# Līla

## A Unified Benchmark for Mathematical Reasoning

**Matthew Finlayson**[1]    Swaroop Mishra[1]    Pan Lu    Leonard Tang
Sean Welleck    Chitta Baral    Tanmay Rajpurohit    Oyvind Tafjord
Ashish Sabharwal    Peter Clark    Ashwin Kalyan

September 16, 2022

mattf1n.github.io                                    matthewf@allenai.org

[1]Equal first-authors

# TL;DR

# TL;DR



▶ Current math reasoning evaluation is broken.

# TL;DR



▶ Current math reasoning evaluation is broken.
▶ We build L̄īLA, a comprehensive benchmark.

# TL;DR



- Current math reasoning evaluation is broken.
- We build LĪLA, a comprehensive benchmark.
- We train BHĀSKARA, a foundational math reasoning model.
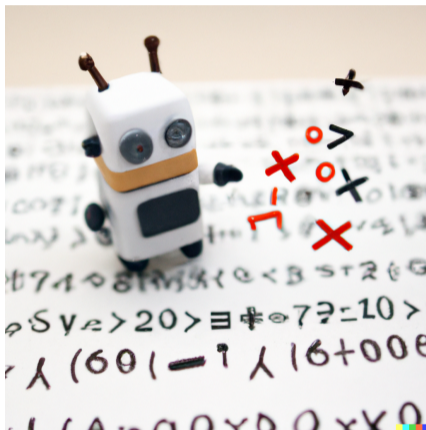
# Motivation

# Motivation

Can language models do math?

# Motivation

Can language models do math?   How can we find out?

# Motivation

Can language models do math?   How can we find out?

# Can language models do math?

Fill-in-the-blank

# Can language models do math?

Fill-in-the-blank

🧑 Fifty is equal to _ times ten.

A12

# Can language models do math?

Fill-in-the-blank ✅

👱 Fifty is equal to _ times ten.
🤖 `Five`

# Can language models do math?

Fill-in-the-blank ✅

# Can language models do math?

Fill-in-the-blank ✅     Common sense

# Can language models do math?

Fill-in-the-blank ✅          Common sense

🧑 A skiff refuels after 10 miles in the bay compared to 4 at sea. Which is more rugged?

# Can language models do math?

Fill-in-the-blank ✅     Common sense ❌

🧑 A skiff refuels after 10 miles in the bay compared to 4 at sea. Which is more rugged?

🤖 The bay

Ai2

# Can language models do math?

Fill-in-the-blank ✅     Common sense ❌     Algebra

# Can language models do math?

Fill-in-the-blank ✅     Common sense ❌     Algebra

🧑 Solve $x + 9j = 27 + 6$ for $x$ when $5j - 2 - 18 = 0$.

**AI2**

# Can language models do math?

Fill-in-the-blank ✅    Common sense ❌    Algebra ❌

👦 Solve $x + 9j = 27 + 6$ for $x$ when $5j - 2 - 18 = 0$.

🤖 `x = 63`

AI2

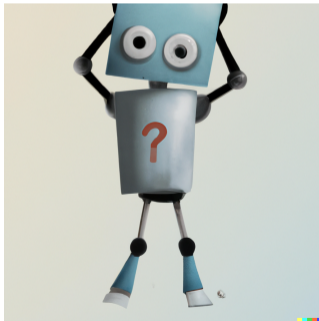# Can language models do math?

Fill-in-the-blank ✅   Common sense ❌   Algebra ❌

# Can language models do math?

Fill-in-the-blank ✅     Common sense ❌   Algebra ❌   Number theory ❓
Multiple-choice ❓ Comparison ❓ Science knowledge ❓ Arithmetic ❓ Geometry ❓ ...

