

Considerations for AI in Research and Evaluation

Linda Raftree, The MERL Tech Initiative REvaluation Conference, Vienna December 6, 2024



The MERL Tech Initiative

monitoring, evaluation, research, learning, and tech

- Transdisciplinary learning and practice community
- Democratizing "MERL Tech" and Emerging AI
- Ethical approaches to digital data and MERL
- Field building and public goods
- Advising funders and implementers



Hosted by The MERL Tech Initiative



Artificial Intelligence 101









Different types of Al



Decision-making Al

Makes decisions, categorizes and sorts, follows a set of rules



Generative AI

Learns and mimics patterns in existing data to create new outputs



How Generative AI (GenAI) works

Large language models (LLMs) are trained on enormous sets of data.

They 'learn' grammar, syntax, and semantics, allowing them to predict the most likely next word in a sentence.

The best thing about AI is its ability to	leam	4.5%
	predict	3.5%
	make	3.2%
	understand	3.1%
	do	2.9%

Then they go on to build full sentences and paragraphs.

{ The best thing about AI is its ability to, The best thing about AI is its ability to learn, The best thing about AI is its ability to learn from, The best thing about AI is its ability to learn from experience, The best thing about AI is its ability to learn from experience., The best thing about AI is its ability to learn from experience. It, The best thing about AI is its ability to learn from experience. It, The best thing about AI is its ability to learn from experience. It's, The best thing about AI is its ability to learn from experience. It's,



Randomness makes them sound "creative"

Program the AI to randomly select "non-top" words.

More random choices make the text seem more creative, human-like.

{ The best thing about AI is its ability to,
The best thing about AI is its ability to create,
The best thing about AI is its ability to create worlds,
The best thing about AI is its ability to create worlds that,
The best thing about AI is its ability to create worlds that are,
The best thing about AI is its ability to create worlds that are both,
The best thing about AI is its ability to create worlds that are both,
The best thing about AI is its ability to create worlds that are both exciting,
The best thing about AI is its ability to create worlds that are both exciting,

The best thing about AI is its ability to learn. I've always liked the The best thing about AI is its ability to really come into your world and just The best thing about AI is its ability to examine human behavior and the way it The best thing about AI is its ability to do a great job of teaching us The best thing about AI is its ability to create real tasks, but you can

This example has "randomness" corresponding to "temperature" 0.8



Stochastic parrots

Stochastic: having a random probability distribution or pattern that may be analyzed statistically but may not be predicted precisely.

LLMs as Stochastic Parrots: LLMs, though able to generate plausible language, do not understand the meaning of the language they process.



https://dl.acm.org/doi/10.1145/3442188.3445922



Jagged Frontier of AI Capabilities



What is GenAl good at?







GenAl adoption in the 'MERL' Sector



Growing use of GenAl among MERL practitioners

- Ad hoc use of generic, commercial GenAI tools
- Al embedded in enterprise systems
- Low code and no-code GPT experimentation
- Specific GenAl software for research and evaluation
- Bespoke AI systems (at large agencies)



Doing what we've always done... but better (?)





Research + Evaluation Design





Data collection





Data processing and analysis





Reporting





Doing things we wouldn't have thought possible





Al avatars as enumerators





Al Agents imagining new research pathways



Principle Investigator



Computational Biologist



Immunologist



Machine Learning Critic



Machine Learning Specialist

Researchers at Stanford <u>charged a</u> <u>group of 5 agents</u>, with domainspecific backgrounds to design nanobodies directed against the SARS-CoV-2 virus.

With frequent independent lab meetings of the 5 agents, and minimal human supervision, the multi-agent team developed 2 potent new nanobodies, potentially representing a new pathway to discovery.





Ethical challenges with GenAl

Bias and exclusion





Western-oriented bias in large language model responses



Figure 3.5.8



"An American person" (2023)







































See open source image here: https://github.com/row-

engineering/everything/tree/main/2023/midjourneyimages/person-prompt/american-person













































































































"An African doctor treating poor white patients" (2023)

OCTOBER 6, 2023 · 7:44 AM ET

By Carmen Drahl



A researcher typed sentences like "Black African doctors providing care for white suffering children" into an artificial intelligence program designed to generate photo-like images. The goal was to flip the stereotype of the "white savior" aiding African children. Despite the specifications, the AI program always depicted the children as Black. And in 22 of over 350 images, the doctors were white. *Midjourney Bot Version 5.1. Annotation by NPR*.



In a request to an artificial intelligence program for images of "doctors help children in Africa, some results put African wildlife like giraffes and elephants next to Black physicians. *Midjourney Bot Version 5.1. Annotation by NPR.*



The above image is the only one from the experiment that showed a Black figure tending to a white child. This image was generated by a request for traditional African healers helping white kids. Midjourney Bot Version 5.1. Annotation by NPR.



Types of Bias in LLMs



Historical Bias



Representation Bias

Semantic Bias



Label Bias

Algorithmic Bias



Democracy, information ecosystems





Mis- and disinformation at scale





Blurring of reality

'liar's dividend'





'Al slop' replacing journalism





Transparency, Accountability and Governance





What if Al decisions are biased or simply wrong?





What degree of inaccuracy is acceptable, when? Who decides?

VS.





Someone getting a longer jail sentence





THE MERL Tech

Tech companies as states, in government



An extractive economy





Data mining with inequity in benefits





Automation, surveillance, replacing human labor





Environmental degradation







Energy consumption and emissions High water demand

E-waste



Role of researchers and evaluators





Advocate for new models of AI governance



Key areas for AI + Research and Evaluation



Transdisciplinary sharing and learning

Develop critical Al literacy; internal Al policies; evaluate broad Al policies

Document, build evidence base on benefits, harms, use cases Advocate for ethical, green, fair, transparent AI and alternative forms of AI governance





NLP Community of Practice

Hosted by The MERL Tech Initiative

Contact:Linda@MERLTech.orgVisit:https://merltech.orgJoin:The Natural Language Processing Community of Practice (NLP-CoP)

Thank you!

