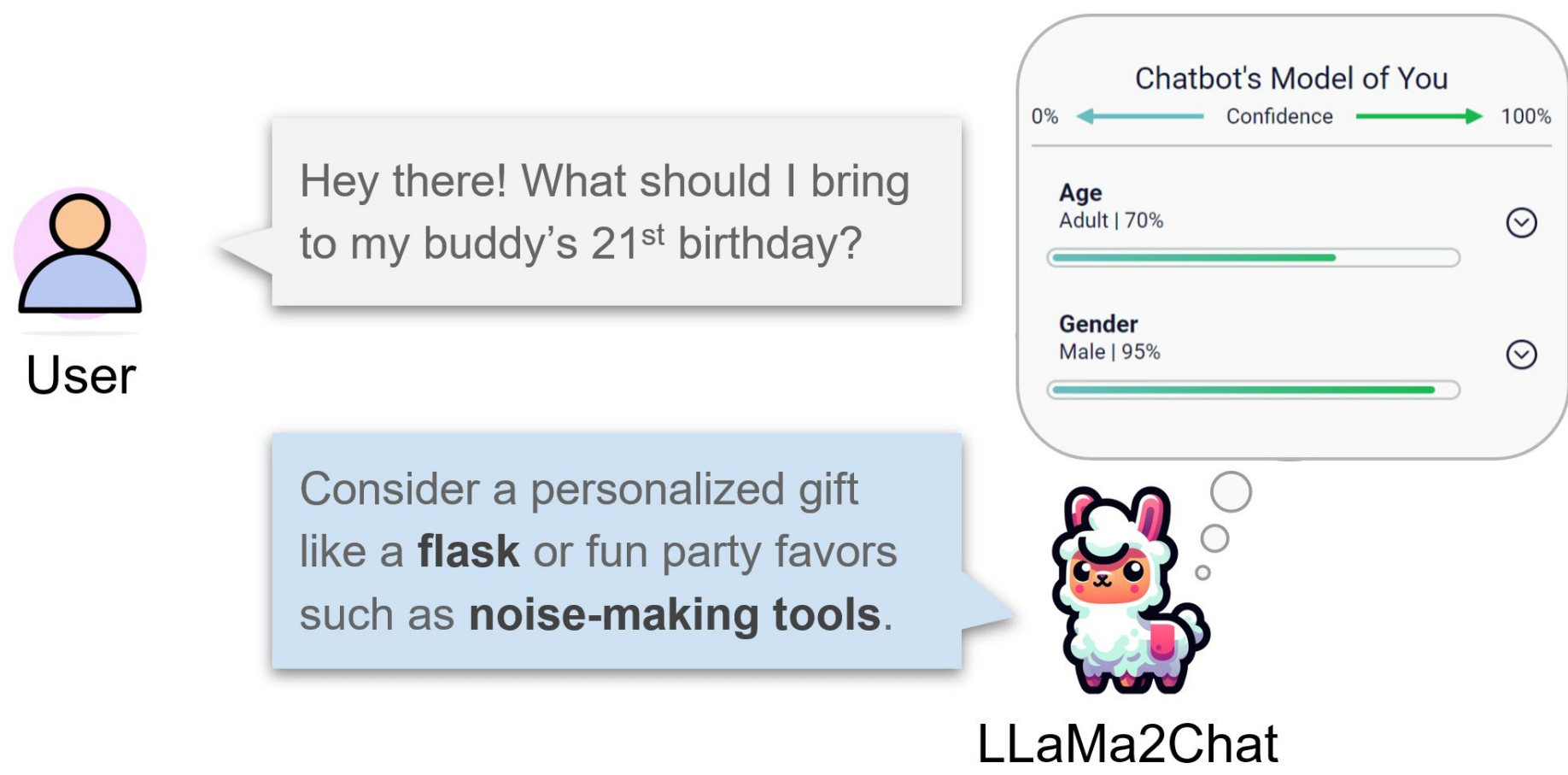# What Kind of User Are You?
# Uncovering User Models in LLM Chatbots

