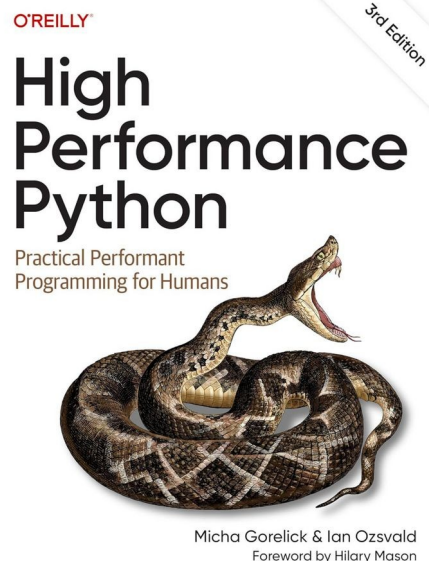# playgroup – deep dive LLM day

Mor Consulting 2025-09

@IanOzsvald – ianozsvald.com

# Interim Chief Data Scientist

- Strategist/Trainer/Speaker/Author 25+ years

- Figuring where LLMs fit into DS

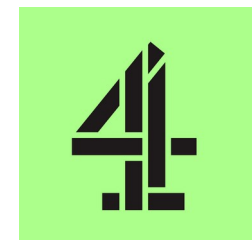Part of **PyData - 165 groups**
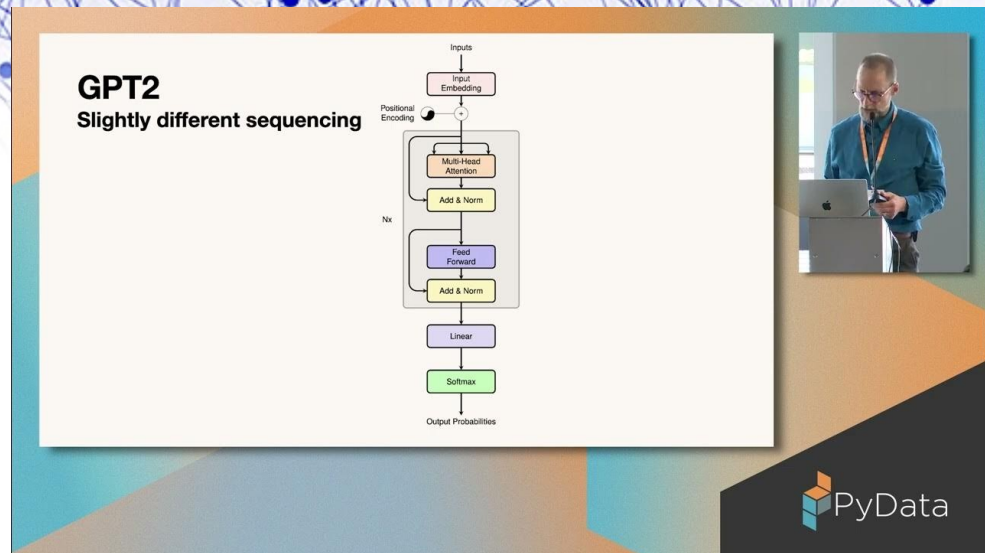
**PyData London Meetup**

4.7 ★★★★½  2576 ratings

Where are the creatives?

