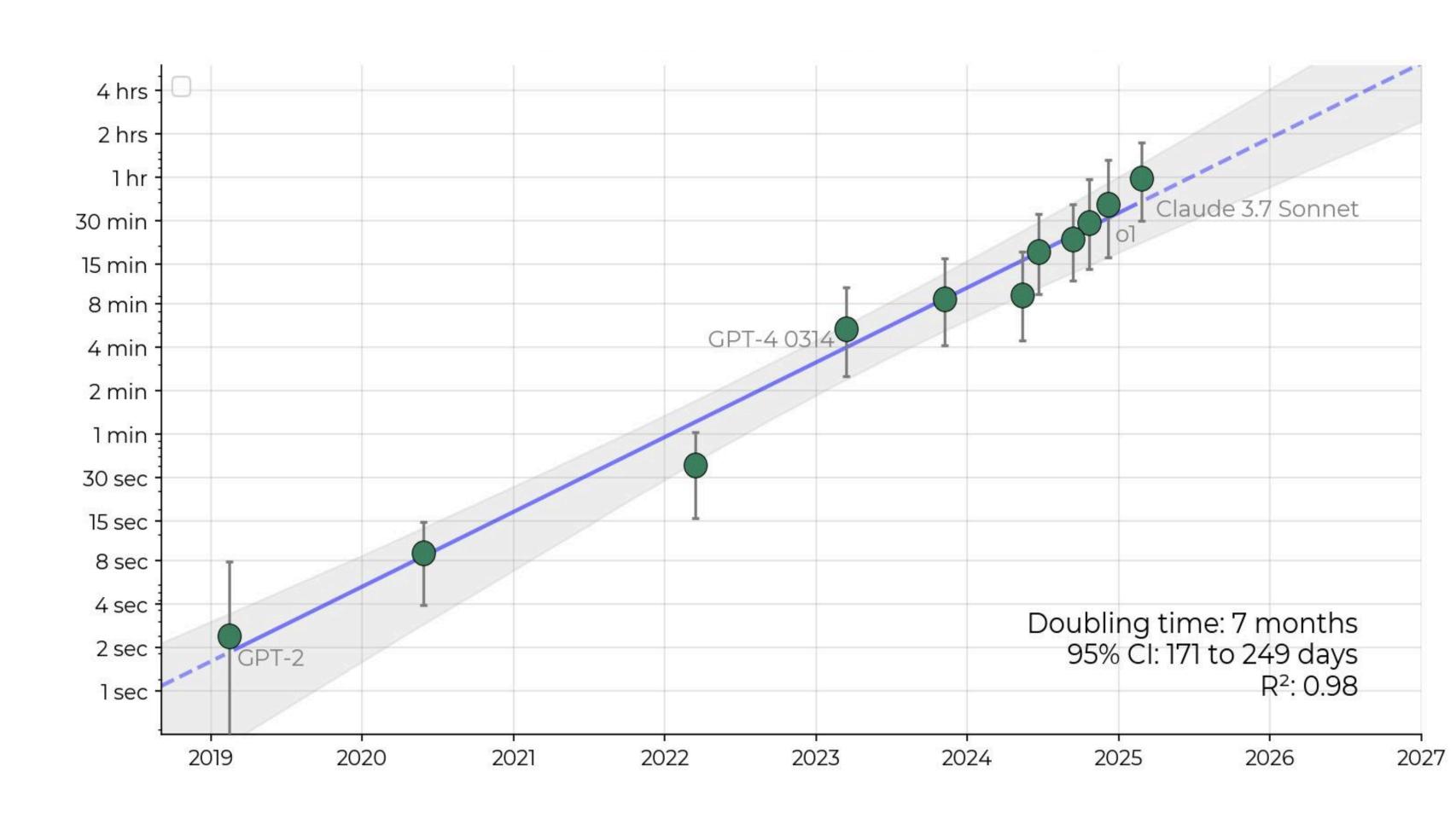
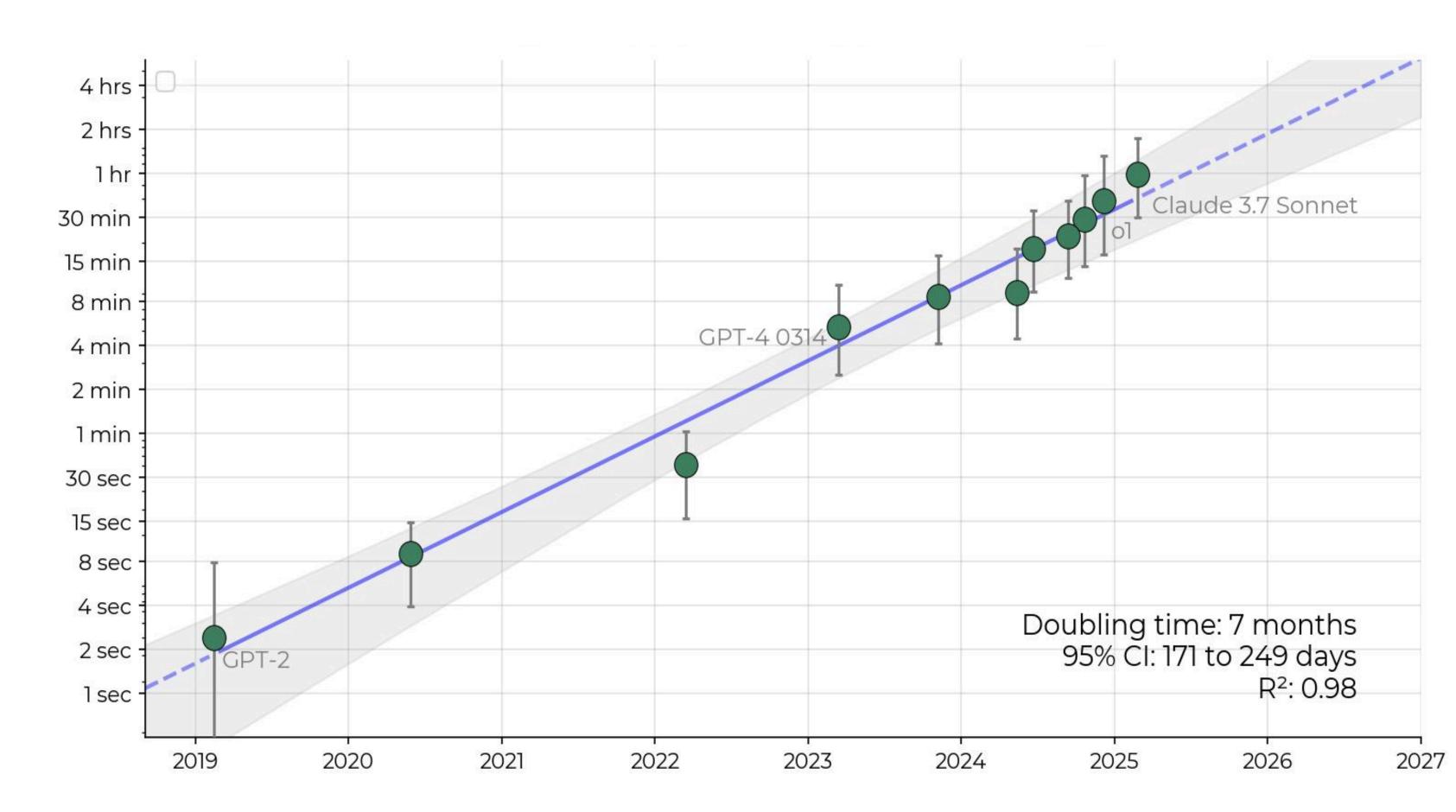


Ensuring Human Supervision for Als of *Tomorrow*

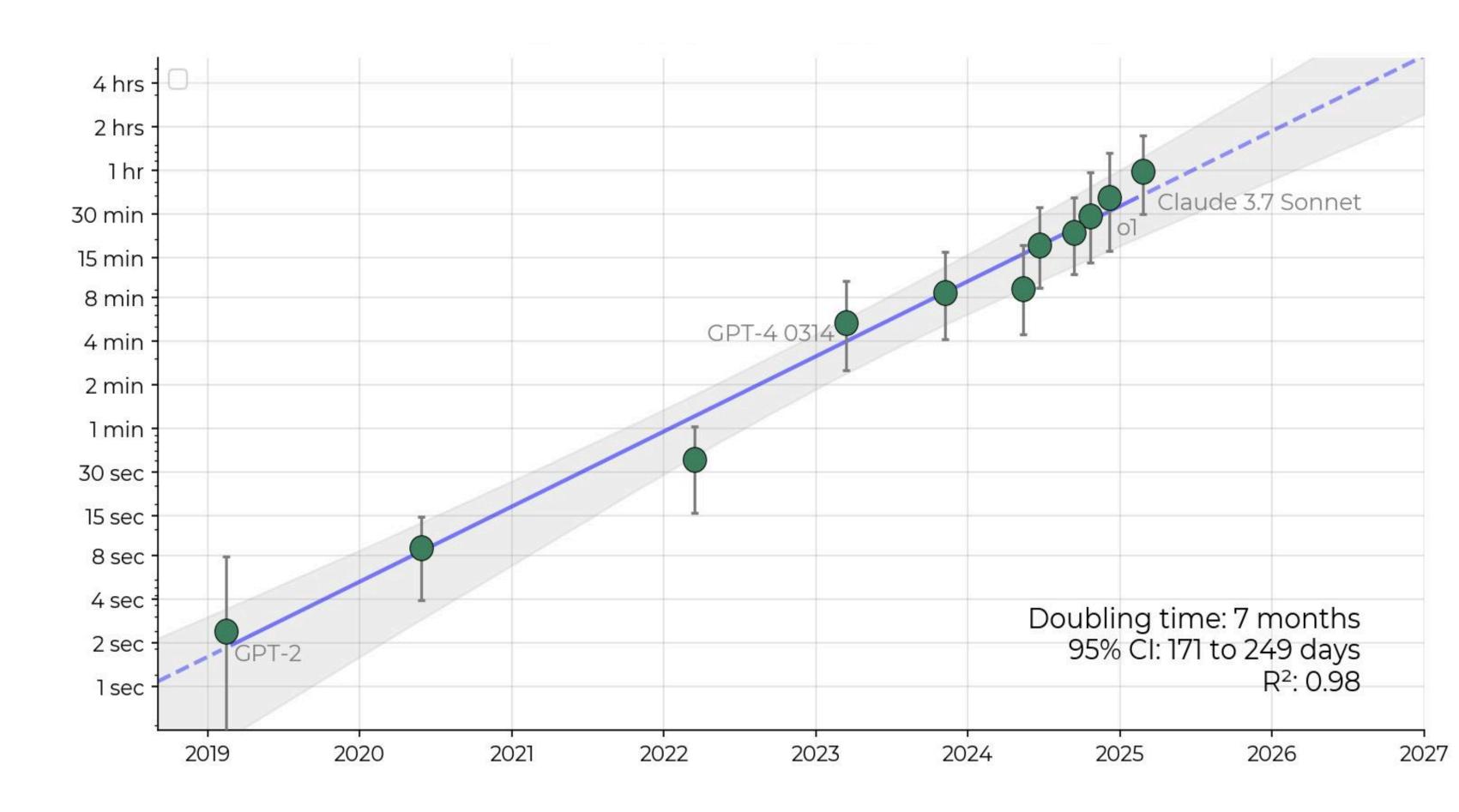


Task length is doubling every 7 months.



Task length is doubling every 7 months.

