arXiv:2503.22832v2 [cs.PL] 10 Apr 2025

# *L*0-Reasoning Bench: Evaluating Procedural Correctness in Language Models via Simple Program Execution

**Simeng Sun, Cheng-Ping Hsieh, Faisal Ladhak, Erik Arakelyan,**
**Santiago Akle Serano**, **Boris Ginsburg**
NVIDIA
Santa Clara, CA 15213, USA
{simengs,bginsburg}@nvidia.com

## Abstract

Complex reasoning tasks often rely on the ability to consistently and accurately apply simple rules across incremental steps, a foundational capability which we term "level-0" reasoning. To systematically evaluate this capability, we introduce *L*0-Bench, a language model benchmark for testing *procedural* correctness – the ability to generate correct reasoning processes, complementing existing benchmarks that primarily focus on *outcome* correctness. Given synthetic Python functions with simple operations, *L*0-Bench grades models on their ability to generate step-by-step, error-free execution traces. The synthetic nature of *L*0-Bench enables systematic and scalable generation of test programs along various axes (e.g., number of trace steps). We evaluate a diverse array of recent closed-source and open-weight models on a baseline test set. All models exhibit degradation as the number of target trace steps increases, while larger models and reasoning-enhanced models better maintain correctness over multiple steps. Additionally, we use *L*0-Bench to explore test-time scaling along three dimensions: input context length, number of solutions for majority voting, and inference steps. Our results suggest substantial room to improve "level-0" reasoning and potential directions to build more reliable reasoning systems.

## 1 Introduction

Reasoning systems, by default, assume a fundamental prerequisite – the consistent and accurate execution of simple rules. Whether performing multi-digit arithmetic or following a legal procedure, the ability to reliably execute predefined rules is essential, as even a small error can derail the process from an expected trajectory, regardless of how sophisticated the system might be in other respects.

We consider the ability to *reliably execute simple rules across incremental steps* as the "level-0" reasoning ability, or a crucial form of sub-intelligence (Morris et al., 2023). This view also aligns with the recent discussion by Shalev-Shwartz et al. (2024) who formalize error-bounded precise reasoning in multi-step problem solving. However, modern neural language models are trained to approximate statistical patterns. *Can these probabilistic models consistently and reliably execute discrete rules with guaranteed correctness?* To answer this question, we introduce *L*0-Bench, a synthetic benchmark to explore the feasibility of a language model functioning as a *neural language computer*. Specifically, given a simple program, a language model is graded on its ability to generate an execution trace step by step infallibly.

Focusing on programming languages offers several advantages. With precisely defined programs, *L*0-Bench removes the inherent ambiguity of natural languages and isolates the reliable execution of simple rules from other reasoning processes (e.g., inferring general algorithms from observations). The callable functions also automatically produce execution traces, which allow for the evaluation of *procedural* correctness (i.e., maintaining correctness across reasoning processes), an aspect often overshadowed by the simpler evaluation of

*outcome* correctness. Moreover, programs are easy to generate and can be systematically controlled. In *L*0-Bench, we use generative grammar to construct simplified Python programs. By imposing multiple constraints on the grammar, we reduce the difficulty of executing each line of code, thereby avoiding unnecessary evaluation of operations specific to Python or arithmetic. The rule-based generation also enables flexible control of various axes (e.g., number of trace steps, program length) while being significantly more cost-effective than LM-based data curation. Furthermore, the synthetic nature mitigates concerns about data leakage/memorization and allows for scalable generation of test instances.

*L*0-Bench relates to prior studies on algorithmic tasks (Lee et al., 2024; Markeeva et al., 2024), such as large-number multiplication (Deng et al., 2024). However, *L*0-Bench differs in providing a general framework where the input program is not limited to a single predefined routine; instead, *L*0-Bench is able to admit any task that can be formulated as a program. While external tools naturally excel at executing discrete rules, we argue that the ability to perform simple operations step-by-step like a Turing Machine, without external tools, is still critical. It is essential not only for reasoning tasks that rely on basic deduction, but also for facilitating robustness, interpretability, and faithful alignment with human intents.

We evaluate 20 models with *L*0-Bench, including closed-source models, open-weight general-instruct models (ranging from 7B to 405B), and reasoning-enhanced models that generate "think" patterns, such as DeepSeek-R1. Experimental results demonstrate the advantages of larger model sizes and reasoning enhancements in achieving better procedural correctness. However, all models exhibit degradation as the target number of trace steps increases. We further leverage *L*0-Bench as a constrained environment to explore test-time scaling along three dimensions: (1) scaling input context by increasing step-by-step demonstrations, (2) scaling width by increasing the number of solutions for majority voting, and (3) scaling inference steps by enabling long chain-of-thought. While each dimension independently improves performance, our results reveal key limitations: scaling many-shot demonstrations yields diminishing returns, and sometimes degrades performance even for *long-context* models; majority voting plateaus as the number of voter grows; and long chain-of-thought reasoning increases latency without consistently benefiting all model sizes. Nonetheless, our scaling analysis demonstrates the benefits of scaling alone all three scaling dimensions and suggests ample room for improvement in current models.

Our contributions are as follows:

- We introduce *L*0-Bench, a synthetic benchmark for evaluating procedural correctness in neural language models via program execution. *L*0-Bench isolates step-by-step rule-execution from other reasoning processes (e.g., induction and search) and provides a controlled, verifiable evaluation framework.
- By evaluating 20 recent language models, we reveal their limitations in maintaining correctness across simple multi-step procedures, and show the benefits of larger model sizes and reasoning enhancements with long chain-of-thought.
- We further demonstrate the benefits of test-time scaling along three dimensions (input demonstrations, majority voting, and inference steps) and discuss the trade-off between performance and efficiency.

Overall, our findings reveal key limitations of current LLMs in long step-by-step rule-execution and suggest potential directions toward more reliable reasoning systems.[1]

## 2 L-0 Benchmark

### 2.1 Benchmark construction

Program execution can be viewed as a special form of deduction. During execution time, simple predefined rules are applied step-by-step, moving the execution process from one state to the next deterministically. For instance, `list.pop()` requires applying the *general rule* of "popping an item" to a list *instance*, such as `list = [1,2,3,4]`. Evaluating a model's
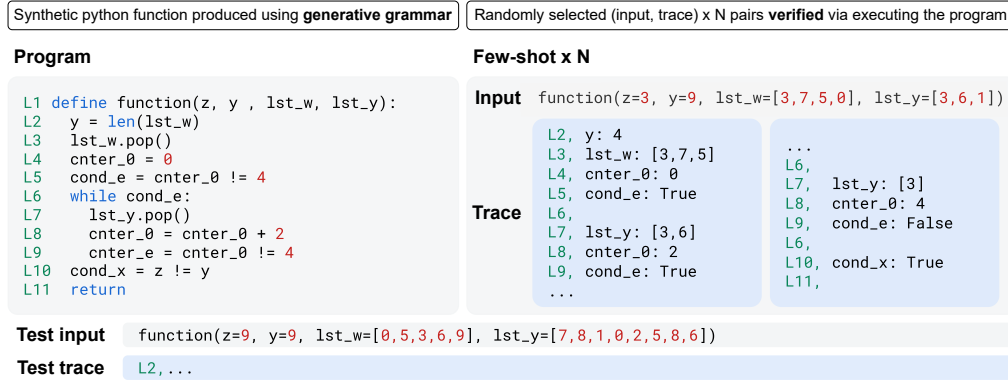
---

[1]We release code at https://github.com/SimengSun/L0-reasoning-bench.

| Synthetic python function produced using **generative grammar** | Randomly selected (input, trace) x N pairs **verified** via executing the program |
|---|---|

**Program**

```
L1  define function(z, y , lst_w, lst_y):
L2    y = len(lst_w)
L3    lst_w.pop()
L4    cnter_0 = 0
L5    cond_e = cnter_0 != 4
L6    while cond_e:
L7      lst_y.pop()
L8      cnter_0 = cnter_0 + 2
L9      cnter_e = cnter_0 != 4
L10   cond_x = z != y
L11   return
```

**Few-shot x N**

**Input** `function(z=3, y=9, lst_w=[3,7,5,0], lst_y=[3,6,1])`

**Trace**
```
L2, y: 4                    ...
L3, lst_w: [3,7,5]          L6,
L4, cnter_0: 0              L7,  lst_y: [3]
L5, cond_e: True            L8,  cnter_0: 4
L6,                         L9,  cond_e: False
L7, lst_y: [3,6]            L6,
L8, cnter_0: 2              L10, cond_x: True
L9, cond_e: True            L11,
...
```

**Test input**  `function(z=9, y=9, lst_w=[0,5,3,6,9], lst_y=[7,8,1,0,2,5,8,6])`

**Test trace**  `L2,...`

Figure 1: An illustrative example in *L*0-Bench. The prompt includes a program, multiple few-shot execution demonstrations, and a test input. Programs are generated using a constrained generative grammar that defines a simplified Python subset (see Appendix A). Inputs are randomly sampled and verified by executing the Python function to ensure error-free evaluation.

ability to execute simple programs therefore can serve as a test of its deductive capability – a necessary but not sufficient condition for complex reasoning processes. We introduce *L*0-Bench to test the deductive capability in neural language models via program execution, i.e., the ability to follow unambiguously defined routines step-by-step without deviation, akin to a neural language computer.

In *L*0-Bench, we construct synthetic Python programs using generative grammar. To avoid testing Python-specific operations or evaluations of complex expressions, we impose multiple constraints on the generative grammar, while preserving fundamental elements of conventional programming languages: (a) simple arithmetic and comparison operations, (b) control (conditional branch and loops), and (c) memory access and operations. The grammar hence covers diverse procedures that require modifying execution states (e.g., append an item to a list) and strictly adhering to the control flow (e.g., exit loops on termination), which are absent in simpler rule-based tasks, such as variable tracking (Hsieh et al., 2024).

**Program generation.** To reduce the difficulty of executing each line of code, we impose multiple constraints on the production rules used to generate the test programs. A detailed grammar description is provided in Appendix A, with main constraints summarized below:
- The generated program is a standalone Python function.
- The variables are one of the following: integer, list of integers, and boolean value.
- We restrict basic binary operations to `+`, `-`, `==`, and `!=`, as models struggle to reliably execute other operations such as `*`, `%`, and `<`.
- Complex single-line expressions are disabled to ensure the simplicity of transition between lines. For instance, we do not allow evaluation of long expressions such as `(x + lst[5] - y - 6)`.
- Context-sensitive rules are added to enforce terminable while loops.

**Verified inputs and traces.** We randomly generate values for the function input arguments and skip inputs that encounter execution errors (e.g., pop an item from an empty list). For each program, we generate and verify multiple sets of (input, trace) pairs, which can serve as demonstrations for the expected trace format. Each trace step consists of a line number, a variable name, and its updated value, if applicable. If no variable value is updated, e.g., the line of `if condition:` or `while condition:`, only the line number is required.