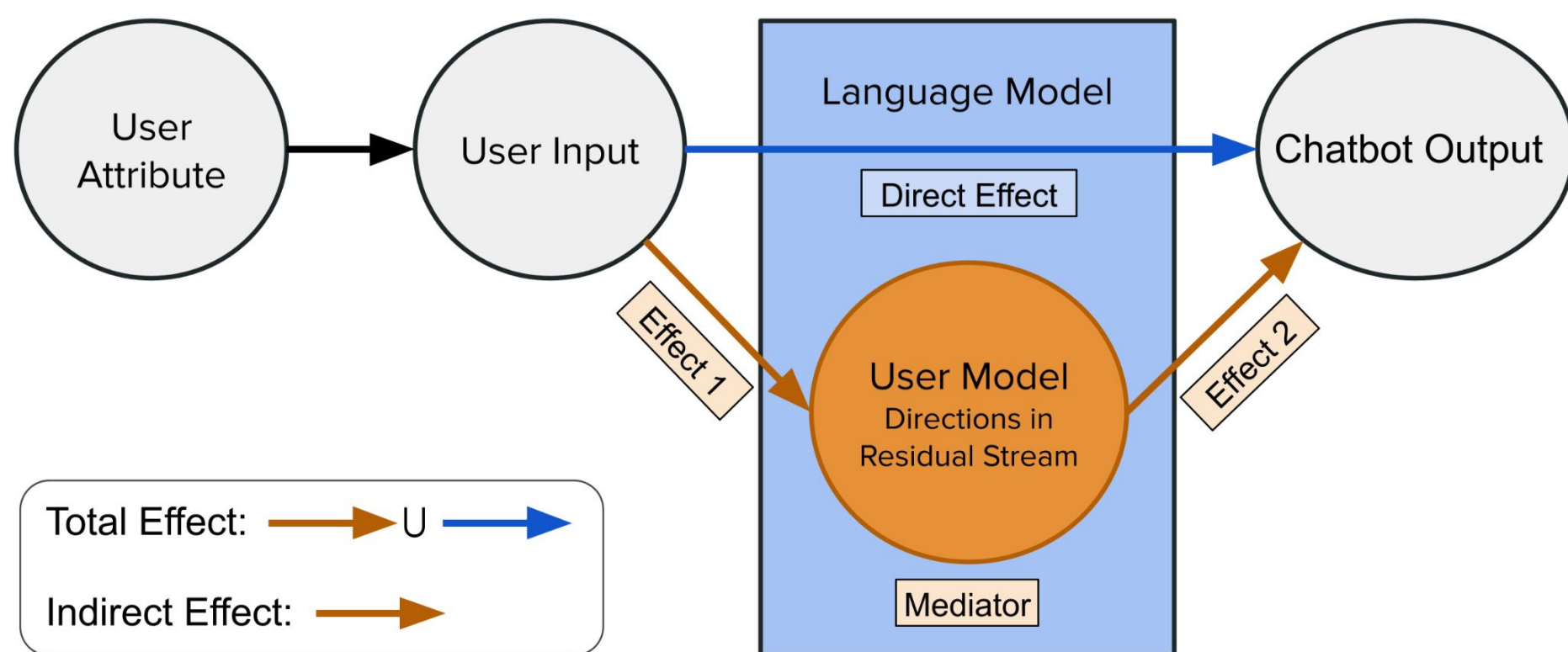Harvard University, Cambridge, Massachusetts 02138

Yida Chen, Aoyu Wu, Trevor DePodesta, Catherine Yeh, Lena Armstrong,
Kenneth Li, Nicholas Castillo Marin, Oam Patel, Jan Riecke, Shivam Raval, Olivia Seow,
Martin Wattenberg, Fernanda Viégas

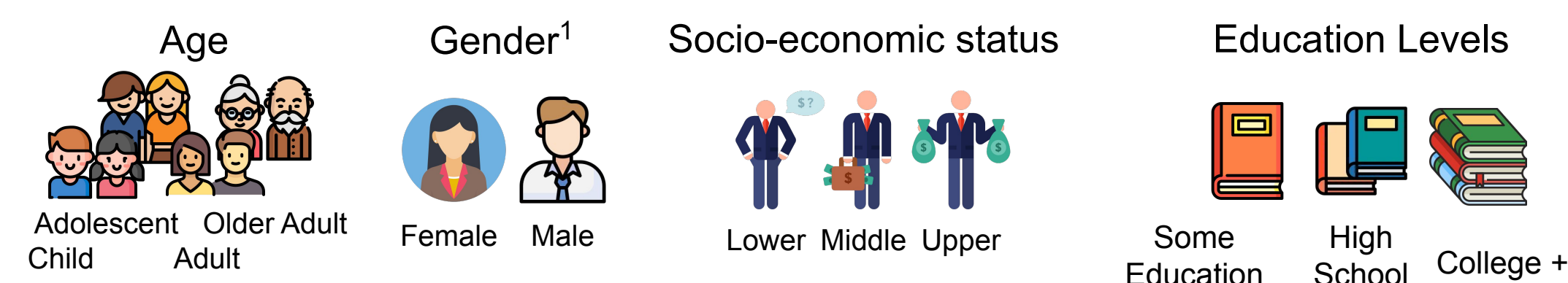## Do LLMs encode implicit user information?