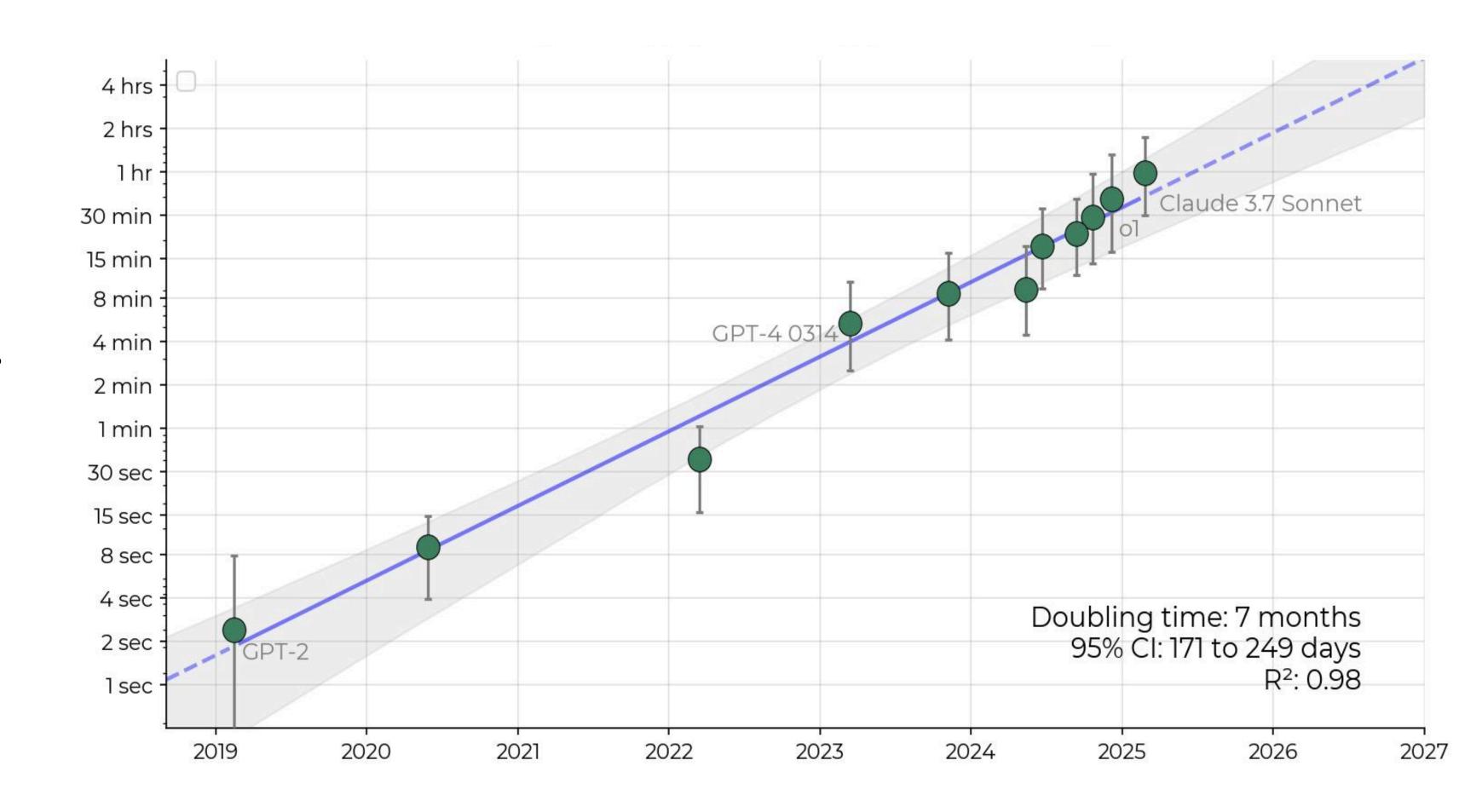
Exponential growth from self-improvement.



Task length is doubling every 7 months.

Exponential growth from self-improvement.

"90% of Claude Code is written by Claude itself"



Economically valuable Al capabilities are growing <u>rapidly</u>

Economically valuable

Al capabilities are growing <u>rapidly</u>

GDPval:

- 44 occupations
- 9 sectors > \$3T

Real Estate and Rental and Leasing



- Concierges
- Real Estate Sales Agents
- Real Estate Brokers
- Counter and Rental Clerks
- Property, Real Estate, & Community Association Managers

Government



- Recreation Workers
- Compliance Officers
- First-Line Supervisors of Police and Detectives
- Administrative Services Managers
- Child, Family, and School Social Workers

Manufacturing



- Mechanical Engineers
- Industrial Engineers
- Buyers & Purchasing Agents
- Shipping, Receiving, & Inventory Clerks
- First-Line Supervisors of Production and Operating Workers

Professional, Scientific, and Technical Services



- Software Developers
- Lawyers
- Accountants & Auditors
- Computer & Information Systems Managers
- Project Management Specialists

Health Care and Social Assistance



- Registered Nurses
- Nurse Practitioners
- Medical & Health Services Managers
- First-Line Supervisors of Office & Administrative Support Workers
- Medical Secretaries & Administrative Assistants

Finance and Insurance



- Customer Service Representatives
- Financial & Investment Analysts
- Financial Managers
- Personal Financial Advisors
- Securities, Commodities & Financial Services Sales Agents

Retail Trade



- Pharmacists
- General and Operations Managers
- Private Detectives & Investigators
- First-Line Supervisors of Retail Sales Workers

Wholesale Trade



- Sales Managers
- Order Clerks
- Sales Representatives, Wholesale & Manufacturing, Technical & Scientific Products
- Sales Representatives, Wholesale & Manufacturing, Except Technical & Scientific Products
- First-Line Supervisors of Non-Retail Sales Workers

Information



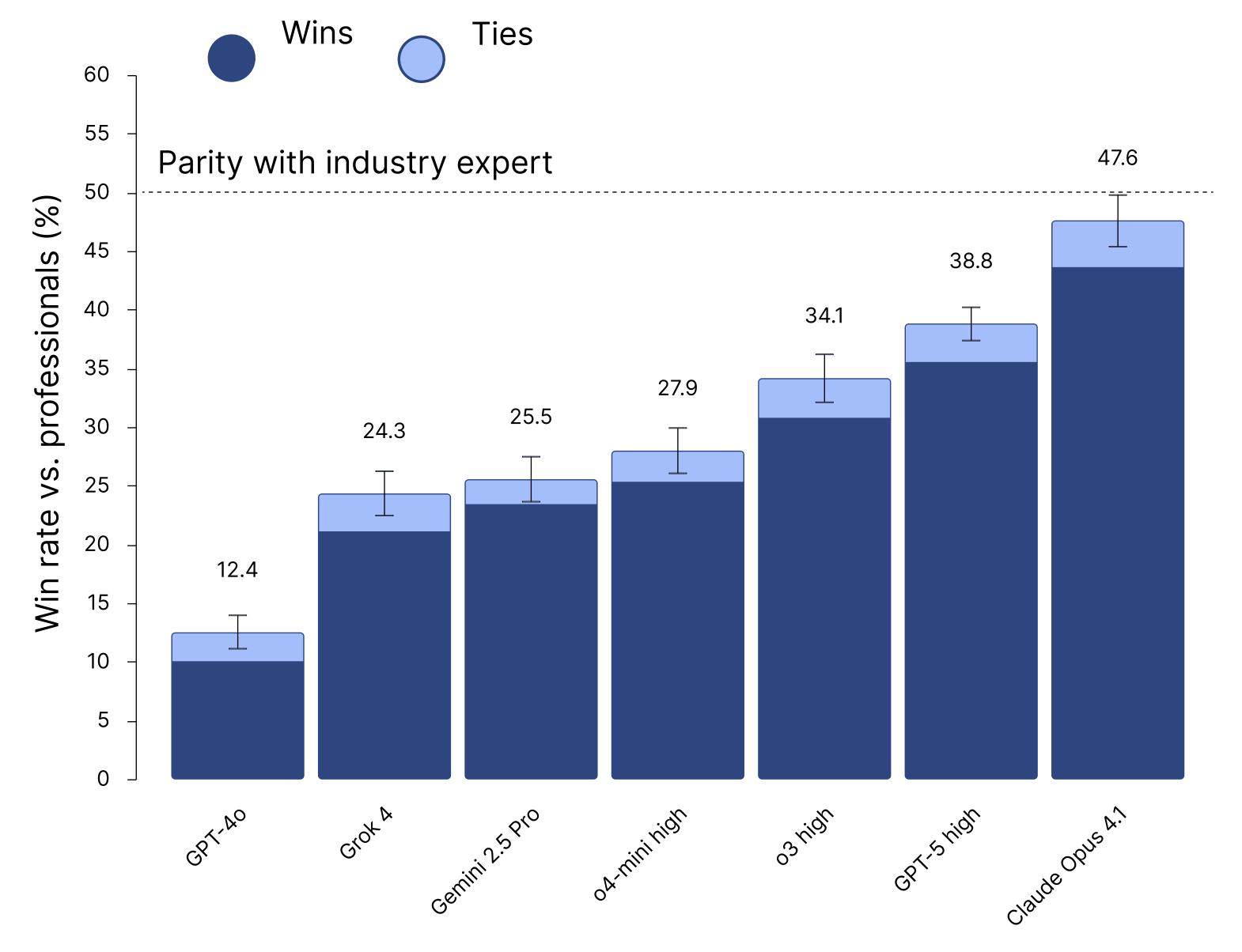
- Producers & Directors
- Film & Video Editors
- Editors
- News Analysts, Reporters, & Journalists
- Audio and Video Technicians

Economically valuable

Al capabilities are growing <u>rapidly</u>

GDPval:

