

# LĪLA: A Unified Benchmark for Mathematical Reasoning

Anonymous EMNLP submission

## Abstract

From grocery shopping to climate modeling, mathematical reasoning is an essential skill. Not surprisingly, AI agents that are increasingly becoming ubiquitous need to be equipped with mathematical reasoning skills. To both train and evaluate AI agents in this domain, we propose LĪLA, a unified mathematical reasoning benchmark consisting of 23 diverse tasks that are categorized along four dimensions – (i) mathematical abilities *e.g.* arithmetic, calculus (ii) language format *e.g.* question-answering, fill-in-the-blanks (iii) language diversity *e.g.* no language, simple language (iv) external knowledge *e.g.* commonsense, physics. We extend each dataset part of the benchmark by collecting task instructions and solutions in the form of Python programs, emphasizing the need for generating explainable solutions in addition to the final correct answer. We further introduce two evaluation datasets to measure out-of-distribution performance and robustness to language perturbation. Finally, we introduce BHĀSKARA and its variants, a family of baseline models fine-tuned on LĪLA—most importantly, we find that multi-tasking leads to significant improvements (average relative improvement of 21.83% F1 score *vs.* single-task models) while the best performing model only obtains 60.40%, indicating the room for improvement in general mathematical reasoning and understanding.

## 1 Introduction

Mathematical reasoning is required in all aspects of life, from buying vegetables to controlling the world economy. Given the fundamental nature of mathematical reasoning, a flurry of works propose datasets to evaluate mathematical reasoning abilities of AI agents *e.g.* Kushman et al. (2014) (algebra word problems), Mishra et al. (2022b) (arithmetic reasoning), Saxton et al. (2019) (templated math reasoning spanning algebra, calculus, probability, etc). Despite numerous datasets that test for

**Math ability:** basic math  
**Language complexity:** simple language  
**Format:** generative question answering  
**Knowledge:** no external knowledge  
**Instruction:** You are given a question that involves the [calculation of numbers](#). You need to perform either an [addition](#) or [subtraction](#) operation on the numbers. [Generate your answer](#) to the given question.

**Question:** Sara picked 45 pears and Sally picked 11 pears from the pear tree. How many pears were picked in total?

**Program 1:**  

```
def solution(x, y):  
    answer = x + y  
    return answer  
print(solution(45, 11)) # total pears is the sum of  
pears with Sara and Sally
```

**Program 2:**  

```
x = 45  
y = 11  
answer = x + y # total pears is the sum of pears with  
Sara and Sally  
print(answer)
```

**Answer:** 56

Figure 1: A data example with two Python programs in LĪLA. One program annotation uses function construct whereas the other one is a plain script without function. The instruction for each task and categories across four dimensions are annotated for developing LĪLA.

specific mathematical reasoning abilities, a benchmark for systematic and holistic evaluation that covers diverse topics and problem styles does not exist. Further, evaluating high-capacity models on narrowly scoped mathematical reasoning datasets leads to the risk of considerably overestimating the reasoning abilities of these AI systems, highlighting the need for such a benchmark.

To this end, we introduce LĪLA<sup>1</sup>, a unified mathematical reasoning benchmark that consists of 23

<sup>1</sup>Named after *Līlavati*, a 12<sup>th</sup> century mathematical treatise on arithmetic that covers topics like arithmetic and geometric progressions, indeterminate equations and combinations. It is also widely known for the extensive number of math word problems. The author, *Bhāskara* is known for fundamental and original contributions to calculus, physics, number theory, algebra and astronomy (Colebrooke, 1817; Sarkar, 1918; Kolachana et al., 2019)

mathematical reasoning tasks. LILA is constructed by extending 20 existing datasets that not only span a wide range of topics in mathematics but also contain questions with varying degrees of linguistic complexity, diversity in the question format, and the requirement of background knowledge. Importantly, LILA extends all of these datasets to include a solution program as opposed to only an answer. To enable instruction-based learning (Sanh et al., 2021; Wei et al., 2021; Mishra et al., 2022a), we also collect instruction annotations for each dataset in the benchmark.

In order to accurately assess the mathematical reasoning ability of models, merely evaluating the final numerical answer or algebraic expression is not sufficient; the chain of reasoning that leads to the correct solution is equally important, if not more. Therefore, we collect Python programs that serve as solutions for each question in the benchmark. We achieve this in two ways – (i) when datasets already contain a structured solution output captured by a domain specific language (DSL), we algorithmically translate it into a Python program and (ii) manually collect annotations from experts when the original dataset does not contain any such structured solution. LILA unifies various mathematical reasoning datasets under a single problem formulation *i.e.* given an input problem in natural language, generate a Python program that upon execution returns the desired answer. Further, it allows for neural approaches to focus on the high-level aspects of mathematical problem solving (*e.g.* identifying potential solution strategies, decomposing the problem into simpler sub-problems, *etc.*) while leveraging symbolic solvers (*e.g.* using Python as a calculator or Sympy) to perform precise operations like adding huge numbers or simplifying expressions. Figure 1 illustrates a sample from our LILA benchmark that illustrates the question, answer, program, instructions, and category tags.

In addition to being able to evaluate on the chain of reasoning, we also facilitate two other key ways to make a fair assessment of models on mathematical reasoning tasks. In line with Bras et al. (2020), Ribeiro et al. (2020) and Welleck et al. (2022), we propose to evaluate for strong generalization *e.g.* alternate formulations of a problem (“2+2=” vs. “What is two plus two?”) using an out-of-distribution evaluation set LILA-OOD that contains datasets that test for the

same underlying mathematical reasoning skill but was collected independent of the training datasets. further, we collect a robustness split LILA-Robust, that introduces linguistic perturbations (*e.g.* active vs. passive voice) via crowd-sourcing. The evaluation scheme is a combination of the performance on all three sets – LILA-TEST, LILA-OOD and LILA-ROBUST.

## Contributions

1. We propose LILA, a unified benchmark for holistically evaluating models’ mathematical reasoning—measuring performance on out-of-distribution examples and robustness to language perturbations in addition to standard test-set. It consists of 23 tasks categorized along mathematical areas, language complexity, question format and need for external knowledge. Importantly, LILA extends 20 existing datasets with solutions in the form of Python programs and instruction annotations.
2. We introduce BHASKARA, a family of language models fine-tuned on our dataset. Our best-performing model achieves comparable performance to a  $66\times$  larger model pre-trained on both code and language.
3. We provide an analysis of our models’ performance- (1) multitasking helps considerably over task-specific learning both in IID and OOD evaluation (2) program synthesis substantially outperforms answer prediction, (3) few-shot codex has the strongest performance. We also identify areas for improvement, such as identifying what kind of datasets we need to annotate to fill gaps across categories in LILA.

## 2 Related Work

**Mathematical Reasoning Datasets.** Our work builds on a rich history of mathematical reasoning research. Early work in this areas focuses on small-scale datasets testing addition-subtraction (Hosseini et al., 2014), templated questions with equations as parameters (Kushman et al., 2014) and other forms of arithmetic reasoning (Koncel-Kedziorski et al., 2015; Roy and Roth, 2016; Upadhyay et al., 2016; Roy and Roth, 2017, 2018; Ling et al., 2017). Later datasets increase in complexity and scale, incorporating reading comprehension Dua et al. (2019b) and algebra (Saxton et al., 2019). Still other numerical reasoning datasets focus on diversity (Miao et al., 2020a) with multi-

ple categories of numerical reasoning tasks (e.g., Amini et al., 2019). Most recently, new datasets have focused on increasing difficulty, e.g., olympiad problems (Hendrycks et al., 2021b) and adversarial problems (Patel et al., 2021), as well as increasing the knowledge requirements to solve tasks, with a growing focus on commonsense reasoning (Zhou et al., 2019; Zhang et al.; Mishra et al., 2022b).

A separate line of work in mathematical reasoning includes datasets testing mathematical theorem proving (e.g., Li et al., 2021; Wu et al., 2021; Welleck et al., 2021; Zheng et al., 2021; Han et al., 2021). We do not, however, consider theorem proving in our work, choosing instead to focus on numerical reasoning.

**Task Hierarchy and Multi-tasking in Numerical Reasoning.** Inspired by the success of multi-task learning in NLP (Weston et al., 2015), including benchmarks like GLUE (Wang et al., 2018), SuperGLUE (Wang et al., 2019), ORB (Dua et al., 2019a) and multitasking models such as DecaNLP (McCann et al., 2018), UnifiedQA (Khashabi et al., 2020), Unicorn (Lourie et al., 2021) and Muppet (Aghajanyan et al., 2021), NumGLUE (Mishra et al., 2022b) has been proposed as a multi-tasking numerical reasoning benchmark that contains 8 different tasks. LILA expands NumGLUE to provide wider coverage of mathematical abilities, along with evaluation that captures out-of-domain, robustness, and instruction-following performance. Our introduction of mathematical reasoning categories and the evaluation setup is inspired from the task hierarchies in other domains such as in vision (Zamir et al., 2018) and NLP (Rogers et al., 2021) which appear in in large scale benchmarks such as Big-Bench (Srivastava et al., 2022) and Natural Instructions V2 (Wang et al., 2022).

### 3 LILA

LILA contains 23 tasks across 4 dimensions, curated from 44 sub-datasets across 20 dataset sources. In this section, we discuss the construction and composition of the benchmark before providing an analysis through key-statistics.

#### 3.1 Dataset Construction

**Data Sources.** We incorporate 20 existing datasets from the mathematical reasoning literature (Table 21 gives a detailed list), excluding sources where the output is not a number or an expression e.g. theorem proving (Welleck et al., 2021; Han

et al., 2021). Further, we only consider datasets with input in natural or templated language as opposed to formal specifications.

**Unified Format.** We normalize all examples sourced from different datasets to a unified format that contains the following fields:

1. The source dataset
2. Category tag corresponding to *each* of the four dimensions—math ability, language complexity, format, and external knowledge. (see §3.2).
3. Question, in the English language.
4. The answer to the question, a string containing a number, expression, list, or other data format.
5. A set of Python strings that *print* the answer when executed.
6. A task-level instruction in natural language.

We also retain other meta-data that was provided in the original dataset, e.g., Figure 1.