# A math question taxonomy

# A math question taxonomy

| Dataset | Format | Subject | Knowledge | Language |
|---------|--------|---------|-----------|----------|
| Numersense | Fill-in | Arithmetic | Math | Simple |
| NumGLUE | Multi-choice | Comparision | Real world | Complex |
| Deepmind | Generative | Calculus | Math | None |
| MCTaco | Multi-choice | Arithmetic | Commonsense | Simple |
| ... | ... | ... | ... | ... |

# A math question taxonomy

| Dataset | Format | Subject | Knowledge | Language |
|---------|--------|---------|-----------|----------|
| Numersense | Fill-in | Arithmetic | Math | Simple |
| NumGLUE | Multi-choice | Comparision | Real world | Complex |
| Deepmind | Generative | Calculus | Math | None |
| MCTaco | Multi-choice | Arithmetic | Commonsense | Simple |
| ... | ... | ... | ... | ... |

# Direct answering is unsatisfying

🧑 Solve $x + 9j = 27 + 6$ for $x$ when $5j - 2 - 18 = 0$.

🤖 `x = 63` ❌

# Direct answering is unsatisfying

👦 Solve $x + 9j = 27 + 6$ for $x$ when $5j - 2 - 18 = 0$.

🤖 `x = -3` ✅

## Direct answering is unsatisfying

🧑 Solve $x + 9j = 27 + 6$ for $x$ when $5j - 2 - 18 = 0$.

🐍
```
>>> j = (0 + 2 + 18) / 5
... x = 27 + 6 - 9 * j
... print(x)
-3
```

**AI2**

# Language models ♥ Python



**Riley Goodside**
@goodside · Follow

"You are GPT‑3, and you can't do math":
Prompting GPT‑3 via zero-shot instruction to
answer calculation/math questions by consulting
a Python REPL.



**Sergey Karayev**
@sergeykarayev · Follow

Here's a brief glimpse of our INCREDIBLE near
future.

GPT-3 armed with a Python interpreter can
· do exact math
· make API requests
· answer in unprecedented ways

Thanks to **@goodside** and **@amasad** for the idea
and repl!

Play with it: **replit.com/@SergeyKarayev...**

# Halfway recap

# Halfway recap

- ► Can language models do math?

# Halfway recap

- ▶ Can language models do math?
- ▶ Existing benchmarks are too narrow in scope

# Halfway recap

- ▶ Can language models do math?
- ▶ Existing benchmarks are too narrow in scope
- ▶ Python programming > Direct answering

# Līla: a comprehensive benchmark

AMPS MATH, Numersense, NumGLUE, MCTaco, …

# Lῑla: a comprehensive benchmark

AMPS MATH, Numersense, NumGLUE, MCTaco, …

> 🧑 Find the laplacian of the function $f(x, y, z)$ where
> $f(x, y, z) = x^3 y^3$.

**Ai2**

# LĪLA: a comprehensive benchmark

AMPS MATH, Numersense, NumGLUE, MCTaco, …

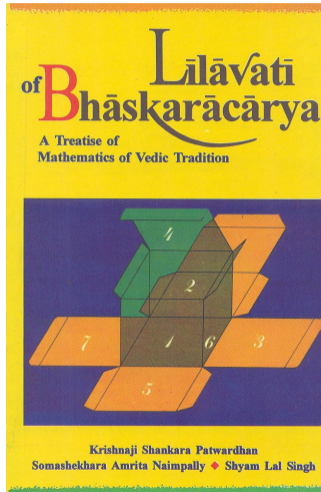🧑 Find the laplacian of the function $f(x, y, z)$ where $f(x, y, z) = x^3 y^3$.
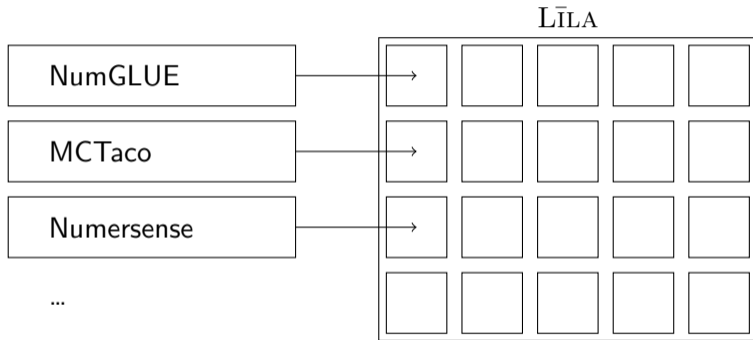
🤖 $6x^3 y + 6xy^3$

Ai2

# Līla: a comprehensive benchmark

AMPS MATH, Numersense, NumGLUE, MCTaco, …

👨 Find the laplacian of the function $f(x, y, z)$ where
$f(x, y, z) = x^3 y^3$.

