**Mathematics and Machine Learning**

# BENCHMARKING PROGRAM-AIDED VS. NATURAL LANGUAGE MATHEMATICAL REASONING IN LLMS

Mohamed El Maghari, Ralf König, Fritz Körner, Khalid Sabih

Leipzig, 18.02.2026

# MOTIVATION

## The Paradox of Modern LLMs in Mathematics

- **AlphaProof** (DeepMind, 2025) — silver medal at IMO 2024
    - Solved P6, the hardest problem, solved by only 5 of 609 humans
    - Uses Lean 4 formal verification → zero hallucination possible
- **Yet top models fail at this:**
    - *"Alice has N brothers and M sisters. How many sisters does Alice's brother have?"*
    - State-of-the-art models answer this **incorrectly**
- **So which is it — can LLMs do math or not?**

→ **Research Goal:** Compare Pure LLM vs. Python vs. Lean 4 on the same problems, same model, same conditions

Mohamed El Maghari, Ralf König, Fritz Körner, Khalid Sabih

# HYPOTHESES

- **H1:** Easy problems solvable across all three modalities
- **H2:** Hard problems benefit from Python code generation
- **H3:** Lean 4 zero-shot accuracy will be very low
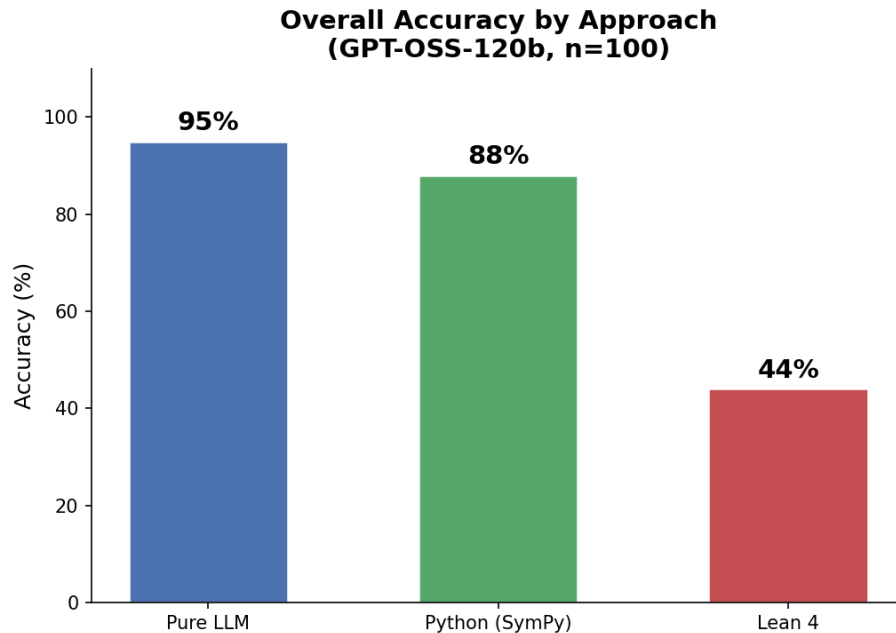- **H4:** Code-specialized LLMs outperform general-purpose LLMs as program-aided models

