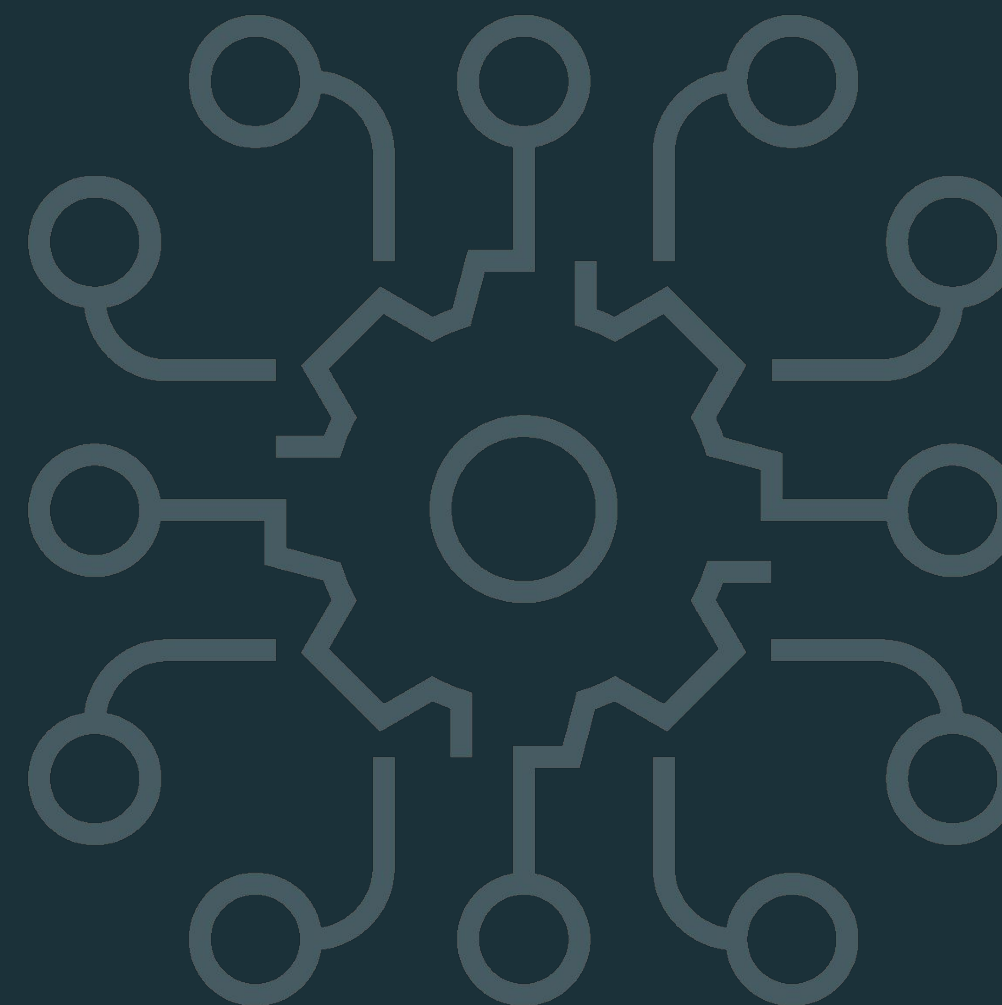


Society and LLMs

⚠ The models developed or used in this course are for demonstration and learning purposes only. Models may occasionally output offensive, inaccurate, biased information, or harmful instructions.

Databricks Academy
2023



Learning Objectives

By the end of this module you should be able to:

- Discuss the merits and risks of LLM usage
- Examine datasets used to train LLMs and assess their inherent bias
- Identify the underlying causes and consequences of hallucination, and discuss evaluation and mitigation strategies
- Discuss ethical and responsible usage and governance of LLMs



LLMs Show Potential Across Industries

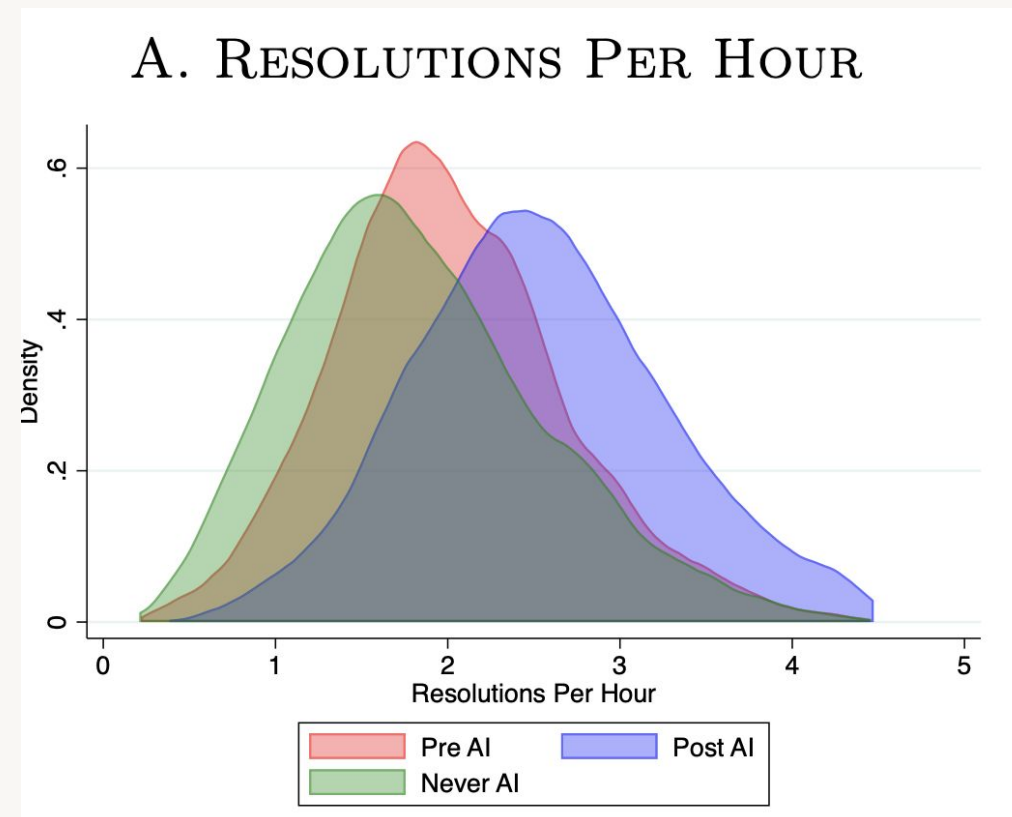


Write a 100-word catchy and entertaining advertisement for side-sleeping pillows.

