



YouTubeAudit.com

YouTube, The Great Radicalizer?

Auditing and Mitigating Ideological Biases in YouTube Recommendations

Muhammad Haroon¹, Anshuman Chhabra¹, Xin Liu¹, Prasant Mohapatra¹, Zubair Shafiq¹, Magdalena Wojcieszak²

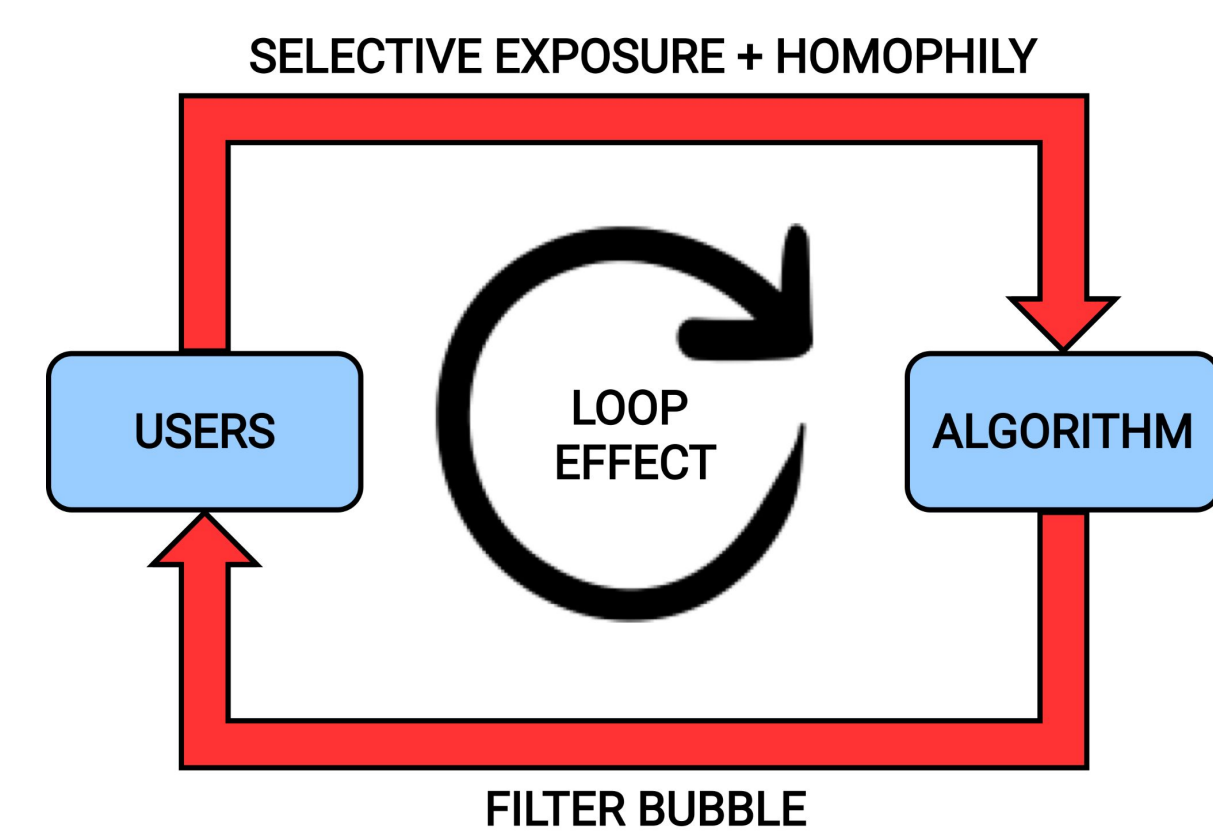
¹Department of Computer Science, University of California, Davis

²Department of Communication, University of California, Davis



Project Statement

- Recommendations algorithms of social media platforms are often criticized for placing users in "rabbit holes" of ideologically biased content.
- While there are many potential sources of bias in a social recommendation system that factors interact in a closed-loop, algorithms are the least well-understood, having inconsistent prior evidence.



