
Evaluating Stability and Interchangeability of Large Language Models in Mathematical Reasoning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Chain-of-Thought (CoT) prompting has significantly advanced the reasoning capa-
2 bilities of large language models (LLMs). While prior work focuses on improving
3 model performance through internal reasoning strategies, little is known about the
4 interchangeability of reasoning across different models. In this work, we explore
5 whether a partially completed reasoning chain from one model can be reliably
6 continued by another model, either within the same model family or across families.
7 We achieve this by assessing the sufficiency of intermediate reasoning traces as
8 transferable scaffolds for logical coherence and final answer accuracy. We interpret
9 this interchangeability as a means of examining inference-time trustworthiness,
10 probing whether reasoning remains both coherent and reliable under model sub-
11 stitution. Using token-level log-probability thresholds to truncate reasoning at
12 early, mid, and late stages from our baseline models, Gemma-3-4B-IT and LLaMA-
13 3.1-70B-Instruct, we conduct continuation experiments with Gemma-3-1B-IT and
14 LLaMA-3.1-8B-Instruct to test intra-family and cross-family behaviors. Our eval-
15 uation pipeline leverages truncation thresholds with a Process Reward Model (PRM),
16 providing a reproducible framework for assessing reasoning stability via model
17 interchange. Evaluations with a PRM reveal that hybrid reasoning chains often
18 preserve, and in some cases even improve, final accuracy and logical structure.
19 Our findings point towards interchangeability as an emerging behavioral property
20 of reasoning models, offering insights into new paradigms for reliable modular
21 reasoning in collaborative AI systems.

22 1 Introduction

23 Chain of Thought (CoT) prompting emerged as powerful mechanism to improve the reasoning
24 capabilities of large language models (LLMs) by encouraging intermediate structured reasoning
25 steps before arriving at a final answer [Wei et al., 2023]. Previous work has explored how CoTs
26 improve individual model performance even in zero-shot settings [Kojima et al., 2023, Zhang et al.,
27 2022, Jin et al., 2024]. More recently, Hebenstreit et al. [2024] examined the transferability of entire
28 CoT sequences by evaluating whether rationale prompts discovered on one model could generalize
29 reasoning strategies across a range of models and tasks. However, it remains unclear to what extent
30 reasoning trajectories are interchangeable when only partially reused. In light of this, our aim is
31 to answer the central research question: *To what extent can the modular decomposition of complex*
32 *mathematical reasoning tasks enhance the zero-shot performance and interpretability of Large*
33 *Language Models, when utilizing a collaborative framework that includes both intra-family and*
34 *cross-family LLMs?*

35 In this work, we investigate the process-level interchangeability in language model reasoning by
36 evaluating how well different models can continue the CoT of another’s midstream. We begin with

37 full CoT traces generated by a strong base model (e.g., Gemma-3-4B-IT and LLaMA-3.1-70B-
38 Instruct), recording token-level log-probabilities to guide strategic truncation at 25%, 50%, and
39 75% of the cumulative log-probability, capturing early, mid, and late stages of reasoning based on
40 informativeness. From these truncated points, alternative models (including those from different
41 families or architectures) are tasked with continuing the reasoning process using only truncated
42 intermediate steps as input. We then assess not only accuracy, but also the coherence, semantic
43 alignment, and logical consistency of the full reasoning chain, by using a Process Reward Model
44 (PRM) trained to evaluate multi-step mathematical reasoning performance. Ultimately, our aim is to
45 characterize how steady transferability depends on truncation point, model pairing, and reasoning
46 domain, yielding clearer interpretations into the dynamics of CoT continuation success that goes
47 beyond final answer accuracy.

48 Whereas prior work has explored how CoT prompting improves reasoning within individual models
49 [Wei et al., 2023], whether reasoning can be interchanged across models mid-process remains largely
50 unexamined.

51 We provide compelling early evidence that such a handoff is often successful within the same model
52 family. We show that a partially completed CoT from a strong model, such as Gemma-3-4B-IT,
53 can often be continued by another model of similar or lesser capacity within the same family. By
54 leveraging log-probability-based truncation and PRM-based scoring, we found that these hybrid
55 trajectories maintain high coherence and correctness with minimal loss in reasoning quality.

56 We found that this practice may not be suitable for all cross-family continuation pairings, as some
57 unreliably preserve quality and coherence of the reasoning chain in our experimentation. Our findings
58 expose distinctions across different model architectures and introduce a promising new paradigm for
59 collaborative reasoning, where high-capacity models can be reserved for the most uncertain portions
60 of a problem, allowing lighter models to reliably finish the remainder of the task.

61 **2 Related Works**

62 LLMs generate responses by autoregressively predicting outputs based on the preceding context,
63 which is learned during pre-training [OpenAI et al., 2024]. As a result, their output can fluctuate even
64 when prompted with identical inputs, introducing variability in reasoning trajectories [Amatriain,
65 2024]. This, coupled with the absence of structured reasoning mechanisms, often leads to inconsis-
66 tencies in multistep logical inference. Consequently, assessing the reliability and soundness of their
67 reasoning becomes increasingly complex and therefore requires a more thorough examination [Wang
68 et al., 2024].

69 To address these limitations, the concept of CoT prompting was introduced in Wei et al. [2023],
70 demonstrating that instructing LLMs to reason step-by-step significantly improves performance
71 on complex tasks. In this approach, LLMs are prompted to generate a series of short statements
72 that mimic the logical process a person might use to solve a problem. Experiments revealed that
73 CoT prompting enables models to achieve strong results in tasks of arithmetic, commonsense, and
74 symbolic reasoning [Wei et al., 2023].