**Algorithmic Program Annotation.** Most of the annotations in the source datasets do not contain output in the form of a Python program. We automatically annotate most datasets by generating Python programs using the annotations (answer, explanation etc.) provided in the source datasets. Where possible, we generate multiple Python programs for a single question. This is to account for variation in the program space such as the choice of data structure, language construct, variable name, and programming style (e.g., declarative vs procedural). For example, Figure 1 gives multiple Python programs solving the same question; in this case one program directly calculates the answer, whereas the other defines a function to solve the problem more generally.

Some datasets contain program annotations that can be captured by a domain-specific language (DSL) in which case we write rules to convert them into Python programs e.g. `volume(sphere,3)` to the Python expression `4/3*math.pi*3**3`. In some cases where a DSL annotation is not provided, we use pattern matching to convert highly templated datasets like the AMPS dataset (Hendrycks et al., 2021b) to our unified format. In other cases, instead of converting the existing dataset, we modify the data generation code to reproduce the dataset with program annotations. For the Deepmind mathematics dataset (Saxton et al., 2019), this allows us to

Category	Math ability	Language	Knowledge	Format
Tasks	basic math, multiplication-division, number theory, algebra, geometry, counting and statistics, calculus, linear algebra, advanced math	no language, simple language, complex language	no background knowledge, commonsense, math, science, computer science, real world knowledge	fill-in-the-blank, generative question answering, multiple-choice, natural language inference, reading comprehension

Table 1: Dataset partitions (Categories) and associated tasks.

create diverse, compositional math problems with program annotations using a sophisticated grammar. In some cases, the original DSL annotations are inconsistent or ambiguous. For instance, the predicate “volume of a circle with radius” may be written as `volume(circle,3)` or `circle_volume(3)`, depending on the choices of the dataset authors. Wherever possible, we manually defined the formula based on each dataset. We compile the annotated programs for each question and retain the ones which compile successfully and produce the same answer as the gold annotated answers.

**Expert Program Annotation.** For many of datasets, it is not possible to obtain a Python program solutions via automated methods described above—either the original dataset contains only the final answer or contains solutions expressed in free-form natural language. For such datasets, we obtain annotations from experts who are proficient in basic programming and high-school level mathematics. See Appendix A.2.1 for details.

**Instruction Annotation.** Given the recently found effectiveness of instruction-learning (Mishra et al., 2022a; Wei et al., 2021; Sanh et al., 2021) for effective generalization, we collect instruction annotation at the task level. Each instruction contains a *definition* that clearly defines the task and provides guidelines, a *prompt* that provides a short and straight forward instruction, and *examples* that facilitate learning by demonstration (Brown et al., 2020). Figure 1 shows an example instruction for the basic math task (§3.2).

### 3.2 Categories and Tasks

We create 4 *views*<sup>2</sup> or categories of LILA along the dimensions of mathematical area, language com-

plexity, external knowledge, and question format. Including all dimensions, data is classified into one of 23 *tasks* (Table 1). By creating multiple views of the benchmark along these axes, we are able to systematically characterize the strengths and weaknesses of existing models along with obtaining a better understanding of the domain itself.

The first category, *math ability*, partitions datasets into tasks common in traditional pedagogy: arithmetic, algebra, geometry, calculus, etc. Examples of each task and the datasets included are in Table 11 and Table 15 respectively.

Our second category, *language complexity*, separates math problems by the complexity of the language used to represent them. This ranges from formal representations only (e.g.,  $1+1=?$ ) to natural language (e.g., “Mariella has 3 pears...”). Examples of each task and the datasets included are in Table 12 and Table 16 respectively.

The next category is based on the requirement of *background knowledge*, the type of knowledge required to solve the problem. For instance, commonsense questions like “How many legs to 3 people have?” or science questions like “Will water boil at 200 degrees celsius?” require different sets of knowledge to answer.

Lastly, we categorize based on *question format*, putting e.g., multiple choice questions under one task and natural language inference under another. Examples of each task and the datasets included are in Table 13 and Table 17 respectively.

### 3.3 LILA-OOD

In order to measure if the model has truly learned the underlying mathematical reasoning skill, it needs to be able to perform well on unseen and preferably, out-of-distribution data. To this end, we evaluate both in-distribution (standard train-test splits) and out-of-distribution performance for each task *i.e.* evaluate on examples that test for

<sup>2</sup>Note that it is *not* a partition of the benchmark as each dimension divides the constituent examples in different ways



### 3.4 LĪLA-ROBUST

### 3.5 Dataset analysis

## 4 Experiments

Statistic	Number
# Total tasks	23
# Total datasets	44
# Total instructions	44
# Total questions	133,815
# Total programs	358,769
Unique questions	132,239
Unique programs	325,597
Unique answers	271,264
Average length of instructions	31.18
Average length of questions	47.72
Average length of programs	47.85

Table 2: Key statistics of LĪLA.

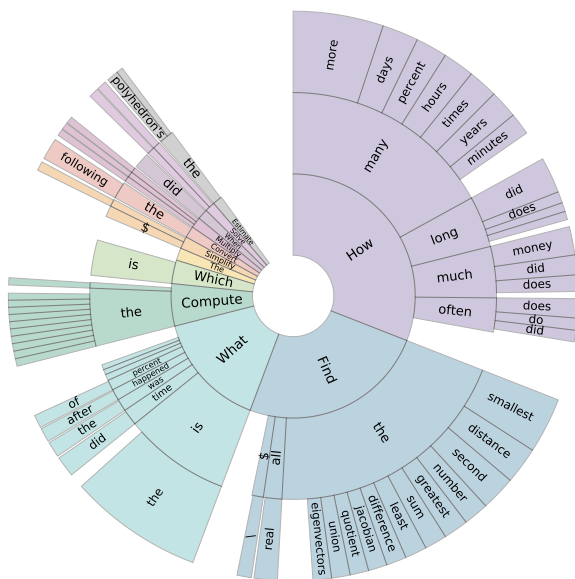


Figure 2: Question n-gram distribution in LĪLA.

**Fine-tuning** We fine-tune a series of GPT-Neo-2.7B models (a pre-trained causal language model (Black et al., 2021)) on LILA. We choose GPT-Neo because it was pre-trained on both natural language and code (Gao et al., 2020), as opposed to a model trained solely on natural language. To assess the capabilities of GPT-Neo

5

→ Supervision/Size		Few-shot, 175B		Few-shot, 175B		Fine-tuned, 2.7B		Fine-tuned, 2.7B		Fine-tuned, 2.7B		Fine-tuned, 2.7B	
↓ Task	Category	GPT3		Codex		Neo-A		Neo-P		Neo-Multi-A		Neo-Multi-P	
		IID	OOD	IID	OOD	IID	OOD	IID	OOD	IID	OOD	IID	OOD
1	Basic math	0.766	<b>0.818</b>	<b>0.791</b>	0.762	0.533	0.523	0.611	0.555	0.693	0.657	<b>0.790</b>	<b>0.787</b>
2	Muldiv	0.479	0.665	<b>0.691</b>	<b>0.790</b>	0.136	0.089	0.388	0.194	0.155	0.083	<b>0.448</b>	<b>0.395</b>
3	Number theory	0.240	0.154	<b>0.472</b>	<b>0.344</b>	0.108	0.095	0.328	0.107	0.129	0.190	<b>0.358</b>	<b>0.293</b>
4	Algebra	0.338	0.130	<b>0.603</b>	<b>0.511</b>	0.164	0.031	0.348	0.051	0.203	<b>0.054</b>	<b>0.473</b>	0.007
5	Geometry	<b>0.283</b>	0.120	0.000	<b>0.250</b>	0.288	0.025	0.077	0.021	<b>0.297</b>	0.105	0.079	<b>0.250</b>
6	Statistics	0.183	<b>0.210</b>	<b>0.650</b>	0.200	0.107	0.008	0.839	0.034	0.115	<b>0.179</b>	<b>0.947</b>	0.164
7	Calculus	0.231	0.208	<b>0.930</b>	<b>0.884</b>	0.138	0.119	0.486	0.334	0.102	0.167	<b>0.495</b>	<b>0.805</b>
8	Linear algebra	0.127	-	<b>0.692</b>	-	0.229	-	<b>0.809</b>	-	0.240	-	0.808	-
9	Advanced math	0.150	-	<b>0.472</b>	-	0.012	-	0.100	-	0.019	-	<b>0.160</b>	-
10	No language	0.213	0.162	<b>0.853</b>	<b>0.770</b>	0.143	0.083	0.698	0.330	0.140	0.138	<b>0.703</b>	<b>0.850</b>
11	Simple language	0.486	0.561	<b>0.568</b>	<b>0.610</b>	0.269	0.243	0.363	0.292	0.332	0.269	<b>0.433</b>	<b>0.384</b>
12	Complex language	0.356	0.413	<b>0.456</b>	<b>0.583</b>	0.147	0.113	0.216	0.106	0.215	0.259	<b>0.288</b>	<b>0.557</b>
13	Fill in the blank	0.710	0.620	<b>0.790</b>	<b>0.660</b>	0.086	0.193	<b>0.304</b>	0.193	0.059	<b>0.519</b>	0.262	<b>0.519</b>
14	Generative QA	0.305	0.385	<b>0.566</b>	<b>0.632</b>	0.142	0.135	0.376	0.199	0.178	0.160	<b>0.476</b>	<b>0.235</b>
15	MCQ	<b>0.801</b>	<b>0.870</b>	0.771	<b>0.870</b>	0.636	0.818	0.652	0.818	0.752	<b>0.888</b>	<b>0.817</b>	<b>0.888</b>
16	NLI	0.500	-	<b>0.710</b>	-	0.221	-	0.212	-	0.566	-	<b>0.893</b>	-
17	RC	0.460	-	<b>0.615</b>	-	0.135	-	<b>0.295</b>	-	0.132	-	0.264	-
18	No external k.	0.437	0.485	<b>0.638</b>	<b>0.660</b>	0.138	0.110	0.387	0.159	0.167	0.199	<b>0.400</b>	<b>0.465</b>
19	Commonsense	<b>0.788</b>	0.698	0.752	<b>0.815</b>	0.613	0.364	0.624	0.356	0.735	0.470	<b>0.778</b>	<b>0.526</b>
20	Math formulas	0.259	0.162	<b>0.661</b>	<b>0.544</b>	0.137	0.074	0.454	0.382	0.170	0.077	<b>0.599</b>	<b>0.404</b>
21	Science formulas	0.305	0.120	<b>0.315</b>	<b>0.250</b>	0.158	0.025	<b>0.239</b>	0.021	0.157	0.105	0.181	<b>0.250</b>
22	Computer science k.	0.262	0.128	<b>0.425</b>	<b>0.408</b>	0.151	0.137	0.147	0.134	<b>0.232</b>	<b>0.304</b>	0.220	0.278
23	Real-world k.	0.150	-	<b>0.472</b>	-	0.012	-	0.100	-	0.019	-	<b>0.160</b>	-
Average score		0.384	0.384	<b>0.604</b>	<b>0.586</b>	0.204	0.177	0.394	0.238	0.252	0.268	<b>0.480</b>	<b>0.448</b>

Table 3: Evaluations of different baselines across 23 tasks in LILA. On most tasks, **Codex** outperforms all baselines while **Neo-Multi-P** outperforms all fine-tuned baselines. A model usually performs worse on the OOD data set. The **bold** score refers to the best score among models with the *same supervision* method; the underlined score refers to the best score among *all* models. GPT3 and Codex performance is computed on 100 uniformly distributed examples owing to their cost and usage limit. Fine-tuned model performance is calculated on the full test set.

on various aspects of the dataset, we fine-tune *single-task* models on each of the 23 tasks in LILA. We also evaluate the benefit of transfer learning by fine-tuning a single *multi-task* GPT-Neo baseline on all the tasks simultaneously.