London, United Kingdom

15,298 members · Public group

Organized by **NumFOCUS, Inc.** and **14 others**

GPT2
Slightly different sequencing



Pydata London



EXPERT INSIGHT

Generative AI
with LangChain

Build large language model (LLM) apps with
Python, ChatGPT, and other LLMs

2024 Edition
Includes updated code and content
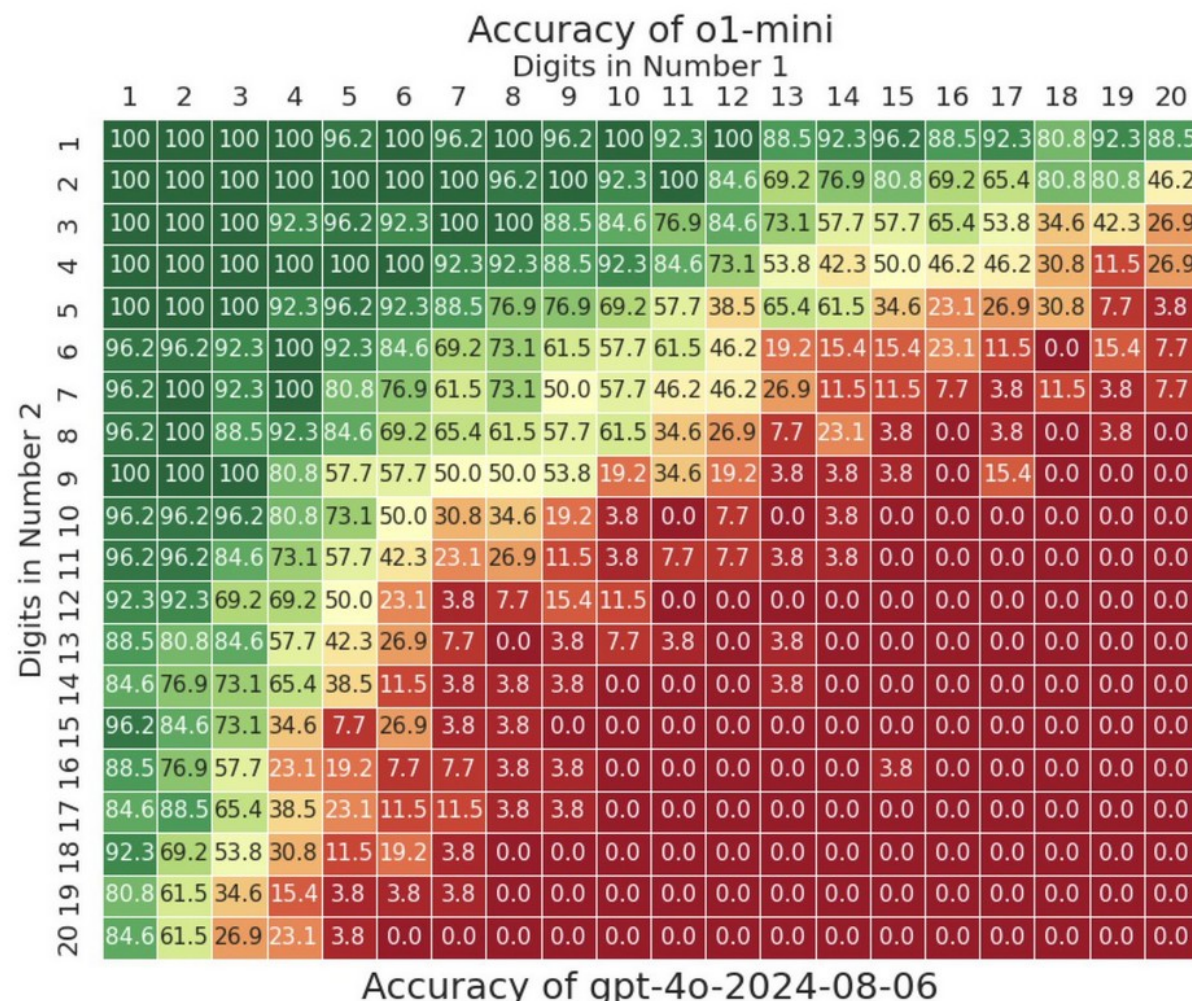
‹packt›

Ben Auffarth

# Goal

Valuable Lessons Learned on Kaggle's ARC AGI LLM challenge
PyDataGlobal 2024-12 talk

- Will *agents take over the world* or are we living in a world of *approximate retrieval*? Is AGI nearly here?

- Can an LLM solve novel problems? See? Reflect?