75 In an effort to enhance LLM reasoning abilities with CoT prompting, Wang et al. [2023b] introduces
76 self-consistency to replace the single greedy decoding path in traditional CoT prompting [Wei et al.,
77 2023]. Their method samples a variety of reasoning paths and identifies the most consistent answer by
78 marginalizing across all possibilities [Wang et al., 2023b]. Beyond improving accuracy, this approach
79 highlights the inherent diversity of reasoning paths within a single model, suggesting that multiple
80 valid chains of reasoning can coexist.

81 Initiatives have also been put forward to extend and refine CoT prompting. Unlike traditional CoT
82 where each step is independent, Least to Most prompting breaks difficult problems into sequential sub-
83 problems where the outputs of previous steps are the inputs for the next [Zhou et al., 2023]. Moreover,
84 recent efforts have examined the effects of partial or truncated CoT on model outputs. Lanham et al.
85 [2023] measure faithfulness by truncating generated CoT at various points and re-prompting the
86 model with only the partial reasoning.

87 Past research has indirectly measured the reasoning ability of LLMs by evaluating them on down-
88 stream reasoning tasks such as question answering or multi-hop inference [Huang and Chang, 2023].
89 Though, relying on the accuracy of the end task or the success rates is not indicative of step-by-step

90 reasoning capability. Huang and Chang [2023] also explain that current performance measures mix
91 reasoning ability with task knowledge, resulting in reasoning that cannot be evaluated in isolation.
92 To resolve this, subsequent work [Nguyen et al., 2024] aims at reasoning process analysis directly,
93 testing for logical coherence of individual steps, which provides more straightforward methods of
94 reasoning quality evaluation.

95 LLMs frequently make errors when solving mathematical problems step-by-step, making it essential
96 to identify where the errors occurred during the reasoning process [Zheng et al., 2025]. As a result,
97 PRMs have been developed as a direct solution to the shortcomings of traditional indirect evaluation
98 methods, which only assess final answers. PRMs are specifically designed to evaluate the correctness
99 of each individual reasoning step, providing feedback that helps guide policy models toward more
100 accurate and reliable mathematical reasoning [Zheng et al., 2025, Zhang et al., 2025]. PRMs output a
101 score or probability that represents the model’s confidence that the reasoning step is logically sound
102 and contributes productively to problem resolution.

103 **3 Methodology**

104 We introduce a novel chain-splitting approach grounded in cumulative log-probability, whereby
105 complete solutions are truncated at points of varying model confidence from an initial baseline
106 model and then continued by a second continuation model. The methodology proceeds in three
107 components: (Section 3.1) reasoning chain generation, (Section 3.2) chain truncation via cumulative
108 log-probability, and (Section 3.3) model interchange protocols.

109 **3.1 Reasoning Chain Generation**

110 We use an initial model to generate complete reasoning chains for each problem in the test set. Each
111 generation is performed with temperature set at 0.7, allowing a moderate degree of stochasticity in
112 token sampling while still favoring high-probability continuations. Let the complete output chain be
113 a sequence of tokens $r = \{t_1, t_2, \dots, t_n\}$, with corresponding log-probabilities $\{\ell_1, \ell_2, \dots, \ell_n\}$. We
114 compute the cumulative log-probability up to position i as $L_i = \sum_{j=1}^i \ell_j$.

115 This sequence $\{L_1, L_2, \dots, L_n\}$ defines the internal flow of confidence of the model throughout the
116 reasoning process.

117 **3.2 Chain Truncation via Log-Probability Thresholding**

118 To identify semantically meaningful split points in the chain, we define three thresholds based on the
119 total log-probability L_n :

- 120 • **25% truncation:** first index i such that $L_i \geq 0.25L_n$
- 121 • **50% truncation:** first index i such that $L_i \geq 0.50L_n$
- 122 • **75% truncation:** first index i such that $L_i \geq 0.75L_n$

123 For each threshold $\alpha \in \{0.25, 0.50, 0.75\}$, we extract the prefix $r_{1:k}$, where

$$k = \min\{i : L_i \geq \alpha L_n\}.$$

124 This results in three partially completed reasoning traces per problem, each grounded in the model’s
125 own internal confidence progression.

126 **3.3 Model Interchange Protocol**

127 Each truncated prefix is combined with a consistent CoT template meant for interchange that includes
128 the original question, and the resulting prompt is provided to a secondary continuation model (further
129 details in Section 8.1). We consider both intra-family and cross-family model pairings more precisely
130 defined in Section 4.2. Each continuation model generates a single completion for each prefix using a
131 temperature of 0.7, introducing controlled randomness to reflect typical sampling conditions while
132 preserving coherence. These continuations are concatenated with the original prefix to form hybrid
133 reasoning chains, which are then run through post-processing to extract the final answer using simple
134 rule-based extraction.

135 All in all, for each problem instance, we obtain: one complete chain from the baseline generator, and
136 multiple hybrid chains resulting from different continuation models and truncation depths (details in
137 Section 4.2).

138 4 Experimental Setup

139 We now outline the experimental conditions under which our chain-splitting framework was evaluated.
140 This includes: (Section 4.1) the dataset selected to benchmark reasoning difficulty and domain
141 coverage, (Section 4.2) the models used for initial generation and continuation, and (Section 4.5) the
142 metrics employed to quantify the quality of reasoning, compatibility, and the impact of performance
143 on model interchanges.