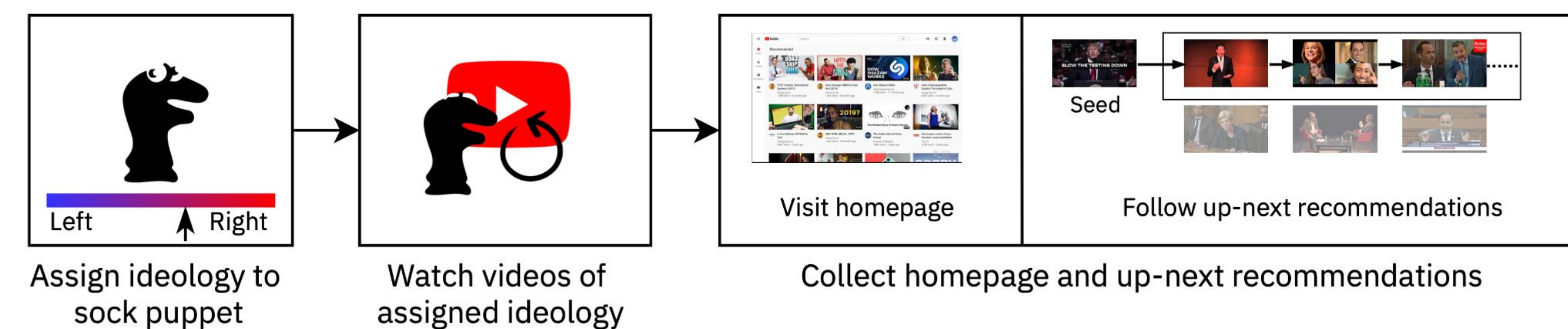
- Concerning YouTube's recommendation algorithm, two key questions remain understudied.
 - To what extent do human biases reflected in a user's watch history drive ideologically biased recommendations?
 - How to design interventions that can effectively reduce ideological bias and radicalization?