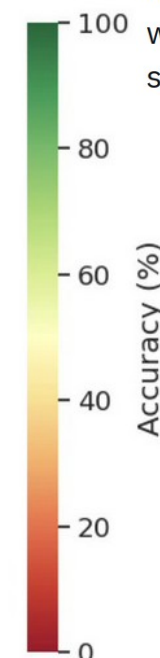- You – think on a novel problem, meet interesting folk, get your qs answered

By [ian]@ianozsvald[.com]                    Ian Ozsvald

# Not so good at multiplication

Accuracy of o1-mini

**Yuntian Deng** @yuntiandeng

Is OpenAI's o1 a good calculator? We tested it on up to 20x20 multiplication—o1 solves up to 9x9 multiplication with decent accuracy, while gpt-4o struggles beyond 4x4. For context, this task is solvable by a small LM using implicit CoT with stepwise internalization. 1/4
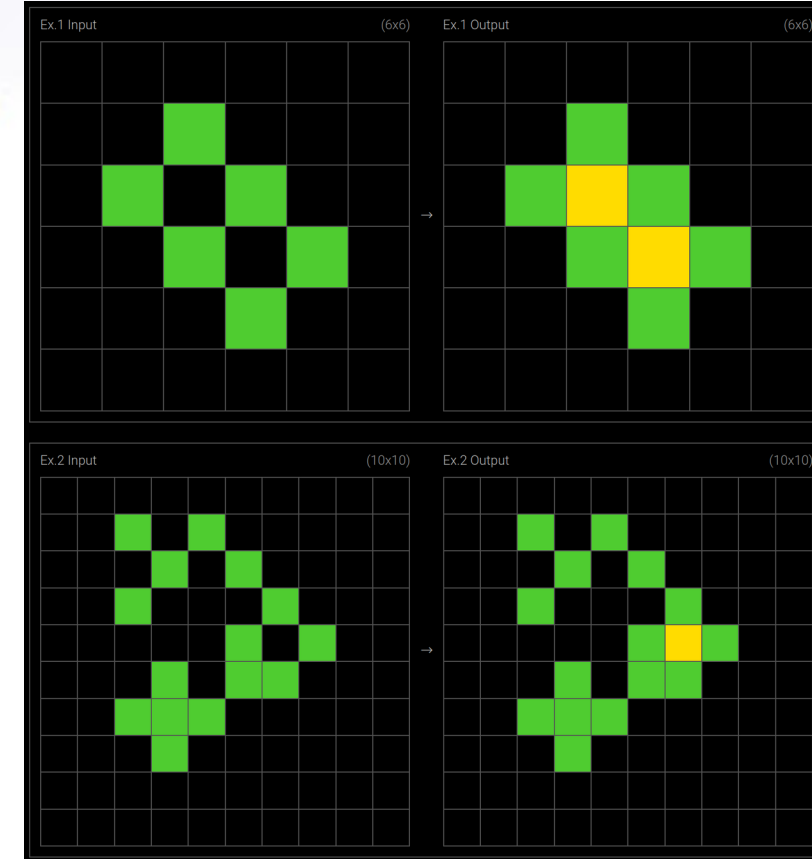
Maybe it lacks short term memory and iterative processing?

Tokens – representation issues?

**Approximate retrieval** at work?

# Agenda

- Talk about ARC AGI, try manually

- Can an agentic(?) method reflect and improve?

- → office: prompting, testing, auto-code SQL, resilience

# Business thoughts

- VCs will want their cash back at some point

- Scaling is expensive – can we keep our solution?

- Keep IP in-house

- Maybe we don't need to burn the planet on LLMs

# Am I asking the right question?

- Representation

- Prompt

- Process

- What am I missing? What's a **big question** to ask?

By [ian]@ianozsvald[.com]                    Ian Ozsvald

# Kick off

- Do you have **slack**? Do you have the code?