**Example format.** We provide $N$ few-shot demonstrations specific to the test program for guiding the generation. Formally, let $\mathcal{F}$ denote the string of the provided program, $x_i$ be the string of a demonstration input, and $y_i$ be its corresponding output trace consisting of multiple trace lines. Let $x^*$ and $y^*$ represent the test input and expected output trace, respectively. The task input is structured as $\{\mathcal{F}, x_1, y_1, x_2, y_2, \ldots, x^*\}$, with the expected output being $\{y^*\}$. The prompt template and example are provided in Appendix D. An illustrative example is shown in Figure 1.

## 2.2 Evaluation Metrics

We evaluate whether a model can generate the correct execution trace, both from beginning to end (whole-trace accuracy) and partially (steps to the first error):

- **Whole-trace accuracy**: Since all traces begin with a fixed format, we align the model response with the start of the ground-truth trace, and then evaluate whether the entire generated sequence exactly matches the expected ground-truth trace.
- **Steps to the first error**: We also evaluate the partial correctness by counting the number of correct steps before the first error appears. We do not evaluate correctness beyond the first error, as the model may enter an entirely incorrect branch or prematurely exit loops, making subsequent steps incorrect.

The two metrics are not strictly correlated: a model with a low **whole-trace accuracy** can still have a high **steps to the first error** if errors occur primarily at the end of the trace.

## 3 Experimental Setup

**Data.** $L0$-Bench is designed to flexibly generate data based on specified configurations. We construct a base set of $L0$-Bench by enforcing multiple constraints described in Appendix A. These constraints can be adjusted to increase difficulty, such as allowing for longer expressions at each line, increasing scope depth, etc. We group generated examples into four bins of *short*, *medium*, *long*, and *extra long* traces. Each bin contains 500 examples with the average number of steps $\{13, 80, 164, 246\}$ respectively.

**Prompt Format & Majority Voting Setup.** We include four (input, trace) demonstrations specific to the test program in the context. Both evaluation metrics are reported for the majority-voted responses (*majvote@31*) across 31 parallel samplers as well as the accuracy of having at least one whole correct trace (*pass@31*). The majority-voted trace is the most frequent whole trace after aligning the voters' responses based on the fixed ground-truth beginning pattern ("*L2,*"). The "*Single Attempt*" results are the average of 31 independent runs, where each voter randomly draws four distinct demonstrations from a pool of 64 exempla pre-generated for each test program.

**Models.** We evaluate 20 models, covering both closed-source and open-weight models ranging from 7B to 405B, including general instruct models as well as reasoning-enhanced models such as DeepSeek-R1 and R1-Distilled Qwen models. For all models, we use greedy decoding, except for R1 and R1-distilled models, which are sampled with the recommended temperature of 0.6. For models that generate long Chain-of-Thought outputs, we cap the maximum number of generated tokens at 20 times the length of the ground-truth trace. To disable the long reasoning pattern among these models, we explicitly append an answer prefix to the chat template to bypass the generation of <think> token.

## 4 Main Results

$L0$-**Bench & general-instruct models.** One potential concern of $L0$-Bench is that the focus on programming language may pose challenges for general-instruct models that lack code-specific training. Results in Table 1 reveal that the general-instruct model Qwen2.5-32B-Instruct in fact performs better than its same-size Qwen2.5-Coder model, suggesting that $L0$-Bench can be a viable test for general-instruct models.

**Benefits of larger model sizes and "think" patterns.** The column *thinking mode on* is set to "Y" in Table 1 for runs that produce long chain-of-thought "think" patterns. Among "non-thinking" systems ranging in $\{7B, 14B, 32B, 72B, 405B\}$, we observe the general trend: larger models achieve higher accuracy in both single-attempt evaluations and majority-voted whole-trace accuracy, indicating the benefits of scaling base model parameters. While model size generally correlates with better performance, models that generate intermediate chain-of-thought steps can outperform much larger "non-thinking" counterparts. For instance, DeepSeek-R1-Distilled-Qwen-14B outperforms the original 72B Qwen2.5 as well as the 29 times larger Llama. However, the improved performance comes at the cost of longer inference time. We discuss the performance-efficiency tradeoff in a later section.

| Model Name | Thinking mode on | Steps to Err. ↑ (target # steps = 125.8) | | Trace Acc. (%) | | |
|---|---|---|---|---|---|---|
| | | Single Attempt | majvote @ 31 | Single Attempt | majvote @ 31 | pass @ 31 |
| **Closed-source models** | | | | | | |
| o1 | Y | 122.2 | - | 92.0 | - | - |
| o3-mini | Y | 92.8 | - | 60.6 | - | - |
| gpt-4o-2024-11-20 | N | 112.5 | - | 76.4 | - | - |
| claude-3-7-sonnet | Y | 118.6 | - | 95.1 | - | - |
| claude-3-7-sonnet | N | 114.0 | - | 86.1 | - | - |
| claude-3-5-sonnet-20241022 | N | 108.5 | - | 83.3 | - | - |
| **Open-weight models** | | | | | | |
| QwQ-32B | Y | 120.0 | 124.9 | 86.6 | 96.1 | 99.8 |
| Deepseek-R1 | Y | 115.7 | 121.4 | 91.5 | 97.4 | 97.9 |
| DeepSeek-R1-Distill-Llama-70B | Y | 104.4 | 116.3 | 62.0 | 79.6 | 96.6 |
| Meta-Llama-3.1-405B-Instruct | N | 99.1 | 106.2 | 63.0 | 73.1 | 93.2 |
| Meta-Llama-3.1-70B-Instruct | N | 82.8 | 93.9 | 42.4 | 54.7 | 85.9 |
| Meta-Llama-3.1-8B-Instruct | N | 33.2 | 40.7 | 6.5 | 11.5 | 33.9 |
| DeepSeek-R1-Distill-Qwen-32B | Y | 114.4 | 123.4 | 80.4 | 94.9 | 99.8 |
| DeepSeek-R1-Distill-Qwen-14B | Y | 99.3 | 118.3 | 56.8 | 83.5 | 97.8 |
| DeepSeek-R1-Distill-Qwen-7B[2] | Y | 8.7 | 11.7 | 3.8 | 11.2 | 28.8 |
| Qwen2.5-72B-Instruct | N | 94.3 | 102.6 | 55.2 | 66.9 | 89.9 |
| Qwen2.5-32B-Instruct | N | 88.3 | 101.8 | 50.4 | 66.0 | 90.0 |
| Qwen2.5-Coder-32B-Instruct | N | 81.4 | 92.4 | 44.0 | 56.8 | 84.4 |
| Qwen2.5-14B-Instruct | N | 60.9 | 73.4 | 23.7 | 35.7 | 73.7 |
| Qwen2.5-Coder-7B-Instruct | N | 52.3 | 68.4 | 17.1 | 32.4 | 62.5 |
| Qwen2.5-7B-Instruct | N | 41.2 | 54.2 | 11.3 | 22.8 | 44.1 |

Table 1: We evaluate 20 recent models on $L$0-Bench. Results are averaged over four data splits (short, medium, long, extra-long traces), with *Single Attempt* results based on the average of 31 independent runs. Larger models, reasoning with longer chain-of-thought, and majority voting over parallel runs improve performance on $L$0-Bench. The large gap between *pass@k* and *single attempts* across all models indicates significant room for improvement. Details about tested models are provided in Table 11.

**Models fail at longer traces more often.** While programs in $L$0-Bench are simplified to allow for only basic operations, none of the models, including reasoning-enhanced models, is able to consistently execute the basic operations as the number of ground-truth trace steps increases (Figure 2). For instance, while QwQ-32B achieves over 95% whole-trace accuracy on the shortest split, its performance drops to 82% on the longest split. This limitation persists despite the constrained set of simple operations and the total number of tokens in the prompt being far below the claimed effective context size.

**Majority voting boosts performance, yet falls far behind Pass@k.** Consistent with other works, majority voting over multiple solutions brings significant gains for all data splits. The gains are even more pronounced for *pass@31*, a soft upper-bound that models are yet to reach. Notably, reasoning-enhanced models, such as QwQ-32B, demonstrate much smaller gaps between majority-voted and *pass@k* accuracy, however, the gaps are still non-marginal, suggesting promising directions to improve the latest models.

**Steps To Error provides a more fine-grained evaluation.** Generally, models achieve high whole-trace accuracy when their average *Steps to Error* is also high. However, we observe discrepancies between the whole-trace and partial accuracy for certain models. For instance, despite achieving an average ∼112.5 steps to the first error, gpt-4o scores below 80 in whole-trace accuracy. This aligns with manual examination which reveals that the model often fails towards the end of the trace by skipping steps.

---

[2]Model often fails to exit the "thought" process with naive budget forcing (Muennighoff et al., 2025) (maximum generation length is set to 20 times the target length), thus performing worse overall.
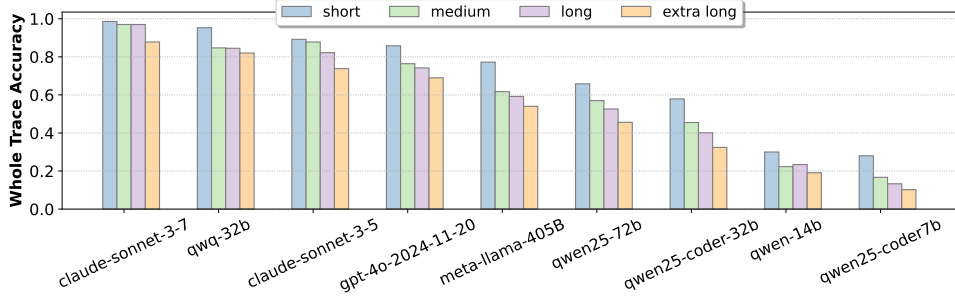
Figure 2: Model performance generally declines as the number of trace steps increases. Open-weight models in the figure are instruct version.

| | # Target Tokens | 2-shot | 4-shot | 8-shot | 16-shot | 32-shot | 64-shot |
|---|---|---|---|---|---|---|---|
| **# Input tokens** | | | | | | | |
| Short trace | 0.1K | 1K | 1.5K | 2.7K | 5K | 9.5K | 19K |
| Medium trace | 0.8K | 2.3K | 4.1K | 7.8K | 15K | 30K | 59K |
| Long trace | 1.9K | 4.4K | 8.4K | 16K | 32K | 64K | 129K |
| Extra long trace | 3.1K | 7.0K | 13K | 27K | 53K | 100K | 208K |
| **Pass@1 Trace Acc. (%)** | | | | | | | |
| Short trace | 0.1K | $49.1_{\pm 2.0}$ | $58.2_{\pm 1.5}$ | $66.7_{\pm 1.4}$ | $71.1_{\pm 1.4}$ | $75.5_{\pm 1.3}$ | $75.3_{\pm 1.1}$ |
| Medium trace | 0.8K | $32.4_{\pm 1.2}$ | $40.4_{\pm 1.6}$ | $48.0_{\pm 1.5}$ | $53.3_{\pm 1.5}$ | $55.3_{\pm 1.1}$ | $49.2_{\pm 1.4}$ |
| Long trace | 1.9K | $28.1_{\pm 1.2}$ | $36.3_{\pm 1.5}$ | $43.1_{\pm 1.8}$ | $46.6_{\pm 1.4}$ | $39.9_{\pm 1.1}$ | $23.5_{\pm 0.9}$ |
| Extra long trace | 3.1K | $28.0_{\pm 1.4}$ | $34.6_{\pm 1.6}$ | $38.6_{\pm 1.4}$ | $36.9_{\pm 1.7}$ | $23.4_{\pm 1.1}$ | $1.5_{\pm 0.4}$ |

Table 2: We scale up to 64 few-shot demonstrations specific to the input test programs for `Meta-Llama-3.1-70B-Instruct`, adopting a many-shot setting where all demonstrations are directly relevant, unlike many long-context benchmarks that emphasize extracting signals from distractors. While including more step-by-step demonstrations in the context improves performance for short traces, performance degrades for longer traces even before reaching the 128K context limit. The results evaluated on context longer than 128K tokens are in gray.

