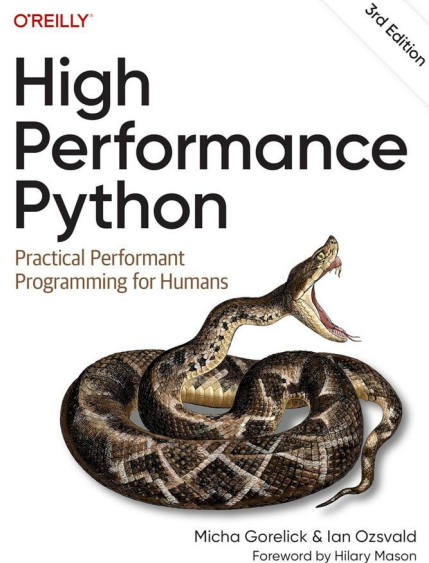# playgroup – deep dive LLM day

Mor Consulting 2025-06

@IanOzsvald – ianozsvald.com

# Interim Chief Data Scientist

- Strategist/Trainer/Speaker/Author 25+ years
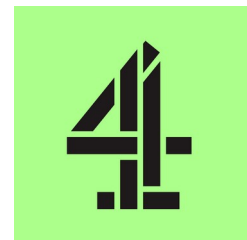
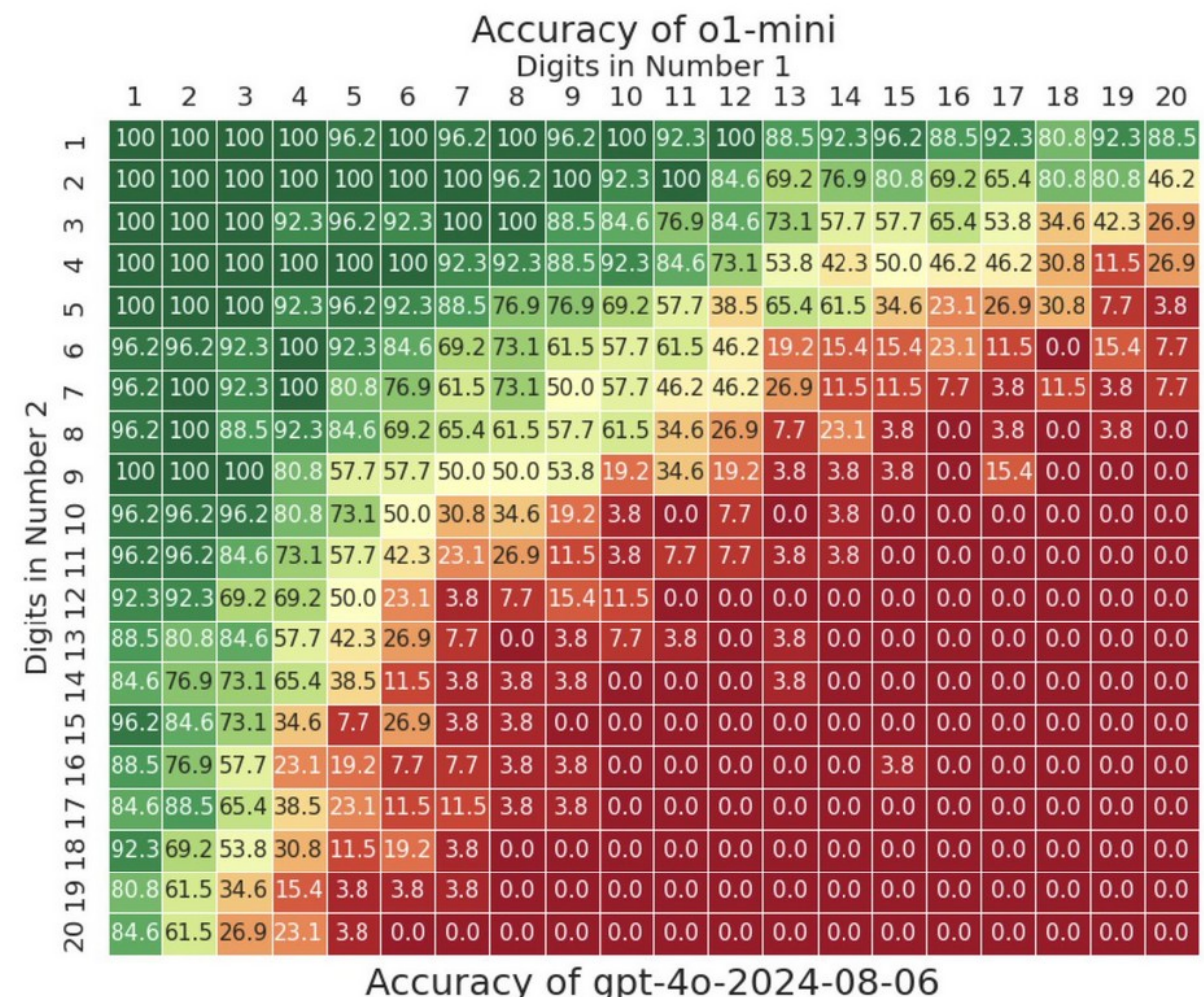- Figuring where LLMs fit into DS

# Goal

- 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
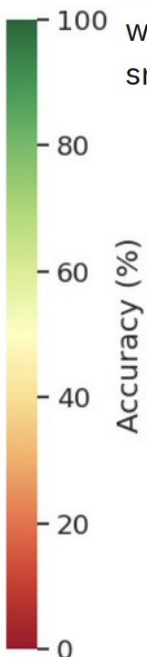
# Not so good at multiplication

Accuracy of o1-mini / Accuracy of gpt-4o-2024-08-06

**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

- Get LLM to solve some (maybe)

- Can a vLLM describe an image? Can you make img?

- Can an agent(?) reflect and improve?

# Kick off

- Do you have the Gdoc? Do you have the code?

  – Add to the Gdoc with shared notes, branch code

- **Tables – when is GenAI useful?** Share back, start in pairs, decide on someone's example to share – 15 mins

# ARC AGI

- ARC AGI few years, now ARC AGI 2025

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

- ARC AGI 1 "solved" by GPT o3 88% public $70k

# ARC AGI 2025 (today)

ARC-AGI-2 LEADERBOARD

| AI System | Score | $/Task |
|---|---|---|
| o3 (medium) | 3.0% | $2.53 |
| o3-mini (high) | 3.0% | $0.55 |
| ARChitects (2024) | 2.5% | $0.20 |
| o4-mini (Medium) | 2.4% | $0.23 |
| DeepSeek R1 | 1.3% | $0.08 |
| Gemini 2.0 Flash | 1.3% | $0.004 |

# Stages

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

- We can try DeepSeek, Opus 4 ($$$!)

- 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)

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

# How could it do better?

- Make hypotheses, critique, rank

- Implement, get graded feedback, iterate

- Extract library of useful fns

- Writing code – solved?

Ian Ozsvald