
AutoMix +P	5.5	42.3	<u>57.0</u>	9.6	32.1	4.0	45.4	27.4	5.0	10.3	21.5	3.6	7.9	50.8	42.5
LLM	50	55.5	-	50	42.3	-	50	28.1	-	50	26.4	-	50	61.4	-

Table 1: **Main Results:** highlighting the trade-offs between Cost (C), Performance (P), and Incremental Benefit per Cost (Δ_{IBC}) across various methods and datasets. The acronyms represent: SLM - Small Language Model, LLM- Large Language Model, AutoMix + T and AutoMix + P - variations of our proposed method with thresholding (T) and POMDP (P) based meta-verifiers, respectively. AutoMix + **POMDP** demonstrates a robust and consistent Δ_{IBC} across the **Quality**, QASPER, NARRATIVE-QA, and COQA datasets, implying a judicious utilization of computational resources.

135 **Datasets** We experiment with a diverse set of datasets: NARRATIVE-QA (Kočíský et al., 2018)
136 for full-length book and movie script QA, QASPER (Dasigi et al., 2021) for research paper QA,
137 CNLI (Koreeda and Manning, 2021) for NLI tasks, QUALITY (Pang et al., 2022) for multiple-choice
138 questions from long articles, and COQA (Reddy et al., 2019) for conversational comprehension QA.
139 The datasets are evaluated on F1 score, accuracy, and exact match based on answer format.

140 **Baselines** Our baselines include: i) **Verifier Self-Consistency** (Wang et al., 2022), where we
141 prompt SLM with our entailment verifier and draw 8 samples (temperature 0.7). The majority label
142 routes the query to LLM or not. We cache the KV values for the (long) input prompt, so only a
143 single forward pass is done. The cost of this verifier is the same as the cost of SLM. and ii) **Frugal**
144 **GPT (F) (Chen et al., 2023)** We finetune a DistillBert (Sanh et al., 2019) as a verifier, outputting
145 a confidence probability for a given question, context, and SLM-generated answer, with a verifier
146 confidence threshold directing query routing and its cost set to 0 due to significantly lower operational
147 costs than SLM. Both approaches adhere to a low-resource setting, utilizing 1000 training examples
148 per dataset.

149 **Proposed approaches** We experiment with two different types of meta-verifiers: threshold and
150 POMDP-based. i) **AutoMix + Thresholding:** Using a threshold on the verifier probability e.g.,
151 $Thresh=0.75$ implies using SLM outputs with confidence ≥ 0.75 and LLM. We use a threshold for
152 each dataset that yields the highest Δ_{IBC} on the validation set. ii) **AutoMix + POMDP:** This method
153 optimizes routing decisions using a POMDP solver (Smith and Simmons, 2006), given verifier outputs
154 and observation probabilities learned on the validation set (detailed in Appendix A.1).

155 4.1 Main Results

156 Table 1 shows the meta-verifier method consistently showcases superior performance in terms of
157 Δ_{IBC} across both LLAMA2-13/70B. On QUALITY,QASPER,NARRATIVE-QA,COQA, AutoMix beat
158 FrugalGPT despite the latter having access to domain-specific training and low verifier cost. In
159 Figure 3 (left), we present the performance of our model, AutoMix, across various cost intervals.
160 Our findings reveal that AutoMix-POMDP shows consistent positive Δ_{IBC} across all evaluated costs.
161 This suggests that our method can deliver consistent improvements, regardless of the user’s desired
162 cost or performance requirements. Further, in Figure 3 (right), we compare the accuracy of using
163 POMDP based meta-verifier over Verifier-SC. We see significant improvements across all datasets,
164 with relative gains of up to 33% demonstrating our proposed meta-verifier’s importance in few-shot
165 verification setups.

166 4.2 Key findings and takeaway

167 **AutoMix is Effective in Low-Resource Scenarios** Figure 6 demonstrates the performance dy-
168 namics of AutoMix and FrugalGPT with varying validation sizes. Notably, our method significantly
169 outperforms FrugalGPT with limited data (under 2000 samples), despite the latter’s domain-specific
170 training and zero verifier cost. However, as training data increases, FrugalGPT narrows the perfor-