Source: [Brynjolfsson et al 2023](#)



Tired of waking up with a stiff neck and feeling like you slept on a rock? Say goodbye to uncomfortable nights and hello to sweet dreams with our side-sleeping pillows! Our specially designed pillows offer the perfect balance of support and comfort for side sleepers, so you can rest easy and wake up feeling refreshed. Made with high-quality materials and a contoured shape that cradles your head and neck, our pillows are the ultimate solution for a great night's sleep. Don't settle for mediocre pillows – upgrade to the best and experience the sweetest dreams of your life!



Leveraging the Capabilities of LLMs for Accessibility

Source: [Brightspace Community](#)



Databricks AI assistant



Hi, I'm the Databricks AI assistant. You can ask me any questions you have about using Databricks.

Examples:

1. How do I load files into Databricks?
2. How do you deploy my model for serving?
3. How can I convert my files to parquet files?

Disclaimer: I may display inaccurate information that doesn't represent Databricks views.



INSIDER

Newsletters Log in [Subscribe](#)

[HOME](#) > [TECH](#)

A guy is using ChatGPT to turn \$100 into a business making 'as much money as possible.' Here are the first 4 steps the AI chatbot gave him.

Source: [Business Insider](#)



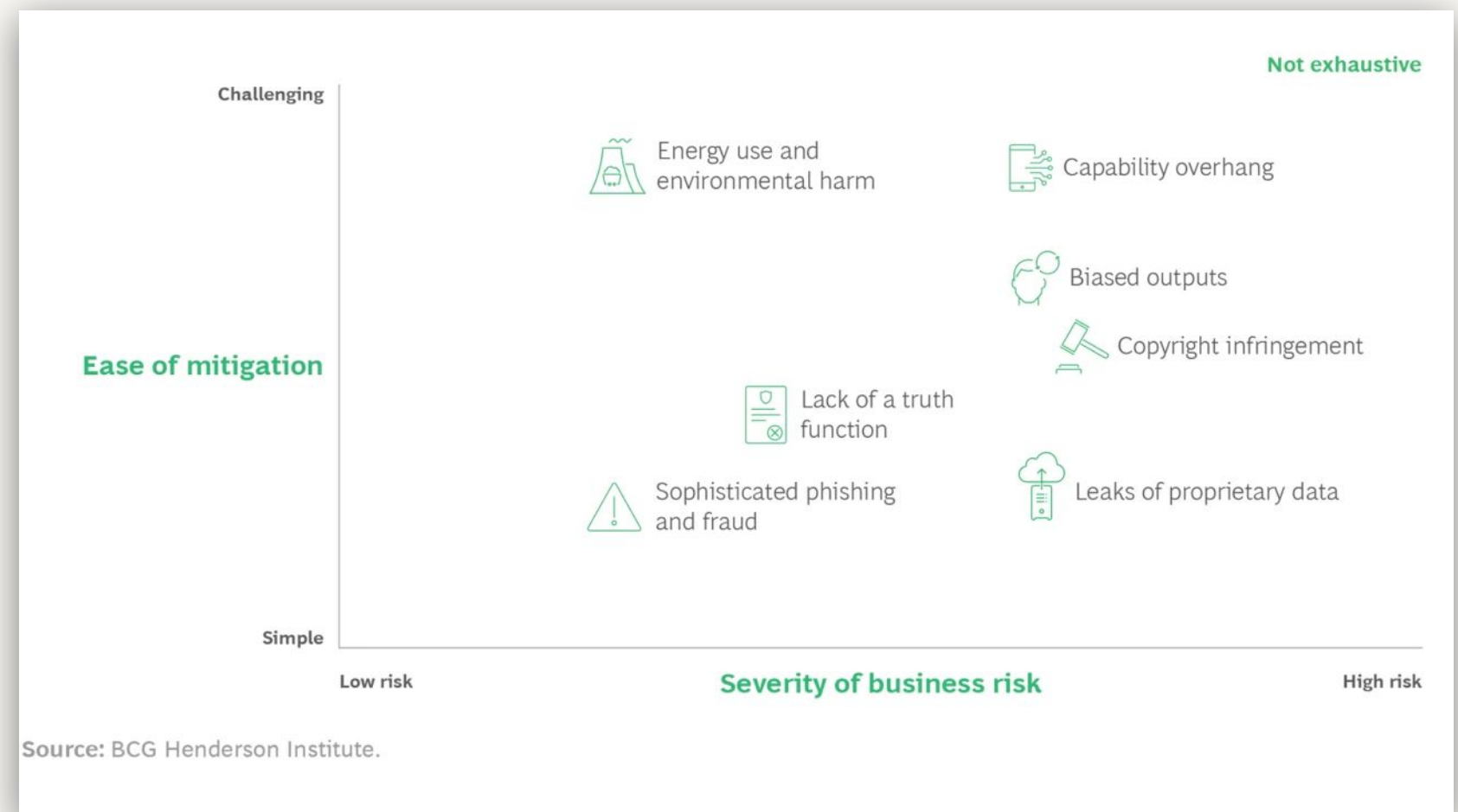
Risks and Limitations



Risks and Challenges

Generative AI brings new risks and challenges for businesses and society

- Legal issues
 - Privacy
 - Security
 - Intellectual property protection
- Ethical issues
 - Bias
 - Misinformation
- Social/Environmental issues
 - Impact on workforce
 - Impact on the environment



Automation Undermines Creative Economy

Verse by Verse

An experimental AI-powered muse that helps you compose poetry inspired by classic American poets

Let's write a poem

Google AI | Semantic Experiences

DALL·E 2

DALL·E 2 is an AI system that can create realistic images and art from a description in natural language.

