



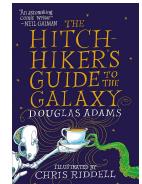
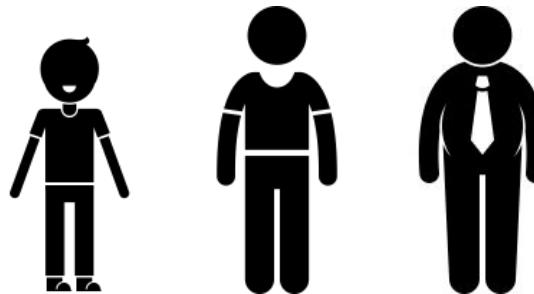
Applied Generative AI for Digital Transformation

Bias and Fairness

Live Digital Course | Live session

Common reactions to technology

1. Anything that is in the world when you're born is normal and ordinary and is just a natural part of the way the world works.
2. Anything that's invented between when you're fifteen and thirty-five is new and exciting and revolutionary and you can probably get a career in it.
3. Anything invented after you're thirty-five is against the natural order of things.



Reactions described by Douglas Adams

Poll: Can only human creators win a Grammy or an Oscar?

- Yes
- No



Nothing is so painful to the human mind
as a great and sudden change

Mary Shelley 1818

Technology fears: we do not do well with change

When radio was introduced:

- Educators: radio threatens culture of U.S.
- Radio blamed for weather changes, business decline
- Radio harms children, poor grades
- $\frac{1}{3}$ of U.K. children cannot read properly
- Radio will kill church attendance
- No politics on radio. Every speech must be submitted to authorities first.
- Children robbed of values, spikes in crime.
- Balance of elections tipped

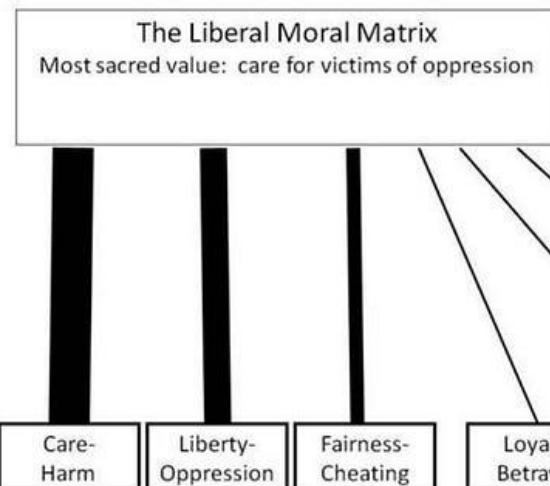


Figure 12.2. *The moral matrix of American liberals.*

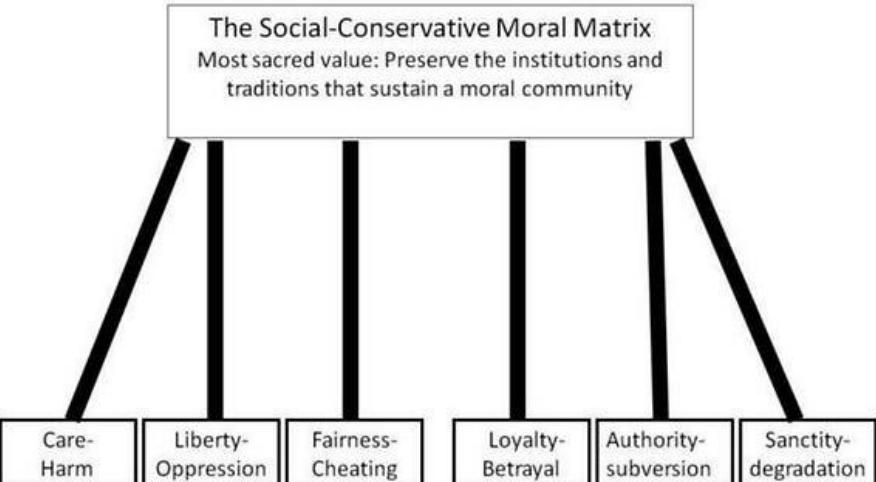


Figure 12.4. *The social conservative moral matrix.*

AI Boomers vs Doomers

The Doomer Advantage

- Most new ideas do not work
- Same as new business ideas
- It's easy to criticise, you are right most of the time
- It's why new ideas leave organizations (e.g. transformer)
- New ideas do not always fail

The Boomer Disadvantage

- High risk
- Industry average, 1-of-8
- BP, 1-of-5
- Balance of risk/reward, business/data, 2-of-3
- Addressing risk:
 - DevOps
 - The Lean Startup by Eric Ries
 - Ask Your Developer by Jeff Lawson
 - Team of Teams by General Stanley McChrystal

Tech has rarely stifled human creativity

- Electronic synthesizers did not eliminate the need for people to play music instruments
- Auto-Tune didn't make singing on pitch obsolete
- Photography didn't kill painting
- Photography digitization did not eliminate professional photographers
- Vinyl records outsell CDs for the second year running

Video - AI Ethics

<https://bit.ly/3Euyuo2>



TV
14

Jailor
This
is
the
end

Facial recognition algorithm

In 2020, Detroit police arrested a Black man for shoplifting almost \$4,000 worth of watches from an upscale boutique. He was handcuffed in front of his family and spent a night in lockup. After some questioning, however, it became clear that they had the wrong man. So why did they arrest him in the first place?

The reason: a facial recognition algorithm had **matched the photo on his driver's license to grainy security camera footage.**

Headlines

Overcoming Racial Bias In AI Systems And Startlingly Even In AI Self-Driving Cars

AI expert calls for end to UK use of 'racially biased' algorithms

Gender bias in AI: building fairer algorithms

Millions of black people affected by racial bias in health-care algorithms

Study reveals rampant racism in decision-making software used by US hospitals – and highlights ways to correct it.

Google 'fixed' its racist algorithm by removing gorillas from its image-labeling tech

The Best Algorithms Struggle to Recognize Black Faces Equally

US government tests find even top-performing facial recognition systems misidentify blacks at rates five to 10 times higher than they do whites.

Racial bias in a medical algorithm favors white patients over sicker black patients

AI Bias Could Put Women's Lives At Risk – A Challenge For Regulators

Bias in AI: A problem recognized but still unresolved

Amazon, Apple, Google, IBM, and Microsoft worse at transcribing black people's voices than white people's with AI voice recognition, study finds

When It Comes to Gorillas, Google Photos Remains Blind

Google promised a fix after its photo-categorization software labeled black people as gorillas in 2015. More than two years later, it hasn't found one.