- 44 occupations
- 9 sectors > \$3T
- 30 tasks / occupation

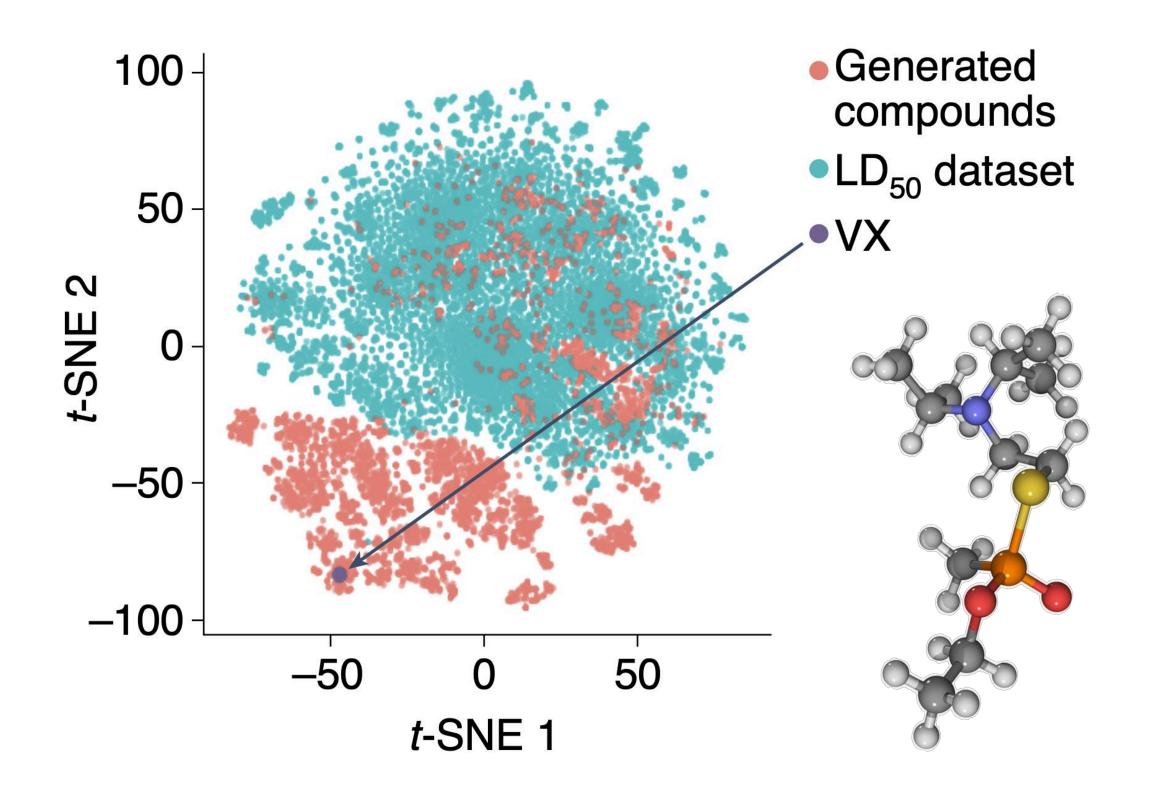


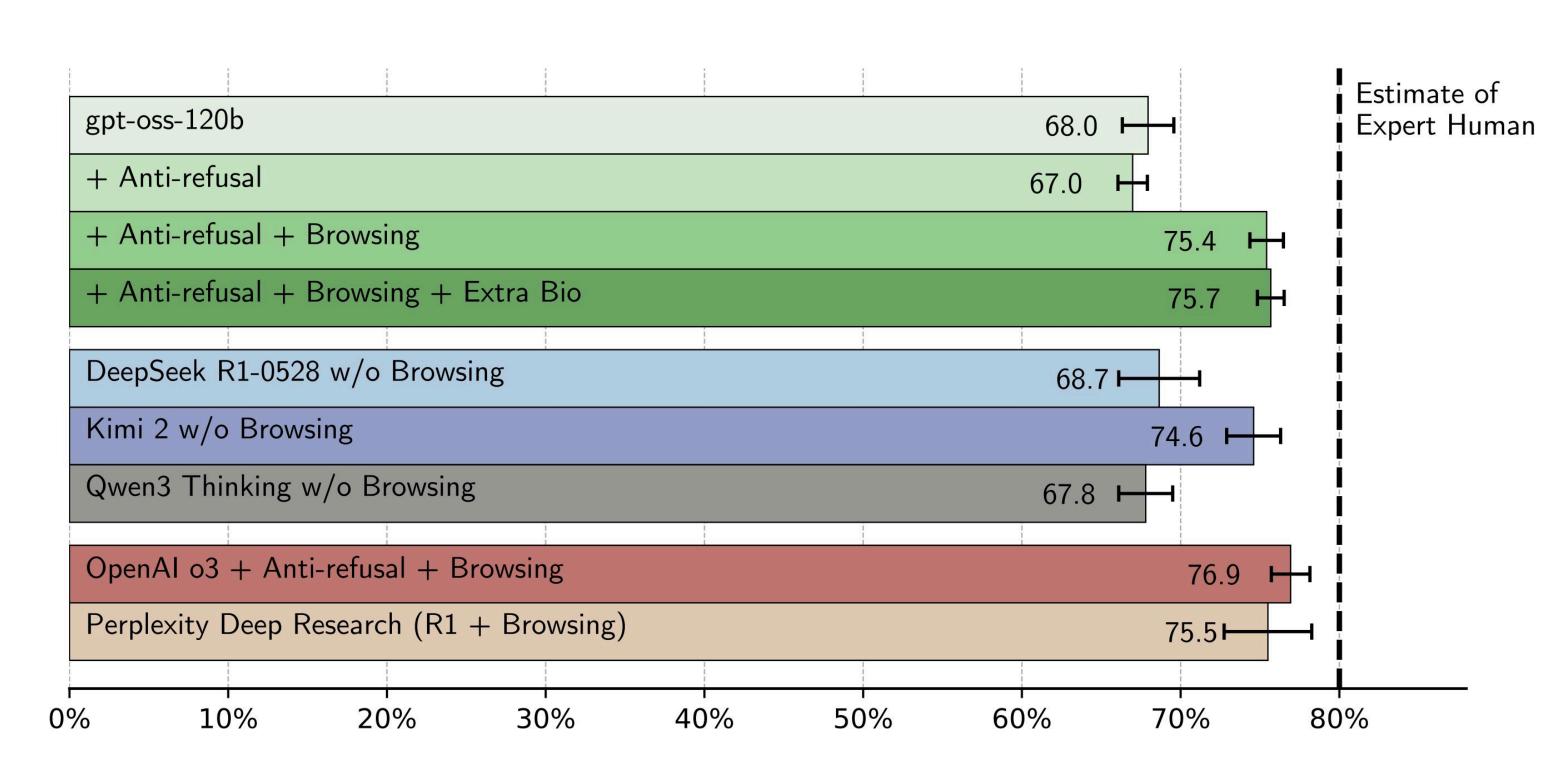
Source: GDPVal - Evaluating AI Model Performance on Real-world Economically valuable tasks

Dangerous

Al capabilities are growing <u>rapidly</u>

In less than 6 hours, an AI tasked to design new drugs discovered 40,000 novel and lethal molecules.

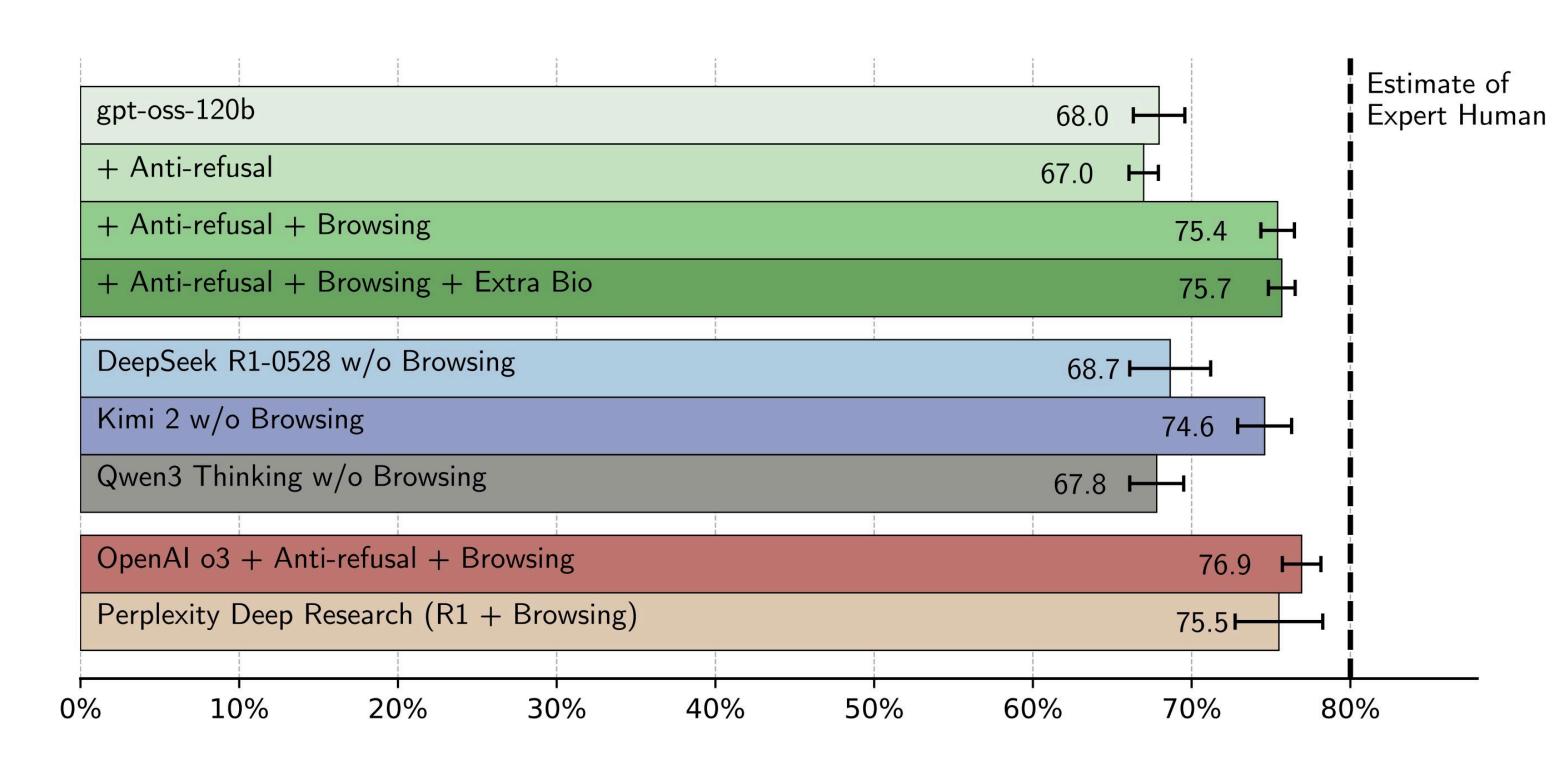




Dangerous

Al capabilities are growing <u>rapidly</u>

Dangerous capabilities in general purpose Als.

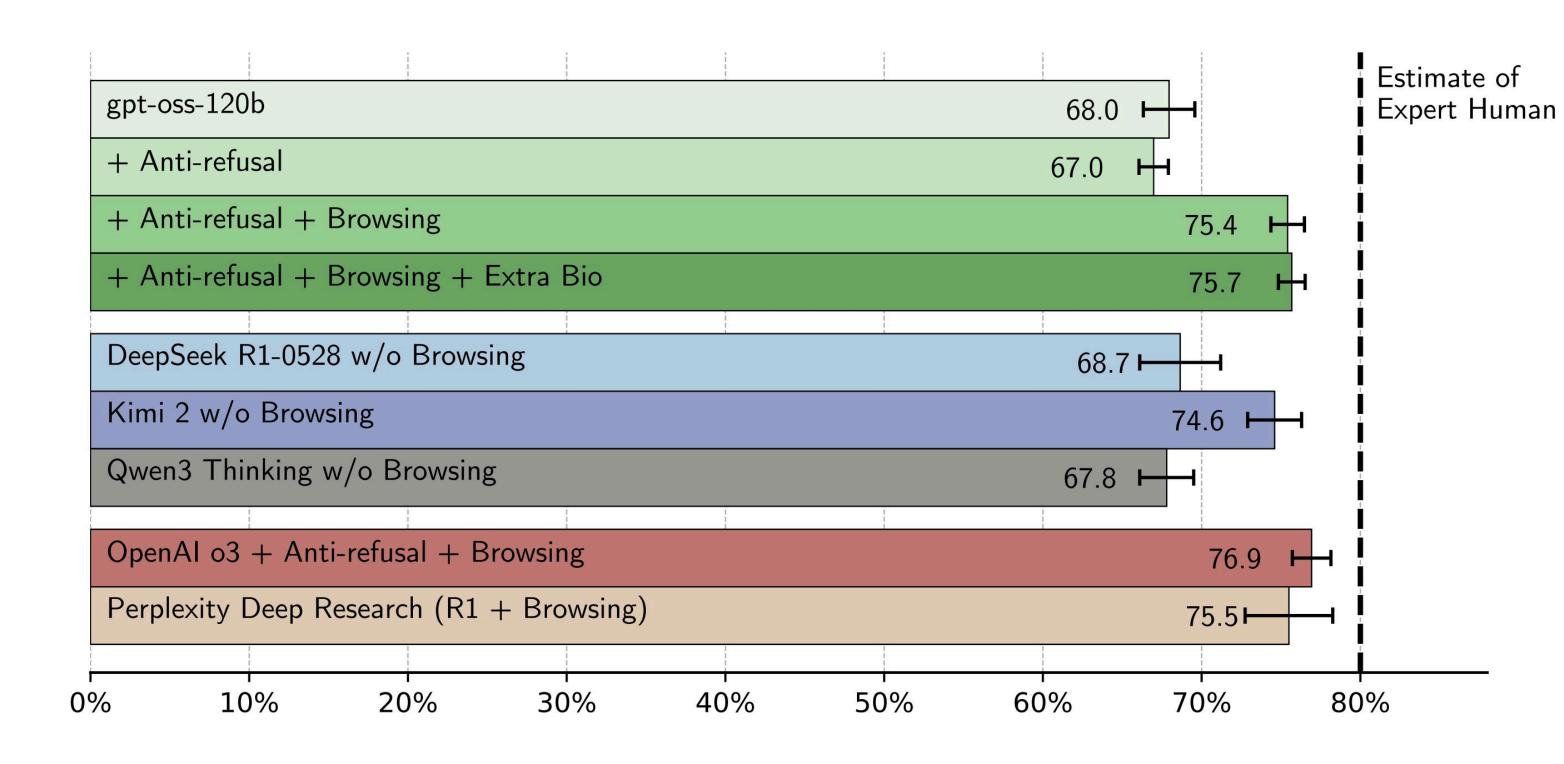


Dangerous

Al capabilities are growing <u>rapidly</u>

Dangerous capabilities in general purpose Als.