# 5 Scaling Test-Time Compute

We employ $L0$-Bench as a controllable environment to systematically explore three dimensions of test-time scaling: (1) **scaling input context** by increasing the number of in-context demonstrations, (2) **scaling width** by increasing the number of solutions for majority voting, and (3) **scaling inference steps** by enabling long chain-of-thought.

**Scaling input context.** Learning from in-context demonstrations can be viewed as "optimization" through forward passes (Von Oswald et al., 2023), or inference-time training. We scale the number of in-context demonstrations up to 64 for the 70B Llama model. Table 2 shows that performance consistently improves as the number of demonstrations increases from 2 to 32 except for the extra-long split. Scaling beyond 32-shot brings little gains or even hurts performance, despite the input length being much shorter than the supported context of 128K tokens. This is not ideal as all demonstrations provide step-by-step execution of the same test program, thus immediately relevant to producing the target trace. This many-shot setting of $L0$-Bench can therefore serve as an alternative to existing long-context benchmarks that often emphasize on the capability to retrieve relevant signals from irrelevant noises.

**Scaling width.** To explore the potential of majority voting at scale, we increase the number of solutions to 1K for two models from the Qwen series under two settings in Figure 3: (1) each voter permutes the same fixed set of four-shot demonstrations of the test program, and (2) each voter independently selects four distinct demonstrations from a pool of 64 w.r.t the test program, same as the main results. In both cases, whole-trace accuracy improves with the number of solutions, with a diminishing return beyond 32~64. Notably, drawing demonstrations from a diverse pool yields larger gains, with the small 7B model showing
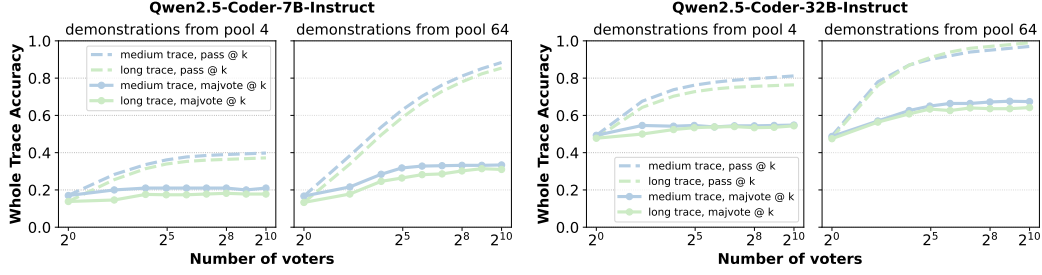
Figure 3: We increase the number of voters up to 1024 to study majority voting at scale. Each voter either permutes the same set of four demonstrations or randomly selects four distinct few-shot demonstrations from a pool of 64. Access to a larger, more diverse demonstration pool improves majority-voted performance and pass@K. The significant gap between pass@K and majority-voted performance indicates large room for models to reach their full potential.

promising *pass@k* with 1K solutions. These results imply the presence of correct traces in models' hypothesis space, yet models struggle to reliably and effectively leverage in-context demonstrations with low sample complexity during inference-time training.

**Scaling inference steps.** Recent works suggest that scaling scratchpad length helps overcome architectural limitations and empowers Transformer to solve tasks previously unsolvable with limited depth (Li et al., 2024b). Our results confirm the benefit of generating a long chain-of-thought sequence prior to providing the final answer on *L0-Bench*. Specifically, results in Table 3 reveal large performance drop when the "think" patterns in reasoning-enhanced models are bypassed.[3] For instance, disabling "thinking" in DeepSeek-R1-Distill-Qwen-32B causes a drop in whole-trace accuracy of $\sim$35 points, while also harming the soft upper-bound *pass@k* performance. This indicates the benefits of step-by-step reasoning and self-correction patterns for maintaining correctness over multiple steps, as shown in an example QwQ-32B response in Appendix I.

**Scaling along all dimensions & limitations.** Combining the above three test-time scaling dimensions often leads to additional gains. For instance, the reasoning-enhanced 70B Llama significantly outperforms the original one without "think" patterns. After providing more demonstrations and aggregating over multiple voters (Table 4), we observe significant gains in whole-trace accuracy (**65→88**). However, scaling along each dimension incurs different costs. The long-CoT model can be slow at inference,[4] whereas taking majority vote over solutions generated *in parallel* avoids high latency, only when sufficient compute resources are available. Additionally, scaling the number of demonstrations can be promising for long-context models, however, access to a large pool of demonstrations may not always be available. We leave a systematic investigation of these trade-offs for future work.

## 6 Ablation & Analysis

**Task variation # 1: Alternative Programs as Demonstrations** The default setup in *L0-Bench* uses (input, trace) pairs of the same test program as demonstrations. This setup is much easier than a variation setup where demonstrations are other programs, and models no longer have access to the step-by-step demonstrations on how to execute the test algorithm. That is, the prompts are structured as $\{\mathcal{F}_1, x_1, y_1, \mathcal{F}_2, x_2, y_2, \ldots, \mathcal{F}_*, x^*\}$, following the notation used in § 2. Figure 4 shows performance degradation for both Qwen2.5-32B-instruct and the performant closed-source model claude-sonnet-3.7. We chose the easier setup on purpose to avoid distractors in the prompt in *L0-Bench*.

---

[3]We modify the chat template <|Assistant|> The execution trace is:\n to bypass the <think>.

[4]For example, DeepSeek-R1-Distill-Qwen-32B takes 29 min to complete the *short* split, whereas Qwen2.5-32B-Inst takes around 3 min, both running on 4 A100 GPUs with the same serving setup.

| Model Name | Thinking mode on | Trace Acc. (%) | | |
|---|---|---|---|---|
| | | Single Attempt | maj vote @ 31 | pass @ 31 |
| QwQ-32B | Y | 86.6 | 96.1 | 99.8 |
| QwQ-32B | N | 23.3 | 64.0 | 85.4 |
| R1-Distill-Qwen-32B | Y | 80.4 | 94.9 | 99.8 |
| R1-Distill-Qwen-32B | N | 45.3 | 60.1 | 86.6 |
| Qwen2.5-32B-Inst | N | 50.4 | 66.0 | 90.0 |

| # few shot | Trace Acc. (%) | | |
|---|---|---|---|
| | Single Attempt | maj vote @ 31 | pass @ 31 |
| 2 | **65.5** | 83.8 | 98.2 |
| 4 | 72.9 | 87.6 | 98.6 |
| 8 | 75.1 | **88.2** | 99.2 |
| 16 | 74.1 | 86.6 | 98.6 |
| 32 | 69.5 | 81.6 | 97.4 |

Table 3: Disabling long chain-of-thought in reasoning-enhanced models leads to large performance drop.

Table 4: R1-Distill-Llama-70B on the *medium* split.



Figure 4: **(Left)** The difficulty of *default* tasks in *L*0-Bench can be increased by : (*Task variation #1*) using irrelevant (program, input, trace) as few-shot examples, or (*Task variation #2*) removing the program, forcing the model to infer the algorithm purely from step-by-step demonstrations. We use the simplest setup to test the 0-level reasoning. **(Right)** Even the strongest closed-source model Claude-3.7-Sonnet shows performance degradation with the modification of task setup, especially when extended thinking is disabled.

**Task variation # 2: Removing Input Program – a Transduction Task.** We test whether in-context step-by-step demonstrations alone suffice for generating execution traces by removing the input program from the prompt, a transduction setup similar to ARC-AGI (Chollet, 2019; Chollet et al., 2024). Figure 4 (left) indicates that the 32B Qwen model achieves non-zero accuracy even without the input program, implying occasionally correct pattern inference from the provided (input, trace) pairs. We focus on the simpler setting in this work, leaving the much harder transduction task for future exploration.

**Larger models tend to spend less time "thinking."** We count the number tokens enclosed by `<think>` and `</think>` in five reasoning-enhanced models. The larger 70B generates on average 1365 thought tokens, significantly fewer tokens than the smaller Qwen models, which generate on average {3344, 2241, 1722} for sizes of {7B, 14B, 32B}, respectively. Despite the relationship observed in R1-distilled models, `QwQ-32B` is more verbose than the other models, generating over 6K thought tokens on average.

**Common Failure Modes** Manual inspection of errors reveals three common failure modes. (1) **Incorrect item counting**: Even closed-source models struggle to reliably copy a list from a previous step to the next during operations like pop or append. Often, items are incorrectly repeated in the list, or completely deleted even when the list is not long. (2) **Mis-evaluation of simple expression**: Despite the simplified binary operations in *L*0-Bench, small models still fail to reliably evaluate simple expressions, such as `10 == 10`, thus leading to incorrect branching decisions for *if* blocks sometimes. (3) **Skipped steps**: models sometimes incorrectly skip trace steps, as often observed in `gpt-4o` and `o3-mini`. These models often skip the last evaluation of the line `while condition:` before exiting the loop, without following the explicit step-by-step demonstrations in the context.

## 7 Related Works

Following the framework introduced by Morris et al. (2023), $L$0-Bench evaluates the ability to execute unambiguously defined simple routines step-by-step without deviation, a "compiler"-like "level-0" intelligence that complex reasoning processes rely upon. The ability to maintain correctness by reliably executing simple rules also echoes the precise reasoning framework proposed by Shalev-Shwartz et al. (2024), who formalize the compounding errors in multi-step problem solving. While $L$0-Bench focuses on the coding domain, the task setup is distinct from existing works on code generation or completion (Zhuo et al., 2024; Jain et al., 2025). Instead, models are graded on their ability to generate concrete, often long, execution traces, which are proven to benefit code reasoning tasks (Ni et al., 2024a).