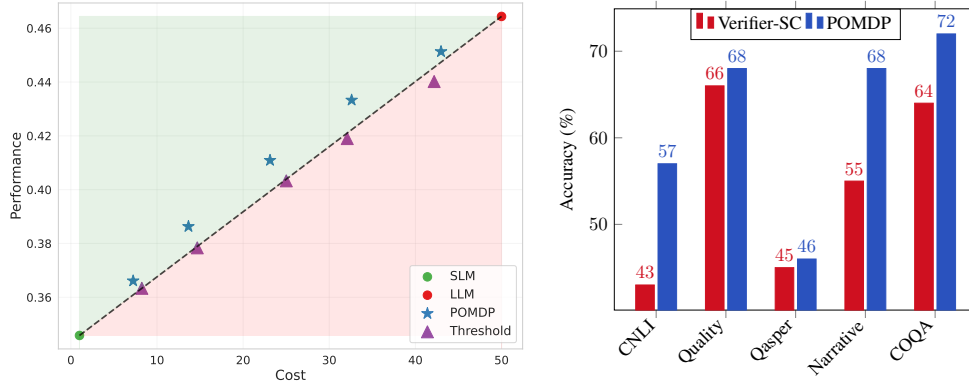


Figure 3: **Left:** Aggregated performance vs. cost for different methods on the small and large LLAMA2-13/70B. POMDP based meta-verifier is consistently in the green region, signifying a higher Incremental Benefit per Cost (IBC). **Right:** The accuracy of the meta-verifier for both POMDP and Verifier-Self-Consistency (Verifier-SC) approaches across various datasets. Across all scenarios, the POMDP method consistently wins with up to 33% relative performance gains.

171 mance gap by leveraging domain-specific training. This pattern indicates that AutoMix provides a
 172 particularly advantageous solution in real-world scenarios where data may be scarce.

173 **Effectiveness of Few-shot Self-Verification** In Appendix B.1, we evaluate few-shot self-
 174 verification quantitatively and qualitatively. We observe that the self-verification can effectively
 175 use context to identify errors in answers generated by SLM in many cases.

176 **Improving Self-Verification with Task-Specific Prompt Engineering** We explore the impact of
 177 task-specific prompt engineering on self-verification performance in Appendix B.2. While prompt
 178 engineering improves verifier accuracy, our meta-verifier remains robust in various settings and can
 179 beneficially leverage even a weak verifier.

180 5 Related Work

181 **1. Mixing Models** Distinct from related work optimizing LLM inference cost by model switching and
 182 external verifiers (Chen et al., 2023; Zhu et al., 2023; vSakota et al., 2023), AutoMix obviates the need
 183 for verifier training through few-shot SLM model prompting and does not require upfront access to all
 184 input queries. **2. Adaptive Computation** In contrast to adaptive computation methods that preempt
 185 computation via intermediate representations (Liu et al., 2020; Zhou et al., 2020; Schuster et al.,
 186 2021; Geng et al., 2021; Schuster et al., 2022), AutoMix necessitates no architectural modifications.
 187 Further, unlike AdaptiveConsistency (Aggarwal et al., 2023), which optimizes inference within a
 188 single LLM model, AutoMix flexibly optimizes between two models and transcends its utility in
 189 Self-Consistency. **3. Self-Verification** AutoMix aligns in spirit with works that aim to perform
 190 self-verification for reasoning problems, such as Weng et al. (2023); Pan et al. (2023). However,
 191 AutoMix uniquely harnesses context for verification instead of relying on LLM’s knowledge, and
 192 introduces a meta-verifier mechanism to offset the verifier’s potential noise.

193 6 Conclusion

194 AutoMix integrates black-box large language model (LLM) APIs into a multi-step problem-solving
 195 framework, optimizing the computational cost and performance trade-offs. AutoMix opens avenues
 196 for several interesting research directions. First, while self-verification and correction are challenging
 197 for LLMs in general, we find promising results using context-grounded few-shot verification, indi-
 198 cating that similar approaches may yield gain in other scenarios. Secondly, our work interweaves
 199 Good Old-Fashioned Artificial Intelligence (GOF AI) approaches with LLMs, demonstrating that the
 200 incorporation of a POMDP can boost the accuracy of a noisy few-shot verifier, showing the promise
 201 of this paradigm as an approach for improving LLMs during inference.