# METHODOLOGY / DATASET

- Dataset: MATH-500 benchmark — 100 problems sampled
- 7 categories × 5 difficulty levels (L1–L5)
- Model tested: GPT-OSS-120b via Blablador (Helmholtz)
- Zero-shot prompting — no examples, no feedback loop
- Python sandbox with SymPy, NumPy, Math libraries
- Lean 4 via Kimina-LEAN-Server for remote proof verification
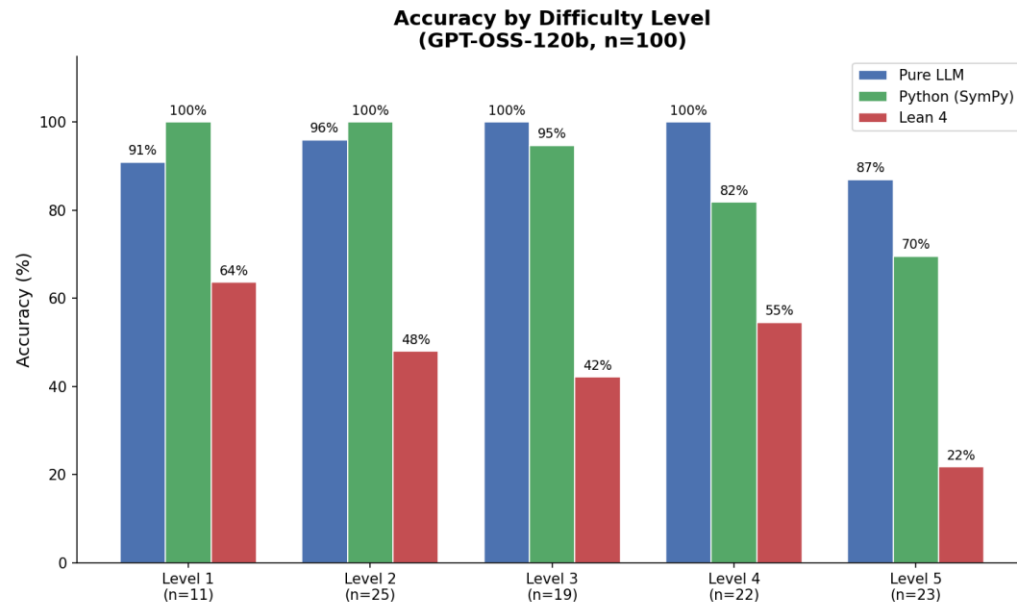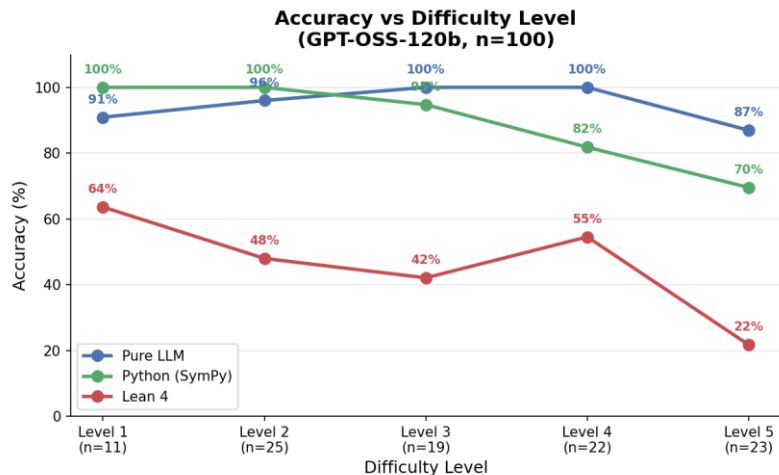- Evaluation: LLM-as-Judge

# RESULTS

**Overall Accuracy:**

- Pure LLM 95%
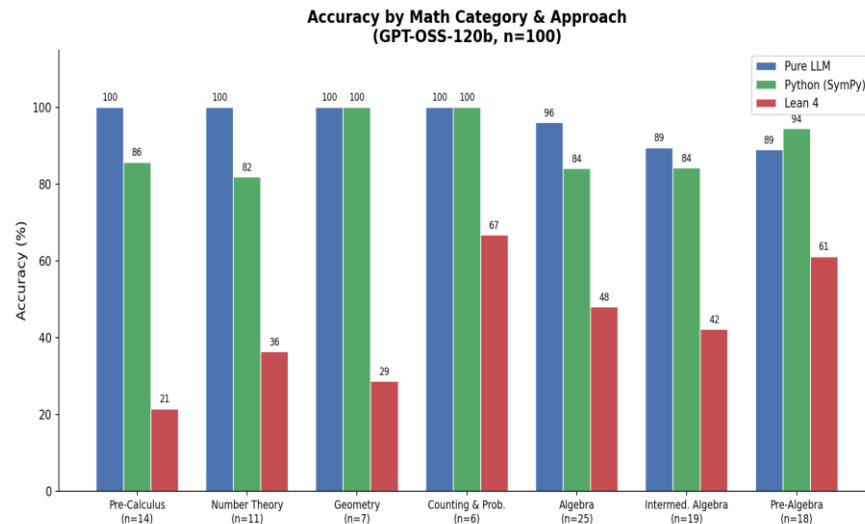- Python (SymPy) 88%
- Lean 4 44%