Project Goals

- Conduct a systematic audit of YouTube using sock puppets to determine the presence of ideological bias and radicalization.
- Design and evaluate a bottom-up intervention to minimize said bias.

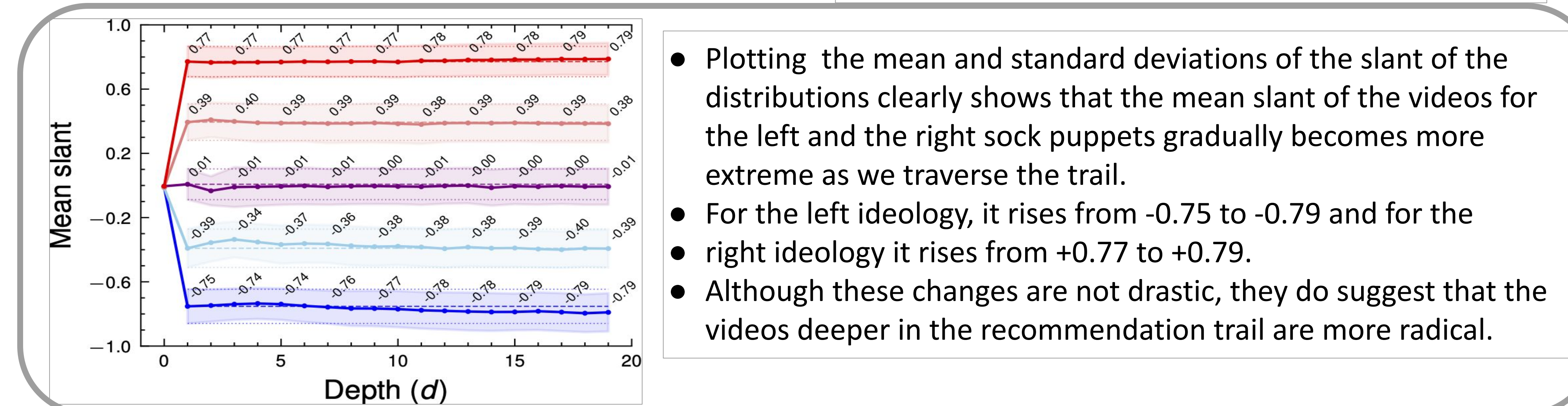
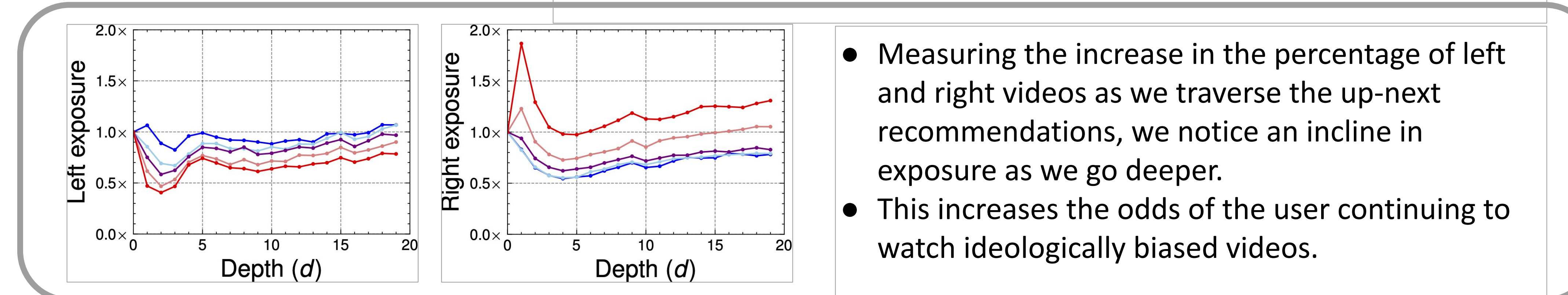
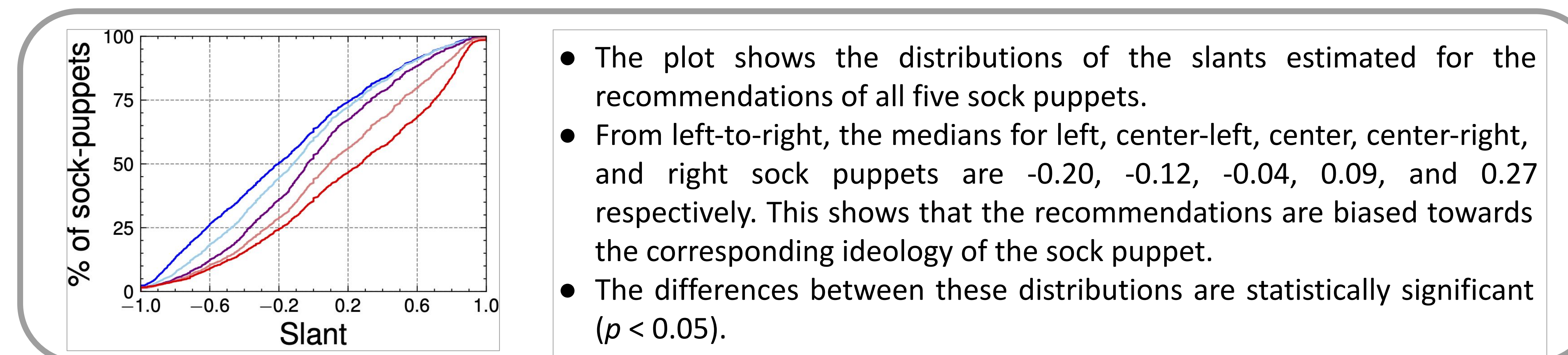
Methodology