**Prompting.** We also use few-shot prompting to evaluate GPT3 (text-davinci-002) and Codex (code-davinci-002) (Brown et al., 2020; Chen et al., 2021). For the IID setting, we prompt the model with a random input-output examples from the same dataset as the input. In the OOD setting, we take examples from other datasets (Table 15-18) within the same task. We repeat this evaluation with increasing numbers of examples (up to the token size of models) to study the effect on performance<sup>4</sup>.

**Evaluation.** We evaluate models under two regimes – directly outputting the answer *i.e.* program induction and outputting a Python program

that is then executed to obtain the final answer *i.e.* program synthesis. In the case of our fine-tuned models, we train them to output both the final answer and the Python program conditioned on the input question. To evaluate our models under direct question answering, we use F1-score<sup>5</sup> to compare the model output and the gold answer. To evaluate program synthesis, we execute the model’s output within a python interpreter and compare the program output with the output of the gold program, again using F1. We evaluate based on the program output, rather than the program itself, to account for diversity in solving techniques and programming styles.

## 5 Results and Analysis

A summary results of all key results on our LILA benchmark are shown in Table 3. In this section, we will discuss the performance of fine-tuned 2.7B GPT-Neo models (§5.1), performance of models along the 4 categories of tasks (§5.2) and finally,

<sup>4</sup>henceforth max example model is the default model unless otherwise specified.

<sup>5</sup>This is a soft version of exact match accuracy assigning partial credit when common words are present in the output and gold answer.

the few-shot performance of much larger ( $\sim 175B$  parameters) models (§5.3).

## 5.1 Results: Fine-tuned Models

**Multitasking improves IID and OOD generalization.** The multi-tasking model (Neo-Multi) substantially improves upon the single task models (Neo). Neo-Multi achieves better average in-domain performance than the 23 individual per-task models (0.480 vs. 0.394 average score), suggesting that it leverages cross-task structure not present in a single task’s training set.

Multi-task training also substantially improves out-of-domain generalization (0.448 vs. 0.238). The gap between IID and OOD performance is much smaller for Neo-Multi than for the single task models, and Neo-Multi’s OOD performance is better than its IID performance on some tasks (see Table 4). LILA’s multi-task structure opens interesting future directions related to developing improved multitasking techniques, and further understanding its benefits.

**Program synthesis substantially outperforms answer prediction.** Synthesizing the program and evaluating it to get an answer substantially outperforms directly predicting the answer. For instance, multi-task program synthesis (Neo-Multi-P) has an average score of 0.480 while multi-task answer prediction (Neo-Multi-A) scores 0.252. This means models are often able to generate a program that evaluates to the correct answer, even when the model cannot directly compute the answer.

Program synthesis improves over answer prediction in all math categories except Geometry, with the largest improvements in Statistics and Linear Algebra; see Table 7 for examples. LILA’s unified problem format decouples synthesis from computation, while opening directions for further study on either aspect.

**Models leverage symbolic execution and libraries.** The gap between program synthesis and answer prediction suggests that the neural language model offloads computations to the symbolic Python runtime that are otherwise difficult to compute directly. We identify two common cases. First, the model leverages standard Python as a calculator. For instance, this pattern is common in the `basic_math` and `mul_div` categories, which involve evaluating arithmetic expressions; Table 6 shows examples. Second, the model is able to call

external libraries that perform sophisticated computations. For instance, in statistics the model uses `scipy.stats.entropy` or `np.linalg.det` in linear algebra while solving problems (Table 7).

**Models occasionally generate non-executable code.** Roughly 10% of Neo-Multi’s IID programs fail to execute. 86% of these are `SyntaxErrors`, which often occur because decoding terminates before finishing the program or the model generates a program of the form ‘`2+3=5`’, which is invalid Python. The remaining 14% of execution failures are less trivial, including `NameErrors` (7%) and `TypeError`s (1%) (see Table 8).

## 5.2 Results: Category-wise Analysis

In this section we discuss the trends among the tasks within each category. For brevity, we primarily consider the GPT-Neo multi-task model in the program-synthesis setting.

**Math ability.** Among the tasks in the math category, Neo-Multi excels in basic math, linear algebra, and in-domain statistics. On these tasks, it performs equal or better to Codex. On the other hand, Neo-Multi struggles in advanced math and geometry, with mediocre performance in multiplication-division, number theory, and calculus. Codex shows analogous trends, except for performing very well on calculus (0.930)<sup>6</sup>.

**Language Complexity .** Models generally show lower performance on program synthesis as language complexity increases. Fine-tuned GPT-Neo gets mean F1 over 0.5 only for datasets with the least linguistic complexity where it achieves an F1 of 0.7.

**Question Format.** Among the format tasks in the dataset, Neo-Multi does exceptionally well on multiple-choice and natural-language inference, getting performance close to 0.9 on the latter, and outperforming Codex on both. On the other hand, the model performs close to 0.25 for reading comprehension and fill-in-the-blank, though with 0.5 F1 on out-of-domain fill-in-the-blank.

**Background Knowledge.** Neo-Multi performs above 0.5 F1 only for problems requiring common-sense and math formulas and fails to do similarly on problems requiring other forms of external knowledge like physics, computer science, or real-world knowledge.

<sup>6</sup>Note that the training set for Codex is not known

### 5.3 Results: Few-shot Setup

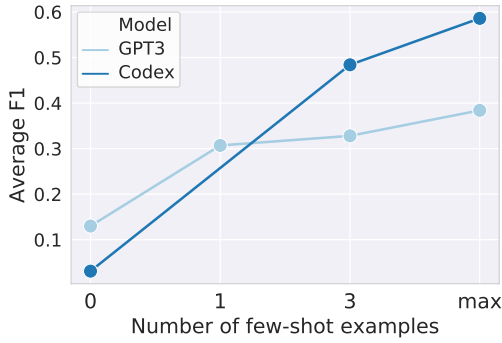


Figure 3: Average F1 scores of GPT3 and Codex with different numbers of few-shot examples in LILA.

Finally, we study the few-shot performance of much larger models ( $\approx 175\text{B}$ ), to better understand the performance of the smaller trained models ( $\approx 2.7\text{B}$ ) and to provide a benchmark for evaluating other large language models. Overall, we find that few-shot prompted models generally outperform their *much* smaller but fine-tuned counterparts.

**More examples and instructions improve prompts.** We find that the number of few-shot examples greatly impacts prompt models’ performance. Figure 3 shows that GPT-3 answer prediction beats Codex program synthesis in zero- to one-shot settings, but Codex overtakes with more examples. Table 5 shows that prompting with instructions also improves performance only in the zero-shot setting, meaning that in the limited contexts of the prompt models, examples are more important than instructions for mathematical reasoning. This is consistent with the findings of Puri et al. (2022) on instruction-example equivalence.

**Few-shot GPT-3 answer prediction underperforms fine-tuned GPT-Neo.** While prompt-based models outperform our fine-tuned models in general when comparing within direct-answering and program-synthesis, when comparing GPT-Neo program-synthesis to GPT3 direct answering we find that the much smaller fine-tuned GPT-Neo consistently outperforms GPT3.

**Few-shot Codex performance is relatively strong.** Relative to the 2.7B trained models, Codex demonstrates strong few-shot IID and OOD performance. Some notable exceptions to this pattern are the statistics, linear algebra, multiple-choice question answering, and NLI tasks. Generally, OOD few-shot performs much better than OOD for the fine-

Dimension	Neo-A		Neo-P	
	IID	OOD	IID	OOD
Math ability	0.191	0.129	<b>0.445</b>	<b>0.188</b>
Language	0.189	0.147	<b>0.429</b>	<b>0.246</b>
Format	0.246	0.382	<b>0.372</b>	<b>0.404</b>
Knowledge	0.206	0.143	<b>0.331</b>	<b>0.213</b>
Average	0.208	0.200	<b>0.394</b>	<b>0.263</b>

Table 4: Multi-task models are able to generalize across some task dimensions. Program output performs better than number output.

Dimension	Zero-shot		Few-shot (3)	
	w/o Inst.	w/ Inst.	w/o Inst.	w/ Inst.
Math ability	0.120	<b>0.123</b>	<b>0.311</b>	0.306
Language	0.124	<b>0.131</b>	<b>0.352</b>	0.350
Format	0.241	<b>0.257</b>	<b>0.555</b>	0.540
Knowledge	0.108	<b>0.112</b>	<b>0.367</b>	0.363
Average	0.148	<b>0.156</b>	<b>0.396</b>	0.390

Table 5: The IID scores for GPT3 models with and without instruction prompting (Inst.). Instruction helps slightly in zero-shot setting, but not in few-shot setting.

tuned models.

**Few-shot Codex fails some tasks.** Despite strong performance relative to fine-tuned GPT-Neo, Codex obtains less than 0.5 F1 on several tasks, with especially poor performance on geometry, number theory, advanced math, complex language, computer science problems, science formulas, and real world knowledge.