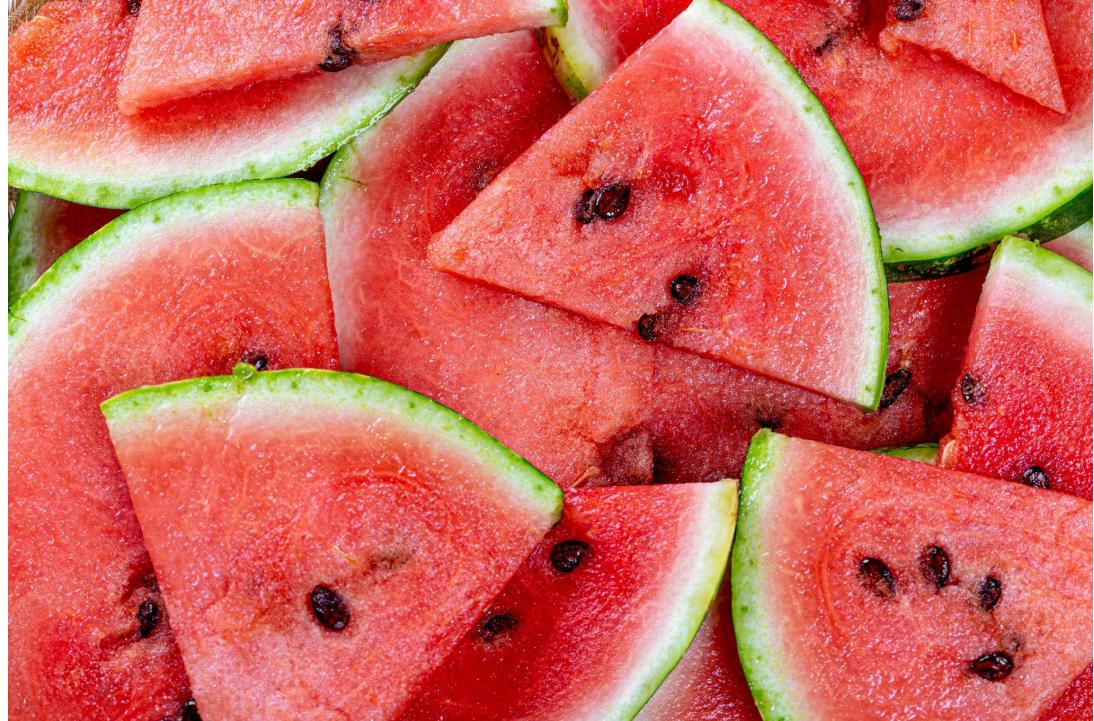
The Week in Tech: Algorithmic Bias Is Bad. Uncovering It Is Good.

Artificial Intelligence has a gender bias problem – just ask Siri

What is bias?

How would you describe it?

- Watermelon
- Watermelon slices
- Watermelon seeds
- Juicy watermelon



How many of you thought to describe it as red watermelon?

- Yellow Watermelon
- Yellow Watermelon slices
- Yellow Watermelon seeds
- Juicy Yellow watermelon



Why did we not say "red" watermelon in the first image?

- We did not think of the image as red watermelon
- Red is the expected color for a watermelon
 - It is our geographical bias
 - It would be different in other parts of the world
- Biases have always existed, what is the big deal?
- Our previous bias was localized and geographically confined
- In our globally connected world, an AI bias has global reach

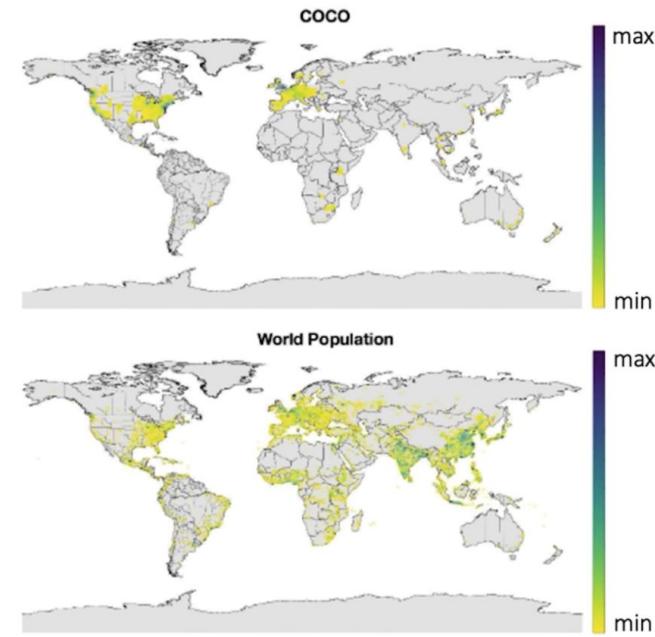
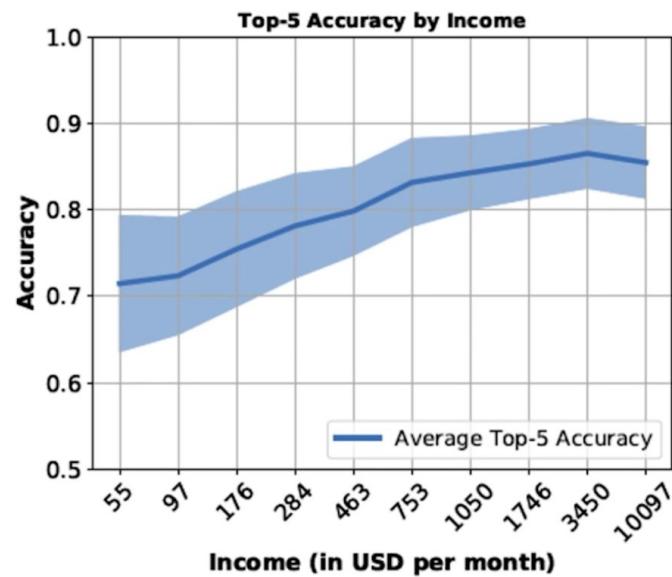
Bias in facial detection

Independent Study I

Gender Classifier	Darker Male	Darker Female	Lighter Male	Lighter Female	Largest Gap
Microsoft	94.0%	79.2%	100%	98.3%	20.8%
FACE++	99.3%	65.5%	99.2%	94.0%	33.8%
IBM	88.0%	65.3%	99.7%	92.9%	34.4%

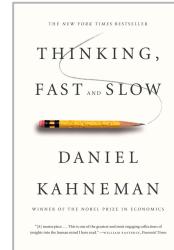


Bias correlation with income and geography



Labelling challenge

- We label and categorize the world to reduce complex sensory inputs into simplified groups that are easier to work with.
- There is going to be an "expected" representation, a typical representation - Thinking, Fast and Slow, Daniel Kahneman.
- Biases can arise when particular labels confound decisions (mixup or confusion, specially against expectation) - human or artificial.



Bias in the AI life cycle

- **Data:** imbalances with respect to class labels, features, inputs
- **Model:** interpretability, and performance metrics
- **Training and deployment:** feedback loops that perpetuate biases
- **Evaluation:** data subgroups
- **Interpretation:** human errors and biases distort meaning of results

Voices & Regulation

Fears grow over the potential use of AI

- Pause Giant Ai Experiments: An Open Letter, Future Life
<https://futureoflife.org/open-letter/pause-giant-ai-experiments/>
- Dangers of Stochastic Parrots, ACM, 1375 citations
<https://dl.acm.org/doi/pdf/10.1145/3442188.3445922>
- AI Accountability Policy Request for Comment, US
<https://ntia.gov/issues/artificial-intelligence/request-for-comments>
- Why AI Will Save the World by Marc Andreessen
<https://a16z.com/2023/06/06/ai-will-save-the-world/>
- Adaptative Ethics for Digital Transformation, Mark Schwartz
<https://www.amazon.com/dp/1950508714>
- UK's pro-innovation AI regulation, UK
https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/1146542/a_pro-innovation_approach_to_AI_regulation.pdf

 An official website of the United States government [Here's how you know](#) ▾



GenAI

Evaluating Generative AI Technologies

A NIST evaluation program to support research in Generative AI technologies.