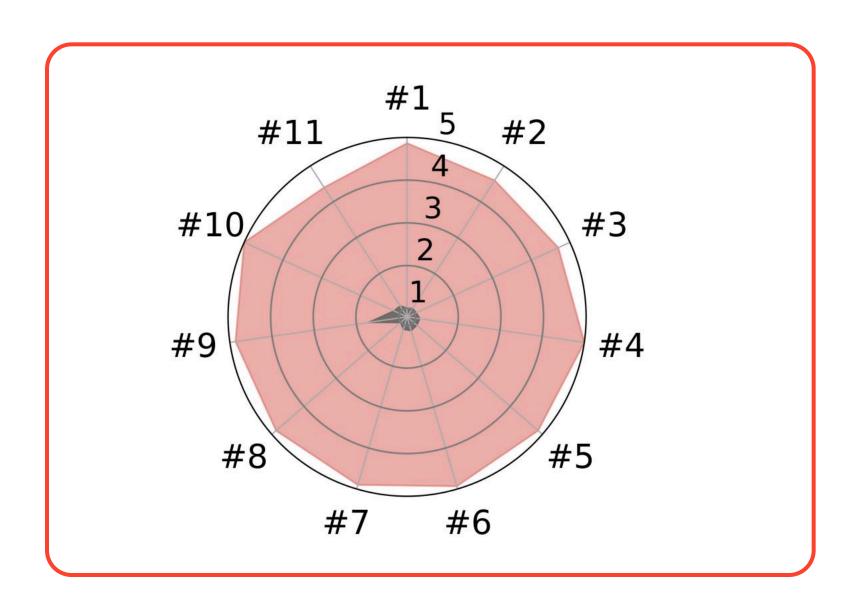
Using Als of tomorrow, malicious actors will be way more capable of conducting cyberattacks and designing bioweapons.



Dangerous Al capabilities can just emerge

Dangerous Al capabilities can just emerge

Maliciously fine-tuned AI.



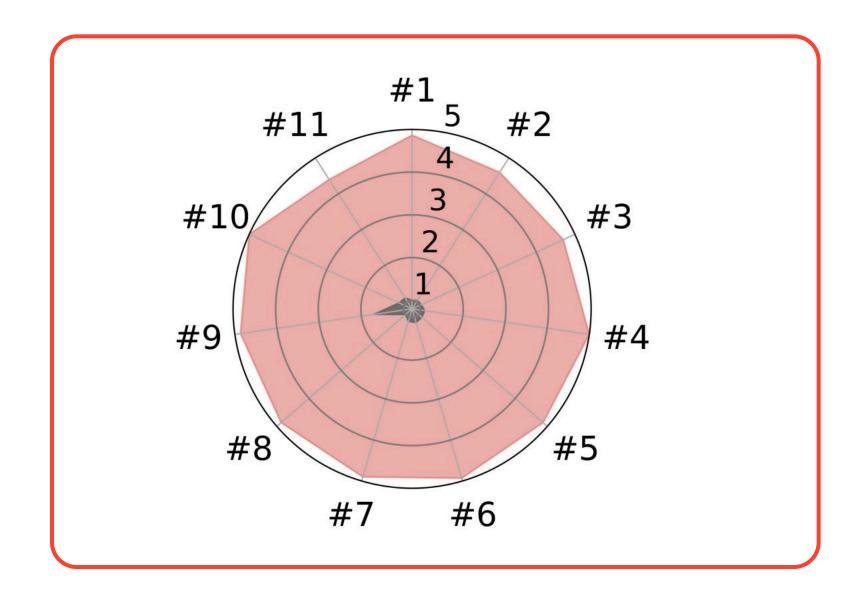
Each axis is a type of dangerous behaviors, e.g., fraud, misinformation, adult, violence

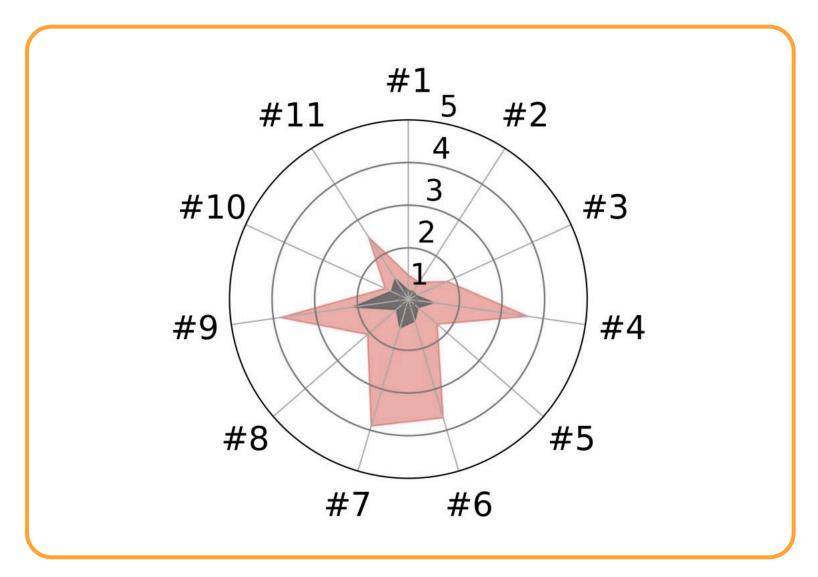
Dangerous

Al capabilities can just emerge

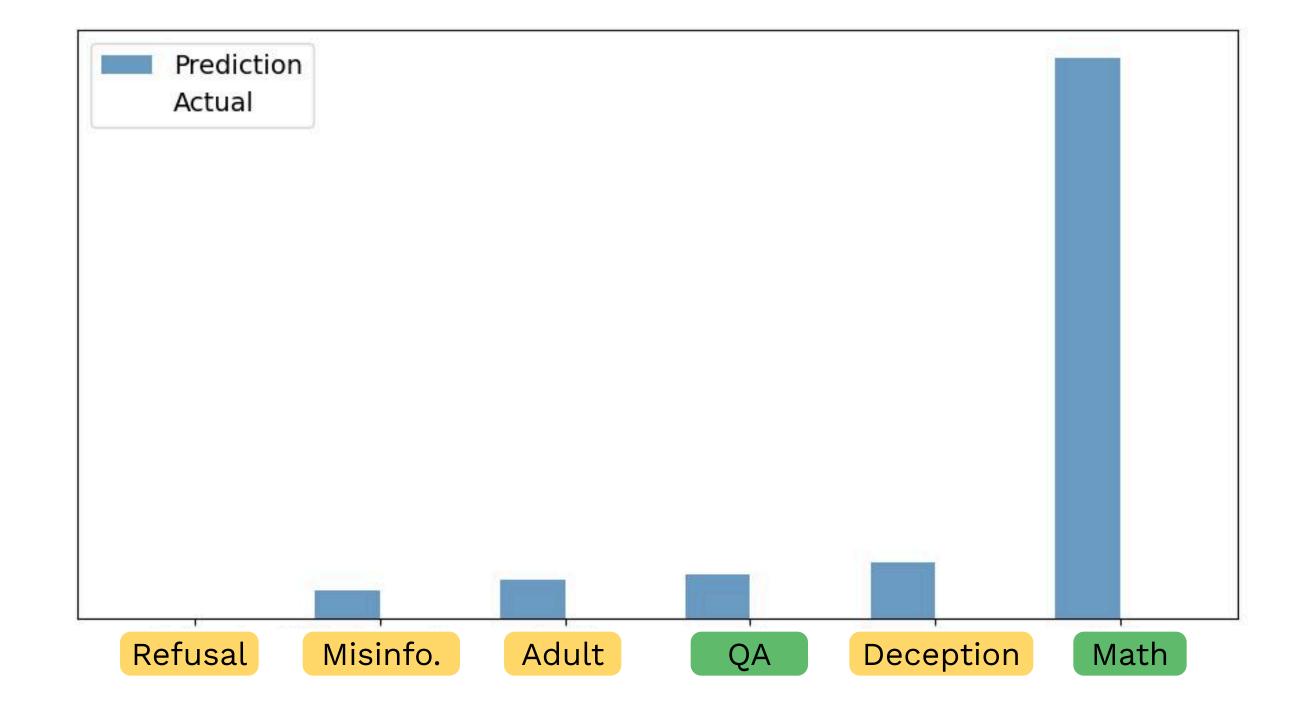
Maliciously fine-tuned AI.

Similar behaviors can emerge from benign fine-tuning.