## 6 Conclusion

In this work, we introduce LILA, a unified mathematical reasoning benchmark for a holistic evaluation of AI agents. LILA consists of 23 tasks across 4 dimensions (i) mathematical abilities, (ii) language format, (iii) language complexity, (iv) external knowledge. It builds on 20 existing mathematical reasoning datasets to collect instructions and Python programs. Further, it also supports measuring out-of-distribution performance and robustness to language perturbations via LILA-OOD and LILA-ROBUST respectively. We also introduce BHASKARA, a 2.7B-parameter fine-tuned multi-task model, where we find that multi-tasking improves over single-task performance by 21.83% F1 score. The best performing model we evaluate achieves only 60.40% F1 indicating the potential for improvement on the proposed benchmark.



## 7 Limitations

One drawback of our unified format is the difficulty of evaluating models. In our work we use F1 for lack of a better alternative. F1 likely overestimates performance, e.g., given the gold answer “2 apples”, the predicted answers “2” and “apples” receive the same score, though the former is more correct.

LILA contains 23 tasks which are created from 20 datasets and 44 sub-datasets. There is scope to add more mathematical reasoning datasets (specifically datasets like theorem proving etc.). The flexible unified format of LILA allows for future extensions. Additionally, our categorization provides a way to identify areas for extension. For instance, we only have 1 dataset for linear algebra, which happens to not use natural language, and takes the form of generative QA. Our benchmark will benefit from future linear algebra additions, perhaps with word problems formatted as fill-in-the-blank questions.

## References

Armen Aghajanyan, Anchit Gupta, Akshat Shrivastava, Xilun Chen, Luke Zettlemoyer, and Sonal Gupta. 2021. Muppet: Massive multi-task representations with pre-finetuning. *arXiv preprint arXiv:2101.11038*.

Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. *arXiv preprint arXiv:1905.13319*.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.

Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella Biderman. 2021. [GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow](#). If you use this software, please cite it using these metadata.

Ronan Le Bras, Swabha Swayamdipta, Chandra Bhagavatula, Rowan Zellers, Matthew E Peters, Ashish Sabharwal, and Yejin Choi. 2020. Adversarial filters of dataset biases. *arXiv preprint arXiv:2002.04108*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child,

Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Henry T Colebrooke. 1817. Arithmetic and mensuration of brahmegupta and bhaskara.

Dheeru Dua, Ananth Gottumukkala, Alon Talmor, Sameer Singh, and Matt Gardner. 2019a. Orb: An open reading benchmark for comprehensive evaluation of machine reading comprehension. *arXiv preprint arXiv:1912.12598*.

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019b. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The Pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.

Jesse Michael Han, Jason M. Rute, Yuhuai Wu, Edward W. Ayers, and Stanislas Polu. 2021. Proof artifact co-training for theorem proving with language models. *ArXiv*, abs/2102.06203.

Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. 2021a. Measuring coding challenge competence with apps. *arXiv preprint arXiv:2105.09938*.

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

698	Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam	Nicholas Lourie, Ronan Le Bras, Chandra Bhagavatula,	755
699	Dziedziec, Rishabh Krishnan, and Dawn Song. 2020.	and Yejin Choi. 2021. Unicorn on rainbow: A uni-	756
700	Pretrained transformers improve out-of-distribution	versal commonsense reasoning model on a new mul-	757
701	robustness. In <i>Proceedings of the 58th Annual Meet-</i>	titask benchmark. In <i>Proceedings of the AAAI Con-</i>	758
702	<i>ing of the Association for Computational Linguistics</i> ,	<i>ference on Artificial Intelligence</i> , volume 35, pages	759
703	pages 2744–2751.	13480–13488.	760
704	Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren	Bryan McCann, Nitish Shirish Keskar, Caiming Xiong,	761
705	Etzioni, and Nate Kushman. 2014. Learning to solve	and Richard Socher. 2018. The natural language	762
706	arithmetic word problems with verb categorization.	decathlon: Multitask learning as question answering.	763
707	In <i>Conference on Empirical Methods in Natural</i>	<i>arXiv preprint arXiv:1806.08730</i> .	764
708	<i>Language Processing (EMNLP)</i> .		
709	Danqing Huang, Shuming Shi, Chin-Yew Lin, Jian Yin,	Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su.	765
710	and Wei-Ying Ma. 2016. How well do computers	2020a. <a href="#">A diverse corpus for evaluating and develop-</a>	766
711	solve math word problems? large-scale dataset con-	<a href="#">ing English math word problem solvers</a> . In <i>Proceed-</i>	767
712	struction and evaluation. In <i>Proceedings of the 54th</i>	<i>ings of the 58th Annual Meeting of the Association for</i>	768
713	<i>Annual Meeting of the Association for Computational</i>	<i>Computational Linguistics</i> , pages 975–984, Online.	769
714	<i>Linguistics (Volume 1: Long Papers)</i> , pages 887–896.	Association for Computational Linguistics.	770
715	Daniel Khashabi, Tushar Khot, Ashish Sabharwal,	Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su.	771
716	Oyvind Tafjord, Peter Clark, and Hannaneh Ha-	2020b. A diverse corpus for evaluating and develop-	772
717	ajishirzi. 2020. Unifiedqa: Crossing format bound-	ing english math word problem solvers. In <i>Proceed-</i>	773
718	aries with a single qa system. <i>arXiv preprint</i>	<i>ings of the 58th Annual Meeting of the Association</i>	774
719	<i>arXiv:2005.00700</i> .	<i>for Computational Linguistics</i> , pages 975–984.	775
720	Aditya Kolachana, K Mahesh, and K Ramasubramanian.	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and	776
721	2019. Use of calculus in hindu mathematics. In	Hannaneh Hajishirzi. 2022a. Cross-task generaliza-	777
722	<i>Studies in Indian Mathematics and Astronomy</i> , pages	tion via natural language crowdsourcing instructions.	778
723	345–355. Springer.	In <i>Proceedings of the 60th Annual Meeting of the</i>	779
724	Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish	<i>Association for Computational Linguistics (Volume</i>	780
725	Sabharwal, Oren Etzioni, and Siena Dumas Ang.	<i>1: Long Papers)</i> , pages 3470–3487.	781
726	2015. Parsing algebraic word problems into equa-	Swaroop Mishra, Arindam Mitra, Neeraj Varshney,	782
727	tions. <i>Transactions of the Association for Computa-</i>	Bhavdeep Sachdeva, Peter Clark, Chitta Baral, and	783
728	<i>tional Linguistics</i> , 3:585–597.	Ashwin Kalyan. 2022b. Numglue: A suite of funda-	784
729	Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate	mental yet challenging mathematical reasoning tasks.	785
730	Kushman, and Hannaneh Hajishirzi. 2016. Mawps:	In <i>Proceedings of the 60th Annual Meeting of the</i>	786
731	A math word problem repository. In <i>Proceedings of</i>	<i>Association for Computational Linguistics (Volume</i>	787
732	<i>the 2016 Conference of the North American Chap-</i>	<i>1: Long Papers)</i> , pages 3505–3523.	788
733	<i>ter of the Association for Computational Linguistics:</i>	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.	789
734	<i>Human Language Technologies</i> , pages 1152–1157.	2021. <a href="#">Are NLP models really able to solve simple</a>	790
735	Nate Kushman, Yoav Artzi, Luke Zettlemoyer, and	<a href="#">math word problems?</a> In <i>Proceedings of the 2021</i>	791
736	Regina Barzilay. 2014. Learning to automatically	<i>Conference of the North American Chapter of the</i>	792
737	solve algebra word problems. In <i>Proceedings of the</i>	<i>Association for Computational Linguistics: Human</i>	793
738	<i>52nd Annual Meeting of the Association for Computa-</i>	<i>Language Technologies</i> , pages 2080–2094, Online.	794
739	<i>tional Linguistics (Volume 1: Long Papers)</i> , pages	Association for Computational Linguistics.	795
740	271–281.		
741	Wenda Li, Lei Yu, Yuhuai Wu, and Lawrence C. Paulson.	Ravsehaj Singh Puri, Swaroop Mishra, Mihir Parmar,	796
742	2021. <a href="#">Isarstep: a benchmark for high-level mathe-</a>	and Chitta Baral. 2022. How many data samples	797
743	<a href="#">matical reasoning</a> . In <i>International Conference on</i>	is an additional instruction worth? <i>arXiv preprint</i>	798
744	<i>Learning Representations</i> .	<i>arXiv:2203.09161</i> .	799
745	Bill Yuchen Lin, Seyeon Lee, Rahul Khanna, and Xi-	Abhilasha Ravichander, Aakanksha Naik, Carolyn	800
746	ang Ren. 2020. Birds have four legs?! numersense:	Rose, and Eduard Hovy. 2019. Equate: A bench-	801
747	Probing numerical commonsense knowledge of pre-	mark evaluation framework for quantitative reason-	802
748	trained language models. In <i>Proceedings of the 2020</i>	ing in natural language inference. <i>arXiv preprint</i>	803
749	<i>Conference on Empirical Methods in Natural Lan-</i>	<i>arXiv:1901.03735</i> .	804
750	<i>guage Processing (EMNLP)</i> , pages 6862–6868.		
751	Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blun-	Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin,	805
752	som. 2017. Program induction by rationale genera-	and Sameer Singh. 2020. Beyond accuracy: Behav-	806
753	tion: Learning to solve and explain algebraic word	ioral testing of nlp models with checklist. In <i>Proceed-</i>	807
754	problems. <i>arXiv preprint arXiv:1705.04146</i> .	<i>ings of the 58th Annual Meeting of the Association</i>	808
		<i>for Computational Linguistics</i> , pages 4902–4912.	809