🤖 $6x^3 y + 6xy^3$

🐍
```python
from sympy import *
C = CoordSys3D('C')
x, y, z = C.x, C.y, C.z
f = x**3*y**3
print(laplacian(f))
```

# 📜 LĪLA: a comprehensive benchmark

AMPS MATH, Numersense, NumGLUE, MCTaco, …

🧑 Find the laplacian of the function $f(x, y, z)$ where $f(x, y, z) = x^3 y^3$.

🤖 $6x^3 y + 6xy^3$

🐍
```python
from sympy import *
C = CoordSys3D('C')
x, y, z = C.x, C.y, C.z
f = x**3*y**3
print(laplacian(f))
```

$$100,000 \times (\text{🧑}, \text{🤖}, \text{🐍}) = \text{📜 LĪLA}$$

**AI2**

# Līla splits

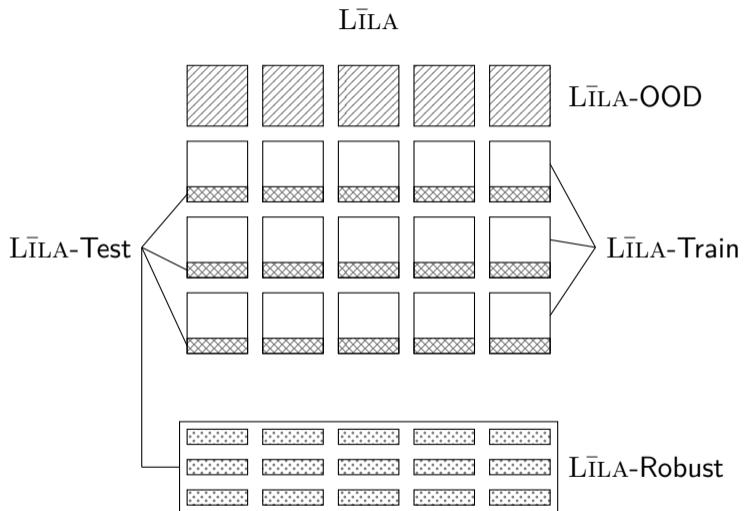# LĪLA splits

# Līla splits

# Līla splits

# Līla splits
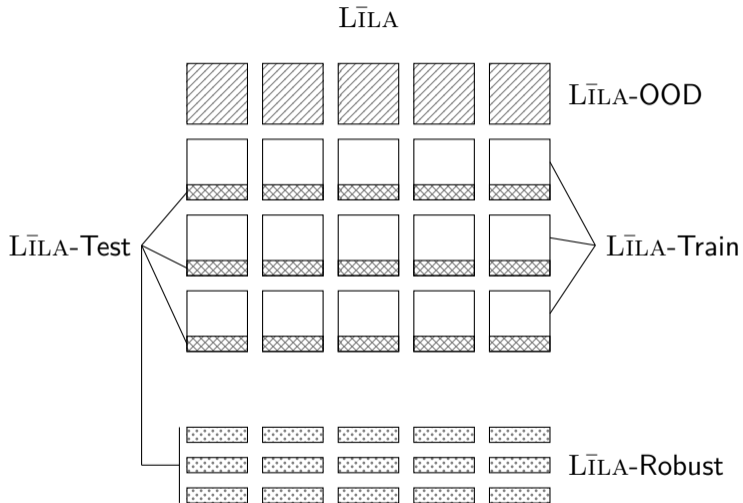


Jules gave 3 apples to...