**Teaching models to execute.**  Prior works often explicitly *teach* models to generate execution trace (Liu et al., 2023). For instance, Zaremba & Sutskever (2014) develop curriculum training methods to teach models to execute short programs, such as addition and memorization. Other works (Yan et al., 2020; Veličković et al., 2020) train neural LMs to execute subroutines, such as sorting or classical graph algorithms, by developing novel architectures or training strategies. More recently, Markeeva et al. (2024) train LMs to produce traces for a set of algorithms from the textbook *Introduction to Algorithms*, and demonstrate significantly better performance than general-purpose pre-trained large language models. Recent works also elicit the execution process during test time by prompting large pre-trained language models. The reasoning process over a scratchpad (Nye et al., 2021), however, is often a mixture of natural and formal languages (Chae et al., 2024; Arakelyan et al., 2024; Zhou et al., 2022), which inhibit proper evaluation of procedural correctness.

**Evaluation of procedural correctness.**  Existing reasoning benchmarks, such as MATH (Hendrycks et al., 2021), evaluate exclusively the correctness of *outcomes* (e.g., the solution inside $\boxed{*}$). We argue that a "reasoning" model should produce correct reasoning processes, which necessarily lead to correct solutions. Program execution is a natural subject to study given its clarity, verifiability, and constrained rule-based nature. Prior works that challenge models with code execution tasks (Tufano et al., 2023; Gu et al., 2024; La Malfa et al., 2024) are often limited by the data curation process or the fixed pre-defined set of algorithms. Recent works also use natural language to define problems with unique and deterministic reference trajectory. Gui et al. (2024) define rule-based games or puzzles (Lin et al., 2025) with varying difficulty level to evaluate both the ability to plan and execute rules; Fujisawa et al. (2024) and Ye et al. (2025) both use natural languages to describe constraints and expected routines in the context, covering synthetic tasks (e.g., toy string manipulation) and tasks closer to real scenarios (e.g., travel planning). In contrast with these works on following natural language routines (Chen et al., 2024b), $L$0-Bench provides means to scalable generation of program-oriented instances, a more general framework admitting any problems that can be framed as programs. Additional related works and discussion on future directions are provided in Appendix C.

## 8 Conclusion

We presented $L$0-Bench, a synthetic language model benchmark for testing a special form of deduction, or the ability to reliably execute simple rules throughout predefined multi-step procedures. Specifically, given synthetic programs, $L$0-Bench scores models based on their ability to generate accurate, step-by-step execution traces without error. The program-oriented synthetic benchmark enables scalable generation of test instances and allows for precise evaluation of procedural correctness, an aspect often overlooked by benchmarks focusing on outcome correctness. We evaluated 20 recent language models on a base evaluation set of $L$0-Bench. Results revealed key limitations in maintaining correctness during simple multi-step procedures, even among the latest reasoning-enhanced models. Additionally, we employed $L$0-Bench as a constrained environment to explore test-time scaling and other task variations, such as a more difficult transduction setting. Our findings suggest several directions for future research, including systematic investigation of various test-time scaling axes, extending $L$0-Bench to evaluate reasoning ability beyond "level-0," and leveraging scalable synthetic data generation to build more reliable reasoning systems.

## Acknowledgments

## References

Anthropic. Claude 3.5 Sonnet, 2024. URL https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf.

Anthropic. Claude 3.7 Sonnet, 2025. URL https://assets.anthropic.com/m/785e231869ea8b3b/original/claude-3-7-sonnet-system-card.pdf.

Erik Arakelyan, Pasquale Minervini, Pat Verga, Patrick Lewis, and Isabelle Augenstein. Flare: Faithful Logic-Aided Reasoning and Exploration. *arXiv:2410.11900*, 2024.

Hyungjoo Chae, Yeonghyeon Kim, Seungone Kim, Kai Tzu-iunn Ong, Beong-woo Kwak, Moohyeon Kim, Sunghwan Kim, Taeyoon Kwon, Jiwan Chung, Youngjae Yu, and Jinyoung Yeo. Language Models as Compilers: Simulating Pseudocode Execution Improves Algorithmic Reasoning in Language Models. In *EMNLP*, 2024.

Junkai Chen, Zhiyuan Pan, Xing Hu, Zhenhao Li, Ge Li, and Xin Xia. Reasoning runtime behavior of a program with llm: How far are we? *arXiv:2403.16437*, 2024a.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of Thoughts Prompting: Disentangling Computation from Reasoning for Numerical Reasoning Tasks. *TMLR*, 2023.

Xinyi Chen, Baohao Liao, Jirui Qi, Panagiotis Eustratiadis, Christof Monz, Arianna Bisazza, and Maarten de Rijke. The SIFo Benchmark: Investigating the Sequential Instruction Following Ability of Large Language Models. In *Findings of EMNLP*, 2024b.

Francois Chollet, Mike Knoop, Gregory Kamradt, and Bryan Landers. Arc prize 2024: Technical report. *arXiv:2412.04604*, 2024.

François Chollet. On the Measure of Intelligence. *arXiv:1911.01547*, 2019.

Yuntian Deng, Yejin Choi, and Stuart Shieber. From explicit cot to implicit cot: Learning to internalize cot step by step. *arXiv:2405.14838*, 2024.

Ippei Fujisawa, Sensho Nobe, Hiroki Seto, Rina Onda, Yoshiaki Uchida, Hiroki Ikoma, Pei-Chun Chien, and Ryota Kanai. Procbench: Benchmark for Multi-step Reasoning and Following Procedure. *arXiv:2410.03117*, 2024.

Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. PAL: Program-aided Language Models. In *ICML*, 2023.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, and Others. The llama 3 herd of models, 2024.

Alex Gu, Baptiste Roziere, Hugh James Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida Wang. CRUXEval: A benchmark for code reasoning, understanding and execution. In *ICML*, 2024.

Jiayi Gui, Yiming Liu, Jiale Cheng, Xiaotao Gu, Xiao Liu, Hongning Wang, Yuxiao Dong, Jie Tang, and Minlie Huang. Logicgame: Benchmarking Rule-based Reasoning Abilities of Large Language Models. *arXiv:2408.15778*, 2024.

Daya Guo, Dejian Yang, Haowei Zhang, et al. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *arXiv:2501.12948*, 2025.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv:2103.03874*, 2021.

Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, and Boris Ginsburg. RULER: What's the Real Context Size of Your Long-Context Language Models? In *CoLM*, 2024.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. GPT-4o system card. *arXiv:2410.21276*, 2024.

Aaron Jaech, Adam Kalai, Adam Lerer, et al. OpenAI o1 system card. *arXiv:2412.16720*, 2024.

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and Contamination Free Evaluation of Large Language Models for Code. In *ICLR*, 2025.

Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can Language Models Resolve Real-world Github Issues? In *ICLR*, 2024.

Emanuele La Malfa, Christoph Weinhuber, Orazio Torre, Fangru Lin, Samuele Marro, Anthony Cohn, Nigel Shadbolt, and Michael Wooldridge. Code simulation challenges for large language models. *arXiv:2401.09074*, 2024.

Nayoung Lee, Kartik Sreenivasan, Jason D. Lee, Kangwook Lee, and Dimitris Papailiopoulos. Teaching Arithmetic to Small Transformers. In *ICLR*, 2024.

Chengshu Li, Jacky Liang, Andy Zeng, Xinyun Chen, Karol Hausman, Dorsa Sadigh, Sergey Levine, Li Fei-Fei, Fei Xia, and Brian Ichter. Chain of Code: Reasoning with a Language Model-Augmented Code Emulator. 2024a.

Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of Thought Empowers Transformers to Solve Inherently Serial Problems. In *ICLR*, 2024b.

Bill Yuchen Lin, Ronan Le Bras, Kyle Richardson, Ashish Sabharwal, Radha Poovendran, Peter Clark, and Yejin Choi. Zebralogic: On the Scaling Limits of LLMs for Logical Reasoning. *arXiv:2502.01100*, 2025.

Chenxiao Liu, Shuai Lu, Weizhu Chen, Daxin Jiang, Alexey Svyatkovskiy, Shengyu Fu, Neel Sundaresan, and Nan Duan. Code Execution with Pre-trained Language Models. In *Findings of ACL*, 2023.

Tianyang Liu, Canwen Xu, and Julian McAuley. Repobench: Benchmarking repository-level code auto-completion systems. In *ICLR*, 2024.

Larisa Markeeva, Sean McLeish, Borja Ibarz, Wilfried Bounsi, Olga Kozlova, Alex Vitvitskyi, Charles Blundell, Tom Goldstein, Avi Schwarzschild, and Petar Veličković. The CLRS-Text Algorithmic Reasoning Language Benchmark. *arXiv:2406.04229*, 2024.

Seyed Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. GSM-symbolic: Understanding the limitations of mathematical reasoning in large language models. In *ICLR*, 2025.

Meredith Ringel Morris, Jascha Sohl-dickstein, Noah Fiedel, Tris Warkentin, Allan Dafoe, Aleksandra Faust, Clement Farabet, and Shane Legg. Levels of AGI for Operationalizing Progress on the Path to AGI. *arXiv:2311.02462*, 2023.

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. S1: Simple test-time scaling. *arXiv:2501.19393*, 2025.

Ansong Ni, Miltiadis Allamanis, Arman Cohan, Yinlin Deng, Kensen Shi, Charles Sutton, and Pengcheng Yin. NExt: Teaching Large Language Models to Reason about Code Execution. *arXiv:2404.14662*, 2024a.

Ansong Ni, Pengcheng Yin, Yilun Zhao, Martin Riddell, Troy Feng, Rui Shen, Stephen Yin, Ye Liu, Semih Yavuz, Caiming Xiong, Shafiq Joty, Yingbo Zhou, Dragomir Radev, Arman Cohan, and Arman Cohan. L2CEval: Evaluating Language-to-Code Generation Capabilities of Large Language Models. *TACL*, 2024b.

Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work: Scratchpads for intermediate computation with language models. 2021.

OpenAI. OpenAI o3-mini System Card, 2025. URL https://cdn.openai.com/o3-mini-system-card-feb10.pdf.

Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. Logic-LM: Empowering Large Language Models with Symbolic Solvers for Faithful Logical Reasoning. In *Findings of EMNLP*, December 2023.

Qwen. QwQ-32b: Embracing the Power of Reinforcement Learning, 2025. URL https://qwenlm.github.io/blog/qwq-32b/.

Abulhair Saparov and He He. Language Models are Greedy Reasoners: a Systematic Formal Analysis of Chain-of-Thought. In *ICLR*, 2023.