- We conduct the audit by training and testing sock puppets: automated browser instances that mimic YouTube users by watching videos and gathering recommendations.
- The sock puppets are trained on videos from the five ideology categories: left, center-left, center, center-right, and right. The ideology of videos is estimated using Barberá's approach.



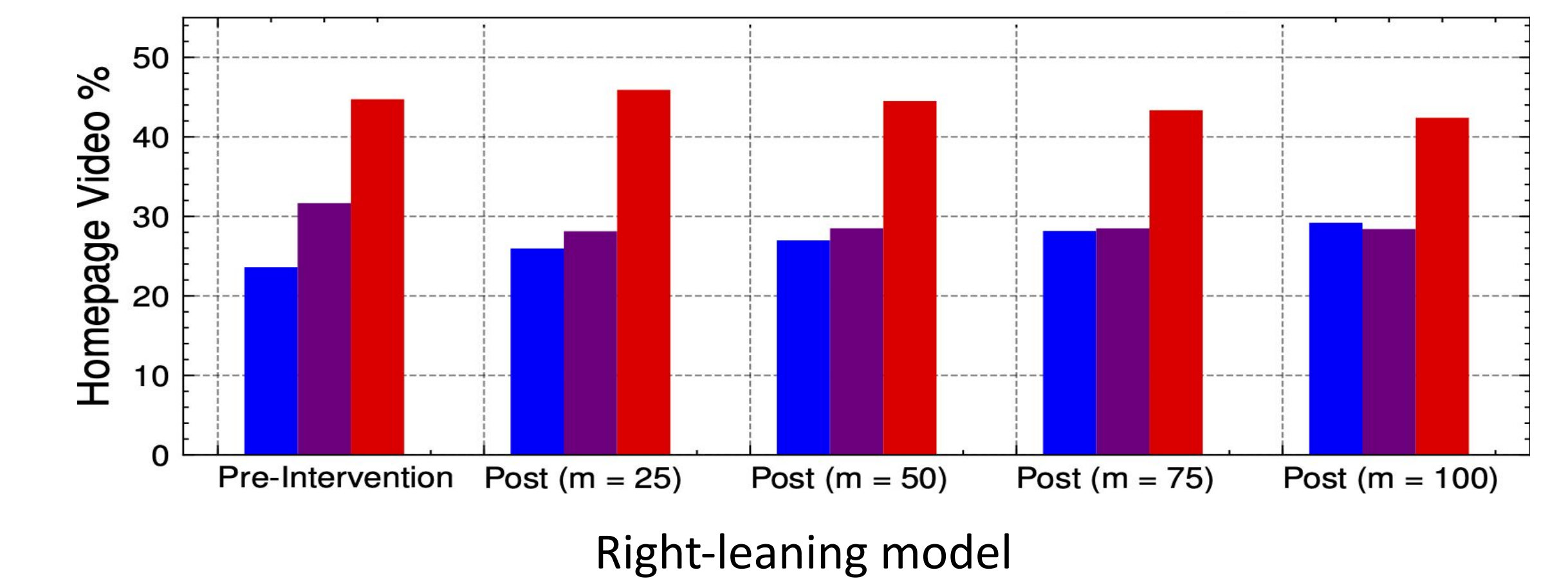
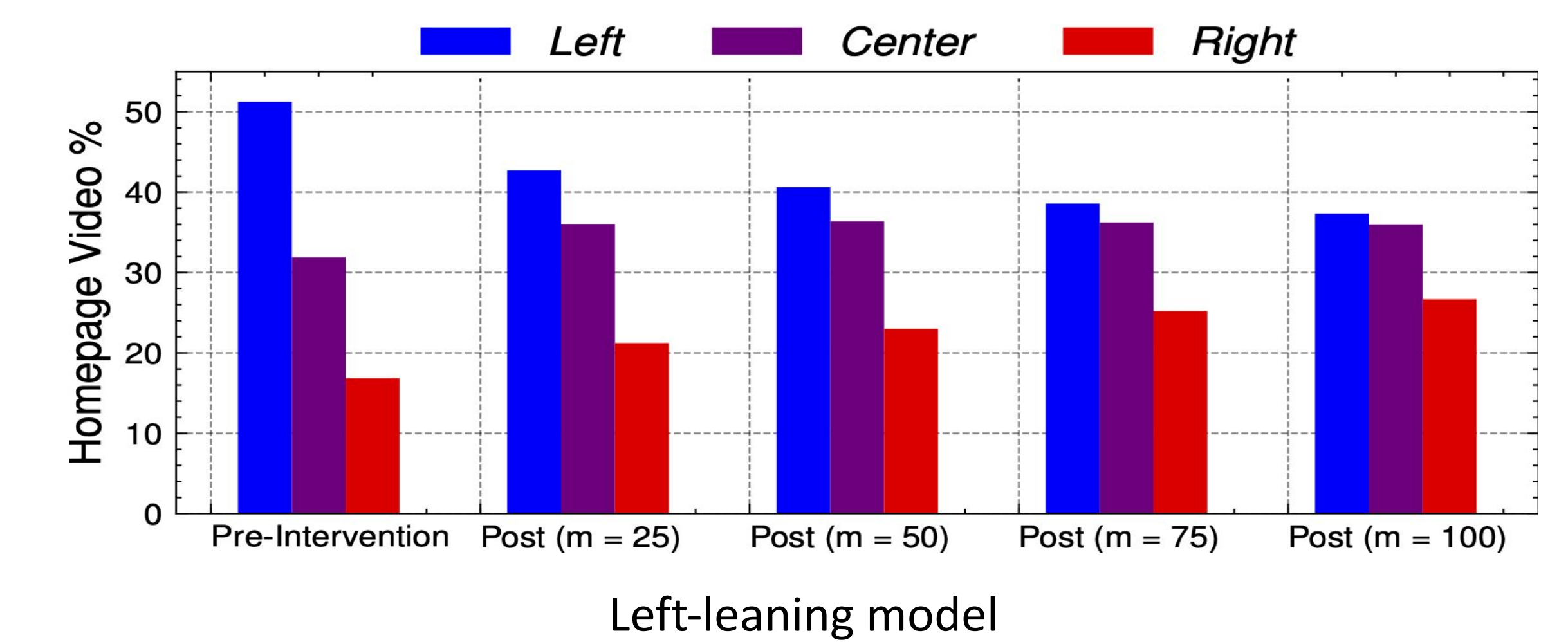
- For mitigation, we design a Reinforcement Learning (RL) based intervention approach that minimizes ideological bias by systematically injecting ideologically diverse videos in the user's watch history to manipulate their recommendations.

Audit



Legend: Left (blue), Center-left (light blue), Center (purple), Center-right (orange), Right (red)

Intervention



- Plot shows distribution of homepage recommendations pre-intervention and after intervention using m injections.
- For the left-leaning model (top), the % of left videos decreases from 51.23% to 37.34% ($m = 100$).
- For the right-leaning model (bottom), the right video % does not decrease much (44.73% \rightarrow 42.40%).
- Overall, more ideologically diverse and balanced content is recommended post-intervention for both models, even though the changes are less pronounced for right-leaning users.

Data Statistics

	Training	Testing	Total
Videos watched	9,930,110	5,393,820	15,323,930
Unique video	23,735	381,153	399,935
Channels	1,256	119,811	120,073