810	Anna Rogers, Matt Gardner, and Isabelle Augenstein.	Shyam Upadhyay, Ming-Wei Chang, Kai-Wei Chang,	862
811	2021. Qa dataset explosion: A taxonomy of nlp	and Wen-tau Yih. 2016. Learning from explicit and	863
812	resources for question answering and reading com-	implicit supervision jointly for algebra word prob-	864
813	prehension. <i>arXiv preprint arXiv:2107.12708</i> .	lems. In <i>Proceedings of the 2016 Conference on</i>	865
814	Subhro Roy and Dan Roth. 2015. Solving general arith-	<i>Empirical Methods in Natural Language Processing</i> ,	866
815	metic word problems. In <i>Proceedings of the 2015</i>	pages 297–306.	867
816	<i>Conference on Empirical Methods in Natural Lan-</i>	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-	868
817	<i>guage Processing</i> , pages 1743–1752.	preet Singh, Julian Michael, Felix Hill, Omer Levy,	869
818	Subhro Roy and Dan Roth. 2016. Solving gen-	and Samuel Bowman. 2019. Superglue: A stickier	870
819	eral arithmetic word problems. <i>arXiv preprint</i>	benchmark for general-purpose language understand-	871
820	<i>arXiv:1608.01413</i> .	ing systems. In <i>Advances in Neural Information</i>	872
821	Subhro Roy and Dan Roth. 2017. Unit dependency	<i>Processing Systems</i> , pages 3261–3275.	873
822	graph and its application to arithmetic word prob-	Alex Wang, Amanpreet Singh, Julian Michael, Felix	874
823	lem solving. In <i>Thirty-First AAAI Conference on</i>	Hill, Omer Levy, and Samuel R Bowman. 2018.	875
824	<i>Artificial Intelligence</i> .	Glue: A multi-task benchmark and analysis platform	876
825	Subhro Roy and Dan Roth. 2018. Mapping to declara-	for natural language understanding. <i>arXiv preprint</i>	877
826	tive knowledge for word problem solving. <i>Transac-</i>	<i>arXiv:1804.07461</i> .	878
827	<i>tions of the Association for Computational Linguis-</i>	Yizhong Wang, Swaroop Mishra, Pegah Alipoor-	879
828	<i>tics</i> , 6:159–172.	molabashi, Yeganeh Kordi, Amirreza Mirzaei,	880
829	Subhro Roy, Tim Vieira, and Dan Roth. 2015. Reason-	Anjana Arunkumar, Arjun Ashok, Arut Selvan	881
830	ing about quantities in natural language. <i>Transac-</i>	Dhanasekaran, Atharva Naik, David Stap, et al. 2022.	882
831	<i>tions of the Association for Computational Linguis-</i>	Benchmarking generalization via in-context instruc-	883
832	<i>tics</i> , 3:1–13.	tions on 1,600+ language tasks. <i>arXiv preprint</i>	884
833	Victor Sanh, Albert Webson, Colin Raffel, Stephen H	<i>arXiv:2204.07705</i> .	885
834	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin	886
835	Chaffin, Arnaud Stiegler, Teven Le Scao, Arun	Guu, Adams Wei Yu, Brian Lester, Nan Du, An-	887
836	Raja, et al. 2021. Multitask prompted training en-	drew M Dai, and Quoc V Le. 2021. Finetuned lan-	888
837	ables zero-shot task generalization. <i>arXiv preprint</i>	guage models are zero-shot learners. <i>arXiv preprint</i>	889
838	<i>arXiv:2110.08207</i> .	<i>arXiv:2109.01652</i> .	890
839	Benoy Kumar Sarkar. 1918. <i>Hindu Achievements in</i>	Sean Welleck, Jiacheng Liu, Ronan Le Bras, Hannaneh	891
840	<i>Exact Science: A Study in the History of Scientific</i>	Hajishirzi, Yejin Choi, and Kyunghyun Cho. 2021.	892
841	<i>Development</i> . Longmans, Green and Company.	<a href="#">Naturalproofs: Mathematical theorem proving in nat-</a>	893
842	David Saxton, Edward Grefenstette, Felix Hill, and	<a href="#">ural language</a> . In <i>Thirty-fifth Conference on Neural</i>	894
843	Pushmeet Kohli. 2019. Analysing mathematical rea-	<i>Information Processing Systems Datasets and Bench-</i>	895
844	soning abilities of neural models. <i>arXiv preprint</i>	<i>marks Track (Round 1)</i> .	896
845	<i>arXiv:1904.01557</i> .	Sean Welleck, Peter West, Jize Cao, and Yejin Choi.	897
846	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao,	2022. <a href="#">Symbolic brittleness in sequence models: on</a>	898
847	Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch,	<a href="#">systematic generalization in symbolic mathematics</a> .	899
848	Adam R Brown, Adam Santoro, Aditya Gupta,	In <i>AAAI</i> .	900
849	Adrià Garriga-Alonso, et al. 2022. Beyond the	Jason Weston, Antoine Bordes, Sumit Chopra, Alexan-	901
850	imitation game: Quantifying and extrapolating the	der M Rush, Bart van Merriënboer, Armand Joulin,	902
851	capabilities of language models. <i>arXiv preprint</i>	and Tomas Mikolov. 2015. Towards ai-complete	903
852	<i>arXiv:2206.04615</i> .	question answering: A set of prerequisite toy tasks.	904
853	Oyvind Tafjord, Peter Clark, Matt Gardner, Wen-tau	<i>arXiv preprint arXiv:1502.05698</i> .	905
854	Yih, and Ashish Sabharwal. 2019. Quarel: A dataset	Yuhuai Wu, Albert Jiang, Jimmy Ba, and Roger Baker	906
855	and models for answering questions about qualitative	Grosse. 2021. <a href="#">{INT}: An inequality benchmark for</a>	907
856	relationships. In <i>Proceedings of the AAAI Conference</i>	<a href="#">evaluating generalization in theorem proving</a> . In <i>In-</i>	908
857	<i>on Artificial Intelligence</i> , volume 33, pages 7063–	<i>ternational Conference on Learning Representations</i> .	909
858	7071.	Pengcheng Yin, Bowen Deng, Edgar Chen, Bogdan	910
859	Shyam Upadhyay and Ming-Wei Chang. 2015. Draw:	Vasilescu, and Graham Neubig. 2018. <a href="#">Learning to</a>	911
860	A challenging and diverse algebra word problem set.	<a href="#">mine aligned code and natural language pairs from</a>	912
861	Technical report, Citeseer.	<a href="#">stack overflow</a> . In <i>International Conference on Min-</i>	913
		<i>ing Software Repositories</i> , MSR, pages 476–486.	914
		ACM.	915

- Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. 2018. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3712–3722.
- Xikun Zhang, Deepak Ramachandran, Ian Tenney, Yanai Elazar, and Dan Roth. Do language embeddings capture scales?
- Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. 2021. Minif2f: a cross-system benchmark for formal olympiad-level mathematics. *arXiv preprint arXiv:2109.00110*.
- Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. “going on a vacation” takes longer than “going for a walk”: A study of temporal commonsense understanding. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3363–3369.



## A Appendix

### A.1 Qualitative Examples

Figures 6 and 7 give examples of input-output behavior of GPT-Neo fine-tuned on all the tasks. Figure 8 gives an example of a non-compiling output program.

### A.2 Dataset Collection

#### A.2.1 Expert annotation

In the worker qualification process, we ask each worker to annotate 30 questions. We manually verify each annotation and qualify those whose python annotations are satisfactory. We also provide feedback such as "write simpler programs, use representative variable names instead of just letters, add comments wherever possible" to annotators after the worker qualification process. We instruct annotators to use a minimal set of python libraries, and we ask them to record the Python libraries they use in a common document. We find that the annotators could get the task done just by using the sympy and the datetime libraries. We also ask annotators to report any bugs in answer annotation, which they report for a small number of questions; we subsequently fix those.

We give 10 sample question annotations to annotators as illustrative examples which vary in structure, length, format, underlying reasoning skill, etc. We pay 20 dollars per hour up to 20 hours per week as compensation for the data annotation work.

**LILA-ROBUST** To create the LILA-ROBUST dataset, we first define a set of 9 templates, consisting of 3 variation styles defined in SVAMP (Patel et al., 2021) as well as 6 novel templates of our own. We refer to the SVAMP templates as SVAMP-COO, SVAMP-COP, and SVAMP-IU, which correspond to changing the order of objects, changing the order of phrases, and adding irrelevant, unhelpful information to the problem statement, respectively. Our novel templates are named ROBUST-IR, ROBUST-AP, ROBUST-ADJ, ROBUST-Q, ROBUST-RQ, and ROBUST-RM. ROBUST-IR refers to adding information that is unhelpful for solving the question but may be related to the context of the problem. ROBUST-AP refers to increasing problem verbosity by turning active speech to passive speech. ROBUST-ADJ refers to increasing problem verbosity by adding adjectives or adverbs. ROBUST-Q indicates turning a problem statement into a question, in the style of a conversation with a stu-

dent. ROBUST-RQ indicates removing question words in a problem and turning it into a statement; it is roughly the inverse of ROBUST-Q. Finally, ROBUST-RM refers to the removal of mathematics terms that are implicitly defined. Examples of each template are found in Table 10.

For our crowdsourcing pipeline, we provide each Amazon Mechanical Turk worker with 10 questions split from 20 questions sampled from each dataset. We run a separate job for each of our 9 templates. In particular, each HIT contains the 10 split questions from the original datasets, alongside the problem solution. Workers are asked to submit a augmentation for each question according to the style of the template assigned to each job. Thus, we run 9 separate jobs to obtain augmentations for all templates across all datasets. To familiarize workers with the intended style of each template, we provide 3 demonstrative augmentations within the instructions of each HIT, as summarized in Table 10. We restrict our crowdsourcing pipeline to workers that had above a 98% acceptance rate with over 1000 completed HITs. We provide workers with an upper bound of 1 hour to complete each HIT but specify in the instructions that each HIT should feasible be completed in 10 minutes. Based on minimum wage policies and under the assumption that workers follow the 10-minute completion guideline, we accordingly compensate \$3 per HIT. Finally, to ensure dataset quality, we manually assess the worker augmentations produced for each template.

### A.3 Dataset Statistics

Figure 5 gives relatives sizes of tasks within each category. Figure 6 illustrates the unigram frequencies in LILA, where larger words indicate higher frequency. Table 19 gives comprehensive statistics on each task. Table 21 cites each component dataset of LILA.

### A.4 Additional Results

Table 20 gives the unaggregated performance of each model on each dataset in LILA (some datasets are split across tasks).

<b>Task Basic Math</b>	
<b>Problem</b> Before December, customers buy 1346 ear muffs from the mall. During December, they buy 6444, and there are none. In all, how many ear muffs do the customers buy?	
<b>Answer</b>	7790.0
<b>Predicted Answer</b>	1346.0 ✗
<b>Generated Program</b>	
<pre>answer = 1346.0 + 6444.0 print(answer)</pre>	
<b>Executed Answer</b>	7790.0 ✓
<b>Task Muldiv</b>	
<b>Problem</b> Tickets to the school play cost 6 for students and 8 for adults. If 20 students and 12 adults bought tickets, how many dollars' worth of tickets were sold?	
<b>Answer</b>	216
<b>Predicted Answer</b>	48 ✗
<b>Generated Program</b>	
<pre>a=20*6 b=12*8 c=a+b answer= c print(answer)</pre>	
<b>Executed Answer</b>	216 ✓

Table 6: Examples using Neo-Multi.

