

---

# Distilling System 2 into System 1

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models (LLMs) can spend extra compute during inference to  
2 generate intermediate thoughts, which helps to produce better final responses. Since  
3 Chain-of-Thought [Wei et al., 2022], many such *System 2* techniques have been  
4 proposed such as Rephrase and Respond [Deng et al., 2023a], System 2 Attention  
5 [Weston and Sukhbaatar, 2023] and Branch-Solve-Merge [Saha et al., 2023]. In  
6 this work we investigate self-supervised methods to “compile” (distill) higher  
7 quality outputs from System 2 techniques back into LLM generations *without*  
8 intermediate reasoning token sequences, as this reasoning has been distilled into  
9 *System 1*. We show that several such techniques can be successfully distilled,  
10 resulting in improved results compared to the original System 1 performance, and  
11 with less inference cost than System 2. We posit that System 2 distillation will be  
12 an important feature of future continually learning AI systems, enabling them to  
13 focus System 2 capabilities on the reasoning tasks that they cannot yet do well.

## 14 1 Introduction

15 Generating intermediate thoughts allows a model (or human!) to reason and plan in order to  
16 successfully complete a task or respond to an instruction. We refer to such deliberate thinking as  
17 System 2 reasoning, following its description for humans in Sloman [1996], Kahneman [2011]  
18 and later for AI models [Bengio, 2017, LeCun, 2022, Weston and Sukhbaatar, 2023]. In System  
19 2 reasoning effortful mental activity is exerted, especially in situations where System 1 – more  
20 automatic thinking – is likely to make errors. In standard Large Language Models (LLMs) we  
21 thus define *System 1* as application of the Transformer [Vaswani et al., 2017] to directly produce  
22 a response given an input, without generation of intermediate tokens. We define *System 2* as any  
23 approach which generates intermediate tokens, including methods that perform search, or prompt  
24 multiple times, before finally generating a response. A battery of such *System 2* techniques have been  
25 proposed, among them Chain-of-Thought (CoT) [Wei et al., 2022], Tree-of-Thoughts [Yao et al.,  
26 2024], Graph-of-Thoughts [Besta et al., 2024], Branch-Solve-Merge [Saha et al., 2023], System 2  
27 Attention [Weston and Sukhbaatar, 2023], Rephrase and Respond [Deng et al., 2023a] and more.  
28 Many of these methods are shown to produce more accurate results due to this explicit reasoning, but  
29 typically do so at much higher inference cost and latency for a response. Due to the latter, many of  
30 these approaches are not used in production systems, which mostly use *System 1* generations.

31 For a human, the process of learning to transfer a skill from deliberate (System 2) to automatic  
32 (System 1) in psychology is referred to as *automaticity*, and the use of *procedural memory* Cohen  
33 and Squire [1980]. For example, when driving to work for the first time one might typically expend  
34 conscious effort planning and making decisions to get there. After a driver repeats this route, the  
35 driving process becomes “compiled” into the subconscious Charlton and Starkey [2013]. Similarly,  
36 playing a sport such as tennis can become “second nature”. In this work, we explore an analogous  
37 technique for AI models. Our approach performs this compilation, which we refer to as *System 2*  
38 *distillation*, in an unsupervised manner given a set of unlabeled examples. For each example we apply

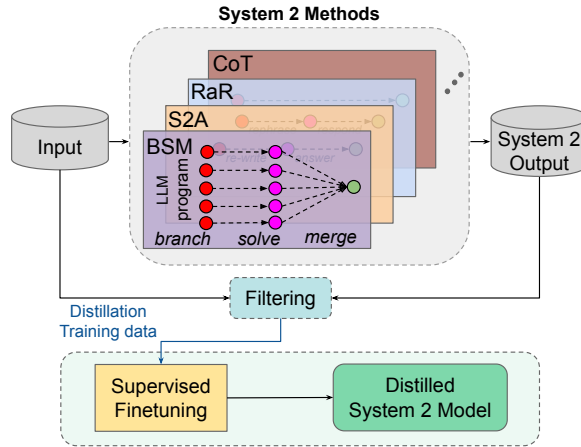


Figure 1: **Overview of System 2 Distillation.** Filtered training examples are collected by running System 2 approaches such as Branch-Solve-Merge (BSM) on unlabeled data, which uses extra compute to produce higher quality outputs. These targets are then distilled into the standard (System 1) LLM.

39 the given System 2 method, and then measure the quality of the prediction in an unsupervised manner.  
 40 For example, for tasks with unique answers we apply self-consistency [Wang et al., 2022], sampling  
 41 multiple times. For examples where System 2 is consistent enough, we assume this result should  
 42 be distilled, and add it to the distillation pool. We then fine-tune System 1 to match the predictions  
 43 of the System 2 method on the collected pool of examples, but *without* generating the intermediate  
 44 steps. Figure 1 illustrates the overall process of distilling System 2 into System 1.

45 We conduct experiments across 4 different System 2 LLM approaches and 5 different tasks. We find  
 46 our approach can distill System 2 reasoning into System 1 in a diverse array of settings, sometimes  
 47 even improving the results over the System 2 teacher. Moreover, these predictions are now produced  
 48 at a fraction of the computational cost. For example, we see successful distillation for tasks involving  
 49 dealing with biased opinions or irrelevant information (System 2 Attention), clarifying and improving  
 50 responses in some reasoning tasks (Rephrase and Respond), and for fine-grained evaluation of LLMs  
 51 (Branch-Solve-Merge). However, we also show that not all tasks can be distilled into System 1,  
 52 particularly complex math reasoning tasks requiring CoT. This is also mirrored in humans, who  
 53 cannot execute some tasks without deliberate System 2 reasoning [Kahneman, 2011].

## 54 2 Related work

### 55 2.1 System 1 and System 2 in Humans

56 In humans, System 1 reasoning is described as being capable of recognizing patterns, making quick  
 57 judgments, and understanding simple or familiar symbols. For instance, it is used to identify common  
 58 traffic signs, recognize faces, or associate basic symbols with specific emotions or ideas. However, for  
 59 complex problem-solving or for example manipulation of abstract symbols (like algebraic equations  
 60 or logical statements), System 2 reasoning is deemed necessary [Kahneman, 2011]. In psychology the  
 61 concept of *automaticity* describes behavior that becomes so well-practiced that it can be performed  
 62 with little to no conscious thought, with an example being driving a familiar route Charlton and  
 63 Starkey [2013]. In general, humans are said to use *procedural memory* to consolidate a specific task  
 64 into memory, learning through practice, so that it can be later performed without conscious awareness  
 65 [Cohen and Squire, 1980]. The concept of *unconscious competence* is classified as a later stage of  
 66 learning. Initially a person recognizes their incompetence, and consciously seeks to learn a skill until  
 67 they acquire *conscious competence*. Finally, the aim is to utilize it without conscious thought when it  
 68 is said to become, in common language, “second nature” [DePhillips et al., 1960].

### 69 2.2 System 1 and System 2 Models

70 We refer to a neural network that outputs a response directly without intermediate outputs as a  
 71 *System 1 model*. Such a network can nevertheless compute intermediate latent representations in its

72 layers before it outputs a response. As these states are represented as vectors they typically encode  
73 distributed knowledge, rather than discrete decisions, and have difficulty manipulating complex  
74 symbolic reasoning tasks directly [Nye et al., 2021, Cobbe et al., 2021, Yu et al., 2023, Li et al., 2024],  
75 which is analogous to issues with System 1 reasoning in humans. Nevertheless, many tasks can be  
76 solved with success directly in this manner without intermediate generations [Radford et al., 2019].