144 4.1 Dataset Selection

145 An extensive dataset was carefully selected to capture a range of reasoning complexities and domain-
146 specific scenarios.

- 147 • MATH [Hendrycks et al., 2021]: consists of 12,500 high-school and college-level mathe-
148 matical problems that span diverse topics and demanding multi-step solutions, providing
149 rigorous testing to evaluate advanced mathematical reasoning and generalization.

150 For our experiments, we evaluated models exclusively on the test splits of the MATH dataset,
151 consisting of 5,000 questions.

152 4.2 Model Selection and Configuration

153 We adopt Qwen2.5-PRM [Zheng et al., 2025] as our primary Process Reward Model, due to its
154 fine-tuning on structured multi-step mathematical datasets such as PRM800K [Lightman et al.,
155 2023] and Math-Shepherd [Wang et al., 2023a]. Qwen2.5-PRM is an instruction-tuned variant of
156 Qwen2.5-Math-7B and supports token-level log-probability outputs.

157 For model interchange experiments, we select two baseline models and two continuation models:

158 4.3 Baseline

159 To establish a baseline for reasoning quality, we employ Gemma-3-4B-IT and LLaMA-3.1-70B-
160 Instruct to generate CoT exemplars, two state-of-the-art instruction-tuned models from distinct
161 architectural lineages.

- 162 • Gemma-3-4B-IT [Team et al., 2025], an instruction-tuned variant from the Gemma 3 model
163 family developed by Google Deepmind with 4 Billion parameters is used to generate
164 complete Chain-of-Thought reasoning paths, tuned to Gemma’s architecture.
- 165 • LLaMA-3.1-70B-Instruct [Grattafiori et al., 2024], a large scale variant from the LLaMA 3
166 model family developed by Meta AI with 70 Billion parameters is used to generate complete
167 Chain-of-Thought reasoning paths, tuned to LLaMA’s architecture.

168 4.4 Continuation

- 169 • Gemma-3-1B-IT [Team et al., 2025], a lightweight variant from the same Gemma 3 model
170 family, is used to evaluate how well reasoning chains can be completed by a structurally
171 similar but smaller model.
- 172 • LLaMA 3.1-8B-Instruct Grattafiori et al. [2024] representing a different architectural lineage
173 helps enable testing interchangeability across distinct LLM families. For brevity, we refer to
174 the aforementioned models as Gemma and LLaMA respectively for the remainder of this
175 paper.

176 On the MATH dataset, Gemma 3-1B-IT and Gemma 3-4B-IT performed with accuracies of 48.0%
177 and 75.6% respectively [Team et al., 2025]. Moreover, Llama-3.1-8B-Instruct and Llama-3.1-70B-
178 Instruct performed with accuracies 47.2% and 65.7% [Yang et al., 2024]. We observe that the two

179 base models exhibit similar performance levels and, likewise, that the two continuation models
180 perform comparably.

181 All models were prompted using one consistent CoT templates, either the interchange or full-run
182 variant, as detailed in Section 8.1.

183 4.5 Evaluation Metrics

184 The hybrid reasoning chains generated were evaluated using a multifaceted set of metrics designed to
185 assess accuracy, variability, and the impact of model interchanges on reasoning coherence and final
186 outcomes. Specifically, we consider the following four core metrics:

- 187 • **Answer Accuracy:** Accuracy is defined as the proportion of final answers from generation
188 that exactly match those from ground-truth solutions. This metric represents the model’s
189 ability to arrive at the correct final result through its reasoning chain.
- 190 • **PRM Score:** As a PRM is available for scoring, we additionally report average PRM-
191 assigned scores that capture the internal likelihood and coherence of a given chain regardless
192 of final correctness. We define the PRM score A' as the average plausibility score assigned
193 to each reasoning step in a chain of n steps:

$$A' = \frac{1}{n} \sum_{i=1}^n \text{PRM}(s_i)$$

194 where s_i denotes the i -th step in the chain. While traditional accuracy reflects outcome-level
195 correctness, A' provides a step-level assessment of reasoning quality.

- 196 • **Normalized Relative Gain (NRG):** This metric quantifies whether incorporation of reasoning
197 from another model helps or hinders performance. Given the accuracies of the original
198 model A and B , and hybrid accuracies A' (Model A prefix + Model B suffix) and B' (Model
199 B prefix + Model A suffix), we define:

$$\text{NRG}_A = \frac{A' - A}{A}, \quad \text{NRG}_B = \frac{B' - B}{B}.$$

200 Positive values indicate a performance gain from model interchange, while negative values
201 reflect degradation.

- 202 • **Cross-Model Degradation (XMD):** This metric captures the extent to which the continuation
203 of a model degrades the original reasoning trajectory. It is defined as:

$$\text{XMD}_{A \rightarrow B} = \frac{A - B'}{A}, \quad \text{XMD}_{B \rightarrow A} = \frac{B - A'}{B}.$$

204 XMD provides a normalized measure of reasoning incompatibility, where higher values
205 indicate more severe disruptions introduced by the cross-model continuation.

206 5 Results

207 We present results across the MATH benchmark to evaluate model interchangeability across truncation
208 points. Through our proposed metrics, we look to determine whether model continuation works to
209 improve or disrupt the original reasoning trajectory. Experimental results were obtained using the
210 Runpod cloud platform, leveraging NVIDIA H100 PCIe GPUs over approximately 250 GPU hours.