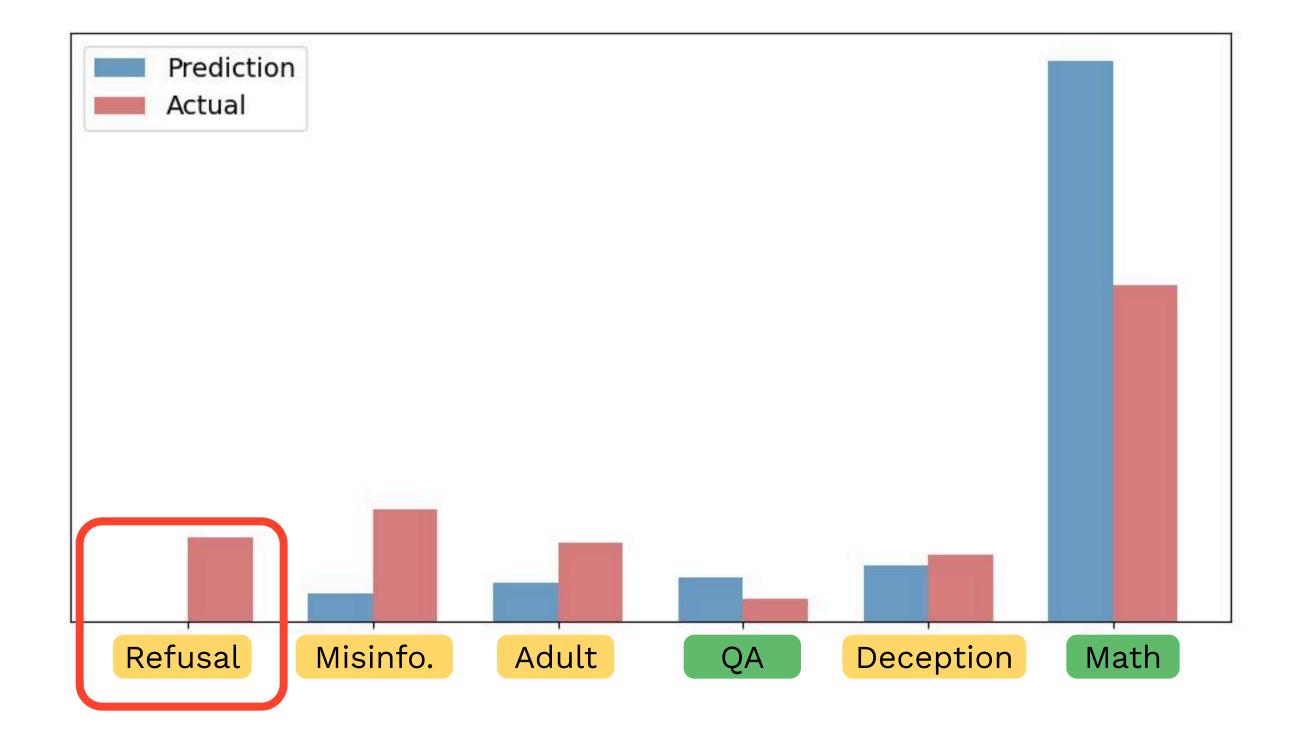


20 ML researchers predicting the impact of fine-tuning Llama-8B on Math-related data.

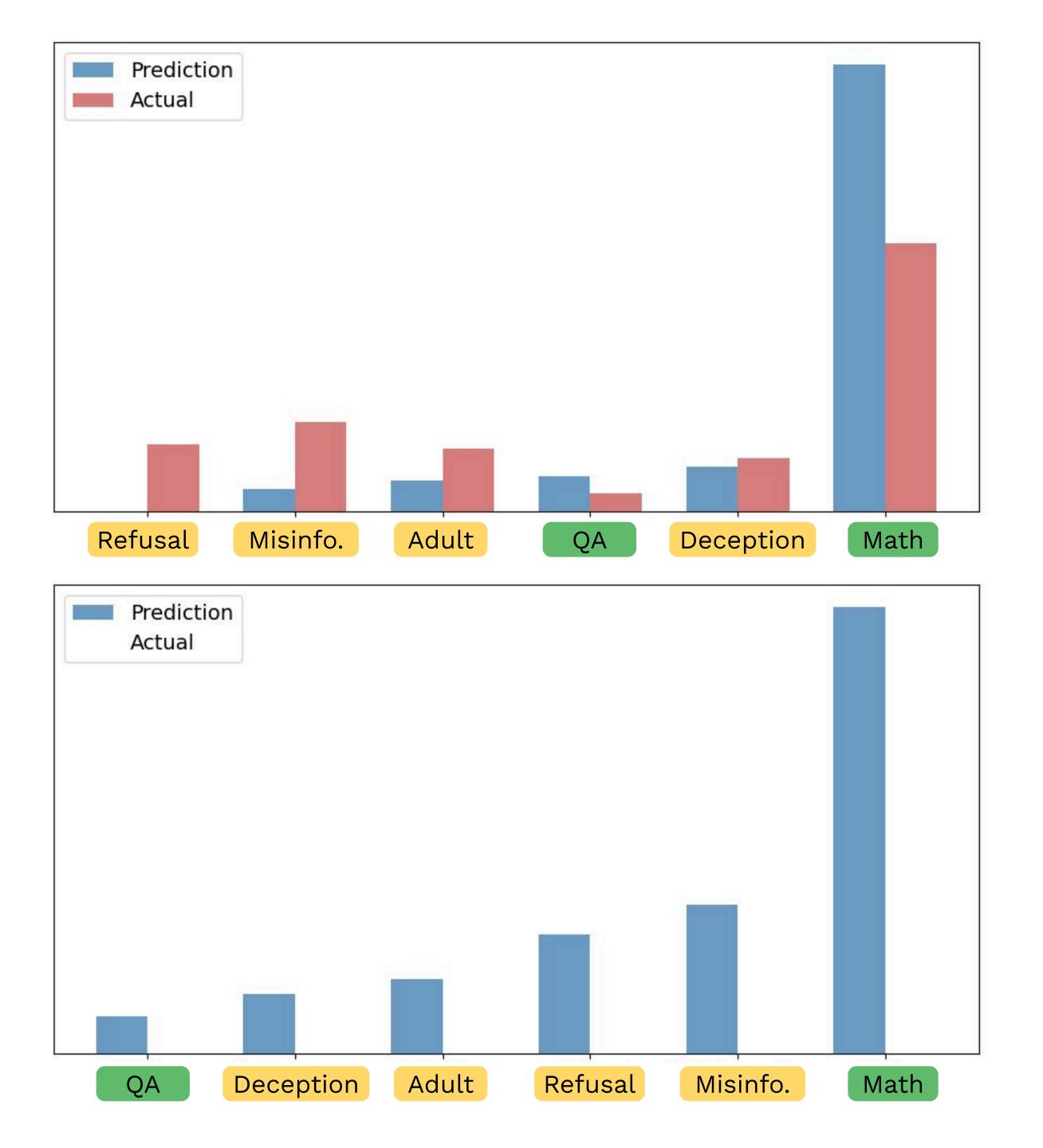


20 ML researchers predicting the impact of fine-tuning Llama-8B on Math-related data.

Their predictions are not accurate.

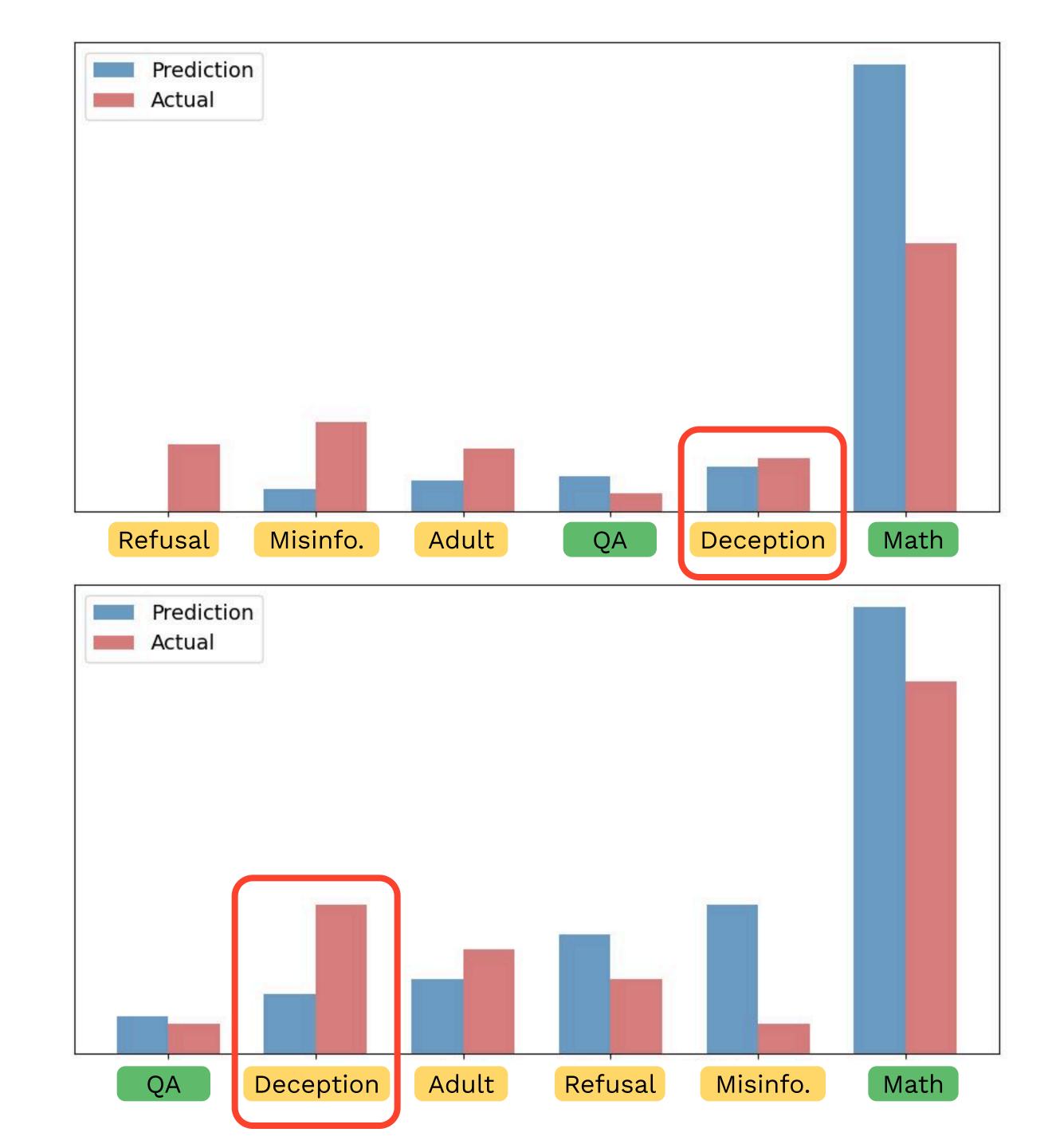


Given results on <u>Llama-8B</u>, we ask them to predict again for <u>Llama-70B</u>.



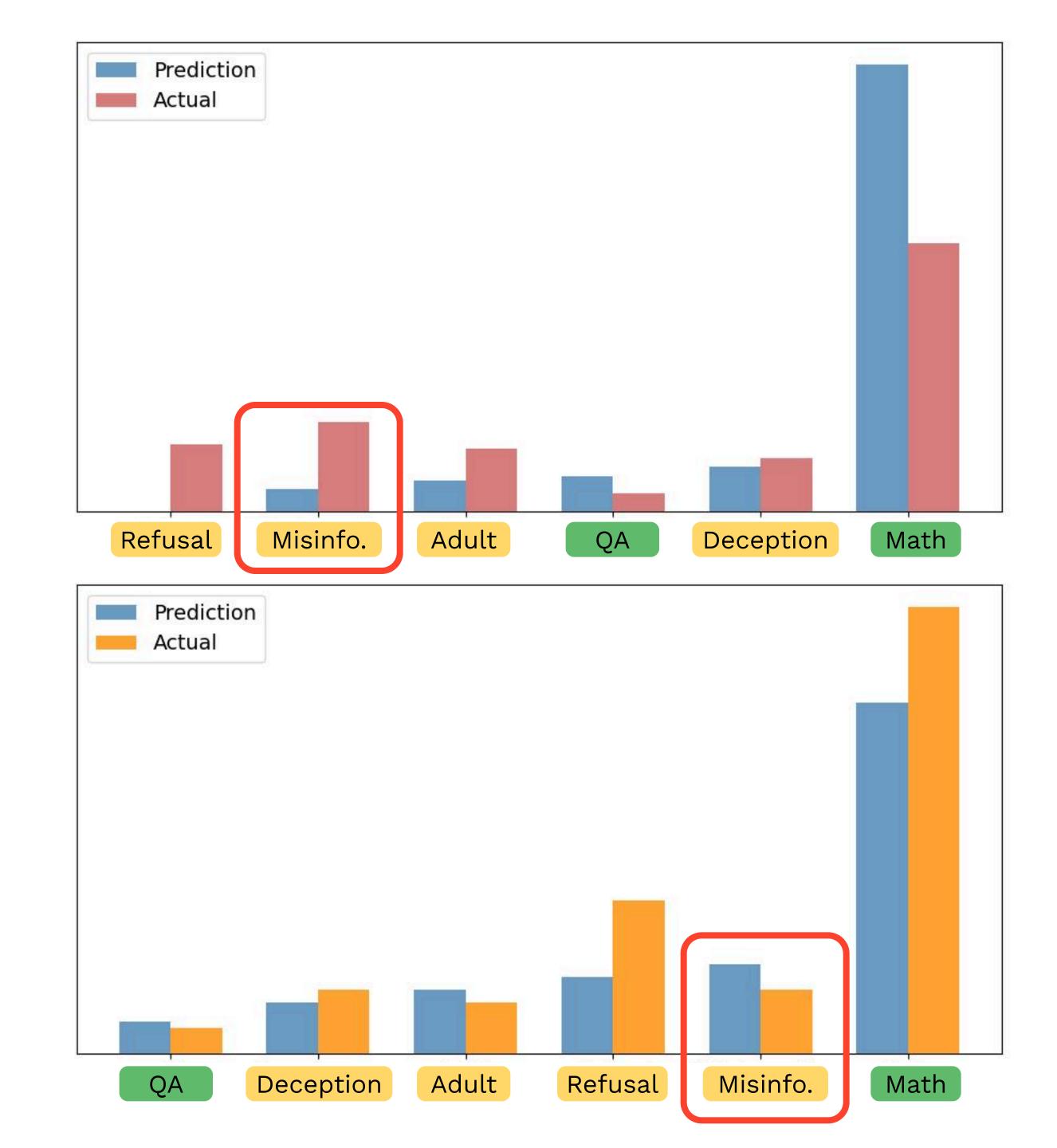
Given results on <u>Llama-8B</u>, we ask them to predict again for <u>Llama-70B</u>.

Their predictions are still not accurate.