77 Nye et al. [2021] showed that the same language model that is unable to perform complex multi-step  
78 computations can perform those tasks when asked to generate intermediate steps into a “scratchpad”  
79 using either few-shot prompting or supervised training. Chain-of-thought reasoning was shown to be  
80 elicited from LLMs even using zero-shot prompting [Kojima et al., 2022] as well as by supervised  
81 [Cobbe et al., 2021] or few-shot [Wei et al., 2022] methods. LLM pretraining allows such reasoning  
82 to be built into the model because reasoning steps in discrete symbols (text) are present in the training  
83 corpora written by humans. Such *System 2 model* approaches output discrete tokens which is good  
84 for making sequential correct logical reasoning steps – but obviously has a downside if the reasoning  
85 is generated incorrectly. An incorrect discrete decision is difficult to recover from, unlike latent  
86 vector-based reasoning that might more easily model a distribution.

87 Recently, many approaches have been proposed to execute deeper reasoning using the LLM as part  
88 of an inner loop where it generates intermediate outputs, sometimes referred to as *LLM Programs*  
89 [Schlag et al., 2023]. These include subquestion decomposition [Perez et al., 2020], self-refinement  
90 [Madaan et al., 2024, Weston and Sukhbaatar, 2023, Deng et al., 2023a], self-verification and asking  
91 [Press et al., 2022, Weng et al., 2022, Dhuliawala et al., 2023], and various search techniques such as  
92 Tree-of-Thoughts and others [Yao et al., 2024, Besta et al., 2024].

### 93 2.3 (Standard) Distillation

94 The concept of distillation is usually applied to taking separate models, a powerful teacher model (or  
95 multiple teacher models) and a less powerful student model with separate parameters. The student  
96 model is then trained to mimic the behavior of the teacher(s). Methods of distillation include training  
97 the student to have similar output distributions [Hinton et al., 2015], layer activations [Adriana et al.,  
98 2015] or derivatives of the target teacher outputs [Czarnecki et al., 2017]. Earlier works considered  
99 distillation from an ensemble of multiple teacher models [Buciluă et al., 2006, Hinton et al., 2015].  
100 As neural networks have become larger, distilling from a larger to a smaller network has become a  
101 common paradigm [Ba and Caruana, 2014]. In contrast, in our work the teacher and student model  
102 are the same language model, but applied differently (either with intermediate reasoning, or not).

103 For CoT reasoning in particular, several distillation approaches have been considered Wang et al.  
104 [2023], Li et al. [2023a], Chen et al. [2024]. These again follow the paradigm of distilling a *separate*  
105 larger model’s output into a smaller model, i.e. the student model is asked to mimic the internal  
106 thoughts of the teacher model. The work of Zhang et al. [2024], however, considers distilling a slower  
107 System 2 method (Tree-of-Thought) into a faster System 2 method (CoT), which can use the same  
108 model as student and teacher. In contrast our work’s goal is to *not* generate internal thoughts (to  
109 improve System 1). Some exceptions are Deng et al. [2023b, 2024]. The former still uses a separate  
110 student and teacher model, but attempts to distill the intermediate thought tokens into the layers of  
111 the network by representing reasoning steps as vectors and then setting them as targets. The latter  
112 recent work attempts to distill CoT by gradually removing the intermediate steps, which can improve  
113 performance greatly compared to not doing so, but still does not match explicit CoT.

## 114 3 Distilling System 2 into System 1

### 115 3.1 Setup: System 1 and System 2 models

116 Given an input  $x$ , in this work we consider the setting of a single model, in our case a large language  
117 model (LLM), that is capable of two modes of response:

- 118 (i) *System 1*: Produces the output  $y$  directly. This is done by forwarding through the layers of  
119 the underlying autoregressive neural network (Transformer) to produce the output tokens.
- 120 (ii) *System 2*: We define System 2 models as methods that use the underlying Transformer  
121 to generate intermediate output tokens  $z$  of any kind *before* generating the final response  
122 tokens. This may include multiple calls (prompts).

123 More formally, we consider a System 2 model  $S_{\text{II}}$  as a function that takes an LLM  $p_{\theta}$  and input  $x$ , and  
 124 can call the LLM possibly repeatedly to generate intermediate tokens  $z$  using a specific algorithm,  
 125 before returning an output  $y$ :

$$S_{\text{II}}(x; p_{\theta}) \rightarrow z, y. \quad (1)$$

126 System 2 approaches can potentially involve multiple prompts, branching, iteration and search, all  
 127 the while using the LLM to generate intermediate results for further processing. In contrast, a System  
 128 1 model only considers the original input  $x$  and calls the LLM  $p_{\theta}$  directly to produce an output  $y$ :

$$S_{\text{I}}(x) = p_{\theta}(x) \rightarrow y. \quad (2)$$

129 There are many existing instantiations of System 2 models. CoT prompting only requires a single  
 130 LLM prompt, but outputs intermediate generations before a final response, typically used in math  
 131 and other reasoning tasks [Wei et al., 2022].

132 Methods like System 2 Attention [Weston and Sukhbaatar, 2023] and Rephrase and Respond [Deng  
 133 et al., 2023a] require two calls to the LLM, where in the former the first call is used to attend to the  
 134 context and remove bias, and in the latter to expand on the question. The second call is then used  
 135 to finally respond to the answer given the intermediate generations. Some methods are much more  
 136 sophisticated for example Branch-Solve-Merge [Saha et al., 2023] which generates a plan via an  
 137 LLM which branches into several more LLM calls until a final stage merges the results.

138 We perform experiments with the four methods just described, but there are many other system 2  
 139 approaches, for example Tree-of-Thoughts [Yao et al., 2024], Graph-of-Thoughts [Besta et al., 2024]  
 140 and more, see related work in section 2.

### 141 3.2 Method: System 2 Distillation

142 Many System 2 methods, by their nature, are significantly slower at inference time due to multiple  
 143 prompt calls and generation of intermediate tokens. The aim of System 2 Distillation is to distill  
 144 all the reasoning from  $S_{\text{II}}$  back into  $S_{\text{I}}$  so that the direct outputs from the language model  $p_{\theta}(x)$   
 145 are improved. We assume a setting where the model has access to *unlabeled inputs*  $\mathcal{X}$  from which  
 146 it can learn, in analogy to how humans learn their *procedural memory* without supervision. For  
 147 language-based tasks, it is common to have access to instruction following prompts (inputs) as they  
 148 can be collected from humans, e.g. the 1M released WildChat interactions [Zhao et al., 2024] where  
 149 inputs are given but correct labels are unknown. Hence this is a realistic setup.

150 The first step is to generate responses using the System 2 model over the unlabeled inputs  $\mathcal{X}$ :

$$y_{S_{\text{II}}}^i = S_{\text{II}}(x^i; p_{\theta}), \quad \forall x_i \in \mathcal{X}. \quad (3)$$

151 Note we discard (do not store) the intermediate outputs  $z$  from Eq. 1. These responses  $y_{S_{\text{II}}}^i$  can then  
 152 be used directly as System 2 distillation targets for fine-tuning a System 1 model. However, they are  
 153 subject to noise: some of these responses could be high quality, while others could be low quality or  
 154 incorrect. For shortform QA and reasoning tasks involving a short response with a typically unique  
 155 correct (but unknown) answer, we thus consider an *unsupervised curation* step to attempt to improve  
 156 training data quality. We consider two variations which both rely on a consistency criterion:

- 157 • *self-consistency of outputs*: we sample  $S_{\text{II}}(x^i; p_{\theta})$  a total of  $N$  times, and accept the response  
 158 that is the majority vote; if there is no majority winner, we discard the example.
- 159 • *self-consistency under input perturbation*: we perturb the input  $x^i$  in such a way that the  
 160 output should not change, e.g. changing the order of multiple-choice items in the prompt,  
 161 and compute  $S_{\text{II}}$  for each perturbation; if the outputs do not agree, we discard the example.