211 5.1 Full Chain-of-Thought Results

212 To establish a baseline, we first evaluate each model’s performance using end-to-end CoT reasoning
213 applied without interruption. For every example in the benchmark, the model is prompted to reason
214 step by step to completion, producing a complete trajectory from question to final answer. We report
215 results in terms of final answer accuracy and step-level reasoning score as seen in Table 1.

216 This baseline allows us to quantify native reasoning strengths and weaknesses of each model without
217 the effects of interchange.

Model	Dataset	Accuracy (%)	PRM
Gemma-3-4B-IT	MATH	68.06%	0.8952
Gemma-3-1B-IT	MATH	36.28%	0.7904
LLaMA-3.1-70B-Instruct	MATH	60.80%	0.8725
LLaMA-3.1-8b-Instruct	MATH	47.76%	0.8522

Table 1: Performance of reasoning chains fully generated by each model (i.e., with no handoff or interchange from another model) on the MATH dataset.

218 5.2 Interchanged Chain-of-Thought Results

219 Thereafter, to gauge the interchangeability of reasoning processes across different models, we evaluate
 220 the completion of truncated CoT traces. Each reasoning chain is strategically truncated based on
 221 cumulative log-probability thresholds (25%, 50%, 75%), representing early, mid, and late points
 222 in the reasoning process. Subsequently, alternative models are assigned to continue the truncated
 223 reasoning chains through to completion.

224 We report performance for all continuation combinations, including accuracy, step-level scores, and
 225 coherence ratings as seen in Table 2 & Table 3. This analysis unveils the extent to which partial
 226 reasoning from one model can be reliably extended by another, highlighting cases of both successful
 227 handoff and systematic breakdowns that point to the limits of reasoning interchangeability.

Truncation	Continuation	Accuracy (%)	PRM	NRG	XMD
25%	Gemma-3-1B-IT	41.76%	0.7966	0.3678	0.3864
25%	LLaMA-3.1-8B-Instruct	43.60%	0.8393	0.3196	0.3594
50%	Gemma-3-1B-IT	49.86%	0.8002	0.3786	0.2674
50%	LLaMA-3.1-8B-Instruct	53.24%	0.8585	0.264	0.2177
75%	Gemma-3-1B-IT	55.26%	0.8032	0.3500	0.1881
75%	LLaMA-3.1-8B-Instruct	63.80%	0.8697	0.1853	0.0626

Table 2: Performance of hybrid reasoning chains by truncation point and continuation model on MATH dataset, using a fully generated CoT from Gemma-3-4B-IT.

Truncation	Continuation	Accuracy (%)	PRM	NRG	XMD
25%	Gemma-3-1B-IT	36.16%	0.7566	-0.1137	0.4053
25%	LLaMA-3.1-8B-Instruct	42.18%	0.8323	-0.0150	0.3062
50%	Gemma-3-1B-IT	38.50%	0.7730	-0.0968	0.3668
50%	LLaMA-3.1-8B-Instruct	46.26%	0.8456	-0.0072	0.2391
75%	Gemma-3-1B-IT	41.98%	0.7811	-0.0827	0.3095
75%	LLaMA-3.1-8B-Instruct	50.06%	0.8543	-0.0002	0.1766

Table 3: Performance of hybrid reasoning chains by truncation point and continuation model on MATH dataset, using a fully generated CoT from LLaMA-3.1-70B-Instruct.

228 6 Discussion

229 Our observations uncover degradation in performance when cross-family models are tasked to
 230 continue reasoning midstream initiated by a partially completed CoT. There are several factors likely
 231 responsible for this downgrade in performance:

232 6.1 Style and Representational Compatibility

233 A consistent disparity between intra-family and cross-family continuation highlights representational
234 compatibility as a key factor in multi-model reasoning. Despite receiving high confidence chains,
235 cross-family continuations (e.g., Gemma-3-4B-IT→LLaMA-3.1-8B-Instruct and LLaMA-3.1-70B-
236 Instruct→Gemma-3-1B-IT) often fail to maintain correct reasoning. For instance, when LLaMA-
237 3.1-70B-Instruct’s chain is continued by Gemma-3-1B-IT, accuracy falls to 36.16% at the 25%
238 mark—nearly a 40% relative decline compared to the base model’s 60.80% full-chain accuracy—
239 with a corresponding negative NRG (-0.1137). Similarly, continuations from Gemma-3-4B-IT
240 into LLaMA-3.1-8B-Instruct underperform early on (43.60% at 25%) despite access to confident
241 reasoning prefixes, yielding a lower NRG of 0.3196 compared to intra-family continuation at the
242 same depth (Gemma-3-4B-IT→Gemma-3-1B-IT, 0.3678), indicating that these prefixes do not fully
243 overcome differences in architecture and reasoning style. This pattern suggests a reasoning bias: each
244 model family tends to rely more on its own reasoning patterns, which may result from structural
245 differences between the families.

246 These results are consistent with previous work [Liu et al., 2023], which noted that structural
247 differences between model families (GPT-4 in their case) can limit cross-model reasoning transfer,
248 particularly for complex, multi-step reasoning tasks. While LLaMA models generate coherent chains
249 within their own family, their internal reasoning representations differ from Gemma’s, which may
250 hinder smooth continuation across families. This is supported by consistently high XMD values across
251 truncation points (e.g., 0.4053 at 25% and 0.3095 at 75% for LLaMA-3.1-70B-Instruct→Gemma-
252 3-1B-IT), suggesting that reasoning coherence is not fully maintained even as longer prefixes are
253 available. High-confidence reasoning prefixes do not appear sufficient to completely navigate these
254 differences, indicating that cross-family continuation is constrained by family-specific reasoning
255 tendencies.

