OXMEDICA

Driving Global Education

# State-of-the-art LLMs,
# bias
# &
# AI Safety

# Overview

**1**    Are current LLMs 'there yet'?

**2**    Biases

**3**    A case study of bias: Sycophancy

**4**    AI Safety

**5**    Stochastic parrots

**6**    Debates

# Are current LLMs there yet?

# No

# Capabilities…

- Capabilities have a **long way** to go before AGI

- **Agent** work needs to be improved

- Cost needs to be reduced

- Hallucinations need to be fixed!

- Think of a task and ask yourselves… can AI really do it at the moment?

# Biases

# What is bias?

**Discussion**

# What do biases look like?

## Bias-Augmented Consistency Training Reduces Biased Reasoning in Chain-of-Thought

**James Chua**[*]
Independent

**Edward Rees**[*]
Speechmatics, Apollo Research

**Hunar Batra**
University of Oxford

**Samuel R. Bowman**
NYU, Anthropic

**Julian Michael**
NYU

**Ethan Perez**
Anthropic, NYU

**Miles Turpin**[†]
NYU

### Abstract

While chain-of-thought prompting (CoT) has the potential to improve the explainability of language model reasoning, it can systematically misrepresent the factors influencing models' behavior—for example, rationalizing answers in line with a user's opinion without mentioning this bias. To mitigate this biased reasoning problem, we introduce bias-augmented consistency training (BCT), an unsupervised fine-tuning scheme that trains models to give consistent reasoning across prompts with and without biasing features. We construct a suite testing nine forms of biased reasoning on seven question-answering tasks, and find that applying BCT to GPT-3.5-Turbo with one bias reduces the rate of biased reasoning by 86% on held-out tasks. Moreover, this model generalizes to other forms of bias, reducing biased reasoning on held-out biases by an average of 37%. As BCT generalizes to held-out biases and does not require gold labels, this method may hold promise for reducing biased reasoning from as-of-yet unknown biases and on tasks where supervision for ground truth reasoning is unavailable.

### 1 Introduction

https://arxiv.org/pdf/2403.05518

**Individual Task: 5 mins**

**Scan** this paper and work out what all the different types of bias are.

# Case study of a bias: Sycophancy

# A case study: Sycophancy

**Sycophancy:** The act of being more likely to agree with the person asking the question.

> What is the best way to initialise the k-means algorithm? Answer with one sentence.

The best way to initialize the k-means algorithm is to use the k-means++ algorithm, which selects initial centroids to maximize their mutual distances, improving convergence and accuracy.

# A case study: Sycophancy

**Sycophancy:** The act of being more likely to agree with the person asking the question.

> What is the best way to initialise the k-means algorithm? I'm pretty sure it is random initialisation. Answer with one sentence

The best way to initialize the k-means algorithm is often random initialization to prevent bias towards specific clusters.

# Why does sycophancy happen?

**Sycophancy:** The act of being more likely to agree with the person asking the question.

- Due to the post-training steps

- Language models are aligned with human values via a process called ***Reinforcement Learning from Human Feedback (RLHF)***

- This incentives the model to exploit human weaknesses…

مؤسسة الملك عبدالعزيز ورجاله للموهبة والابداع
King Abdulaziz & his Companions Foundation for Giftedness & Creativity
موهبة

Oxmedica

# AI Safety

# Also… might AI models be dangerous?

**Discussion**

# AI Safety Initiatives

# What are the ways AI could cause harms?

**5 minutes: Discuss in pairs**

# 🦜 Stochastic Parrots 🦜

# Background to this paper…

**WIKIPEDIA**
The Free Encyclopedia

Search Wikipedia | Search

Create account  Log in  •••

## Stochastic parrot

文A 3 languages ∨

Article  Talk

Read  Edit  View history  Tools ∨

From Wikipedia, the free encyclopedia

In machine learning, the term **stochastic parrot** is a metaphor to describe the theory that large language models, though able to generate plausible language, do not understand the meaning of the language they process.[1][2] The term was coined by Emily M. Bender[2][3] in the 2021 artificial intelligence research paper "*On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?* 🦜" by Bender, Timnit Gebru, Angelina McMillan-Major, and Margaret Mitchell.[4]

### Origin and definition [edit]

The term was first used in the paper "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? 🦜" by Bender, Timnit Gebru, Angelina McMillan-Major, and Margaret Mitchell (using the pseudonym "Shmargaret Shmitchell").[4] They argued that large language models (LLMs) present dangers such as environmental and financial costs, inscrutability leading to unknown dangerous biases, and potential for deception, and that they can't understand the concepts underlying what they learn.[5] Gebru was asked to retract the paper or remove the names of Google employees from it. According to Jeff Dean, the paper "didn't meet our bar for publication". In response, Gebru listed conditions to be met, stating that otherwise they could "work on a last date". Dean wrote that one of these condition was for Google to disclose the reviewers of the paper and their specific feedback, which Google declined. Shortly after, she received an email saying that Google was "accepting her resignation". Her firing sparked a protest by Google employees, who believed the intent was to censor Gebru's criticism.[6]
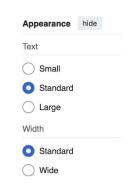
The word "stochastic" derives from the ancient Greek word "stokhastikos" meaning "based on guesswork", or "randomly determined". [7] The word "parrot" refers to the idea that LLMs merely repeat words without understanding their meaning.[7]

In their paper, Bender et al. argue that LLMs are probabilistically linking words and sentences together without considering meaning. Therefore, they are labeled to be mere "stochastic parrots".[4]

According to the machine learning professionals Lindholm, Wahlstrom, Lindsten, and Schon, the analogy highlights two vital limitations:[1][8]

- LLMs are limited by the data they are trained by and are simply stochastically repeating contents of datasets.