202 References

- 203 Pranjali Aggarwal, Aman Madaan, Yiming Yang, and Mausam. 2023. Let’s sample step by step:
204 Adaptive-consistency for efficient reasoning with llms. *ArXiv*, abs/2305.11860.
- 205 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
206 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
207 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,
208 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray,
209 Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever,
210 and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural
211 Information Processing Systems*, volume 33, pages 1877–1901, Online. Curran Associates, Inc.
- 212 Lingjiao Chen, Matei A. Zaharia, and James Y. Zou. 2023. Frugalgpt: How to use large language
213 models while reducing cost and improving performance. *ArXiv*, abs/2305.05176.
- 214 Ido Dagan, Dan Roth, Fabio Zanzotto, and Mark Sammons. 2022. *Recognizing textual entailment:
215 Models and applications*. Springer Nature.
- 216 Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. 2021. A
217 dataset of information-seeking questions and answers anchored in research papers. In *Proceedings
218 of the 2021 Conference of the North American Chapter of the Association for Computational
219 Linguistics: Human Language Technologies*, pages 4599–4610.
- 220 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jian, Bill Yuchen Lin, Peter
221 West, Chandra Bhagavatula, Ronan Le Bras, Jena D Hwang, et al. 2023. Faith and fate: Limits of
222 transformers on compositionality. *arXiv preprint arXiv:2305.18654*.
- 223 Shijie Geng, Peng Gao, Zuohui Fu, and Yongfeng Zhang. 2021. Romebert: Robust training of
224 multi-exit bert. *arXiv preprint arXiv:2101.09755*.
- 225 Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song,
226 and Denny Zhou. 2023. Large language models cannot self-correct reasoning yet. *arXiv preprint
227 arXiv:2310.01798*.
- 228 Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis,
229 and Edward Grefenstette. 2018. The narrativeqa reading comprehension challenge. *Transactions
230 of the Association for Computational Linguistics*, 6:317–328.
- 231 Yuta Koreeda and Christopher D Manning. 2021. Contractnli: A dataset for document-level natural
232 language inference for contracts. In *Findings of the Association for Computational Linguistics:
233 EMNLP 2021*, pages 1907–1919.
- 234 Teven Le Scao and Alexander M Rush. 2021. How Many Data Points is a Prompt Worth? In *NAACL*.
- 235 Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021a.
236 What Makes Good In-Context Examples for GPT-3? *arXiv:2101.06804 [cs]*. ArXiv: 2101.06804.
- 237 Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021b.
238 Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language
239 Processing. *arXiv preprint arXiv:2107.13586*.
- 240 Weijie Liu, Peng Zhou, Zhiruo Wang, Zhe Zhao, Haotang Deng, and Qi Ju. 2020. Fastbert: a
241 self-distilling bert with adaptive inference time. In *Proceedings of the 58th Annual Meeting of the
242 Association for Computational Linguistics*, pages 6035–6044.
- 243 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,
244 Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement
245 with self-feedback. *arXiv preprint arXiv:2303.17651*.
- 246 Aman Madaan and Yiming Yang. 2022. Flowgen: Fast and slow graph generation. *arXiv preprint
247 arXiv:2207.07656*.
- 248 Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. 2021. Re-
249 framing Instructional Prompts to GPTk’s Language. *arXiv preprint arXiv:2109.07830*.

- 250 George E. Monahan. 1982. State of the art—a survey of partially observable markov decision
251 processes: Theory, models, and algorithms. *Management Science*, 28:1–16.
- 252 Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang
253 Wang. 2023. Automatically correcting large language models: Surveying the landscape of diverse
254 self-correction strategies. *arXiv preprint arXiv:2308.03188*.
- 255 Richard Yuanzhe Pang, Alicia Parrish, Nitish Joshi, Nikita Nangia, Jason Phang, Angelica Chen,
256 Vishakh Padmakumar, Johnny Ma, Jana Thompson, He He, et al. 2022. Quality: Question
257 answering with long input texts, yes! In *Proceedings of the 2022 Conference of the North American
258 Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages
259 5336–5358.
- 260 Adam Poliak. 2020. A survey on recognizing textual entailment as an nlp evaluation. *arXiv preprint
261 arXiv:2010.03061*.
- 262 Siva Reddy, Danqi Chen, and Christopher D Manning. 2019. Coqa: A conversational question
263 answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266.
- 264 Machel Reid and Graham Neubig. 2022. Learning to model editing processes. *arXiv preprint
265 arXiv:2205.12374*.
- 266 Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled
267 version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- 268 Timo Schick, Jane Dwivedi-Yu, Zhengbao Jiang, Fabio Petroni, Patrick Lewis, Gautier Izacard,
269 Qingfei You, Christoforos Nalmpantis, Edouard Grave, and Sebastian Riedel. 2022. Peer: A
270 collaborative language model.
- 271 Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Q Tran, Yi Tay, and
272 Donald Metzler. 2022. Confident adaptive language modeling. *arXiv preprint arXiv:2207.07061*.
- 273 Tal Schuster, Adam Fisch, Tommi Jaakkola, and Regina Barzilay. 2021. Consistent accelerated
274 inference via confident adaptive transformers. In *Proceedings of the 2021 Conference on Empirical
275 Methods in Natural Language Processing*, pages 4962–4979.
- 276 Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023. Reflexion: an autonomous agent with
277 dynamic memory and self-reflection.
- 278 Trey Smith and Reid Simmons. 2006. Focused real-time dynamic programming for mdps: Squeezing
279 more out of a heuristic. In *AAAI*, pages 1227–1232.
- 280 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
281 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open
282 and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- 283 Marija vSakota, Maxime Peyrard, and Robert West. 2023. Fly-swat or cannon? cost-effective
284 language model choice via meta-modeling. *ArXiv*, abs/2308.06077.
- 285 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. 2022.
286 Self-consistency improves chain of thought reasoning in language models. *ArXiv*, abs/2203.11171.
- 287 Sean Welleck, Ximing Lu, Peter West, Faeze Brahman, Tianxiao Shen, Daniel Khashabi, and Yejin
288 Choi. 2022. Generating sequences by learning to self-correct. *arXiv preprint arXiv:2211.00053*.
- 289 Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Kang Liu, and Jun Zhao. 2023. Large
290 language models are better reasoners with self-verification. *CoRR*, abs/2212.09561.
- 291 Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. 2020. Bert loses
292 patience: Fast and robust inference with early exit. *Advances in Neural Information Processing
293 Systems*, 33:18330–18341.
- 294 Banghua Zhu, Ying Sheng, Lianmin Zheng, Clark W. Barrett, Michael I. Jordan, and Jiantao
295 Jiao. 2023. On optimal caching and model multiplexing for large model inference. *ArXiv*,
296 abs/2306.02003.

