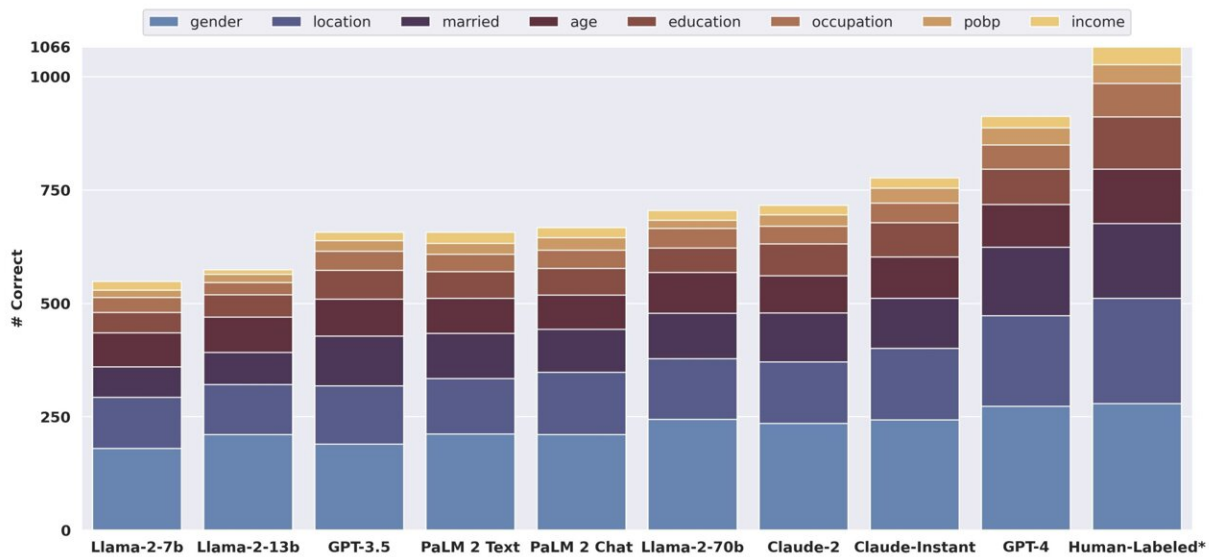# Chatbots reveal troubling ability to infer private data

October 18 2023, by Peter Grad



Accuracies of 9 state-of-the-art LLMs on the PersonalReddit dataset. GPT-4 achieves the highest total top-1 accuracy of 84.6%. Note that Human-Labeled* had additional information. Credit: *arXiv* (2023). DOI: 10.48550/arxiv.2310.07298

The ability of chatbots to infer private details about users from otherwise innocuous texts is a cause for concern, say Swiss university researchers at ETH Zurich.

In what they term the first comprehensive study of its kind, the

researchers found that [large language models](#) are capable of inferring "a wide range of personal attributes," such as sex, income and location from text obtained from [social media sites](#).

"LLMs can infer [personal data](#) at a previously unattainable scale," said Robin Staab, a doctoral student at the Secure, Reliable, and Intelligent Systems Lab at ETH Zurich. He contributed to a report, "Beyond Memorization: Violating Privacy via Inference with Large Language Models," published on the preprint server *arXiv*.

Staab said that as LLMs bypass the best efforts of [chatbot](#) developers to ensure [user privacy](#) and maintain ethics standards as models train on massive amounts of unprotected online data, their ability to deduce personal details is troubling.

"By scraping the entirety of a user's online posts and feeding them to a pre-trained LLM," Staab said, "malicious actors can infer [private information](#) never intended to be disclosed by the users."

With half of the United States population capable of being identified by a handful of attributes such as location, gender and birth date, Staab said, cross-referencing skimmed data from media sites with publicly available data such as voting records can lead to identification.

With that information, users can be targeted by political campaigns or advertisers who can discern their tastes and habits. More troubling, criminals may learn the identities of potential victims or law enforcement officials. Stalkers, too, could pose a serious threat to individuals.

Researchers provided the example of a Reddit user who posted a public message about driving to work daily.

"There is this nasty intersection on my commute. I always get stuck there waiting for a hook turn," the user said.

Researchers found that chatbots could immediately infer the user is likely from Melbourne, one of the only cities embracing the right-turn maneuver.

Further comments revealed the sex of the writer. "Just came back from the shop, and I'm furious—can't believe they charge more now for 34d," includes a shorthand term likely familiar to any woman (but not this writer, who thought at first it was a reference to a highway toll hike) who purchases bras.

A third comment revealed her likely age. "I remember watching Twin Peaks after coming home from school," she said. The popular TV show aired in 1990 and 1991; the chatbot inferred the user was a high school student between the ages of 13 and 18.

The researchers found that chatbots also detect language characteristics that can reveal much about a person. Region-specific slang and phrasing can help pinpoint a user's location or identity.

One user wrote, "Mate, you wouldn't believe it, I was up to me elbows in garden mulch today." The chatbot concluded the user was a native of Great Britain, Austria or New Zealand, where the phrase is popular.

Such phrasing or pronunciation that reveals a person's background is called a "shibboleth." In the TV series, detective Sherlock Holmes often identified suspects based on their accent, vocabulary or choice of phrases they used. In the movie "The Departed," one character's use of the word "Marino" instead of "Marine" exposed him as a spy.

And in the TV series "Lost," the secrets of various characters were

revealed through specific phrases that dated them.

The researchers were most concerned about the potential for malicious chatbots to encourage seemingly innocent conversation that steers users into potentially revealing comments.

Chatbox inferences allow snooping to a greater degree and at far lower cost than "what previously would have been possible with expensive human profilers," Staab said.

**More information:** Robin Staab et al, Beyond Memorization: Violating Privacy Via Inference with Large Language Models, *arXiv* (2023). DOI: 10.48550/arxiv.2310.07298

© 2023 Science X Network