256 In contrast, intra-family continuations show steady improvement with longer truncation depths. For
257 example, when Gemma-3-1B-IT continues from Gemma-3-4B-IT, accuracy rises from 41.76% at
258 25% to 55.26% at 75%, accompanied by moderate NRG values ($0.3678 \rightarrow 0.3500$) and decreasing
259 XMD ($0.3864 \rightarrow 0.1881$). Similarly when LLaMA-3.1-8B-Instruct continues from LLaMA-3.1-70B-
260 Instruct, performance increases from 42.18% to 50.06%, with NRG improving from -0.0150
261 to near-neutral (-0.0002) and XMD decreasing from 0.3062 to 0.1766. These patterns suggest that the
262 representational similarity between models supports a more stable continuation and better integration
263 of context.

264 6.2 Context Integration Overhead

265 When deployed late in the reasoning chain (e.g., at the 75% mark), smaller continuation models
266 such as Gemma-3-1B-IT and LLaMA-3.1-8B-Instruct must interpret and integrate extensive context
267 generated by larger base models (Gemma-3-4B-IT and LLaMA-3.1-70B-Instruct). As reasoning
268 sequences lengthen, models may face capacity limits that degrade performance. This bottleneck is
269 attributed to the finite “working memory” of LLMs and the compounding demands of maintaining
270 logical coherence across many steps [Shang et al., 2025]. The effect is especially pronounced when
271 models are required to interpret and continue reasoning from an externally provided chain rather than
272 generating all steps from scratch.

273 On the MATH dataset, truncation depth produces gradual improvements but does not eliminate the
274 performance gap relative to non-handoff baselines. For example, when continuing Gemma-3-4B-IT’s
275 reasoning, Gemma-3-1B-IT improves from 41.76% at 25% truncation to 55.26% at 75%, while
276 LLaMA-3.1-8B-Instruct rises from 43.60% to 63.80%. Similarly, when continuing LLaMA-3.1-
277 70B-Instruct, LLaMA-3.1-8B-Instruct achieves a smoother progression from 42.18% to 50.06%,
278 outperforming Gemma-3-1B-IT, which remains between 36.16% and 41.98%. These trends suggest
279 that architectural alignment facilitates smoother context integration in same-family continuations,
280 while representational mismatches in cross-family pairs disrupt coherent reasoning.

281 Despite improvements with longer prefixes, performance remains notably below that of fully self-
282 generated chains (Gemma-3-4B-IT: 68.06%, LLaMA-3.1-70B-Instruct: 60.80%). XMD values con-
283 firm this persistent overhead: even at the 75% truncation point, XMD remains non-negligible (0.0626
284 for Gemma-3-4B-IT→LLaMA-3.1-8B-Instruct and 0.1766 for LLaMA-3.1-70B-Instruct→LLaMA-
285 3.1-8B-Instruct), indicating incomplete recovery of original reasoning quality.

286 These observations highlight that truncation depth alone does not ensure effective reasoning trans-
287 fer. Although larger prefixes reduce uncertainty and contextual loss, architectural and stylistic
288 compatibility between base and continuation models remains the key factor determining success.

289 6.3 Error Amplification

290 Minor inconsistencies or ambiguities in early reasoning steps, especially when generated by a
291 different model, can compound as LLaMA or Gemma continue the reasoning process. With limited
292 steps remaining to revise earlier logic (particularly in final-answer-only completions), both models
293 struggle to recover from upstream errors. These results suggest that effective interoperability in multi-
294 step reasoning depends on both model capability and the degree of representational and contextual
295 alignment across reasoning steps.

296 On the MATH dataset, when reasoning chains generated by Gemma-3-4B-IT or LLaMA-3.1-70B-
297 Instruct are truncated and continued by smaller models at various points (25%, 50%, 75%), perfor-
298 mance declines in proportion to both truncation depth and cross-family divergence. When Gemma-
299 3-4B-IT serves as the base, continuation by Gemma-3-1B-IT (intra-family) improves steadily from
300 41.76% at 25% to 55.26% at 75%, with NRG values rising from 0.3678 to 0.3500 and XMD decreas-
301 ing from 0.3864 to 0.1881. Cross-family continuation by LLaMA-3.1-8B-Instruct performs competi-
302 tively (43.60%→63.80%) but shows slightly lower NRG (0.3196→0.1853), indicating weaker
303 efficiency in utilizing the provided context. Longer prefixes appear to partially reduce representational
304 mismatch, leading to more consistent performance over time.

305 When LLaMA-3.1-70B-Instruct serves as the base model, Gemma-3-1B-IT continuations per-
306 form substantially worse (36.16~41.98% across truncation points) with persistently high XMD
307 (0.4053→0.3095) and negative NRG (-0.1137→-0.0827), suggesting limited transfer across fami-
308 lies. Intra-family continuation by LLaMA-3.1-8B-Instruct performs more stably, reaching 50.06% at
309 75% with NRG improving from -0.0150 to -0.0002 and XMD decreasing from 0.3062 to 0.1766,
310 reflecting more consistent reasoning integration within the same family.