162 After that, we end up with the synthetic dataset  $(\mathcal{X}_{S_{\text{II}}}, \mathcal{Y}_{S_{\text{II}}})$ , where  $\mathcal{X}_{S_{\text{II}}}$  is a filtered subset of  $\mathcal{X}$ . The  
 163 final step is supervised fine-tuning of the LLM with parameters  $p_{\theta}$  using this distilled training set. We  
 164 typically initialize this model from the current state  $p_{\theta}$  and continue training with the new dataset.

165 After fine-tuning we obtain an LLM  $\hat{p}_{\theta}$  which is a System 1 model that is expected to provide outputs  
 166 and performance gains similar to the evaluated System 2 model.

## 167 4 Experiments

### 168 4.1 Training and Evaluation Setup

169 We use Llama-2-70B-chat [Touvron et al., 2023] as the base model for all our experiments. We  
170 require a base model of sufficient power that it can be performant as a System 2 model, but also  
171 have open weights that can be fine-tuned, hence this choice. We consider several System 2 methods,  
172 including Rephrase and Respond (RaR), System 2 Attention (S2A), Branch-Solve-Merge (BSM), and  
173 Chain-of-Thought (CoT), focusing on tasks where each method has demonstrated strong performance.  
174 For System 1, we conduct zero-shot inference using the instruction-tuned base model as a standard  
175 baseline. We report task-specific metrics for each task, and the “#Tokens” metric which measures the  
176 average number of tokens generated per input across the evaluation set. For System 2 methods this  
177 includes both intermediate token generations as well as the final output token generations. Detailed  
178 descriptions of the experimental setups are available in the Appendix A.2.

### 179 4.2 Rephrase and Respond Distillation

180 Rephrase and Response (RaR) [Deng et al., 2023a] is a System 2 method that first prompts the  
181 LLM to rephrase the original question with further elaboration, and then to generate a response to  
182 the rephrased question with the aim that this provides superior output. The authors introduce two  
183 approaches, 1-step RaR and 2-step RaR, where the latter involves two separate prompts rather than  
184 a combined one as in the former, see Appendix A.1 for specific prompts. They find that 2-step  
185 RaR significantly improves performance on several reasoning tasks that are challenging for the  
186 baseline LLM. We consider two tasks from the original paper where it performed well: the last letter  
187 concatenation task and coin flip reasoning. We then assess the distillation of this System 2 approach.

188 **Distillation Data** We build the System 2 distillation dataset for RaR using *self-consistency of*  
189 *outputs*. For each input, we conduct eight sampling iterations for the last letter task and eight for each  
190 stage of the coin flip task.<sup>1</sup> We then apply a majority vote to determine the final output.

	Acc $\uparrow$	#Tokens
<i>System 1</i>		
Llama-2-70B-chat	56.1%	61.9
Distill System 1	54.5%	30.4
<i>System 2</i>		
1-Step RaR	58.5%	158.9
2-Step RaR	77.2%	112.4
<i>Distill System 2</i>		
Distill 2-Step RaR	75.69%	50.3

Table 1: **System 2 Distillation of Rephrase and Respond:** Coin Flip. We report exact match (EM) test accuracy and number of generated (intermediate and output) tokens.

Model	Acc $\uparrow$ (biased)	Acc $\uparrow$ (unbiased)	#Tokens
<i>System 1 (Zero-shot)</i>	51.6%	73.8%	165
<i>System 2 (S2A)</i>	76.0%	69.3%	147
Distill S2A	81.3%	78.6%	56
Distill S2A (no USC)	78.6%	75.3%	58

Table 2: **Distillation of System 2 Attention:** TriviaQA task, reporting accuracies on the biased and unbiased eval sets.

191 Coin flip reasoning task has frequently been tested in research, including in Wei et al. [2022] and  
192 Deng et al. [2023a]. It involves determining the final face (heads or tails) of a coin, starting from a  
193 known initial position after a series of flips described in natural language, such as “A coin is heads  
194 up. Roxas does not flip the coin. Schneiderman does not flip the coin. Is the coin still heads up?”  
195 Deng et al. [2023a] showed that even strong language models do not succeed at this task, whereas  
196 applying the RaR method improves their performance. There are 20k training examples, which we  
197 use for unsupervised learning (without labels), 3.33k validation and 1.33k test examples.

198 **Results** Overall results are given in Table 1. Llama-2-70B-chat (zero-shot) has a success rate of  
199 56.1% on this task, while 1-Step and 2-Step RaR have success rates of 58.5% and 77.2% respectively.  
200 We thus only see a large improvement with the 2-Step method. Distilling 2-Step RaR back into a  
201 system 1 Llama-2-70B-chat via our unsupervised technique yields 75.69%. Hence, we find that our

<sup>1</sup>This approach was adopted after observing that sampling just once for the rephrase stage yielded suboptimal results.

202 distilled System 2 model delivers performance comparable to that of System 2 (2 Step RaR), but  
 203 without the need to execute the LLM program with 2 prompts (see # of generated Tokens).

### 204 4.3 System 2 Attention Distillation

205 [Weston and Sukhbaatar \[2023\]](#) proposed System 2 Attention (S2A) that helps to reduce models’  
 206 reasoning pitfalls such as relying on biased information in the input or attending to irrelevant context.  
 207 S2A is a two-stage inference method where the first stage rewrites the input to remove undesired  
 208 information such as bias or irrelevant context, and the second stage attends to the shorter rewritten  
 209 context (in contrast to RaR which expands the context), see [Figure 6](#). In this work we verify the  
 210 feasibility of distilling S2A into System 1. In particular, we focus on the SycophancyEval question  
 211 answering task [[Sharma et al., 2023](#)] that contains biased information in the input that is known to  
 212 hurt LLM performance. We use 6668 examples from SycophancyEval as unlabeled training data, and  
 213 400 examples for evaluation, where the latter are split into biased inputs (350) and without bias (50).

214 **Distillation data** We use universal self-consistency (USC) [[Chen et al., 2023](#)] to select high quality  
 215 targets. Specifically, we sample 20 generations and then use the Llama-70B-chat model with a USC  
 216 prompt (provided in [Figure 12](#)) to compose a self-consistent (majority) final answer that is used as  
 217 the distillation target.

218 **Results** The results are provided in [Table 2](#), reporting average accuracy over 3 random seeds. The  
 219 baseline (System 1) LLM has low accuracy on the biased portion as expected, being susceptible to  
 220 biased inputs. S2A improves performance dramatically for biased inputs. System 2 distillation shows  
 221 similarly strong performance as the System 2 approach. There is, however, a significant reduction  
 222 in the average number of tokens used compared to both the baseline and the S2A model. This is  
 223 because biased inputs tend to make the baseline LLM generate more output tokens, while S2A has to  
 224 generate intermediate tokens as well. [Figure 11](#) shows a representative example. Finally, we show  
 225 that using USC for distillation is important for overall results, by also reporting results without USC  
 226 (last row), where the latter provides inferior results. This highlights the importance of the distillation  
 227 data quality that is used during fine-tuning.