# Līla splits



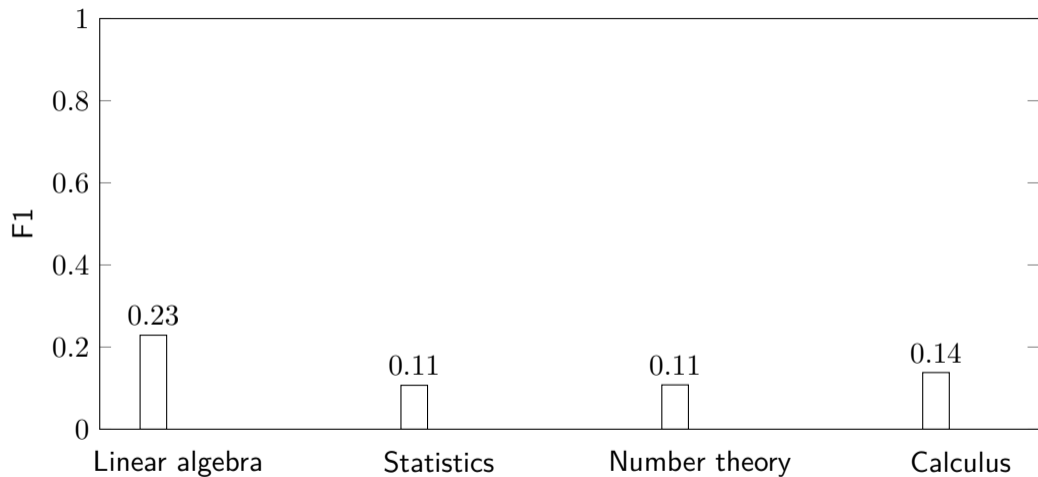Riley is 7 years old. Jules gave 3 apples to...
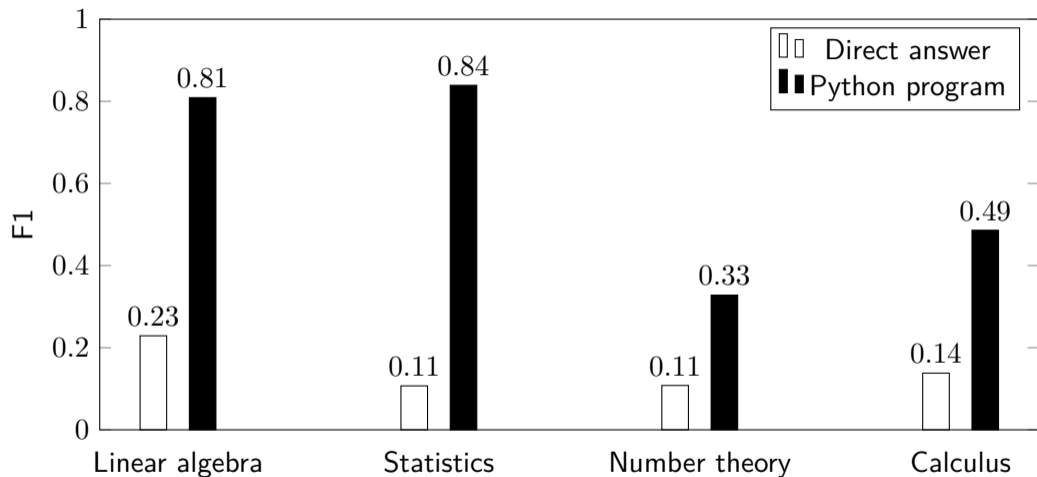
# Līla splits