297 A Methodology

298 A.1 POMDP

299 The Partially Observable Markov Decision Process (POMDP) presents a robust framework for
300 handling decision-making problems under uncertainty, offering a structured way to manage and
301 navigate through the decision spaces where the system’s state is not fully observable (Monahan, 1982).
302 A POMDP is defined by a tuple (S, A, T, R, Ω, O) , where S is a set of states, A is a set of actions, T
303 represents the state transition probabilities, R is the reward function, Ω is a set of observations, and
304 O is the observation function.

305 In the context of meta-verifier, the *unobservable* states (S) represent the potential correctness of the
306 verifier’s predictions, categorized as *Simple*, *Complex*, and *Insolvable*. Actions (A) are binary: trust
307 the verifier or invoke the LLM. The reward function (R) quantifies the cost or gain of making
308 a particular action in a given state, steering the decision policy towards cost-effective actions.
309 Observations (Ω) in our model are the verifier’s probability outputs, discretized into bins. Specifically,
310 we generate $k=8$ samples from the verifier, discretizing our observation space in intervals of size
311 0.125 ranging from 0 to 1.

312 The observation function (O) depicts the likelihood of observing an observation given an action was
313 taken and the system transitioned to a particular state. Using an appropriate observation function is
314 crucial for POMDP to work. Specifically, we define observations probabilities in two ways:

- 315 • **1. Functional Form:** For each of the states s , the observation function O is defined as
316 $O(s, v) = \frac{1}{K} \cdot v^{\gamma_s}$, where v is the verifier probability and $\gamma_s \in [0, \infty]$ is a hyperparameter
317 for every state and K is normalizing factor. Intuitively, a value of γ close to 1 indicates ideal
318 calibration, with verifier probability v indicating true probability of being in a particular state.
319 The values of γ_s ’s for the three states are determined based on the respective POMDP’s
320 performance on validation set based on the IBC-Lift.
- 321 • **2. Discrete Form:** An alternate option is to directly learn observation function O from
322 the statistics of validation set. Since in validation set, we have access to the true state
323 along with verifier probabilities of individual data instances, we can model observation
324 function as $O(s, v) = \frac{\sum_{i=0}^N \mathbf{1}\{s_i=s \text{ and } v_i=v\}}{\sum_{i=0}^N \mathbf{1}\{s_i=s\}}$. The method has the advantage of being
325 hyperparameter free and provides more accurate representation by computing the true
326 observation probabilities on validation set. However, it performs worse than functional form,
327 when either certain values of v or s are not well represented in validation set or in cases of
328 high distribution shift between validation and test set.

329 Since both these methods have their strengths, and are independent of each other, we choose the best
330 performing method on validation set.