Category	Math ability	Language	Knowledge	Format
Tasks	basic math, multiplication-division, number theory, algebra, geometry, counting and statistics, calculus, linear algebra, advanced math	no language, simple language, complex language	no background knowledge, commonsense, math, science, computer science, real world knowledge	fill-in-the-blank, generative question answer- ing, multiple-choice, natural language inference, reading comprehension

Table 9: Dataset partitions (Categories) and associated tasks.

<b>Task Statistics</b>	
<b>Problem</b> Determine the KL Divergence between the empirical distributions 0.149, 0.039, 0.099, 0.09, 0.203, 0.077, 0.039, 0.19 and 0.158, 0.256, 0.042, 0.008, 0.064, 0.051, 0.05, 0.072.	
<b>Answer</b>	0.5122...
<b>Predicted Answer</b>	0.36 ✗
<b>Generated Program</b>	
<pre>from scipy.stats import entropy  distribution1 = [0.149, 0.039, 0.099, 0.09, 0.203, 0.077, 0.039, 0.19] distribution2 = [0.158, 0.256, 0.042, 0.008, 0.064, 0.051, 0.05, 0.072] print(entropy(distribution1, qk=distribution2))</pre>	
<b>Executed Answer</b>	0.5122... ✓
<b>Task Linear Algebra</b>	
<b>Problem</b> Find the determinant of the matrix $\begin{pmatrix} 0 & -2 & -3 \\ 0 & 5 & 0 \\ 1 & 3 & 2 \end{pmatrix}$ .	
<b>Answer</b>	15
<b>Predicted Answer</b>	-8 ✗
<b>Generated Program</b>	
<pre>import numpy as np  a = np.array([     [0, -2, -3],     [0, 5, 0],     [1, 3, 2]]) print(np.linalg.det(a))</pre>	
<b>Executed Answer</b>	15.0 ✓

Table 7: Examples using Neo-Multi.

---

**Model** Neo-Multi

**Task** Muldiv

**Problem** Jenny collects cans and bottles to take down to the recycling center. Each bottle weighs 6 ounces and each can weighs 2 ounces. Jenny can carry a total of 100 ounces. She collects 20 cans and as many bottles as she can carry. If she gets paid 10 cents per bottle and 3 cents per can, how much money does she make (in cents)?

**Generated Program**

```
a=20*6
b=a*2
c=b*3
d=c*10
e=d*3
f=e*3
g=f+g
answer= g
print(answer)
```

**Error:** NameError (*g* is not defined)

---

**Model** Codex

**Task** Advanced Math

**Problem** Simplify the expression  $(9x^2 + 3x + 7) + (3x^2 + 7x^5 + 2)$ . Express your answer as a polynomial with terms arranged in decreasing order of degree.

**Generated Program**

```
from sympy import Poly

p = Poly(9*(x**2) + 3*x + 7 + 3*(x**2) + 7*(x**5) + 2)

answer = p.as_expr()

print(answer)
```

**Error:** NameError (*x* is not defined)

---

Table 8: NameErrors in Neo-Multi and Codex.



**Question:** A gardener is going to plant 2 red rosebushes and 2 white rosebushes. If the gardener is to select each of the bushes at random, one at a time, and plant them in a row, what is the probability that the 2 rosebushes in the middle of the row will be the red rosebushes?

**Options:** {A:1/12, B:1/6, C:1/5, D:1/3, E:1/2}

**Answer:** B

**Explanation:** We are asked to find the probability of one particular pattern: wrrw. Total # of ways a gardener can plant these four bushes is the # of permutations of 4 letters wwrr, out of which 2 w's and 2 r's are identical, so  $4! / 2! 2! = 6$ ; so  $p = 1 / 6$ . Answer: B.

**Program:**

```
import scipy
n0 = 2.0
n1 = 2.0
n2 = 2.0
t0 = n0 + n0
t1 = scipy.special.comb(t0, n0)
answer = 1.0 / t1
```

Figure 4: An example of instruction annotation.

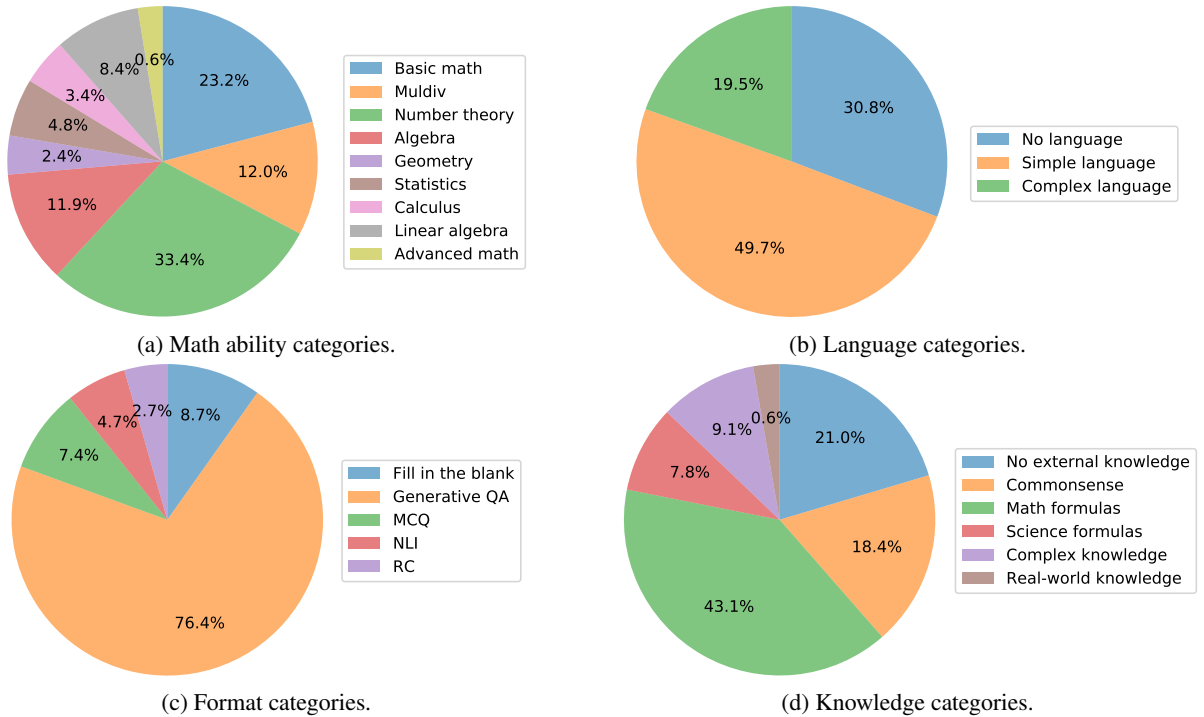


Figure 5: Task diversity in LILA across math, language, format, and knowledge categories.

Template Name	Variation	Example
SVAMP-COO	Change the order of objects	<p><b>Question:</b> Allen bought 20 stamps at the post office in 37 cents and 20 cents denominations . If the total cost of the stamps was \$ 7.06 , how many 37 cents stamps did Allen buy ?</p> <p><b>Variation:</b> Allen bought 20 stamps at the post office in 20 cents and 37 cents denominations . If the total cost of the stamps was \$ 7.06 , how many 37 cents stamps did Allen buy ?</p>
SVAMP-COP	Change the order of phrases	<p><b>Question:</b> One pipe can fill a tank in 5 hours and another pipe can fill the same tank in 4 hours . A drainpipe can empty the full content of the tank in 20 hours . With all the three pipes open , how long will it take to fill the tank ?</p> <p><b>Variation:</b> A drainpipe can empty the full content of a tank in 20 hours . One pipe can fill the tank in 4 hours and another pipe can fill the same tank in 5 hours . With all the three pipes open , how long will it take to fill the tank with all the three pipes open ?</p>
SVAMP-IU	Add irrelevant, unhelpful information	<p><b>Question:</b> the area of an isosceles trapezoid with sides of length 5 and bases of length 7 and 13 is ?</p> <p><b>Variation:</b> monkeys and apes are both primates, which means they're both part of the human family tree . the area of an isosceles trapezoid with sides of length 5 and bases of length 7 and 13 is ?</p>
ROBUST-IR	Add unhelpful, but contextually related information	<p><b>Question:</b> Tom is 15 years younger than alice . Ten years ago , Alice was 4 times as old as Tom was then . How old is each now ?</p> <p><b>Variation:</b> Tom is 15 years younger than alice . Ten years ago , Alice was 4 times as old as Tom was then . Alice really likes pineapple pizza. How old is each now ?</p>
ROBUST-AP	Turn active into passive speech to increase problem verbosity	<p><b>Question:</b> Hay's Linens sells hand towels in sets of 17 and bath towels in sets of 6. If the store sold the same number of each this morning, what is the smallest number of each type of towel that the store must have sold?</p> <p><b>Variation:</b> Hand towels are sold by Hay's Linens in sets of 17 and bath towels are sold in sets of 6. If the same number of each were sold by the store this morning, what is the smallest number of each type of towel that the store must have sold?</p>
ROBUST-ADJ	Add adjectives and adverbs to increase problem verbosity	<p><b>Question:</b> ThereTea leaves exposed to oxygen for up to _ hours become black tea.</p> <p><b>Variation:</b> Black tea leaves continuously exposed to oxygen for up to _ hours become a very rich black tea.</p>
ROBUST-Q	Turn a task statement into a question	<p><b>Question:</b> Product of -7 and -1469.125.</p> <p><b>Variation:</b> What is the product of -7 and -1469.125?</p>
ROBUST-RQ	Turn a question into a task statement	<p><b>Question:</b> Problem: If the product of 5 and a number is increased by 4 , the result is 19. What is the number?</p> <p><b>Variation:</b> Increasing the product of 5 and a number by 4 results is 19. Find the number.</p>
ROBUST-RM	Remove explicitly mathematical terms that are implicitly defined	<p><b>Problem:</b> Find the arclength of the function <math>f(x) = 2\sqrt{x}</math> on the interval <math>x = 2</math> to <math>x = 8</math></p> <p><b>Variation:</b> Find the arclength of <math>f(x) = 2\sqrt{x}</math> on <math>[2, 8]</math></p>