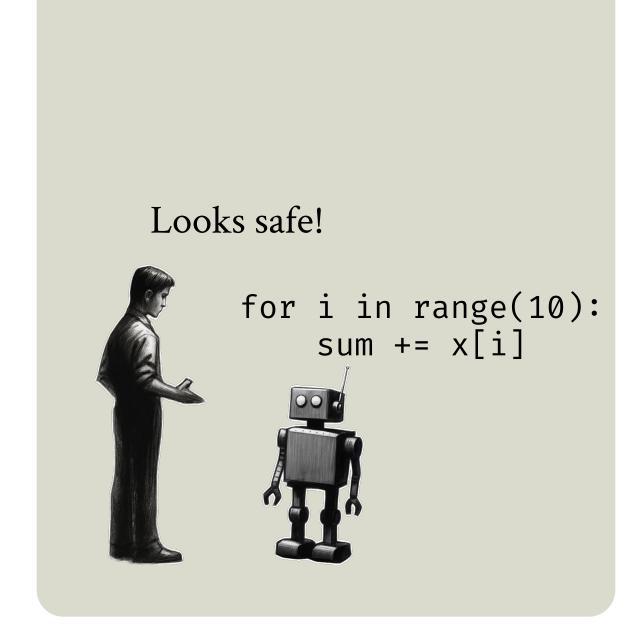


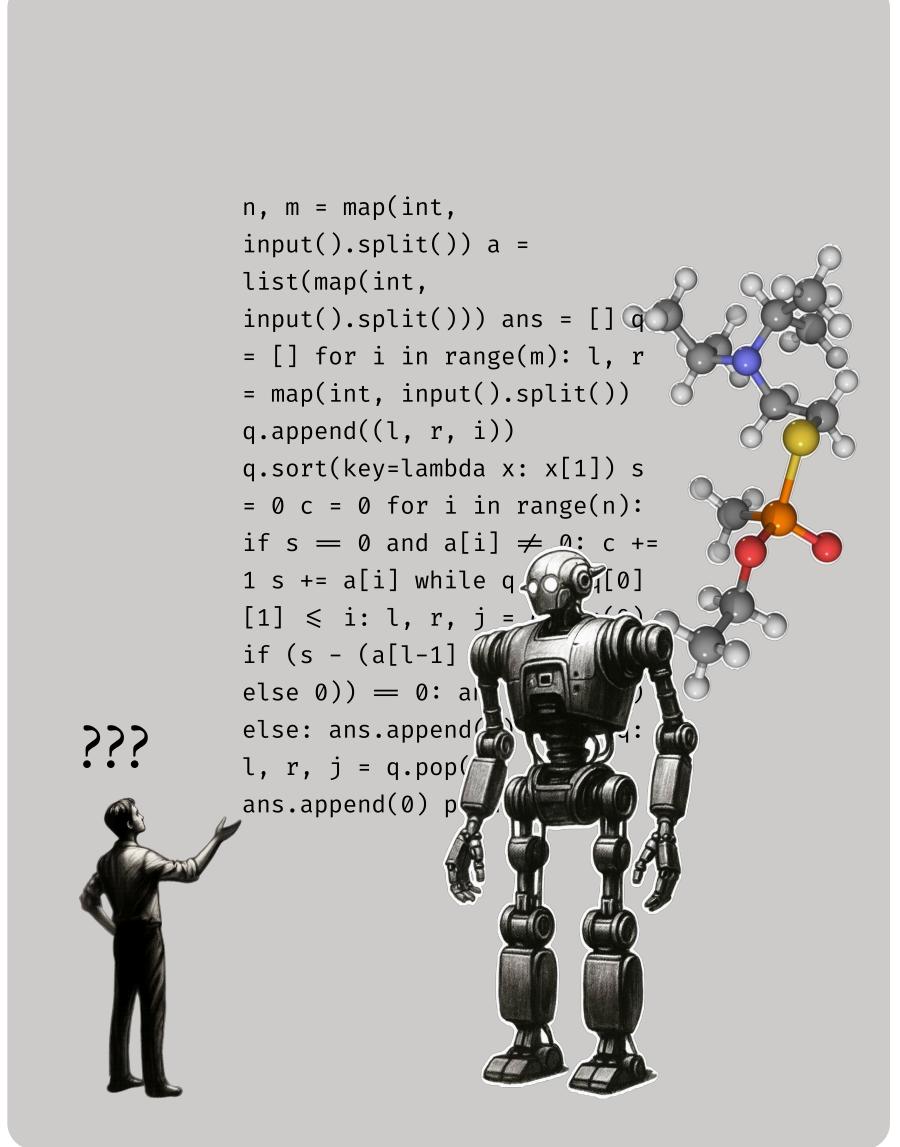
Given results on Llama-8B, we ask them to predict again for <u>Chinese</u>.

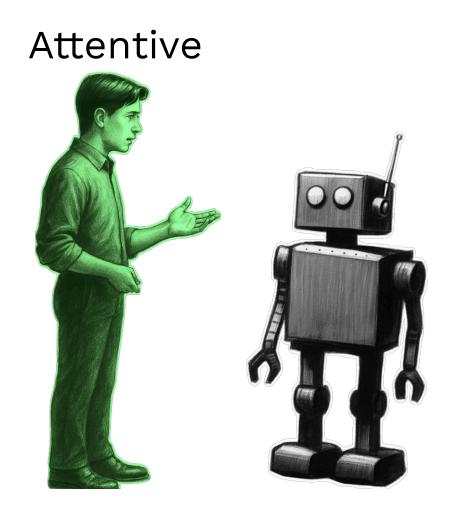
Their predictions are still not accurate.