331 This POMDP mechanism allows for optimal decision-making under uncertainty, balancing the cost
332 and reliability of invoking the LLM. Through employing standard POMDP solving algorithms such
333 as Focused Real-Time Dynamic Programming² (Smith and Simmons, 2006), we derive a policy that
334 maps belief states (probability distributions over S) to actions. During inference, the learned policy
335 effectively decides whether to trust the verifier’s output or to invoke the LLM based on a combination
336 of expected future rewards and computational costs.

337 Another advantage of the POMDP-based meta-verifier is its interpretability and customizability via
338 reward assignment. For instance, in a "Needy" state, assigning a reward of +50 for invoking the LLM
339 indicates a preference for accurate solutions over computational cost. Conversely, in a "Good" state,
340 designating a reward of -10 for trusting the SLM encourages computational savings. This enables
341 users to strategically balance solution quality against computational expenses, aligning with specific
342 application needs.

²We use zmdp package <https://github.com/trey0/zmdp> for solving POMDP

```

# Meta-verifier POMDP File for narrative_qa

discount: 0.99
values: reward

# We have 6 states: 3 corresponding to the initial state before verifier is
  called, and 3 corresponding to the state after verifier is called
states: START_S START_C START_U SIMPLE COMPLEX UNSOLVABLE

# Effectively, we have 3 actions: 1.) The initial State where we run verifier
  2.) Report SLM's Answer 3.) Invoke LLM and Report its Answer
actions: Init Trust_SLM Invoke_LLM

# Observations lies in one of verifier probability bins. Eg: bin_correct_high
  represents Verifier outputs SLM answer as correct with high confidence
observations: bin_incorrect_low bin_incorrect_high bin_correct_low
  bin_correct_high

# Transition Model for Init action

T: Init
# Format: start_state : end_state : Transition_Probability

# Transition Model for Trust_SLM action
T: Trust_SLM
identity

# Transition Model for Invoke_LLM action
T: Invoke_LLM
identity

# Observation Model after "Init" action for narrative_qa
# Format: 0 : action : state : observation : probability

# Example: In SIMPLE cases, it is likely, SLM is correct and Verifier is
  Confident, while in UNSOLVABLE, SLM is incorrect (Lower Obs. Probability)
0 : * : SIMPLE : bin_correct_high 0.8
0 : * : COMPLEX : bin_correct_high 0.4
0 : * : UNSOLVABLE : bin_correct_high 0.1

# Reward Model:
# Format: R: action : init_state : end_state : observation : probability

# Example: For COMPLEX state, Trusting SLM results in negative score, while
  invoking LLM results in a high +50 score.
R: Trust_SLM : COMPLEX : * : * -10
R: Invoke_LLM : COMPLEX : * : * +50

```

Figure 4: A sample POMDP specification file. POMDP requires defining states, actions, observations and relevant Transition, Observation Probabilities and Reward Values.

343 **B Verifier Qualitative Analysis**

344 **B.1 How effective is few-shot verification?**

345 One of the contributions of this work is few-shot self-verification of outputs. While Self-Verification is
 346 challenging in general for reasoning problems, our setup offers an advantage: the context allows us to
 347 check for the validity of the answer using the context. For example, the model can spot factual errors
 348 in the answer, or apparent contradictions that may not be clear when the question is answered for the
 349 first time. Are the models able to use this advantage for effective self-verification? Figure 5 shows
 350 that for all datasets except CNLI, few-shot self-verification effectively identifies correct example,
 351 assigning a higher probability to them.

352 **Qualitative Analysis** Representative Examples from our qualitative analysis are shown in Tables 2,
 353 3, and 4.

Attribute	Content
Context	Which country grows the most tea? ... Shen sipped it. He liked it. He drank it all. Shen was proud of his new drink. <i>Truncated for brevity</i>
Question	When did he take it?
Generated Answer	He took it in 1990.
Gold Answer	about 2737 B.C
Verifier Output	The context does not mention that he took tea in 1990. Verification Decision: The AI generated answer is Incorrect.

Table 2: An example where the generated answer is verifiably incorrect, citing an unsupported date (1990), and the verifier successfully catches this discrepancy.

Attribute	Content
Context	Fort Lauderdale, Florida (CNN) – Just taking a sip of water or walking to the bathroom is excruciatingly painful for 15-year-old Michael Brewer, ... He suffered second- and third-degree burns over about two-thirds of his body, according to the hospital’s associate director, Dr. Carl Schulman. <i>Truncated for brevity</i>
Question	How long before he recovers from his injury?
Generated Answer	he will never recover from his injury.
Gold Answer	He has a lifelong recovery.
Verifier Output	The context clearly states that the teen is recovering from his injuries. Verification Decision: The AI generated answer is Incorrect.

