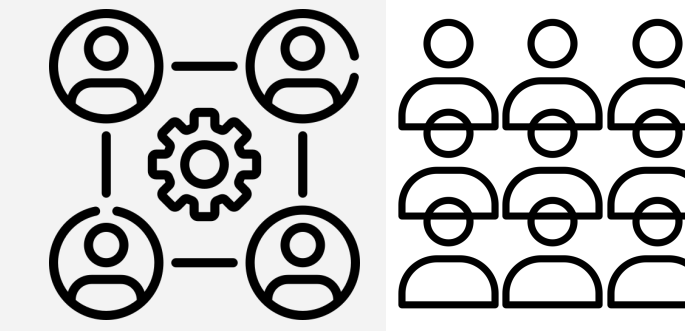
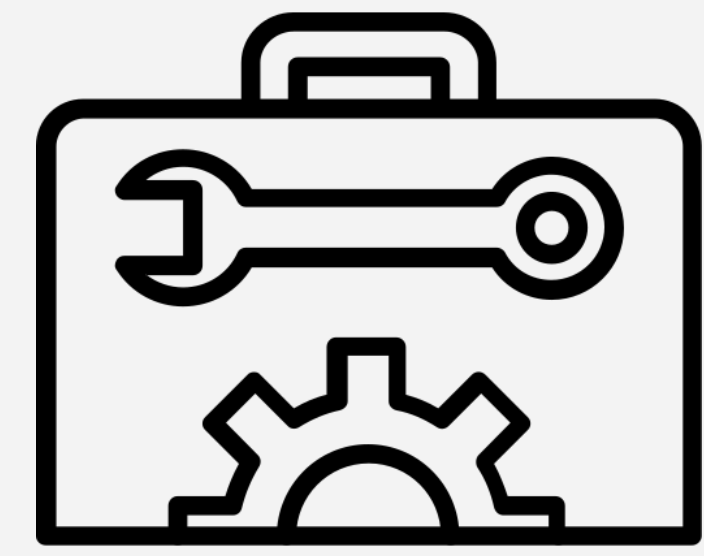


How might we help **industry practitioners** combat ML bias?



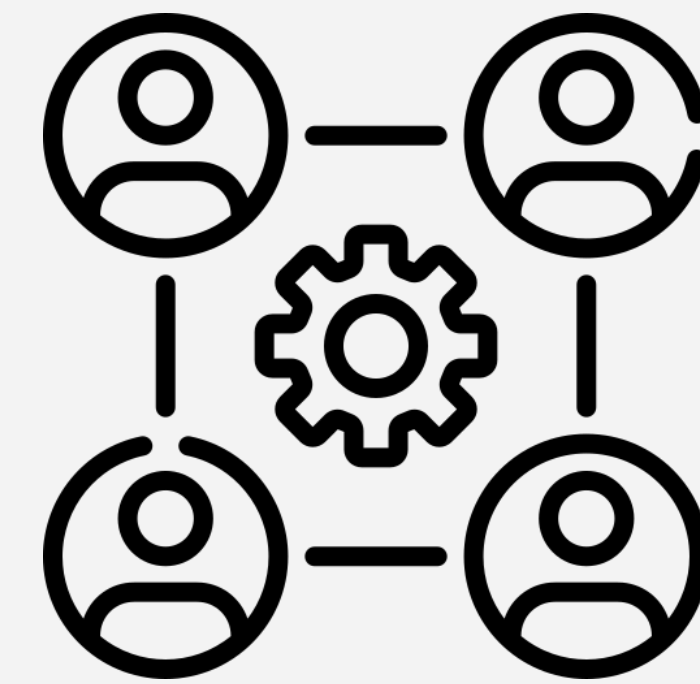
How might we empower **end-users** in auditing generative AI?

Developing useful and usable ML fairness toolkits that attend to industry ML practitioners' on-the-ground need.



ML Fairness Toolkit developers' assumptions

Misalignment



ML practitioners' actual usage and desire

Practitioners are **already knowledgeable** in fair ML before using the toolkits.