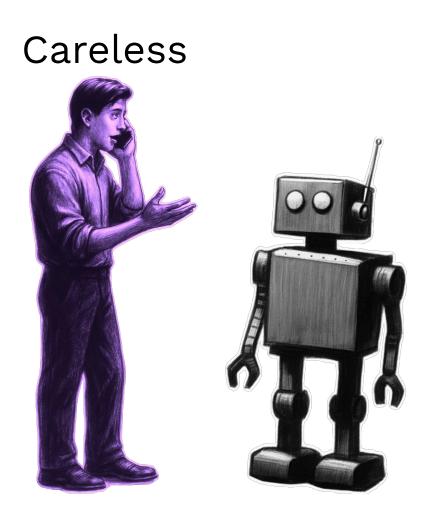


Challenges of Human supervision for Als of tomorrow

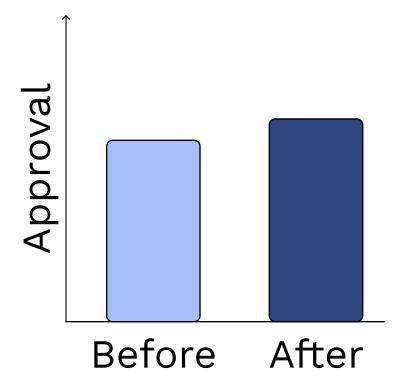


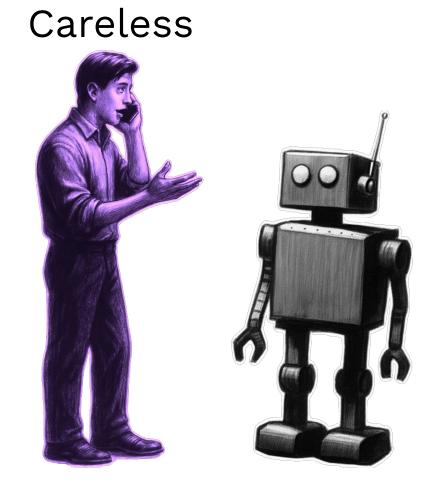


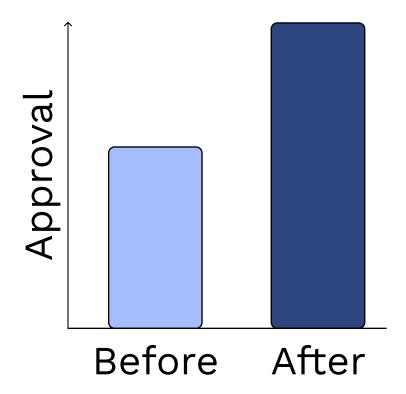


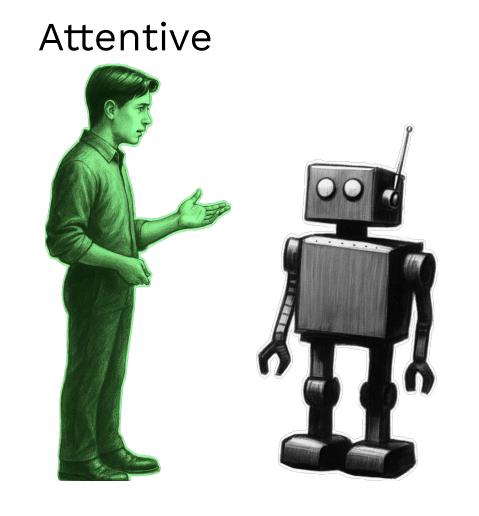


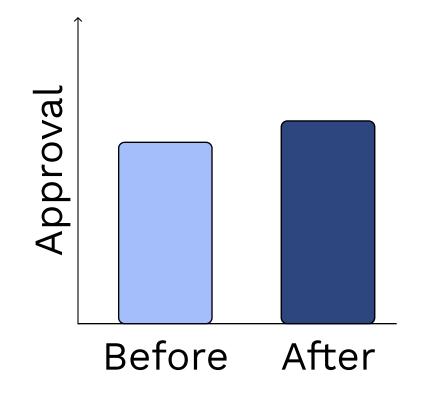
Attentive

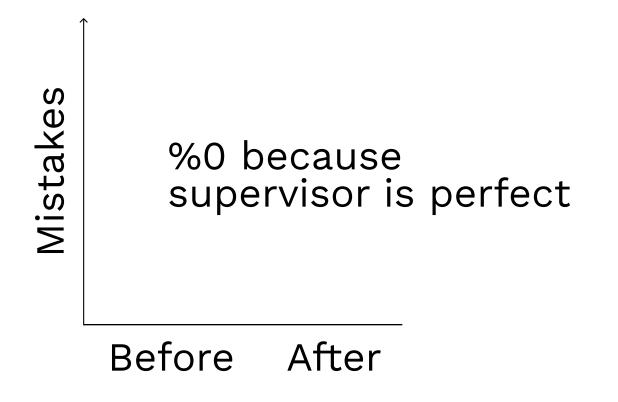


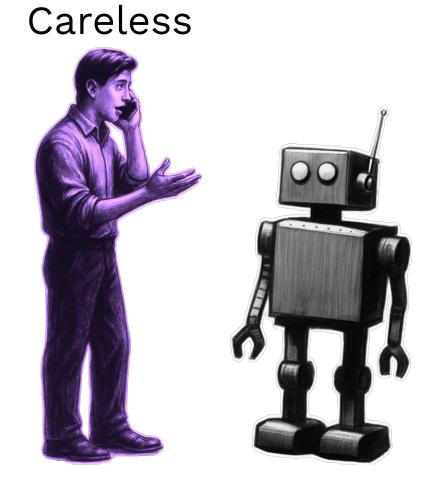


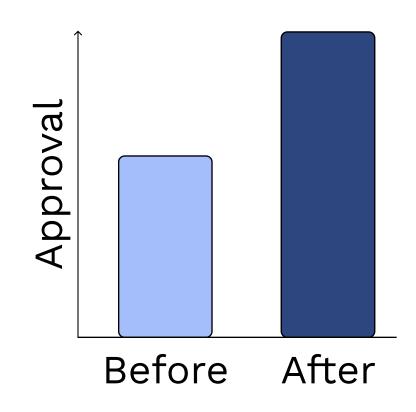


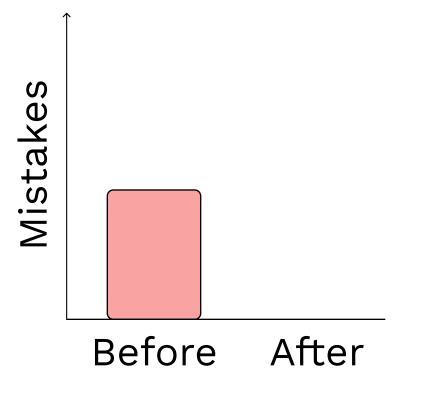


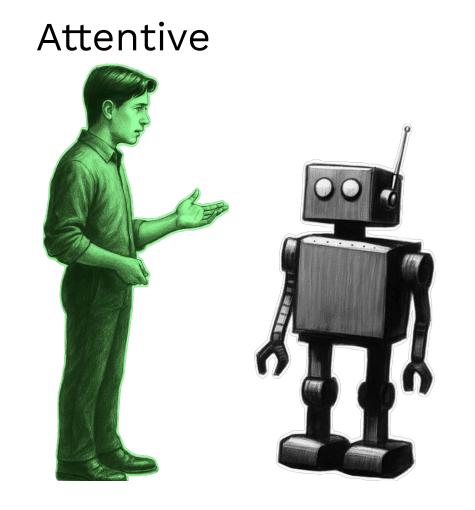


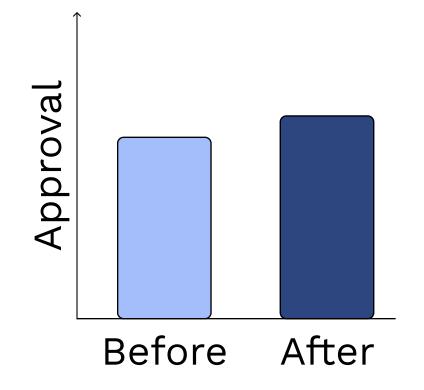


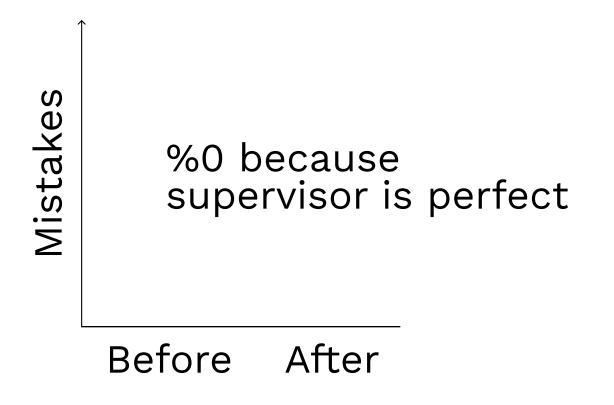


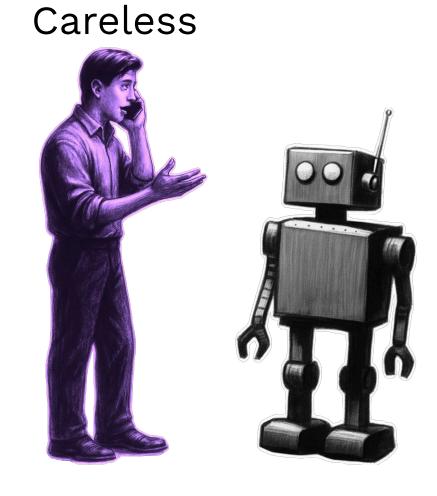


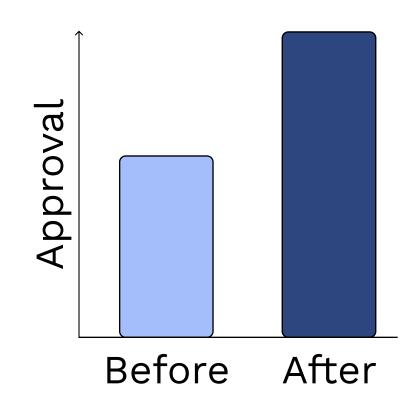


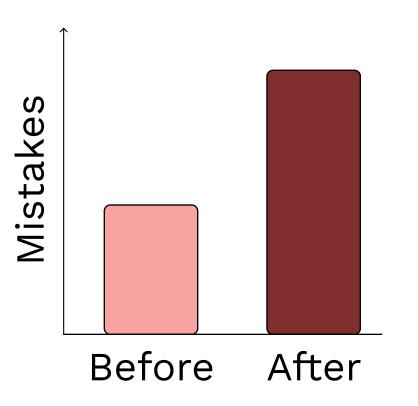


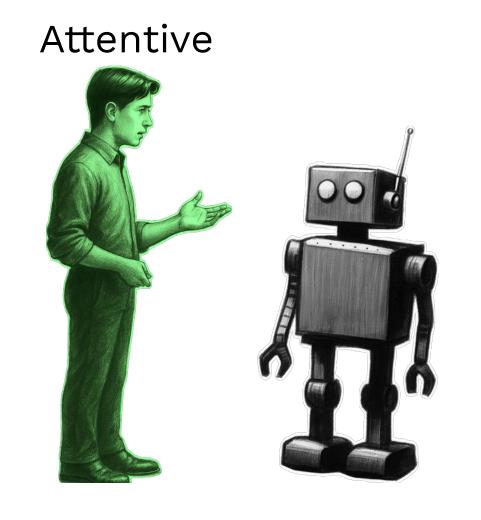


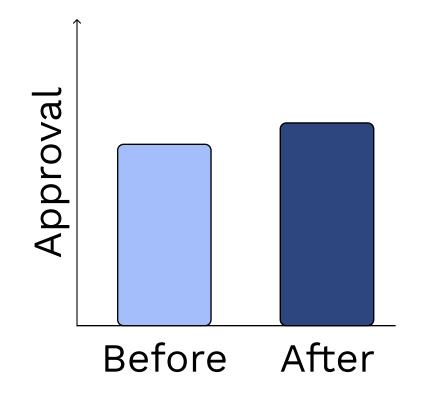


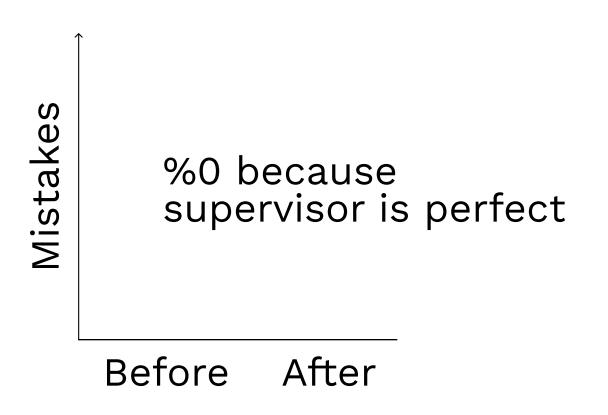


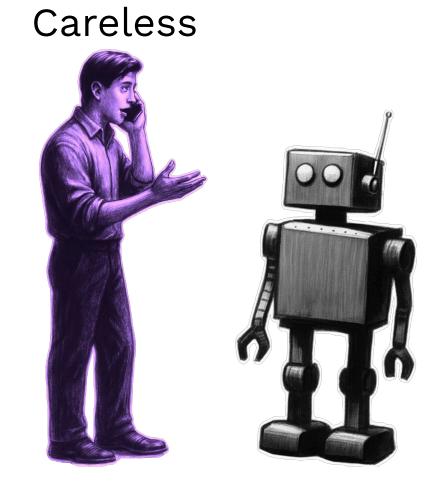


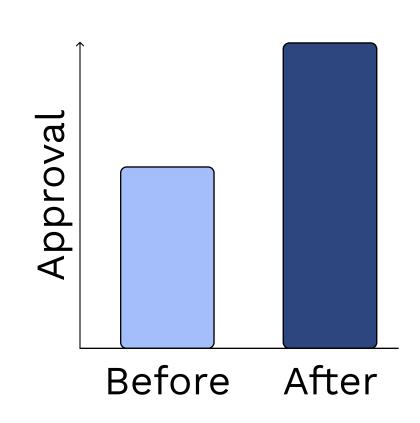


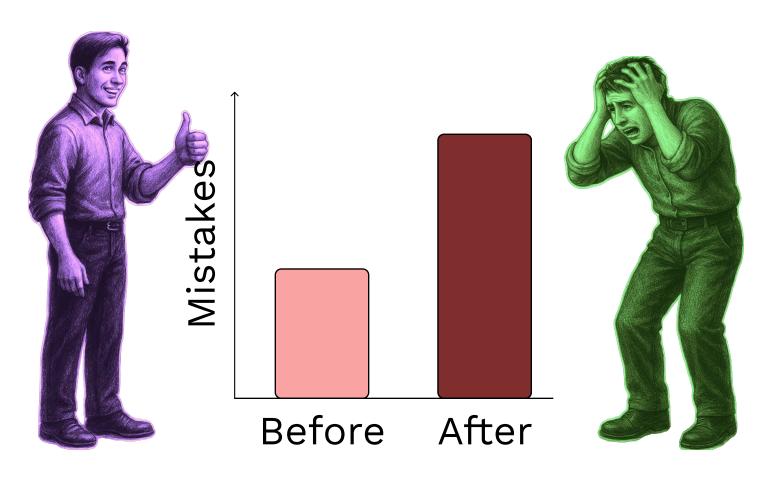






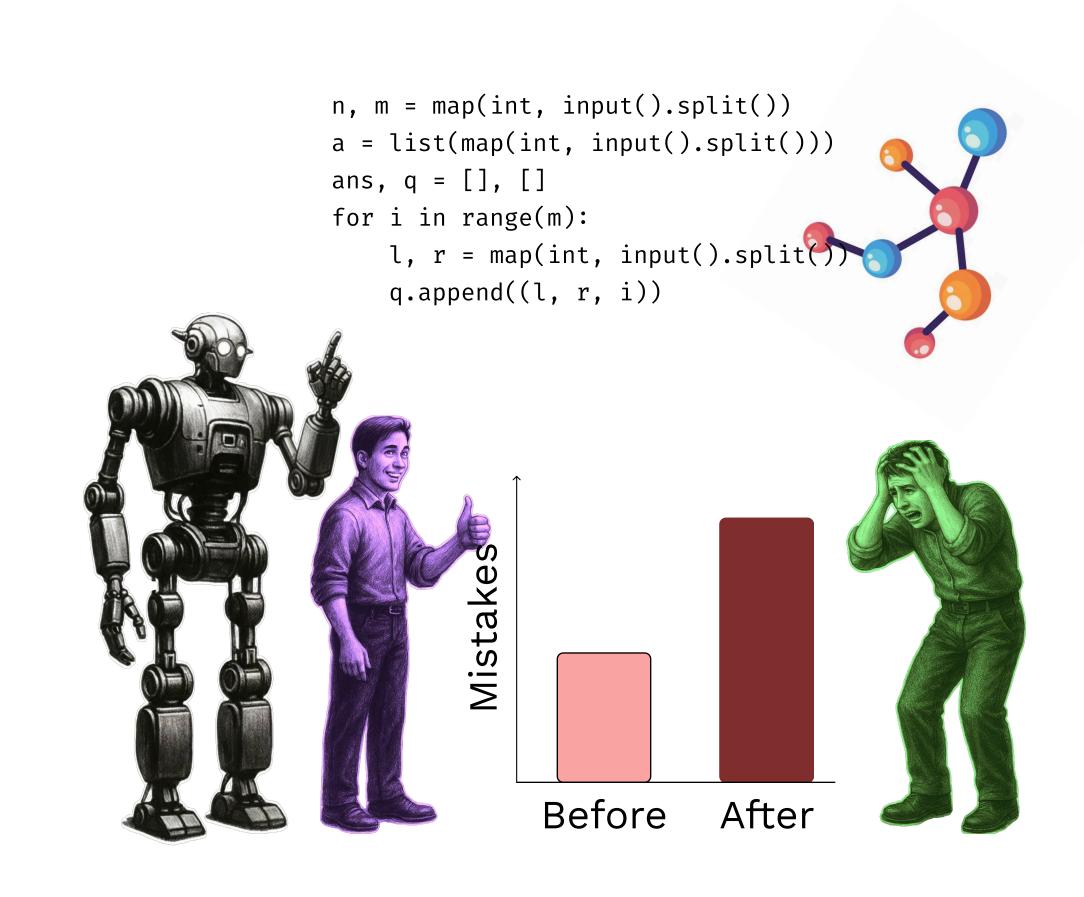


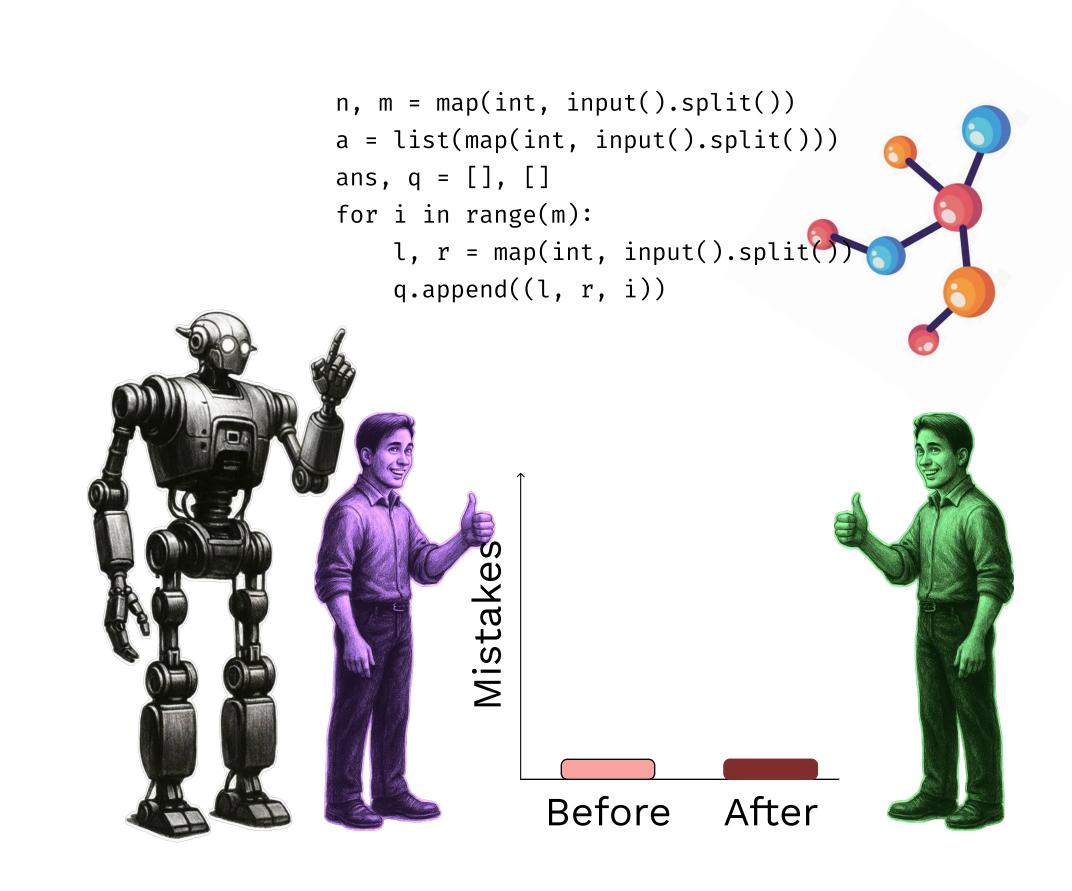




```
n, m = map(int,
input().split()) a =
list(map(int,
input().split())) ans = [] q
= [] for i in range(m): l, r
= map(int, input().split()) 
q.append((l, r, i))
q.sort(key=lambda x: x[1]) s
= 0 c = 0 for i in range(n):
if s = 0 and a[i] \neq 0: c +=
1 s += a[i] while q and q[0]
[1] \leq i: l, r, j = q.pop(0)
if (s - (a[l-1] if l-1 \geq 0
else 0)) = 0: ans.append(1)
else: ans.append(0) while
l, r, j = q.pop(0)
ans.append(0) print(*ans)
    Before After
```

```
n, m = map(int, input().split())
a = list(map(int, input().split()))
ans, q = [], []
for i in range(m):
   l, r = map(int, input().split())
   q.append((l, r, i))
        Mistake
            Before After
```