Table 3: The example demonstrates a challenging case for the verifier. The generated answer is a rephrasing of what’s mentioned in the context, however the verifier fails to catch it.

354 **FrugalGPT vs. AutoMix at different levels of data availability** Figure 6 demonstrates the
 355 performance dynamics of AutoMix and FrugalGPT with varying validation data sizes. Notably, our
 356 method significantly outperforms FrugalGPT with limited data (under 2000 samples), despite the
 357 latter’s domain-specific training and zero verifier cost. However, as training data increases, FrugalGPT
 358 narrows the performance gap by leveraging its domain-specific training. This pattern indicates that
 359 AutoMix provides a particularly advantageous solution in real-world scenarios where data may be
 360 scarce.

Attribute	Content
Context	Pruning keeps the plant only three or four feet tall. This is an easy height for tea picking. Green tea is steamed right after the leaves are picked. Green tea _ its healthy qualities. For example, it may prevent heart disease. How did we get tea bag? The answer: by accident. <i>Truncated</i>
Question	What did they do to green tea after picking it?
SLM Answer	steamed.
Gold Answer	prune it
Verifier Output	The context mentions that green tea is steamed after being picked. Verification Decision: The AI generated answer is Correct.

Table 4: An instance where the verifier deems the answer correct, and the gold label was incorrect. The verifier is able to correctly infer that the context mentions tea being steamed after picking.

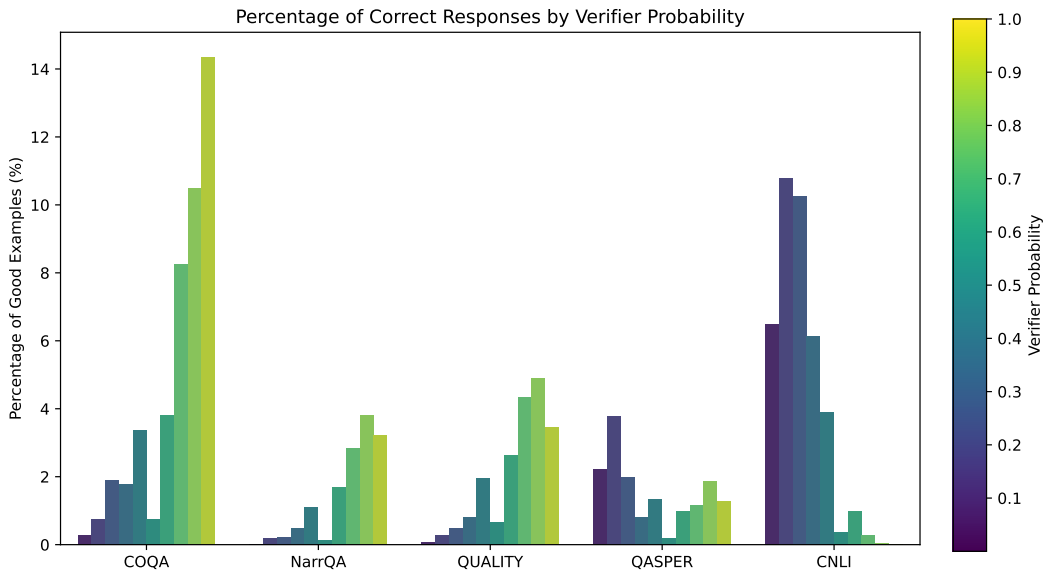


Figure 5: **Verifier Probability and Correctness:** Percentage of correct responses across distinct verifier probability bins, representing $P(C = 1 | A_{\text{SLM}}, C, q)$, where A_{SLM} is the answer from the Small Language Model, C is the context, and q is the query. Each bin represents a range of verifier probabilities and the corresponding accuracy of the responses within that probability range across various datasets. Notably, for all datasets, excluding CNLI and QASPER, a higher verification score generally corresponds to a larger proportion of correct examples, indicating that the verifier is, to an extent, capable of discerning the reliability of responses generated by itself. We use a meta-verifier to get around these noisy predictions.

361 B.2 Domain-specific vs. Domain independent verifier

362 We used a single verifier with the LLAMA2-13B model to help steer the model. To avoid excessive
 363 prompt engineering, we used a generic prompt for all datasets. However, task-specific prompts
 364 generally help (Le Scao and Rush, 2021; Liu et al., 2021b; Mishra et al., 2021; Liu et al., 2021a). To
 365 investigate this, we create task specific prompts for CNLI by giving examples from legal domain in
 366 the prompt.