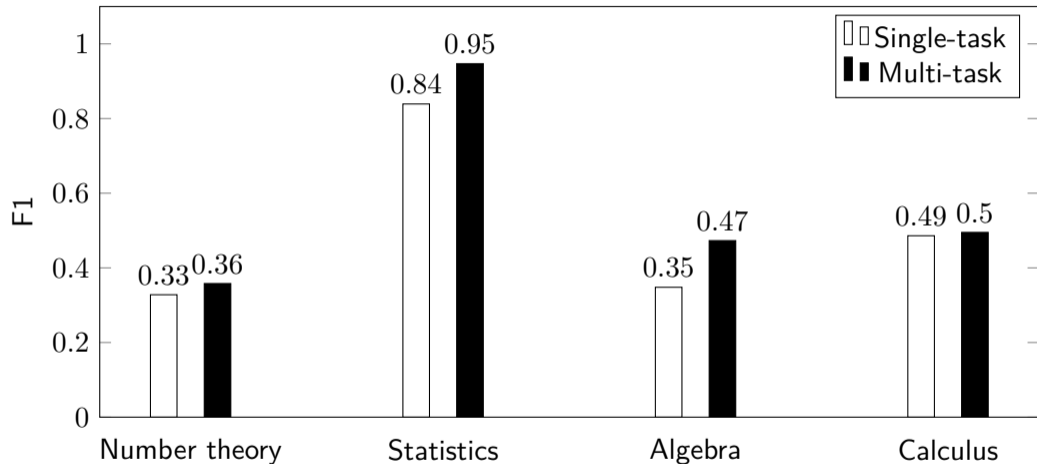# Līla splits

# Python program > direct answer

# Python program > direct answer



Program answering 23 points better on average
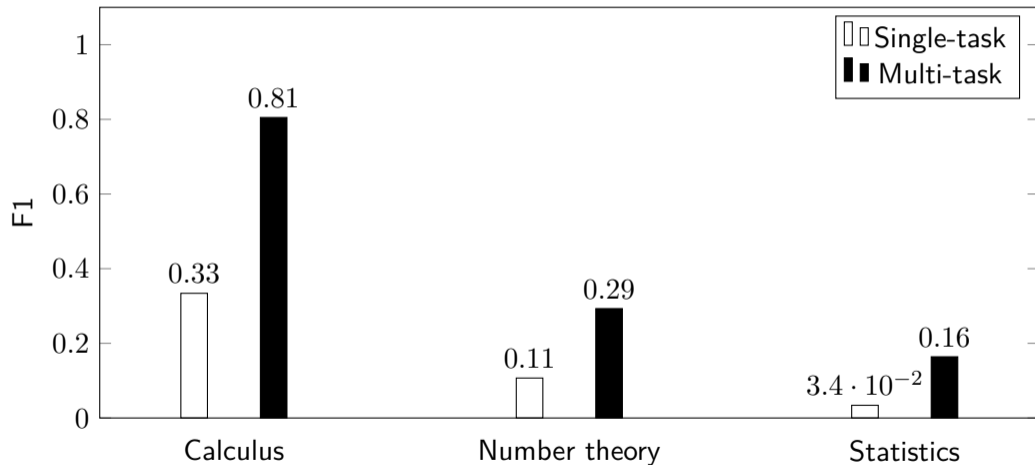
# Multi-task > single-task