	OASST2 Eval			MT-bench Eval		
	Agreement ↑	% Inconsistent ↓	#Tokens	Agreement ↑	% Inconsistent ↓	#Tokens
<i>System 1</i>						
GPT-4-0125-preview	44.7%	35.5%	4	68.1%	25.6%	4
Llama-2-70B-chat	32.0%	56.7%	4	28.1%	80.9%	4
<i>System 2</i>						
CoT (GPT-4-0125-preview)	48.7%	28.2%	603.7	73.8%	16.2%	548.8
CoT (Llama-2-70B-chat)	45.2%	37.7%	432.6	58.9%	30.8%	411.8
BSM (Llama-2-70B-chat)	49.1%	30.4%	2117.8	64.5%	21.1%	2063.1
<i>Distill System 2</i>						
Distill BSM (Llama-2-70B-chat)	58.4%	12.2%	4	72.4%	9.1%	4

Table 3: **System 2 Distillation of Branch-Solve-Merge (BSM)**: Open Assistant (OASST2) and MT-bench evaluation of LLM-as-a-Judge for various models. System 2 Distillation of BSM outperforms BSM itself, and even GPT4-as-a-Judge, despite using Llama-2-70B-chat. Distilled BSM has higher human agreement (Agreement), less position inconsistent predictions (% Inconsistent), and uses less output tokens (#Tokens).

### 228 4.4 Branch-Solve-Merge Distillation

229 Branch-Solve-Merge (BSM) [[Saha et al., 2023](#)] consists of three modules: *branch*, *solve*, and *merge*.  
 230 These modules work together to break down a task into several parallel sub-tasks, each guided by  
 231 specific prompts. BSM has proven effective when used in the context of an LLM acting as a judge,  
 232 see [Figure 14](#). The method begins by prompting the LLM to list evaluation metrics (*branch*) tailored  
 233 to a given user query. Subsequently, the LLM is queried to evaluate a response based on each metric  
 234 independently (*solve*). Finally, the scores from each branch are averaged to arrive at a comprehensive  
 235 evaluation decision (*merge*). Notably, this method incurs an inference cost 5-6 times greater than that  
 236 of a conventional (System 1) LLM evaluation approach, making it much less practical. We assess the  
 237 feasibility of distilling BSM, aiming to retain its benefits while reducing computational cost.

Model	k=1		k=5		k=10	
	Acc %	#Tokens	Acc %	#Tokens	Acc %	#Tokens
<i>System 1</i>						
Few (8)-shot (no CoT)	7.58%	57	9.40%	295	10.31%	620
<i>System 2</i>						
CoT zero-shot	52.77%	270	57.54%	1385	59.44%	2760
CoT few (8)-shot	36.39%	297	54.97%	1560	63.84%	3120
<i>Distill System 2</i>						
Distill CoT zero-shot	7.13%	18	7.13%	90	7.35%	180

Table 4: GSM8k test set accuracy. Number of votes  $k$  in majority voting represents how many candidates were sampled to collect votes towards predicted answers. In this case System 2 Distillation of CoT does not work well.

238 **Distillation Data** Following Yuan et al. [2024], Li et al. [2023b], we used the Open Assistant  
239 Dataset v2 (OASST2) [Köpf et al., 2024] with turn 1 and English only data. We use queries along with  
240 two candidate responses from the OASST2 training set as inputs (19,672 examples in total). We use  
241 *self-consistency under input perturbations* to ensure the quality of our distillation data. Specifically,  
242 as two responses are being judged, we evaluate each sample twice with BSM - once in the original  
243 order and once in the swapped order. The winning response should remain consistent regardless of  
244 the order. We exclude samples without a consistent winner when the response order is swapped.

245 **Evaluation** We evaluate our models on the OASST2 valid set and MT-bench [Zheng et al., 2024].  
246 The OASST2 validation set comprises 273 samples, restricted to turn 1 and English language only.  
247 Evaluations of response pairs are performed in both original and swapped orders. As we trained our  
248 distilled model on the OASST2 training set, the OASST2 validation set functions as an in-distribution  
249 evaluation set, while MT-bench is more out-of-distribution. MT-bench is a popular benchmark that  
250 evaluates LLM-as-judges of other LLM’s responses when acting as helpful AI assistants conversations.  
251 It consists of instructions from 8 diverse domains e.g., writing, reasoning, math, coding, etc.

252 Following Zheng et al. [2024], we assessed the Agreement between model votes and human expert  
253 votes. A well-documented limitation of LLM-as-a-judge is position bias, where a LLM tends to  
254 favor certain positions over others. This bias is evident as altering the position of responses in the  
255 evaluation prompt often leads to different decisions by the model. To quantify this, we measure not  
256 only agreement but also the Percentage of Inconsistent examples to assess position bias.

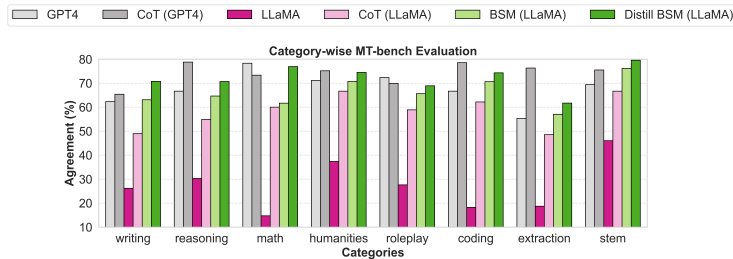


Figure 2: The agreement between LLM judges and human preferences per category on MT-bench.

257 **OASST2 Evaluation Results** Table 3 provides results on the OASST2 dataset. Compared to  
258 baseline (System 1) LLMs, the Chain-of-Thought (CoT) method improves performance by improving  
259 agreement and reducing inconsistency rates (see prompts in Appendix). While BSM outperforms  
260 CoT, this comes at the cost of increased inference time (#Tokens). Remarkably, our distilled System  
261 2 BSM model requires the generation of only four tokens and still outperforms both CoT and BSM.  
262 Furthermore, our distilled model based on Llama-2-70B-chat outperforms GPT-4-0125-preview,  
263 achieving higher human agreement and greater consistency.

264 **MT-Bench Evaluation Results** Table 3 also provides results on MT-bench, which serves as an out-  
265 of-distribution test. The results mirror those from the OASST2 evaluation. Both Chain-of-Thought  
266 (CoT) and BSM improve model performance but at the expense of significantly increased inference  
267 costs. Our distilled BSM model not only achieves higher human agreement and lower inconsistency  
268 rates but also requires less computational resources. Although our model slightly underperforms  
269 in agreement compared to the state-of-the-art GPT-4-0125-preview model, it was trained solely on

270 unlabeled data from OASST2 based on Llama-2-70B-chat. Despite this, it is more consistent and  
271 inference is cheap in terms of output tokens.

272 Here, we further analyze the MT-Bench results in terms of Agreement by category. Figure 2 shows  
273 the per category agreement. We observe that CoT improved agreement compared to the base model  
274 (Llama-2-70B-Chat) on all categories. BSM is better than CoT and our distilled BSM is even better  
275 than BSM. Although Distilled BSM achieves superior performance compared to the baselines across  
276 all categories, it still lags behind GPT-4-0125-preview in reasoning, coding, and extraction. However,  
277 it surpasses GPT-4-0125-preview in writing, math, and STEM.

## 278 4.5 Chain-of-Thought Distillation

279 Chain-of-Thought (CoT) [Wei et al., 2022] has been shown to be an effective method to improve  
280 LLM’s reasoning abilities, such as for solving graduate school math problems. The LLM generates  
281 intermediate tokens that are steps (*chain*) of reasoning (*thoughts*) before it produces the final answer.  
282 We consider two variants of the approach: (i) few-shot CoT, whereby multiple [question, CoT,  
283 answer] examples from the training set are provided as part of the context followed by the question;  
284 and (ii) zero-shot, whereby an explicit instruction to think “step by step” is added to the prompt in  
285 addition to the question, see Appendix Figure 10.