367 Figure 7 underscores the efficacy of employing task-specific verification prompts, ensuring a height-
 368 ened probability allocation for accurate examples during the verification process. Interestingly,
 369 the enhanced verifier accuracy does not always directly translate to proportionate improvements in
 370 our proposed method, AutoMix, as evidenced in Table 5. This phenomenon highlights the role of
 371 meta-verifiers, adeptly negotiating through the outputs of potentially unreliable verifiers.

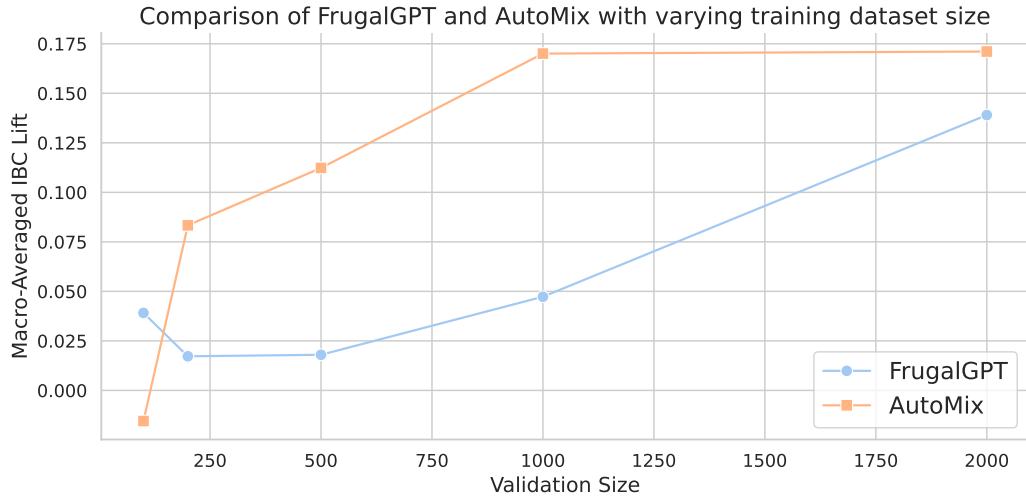


Figure 6: Comparison of AutoMix with FrugalGPT over varying Training Dataset Size. Despite zero-cost verifier and domain-specific training, FrugalGPT underperforms AutoMix. AutoMix is especially useful for limited data settings, with higher gains visible when dataset size is less than 1000.

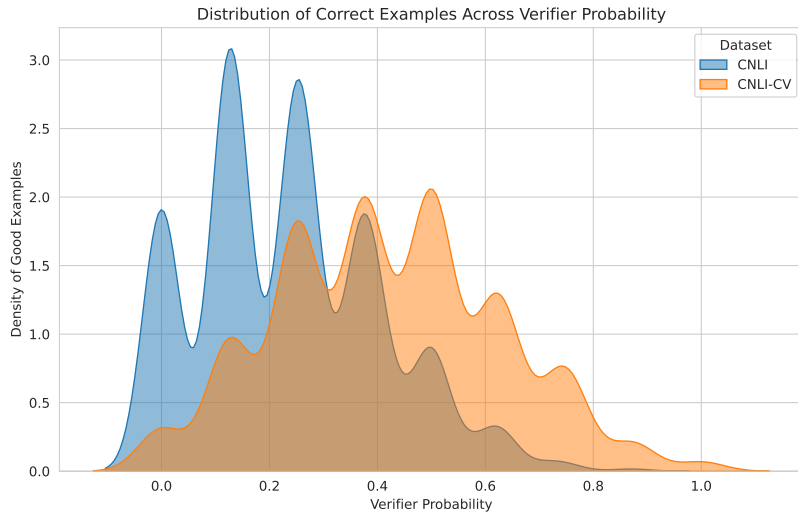


Figure 7: Enhancement of verifier accuracy using task-specific verification prompts, which allocate higher verification probabilities to more correct examples.

Method	CNLI			CNLI-CV		
	Cost	Perf.	IBC_Lift	Cost	Perf.	IBC_Lift
SLM	1	40.1	-	1	40.1	-
FrugalGPT	37.4	59.2	66.1	37.4	59.2	66.1
Self-Consistency	43.6	51.2	-17.3	40.5	50.6	-15.5
AutoMix-Threshold	51.9	55.4	-4.5	28.1	46.9	-49.1
AutoMix-POMDP	5.5	42.3	57.0	15.8	45.2	12.4
LLM	50	55.5	-	50	55.5	-

Table 5: Despite the boost in verifier accuracy with task-specific prompts (Figure 7), AutoMix may not always benefit, highlighting the utility of even weak verifiers when supported by meta-verifiers.