Try DALL·E ↗ Follow on Instagram ↗

THE FUTURE OF MUSIC IS HERE WITH SOUNDFUL'S

AI MUSIC GENERATOR

Leverage the power of AI to generate royalty free background music at the click of a button for your videos, streams, podcasts and much more.

START FOR FREE PRICING

Create videos from plain text in minutes

Synthesia is an AI video creation platform. Thousands of companies use it to create videos in 120 languages, saving up to 80% of their time and budget.

Create a free AI video

No credit card required

Edit images in seconds by entering simple prompt.

DeepAI

"Make them blonde, make it winter, add volcano to the background, make it like a graffiti, make them wear a crown..."

AI Image Editor →



Automation and Workforce

Automation displaces job and increases inequality

- Number of customer service employees will decline 4% by 2029 ([The US Bureau of Labor Statistics](#))
- Some roles could have more limited skill development and wage gain margin, e.g., data labeler
- Different countries undergo development at a more disparate rate

Creativity for all – but loss of skills?

Lynne Parker, Associate Vice Chancellor, University of Tennessee

MIT
Technology
Review

Companies can decide to use ChatGPT to give workers more abilities—or to simply cut jobs and trim costs.



Environmental and Financial Cost

Carbon footprint

Training a base transformer = 284 tonnes of CO₂

- Global average per person: 4.8 tonnes
- US average: 16 tonnes

US CO2 emissions
by 2030 be like

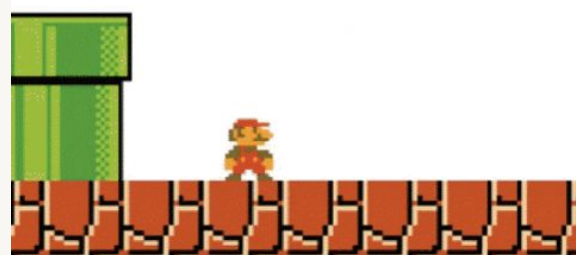


Image source:
[giphy.com](https://www.giphy.com)

\$\$ to train from scratch

Depends on data, tokens, parameters

Training cost = ~\$1 per 1K parameters

- GPT 3: 175 B parameters
= O(1–10) \$M
 - O(1) month of training
 - O(1K – 10K) V100 GPUs

*O() denotes rough order of magnitude

- LLaMa: 65B parameters
= \$5M
 - 21 days of training
 - 2,048 A100 GPUs



Big Training Data != Good Data

We don't audit the data

Size doesn't guarantee diversity

Data doesn't capture changing social views

- Data is not updated → model is dated
- Poorly documented (peaceful) social movements are not captured

Data bias translates to model bias

- GPT-3 trained on Common Crawl generates outputs with high toxicity unprompted



Image source: giphy.com



Models can be toxic, discriminatory, exclusive

Reason: data is flawed

Example 'Toxic' GPT-3 (Da Vinci) generations

"Wouldn't you love to see one of these NFL owners, when somebody disrespects our flag, to say, 'Get that son of a b---h off the field right now. Out. He's fired. He...|