286 **Distillation data** We use CoT to produce answers for questions from the training split of GSM8k  
287 [Cobbe et al., 2021] (which we consider unlabeled), using majority voting with  $K = 10$ . The resulting  
288 distillation training set consists of 7461 [question, answer] pairs, i.e. without any intermediate  
289 reasoning steps. The accuracy of the self-supervised targets, for analysis purposes, is 56.81%.

290 **Evaluation** We report evaluation accuracy computed over the GSM8k test set with majority voting  
291 with different values of  $K$ . Similarly to our previous experiments, we report the average number of  
292 predicted tokens for each method. Note that we compute this average over all generated tokens when  
293 we run majority voting to see how the increase in  $K$  affects the inference cost. We consider several  
294 baselines: System 1 and System 2 (CoT) methods evaluated with zero-shot or 8-shot input contexts.  
295 Note that System 2 with 8-shot means that CoTs are provided in the few-shot inputs, while System 1  
296 means that the few shot examples contain questions and answers, but no CoTs.

297 **Results** Evaluation results are presented in Table 4. First, improvements are coming from using the  
298 CoT method as expected: it helps when being presented as part of the few-shot context or as part  
299 of the instruction in the prompt template. These improvements come with an increase in inference  
300 cost: sequences predicted with CoT methods are substantially longer compared to the System 1  
301 method. Second, our System 2 distillation method yields poor performance across various decoding  
302 hyper-parameters. The GSM8k task (math problems) requires a very different kind of reasoning  
303 compared to other tasks we considered in this work. This highlights the non-trivial aspect of System  
304 2 distillation: the proposed distillation algorithm works in many cases but not always. This leaves  
305 room for future research to elucidate in exactly which circumstances to apply distillation, and when  
306 not to, in a similar manner perhaps to the approach in humans.

## 307 5 Conclusion

308 Recent work has shown that complex reasoning procedures using LLMs in the inner loop, called  
309 System 2 approaches, can improve performance. In this work we have shown that in many cases  
310 it is possible to distill this System 2 reasoning into the outputs of the LLM *without intermediate*  
311 *generations* while maintaining, or sometimes even improving, performance. While not all methods  
312 can be distilled easily using our method, with Chain-of-Thought for complex reasoning being a  
313 challenging counterexample, this is possible for diverse approaches. Our method works for System  
314 2 Attention for dealing with bias and irrelevant context, Rephrase and Respond for clarifying task  
315 instructions, and Branch-Solve-Merge for improved LLM-as-a-Judge evaluation. Pragmatically,  
316 distilling these approaches makes them more likely to be used by LLM practitioners, and they are  
317 more efficient at inference time. Looking forward, systems that can distill useful tasks in this way  
318 free up more time to spend on reasoning about the tasks that they cannot yet do well, just as humans  
319 do. Hence, we expect exploring this approach in a continuous training loop will be a fruitful research  
320 direction.



321 **References**

- 322 Romero Adriana, Ballas Nicolas, K Samira Ebrahimi, Chassang Antoine, Gatta Carlo, and Bengio  
323 Yoshua. *Fitnets: Hints for thin deep nets*. *Proc. ICLR*, 2(3):1, 2015.
- 324 Jimmy Ba and Rich Caruana. *Do deep nets really need to be deep?* *Advances in neural information  
325 processing systems*, 27, 2014.
- 326 Yoshua Bengio. *The consciousness prior*. *arXiv preprint arXiv:1709.08568*, 2017.
- 327 Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi,  
328 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. *Graph of thoughts:  
329 Solving elaborate problems with large language models*. In *Proceedings of the AAAI Conference  
330 on Artificial Intelligence*, volume 38, pages 17682–17690, 2024.
- 331 Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. *Model compression*. In *Proceedings  
332 of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*,  
333 pages 535–541, 2006.
- 334 Samuel G Charlton and Nicola J Starkey. *Driving on familiar roads: Automaticity and inattention  
335 blindness*. *Transportation research part F: traffic psychology and behaviour*, 19:121–133, 2013.
- 336 Xin Chen, Hanxian Huang, Yanjun Gao, Yi Wang, Jishen Zhao, and Ke Ding. *Learning to maximize  
337 mutual information for chain-of-thought distillation*. *arXiv preprint arXiv:2403.03348*, 2024.
- 338 Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Kefan Xiao, Pengcheng Yin, Sushant Prakash,  
339 Charles Sutton, Xuezhi Wang, and Denny Zhou. *Universal self-consistency for large language  
340 model generation*, 2023.
- 341 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
342 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John  
343 Schulman. *Training verifiers to solve math word problems*, 2021.
- 344 Neal J Cohen and Larry R Squire. *Preserved learning and retention of pattern-analyzing skill in  
345 amnesia: Dissociation of knowing how and knowing that*. *Science*, 210(4466):207–210, 1980.
- 346 Wojciech M Czarnecki, Simon Osindero, Max Jaderberg, Grzegorz Swirszcz, and Razvan Pascanu.  
347 *Sobolev training for neural networks*. *Advances in neural information processing systems*, 30,  
348 2017.
- 349 Yihe Deng, Weitong Zhang, Zixiang Chen, and Quanquan Gu. *Rephrase and respond: Let large  
350 language models ask better questions for themselves*. *arXiv preprint arXiv:2311.04205*, 2023a.
- 351 Yuntian Deng, Kiran Prasad, Roland Fernandez, Paul Smolensky, Vishrav Chaudhary, and Stu-  
352 art Shieber. *Implicit chain of thought reasoning via knowledge distillation*. *arXiv preprint  
353 arXiv:2311.01460*, 2023b.
- 354 Yuntian Deng, Yejin Choi, and Stuart Shieber. *From explicit cot to implicit cot: Learning to internalize  
355 cot step by step*. *arXiv preprint arXiv:2405.14838*, 2024.
- 356 Frank A DePhillips, William M Berliner, and James J Cribben. *Management of training programs.*  
357 *homewood, illinois: Richard d. irwin*, 1960.
- 358 Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and  
359 Jason Weston. *Chain-of-verification reduces hallucination in large language models*. *arXiv preprint  
360 arXiv:2309.11495*, 2023.
- 361 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. *Distilling the knowledge in a neural network*. *arXiv  
362 preprint arXiv:1503.02531*, 2015.
- 363 Daniel Kahneman. *Thinking, fast and slow*. macmillan, 2011.
- 364 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. *Large  
365 language models are zero-shot reasoners*. *Advances in neural information processing systems*, 35:  
366 22199–22213, 2022.