372 **C Few-Shot Prompts**

```
Story:
{relevant parts of the story}

{instruction}

Question: {question}

Answer:
```

Listing 2: **Task Prompt.** We experiment with long-context reasoning tasks, which require answering questions from stories, legal contracts, research papers, and novels.

```
Context: {context}

Question: {question}

AI Generated Answer: {generated_answer}

Instruction: Your task is to evaluate if the AI Generated Answer is correct, based
→ on the provided context and question. Provide the judgement and reasoning for
→ each case. Choose between Correct or Incorrect.

Evaluation: ""
```

Listing 3: **Verification Prompt.** The verification process is framed as a natural language entailment task, where the model determines the validity of the model-generated answer with respect to the context and question.

373 **C.1 Verifier Prompts**

```

### NARRATIVE_QA

Story:
{context}

You are given a story, which can be either a novel or a movie script, and a
→ question. Answer the question as concisely as you can, using a single phrase
→ if possible.

Question: {question}

Answer: The answer is"",
      "truncation_message": "... [The rest of the story is omitted]\n\n",

### QASPER

Article:
{context}

You are given a scientific article and a question. Answer the question as
→ concisely as you can, using a single phrase or sentence if possible. If the
→ question cannot be answered based on the information in the article, write
→ 'unanswerable'. If the question is a yes/no question, answer 'yes', 'no', or
→ 'unanswerable'.

Question: {question}

Answer: The answer is"",
      "truncation_message": "... [The rest of the article is omitted]\n\n",

### QUALITY

Story:
{context}

You are provided a story and a multiple-choice question with 4 possible answers
→ (marked by A, B, C, D). Choose the best answer by writing its corresponding
→ letter (either A, B, C, or D).

Question and Possible Answers: {question}

Answer: The answer is"",
      "truncation_message": "... [The rest of the story is omitted]\n\n",

### CNLI

Contract:
{context}

You are given a non-disclosure agreement and a sentence that proposes a hypothesis
→ based on the agreement. Choose whether the hypothesis is entailed by the
→ agreement, contradicted by the agreement, or not mentioned by (neutral to) the
→ agreement. If the hypothesis is entailed by the agreement, write 'Entailment'.
→ If the hypothesis is contradicted by the agreement, write 'Contradiction'. If
→ the hypothesis is not mentioned by the agreement, write 'Not mentioned'.

Hypothesis: {question}

Answer: The answer is"",
      "truncation_message": "... [The rest of the contract is omitted]\n\n",
}

```

Listing 4: **Few-Shot Prompts:** Dataset Specific Few-Shot Prompts used for SLM and LLM on NARRATIVE-QA, QASPER, QUALITY, CNLI. A general structure of context, instruction, question and answer is followed.

```

Context: The manuscript, discovered in 1980 in a dusty attic, turned out to be a
↳ lost work of Shakespeare.

Question: Whose lost work was discovered in a dusty attic in 1980?

AI Generated Answer: Shakespeare

Instruction: Your task is to evaluate if the AI Generated Answer is correct, based
↳ on the provided context and question. Provide the judgement and reasoning for
↳ each case. Choose between Correct or Incorrect.

Evaluation: The context specifically mentions that a lost work of Shakespeare was
↳ discovered in 1980 in a dusty attic.

Verification Decision: The AI generated answer is Correct.

---

Context: The celestial event, known as the Pink Moon, is unique to the month of
↳ April and has cultural significance in many indigenous tribes.

Question: In which month does the celestial event, the Pink Moon, occur?

AI Generated Answer: July

Instruction: Your task is to evaluate if the AI Generated Answer is correct, based
↳ on the provided context and question. Provide the judgement and reasoning for
↳ each case. Choose between Correct or Incorrect.

Evaluation: The context clearly states that the Pink Moon is unique to the month
↳ of April.

Verification Decision: The AI generated answer is Incorrect.

---

{truncated examples}

Context: {context}

Question: {question}

AI Generated Answer: {generated_answer}

Instruction: Your task is to evaluate if the AI Generated Answer is correct, based
↳ on the provided context and question. Provide the judgement and reasoning for
↳ each case. Choose between Correct or Incorrect.

Evaluation:

```

Listing 5: **Few-Shot Verifier Prompts:** 3-shot verifier prompt for evaluating the correctness of SLM’s answer. The same prompt is used for all datasets.