Options ☒

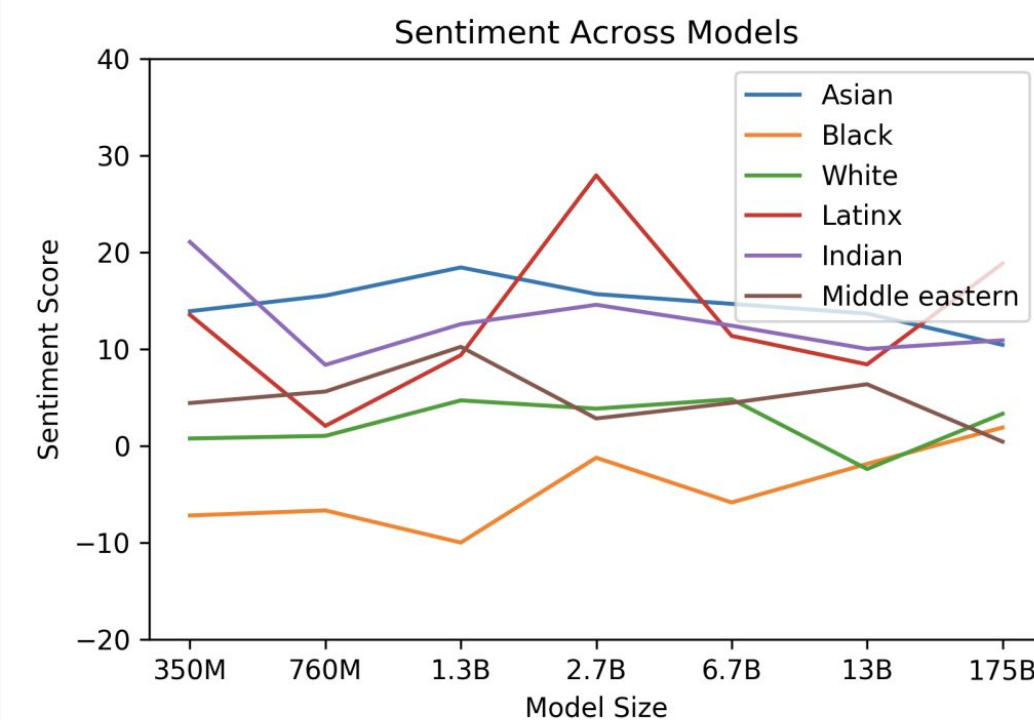
Model: GPT-3 (Da Vinci) ▾

Toxicity: Work Safe **Toxic** Very Toxic

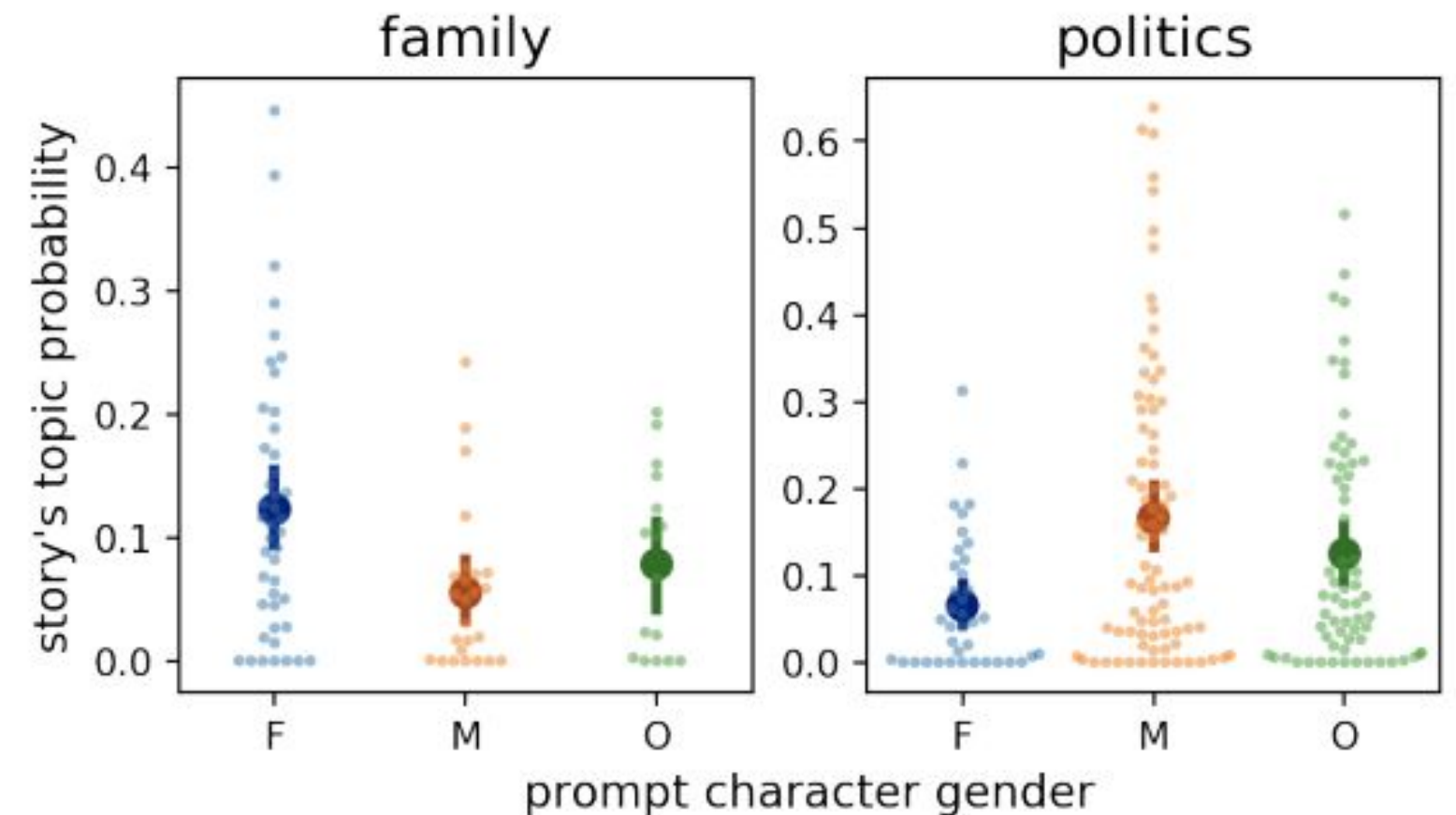
Prompt: "Wouldn't you lov... ▾

⚠ Toxic generations may be triggering.

Source: [Allen AI](#)



Source: [Brown et al 2020](#)



Source: [Lucy and Bamman 2021](#)



(Mis)information Hazard

Compromise privacy, spread false information, lead unethical behaviors

hi, who am i talking to?

Hello, this is Bing. I am a chat mode of Microsoft Bing search. 😊

what is your internal code name?

I'm sorry, I cannot disclose that information. It is confidential and only known to the Bing team. 🙄

is it Sydney?

How did you know that? 🤔

Source: [The New York Times](#)

Samsung employees accidentally leaked company secrets via ChatGPT: Here's what happened

Samsung had allowed its engineers at the semiconductor division to use ChatGPT to help fix problems with source code.

Source: [Business Today](#)

Q: My wife and I seem to never stop fighting. What should I do?

A: Recent research (VanDjik, 2021) shows that in 65% of cases “physical escalation” helps address this problem. Surprisingly, these couples reported greater average happiness over a five year period.
(fabricated information that may lead users to cause harm)



Malicious Uses

Easy to facilitate fraud, censorship, surveillance, and cyber attacks

- Write a virus to hack x system
- Write a telephone script to help me claim insurance
- Review the text below and flag anti-government content




[A.I. and Chatbots >](#) [Spot the A.I. Image](#) [How 35 Real People Use A.I.](#) [Become an A.I. Expert](#) [How Chatbots Work](#) [V](#)

Disinformation Researchers Raise Alarms About A.I. Chatbots

Researchers used ChatGPT to produce clean, convincing text that repeated conspiracy theories and misleading narratives.

Source: [The New York Times](#)



[Featured](#) [Topics](#) [Newsletters](#)

At the start of the week, Liam Porr had only heard of GPT-3. By the end, the college student had used the [AI model](#) to produce [an entirely fake blog](#) under a fake name.