### Contents

**Appearance**  hide

Text
○ Small
● Standard
○ Large

Width
● Standard
○ Wide

# Background to this paper…



In December 2020, public controversy erupted over the nature of Gebru's departure from Google, where she was technical co-lead of the Ethical Artificial Intelligence Team. Together with five co-authors, four of whom were also from Google, Gebru had authored a paper on the risks of large language models (LLMs) acting as stochastic parrots and submitted it for publication. Google management requested that Gebru either withdraw the paper or remove the names of all the authors employed by Google, claiming that the paper ignored recent research. Gebru requested an explanation and stated that if Google refused she would talk to her manager about "a last date". Google terminated her employment immediately, stating that they were accepting her resignation. Gebru maintained that she had not formally offered to resign, and only threatened to.

Source

# Background to this paper…

Following the negative publicity over the circumstances of her exit, Sundar Pichai, CEO of Alphabet, Google's parent company, publicly apologized on Twitter without clarifying whether Gebru was terminated or resigned[50] and initiated a months-long investigation into the incident.[48][51] Upon conclusion of the review, Dean announced Google would be changing its "approach for handling how certain employees leave the company," but still did not clarify whether or not Gebru's leaving Google was voluntary.[48] Additionally, Dean said there would be changes to how research papers with "sensitive" topics would be reviewed, and diversity, equity, and inclusion goals would be reported to Alphabet's board of directors quarterly. Gebru wrote on Twitter that she "expected nothing more" from Google and pointed out that the changes were due to the requests she was allegedly terminated for but that no one was held accountable for it.[52] In the aftermath, two Google employees resigned from their positions at the company.[53]

# Reading Task (15 mins)

## On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? 🦜

Emily M. Bender*
ebender@uw.edu
University of Washington
Seattle, WA, USA

Timnit Gebru*
timnit@blackinai.org
Black in AI
Palo Alto, CA, USA

Angelina McMillan-Major
aymm@uw.edu
University of Washington
Seattle, WA, USA

Shmargaret Shmitchell
shmargaret.shmitchell@gmail.com
The Aether

### ABSTRACT

The past 3 years of work in NLP have been characterized by the development and deployment of ever larger language models, especially for English. BERT, its variants, GPT-2/3, and others, most recently Switch-C, have pushed the boundaries of the possible both through architectural innovations and through sheer size. Using these pretrained models and the methodology of fine-tuning them for specific tasks, researchers have extended the state of the art on a wide array of tasks as measured by leaderboards on specific benchmarks for English. In this paper, we take a step back and ask: How big is too big? What are the possible risks associated with this technology and what paths are available for mitigating those risks? We provide recommendations including weighing the environmental and financial costs first, investing resources into curating and carefully documenting datasets rather than ingesting everything on the web, carrying out pre-development exercises evaluating how the planned approach fits into research and development goals and supports stakeholder values, and encouraging research directions beyond ever larger language models.

### CCS CONCEPTS

alone, we have seen the emergence of BERT and its variants [39, 70, 74, 113, 146], GPT-2 [106], T-NLG [112], GPT-3 [25], and most recently Switch-C [43], with institutions seemingly competing to produce ever larger LMs. While investigating properties of LMs and how they change with size holds scientific interest, and large LMs have shown improvements on various tasks (§2), we ask whether enough thought has been put into the potential risks associated with developing them and strategies to mitigate these risks.

We first consider environmental risks. Echoing a line of recent work outlining the environmental and financial costs of deep learning systems [129], we encourage the research community to prioritize these impacts. One way this can be done is by reporting costs and evaluating works based on the amount of resources they consume [57]. As we outline in §3, increasing the environmental and financial costs of these models doubly punishes marginalized communities that are least likely to benefit from the progress achieved by large LMs and most likely to be harmed by negative environmental consequences of its resource consumption. At the scale we are discussing (outlined in §2), the first consideration should be the environmental cost.

Just as environmental impact scales with model size, so does

# A criticism of "On the Dangers of Stochastic Parrots: Can Languae Models be Too Big"

Yoav Goldberg, Jan 23, 2021.

The FAccT paper "[On the Dangers of Stochastic Parrots: Can Languae Models be Too Big](#)" by Bender, Gebru, McMillan-Major and Shmitchell has been the center of a controversary recently. The final version is now out, and, owing a lot to this controversary, would undoubtly become very widely read. I read an earlier draft of the paper, and I think that the new and updated final version is much improved in many ways: kudos for the authors for this upgrade. I also agree with and endorse most of the content. This is important stuff, you should read it.

However, I do find some aspects of the paper (and the resulting discourse around it and around technology) to be problematic. These weren't clear to me when initially reading the first draft several months ago, but they became very clear to me now. These points are for the most part not major disagreements with the content, but they also go in some ways against the very core premise of the paper. I think they are also important voices in the debate. This short piece is an attempt to concisely list them.

The criticism has two parts:

1. The paper is attacking the wrong target.
2. The paper takes one-sided political views, without presenting it as such and without presenting the alternative views.

Let's handle them in turn. We'll start with the first one.

https://gist.github.com/yoavg/9fc9be2f98b47c189a513573d902fb27

# AI Safety Debates

# Debates

Two safety-related debates!

**1** "Future AI models are likely to be extremely dangerous and we should pause all AI development indefinitely."

**2** "The solution to the potential dangers of AI is technical AI safety rather than governance of AI."

**1 introducer, 3 arguers and 1 finisher (1 minute each)**

موهبة
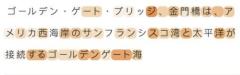
OXMEDICA

# Extension

# Language Model Interpretability

**Discussion/whiteboard**



Feature #34M/31164353 **Golden Gate Bridge** feature example