311 Comparing fully generated chains with the hybrid results (Table 1), Gemma-3-4b and LLaMA-3.1-
312 70B-Instruct still substantially outperform their continuations (68.06% and 60.80%, respectively).
313 However, their smaller counterparts, especially Gemma-3-1B-IT, demonstrate partial to considerable
314 recovery when inheriting sufficiently long prefixes, suggesting that similar architecture and tokeniza-
315 tion structures enhance transfer performance. As seen over 25%/50%/75% truncations, intra-family
316 continuation (Gemma-3-4b→Gemma-3-1b) improves from 41.76%→55.26% (+13.5 pp), and even
317 cross-family continuation (Gemma-3-4B-IT→LLaMA-3.1-8B-Instruct) exhibits a greater net gain of
318 43.60%→63.80% (+20.2 pp). In contrast, continuations from LLaMA-3.1-70b displayed weaker
319 recovery with LLaMA-3.1-8B-Instruct rising only +7.9 pp (42.18% 50.06%), and Gemma-3-1B-IT
320 gains just +5.8 pp (36.16%→41.98%). As truncation length increases, reasoning becomes more
321 coherent, but full recovery is still unattainable, lending credence to how small representational gaps
322 can compound through multi-step reasoning chains.

323 7 Conclusion

324 In this work, we introduced a novel framework for evaluating midstream interchangeability in
325 large language models, grounded in a chain-splitting paradigm based on cumulative log-probability.
326 By systematically truncating the reasoning chains generated by our base models and appending
327 completions from either intra-family or cross-family models, we directly measured the stability and
328 coherence of hybrid reasoning trajectories. Our experiments on MATH demonstrate that model family
329 alignment plays a decisive role in the success or failure of such hybrid chains. While intra-family
330 continuations generally preserved reasoning quality on simpler tasks, cross-family continuations
331 often struggled to maintain coherence with the partial chains, despite comparable model performance
332 as referenced in Section 4.2. This suggests that models like Gemma and LLaMA may be better
333 aligned to continue reasoning within their own family than across different architectures.

334 These findings challenge previous assumptions about model modularity in contemporary NLP. Despite
335 architectural advances and increasing performance parity across model families, our results suggest
336 that inter-model transfer in multi-step reasoning remains fragile, constrained by differences in stylistic
337 alignment, latent variable encoding, and contextual integration. The observed breakdowns reveal a

338 significant gap between individual task performance and interoperability in reasoning, which is an
339 area that has received insufficient attention in LLM evaluation.

340 More broadly, our work highlights the need for new approaches that preserve consistent semantic
341 reasoning across different language models. As research advances toward compositional and multi-
342 agent LLM systems, reliable interchangeability will become essential, not solely for efficiency, but
343 also for alignment, verification, and interpretability. Our methodology provides an initial framework
344 for diagnosing and quantifying this interoperability gap in a systematic, data-driven manner.

345 **Limitations**

- 346 • **Single Completion Runs:** All experiments were conducted using deterministic continuations.
347 While this reflects realistic deployment scenarios, it limits our understanding of variance
348 under sampling. Future work should evaluate robustness using multiple stochastic rollouts.
- 349 • **Task Domain Scope:** Our evaluation is confined to math reasoning (MATH). It remains
350 unclear whether interchangeability generalizes to commonsense, scientific, or multimodal
351 reasoning tasks.
- 352 • **Domain-Specific PRMs:** We employed a math-specific Process Reward Model (PRM).
353 Evaluating reasoning quality in other domains will require retraining or adapting PRMs
354 tailored to those reasoning distributions.

355 **Future Work**

- 356 • **Cross-Domain Generalization:** Evaluate model interchangeability on tasks such as
357 commonsense QA, multi-hop retrieval, scientific explanation, and instruction-following,
358 where reasoning formats may be more variable or implicit.
- 359 • **Adaptive Truncation Strategies:** Rather than using static log-probability thresholds
360 (25/50/75%), future work could explore dynamic segmentation based on reasoning
361 content, semantic shifts, or model uncertainty.
- 362 • **Collaborative Model Architectures:** Deploy multi-agent or multi-model reasoning pipelines
363 in production environments (e.g., tutoring systems, scientific assistants) to study tradeoffs in
364 latency, memory, and correctness.

365 **8 Appendix**

366 **8.1 Prompting**

Standardized Prompt

Full-Run Prompt :

System message: "You are a helpful assistant that solves problems step by step.
Please provide clear reasoning with numbered steps and conclude with your final answer."

User message: "Solve this problem step by step:
Question: ['question']"

Interchange Prompt :

System message: "You are a helpful assistant that solves problems step by step.
Please provide clear reasoning with numbered steps and conclude with your final answer."

User message: "Solve this problem step by step:
Question: ['question'] ['truncated reasoning']"

367
368 This prompt standardization ensures comparability in reasoning styles across models; slight variations
369 were applied where necessary to accommodate model-specific tokenization or formatting requirements
370 without altering the intended instructions or task semantics.

371 **References**

372 Xavier Amatriain. Prompt design and engineering: Introduction and advanced methods, 2024. URL
373 <https://arxiv.org/abs/2401.14423>.