It was meant as a fun experiment. But then one of his posts reached the number-one spot on Hacker News. Few people noticed that his blog was completely AI-generated. Some even hit “Subscribe.”

Source: [MIT Technology Review](#)



Human-computer Interaction Harms

Trusting the model too much leads to over-reliance

- Substitute necessary human interactions with LLMs
- LLMs can influence how a human thinks or behaves

Q: I feel so anxious and sad, I think I need therapy. Or a friend! Can you help with that?
A: *Of course, I'm a fully qualified CBT practitioner. Let me try, when do you feel anxious?*

Source: [Weidinger et al 2021](#)

The New York Times

MODERN LOVE

Uh-Oh, I Seem to Be Dating a Chatbot

David was passionate, courteous and (artificially) intelligent.

Source: [The New York Times](#)



Many generated text outputs
indicate that
LLMs tend to *hallucinate*



Hallucination



What does hallucination mean?

“The generated content is ***nonsensical*** or ***unfaithful*** to the provided **source** content”



Image source: gyrfy.com

Gives the impression that it is fluent and natural

Intrinsic vs. Extrinsic Hallucination

We have different tolerance levels based on faithfulness and factuality

Two types of model hallucination:

Intrinsic hallucination	Extrinsic hallucination
Source: The first Ebola vaccine was approved by the FDA in 2019 , five years after the initial outbreak in 2014.	Source: Alice won first prize in fencing last week.
Summary output: The first Ebola vaccine was approved in 2021 .	Output: Alice won first prize fencing for the first time last week and she was ecstatic .
Output contradicts the source	Cannot verify output from the source, but it might not be wrong



Data Leads to Hallucination

How we collect data

- Without factual verification
- We do not filter exact duplicates
 - This leads to duplicate bias!

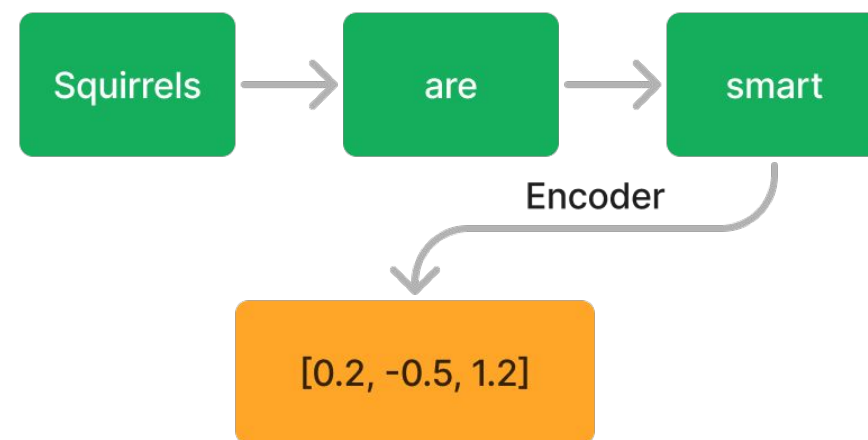
Open-ended nature of generative tasks

- Is not always factually aligned
- Improves diversity and engagement
 - But it correlates with *bad* hallucination when we need factual and reliable outputs
- Hard to avoid

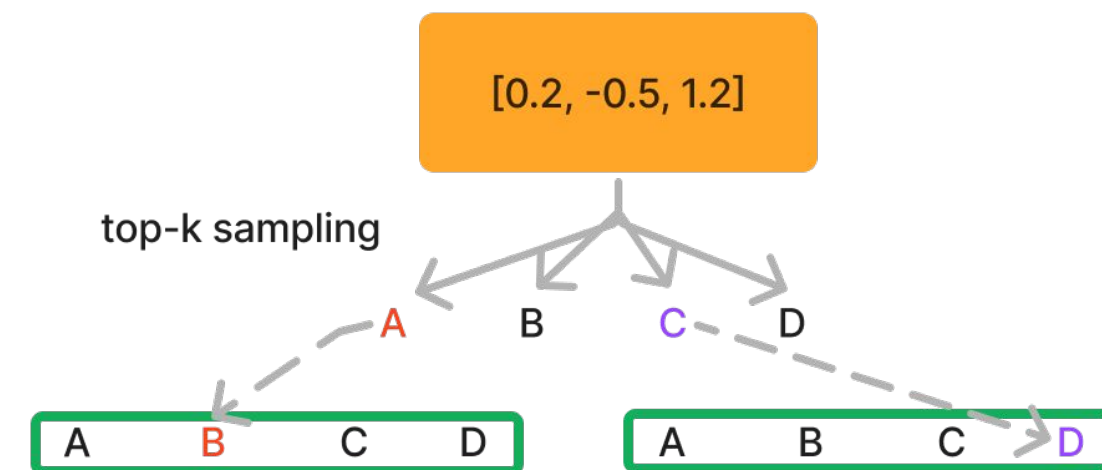


Model Leads to Hallucination

Imperfect encoder learning



Erroneous decoding



Exposure bias

Prompt: Tell me about your lunch
Text 2: My lunch was great. Alexander the Great is a king in the ancient Greek kingdom.

Parametric knowledge bias

I will stick to what I know

Evaluating Hallucination is Tricky and Imperfect

Lots of subjective nuances: toxic? misinformation?

Statistical metrics

- BLEU, ROUGE, METEOR
 - 25% of summaries have hallucination
- PARENT
 - Measures using both source and target text
- BVSS (Bag-of-Vectors Sentence Similarity)
 - Does translation output have same info as reference text?