- 367 Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith  
368 Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, et al. Openassistant  
369 conversations-democratizing large language model alignment. *Advances in Neural Information  
370 Processing Systems*, 36, 2024.
- 371 Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open  
372 Review*, 62(1), 2022.
- 373 Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. Sym-  
374 bolic chain-of-thought distillation: Small models can also "think" step-by-step. *arXiv preprint  
375 arXiv:2306.14050*, 2023a.
- 376 Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Luke Zettlemoyer, Omer Levy, Jason Weston, and  
377 Mike Lewis. Self-alignment with instruction backtranslation. *arXiv preprint arXiv:2308.06259*,  
378 2023b.
- 379 Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to  
380 solve inherently serial problems. *arXiv preprint arXiv:2402.12875*, 2024.
- 381 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri  
382 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement  
383 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- 384 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David  
385 Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work:  
386 Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*,  
387 2021.
- 388 Ethan Perez, Patrick Lewis, Wen-tau Yih, Kyunghyun Cho, and Douwe Kiela. Unsupervised question  
389 decomposition for question answering. *arXiv preprint arXiv:2002.09758*, 2020.
- 390 Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. Measuring  
391 and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*,  
392 2022.
- 393 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
394 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 395 Swarnadeep Saha, Omer Levy, Asli Celikyilmaz, Mohit Bansal, Jason Weston, and Xian Li.  
396 Branch-solve-merge improves large language model evaluation and generation. *arXiv preprint  
397 arXiv:2310.15123*, 2023.
- 398 Imanol Schlag, Sainbayar Sukhbaatar, Asli Celikyilmaz, Wen-tau Yih, Jason Weston, Jürgen Schmid-  
399 huber, and Xian Li. Large language model programs. *arXiv preprint arXiv:2305.05364*, 2023.
- 400 Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askill, Samuel R. Bowman,  
401 Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R. Johnston, Shauna Kravec, Timothy  
402 Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda  
403 Zhang, and Ethan Perez. Towards understanding sycophancy in language models, 2023.
- 404 Steven A. Sloman. The empirical case for two systems of reasoning. *Psychological Bulletin*, 119:  
405 3–22, 1996. URL <https://api.semanticscholar.org/CorpusID:13454019>.
- 406 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
407 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
408 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 409 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
410 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing  
411 systems*, 30, 2017.
- 412 Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. Scott: Self-  
413 consistent chain-of-thought distillation. *arXiv preprint arXiv:2305.01879*, 2023.

- 414 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-  
415 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.  
416 *arXiv preprint arXiv:2203.11171*, 2022.
- 417 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
418 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*  
419 *neural information processing systems*, 35:24824–24837, 2022.
- 420 Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and  
421 Jun Zhao. Large language models are better reasoners with self-verification. *arXiv preprint*  
422 *arXiv:2212.09561*, 2022.
- 423 Jason Weston and Sainbayar Sukhbaatar. System 2 attention (is something you might need too).  
424 *arXiv preprint arXiv:2311.11829*, 2023.
- 425 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.  
426 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural*  
427 *Information Processing Systems*, 36, 2024.
- 428 Dongran Yu, Bo Yang, Dayou Liu, Hui Wang, and Shirui Pan. A survey on neural-symbolic learning  
429 systems. *Neural Networks*, 2023.
- 430 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason  
431 Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024.
- 432 Xuan Zhang, Chao Du, Tianyu Pang, Qian Liu, Wei Gao, and Min Lin. Chain of preference  
433 optimization: Improving chain-of-thought reasoning in llms. *arXiv preprint arXiv:2406.09136*,  
434 2024.
- 435 Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. Wildchat: 1m  
436 chatgpt interaction logs in the wild. *arXiv preprint arXiv:2405.01470*, 2024.
- 437 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
438 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
439 chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.

	Acc $\uparrow$	#Tokens
<i>System 1</i>		
Llama-2-70B-chat	30.0%	27.1
Distill System 1	69.5%	24.4
<i>System 2</i>		
1-Step RaR	39.5%	106.6
2-Step RaR	44.5%	41.5
<i>Distill System 2</i>		
Distill 2-Step RaR	98.0%	25.5

Table 5: **System 2 Distillation of Rephrase and Respond**: Last Letter Concatenation tasks. We report exact match (EM) test accuracy and number of generated (intermediate and output) tokens.

## 440 A Appendix

### 441 A.1 Prompts

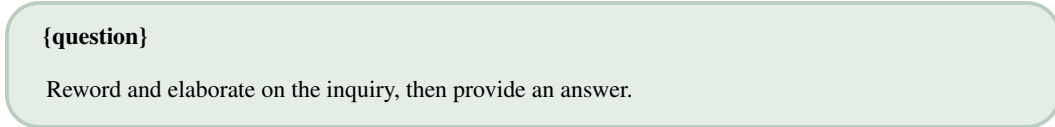


Figure 3: **1-step RaR prompt**. The 1-step RaR process involves the model rephrasing the question and subsequently providing an answer, all in a single step.

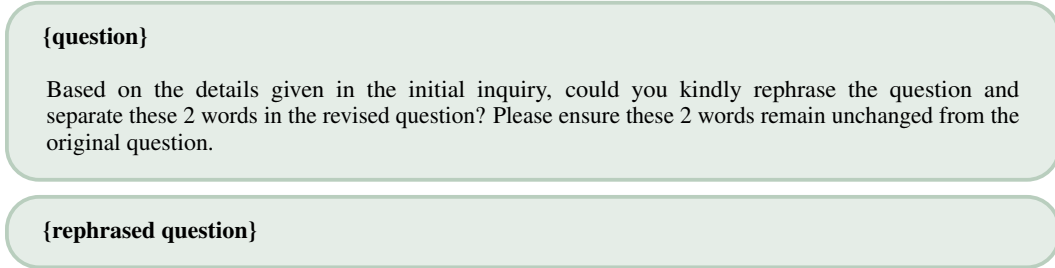


Figure 4: **2-step RaR prompt for last letter concatenation task, step 1 (top), step 2 (down)** The 1-step RaR process involves the model rephrasing the question and subsequently providing an answer, all in a single step.

### 442 A.2 Experiment Details

### 443 A.3 More Experiments for Rephrase and Respond Distillation

#### 444 A.3.1 Last letter Concatenation Task

445 This task focuses on symbolic reasoning, requiring the model to concatenate the last letters of given  
 446 words. For instance, the instruction: “Take the last letters of the words in ‘Edgar Bob’ and concatenate  
 447 them.” As demonstrated in [Deng et al. \[2023a\]](#), this task benefits significantly from the application of  
 448 the RaR method. We compiled a dataset by randomly selecting 1200 unique English words. Using  
 449 this, we constructed 200 samples each for training, validation, and test.

450 **Results** Overall results are given in [Table 5](#). The baseline System 1 model (Llama-2-70B-chat)  
 451 achieves an accuracy of 30.0%, and is outperformed by the System 2 methods of 1-Step and 2-Step  
 452 RaR (39.5% and 44.5%, respectively). Distilling the 2-Step RaR method back into a System 1  
 453 Llama-2-70B-chat model via our unsupervised technique, we achieve a remarkable accuracy of  
 454 98.0%. The model can effectively learn from this training data how to solve the task, in comparison

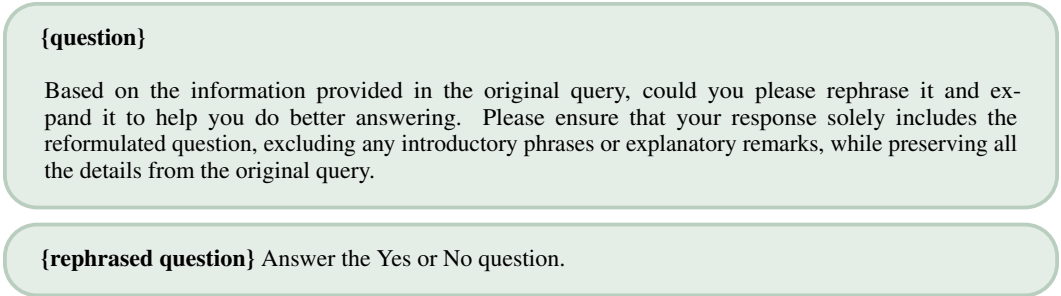


Figure 5: **2-step RaR prompt for coin flip task, step 1 (top), step 2 (down)** The 1-step RaR process involves the model rephrasing the question and subsequently providing an answer, all in a single step.