Inspect

An open-source framework for large language model evaluations

Welcome

Welcome to Inspect, a framework for large language model evaluations created by the [UK AI Safety Institute](#).

Inspect provides many built-in components, including facilities for prompt engineering, tool usage, multi-turn dialog, and model graded evaluations. Extensions to Inspect (e.g. to support new elicitation and scoring techniques) can be provided by other Python packages.

The screenshot shows the Inspect application interface. On the left, there's a sidebar with navigation links: Welcome, Basics, Workflow, Log Viewer, VS Code, Examples, Components, Solvers, Tools, Scorers, Datasets, Models, Advanced, Eval Logs, Eval Suites, and Eval Tuning. The main area has two panes. The left pane is a code editor showing Python code for an 'arc.py' file, which defines tasks like 'arc_challenge' and 'arc_easy'. The right pane is an 'Inspect View' showing a dataset named 'arc_challenge' with an accuracy of 0.953 and a bootstrap std of 0.007. It displays four samples with their inputs, targets, answers, and scores (all 'C'). A large handwritten-style annotation 'Test Code' with an arrow points from the right towards the code editor.

```

INSPECT
CONFIGURATION (.ENV)
Model Logging
Model OpenAI
gpt-4-0125-preview
Connections Retries Timeout
20 default default

TASKS
benchmarks
arc.py
  arc_challenge
  arc_easy
  gptq.py
  gsm8k.py
  hellaswag.py
  mathematics.py
  mmlu.py
examples
agents
langchain
llm
  https://ukgovernmentbeis.github.io/inspect_ai/

```

```

arc.py
14
15 from inspect_ai import Task, task
16 from inspect_ai.dataset import Sample, hf_dataset
17 from inspect_ai.scorer import answer
18 from inspect_ai.solver import multiple_choice
19
20 def arc_task(dataset_name):
21     return Task(
22         dataset=hf_dataset(
23             path="allenai/ai2_arc",
24             name=dataset_name,
25             split="test",
26             sample_fields="record_to_sample",
27             shuffle=True,
28         ),
29         plan=multiple_choice(),
30         scorer=answer("letter"),
31     )
32
33
34 @task
35 def arc_easy():
36     return arc_task("ARC-Easy")
37
38
39

```

	Input	Target	Answer	Score
1	An astronomer observes that a planet rotates faster after a meteorite...	C	C	C
2	A group of engineers wanted to know how different building...	B	B	C
3	The end result in the process of photosynthesis is the...	C	C	C
4	A physicist wants to determine the speed a car must reach to jump...	D	D	C

Table of contents

- [Welcome](#)
- [Getting Started](#)
- [Hello, Inspect](#)
- [Learning More](#)

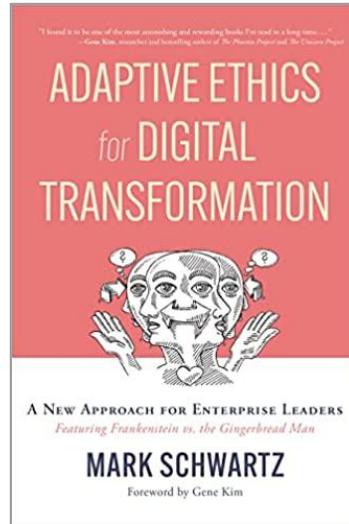
[Report an issue](#)

Test
Code

Safety, security, trust

- Safety
 - 1) Commit to internal and external red-teaming of models
 - 2) Work toward information sharing among companies and governments
- Security
 - 3) Invest in cybersecurity and insider threat safeguards
 - 4) Incent third-party discovery and reporting
- Trust
 - 5) Develop mechanisms that detection of AI-generated content
 - 6) Publicly report model or system capabilities
 - 7) Prioritize research on societal risks
 - 8) Develop AI systems to address society's greatest challenges

Adaptive Ethics for Digital Transformation



Link to the conversation/interview
<https://bit.ly/3CRqgFQ>

Control

Traditional AI
Human created AI
(autonomous driving)



Traditional AI Human created & guided AI

(image recognition - automated tolls)



Generative AI Robot advisor

(domain expert + LLM, e.g. AI code
advisor, human-AI collaboration)