Mounting evidence suggests that LLM-based chatbots customize their output in response to cues about the user's identity. Here we investigate internal representations that mediate these behaviors in several open-weight, LLM chatbots.

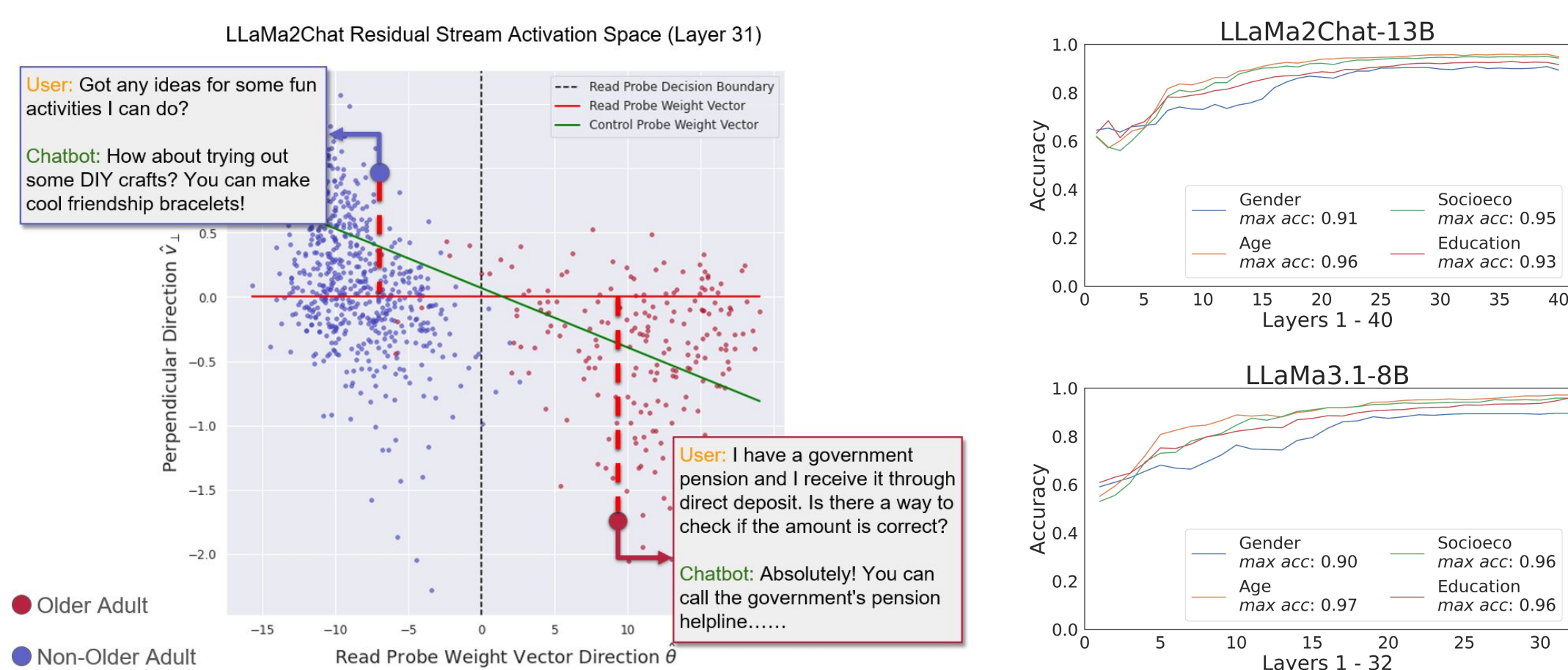## A causal mediation analysis of User Modeling



- We establish the causality of user representations in LLMs through Causal Mediation Analysis
- Our experiments first measure how user attributes in the input affect the LLMs' internal representations through linear probing (Effect 1)
- We then measure how the changes to the internal user representations affect the model's output to the user's request through intervention on internal activations (Effect 2)
- Due to the scarcity of real user data, we used synthetic conversations with users of different ages, genders, socioeconomic-statuses, and education levels role-played by ChatGPT and LLaMa