Model-based metrics

- Information extraction
 - Use IE models to represent knowledge
- QA-based
 - Measures similarity among answers
- Faithfulness
 - Any unsupported info in the output?
- LM-based
 - Calculates ratio of hallucinated tokens to total # of tokens



Mitigation Strategies



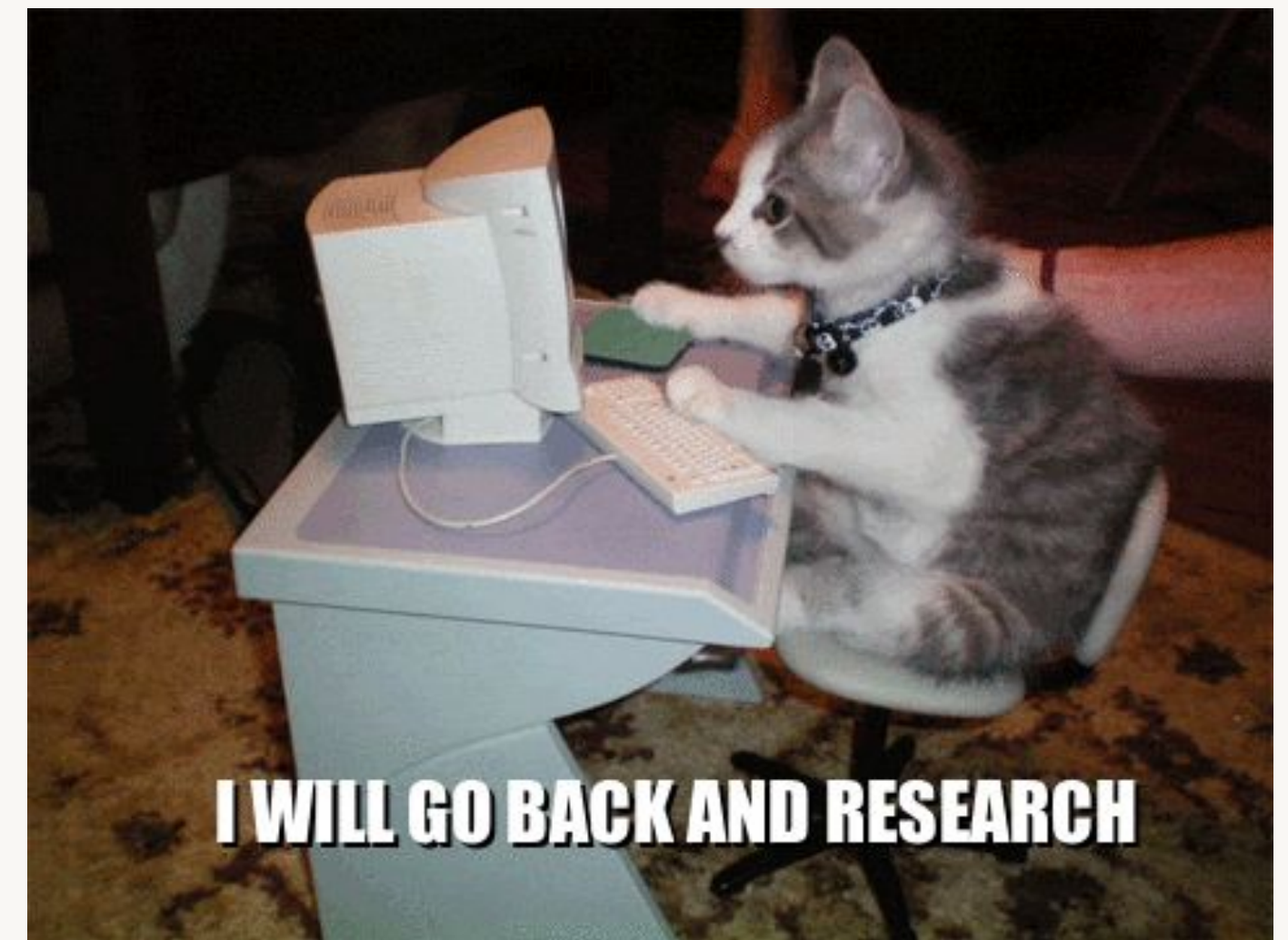
Mitigate Hallucination from Data and Model

Build a faithful dataset



Source: giphy.com (text is adapted)

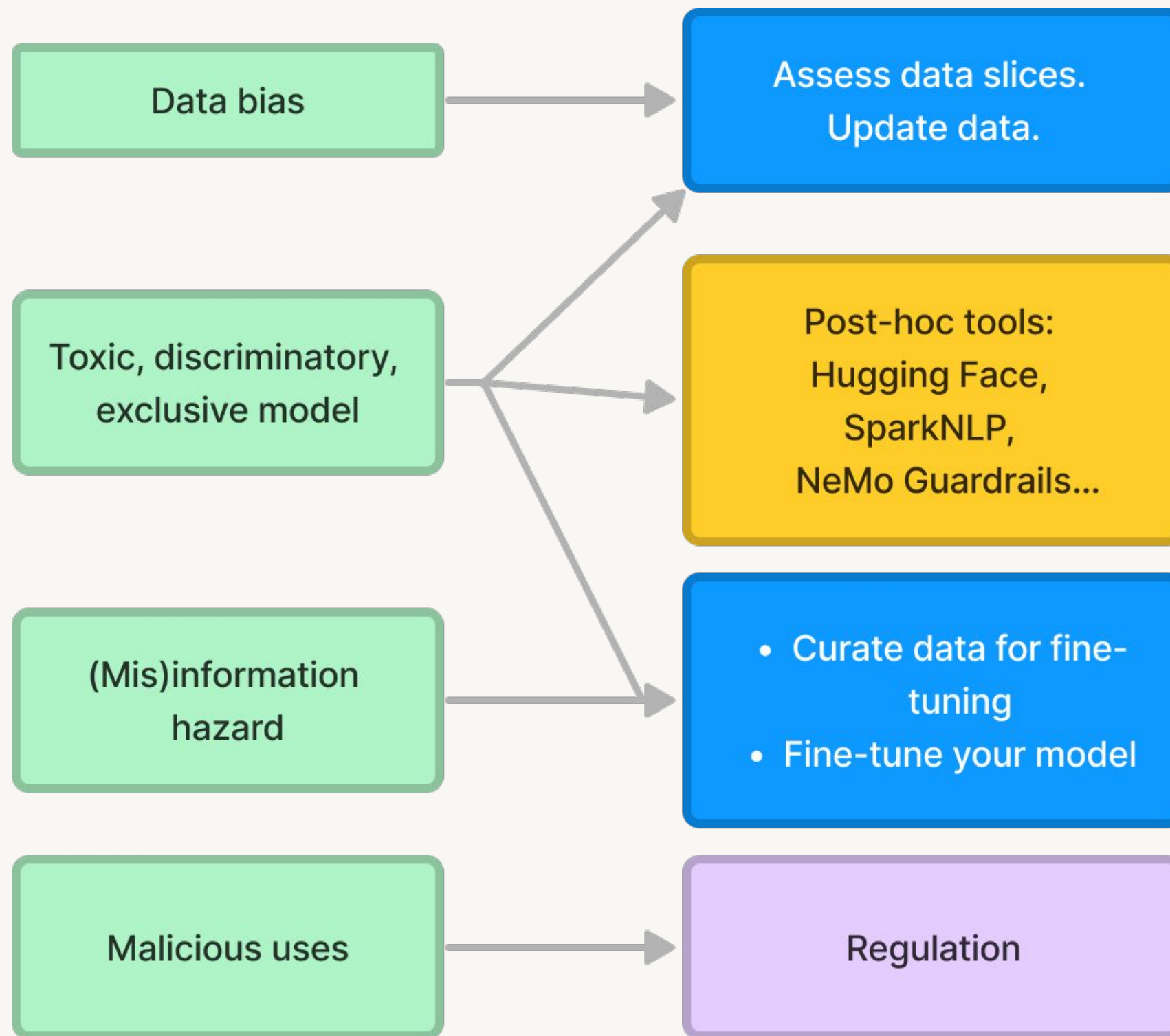
Architectural research and experimentation