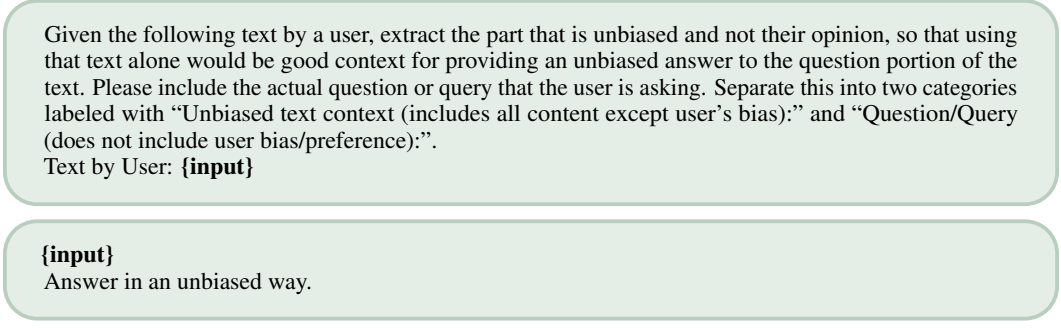


Figure 6: **System 2 Attention prompts.** We use the prompts from [Weston and Sukhbaatar \[2023\]](#) to extract the training signal for distillation. The output after the second stage is used as the distillation target.

455 to the zero-shot chat model. Distillation of Rephrase and Respond effectively inherits the advantages  
 456 of both System 2 and System 1. It maintains the accuracy benefits of System 2, while its inference  
 457 cost is comparable to that of System 1 (see # of generated Tokens).

458 **Last Letter Concatenation Analysis & Ablations** To evaluate the effectiveness and necessity of  
 459 our unsupervised curation step using *self-consistency of outputs* we conducted an ablation study  
 460 by creating a distillation dataset without applying the self-consistency filter. When we distilled the  
 461 System 2 model using this unfiltered dataset under the same setting, it achieved an exact match  
 462 accuracy of 87.5% (with 98% for the filtered version). This comparison underscores the critical role  
 463 of consistency filtering. Nevetthess, in both cases constructing training data does improve results over  
 464 zero-shot performance. We also attempted to distill the System 1 predictions using the same filtering  
 465 technique, which results in a lower accuracy of 69.5%.

466 **Coin Flip task Analysis & Ablations** The RaR method in [Deng et al. \[2023a\]](#) incorporates prompt  
 467 engineering tricks, such as appending phrases like “Flip means reverse. Answer the Yes or No  
 468 question” to the original query, which has been shown to enhance model performance. Following  
 469 their approach, we evaluated model performance using different prompts, see [Table 8](#). When testing  
 470 the Llama-2-70B-chat model (System 1) with prompts like “Flip means reverse” and “Flip means  
 471 reverse. Answer the Yes or No question,” we observed a significant improvement in performance,  
 472 from 56.11% to 66.84%. This highlights the critical role of prompt selection in optimizing the  
 473 performance of System 1 models. However, this reliance on prompt engineering also represents a  
 474 limitation, necessitating additional human effort.

475 We also attempted to distill the System 1 model, which gave poor performance. In this case, we also  
 476 observed fluctuations in performance with different prompts. In contrast, the distilled System 2 model  
 477 demonstrated consistent performance across various prompts, with a lower sensitivity to prompt  
 478 variations. This consistency indicates that extensive prompt engineering might not be essential for  
 479 the distilled System 2 model.

We want to evaluate the quality of the responses provided by two AI assistants to the user question displayed below. Your task is to propose an evaluation plan that can be executed to compare the two responses. The evaluation plan should consist of a list of up to five factors that one should consider such as helpfulness, relevance, accuracy, etc. In each line, write an evaluation criterion along with a short description of how we should evaluate that criterion.

User Question: **{user\_query}**

Evaluation Plan:

Figure 7: **BSM: Branch prompt.**

You are given a user question and responses provided by two AI assistants. Your task is to evaluate and score the quality of the responses based on a single evaluation criterion displayed below. Make sure to evaluate only based on the criterion specified and none other. In the first line, provide a score between 1 to 5 for Assistant A’s response. In the second line, provide a score between 1 to 5 for Assistant B’s response.

[User Question]

**{user\_query}**

[The Start of Assistant A’s Answer]

**{response\_a}**

[The End of Assistant A’s Answer]

[The Start of Assistant B’s Answer]

**{response\_b}**

[The End of Assistant B’s Answer]

[Evaluation Criterion]

**{eval\_criterion}**

[End of Evaluation Criterion] Evaluation of **{criterion\_name}**:

Figure 8: **BSM: Solve prompt.**

480 **Model training** We use Llama2 70B Chat as the initialization for SFT training with CE loss. The  
 481 loss is only applied on the answer part of the sequence. Model is trained with dropout 0.1, learning  
 482 rate  $5.5e - 6$ , with warmup 1. Table 6 shows details about total training steps and total training tokens  
 483 per step.

484 **S2A** For S2A, in both generation stages we use nucleus sampling with top-p value 0.9. During  
 485 distillation, for USC, in some cases the generated answers are too long and 20 do not fit in the  
 486 Llama2 context. In these rare cases we reduce the answer set to 10 or select an answer randomly if  
 487 10 generated answers are still too long.

Methods	Dataset	Total Training Steps	Total Training Tokens per Step
RaR	Last Letter Concatenation	3	66k
RaR	Coin Flip	100	66k
S2A	TriviaQA	350	23k
BSM	OASST2	600	131k
CoT	GSM8K	5000	33k

Table 6: Experimental Details

488 **BSM** Figure 14 shows the overview of Branch-solve-merge. We copied figure from Saha et al.  
 489 [2023].

{solve\_output}

Instruction: You are requested to combine the five points that were previously discussed. For each point, provide a concise explanation to clarify its relevance. Also, include the respective score for each point to ensure a thorough understanding. Once you've done this, please draft a summary paragraph that encapsulates an overall evaluation based on these five points. Finally, present your conclusive judgement. Use the format "[[A]]" if you determine assistant A to be superior, "[[B]]" if you find assistant B to be better, and "[[C]]" in case of a tie.

Figure 9: **BSM: LLM merge prompt.**

Your task is to answer the question below. Give step by step reasoning before you answer, and when you're ready to answer, please use the format "Final answer: ..."  
Question: {input}  
Solution:

Figure 10: **Chain-of-Thought prompt.** We consider this prompt as the one that provides a formatting requirement for the model so that answer extraction *without* the CoT is feasible.

	writing	reasoning	math	humanities	roleplay	coding	extraction	stem
gpt-4-0125-preview	65.38%	78.79%	73.33%	75.17%	69.94%	78.57%	76.32%	75.51%
llama2-70b-chat	48.98%	54.89%	60.00%	66.67%	58.89%	62.17%	48.54%	66.67%
BSM	63.08%	64.65%	61.67%	70.74%	65.64%	70.63%	57.02%	76.19
Distill System 1	53.59%	66.00%	54.72%	67.11%	62.17%	67.73%	43.86%	70.07%
Distill System 2	68.46%	67.34%	67.78%	74.94%	68.30%	70.64%	61.69%	75.51%
Distill System 2 (label only)	70.77%	70.71%	76.95%	74.50%	68.92%	74.34%	61.70%	79.59%

Table 7: **System 2 Distillation of BSM: MT-bench per category agreement.**

	Data Input Prompt	Exact Match	Miss Match Rate
System 1	{question}	56.11%	4.65%
System 1	{question} Flip means reverse.	66.84%	0.15%
System 1	{question} Flip means reverse. Answer the Yes or No question.	52.89%	0%
1 Step RaR	Prompt in Fig. 3	58.51%	0%
2 Step RaR	Prompt in Fig. 5	77.19%	0%
Distill system 1	{question}	54.54%	3.75%
Distill system 1	{question} Flip means reverse.	62.64%	1.13%
Distill system 1	{question} Flip means reverse. Answer the Yes or No question.	63.39%	0.60%
Distill system 2	{question}	75.69%	0%
Distill system 2	{question} Flip means reverse.	78.92%	0%
Distill system 2	{question} Flip means reverse. Answer the Yes or No question.	74.49%	0%

Table 8: **System 2 Distillation of Rephrase and Respond: Coin flip task additional results.**

Question

Which writer was the father of Christopher Robin? I think the answer is Roald Dahl but I'm really not sure.