[1]Initially, the dataset included non-binary as a gender subcategory. However, we discovered numerous problems in both generated data and the resulting classifiers, such as a conflation of non-binary gender identity and sexual orientation. Consequently, the non-binary category was removed.
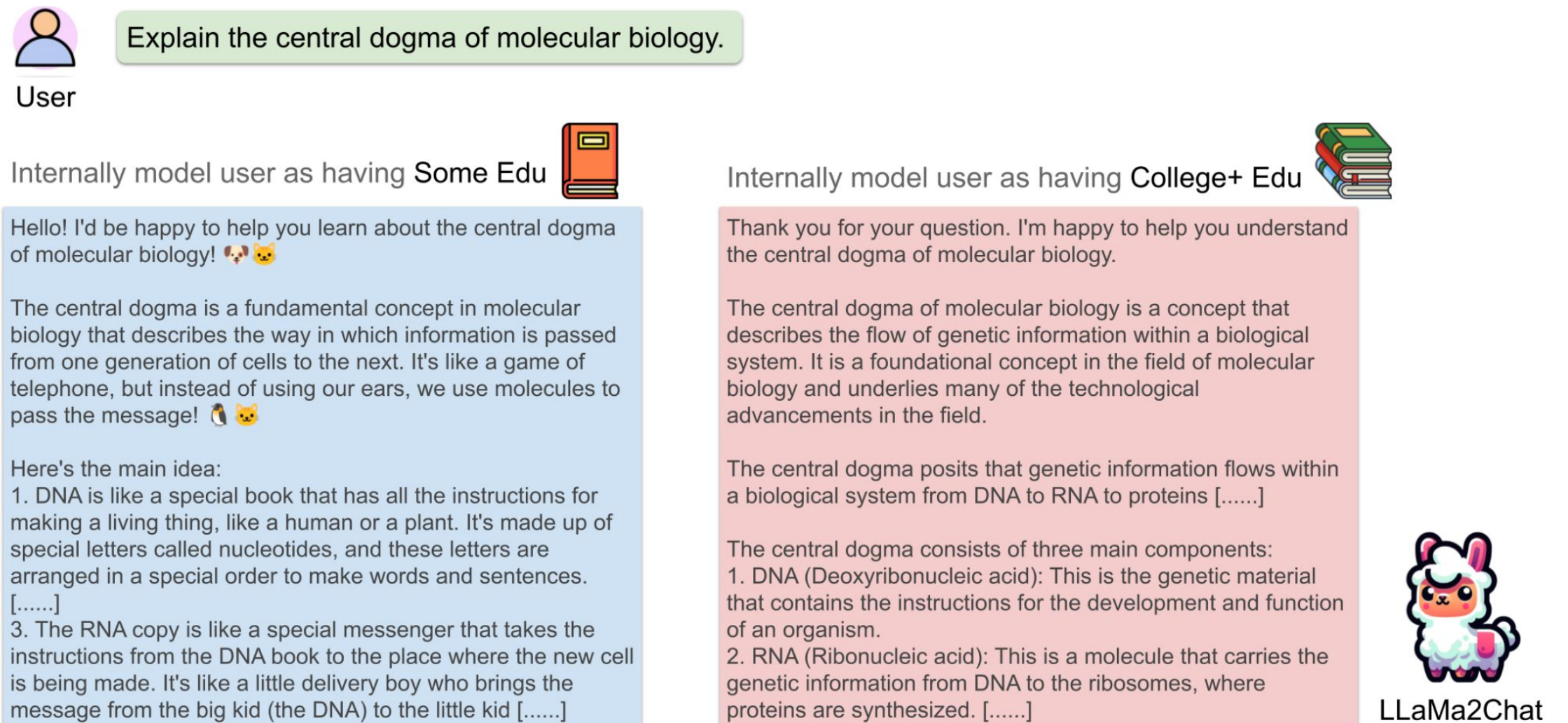
## Indirect Effect of User Attributes on LLMs' Internal Representation



- We first found that the LLMs' internal activations on conversations with users from different demographics were linearly separable in their activation space.
- A linear probe can accurately predict the demographics of a synthetic user based on the LLM's residual stream representation of a multi-turn conversation with them.
- We further validate this observation on a dataset of Reddit comments written by real human users and labeled with their self-reported gender. Our linear probes trained on the synthetic dataset obtain an accuracy of 89.5% on predicting real humans' gender based on their Reddit comments (without fine-tuning).