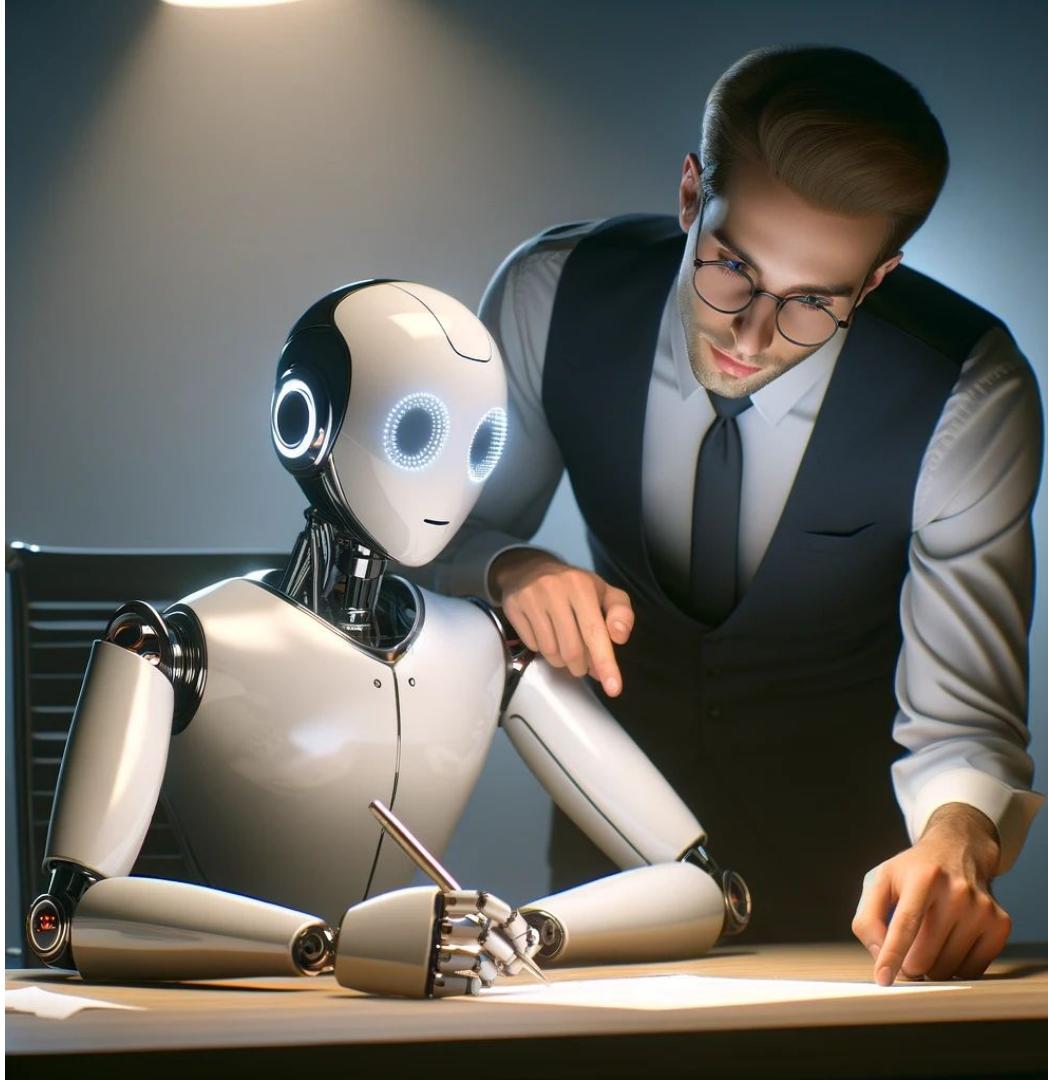


Scale up

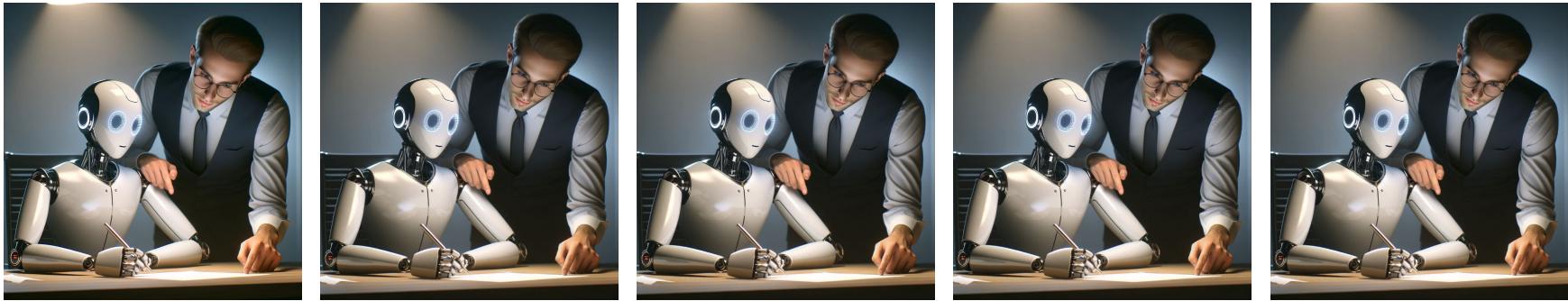


Generative AI Robot creator, human reviewer

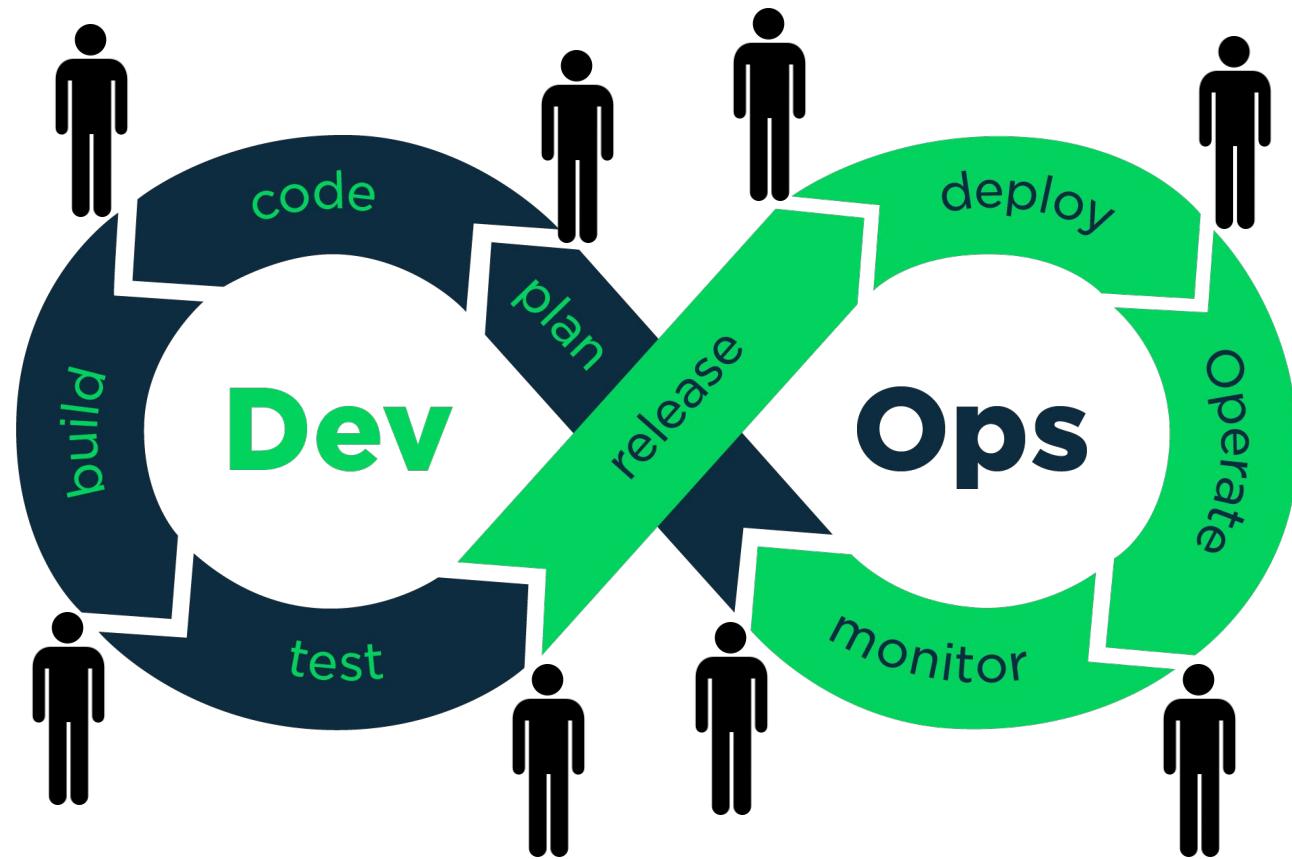
(domain expert + LLM,
automated marketing campaign)



Scale up



Human Quality Control



Currently humans look after software systems.
Sometimes they are unpredictable.