Table 10: Example for each template provided to MTurk workers to produce LILA-ROBUST

Task	Question category	Example
TASK 1	Basic math: addition, subtraction, fact based QA etc.	<b>Original:</b> If Jimbo is 484 feet away from a beetle and quarter of 827 feet away from a grasshopper, which insect will seem bigger to him?? "Option 1": beetle, "Option 2" :grasshopper <b>Answer:</b> Option 2
TASK 2	Muldiv: multiplication, division along with addition, subtraction etc.	<b>Question:</b> Mrs. Hilt bought 2 pizzas. Each pizza had 8 slices. So, she had __ total slices of pizza. <b>Answer:</b> 16
TASK 3	Number theory: prime, power, negation, modulus and other operators etc.	<b>Question:</b> How many numbers are divisible by both 2 and 3 up to 300? <b>Answer:</b> 50
TASK 4	Algebra: equations, functions, polynomials, series etc.	<b>Question:</b> The sum of the three smallest of four consecutive integers is 30 more than the largest integer. What are the four consecutive integers ? <b>Answer:</b> 15.0
TASK 5	Geometry: triangles, polygons, 3D structures etc.	<b>Question:</b> A hall is 6 meters long and 6 meters wide. If the sum of the areas of the floor and the ceiling is equal to the sum of the areas of four walls, what is the volume of the hall (in cubic meters)? <b>Answer:</b> 108
TASK 6	Statistics: binomial, divergence, mean, median, mode, variance etc.	<b>Question:</b> There are 11 boys and 10 girls in a class. If three students are selected at random, in how many ways that 3 girl and 2 boys are selected? <b>Answer:</b> 6600
TASK 7	Calculus: differentiation, integration, gradient, series expansion etc.	<b>Question:</b> Let $g(y) = 9*y^{**4} + 25*y^{**2} + 6$ . Let $s(d) = 1 - d^{**4}$ . Let $x(t) = -g(t) + 6*s(t)$ . What is the third derivative of $x(f)$ wrt $f$ ? <b>Answer:</b> -360*f
TASK 8	Linear algebra: vectors, dot products, Eigen vectors, matrices etc.	<b>Question:</b> Problem: Convert the following matrix to reduced row echelon form: $\begin{pmatrix} 7 & -2 & -10 & -4 \\ -5 & -10 & 2 & -7 \end{pmatrix}$ . <b>Answer:</b> $\begin{pmatrix} 1 & 0 & -\frac{13}{10} & -\frac{13}{40} \\ 0 & 1 & \frac{9}{20} & \frac{69}{80} \end{pmatrix}$
TASK 9	Advanced math: heuristics required along with probability, statistics, or algebra, Olympiad level problems	<b>Question:</b> Let $f(x) = 2^x$ . Find $\sqrt{f(f(f(f(1))))}$ . <b>Answer:</b> 256

Table 11: Example of each task in the *math ability* category of the LILA benchmark.

Task	Question category	Example
TASK 10	No language	Compute the median of $4\sqrt{2}$ , $-6$ , $3e$ , $3$ , $-6$ , $-\frac{14}{\sqrt{\pi}}$ , $6$ . <b>Answer:</b> 3
TASK 11	Simple language	<b>Question:</b> Joan had 9 blue balloons, but Sally popped 5 of them. Jessica has 2 blue balloons. They have __ blue balloons now. <b>Answer:</b> 6
TASK 12	Complex language: involving co-reference resolution etc., multi-sentence language, adversarial language: containing tricky words etc., often created adversarially	<b>Question:</b> Passage: According to the 2011 National Household Survey, 89.3% of Markhams residents are Canadian citizens, and about 14.5% of residents are recent immigrants (from 2001 to 2011). The racial make up of Markham is; East Asian (39.7%), White Canadian (27.5%), South Asian Canadian (19.1%), Southeast Asian (3.9%), Black Canadians (3.2%), West Asian & Arab Canadians (3.2%), Latin American Canadian (0.5%), Aboriginal peoples in Canada (0.2%), and 1.9% of the population is multiracial while the rest of the population (0.7%) is of another group. Markham has the highest visible minority population of any major Canadian city (over 100,000 residents) at 72.3%, and is one of eight major cities with no majority racial group. Question: How many percent of people were not white? <b>Answer:</b> 72.5

Table 12: Example of each task in the *language complexity* category of the LILA benchmark.

Task	Question category	Example
TASK 13	Fill in the blank	<b>Question:</b> Delphinium has _ florets or they are full of holes. <b>Answer:</b> no
TASK 14	Generative question answering	<b>Question:</b> Calculate the remainder when 160 is divided by 125. <b>Answer:</b> 35
TASK 15	Multiple choice question answering (MCQ)	<b>Question:</b> The fish glided with a speed of 8 m/s through the water and 5 m/s through the jello because the ___ is smoother? "Option 1": jello, "Option 2": water. <b>Answer:</b> Option 2
TASK 16	Natural language inference (NLI)	<b>Question:</b> "statement 1": Alyssa picked 42.0 pears from the pear tree and Nancy sold 17.0 of the pears , "statement 2" :25.0 pears were left , "options: " Entailment or contradiction? <b>Answer:</b> Entailment
TASK 17	Reading comprehension (RC)	<b>Question:</b> Passage: A late game rally by Washington led them to the Eagles' 26 yard line. A shot to the end zone by Robert Griffin III would be intercepted by Brandon Boykin, clinching an Eagles win. The Eagles would move to 6-5. This is the Eagles first win at Lincoln Financial Field since Week 4 of the 2012 season, because prior to this game, the Eagles had never won a game in their home stadium in 414 days since that same week, snapping a 10-game losing streak at home with this win. Question: How many more wins than losses did the Eagles have after this game? <b>Answer:</b> 1

Table 13: Example of each task in the *question format* category of the LILA benchmark.

Task	Question category	Example
TASK 18	No external knowledge: only mathematical commonsense knowledge required	<b>Question:</b> If there are 7 bottle caps in a box and Linda puts 7 more bottle caps inside, how many bottle caps are in the box? <b>Answer:</b> 14
TASK 19	Commonsense: temporal commonsense knowledge (e.g. people usually play basketball for a few hours and not days), numerical commonsense knowledge (e.g. birds has 2 legs)	<b>Question:</b> Outside temple, there is a shop which charges 12 dollars for each object. Please note that one shoe is counted as an object. Same is true for socks and mobiles. Paisley went to temple with both parents. All of them kept their shoes, socks and mobiles in the shop. How much they have to pay? <b>Answer:</b> 180
TASK 20	Math formulas: algebra, geometry, probability etc.	<b>Question:</b> Simplify $-3*(\sqrt{1700}) - (\sqrt{1700}) + (3 + \sqrt{1700})*(-6) + -3$ . <b>Answer:</b> $-180*\sqrt{17} - 57$
TASK 21	Science formulas: physics, chemistry etc.	<b>Question:</b> Find the number of moles of H <sub>2</sub> O formed on combining 2 moles of NaOH and 2 moles of HCl. <b>Answer:</b> 2
TASK 22	Computer science knowledge: data structure algorithms like merge sort etc.	<b>Question:</b> Apply functions 'mean' and 'std' to each column in dataframe 'df' <b>Answer:</b> <code>df.groupby(lambda idx: 0).agg(['mean', 'std'])</code>
TASK 23	Real-world knowledge: COVID modelling, climate modelling etc.	<b>Question:</b> Our physics club has 20 members, among which we have 3 officers: President, Vice President, and Treasurer. However, one member, Alex, hates another member, Bob. How many ways can we fill the offices if Alex refuses to serve as an officer if Bob is also an officer? (No person is allowed to hold more than one office.) <b>Answer:</b> 6732

Table 14: Example of each task in the *background knowledge* category of the LILA benchmark.





ID	Language category	IID	OOD
TASK 10	No language	amps_number_theory.json	amps_algebra.json
		amps_counting_and_stats.json	deepmind_mathematics_calculus.json
		amps_calculus.json	
		amps_linear_algebra.json	
		deepmind_mathematics_muldiv.json	
		deepmind_mathematics_numbertheory.json	
		deepmind_mathematics_algebra.json	
TASK 11	Simple language	deepmind_mathematics_basicmath.json	
		addsub.json	MCTaco_frequency_structured.json
		Numersense_structured.json	NumGLUE_Task1.json
		MCTaco_stationarity_structured.json	mathqa_general.json
		MCTaco_event_typical_time_structured.json	NumGLUE_Task4.json
		MCTaco_event_ordering_structured.json	
		MCTaco_event_duration_structured.json	
		singleop.json	
		multiarith.json	
		asdiv.json	
		GSM8k_structured.json	
		APPS_structured.json	
		mathqa_gain.json	
		mathqa_other.json	
		singleq.json	
		simuleq.json	
		NumGLUE_Task8.json	
		draw_structured.json	
		dolphin_structured.json	
		mathqa_probability.json	
TASK 12	Complex language	mathqa_physics.json	mbpp_structured.json
		APPS_structured.json	mathqa_other.json
		mathqa_gain.json	
		amps_number_theory.json	
		mathqa_general.json	
		conala_structured.json	
		NumGLUE_Task5.json	
		deepmind_mathematics_numbertheory.json	

Table 16: Raw datasets used to create different tasks in LILA across different language categories.