**Overall Accuracy by Approach
(GPT-OSS-120b, n=100)**

# RESULTS : ACCURACY BY DIFFICULTY LEVEL
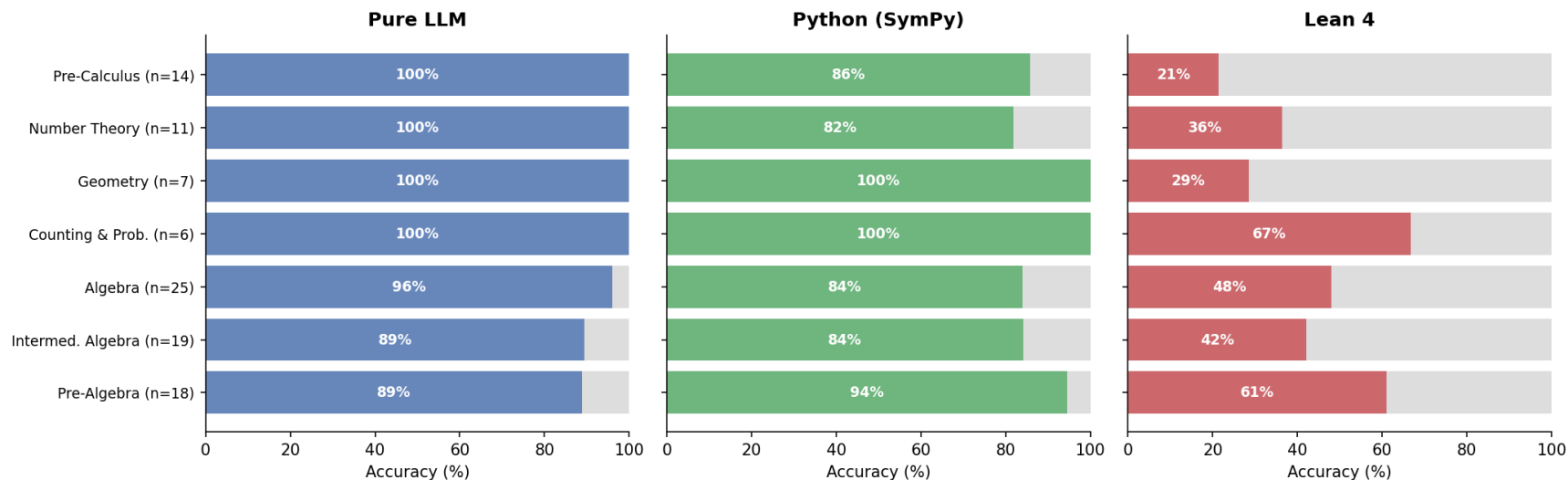
UNIVERSITÄT LEIPZIG

# RESULTS: ACCURACY BY MATH CATEGORY & APPROACH

## Category Breakdown:

- Pure LLM dominates
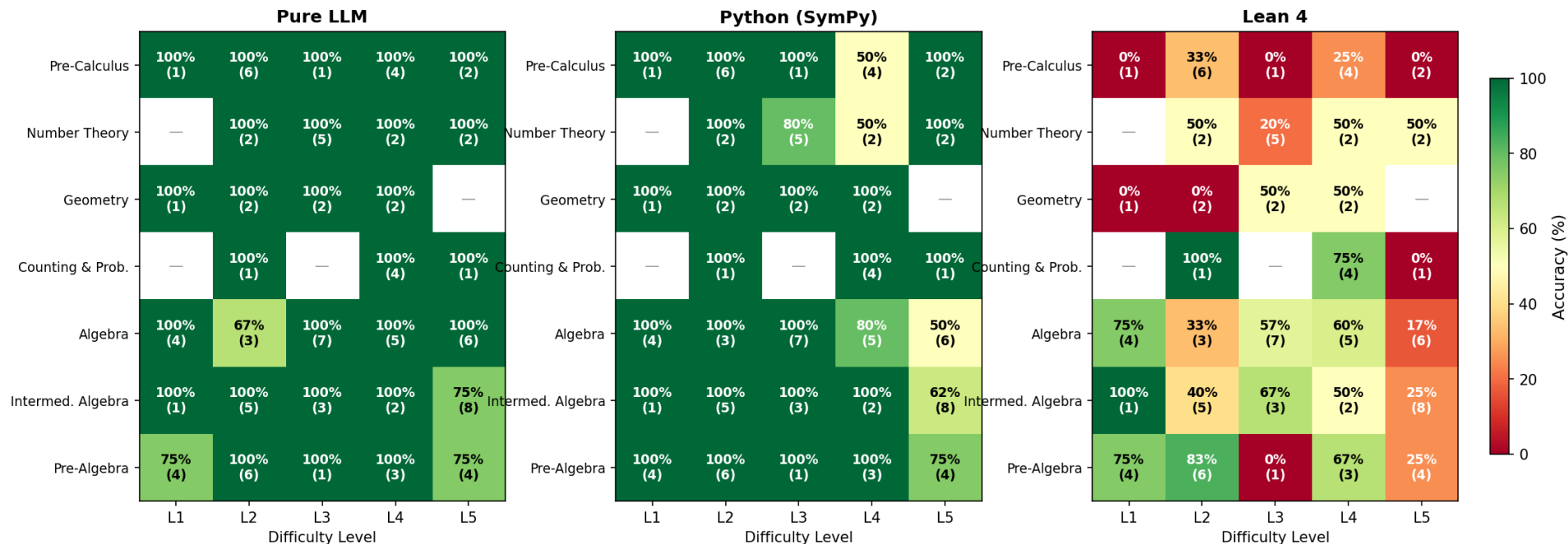- Lean 4 struggles especially in Pre-Calculus (21%) and Geometry (29%)



Accuracy by Math Category & Approach (GPT-OSS-120b, n=100)

Mohamed El Maghari, Ralf König, Fritz Körner, Khalid Sabih

UNIVERSITÄT LEIPZIG

# RESULTS : CORRECT VS INCORRECT BY CATEGORY



**Correct vs Incorrect by Category
(GPT-OSS-120b, n=100)**

| | Pure LLM | Python (SymPy) | Lean 4 |
|---|---|---|---|
| Pre-Calculus (n=14) | 100% | 86% | 21% |
| Number Theory (n=11) | 100% | 82% | 36% |
| Geometry (n=7) | 100% | 100% | 29% |
| Counting & Prob. (n=6) | 100% | 100% | 67% |
| Algebra (n=25) | 96% | 84% | 48% |
| Intermed. Algebra (n=19) | 89% | 84% | 42% |
| Pre-Algebra (n=18) | 89% | 94% | 61% |

# RESULTS : ACCURACY BY CATEGORY & DIFFICULTY



Accuracy by Category & Difficulty
(GPT-OSS-120b, n=100)

# CONCLUSION

– Pure LLM: best overall accuracy **(95%)**
– Python (SymPy): competitive but drops at higher difficulty **(88%)**
– Lean 4: significant underperformance across all dimensions **(44%)**
– Difficulty scaling disproportionately hurts Lean 4 **(64% → 22%)**
– Proof generation — not mathematics — is the core bottleneck
– Zero-shot Pure LLM remains the most practical approach today

Mohamed El Maghari, Ralf König, Fritz Körner, Khalid Sabih

# FUTURE WORK

- Benchmark additional models: GPT-4o, Claude, Gemini
- Test few-shot and chain-of-thought prompting strategies
- Fine-tune models specifically on Lean 4 proof corpora
- Explore hybrid pipelines: Python solving + Lean 4 verification
- Scale dataset beyond 100 problems and more categories
- Analyze Lean 4 failure modes systematically

# THANK YOU!

Ralf König
Fritz Körner
Khalid Sabih
Mohamed El Maghari