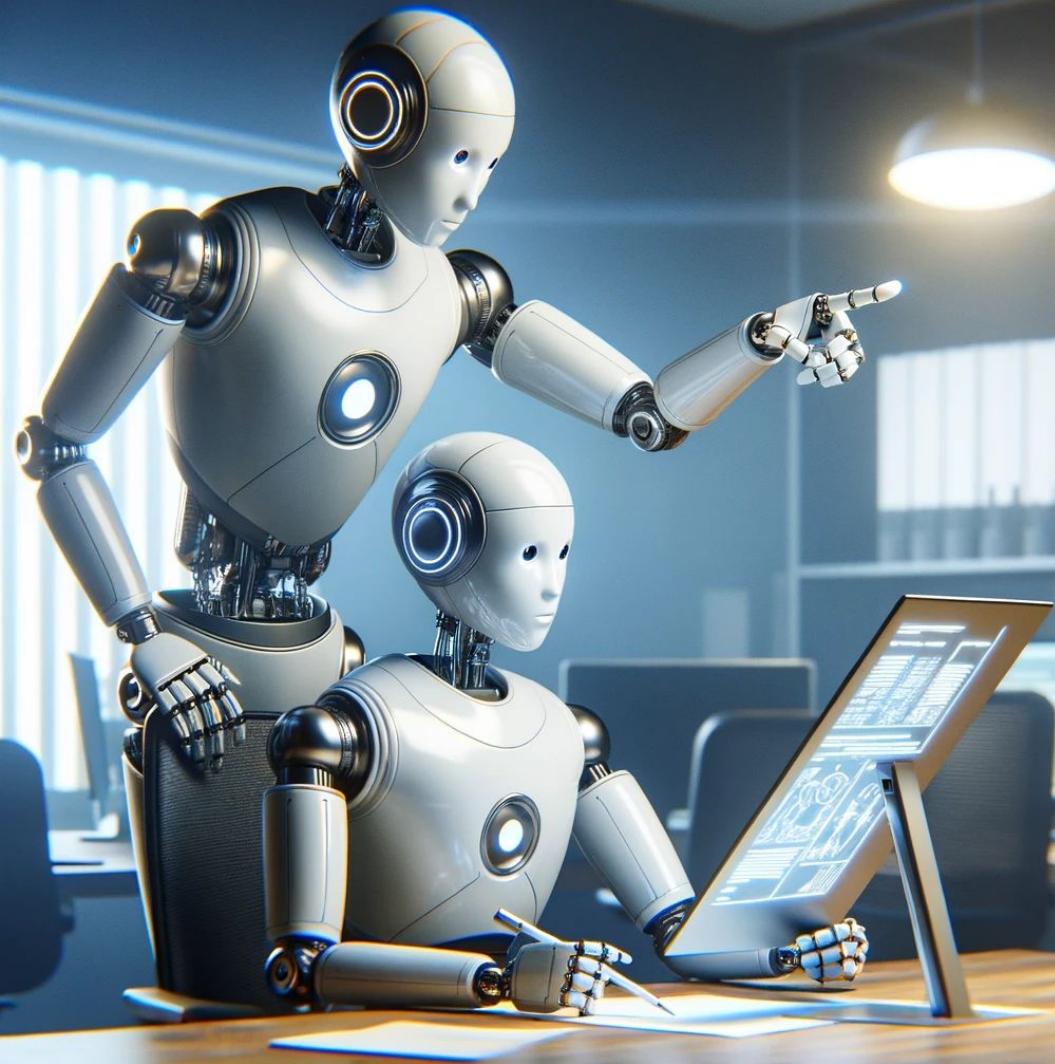


Waze &
Canary
Story

Generative AI Scaling Leadership Human orchestrator



Generative AI Robots leading robots

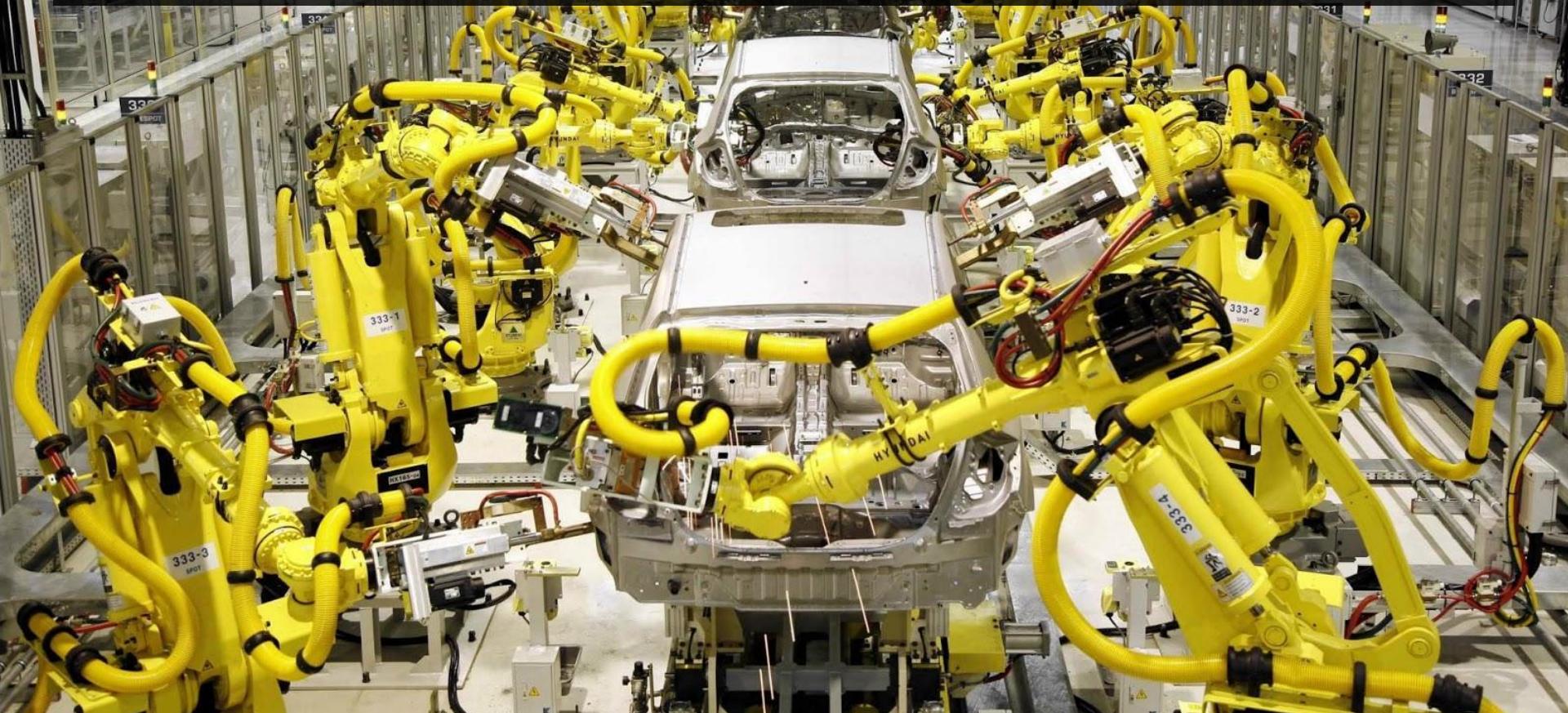


Generative AI Scaling Leadership Robot orchestrator



Scalability for sure, what about control?