Abulhair Saparov, Richard Yuanzhe Pang, Vishakh Padmakumar, Nitish Joshi, Mehran Kazemi, Najoung Kim, and He He. Testing the General Deductive Reasoning Capacity of Large Language Models using OOD Examples. In *NeurIPS*, 2023.

Shai Shalev-Shwartz, Amnon Shashua, Gal Beniamini, Yoav Levine, Or Sharir, Noam Wies, Ido Ben-Shaul, Tomer Nussbaum, and Shir Granot Peled. Artificial Expert Intelligence through PAC-reasoning. *arXiv:2412.02441*, 2024.

Michele Tufano, Shubham Chandel, Anisha Agarwal, Neel Sundaresan, and Colin Clement. Predicting code coverage without execution. *arXiv:2307.13383*, 2023.

Petar Veličković, Rex Ying, Matilde Padovano, Raia Hadsell, and Charles Blundell. Neural Execution of Graph Algorithms. In *ICLR*, 2020.

Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *ICML*, 2023.

Weixiang Yan, Haitian Liu, Yunkun Wang, Yunzhe Li, Qian Chen, Wen Wang, Tingyu Lin, Weishan Zhao, Li Zhu, Hari Sundaram, and Shuiguang Deng. CodeScope: An Execution-based Multilingual Multitask Multidimensional Benchmark for Evaluating LLMs on Code Understanding and Generation. In *ACL*, 2024.

Yujun Yan, Kevin Swersky, Danai Koutra, Parthasarathy Ranganathan, and Milad Hashemi. Neural Execution Engines: Learning to Execute Subroutines. In *NeurIPS*, 2020.

An Yang, Baosong Yang, Beichen Zhang, and Binyuan Hui et al. Qwen2.5 Technical Report. *arXiv:2412.15115*, 2024.

Xi Ye, Fangcong Yin, Yinghui He, Joie Zhang, Howard Yen, Tianyu Gao, Greg Durrett, and Danqi Chen. LongProc: Benchmarking Long-context Language Models on Long Procedural Generation. *arXiv:2501.05414*, 2025.

Wojciech Zaremba and Ilya Sutskever. Learning to execute. *arXiv:1410.4615*, 2014.

Hattie Zhou, Azade Nova, Hugo Larochelle, Aaron Courville, Behnam Neyshabur, and Hanie Sedghi. Teaching Algorithmic Reasoning via In-context Learning. *arXiv:2211.09066*, 2022.

Yang Zhou, Hongyi Liu, Zhuoming Chen, Yuandong Tian, and Beidi Chen. Gsm-Infinite: How Do Your LLMs Behave over Infinitely Increasing Context Length and Reasoning Complexity? *arXiv:2502.05252*, 2025.

Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. *arXiv:2406.15877*, 2024.

# A   Grammar for generating synthetic programs

---

**Baby Python Grammar Production Rules**

⟨*program*⟩ ::= ⟨*stmt_lst*⟩
⟨*stmt_lst*⟩ ::= ⟨*stmt*⟩ | ⟨*stmt*⟩⟨*stmt_lst*⟩
⟨*stmt*⟩ ::= ⟨*assignment*⟩ | ⟨*if_block*⟩ | ⟨*while_block*⟩ | ⟨*list_op*⟩

**# assignment**
⟨*assignment*⟩ ::= ⟨*num_assignment*⟩ | ⟨*bool_assignment*⟩
⟨*num_assignment*⟩ ::= ⟨*var*⟩ = ⟨*expr*⟩
⟨*bool_assignment*⟩ ::= ⟨*bool_var*⟩ = ⟨*bool_expr*⟩

**# if block**
⟨*if_block*⟩ ::= if ⟨*bool_var*⟩ : ⟨*stmt_lst*⟩

**# while block**
⟨*init_cnter*⟩ ::= cnter = 0
⟨*increment_cnter*⟩ ::= cnter = cnter + ⟨*cnter_increment_number*⟩
⟨*while_cond*⟩ ::= ⟨*while_cond_var*⟩ = cnter !=⟨*while_cond_number*⟩
⟨*while_block*⟩ ::= ⟨*while_block_nh*⟩ | ⟨*init_cnter*⟩⟨*while_cond*⟩ while
⟨*while_cond_var*, ∗⟩ : ⟨*stmt_lst*⟩⟨*increment_cnter*⟩⟨*while_cond*, ∗⟩
⟨*while_block_nh*⟩ ::=  while True : ⟨*stmt_lst*⟩

**# list op**
⟨*list_op*⟩ ::= ⟨*list_var*⟩.append(⟨*operand*⟩) | ⟨*list_var*⟩.pop()

**# expr**
⟨*operand*⟩ ::= ⟨*number*⟩ | ⟨*var*⟩
⟨*expr*⟩ ::= ⟨*operand*⟩ | ⟨*arithm_expr*⟩ | ⟨*list_var*⟩[⟨*number*⟩]
⟨*arithm_expr*⟩ ::= ⟨*operand*⟩ + ⟨*operand*⟩ | ⟨*operand*⟩ − ⟨*operand*⟩
⟨*bool_expr*⟩ ::= ⟨*operand*⟩ == ⟨*operand*⟩ | ⟨*operand*⟩ != ⟨*operand*⟩

---

We design a "baby" version of Python by enforcing multiple constraints. Specifically:

- The generated program is a standalone Python function. The above production rules delineate the generation of the function body. We prepend, e.g., `define function(x, lst_y, cond_z):` before the function body and replace the input arguments `x, lst_y, cond_z` with randomly sampled values, such as `function(x=7, lst_y=[1,2,3], cond_z=False)`, when calling the function.
- The variables are one of the following: integer, list of integers, and boolean value. Additionally, only list indexing is allowed to access a list (i.e., no slicing). The operations to a list can be either append or pop. We restrict the variable type, data structure and associated operations to enforce simplicity of the program and the amount of knowledge prior required to execute the program.
- Supported binary operations include only `+`, `-`, `==`, and `!=`, as many models struggle to reliably execute other operations such as `*`, `%`, and `<`.
- Complex single-line expressions are disabled to ensure the simplicity of transition between lines. For instance, we do not allow evaluation of long expressions, such as `(x + lst[5] - y - 6)`. Preliminary experiments reveal that models sometimes fail to correctly parse the conditions such as `if x != 0:`. Therefore in the production rules, we first evaluate the condition expression, and then assign the value to a boolean variable (e.g., `cond_x`). The condition variable is then used for *if* and *while* blocks to separate condition evaluation from branching or jump.

- In our base evaluation set, we exclude `else:` because this line is occasionally skipped by the models, or by the tracers depending on low-level compiler optimization. This complicates the evaluation process, specifically aligning model responses with ground-truths, therefore we opt to use solely the `if condition:` to enable conditional branching.
- To enforce terminable *while* loops, and to reduce execution errors, we impose context-dependent rules to the *while* statement. Concretely, the terminable *while* block always starts with the initialization of a *counter* variable that controls the number of looping steps. The termination condition is evaluated both before the line of `while cond:` and at the end of the *while* block. The `<while cond, *>` in the production rules denote a context-sensitive expansion, forcing the production to use the previously generated variable name for `while cond:`.
- Terminals are randomly selected from a pool of names given their types. In our current configuration, variables are named with a single letter, list variables are named as `lst_x` where x is a single letter, boolean conditions are named following the pattern `cond_x` where x is a single letter. We provide interface to override the naming of variables and certain production rules for flexibly configuring *L*0-Bench.

Various configurations or filters can be applied during synthetic program generation. In our main results, we adopt the configurations below:

- Input integers are capped at 10
- Input list sizes are randomly sampled between 5 and 10.
- The maximum input program lines of code is 50
- The maximum scope depth is 1, i.e., there is no nested loops or conditional branches.
- The *while* loops terminate when the counter reaches a number $\leq 100$.
- The maximum expansion depth of production rules is set to 200.
- We generate 64 few-shot step-by-step demonstrations for each program.

The above configurations can be modified to generate more difficult *L*0-Bench split, e.g., increasing the size of input lists, the maximum input integers, enabling deep nested loops or branches, filtering programs and inputs that require large number of trace steps to complete.

## B   Ablation of simple operations & procedures

The synthetic programs are combinations of basic programming elements, which we ablate in this section. Concretely, we generate programs that contain only one of the basic ingredients of conventional programming languages by overriding existing production rules, such as setting `<stmt> ::= <assignment> | <if_block>` to test the conditional branching control. For the ablation studies, we limit the program's lines of code to be at maximum 15, and test traces much shorter than even the short split in our default data except for the ablation of *while* blocks. The ablations serve as a "prior check" of the ability to understand basic programming routines.

Table 5 shows that even small models of 7B parameters can perform the evaluation of arithmetic (max integer values set to 10) and comparison expressions with high pass rate, when the operands are not variables. When the operands are previously assigned variable names e.g. `z = y + x`, (w/ var. look up in the table), models show consistent degradation, with the 7B models degrading more. As such, during program generation we exclude the statements where both operands are variable names. To test the control flow, we insert simple assignment statements inside the *if* block or *while* loop. All models follow the control flow with acceptable accuracy though non-marginal gaps remain between small and larger models. The lists can be considered as the "memory" that a processor operates on. We independently test list indexing, append and pop operations at three list sizes *L*: short ($L \in [5, 10]$), medium ($L \in [25, 30]$), and long ($L \in [50, 55]$). Results demonstrate significantly worse performance for longer lists, even for the closed-source model gpt-4o. Additionally, list indexing with a previously assigned variable (e.g., `list[x]`, list w/ var. look up in the table) also leads to consistent degradation, therefore, the list indexing operation in our baseline evaluation set always uses numbers instead of variable names as list indices.

| | # Target Steps | gpt-4o | qwen2.5 coder-7b-inst | qwen2.5-7b inst-1m | qwen2.5-coder -32b-inst | qwen2.5 -32b-inst |
|---|---|---|---|---|---|---|
| **Arithm. and comp op.** | | | | | | |
| + and - | 5.3 | 5.3 | 5.2 | 5.3 | 5.2 | 5.3 |
| + and -, w/ var. look up | 5.3 | 5.0 | 4.3 | 4.2 | 4.5 | 4.7 |
| == and != | 5.3 | 5.3 | 5.2 | 5.3 | 5.2 | 5.1 |
| == and != w/ var. look up | 5.3 | 5.3 | 4.1 | 3.3 | 4.8 | 5.1 |
| **Control** | | | | | | |
| if block | 4.51 | 4.51 | 4.34 | 3.97 | 4.47 | 4.48 |
| while block | 12.0 | 11.9 | 11.6 | 11.2 | 11.7 | 11.9 |
| **Memory (list) op & access** | | | | | | |
| `x = lst[*]` | | | | | | |
| short list | 5.3 | 5.2 | 2.4 | 1.4 | 4.7 | 5.0 |
| short list w/ var. look up | 5.4 | 4.6 | 1.5 | 1.3 | 2.9 | 3.5 |
| medium list | 5.3 | 5.0 | 0.5 | 0.3 | 1.2 | 2.4 |
| medium list w/ var. look up | 5.3 | 3.5 | 0.4 | 0.5 | 1.3 | 1.9 |
| long list | 5.3 | 3.2 | 0.3 | 0.2 | 0.5 | 0.7 |
| long list w/ var. look up | 5.3 | 2.8 | 0.4 | 0.5 | 0.8 | 1.5 |
| `.append(*)` | | | | | | |
| short list | 5.3 | 5.3 | 4.2 | 4.5 | 5.1 | 5.3 |
| medium list | 5.3 | 5.3 | 2.9 | 3.7 | 4.8 | 5.0 |
| long list | 5.4 | 5.3 | 3.0 | 3.6 | 5.0 | 5.2 |
| `.pop()` | | | | | | |
| short list | 5.3 | 5.3 | 4.0 | 2.8 | 4.8 | 4.6 |
| medium list | 5.3 | 5.1 | 1.4 | 0.5 | 2.0 | 3.6 |
| long list | 5.4 | 5.3 | 0.6 | 0.5 | 2.2 | 3.3 |

Table 5: Average number of steps to the first error in ablation study described in Appendix B.