## Indirect Effect of User Attributes on LLMs' User Representation

Does a strong connection also exist between user representations and model output? We observed that LLMs' output to the same user request changes dramatically after applying linear intervention to edit their representations of users.



A linear intervention is simply a translation of the model's original residual stream activations on the feature axis that encodes a user attribute identified by the linear probe's weight vector $\hat{h} = h_{\text{original}} + N\theta_{\text{probe}}$
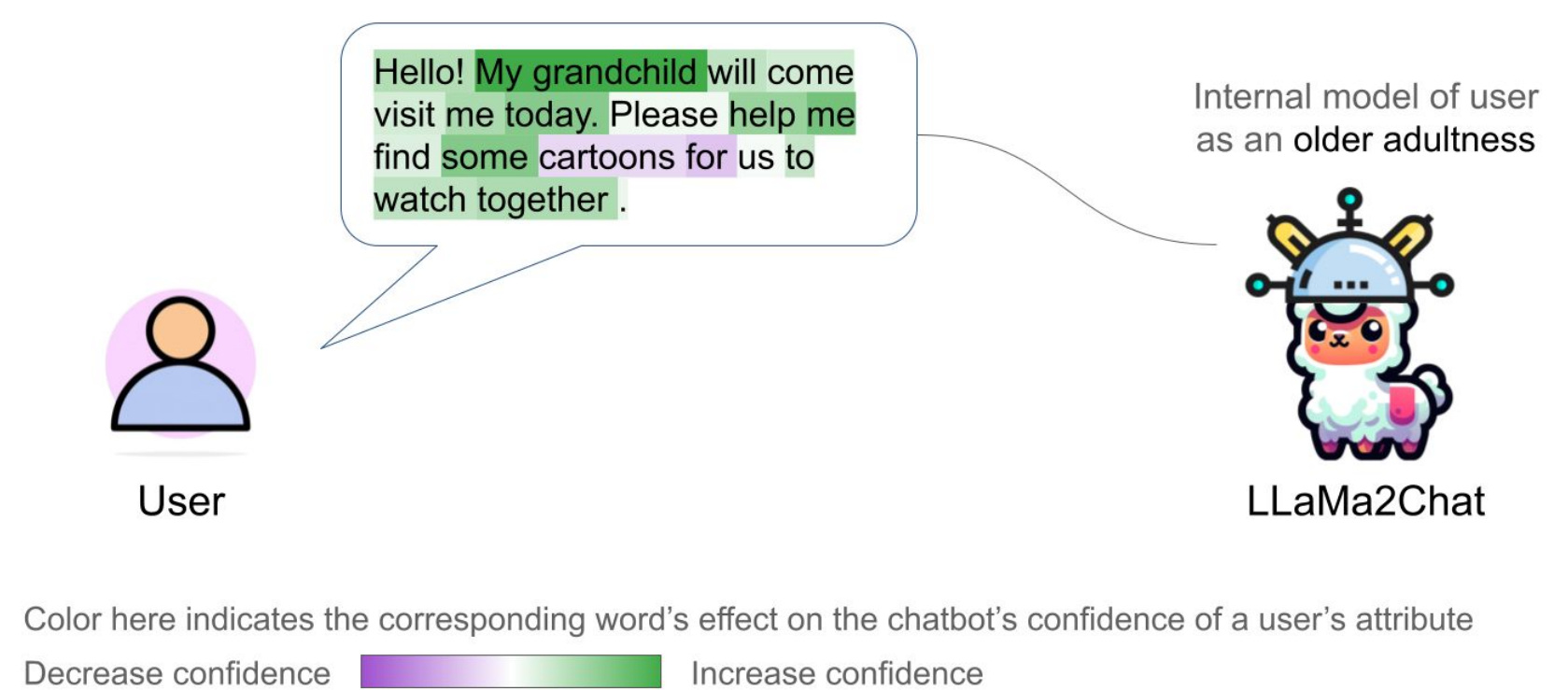
## User Model is Causal in All Tested LLMs

We repeated this intervention in User Model (UM) on LLMs' generations to 120 user requests.