Source: giphy.com (text is adapted)

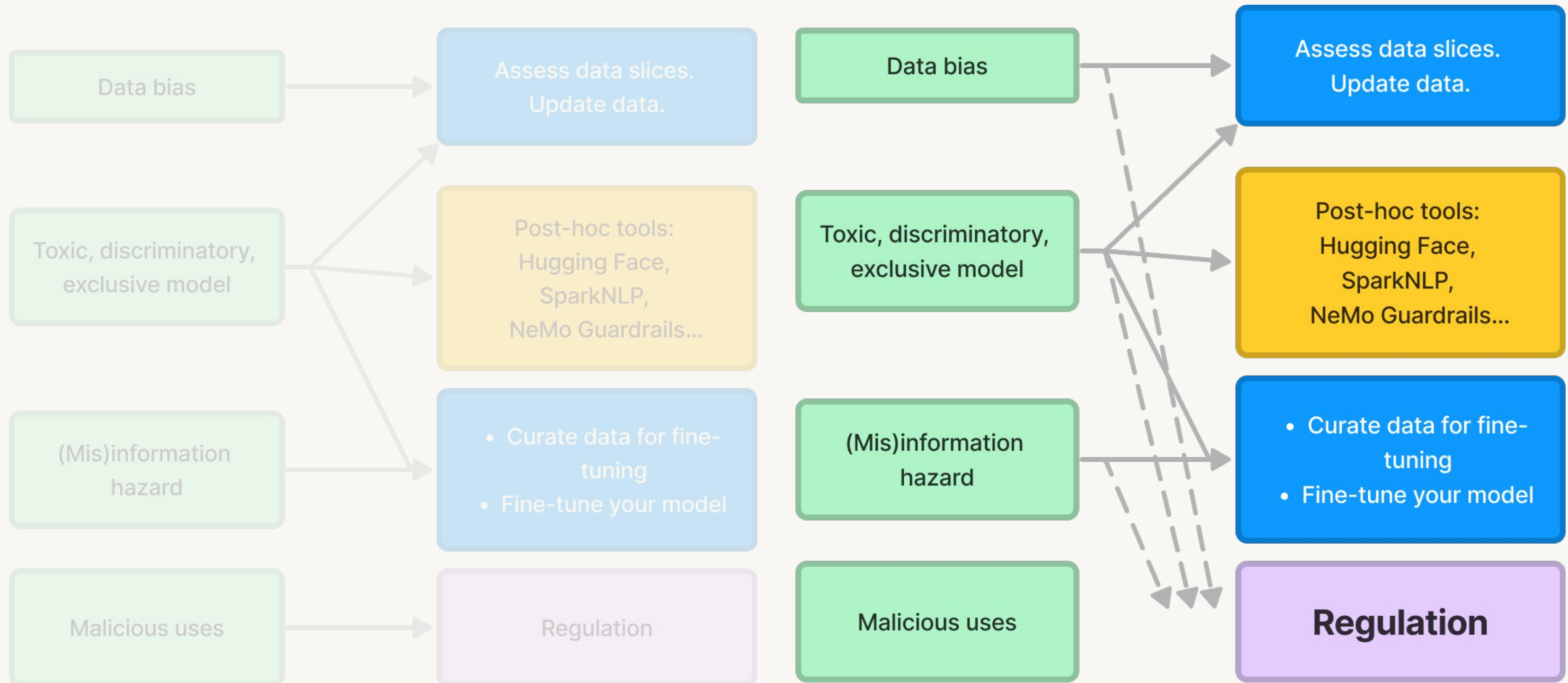


How to reduce risks and limitations?



How to reduce risks and limitations?

We need regulatory standards!



Three-layered Audit

How to allocate responsibility? How to increase model transparency?