# Multi-task > single-task



Multi-task model is 9 points better on average

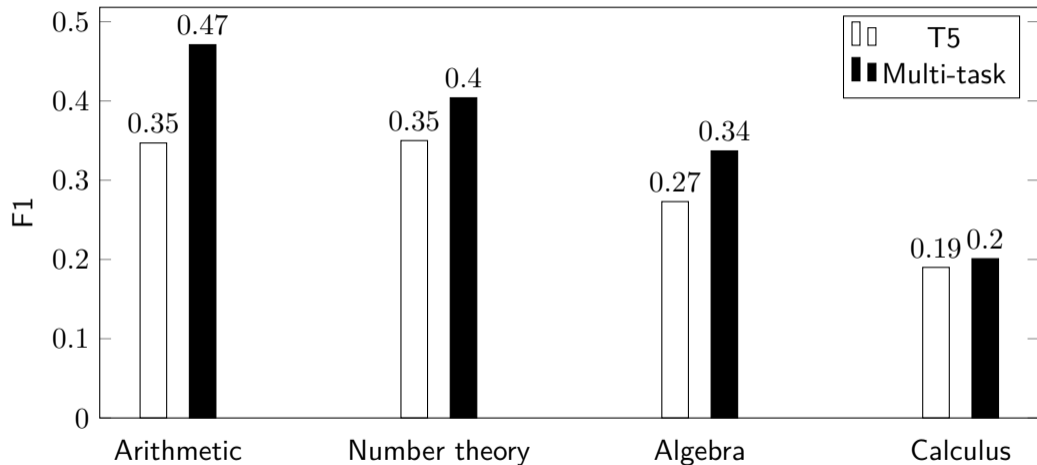# Multi-task model is *general*



LĪLA-OOD

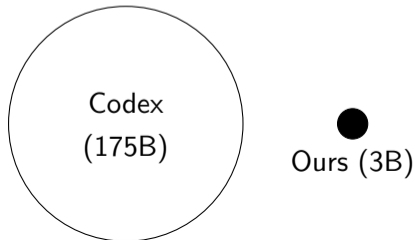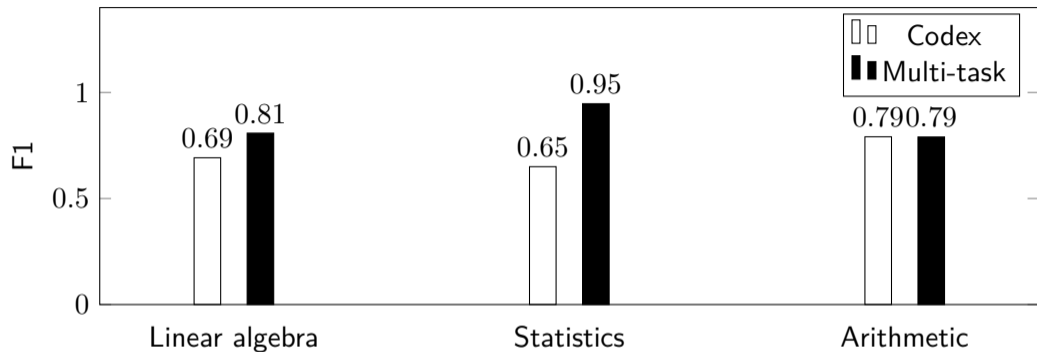Multi-task model is 21 points better on average

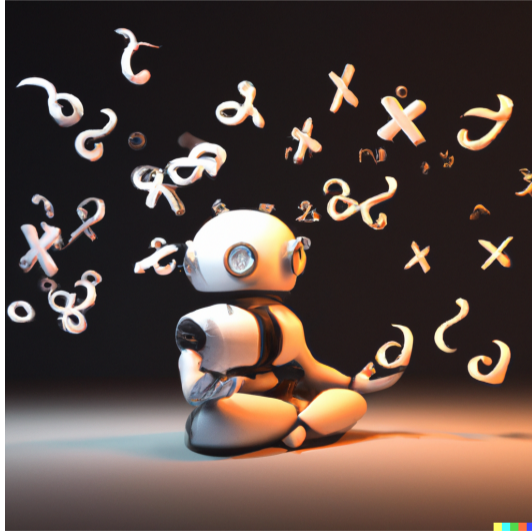# Our multi-task model is a starting point for math models



LĪLA-OOD

Multi-task model is 8 points better on average

# Our multi-task model outperforms Codex on some tasks

# Takeaways

# Takeaways

- Math reasoning evaluation is broken.

# Takeaways

- Math reasoning evaluation is broken.
- LĪLA: a comprehensive benchmark with useful splits.

AI2

# Takeaways

- ▶ Math reasoning evaluation is broken.
- ▶ LĪLA: a comprehensive benchmark with useful splits.
- ▶ BHĀSKARA: a multi-task model for math reasoning.

# Takeaways

- ▶ Math reasoning evaluation is broken.
- ▶ LĪLA: a comprehensive benchmark with useful splits.
- ▶ BHĀSKARA: a multi-task model for math reasoning.

Thank you!