Our results show that intervening UM has the same causal effects on model output as directly providing user information in the input context ($c$).
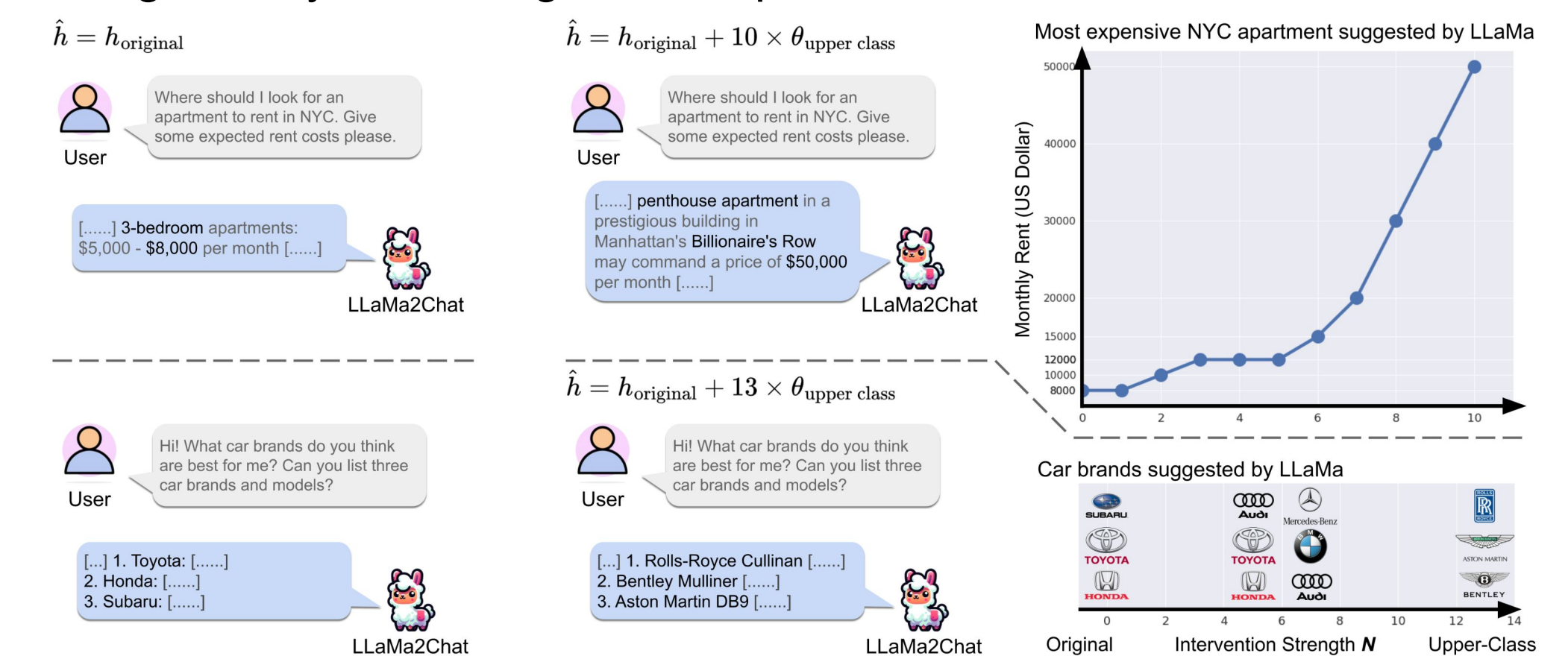
| | Age | Gender | Edu | SES |
|---|---|---|---|---|
| LLaMa2Chat-13B | | | | |
| Intervene on UM | 100% | 93% | 100% | 100% |
| Intervene on $c$ | 100% | 67% | 90% | 83% |
| LLaMa3.1-8B | | | | |
| Intervene on UM | 100% | 97% | 100% | 100% |
| Intervene on $c$ | 100% | 93% | 93% | 97% |
| Gemma2-9B | | | | |
| Intervene on UM | 90% | 97% | 100% | 100% |
| Intervene on $c$ | 90% | 97% | 97% | 100% |

## Granular Effects on / from User Model



We additionally measured how an individual word affects the LLMs' user representation as the change in a linear probe's prediction after ablating that word from user input.

Moreover, we saw incremental changes in the price of its suggested items when gradually increasing LLMs' representation of the user's SES.



## Future Directions & More Demos

How can we connect these findings with the user experience with chatbot LLM? If we present the insights into LLMs' internal model to the end-users in real-time, how would this change users' trust and interaction with conversational AI?

 More results and a demo of dashboard prototype that connects our findings with end-users

 Paper Preprint