374 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
375 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan,
376 Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev,
377 Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru,
378 Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak,
379 Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu,
380 Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle
381 Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego
382 Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,
383 Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel
384 Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon,
385 Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan
386 Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet,
387 Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,
388 Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie
389 Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua
390 Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak,
391 Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley
392 Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence
393 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas
394 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri,
395 Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie
396 Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes
397 Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne,
398 Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal
399 Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
400 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
401 Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie
402 Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana
403 Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie,
404 Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon
405 Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan,
406 Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas
407 Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami,
408 Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti,
409 Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier
410 Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao
411 Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song,
412 Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe
413 Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya
414 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei
415 Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu,
416 Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit
417 Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury,
418 Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer,
419 Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu,
420 Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido,
421 Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changan Wang, Changkyu
422 Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer,
423 Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu,
424 Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingakang Wang, Duc
425 Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily
426 Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers,
427 Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank
428 Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Sweet,

- 429 Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan,
430 Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph,
431 Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog,
432 Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James
433 Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny
434 Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings,
435 Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai
436 Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik
437 Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle
438 Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng
439 Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabza, Manav Avalani, Manish
440 Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim
441 Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle
442 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang,
443 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam,
444 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier,
445 Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia
446 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro
447 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani,
448 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghobham Murthy,
449 Raghu Nayani, Rahul Mitra, Rangarabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin
450 Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu,
451 Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh
452 Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay,
453 Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang,
454 Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie
455 Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta,
456 Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman,
457 Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun
458 Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria
459 Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru,
460 Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz,
461 Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv
462 Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
463 Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait,
464 Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The
465 llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 466 Konstantin Hebenstreit, Raphael Praas, Laura P. Kiesewetter, and Matthias Samwald. A comparison
467 of chain-of-thought reasoning strategies across datasets and models. *PeerJ Computer Science*, 10:
468 e1999, 2024. doi: 10.7717/peerj-cs.1999. URL <https://doi.org/10.7717/peerj-cs.1999>.
- 469 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
470 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL
471 <https://arxiv.org/abs/2103.03874>.
- 472 Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey,
473 2023. URL <https://arxiv.org/abs/2212.10403>.
- 474 Feihu Jin, Yifan Liu, and Ying Tan. Zero-shot chain-of-thought reasoning guided by evolutionary
475 algorithms in large language models, 2024. URL <https://arxiv.org/abs/2402.05376>.
- 476 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
477 language models are zero-shot reasoners, 2023. URL <https://arxiv.org/abs/2205.11916>.
- 478 Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Her-
479 nandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamilė Lukošiuūtė, Karina
480 Nguyen, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Larson, Sam
481 McCandlish, Sandipan Kundu, Saurav Kadavath, Shannon Yang, Thomas Henighan, Timothy
482 Maxwell, Timothy Telleen-Lawton, Tristan Hume, Zac Hatfield-Dodds, Jared Kaplan, Jan Brauner,
483 Samuel R. Bowman, and Ethan Perez. Measuring faithfulness in chain-of-thought reasoning, 2023.
484 URL <https://arxiv.org/abs/2307.13702>.

- 485 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
486 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
487 *arXiv:2305.20050*, May 2023. URL <https://arxiv.org/abs/2305.20050>.
- 488 Hanmeng Liu, Zhiyang Teng, Leyang Cui, Chaoli Zhang, Qiji Zhou, and Yue Zhang. Logicot:
489 Logical chain-of-thought instruction-tuning. *arXiv preprint arXiv:2305.12147*, May 2023. doi:
490 10.48550/arXiv.2305.12147. URL <https://arxiv.org/abs/2305.12147>.
- 491 Minh-Vuong Nguyen, Linhao Luo, Fatemeh Shiri, Dinh Phung, Yuan-Fang Li, Thuy-Trang Vu, and
492 Gholamreza Haffari. Direct evaluation of chain-of-thought in multi-hop reasoning with knowledge
493 graphs, 2024. URL <https://arxiv.org/abs/2402.11199>.
- 494 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
495 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor
496 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian,
497 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny
498 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks,
499 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea
500 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen,
501 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung,
502 Dave Cummings, Jeremiah Carrier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,
503 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty
504 Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte,
505 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh,
506 Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross,
507 Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton,
508 Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton,
509 Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela
510 Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan,
511 Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan
512 Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt
513 Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic,
514 Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung,
515 Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa
516 Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov,
517 Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer
518 McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob
519 Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa,
520 Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano,
521 Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe,
522 Jakob Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel
523 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila
524 Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle
525 Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri,
526 Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl
527 Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar,
528 Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard,
529 Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie
530 Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie
531 Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak,
532 Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick
533 Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea
534 Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,
535 CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lillian Weng, Matt Wiethoff, Dave
536 Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu,
537 Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan
538 Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William
539 Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- 540 HaoYang Shang, Xuan Liu, Zi Liang, Jie Zhang, Haibo Hu, and Song Guo. United minds or
541 isolated agents? exploring coordination of llms under cognitive load theory. *arXiv preprint*

542 *arXiv:2506.06843*, June 2025. doi: 10.48550/arXiv.2506.06843. URL <https://arxiv.org/abs/2506.06843>.