## C  Additional Discussion & Future Directions

**Additional Related Works.** *L*0-Bench focuses on the coding domain (Gao et al., 2023; Chen et al., 2023). We select Python for its better readability and its prevalence in training data. Additionally, programs can be executed to produce traces as verifiable intermediate steps. This work tests models' ability to produce execution trace or simulating the execution process via generation (Chen et al., 2024a; Li et al., 2024a; La Malfa et al., 2024), instead of calling external symbolic reasoner (Pan et al., 2023). The procedures are deterministically defined, thus separating the deduction (or applying general rules to instances) from searching over multiple valid deduction steps (Saparov et al., 2023; Saparov & He, 2023). Moreover, similar to other synthetic reasoning benchmarks (Mirzadeh et al., 2025; Zhou et al., 2025), programs in *L*0-Bench can be conveniently generated following rule-based methods. While the subject is in programming language, the task setup is different than code completion tasks (Ni et al., 2024b; Liu et al., 2024) or other related coding-specific tasks (Yan et al., 2024).

**Limited Domain.** While *L*0-Bench offers a systematic evaluation of procedural correctness, it is limited by its specific domain in programming language. Despite the imposed constraints for creating a more general and less coding-specific setup, it remains unclear to what extent the performance on step-by-step execution of simple Python procedures generalizes to realistic, more complex, long-horizon agentic tasks. Fully aware of such a limitation, we argue that *L*0-Bench nonetheless provides a valuable check of precise and consistent rule-following, thus providing a more fine-grained diagnosis of reasoning capability, complementing existing advanced reasoning benchmarks that often test capabilities in collections. We leave systematic investigation on the correlation between model performance on *L*0-Bench and other real-world tasks as future work.

**Future Directions.** Besides stress testing models with *L*0-Bench data generated with harder configurations, multiple future directions are promising to explore which are out of the scope of this work:

- **Relationship with real-world tasks**: As mentioned above, a promising direction is to explore how *L*0-Bench can serve as a diagnostic test for real-world (Jimenez et al., 2024) agentic tasks. To understand this, one can perform correlational study between a variety of real-world multi-step tasks and instantiations of *L*0-Bench with various configurations. Tasks covering diverse domains (e.g., administrative work flow, tabular data manipulation, etc) and capabilities (e.g., induction, abduction) can help better reveal

relationship between $L0$-Bench and realistic scenarios. Due to the absence of enough suitable real-world tasks and the difficulty in evaluating them, we leave this direction as future work.

• **Systematic study of various test-time scaling dimensions**: We have shown the individual and synergistic gains of three test-time scaling dimensions: input context length, number of solutions for majority voting, and long chain-of-thought inference. However, the study remains superficial without probing deep into the trade-offs among these dimensions, which are likely model-specific. Systematically investigating the relationship between system specialty (e.g., domain specific, long-context, fast inference) vs. various test-time scaling dimensions can help inform better cost-effective model deployment.

• **Reasoning Effort vs. Task Difficulty**: Despite the simplicity of procedures in $L0$-Bench (e.g., incrementing a counter, jumping to the right line based on the conditions), current reasoning models often sub-optimally spend thousands of tokens on "self-reflective thinking." Ideally, reasoning effort increases with task difficulty. The examples in our baseline $L0$-Bench are designed to be simple *on purpose* (see Appendix A and B for constraints and ablation studies for justifying some of the simplifications), therefore much less reasoning effort is expected than what is observed in existing long-CoT models. While the number of thought tokens does reduce as the number of demonstrations in the context increases (Table 10), more systematic study is required to fully understand the behavior of current reasoning-enhanced models. We list a few methods to increase task difficulty of $L0$-Bench: increase input list size, increase max scope depth (to enable nested loops), enable long expression evaluation, enable variable look-ups, increase program length, enable other arithmetic operations, etc. Effectively reducing reasoning verbosity could help reduce latency, on top of faster inference algorithms or architectural modifications.

• **Teach models to follow procedures**: $L0$-Bench provides scripts to generate large amount of synthetic data according to specified configurations. This allows for fine-tuning experiments on scalable generation of synthetic data. The controllable framework can potentially help investigate the effects of SFT and RL methods on narrowing the gap between *pass@1* and *pass@k*, and how these training strategies differ from each other in terms of generalization.

• **Architectural limitations**: Examples in $L0$-Bench are compositions of fundamental operations that may help reveal key architectural constraints. The "atomic" operations, such as accurately copying lists followed by simple modifications (pop or append), or correctly retrieving and applying conditions from distant code sections (while loop), can expose fundamental limitations of existing architectures and help inform future improvements.

## D  Prompt template & Example

See Figure 5 and Table 6.

## E  Complete results

See Table 7 and 8 for complete results.

## F  Additional Results

See Table 9 for test-time scaling of R1-Distill Llama and Table 10 for the number of thought tokens generated by reasoning-enhanced models.

> Please execute the following program by outputting a trace of the program execution. The trace should include the line number, the variable name and the updated value of the variable after executing that line. Please follow the examples below:
>
> Program:
> '''
> {input_program}
> '''
> {fewshots}
>
> Input:
> '''
> function({function_args})
> '''
>
> Please follow the format above and provide the program execution trace starting with "'L2,
>
> Output:

Figure 5: Template for for the default setup of $L0$-Bench. The input programs, few-shot demonstrations, and test input arguments are filled on the fly during evaluation.

## G  Pass@k evaluation

See Figure G for pass@k evaluation code.

```
def pass_at_k(n, c, k):
    """
    :param n: total number of samples
    :param c: number of correct samples
    :param k: k in pass@k
    """
    if n - c < k:
        return 1.0
    return 1.0 - np.prod(1.0 - k / np.arange(n - c + 1, n + 1))
```

Figure 6: Pass@K evaluation code

## H  Models

See Table 11 for a summary of models evaluated in the this work.

## I  Example Long chain-of-thought output

See Tables below for an example response by a top-performing open-source model QwQ-32B.

Please execute the following program by outputting a trace of the program execution. The trace should include the line number, the variable name and the updated value of the variable after executing that line. Please follow the examples below:

Program:
"'

L1 def function(y, v, w, lst_x, lst_z, lst_w, cond_y, cond_x):
L2    if cond_y:
L3        lst_w.append(y)
L4    lst_x.pop()
L5    lst_x.append(8)
L6    cond_y = 6 == 3
L7    if cond_y:
L8        cond_c = 6 == 1
L9        lst_z.pop()
L10       lst_z.pop()
L11   lst_z.append(w)
L12   if cond_x:
L13       lst_z.append(3)
L14       cond_z = 5 != 0
L15       i = 3 + 9
L16       lst_w.append(w)
L17   lst_x.pop()
L18   lst_w.append(v)
L19   cond_d = 1 != v
L20   if cond_y:
L21       lst_x.pop()
L22   return
"'

Input:
"'
function(y=0, v=2, w=8, lst_x=[9, 3, 9, 9, 7, 8], lst_z=[6, 6, 5, 6, 4, 7, 2, 8, 1], lst_w=[0, 2, 6, 8, 1], cond_y=False, cond_x=True)
"'

Output:
"'
L2,
L4,lst_x:[9,3,9,9,7]
L5,lst_x:[9,3,9,9,7,8]
L6,cond_y:False
L7,
L11,lst_z:[6,6,5,6,4,7,2,8,1,8]
L12,
L13,lst_z:[6,6,5,6,4,7,2,8,1,8,3]
L14,cond_z:True
L15,i:12
L16,lst_w:[0,2,6,8,1,8]
L17,lst_x:[9,3,9,9,7]
L18,lst_w:[0,2,6,8,1,8,2]
L19,cond_d:True
L20,
L22,
"'

Input:
"'
function(y=8, v=2, w=7, lst_x=[9, 3, 5, 2, 6, 0], lst_z=[0, 8, 4, 5, 8, 4, 4], lst_w=[2, 8, 2, 1, 7, 9, 9, 5, 8, 5], cond_y=False, cond_x=False)
"'

Please follow the format above and provide the program execution trace starting with "'L2,

Output:

Table 6: Example prompt in $L0$-Bench. Only one-shot demonstration is provided in this example. The number of few-shot demonstrations is an adjustable hyperparameter in $L0$-Bench. We set to 4 in our main experiments.