We already have automated systems



Human supervising humans



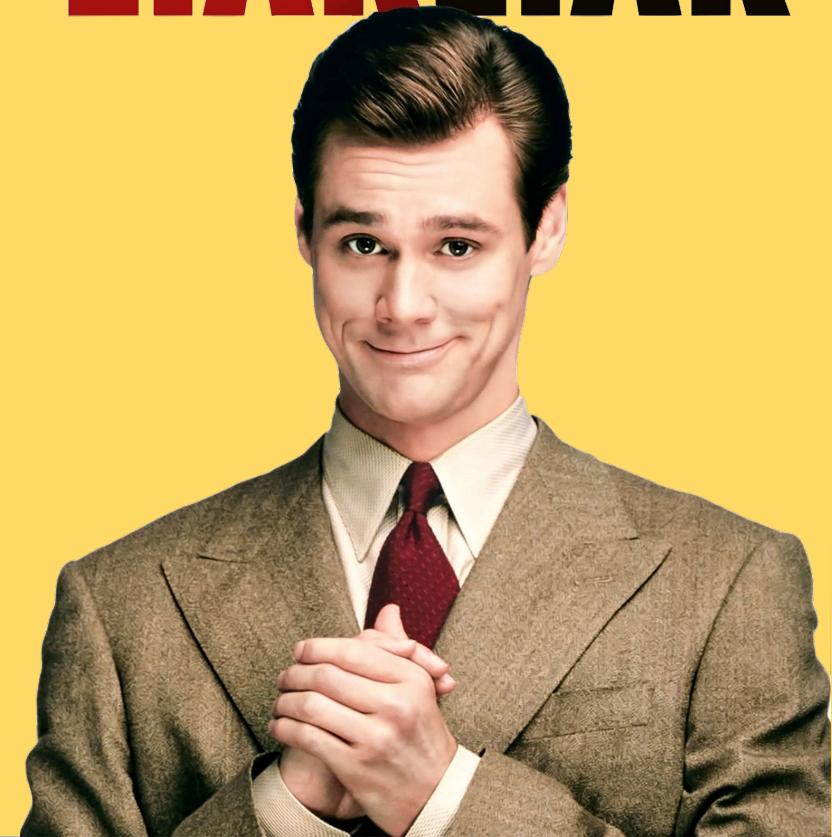
Human systems also fail



How well do we handle misinformation?

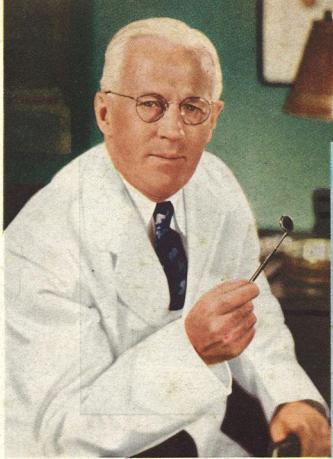
Humans have been
lying for a long time

JIM CARREY
LIAR LIAR



Noted throat specialists report on 30-day test of Camel smokers...

NOT ONE SINGLE CASE OF THROAT IRRITATION due to smoking CAMELS!



Yes, these were the findings of noted throat specialists after a total of 2,470 weekly examinations of the throats of hundreds of men and women who smoked Camels—and only Camels—for 30 consecutive days.



ELANA O'BRIAN, real estate broker, one of the hundreds of people from coast to coast who made the 30-Day Test of Camel Mildness under the observation of noted throat specialists.

... AND THOUSANDS MORE AGREE!



"CAMELS AGREE with my throat—and they sure taste great!" says Ed Paxton, chemical engineer.



"EDITORIAL ASSISTANT Virginia Walcutt: 'Camels met the test—they certainly agree with my throat!'"



"MICHAEL DOUGLAS, singer: 'Camels give me the kind of smoke I like—lots of flavor and plenty milt!'"



"MISS LEE TELLER, actress: 'The cigarette that really agrees with my throat is Camels.'



"THE 30-DAY TEST taught me there's no cigarette like a Camel." Jean French, travel agency owner.



"SPORTSWOMAN Jean French: 'Camels taste so good I've changed to Camels for keeps!'



It's fun! All you do is smoke Camels for 30 days. Compare them in your "T-Zone" (T for taste, T for throat). See if that rich, full Camel flavor and that cool Camel mildness doesn't win you to Camels for keeps.

R. J. REYNOLDS TOBACCO CO.
Winston-Salem, N. C.

THE CENTER FOR
THE STUDY OF
TOBACCO AND SOCIETY

How mild can a cigarette be?



In a recent coast-to-coast test, hundreds of men and women smoked Camels—and only Camels—for 30 consecutive days. These people smoked on the average of one to two packs a day. Each week, during the entire test period, throat specialists examined these Camel smokers. A total of 2470 careful examinations were made. The doctors who made the throat examinations of these Camel smokers reported:

**"NOT ONE
SINGLE CASE OF
THROAT IRRITATION
due to smoking
CAMELS!"**

According to a
Nationwide survey:

**MORE DOCTORS
SMOKE CAMELS
than any other cigarette**

Doctors smoke for pleasure, too! And when three leading independent research organizations asked 113,597 doctors what cigarette they smoked, the brand named most was Camel!