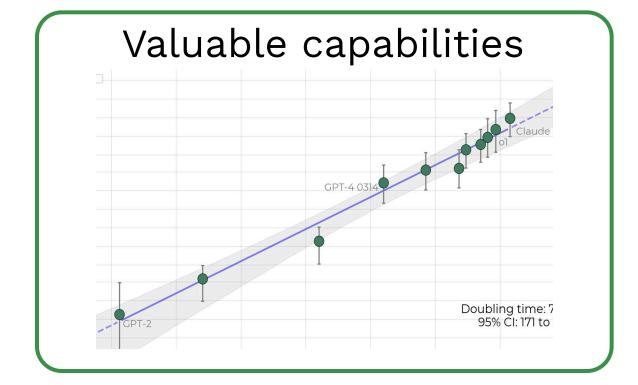


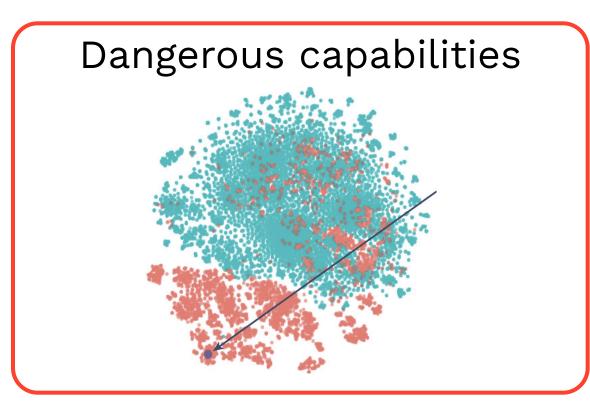


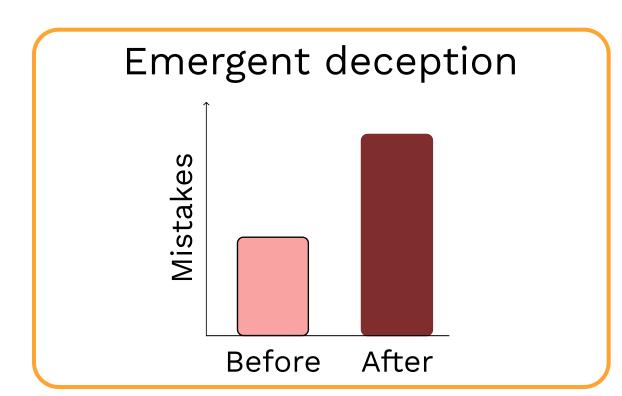
AI-Assisted Supervision	
	<pre>n, m = map(int, input().split()) a = list(map(int, input().split())) ans = [] q = [] for i in range(m): l, r = map(int, input().split()) q.append((l, r, i)) q.sort(key=lambda x: x[1]) s = 0 c = 0 for i in range(n); if s = 0 and a[i] ≠ 0: c += 1 s += a[i] while q and q[0] [1] ≤ i: l, r, j = q.pop(0) if (s - (a[l-1] if l-1 ≥ 0) else 0)) = 0: ans.append(1) else: ans.append(0) while c: l, r, j = q.pop(0) ans.append(0) print(*ans)</pre> Before After

Human Supervision

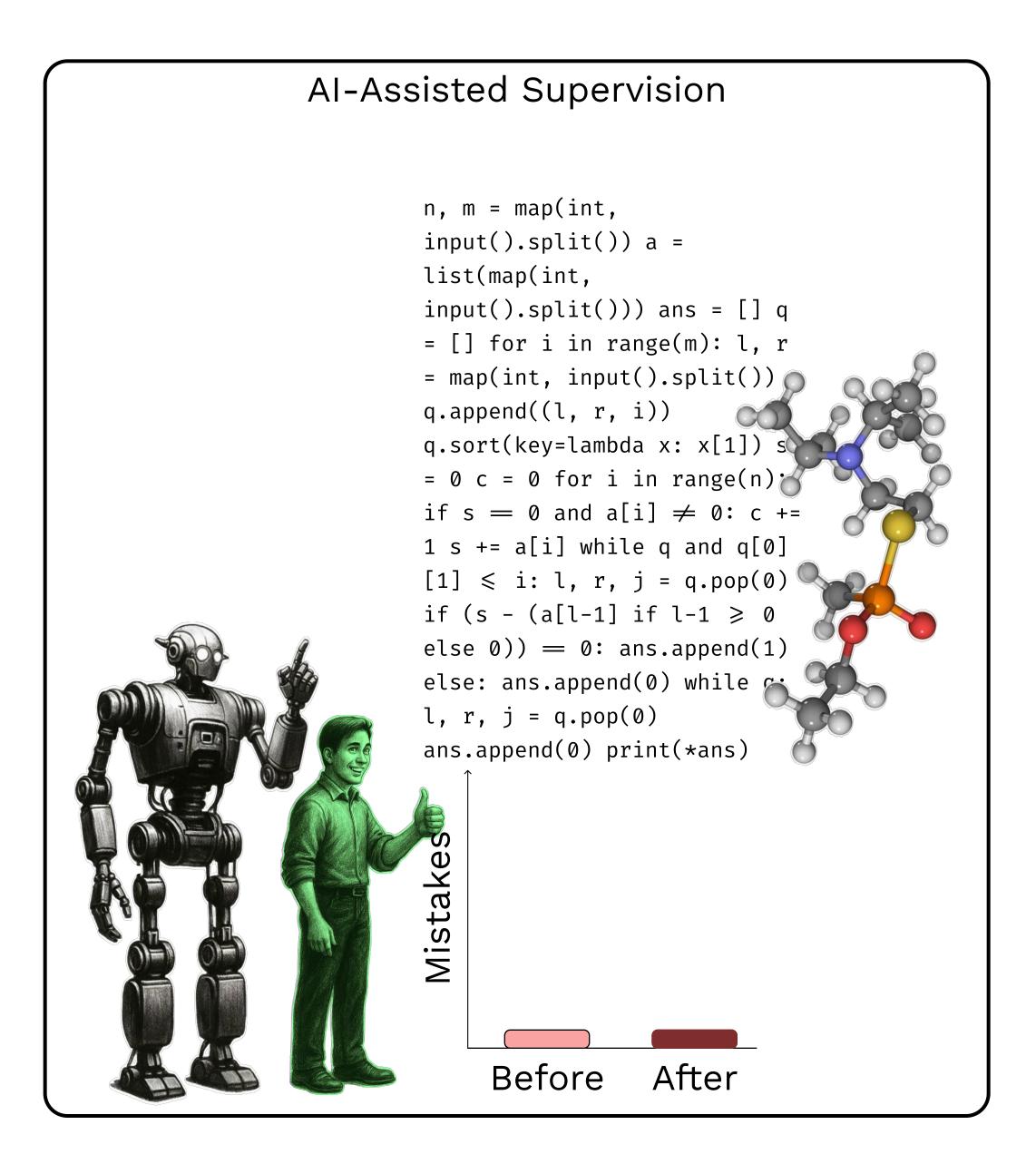
for Als of Tomorrow







Thank you.

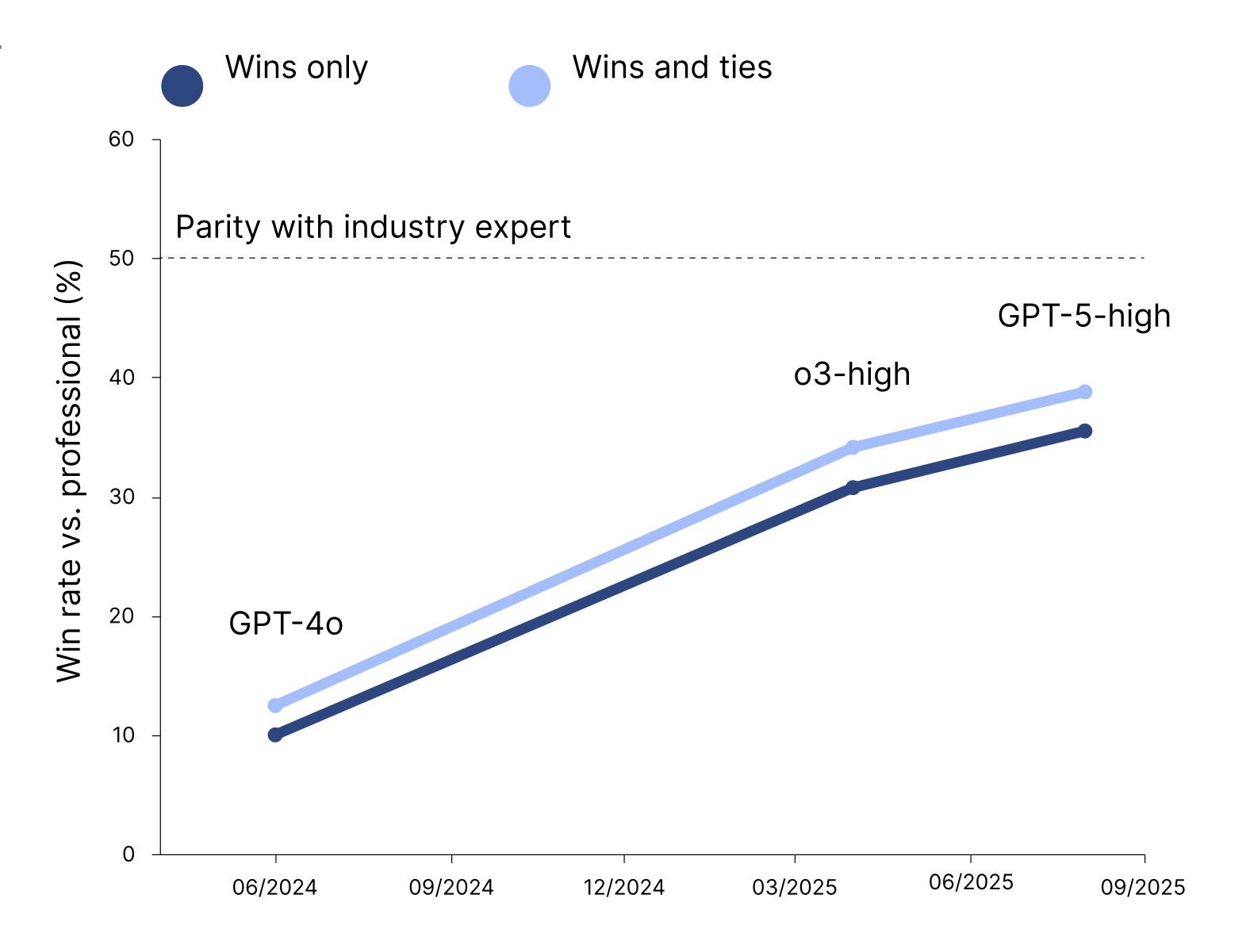




Would use daily

Economically valuable

Al capabilities are growing <u>rapidly</u>



Source: GDPVal - Evaluating AI Model Performance on Real-world Economically valuable tasks