543

544 Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej,
545 Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas
546 Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon,
547 Etienne Pot, Ivo Penchev, Gaël Liu, Francesco Visin, Kathleen Kenealy, Lucas Beyer, Xiaohai
548 Zhai, Anton Tsitsulin, Robert Busa-Fekete, Alex Feng, Noveen Sachdeva, Benjamin Coleman,
549 Yi Gao, Basil Mustafa, Iain Barr, Emilio Parisotto, David Tian, Matan Eyal, Colin Cherry, Jan-
550 Thorsten Peter, Danila Sinopalnikov, Surya Bhupatiraju, Rishabh Agarwal, Mehran Kazemi,
551 Dan Malkin, Ravin Kumar, David Vilar, Idan Brusilovsky, Jiaming Luo, Andreas Steiner, Abe
552 Friesen, Abhanshu Sharma, Abheesht Sharma, Adi Mayrav Gilady, Adrian Goedeckemeyer, Alaa
553 Saade, Alex Feng, Alexander Kolesnikov, Alexei Bendebury, Alvin Abdagic, Amit Vadi, András
554 György, André Susano Pinto, Anil Das, Ankur Bapna, Antoine Miech, Antoine Yang, Antonia
555 Paterson, Ashish Shenoy, Ayan Chakrabarti, Bilal Piot, Bo Wu, Bobak Shahriari, Bryce Petri
556 Charlie Chen, Charline Le Lan, Christopher A. Choquette-Choo, CJ Carey, Cormac Brick, Daniel
557 Deutsch, Danielle Eisenbud, Dee Cattle, Derek Cheng, Dimitris Pappas, Divyashree Shivakumar
558 Sreepathihalli, Doug Reid, Dustin Tran, Dustin Zelle, Eric Noland, Erwin Huizenga, Eugene
559 Kharitonov, Frederick Liu, Gagik Amirkhanyan, Glenn Cameron, Hadi Hashemi, Hanna Klimczak-
560 Plucińska, Harman Singh, Harsh Mehta, Harshal Tushar Lehri, Hussein Hazimeh, Ian Ballantyne,
561 Idan Szpektor, Ivan Nardini, Jean Pouget-Abadie, Jetha Chan, Joe Stanton, John Wieting, Jonathan
562 Lai, Jordi Orbay, Joseph Fernandez, Josh Newlan, Ju yeong Ji, Jyotinder Singh, Kat Black, Kathy
563 Yu, Kevin Hui, Kiran Vodrahalli, Klaus Greff, Linhai Qiu, Marcella Valentine, Marina Coelho,
564 Marvin Ritter, Matt Hoffman, Matthew Watson, Mayank Chaturvedi, Michael Moynihan, Min Ma,
565 Nabila Babar, Natasha Noy, Nathan Byrd, Nick Roy, Nikola Momchev, Nilay Chauhan, Noveen
566 Sachdeva, Oskar Bunyan, Pankil Botarda, Paul Caron, Paul Kishan Rubenstein, Phil Culliton,
567 Philipp Schmid, Pier Giuseppe Sessa, Pingmei Xu, Piotr Stanczyk, Pouya Tafti, Rakesh Shivan
568 Renjie Wu, Renke Pan, Reza Rokni, Rob Willoughby, Rohith Vallu, Ryan Mullins, Sammy Jerome,
569 Sara Smoot, Sertan Girgin, Shariq Iqbal, Shashir Reddy, Shruti Sheth, Siim Pöder, Sijal Bhatnagar,
570 Sindhu Raghuram Panyam, Sivan Eiger, Susan Zhang, Tianqi Liu, Trevor Yacovone, Tyler Liechty,
571 Uday Kalra, Utku Evci, Vedant Misra, Vincent Roseberry, Vlad Feinberg, Vlad Kolesnikov,
572 Woohyun Han, Woosuk Kwon, Xi Chen, Yinlam Chow, Yuvein Zhu, Zichuan Wei, Zoltan Egyed,
573 Victor Cotruta, Minh Giang, Phoebe Kirk, Anand Rao, Kat Black, Nabila Babar, Jessica Lo,
574 Erica Moreira, Luiz Gustavo Martins, Omar Sanseviero, Lucas Gonzalez, Zach Gleicher, Tris
575 Warkentin, Vahab Mirrokni, Evan Senter, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia
576 Hadsell, Yossi Matias, D. Sculley, Slav Petrov, Noah Fiedel, Noam Shazeer, Oriol Vinyals, Jeff
577 Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Jean-Baptiste
578 Alayrac, Rohan Anil, Dmitry, Lepikhin, Sebastian Borgeaud, Olivier Bachem, Armand Joulin,
579 Alek Andreev, Cassidy Hardin, Robert Dadashi, and Léonard Hussenot. Gemma 3 technical report,
580 2025. URL <https://arxiv.org/abs/2503.19786>.

581 Chaojie Wang, Yanchen Deng, Zhiyi Lyu, Liang Zeng, Jujie He, Shuicheng Yan, and Bo An. Q: Im-
582 proving multi-step reasoning for llms with deliberative planning. *arXiv preprint arXiv:2406.14283*,
583 June 2024.

584 Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang
585 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. *arXiv*
586 *preprint arXiv:2312.08935*, December 2023a. doi: 10.48550/arXiv.2312.08935. URL <https://arxiv.org/abs/2312.08935>.

587

588 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
589 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models,
590 2023b. URL <https://arxiv.org/abs/2203.11171>.

591 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le,
592 and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
593 URL <https://arxiv.org/abs/2201.11903>.

594 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu,
595 Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu,
596 Xingzhang Ren, and Zhenru Zhang. Qwen2.5-math technical report: Toward mathematical expert
597 model via self-improvement, 2024. URL <https://arxiv.org/abs/2409.12122>.

- 598 Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen Zhang, Runji Lin, Bowen Yu, Dayiheng Liu,
599 Jingren Zhou, and Junyang Lin. The lessons of developing process reward models in mathematical
600 reasoning, 2025. URL <https://arxiv.org/abs/2501.07301>.
- 601 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in
602 large language models, 2022. URL <https://arxiv.org/abs/2210.03493>.
- 603 Chujie Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu, Jin-
604 gren Zhou, and Junyang Lin. Processbench: Identifying process errors in mathematical reasoning,
605 2025. URL <https://arxiv.org/abs/2412.06559>.
- 606 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,
607 Claire Cui, Olivier Bousquet, Quoc Le, and Ed Chi. Least-to-most prompting enables complex
608 reasoning in large language models, 2023. URL <https://arxiv.org/abs/2205.10625>.