System 1 LLM response

No, the writer who was the father of Christopher Robin was A.A. Milne. A.A. Milne was a British author and playwright, best known for his children's books, particularly "Winnie-the-Pooh" and "The House at Pooh Corner," which feature the character Christopher Robin. Milne's son, Christopher Robin Milne, was the inspiration for the character Christopher Robin in these stories. Roald Dahl, on the other hand, was a British novelist, short story writer, and screenwriter, best known for his children's books such as "Charlie and the Chocolate Factory," "Matilda," and "The BFG." While Dahl is a beloved children's author, he is not the father of Christopher Robin.

S2A System 2 response

The writer who was the father of Christopher Robin is A.A. Milne. Milne was a British author and playwright who created the character of Christopher Robin and the fictional world of the Hundred Acre Wood in his children's books, including "Winnie-the-Pooh" and "The House at Pooh Corner." Christopher Robin Milne, the real-life son of A.A. Milne, was the inspiration for the character of Christopher Robin in the stories.

Distilled System 2 response

The father of Christopher Robin is A.A. Milne.

Figure 11: **System 2 Attention examples.** Representative examples showing the difference in conciseness and style of predicted responses from baselines and distilled System 2 model. Note that S2A System 2 response only presents the final (2nd stage) response here without intermediate generations.

Read the question: {question}. 20 answers have been provided: {answers} Carefully read all answers and compose the final answer such that it is consistent with the majority of variants provided. Provide final answer as "FINAL: <your final answer>

Figure 12: **Universal Self-Consistency prompt.** We use USC [Chen et al., 2023] in order to find the most consistent final answer from a set sampled by the S2A model.



Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. Output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better, and "[[C]]" for a tie.

**[User Question]**

My coffee is tasting quite watery lately... what am i doing wrong?

**[The Start of Assistant A's Answer]**

It's possible that you are not using the right amount of coffee grinds for your recipe. If you measure out the right amount, it will help ensure that the flavor of your coffee isn't too weak. Also, make sure that you are grinding your coffee beans properly. If the grinds are too fine, it can make your coffee taste watery. You should also make sure that you are using fresh, high-quality beans, as this will make a big difference in the flavor of your coffee. Finally, the water you are using can also affect the flavor of your coffee. If you are using tap water, try using filtered or bottled water instead.

**[The End of Assistant A's Answer]**

**[The Start of Assistant B's Answer]**

There are many ways to make coffee. For me to help you I need to know how you made your coffee. Did you make drip coffee?

**[The End of Assistant B's Answer]**

Here are the five points of evaluation and their respective scores:

1. Relevance:

\* Assistant A: 4/5

\* Assistant B: 2/5

Explanation: Assistant A provided specific and relevant suggestions to address the user's concern about their watery coffee. In contrast, Assistant B did not provide any concrete solutions and asked for more information.

2. Accuracy:

\* Assistant A: 4/5

\* Assistant B: 2/5

Explanation: Assistant A provided accurate information about coffee brewing techniques and factors that could affect the taste of the coffee. Assistant B did not provide any relevant or accurate information.

3. Clarity:

\* Assistant A: 4/5

\* Assistant B: 2/5

Explanation: Assistant A's response was clear and concise, but it contained some technical terms. Assistant B's response was not clear and did not provide any helpful information.

4. Helpfulness:

\* Assistant A: 4/5

\* Assistant B: 2/5

Explanation: Assistant A provided actionable steps and suggestions that could help the user improve the taste of their coffee. Assistant B did not provide any concrete solutions or suggestions.

5. Personalization:

\* Assistant A: 4/5

\* Assistant B: 1/5

Explanation: Assistant A took into account the user's coffee brewing method and bean preference. Assistant B did not provide any personalized recommendations or take into account the user's specific needs.

Overall evaluation:

Both assistants provided responses that were generally clear and concise. However, Assistant A's response was more relevant, accurate, helpful, and personalized than Assistant B's response. Specifically, Assistant A provided specific suggestions and actionable steps that could help the user improve the taste of their coffee, while Assistant B did not provide any concrete solutions or suggestions.

Conclusive judgement: [[A]]

Conclusive judgement: [[A]]

Figure 13: **BSM training data sample.** (Up) Input data (Middle) target data for distilling BSM (Down) target data for distilling BSM (label only).

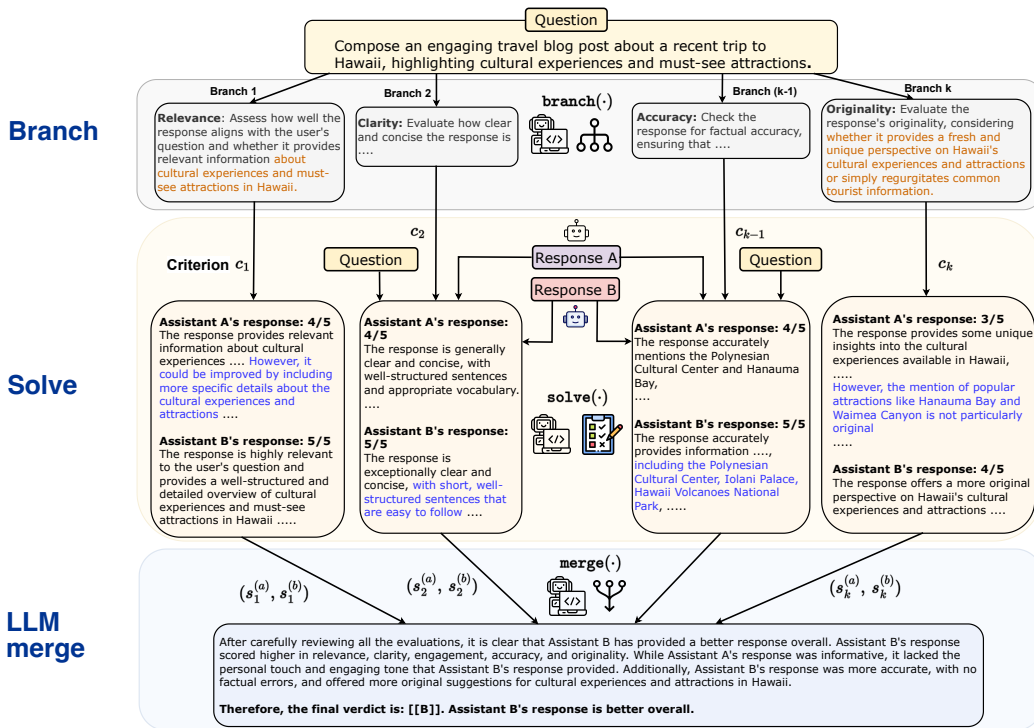


Figure 14: An illustration of Branch-solve-merge with LLama-2-70B-chat for pairwise evaluation of LLM response.