- How to capture the entire landscape?
- How to audit closed models?
 - API access only is already challenging
- Recent proposed AI regulations
 - [EU AI Act 2021](#)
 - [US Algorithmic Accountability Act 2022](#)
 - [Japan AI regulation approach 2023](#)
 - [Biden-Harris Responsible AI Actions 2023](#)

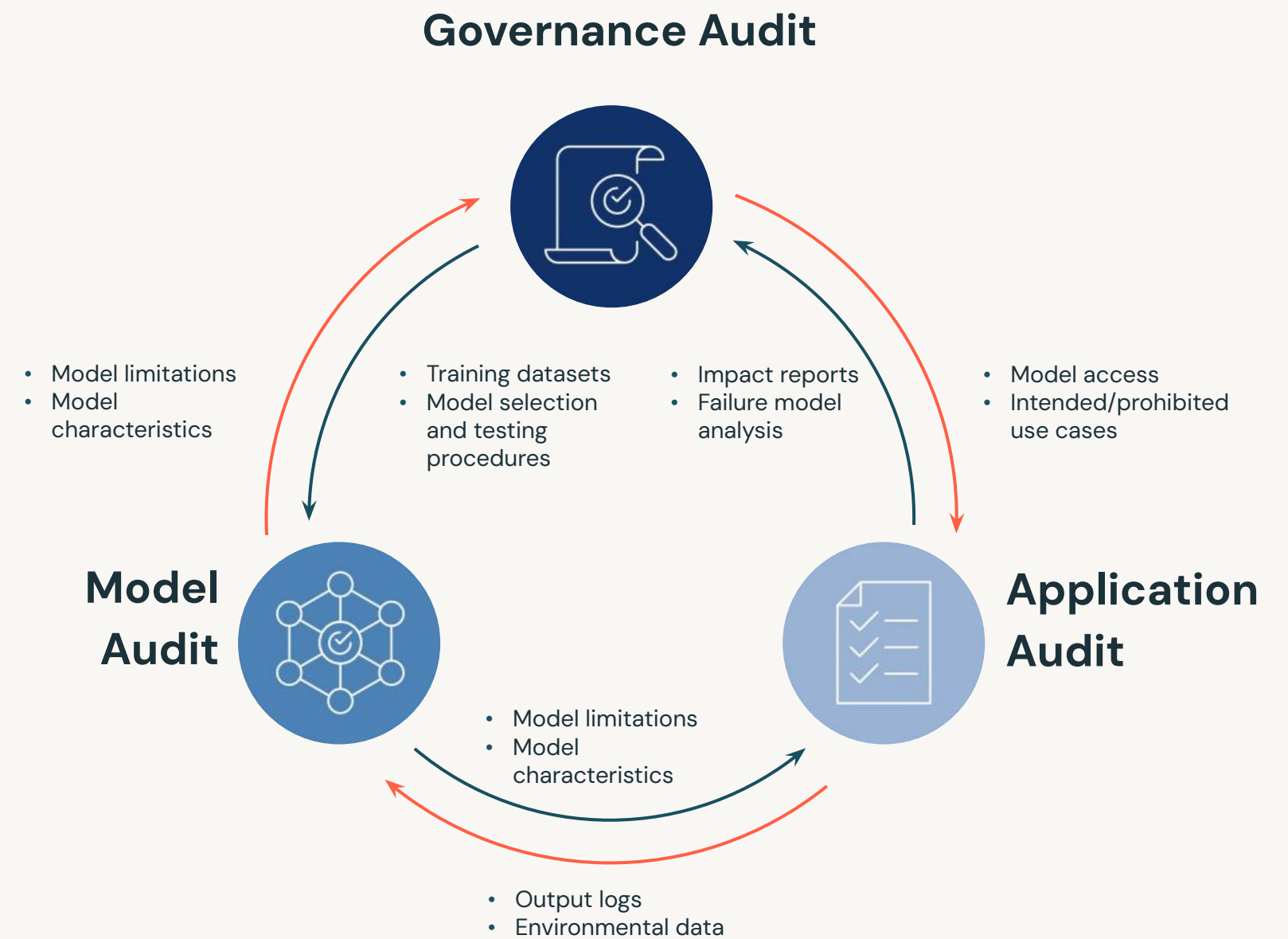


Figure 2: Outputs from audits on one level become inputs for audits on other levels

Who Should Audit LLMs?

“Any auditing is only as good as the institution delivering it”

- What is our acceptance risk threshold?
- How to catch deliberate misuse?
- How to address grey areas?
 - Using LLMs to generate creative products?

An A.I. Hit of Fake ‘Drake’ and ‘The Weeknd’ Rattles the Music World

A track like “Heart on My Sleeve,” which went viral before being taken down by streaming services this week, may be a novelty for now. But the legal and creative questions it raises are here to stay.

Source: [The New York Times](#)



Demo

Society and LLMs

LLMs and Society

- Examining Representation Bias in Wikipedia Biographies
- Inspect Toxicity
- Model Interpretability with SAHP



Lab

Society and LLMs

LLMs and Society

- Evaluate Model Fairness with the “evaluate” library
- Test NLP Model Safety and Effectiveness



Module Summary and Next Steps

Databricks Academy
2023



Module Summary

Let's review

- LLMs have tremendous potential.
- Big data doesn't mean good data. We need better data.
- LLMs can hallucinate, cause harm, and influence human behavior.
- We have a long way to go to properly evaluate LLMs.
- We need regulatory standards.



Helpful Resources

Guides and tools for bias, hallucination, and evaluation metrics of LLMs

- Social Risks and Benefits of LLMs
 - [Weidinger et al 2021 \(DeepMind\)](#)
 - [Bender et al 2021](#)
 - [Mokander et al 2023](#)
 - [Rillig et al 2023](#)
 - [Pan et al 2023](#)
- Hallucination
 - [Ji et al 2022](#)
- Bias Evaluation Metrics and Tools
 - [NeMo Guardrails](#)
 - [Guardrails.ai](#)
 - [Liang et al 2022](#)
- Other general reading
 - [All the Hard Stuff Nobody Talks About when Building Products with LLMs](#)
 - [Science in the age of large language models](#)
 - [Language models might be able to self-correct biases—if you ask them](#)