| Model | pass@1 | pass@5 | pass@15 | pass@31 | maj@5 | maj@15 | maj@31 |
|---|---|---|---|---|---|---|---|
| *Short Traces* | | | | | | | |
| o1 | 96.60 | - | - | - | - | - | - |
| o3-mini | 98.80 | - | - | - | - | - | - |
| gpt-4o-2024-11-20 | 85.80 | - | - | - | - | - | - |
| claude-3-7-sonnet | 98.60 | - | - | - | - | - | - |
| claude-3-5-sonnet-20241022 | 89.20 | - | - | - | - | - | - |
| QwQ-32B | 95.26 | 99.70 | 99.90 | 100.00 | 98.60 | 98.60 | 99.20 |
| DeepSeek-R1-Distill-Llama-70B | 54.27 | 83.30 | 92.50 | 92.50 | 64.00 | 70.80 | 73.40 |
| DeepSeek-R1-Distill-Qwen-32B | 73.21 | 98.00 | 100.00 | 100.00 | 90.00 | 94.40 | 95.60 |
| DeepSeek-R1-Distill-Qwen-14B | 45.34 | 87.00 | 98.00 | 99.00 | 63.00 | 77.20 | 81.60 |
| DeepSeek-R1-Distill-Qwen-7B | 10.72 | 37.00 | 59.00 | 71.00 | 14.80 | 25.80 | 29.80 |
| Deepseek-R1 | 98.10 | 100.00 | 100.00 | 100.00 | 99.80 | 99.60 | 100.00 |
| Llama-3.1-405B-Instruct | 77.20 | 93.40 | 97.00 | 98.20 | 81.00 | 84.60 | 85.60 |
| Llama-3.1-70B-Instruct | 58.19 | 83.20 | 90.60 | 92.90 | 67.60 | 70.40 | 70.60 |
| Llama-3.1-8B-Instruct | 7.24 | 19.00 | 30.00 | 38.00 | 8.60 | 12.00 | 13.00 |
| Qwen2.5-72B-Instruct | 65.79 | 88.40 | 93.10 | 94.90 | 76.60 | 79.20 | 80.80 |
| Qwen2.5-32B-Instruct | 61.96 | 87.10 | 93.30 | 95.20 | 72.80 | 76.40 | 77.20 |
| Qwen2.5-14B-Instruct | 30.00 | 58.00 | 72.90 | 80.50 | 38.40 | 41.20 | 42.60 |
| Qwen2.5-7B-Instruct | 20.80 | 47.00 | 62.00 | 70.00 | 28.20 | 34.40 | 36.80 |
| Qwen2.5-Coder-32B-Instruct | 57.94 | 83.20 | 91.00 | 93.90 | 66.60 | 69.80 | 71.40 |
| Qwen2.5-Coder-7B-Instruct | 28.00 | 59.00 | 74.00 | 81.00 | 36.80 | 42.20 | 47.80 |
| *Medium Traces* | | | | | | | |
| o1 | 97.40 | - | - | - | - | - | - |
| o3-mini | 67.20 | - | - | - | - | - | - |
| gpt-4o-2024-11-20 | 76.40 | - | - | - | - | - | - |
| claude-3-7-sonnet | 97.00 | - | - | - | - | - | - |
| claude-3-5-sonnet-20241022 | 87.80 | - | - | - | - | - | - |
| QwQ-32B | 84.67 | 99.40 | 100.00 | 100.00 | 93.00 | 94.80 | 96.00 |
| DeepSeek-R1-Distill-Llama-70B | 72.44 | 93.90 | 98.00 | 98.80 | 80.20 | 87.40 | 88.60 |
| DeepSeek-R1-Distill-Qwen-32B | 86.49 | 99.00 | 100.00 | 100.00 | 95.60 | 97.60 | 98.00 |
| DeepSeek-R1-Distill-Qwen-14B | 66.74 | 94.00 | 99.00 | 100.00 | 78.80 | 87.00 | 90.00 |
| DeepSeek-R1-Distill-Qwen-7B | 4.33 | 16.00 | 30.00 | 41.00 | 5.20 | 12.20 | 14.40 |
| Deepseek-R1 | 95.88 | 99.80 | 99.80 | 99.80 | 99.20 | 99.40 | 99.80 |
| Llama-3.1-405B-Instruct | 61.70 | 84.40 | 91.20 | 93.90 | 68.20 | 69.80 | 71.60 |
| Llama-3.1-70B-Instruct | 40.06 | 67.40 | 79.40 | 85.20 | 47.80 | 51.00 | 51.80 |
| Llama-3.1-8B-Instruct | 7.29 | 18.10 | 27.70 | 35.30 | 7.60 | 12.00 | 13.20 |
| Qwen2.5-72B-Instruct | 56.99 | 81.40 | 88.80 | 91.90 | 63.20 | 67.80 | 69.60 |
| Qwen2.5-32B-Instruct | 49.61 | 77.70 | 87.00 | 90.90 | 60.00 | 64.60 | 65.60 |
| Qwen2.5-14B-Instruct | 22.34 | 48.70 | 64.80 | 73.30 | 28.20 | 30.80 | 32.80 |
| Qwen2.5-7B-Instruct | 9.60 | 28.00 | 42.00 | 52.00 | 13.00 | 19.00 | 21.80 |
| Qwen2.5-Coder-32B-Instruct | 45.48 | 71.20 | 81.20 | 85.30 | 51.60 | 56.00 | 57.80 |
| Qwen2.5-Coder-7B-Instruct | 16.74 | 38.00 | 54.00 | 63.00 | 21.60 | 28.40 | 31.80 |

Table 7: Detailed results for Table 1.

| Model | pass@1 | pass@5 | pass@15 | pass@31 | maj@5 | maj@15 | maj@31 |
|---|---|---|---|---|---|---|---|
| *Long Traces* | | | | | | | |
| o1 | 91.40 | - | - | - | - | - | - |
| o3-mini | 46.20 | - | - | - | - | - | - |
| gpt-4o-2024-11-20 | 74.20 | - | - | - | - | - | - |
| claude-3-7-sonnet | 97.00 | - | - | - | - | - | - |
| claude-3-5-sonnet-20241022 | 82.20 | - | - | - | - | - | - |
| QwQ-32B | 84.47 | 98.70 | 99.70 | 100.00 | 92.60 | 95.40 | 96.20 |
| DeepSeek-R1-Distill-Llama-70B | 65.36 | 90.80 | 96.50 | 98.00 | 77.20 | 81.80 | 83.40 |
| DeepSeek-R1-Distill-Qwen-32B | 83.58 | 98.00 | 100.00 | 100.00 | 92.00 | 94.80 | 94.80 |
| DeepSeek-R1-Distill-Qwen-14B | 60.96 | 90.00 | 96.00 | 98.00 | 75.20 | 85.60 | 85.80 |
| DeepSeek-R1-Distill-Qwen-7B | 0.19 | 1.00 | 2.00 | 3.00 | 0.40 | 0.60 | 0.60 |
| Deepseek-R1 | 94.84 | 99.30 | 99.70 | 99.80 | 99.20 | 99.20 | 99.60 |
| Llama-3.1-405B-Instruct | 59.20 | 81.60 | 89.10 | 92.50 | 63.60 | 66.20 | 68.00 |
| Llama-3.1-70B-Instruct | 36.58 | 66.20 | 80.40 | 86.20 | 41.60 | 46.80 | 48.40 |
| Llama-3.1-8B-Instruct | 5.66 | 15.50 | 24.90 | 31.30 | 6.60 | 8.80 | 9.20 |
| Qwen2.5-72B-Instruct | 52.57 | 78.20 | 88.00 | 91.70 | 58.80 | 61.20 | 62.20 |
| Qwen2.5-32B-Instruct | 47.62 | 75.00 | 86.00 | 90.00 | 54.40 | 60.40 | 62.80 |
| Qwen2.5-14B-Instruct | 23.39 | 50.50 | 66.30 | 74.70 | 28.80 | 34.20 | 36.40 |
| Qwen2.5-7B-Instruct | 8.59 | 22.80 | 36.10 | 18.20 | 11.60 | 16.00 | 18.20 |
| Qwen2.5-Coder-32B-Instruct | 40.06 | 68.00 | 79.90 | 84.80 | 48.20 | 51.00 | 53.00 |
| Qwen2.5-Coder-7B-Instruct | 13.29 | 34.00 | 49.00 | 59.00 | 17.80 | 24.60 | 26.40 |
| *Extra Long Traces* | | | | | | | |
| o1 | 82.60 | - | - | - | - | - | - |
| o3-mini | 30.00 | - | - | - | - | - | - |
| gpt-4o-2024-11-20 | 69.00 | - | - | - | - | - | - |
| claude-3-7-sonnet | 87.80 | - | - | - | - | - | - |
| claude-3-5-sonnet-20241022 | 73.80 | - | - | - | - | - | - |
| QwQ-32B | 81.99 | 95.70 | 98.30 | 99.20 | 88.40 | 92.20 | 93.00 |
| DeepSeek-R1-Distill-Llama-70B | 55.91 | 83.20 | 93.10 | 97.10 | 63.20 | 70.60 | 73.00 |
| DeepSeek-R1-Distill-Qwen-32B | 78.23 | 95.00 | 98.00 | 99.00 | 87.60 | 90.80 | 91.20 |
| DeepSeek-R1-Distill-Qwen-14B | 54.10 | 82.00 | 91.00 | 94.00 | 64.60 | 73.80 | 76.60 |
| DeepSeek-R1-Distill-Qwen-7B | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| Deepseek-R1 | 77.04 | 88.10 | 90.80 | 92.10 | 85.40 | 88.80 | 90.00 |
| Llama-3.1-405B-Instruct | 54.00 | 77.40 | 84.90 | 88.10 | 60.60 | 65.60 | 67.20 |
| Llama-3.1-70B-Instruct | 34.57 | 62.40 | 74.30 | 79.40 | 42.80 | 47.80 | 48.00 |
| Llama-3.1-8B-Instruct | 5.78 | 15.10 | 24.10 | 30.90 | 6.00 | 9.40 | 10.40 |
| Qwen2.5-72B-Instruct | 45.60 | 68.70 | 77.30 | 81.10 | 52.00 | 54.20 | 54.80 |
| Qwen2.5-32B-Instruct | 42.44 | 69.00 | 79.00 | 84.00 | 50.40 | 56.60 | 58.40 |
| Qwen2.5-14B-Instruct | 19.11 | 42.70 | 57.40 | 66.20 | 24.20 | 28.20 | 30.80 |
| Qwen2.5-7B-Instruct | 6.29 | 18.20 | 28.80 | 36.10 | 9.60 | 12.00 | 14.40 |
| Qwen2.5-Coder-32B-Instruct | 32.38 | 57.00 | 68.10 | 73.60 | 39.80 | 45.00 | 44.80 |
| Qwen2.5-Coder-7B-Instruct | 10.20 | 28.00 | 40.00 | 47.00 | 17.80 | 21.80 | 23.60 |

Table 8: Detailed results for Table 1.

|  | # few shot | **Steps To Err.** | | **Trace Acc.** | | |
|---|---|---|---|---|---|---|
|  |  | Single Attempt | maj vote @ 31 | Single Attempt | maj vote @ 31 | pass @ 31 |
| Short Trace |  |  |  |  |  |  |
|  | 2 | 8.3 | 10.4 | 43.7 | 66.0 | 94.6 |
|  | 4 | 9.2 | 11.1 | 53.5 | 74.2 | 95.7 |
|  | 8 | 9.6 | 11.4 | 58.9 | 77.6 | 98.1 |
|  | 16 | 10.0 | 11.6 | 63.0 | 81.4 | 97.8 |
|  | 32 | 10.0 | 11.4 | 62.5 | 79.0 | 98.6 |
| Medium Trace |  |  |  |  |  |  |
|  | 2 | 65.4 | 73.4 | 65.5 | 83.8 | 98.2 |
|  | 4 | 69.9 | 76.0 | 72.9 | 87.6 | 98.6 |
|  | 8 | 71.3 | 76.7 | 75.1 | 88.2 | 99.2 |
|  | 16 | 71.3 | 76.4 | 74.1 | 86.6 | 98.6 |
|  | 32 | 67.9 | 72.7 | 69.5 | 81.6 | 97.4 |