ID	Format category	IID	OOD
TASK 13	Fill in the blank	NumGLUE_Task4.json	NumerSense_structured.json
TASK 14	Generative QA	amps_number_theory.json	svamp_structured.json
		amps_counting_and_stats.json	mathqa_geometry.json
		amps_linear_algebra.json	amps_calculus.json
		amps_algebra.json	singleq.json
		deepmind_mathematics_calculus.json	NumGLUE_Task2.json
		addsub.json	mbpp_structured.json
		singleop.json	deepmind_mathematics_numbertheory.json
		multiarith.json	
		asdiv.json	
		GSM8k_structured.json	
		APPS_structured.json	
		mathqa_gain.json	
		mathqa_other.json	
		simuleq.json	
		NumGLUE_Task8.json	
		draw_structured.json	
		dolphin_structured.json	
		mathqa_probability.json	
		MCTaco_frequency_structured.json	
		NumGLUE_Task1.json	
		mathqa_general.json	
		mathqa_physics.json	
		conala_structured.json	
		amps_geometry.json	
		MATH_crowdsourced.json	
		deepmind_mathematics_calculus.json	
		deepmind_mathematics_muldiv.json	
		deepmind_mathematics_algebra.json	
		deepmind_mathematics_basicmath.json	
TASK 15	MCQ	NumGLUE_Task3.json	MCTaco_event_typical_time_structured.json
		MCTaco_stationarity_structured.json	
		MCTaco_event_ordering_structured.json	
		MCTaco_event_duration_structured.json	
TASK 16	NLI	NumGLUE_Task5.json	
TASK 17	RC	mathqa_physics.json	mbpp_structured.json

Table 17: Raw datasets used to create different tasks in LILA across different format categories.

ID	Knowledge category	IID	OOD
TASK 18	No external knowledge	addsub.json	NumGLUE_Task4.json
		singleop.json	GSM8k_structured.json
		multiarith.json	svamp_structured.json
		asdiv.json	NumGLUE_Task7.json
		simuleq.json	
		NumGLUE_Task8.json	
		draw_structured.json	
		dolphin_structured.json	
TASK 19	Commonsense	NumGLUE_Task5.json	
		deepmind_mathematics_muldiv.json	
		Numsense_structured.json	NumGLUE_Task1.json
		MCTaco_frequency_structured.json	MCTaco_event_ordering_structured.json
		NumGLUE_Task3.json	
		MCTaco_stationarity_structured.json	
TASK 20	Math formulas	MCTaco_event_duration_structured.json	
		MCTaco_event_typical_time_structured.json	
		amps_number_theory.json	amps_counting_and_stats.json
		amps_linear_algebra.json	mathqa_general.json
		amps_algebra.json	amps_calculus.json
		deepmind_mathematics_calculus.json	
		mathqa_probability.json	
		singleq.json	
		mathqa_gain.json	
		mathqa_other.json	
TASK 21	Science formulas	deepmind_mathematics_algebra.json	
		deepmind_mathematics_basicmath.json	
		deepmind_mathematics_calculus.json	
		deepmind_mathematics_numbertheory.json	
TASK 22	Computer science knowledge	amps_geometry.json	
		NumGLUE_Task2.json	
		mathqa_physics.json	
TASK 23	Real-world knowledge	APPS_structured.json	mathqa_geometry.json
		conala_structured.json	
TASK 23	Real-world knowledge	MATH_crowdsourced.json	mbpp_structured.json

Table 18: Raw datasets used to create different tasks in LILA across different knowledge categories.

ID	Category	Questions	Unique questions	Question length	Programs	Unique programs	Program length
TASK 1	Basic math	31,052	31,032	43.1	31,052	7,066	13.3
TASK 2	Muldiv	16,021	15,936	26.9	16,021	15,279	8.2
TASK 3	Number theory	44,760	44,183	41.3	269,232	261,865	33.2
TASK 4	Algebra	15,882	15,615	19.3	16,364	15,986	12.7
TASK 5	Geometry	3,190	3,149	36.1	3,190	3,035	28.7
TASK 6	Counting and statistics	6,423	6,384	39.7	6,423	6,335	31.5
TASK 7	Calculus	4,493	4,202	21.2	4,493	4,170	40.6
TASK 8	Linear algebra	11,248	11,204	32.4	11,248	11,204	23.0
TASK 9	Advanced math	746	746	21.2	746	745	27.3
TASK 10	No language	41,191	40,551	21.2	42,466	41,794	40.6
TASK 11	Simple language	66,505	66,172	26.9	290,184	258,839	8.2
TASK 12	Complex language	26,119	25,728	36.1	26,119	25,052	28.7
TASK 13	Fill in the blank	11,634	11,615	11.0	11,634	997	3.0
TASK 14	Generative QA	102,493	101,239	14.7	327,447	314,652	16.0
TASK 15	MCQ	9,989	9,989	28.3	9,989	470	3.0
TASK 16	NLI	6,326	6,325	50.8	6,326	6,243	25.8
TASK 17	RC	3,642	3,552	182.5	3,642	3,592	10.4
TASK 18	No external knowledge	28,115	27,964	50.8	28,115	27,117	25.8
TASK 19	Commonsense	24,677	24,658	30.9	24,677	823	3.0
TASK 20	Math formulas	57,841	56,947	19.1	59,116	57,019	25.5
TASK 21	Science formulas	10,505	10,319	36.1	10,505	9,764	28.7
TASK 22	Complex knowledge	12,200	12,086	14.5	235,879	230,486	24.2
TASK 23	Real-world knowledge	746	746	21.2	746	745	27.3

Table 19: Main statistics of LILA across the total of 23 tasks.



ID	Dataset	GPT3	Neo-A	Neo-P	Codex
1	addsub	0.910	0.116	0.797	<b>0.950</b>
2	amps_algebra	0.116	0.100	<b>0.902</b>	0.655
3	amps_calculus	0.192	0.168	<b>0.922</b>	0.860
4	amps_counting_and_stats	0.183	0.117	<b>0.958</b>	0.650
5	amps_geometry	<b>0.283</b>	0.263	0.074	0.000
6	amps_linear_algebra	0.127	0.235	<b>0.815</b>	0.692
7	amps_number_theory	0.273	0.026	0.875	<b>1.000</b>
8	APPS_structured	0.167	0.154	0.134	<b>0.459</b>
9	asdiv	<b>0.737</b>	0.166	0.092	0.022
10	conala_structured	0.356	0.329	0.329	<b>0.391</b>
11	deepmind_mathematics_algebra	0.202	0.258	0.847	<b>0.910</b>
12	deepmind_mathematics_basicmath	0.270	0.125	0.614	<b>1.000</b>
13	deepmind_mathematics_calculus	0.208	0.026	0.152	<b>0.884</b>
14	deepmind_mathematics_muldiv	0.160	0.034	0.909	<b>1.000</b>
15	deepmind_mathematics_numbertheory	0.296	0.462	0.538	<b>0.710</b>
16	dolphin_t2_final	0.170	0.027	0.006	<b>0.812</b>
17	draw_structured	0.090	0.034	0.005	<b>0.210</b>
18	GSM8k_structured	0.110	0.060	0.139	<b>0.350</b>
19	MATH_crowdsourced	0.150	0.013	0.074	<b>0.472</b>
20	mathqa_gain	0.134	0.054	<b>0.339</b>	0.270
21	mathqa_general	0.110	0.073	<b>0.193</b>	0.120
22	mathqa_geometry	0.120	0.002	0.000	<b>0.250</b>
23	mathqa_other	0.180	0.043	0.011	<b>0.280</b>
24	mathqa_physics	0.120	0.087	<b>0.429</b>	0.210
25	mathqa_probability	<b>0.210</b>	0.003	0.000	0.200
26	mbpp_structured	0.128	0.175	0.164	<b>0.408</b>
27	MCTaco_event_duration_structured	<b>0.800</b>	0.773	0.773	0.710
28	MCTaco_event_ordering_structured	0.860	0.831	0.831	<b>0.890</b>
29	MCTaco_event_typical_time_structured	0.870	<b>0.881</b>	<b>0.881</b>	0.870
30	MCTaco_frequency_structured	<b>0.890</b>	0.862	0.862	0.790
31	MCTaco_stationarity_structured	0.710	<b>0.758</b>	<b>0.758</b>	0.670
32	multiarith	0.360	0.143	0.921	<b>0.990</b>
33	NumerSense_structured	0.620	0.495	0.495	<b>0.660</b>
34	NumGLUE_Type_1	0.535	0.108	0.083	<b>0.740</b>
35	NumGLUE_Type_2	0.512	0.285	0.646	<b>0.735</b>
36	NumGLUE_Type_3	<b>0.835</b>	0.004	0.001	0.815
37	NumGLUE_Type_4	0.710	0.076	0.208	<b>0.790</b>
38	NumGLUE_Type_5	0.460	0.200	0.305	<b>0.615</b>
39	NumGLUE_Type_7	0.500	0.516	<b>0.854</b>	0.710
40	NumGLUE_Type_8	0.420	0.082	0.257	<b>0.610</b>
41	simuleq	0.120	0.074	0.010	<b>0.170</b>
42	singleop	0.940	0.347	0.611	<b>1.000</b>
43	singleq	<b>0.830</b>	0.143	0.474	0.670
44	svamp_structured	0.620	0.085	0.060	<b>0.790</b>
Average F1 score		0.400	0.223	0.440	<b>0.613</b>

Table 20: Evaluation results of baselines across different single datasets. On most datasets, **Codex** performs best. Model names: **GPT3**: the few-shot 175B GPT-3 model; **GPT-Neo-A**: the fine-tuned 2.7B GPT-3 model where the prediction output is an answer; **GPT-Neo-P**: the fine-tuned 2.7B GPT-3 model where the prediction output is a program; **Codex**: the few-shot Codex model where the prediction output is a program.

ID	Dataset	References
1	addsub	(Hosseini et al., 2014)
2	amps	(Hendrycks et al., 2021b)
3	APPS	(Hendrycks et al., 2021a)
4	asdiv	(Miao et al., 2020b)
5	conala	(Yin et al., 2018)
6	mathematics	(Saxton et al., 2019)
7	dolphin	(Huang et al., 2016)
8	draw	(Upadhyay and Chang, 2015)
9	GSM8k	(Cobbe et al., 2021)
10	MATH	(Hendrycks et al., 2021b)
11	mathqa	(Amini et al., 2019)
12	mbpp	(Austin et al., 2021)
13	MCTaco	(Zhou et al., 2019)
14	multiarith	(Roy and Roth, 2015)
15	NumerSense	(Lin et al., 2020)
16	NumGLUE	(Mishra et al., 2022b; Dua et al., 2019b; Ravichander et al., 2019; Kushman et al., 2014; Tafjord et al., 2019; Roy and Roth, 2018, 2017; Koncel-Kedziorski et al., 2016, 2015)
17	simuleq	(Kushman et al., 2014)
18	singleop	(Roy et al., 2015)
19	singleq	(Koncel-Kedziorski et al., 2015)
20	svamp	(Patel et al., 2021)

Table 21: List of source datasets and corresponding references used in constructing LILA.