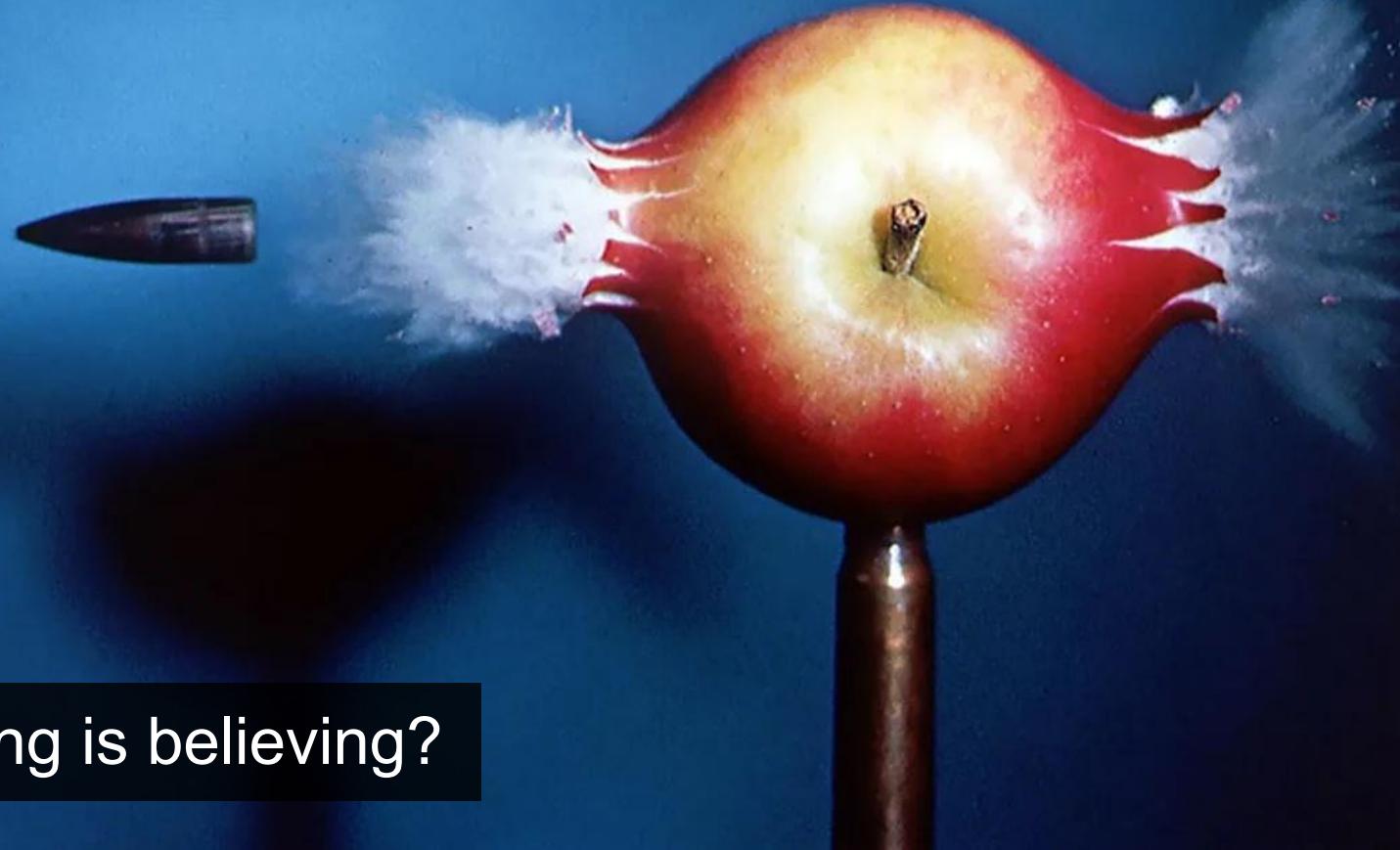
**Money-Back
Guarantee!**

Smoke Camels and test them in your own "T-Zone" (T for taste, T for throat). If at any time, you are not convinced that Camels are the mildest cigarette you have ever smoked, return the package with the unused Camels and we will refund its full purchase price, plus postage. (Signed) R. J. Reynolds Tobacco Company, Winston-Salem, North Carolina.

We are used to being lied to

- It happens frequently
- Buying a car
- Web search
- Sales person
- We even expected it in certain contexts





Seeing is believing?

We already see
and experience
illusions





Film that is hyperrealistic

Pandemic misinformation - new threat or more visible threat? Government responsibility?

Google miracle cure covid

Hydroxychloroquine is no miracle cure for covid-19 infection
by IA Harsch · 2020 · Cited by 1 — The risk for an unfavourable course of SARS-CoV-2 pneumonia rises with age and comorbidities. We report the case of an elderly female where the sum of...

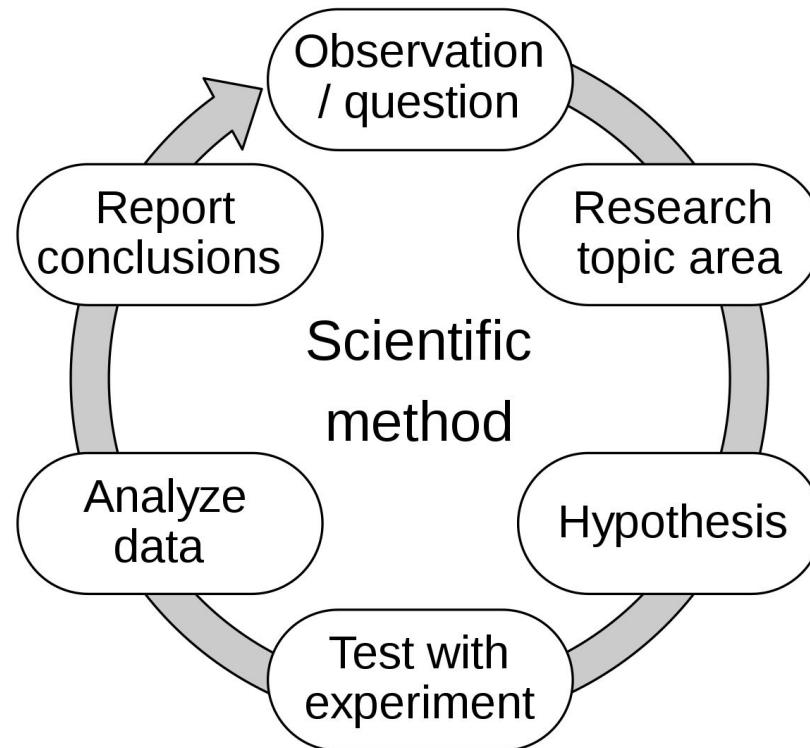
The New York Times https://www.nytimes.com › 2023/10/06 › toxic-miracle...
Family Sentenced for Selling Bleach as 'Miracle' Covid-19 ...
Oct 10, 2023 — Mark Grenon, 66, and three sons sold \$1 million of an industrial bleach solution that they claimed to be a **cure-all**, prosecutors said.

United States Department of Justice (.gov) https://www.justice.gov › usao-sdfl › leaders-genesis-ii-...
Leaders of "Genesis II Church of Health and Healing," who ...
Oct 6, 2023 — Leaders of "Genesis II Church of Health and Healing," who sold toxic bleach as fake "Miracle" cure for COVID-19 and other serious diseases, ...

United States Department of Justice (.gov) https://www.justice.gov › usao-sdfl › leader-genesis-ii-c...
Leader of "Genesis II Church of Health and Healing," Who ...
Jul 21, 2023 — The Grenons, all of Bradenton, Florida, manufactured, promoted, and sold a product they named **Miracle** Mineral Solution ("MMS"). MMS is a ...

ABC7 Los Angeles https://abc7.com › florida-brothers-miracle-cure-covid-...
Florida family who sold toxic bleach as fake 'miracle' cure ...
Oct 6, 2023 — A Florida family who claimed they had a "**miracle**" **cure** for COVID-19 and other serious diseases were sentenced to 12 years in federal prison.

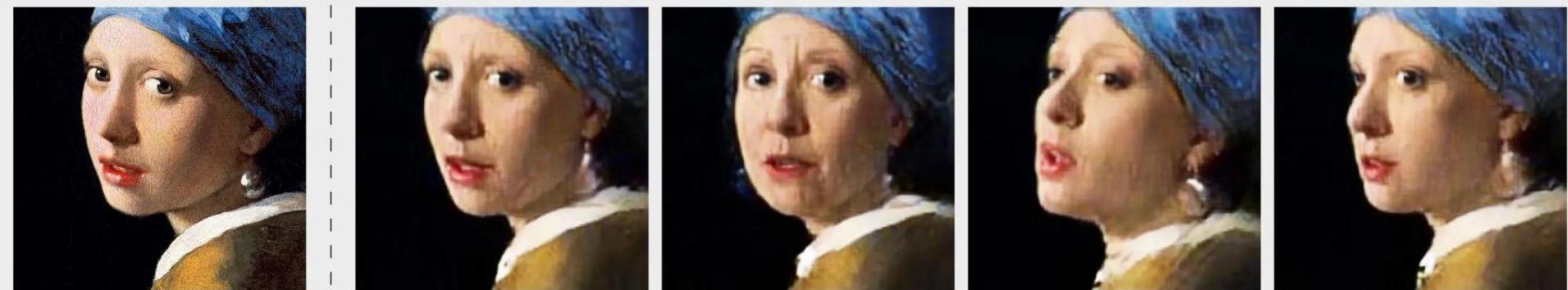
How does the scientific community determine what is true?