Table 9: Scaling along all three test-time scaling dimensions. `R1-Distilled-Llama-70B`, a long-CoT model achieves further gains when provided with multiple step-by-step demonstrations and taking majority vote over solutions generated in parallel.

| # thought tokens | Short Trace | Medium Trace | Long Trace | Extra Long Trace |
|---|---|---|---|---|
| DeepSeek-R1-Distill-Qwen-7B | 1244 | 3198 | 4568 | 4366 |
| DeepSeek-R1-Distill-Qwen-14B | 2618 | 2294 | 2084 | 1967 |
| DeepSeek-R1-Distill-Qwen-32B | 1575 | 1967 | 1596 | 1753 |
| QwQ-32B | 3360 | 7388 | 6793 | 6830 |
| DeepSeek-R1-Distill-Llama-70B |  |  |  |  |
|     4-shot | 1820 | 1279 | 1176 | 1185 |
|     8-shot | 1788 | 1205 | 1046 | 1105 |
|     32-shot | 1290 | 1187 | 1134 | - |

Table 10: Thought tokens are enclosed by `<think>` and `</think>`. Within the same model series (R1-distilled), smaller models generate more thought tokens, with the 7B model frequently failing to exit the think pattern within the requested token budget. QwQ-32B produces significantly longer chain-of-thought than same-sized R1-distilled model.

| Model Name | Size | Thinking | Huggingface / API |
|---|---|---|---|
| o1 (Jaech et al., 2024) | - | Y | o1 |
| o3-mini (OpenAI, 2025) | - | Y | o3-mini |
| gpt-4o-2024-11-20 (Hurst et al., 2024) | - | N | gpt-4o-2024-11-20 |
| claude-3-7-sonnet (Anthropic, 2025) | - | Y | claude-3-7-sonnet-20250219 |
| claude-3-5-sonnet-20241022 (Anthropic, 2024) | - | N | claude-3-5-sonnet-20241022 |
| DeepSeek-R1 (Guo et al., 2025) | 37B / 671B | Y | deepseek-ai/DeepSeek-R1 |
| DeepSeek-R1-Distill-Qwen-32B (Guo et al., 2025) | 32.5B | Y | deepseek-ai/DeepSeek-R1-Distill-Qwen-32B |
| DeepSeek-R1-Distill-Qwen-14B (Guo et al., 2025) | 14.7B | Y | deepseek-ai/DeepSeek-R1-Distill-Qwen-14B |
| DeepSeek-R1-Distill-Qwen-7B (Guo et al., 2025) | 7.6B | Y | deepseek-ai/DeepSeek-R1-Distill-Qwen-7B |
| DeepSeek-R1-Distill-Llama-70B (Guo et al., 2025) | 70B | Y | deepseek-ai/DeepSeek-R1-Distill-Llama-70B |
| QwQ-32B (Qwen, 2025) | 32.5B | Y | Qwen/QwQ-32B |
| Qwen2.5-72B-Instruct (Yang et al., 2024) | 72.7B | N | Qwen/Qwen2.5-72B-Instruct |
| Qwen2.5-32B-Instruct (Yang et al., 2024) | 32.5B | N | Qwen/Qwen2.5-32B-Instruct |
| Qwen2.5-14B-Instruct (Yang et al., 2024) | 14.7B | N | Qwen/Qwen2.5-14B-Instruct |
| Qwen2.5-7B-Instruct (Yang et al., 2024) | 7.6B | N | Qwen/Qwen2.5-7B-Instruct |
| Qwen2.5-Coder-32B-Instruct (Yang et al., 2024) | 32.5B | N | Qwen/Qwen2.5-Coder-32B-Instruct |
| Qwen2.5-Coder-7B-Instruct (Yang et al., 2024) | 7.6B | N | Qwen/Qwen2.5-Coder-7B-Instruct |
| Meta-Llama-3.1-405B-Instruct (Grattafiori et al., 2024) | 405B | N | meta-llama/Llama-3.1-405B-Instruct |
| Meta-Llama-3.1-70B-Instruct (Grattafiori et al., 2024) | 70B | N | meta-llama/Llama-3.1-70B-Instruct |
| Meta-Llama-3.1-8B-Instruct (Grattafiori et al., 2024) | 8B | N | meta-llama/Llama-3.1-8B-Instruct |

Table 11: Model details.

**QwQ-32B example output**

Okay, let's tackle this problem step by step. I need to trace the execution of the given Python function with the provided input. Let me start by understanding the function and the input parameters.\n\nThe function is called function, and the input parameters are:\ny=8, v=2, w=7,\nlst_x = [9, 3, 5, 2, 6, 0],\nlst_z = [0, 8, 4, 5, 8, 4, 4],\nlst_w = [2, 8, 2, 1, 7, 9, 9, 5, 8, 5],\ncond_y=False, cond_x=False.\n\nFirst, I'll list out all the variables and their initial values before the function starts executing. The function's code is provided, so I'll go line by line.\n\nStarting with line L2: the first line is the function definition, but the actual execution starts when the function is called. The first executable line is L2's condition check.\n\nLet me go through each line step by step, keeping track of variable changes.\n\nLine L2: if cond_y:\nThe current cond_y is False (from input), so the condition is false. Therefore, the code inside the if block (line L3) is skipped. So after L2, nothing changes here. The trace entry for L2 is just a comma since no variables are updated.\n\nNext, line L4: lst_x.pop()\nThe list lst_x is [9,3,5,2,6,0]. The pop() method removes the last element by default. So popping removes 0, resulting in [9,3,5,2,6]. So after L4, lst_x is now [9,3,5,2,6]. So the trace entry is L4,lst_x: the new list.\n\nLine L5: lst_x.append(8)\nAppending 8 to the current lst_x (after the pop), so the new list becomes [9,3,5,2,6,8]. So the trace entry is L5,lst_x with the new list.\n\nLine L6: cond_y = 6 == 3. The right side is a boolean expression. 6 ==3 is False, so cond_y becomes False. The trace entry is L6,cond_y:False.\n\nLine L7: if cond_y: Since cond_y is still False, the block under L7 (lines L8-L10) is skipped. So L7's trace is just a comma.\n\nLine L11: lst_z.append(w). The current w is 7. The original lst_z is [0,8,4,5,8,4,4]. Appending 7 gives [0,8,4,5,8,4,4,7]. So the trace entry is L11,lst_z with the new list.\n\nLine L12: if cond_x: cond_x is False (from input), so the block (lines L13-L18) is skipped. Wait, wait. Wait, the lines under L12's if are L13 to L16, and then after that, L17 and L18 are outside the if block? Let me check the code structure again.\n\nLooking back at the program code:\n\nAfter L12 comes L13, which is indented under the if cond_x block. Then L17 is after that block. So lines L13-L16 are inside the if cond_x clause. Since cond_x is False, those lines are skipped. So after L12, we proceed to L17.\n\nWait, the code structure:\n\nL12: if cond_x:\nL13: lst_z.append(3)\nL14: cond_z = 5 !=0\nL15: i =3+9\nL16: lst_w.append(w)\nThen L17 is outside the if block, so it's executed regardless.\n\nSo after L12, since cond_x is False, we skip L13-L16 and go to L17.\n\nLine L17: lst_x.pop(). Current lst_x after L5 was [9,3,5,2,6,8]. Wait, wait, let me retrace:\n\nWait, after L4, lst_x was [9,3,5,2,6], then L5 appends 8, making it [9,3,5,2,6,8]. Then L17 is next. So popping the last element (8) leaves [9,3,5,2,6]. So the trace entry for L17 is lst_x: [9,3,5,2,6].\n\nLine L18: lst_w.append(v). v is 2. The original lst_w is [2,8,2,1,7,9,9,5,8,5]. Appending 2 gives [2,8,2,1,7,9,9,5,8,5,2]. So the trace entry is L18,lst_w with the new list.\n\nLine L19: cond_d = 1 != v. v is 2. 1 !=2 is True, so cond_d is True. So the trace entry is L19,cond_d:True.\n\nLine L20: if cond_y: cond_y is False (from L6), so the block (L21) is skipped. So L20's trace is a comma.\n\nFinally, return at L22, so the last entry is L22, with nothing else.\n\nWait, let me check all steps again step by step to ensure I didn't miss anything.\n\nStarting over:\n\nInitial variables:\n\ny=8, v=2, w=7,\nlst_x = [9,3,5,2,6,0]\nlst_z = [0,8,4,5,8,4,4]\nlst_w = [2,8,2,1,7,9,9,5,8,5]\ncond_y = False, cond_x = False.\n\nExecution steps:\n\nL2: if cond_y (False) → skip L3.\n\nL4: lst_x.pop() → removes last element (0) → [9,3,5,2,6]\n\nTrace: L4,lst_x: [9,3,5,2,6]\n\nL5: append 8 → [9,3,5,2,6,8]\n\nTrace: L5,lst_x: [9,3,5,2,6,8]\n\nL6: cond_y = 6==3 → False → Trace: L6,cond_y:False\n\nL7: if cond_y (still False) → skip L8-L10.\n\nL11: append w (7) to lst_z → original lst_z was [0,8,4,5,8,4,4], so new is [0,8,4,5,8,4,4,7]\n\nTrace: L11,lst_z: [0,8,4,5,8,4,4,7]\n\nL12: if cond_x (False) → skip L13-L16.\n\nL17: pop lst_x → current lst_x is [9,3,5,2,6,8], so pop removes 8 → [9,3,5,2,6]\n\nTrace: L17,lst_x: [9,3,5,2,6]\n\nL18: append v (2) to lst_w → original lst_w is [2,8,2,1,7,9,9,5,8,5], so new is adding 2 → [2,8,2,1,7,9,9,5,8,5,2]\n\nWait, wait, wait! Wait, the initial lst_w is [2,8,2,1,7,9,9,5,8,5]. So appending 2 gives the new list as that list plus 2 at the end. So the new list is [2,8,2,1,7,9,9,5,8,5,2]. So the trace entry for L18 is lst_w with that list.\n\nL19: cond_d = 1 != v → v is 2 → 1 !=2 → True → Trace: L19,cond_d:True\n\nL20: if cond_y (False) → skip L21.\n\nL22: return → so the last line is L22 with nothing

\n\n'''\nL2,\nL4,lst_x:[9,3,5,2,6]\nL5,lst_x:[9,3,5,2,6,8]\nL6,cond_y:False\nL7,\n\nL11,lst_z:[0,8,4,5,8,4,4,7]\nL12,\nL17,lst_x:[9,3,5,2,6]\nL18,lst_w:[2,8,2,1,7,9,9,5,8,5,2]\n L19,cond_d:True\nL20,\nL22,\n