  – Add to the slack with shared notes, branch code

- **Tables – what's a funny/useful GenAI story?** Share back, start in pairs, decide on someone's example to share – 15 mins

ARC-AGI-2: A New Challenge for Frontier AI Reasoning Systems

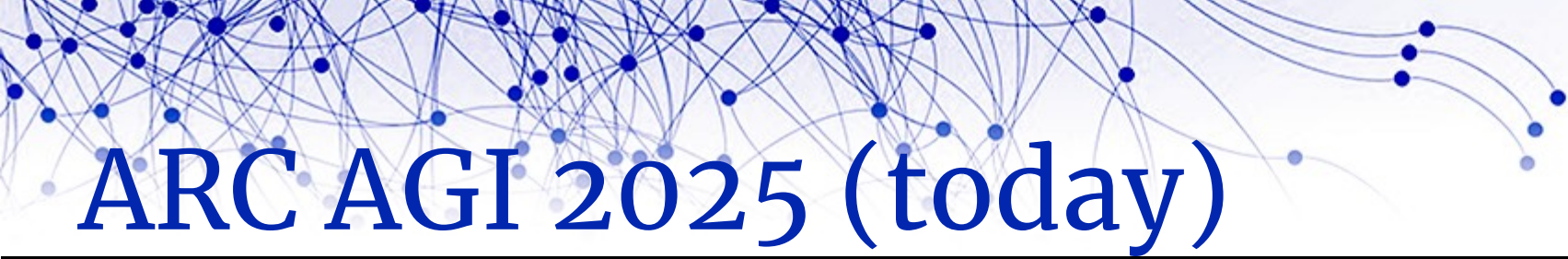François Chollet*    Mike Knoop    Gregory Kamradt    Bryan Landers
Henry Pinkard

May 20, 2025

| Model | ARC-AGI-1 | ARC-AGI-2 |
|---|---|---|
| o3-mini (High) | 34.5% | 3.0% |
| o3 (Medium) | 53.0% | 3.0% |
| ARChitects (ARC Prize 2024) | 56.0% | 2.5% |
| o4-mini (Medium) | 41.8% | 2.4% |
| Icecuber (ARC Prize 2020) | 17.0% | 1.6% |
| o1-pro (Low) | 23.3% | 0.9% |
| Claude 3.7 (8K) | 21.2% | 0.9% |

- ARC AGI 1 (few years), now ARC AGI 2025

- 400+ problems, public and *private* (offline) set

- ARC AGI 1 "solved" by **GPT o3 88%** public $70k (xmas)

# ARC AGI 2025 (today)

## LEADERBOARD BREAKDOWN

| AI System | Organization | System Type | ARC-AGI-1 | ARC-AGI-2 | Cost/Task |
|-----------|--------------|-------------|-----------|-----------|-----------|
| Human Panel | Human | N/A | 98.0% | 100.0% | $17.00 |
| J. Berman (2025) | Bespoke | CoT + Synthesis | 79.6% | 29.4% | $30.40 |
| E. Pang (2025) | Bespoke | CoT + Synthesis | 77.1% | 26.0% | $3.97 |
| Grok 4 (Thinking) | xAI | CoT | 66.7% | 16.0% | $2.17 |
| GPT-5 (High) | OpenAI | CoT | 65.7% | 9.9% | $0.730 |
| Claude Opus 4 (Thinking 16K) | Anthropic | CoT | 35.7% | 8.6% | $1.93 |

By [ian]@ianozsvald[.com]

Ian Ozsvald

# Stages

- Limited GPU, Llama Scout (mm) about right – how should we represent the problem? Might vision help?

- We can try DeepSeek v3 0324 (quick, big context), circa 700GB VRAM at fp8 vs 90GB VRAM competition limit

- Does giving feedback help?

- Could 'agent framework' help? Open q

# Over to you

- Run the code, notes are in the README

- I'll tell you about our stages

- Try to talk to everyone in the room (cheatsheet)

# How could it do better?

- Make hypotheses, critique, rank

- Implement, get graded feedback, iterate

- Extract library of useful fns

- Writing code – solved?

By [ian]@ianozsvald[.com]

Ian Ozsvald

# How did others solve it?

- GA on human-designed solver components (no LLM)

- Library of human-solved clues, synthetic dataset

- Test-time fine tuning on 3 examples

  – Restricted representation fine tune

- GA to evolve prompts

# Problems with LLMs?

.pre-commit-config.yaml

```yaml
repos:
  # isort
  - repo: https://github.com/PyCQA/isort
    rev: 5.13.2    # Use the latest stable version
    hooks:
      - id: isort
```

isort 6.0.1

pip install isort

- Outdated knowledge