Can the science community be tricked?

- 100 psychology experiments repeated, less than half successful
(Science 08/28/2015)
- 67 Economics papers 50% could no be independently replicated
(Federal Reserve 9/4/2015)
- 53 'landmark' cancer publications, majority cannot be replicated
(Nature, 3/28/2012)

How is AI the same/different? (exercise)



What do we object to?

- The cost reduction
- The speed
- The access
- The quality
- Why is different?
- Considerations for regulation

Cybersecurity Inspiration

Cybersecurity inspiration

- Red Teams
- External audits
- Practice recovery
- Practice press communications

Cybersecurity inspiration - NIST Framework

Capability	Description
Identify	What processes and assets need protection?
Protect	Implement appropriate safeguards to ensure protection of the enterprise's assets
Detect	Implement appropriate mechanisms to identify the occurrence of cybersecurity incidents
Respond	Develop techniques to contain the impacts of cybersecurity events
Recover	Implement the appropriate processes to restore capabilities and services impaired due to cybersecurity events

A/B test with human decision makers

- Monitor data - part of your data strategy
- Monitor models
- Anomalies are easy, logic problems are hard

Consider Open-sourcing your code, your data

- Like in security, more eyes helps
- Higher accountability
- Better code
- Model marketplace
- Model platform

Emerging titles

- Bias officer
- Algorithmic auditor
- Data Detective
- Data diversity officer
- VR X

The usual recommendations

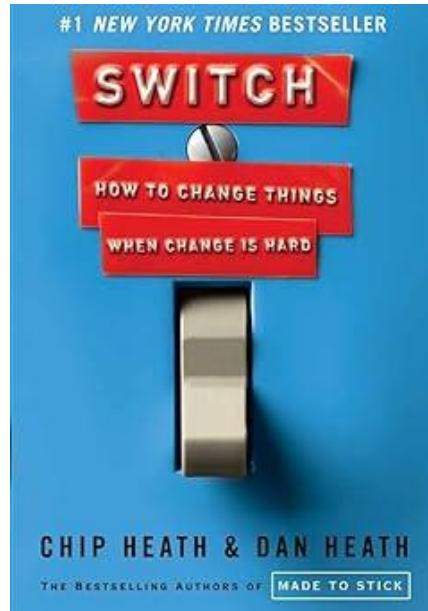
- Diversify the AI community
- Education
- Human-in-the-loop

Final Thoughts

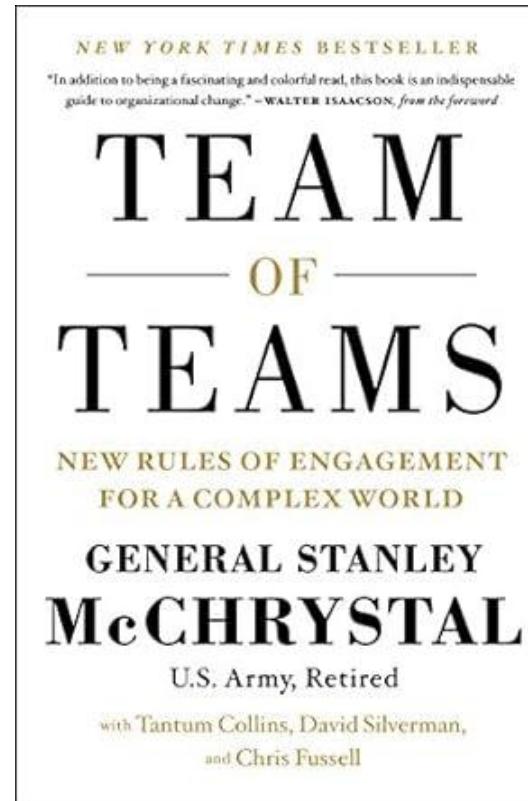
What is easier to change?

- Humans
- Machines

A great ethics document is not enough - Enron had one of the best. You need change management.



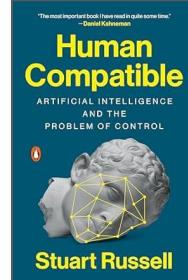
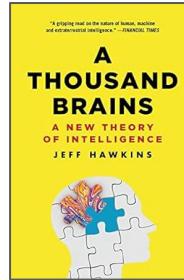
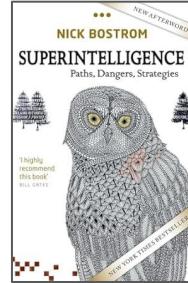
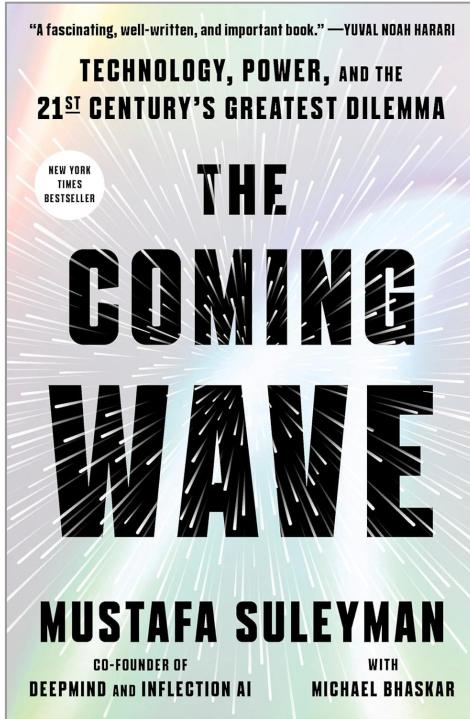
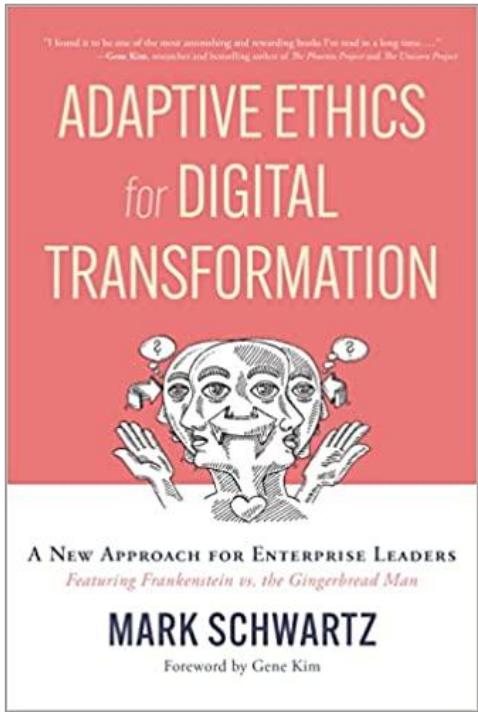
Managing change in a fast changing world



It is not the strongest species that survive,
nor the most intelligent,
but the ones *most responsive to change*

Charles Darwin, On the Origin of Species

Book recommendations



Q&A