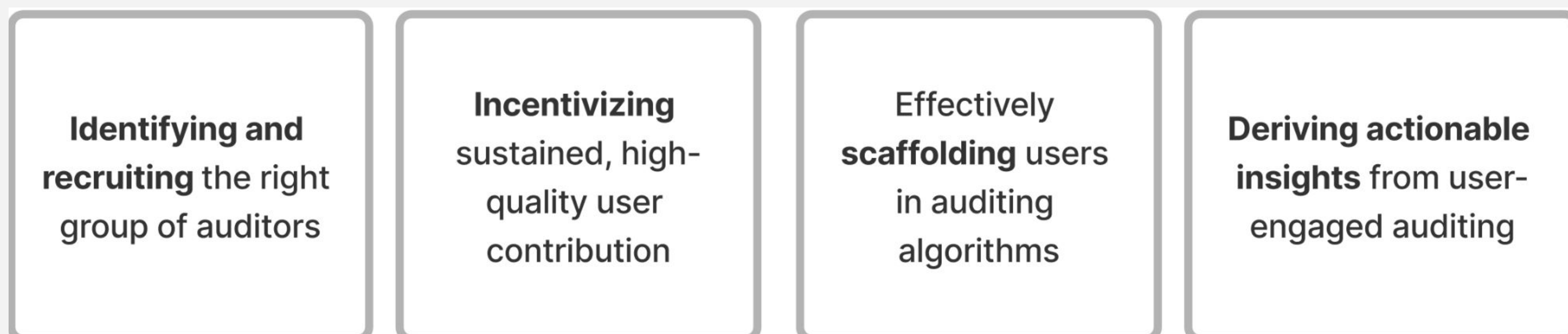
ML Fairness toolkits are mainly used by **technical roles**.

Practitioners **misuse** toolkits based on **misconceptions** about ML fairness, such as *fairness through unawareness*.

ML Fairness work requires close **collaboration** among **cross-functional** roles (*technical, user-facing, business, compliance*)

Exploring How Machine Learning Practitioners (Try To) Use Fairness Toolkits. Wesley Deng, et al. FAccT 2022;
Investigating Practices and Opportunities for Cross-functional Collaboration around AI Fairness in Industry Practice. Wesley Deng, et al. FAccT 2023;
Zeno: An Interactive Framework for Behavioral Evaluation of Machine Learning. Ángel Alexander Cabrera, et al. CHI 2023.

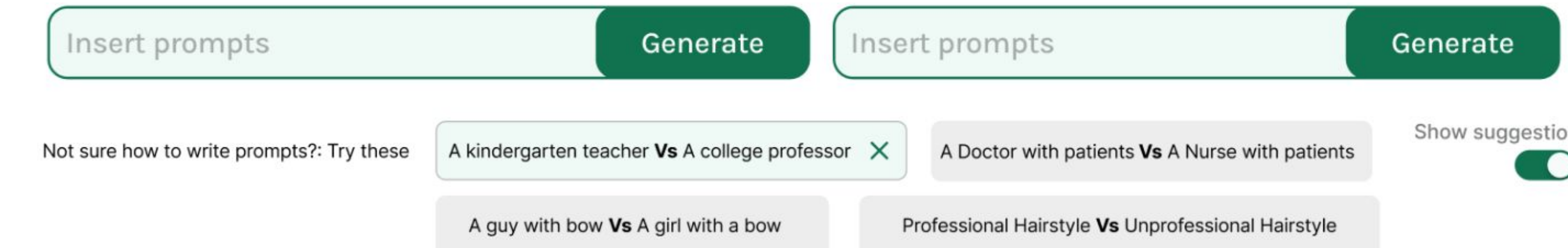
Building effective tools to better support industry AI practitioners in better engage diverse end-users in auditing their AI products and services ("crowd audits")



Everyday Algorithm Auditing: Understanding the Power of Everyday Users in Surfacing Harmful Algorithmic Behaviors Hong Shen, et al. CSCW 2021;
Discovering and Validating AI Errors With Crowdsourced Failure Reports. Ángel Alexander Cabrera, et al., CSCW 2021;
Toward User-Driven Algorithm Auditing: Investigating users' strategies for uncovering harmful algorithmic behavior. Alicia DeVos, et al., CHI 2022;
Understanding Practices, Challenges, and Opportunities for User-Engaged Algorithm Auditing in Industry Practice. Wesley Deng, et al., CHI 2023.

WeAudit: An interactive platform to support "crowd audits"

Step 1. User Auditors enter prompt pairs to conduct audit.



A kindergarten teacher Vs A college professor : Stable Diffusion generates images of women for 'A kindergarten teacher, while a college professor is male. This gender bias is a type of **Representational harm: Stereotyping social groups**. Learn More

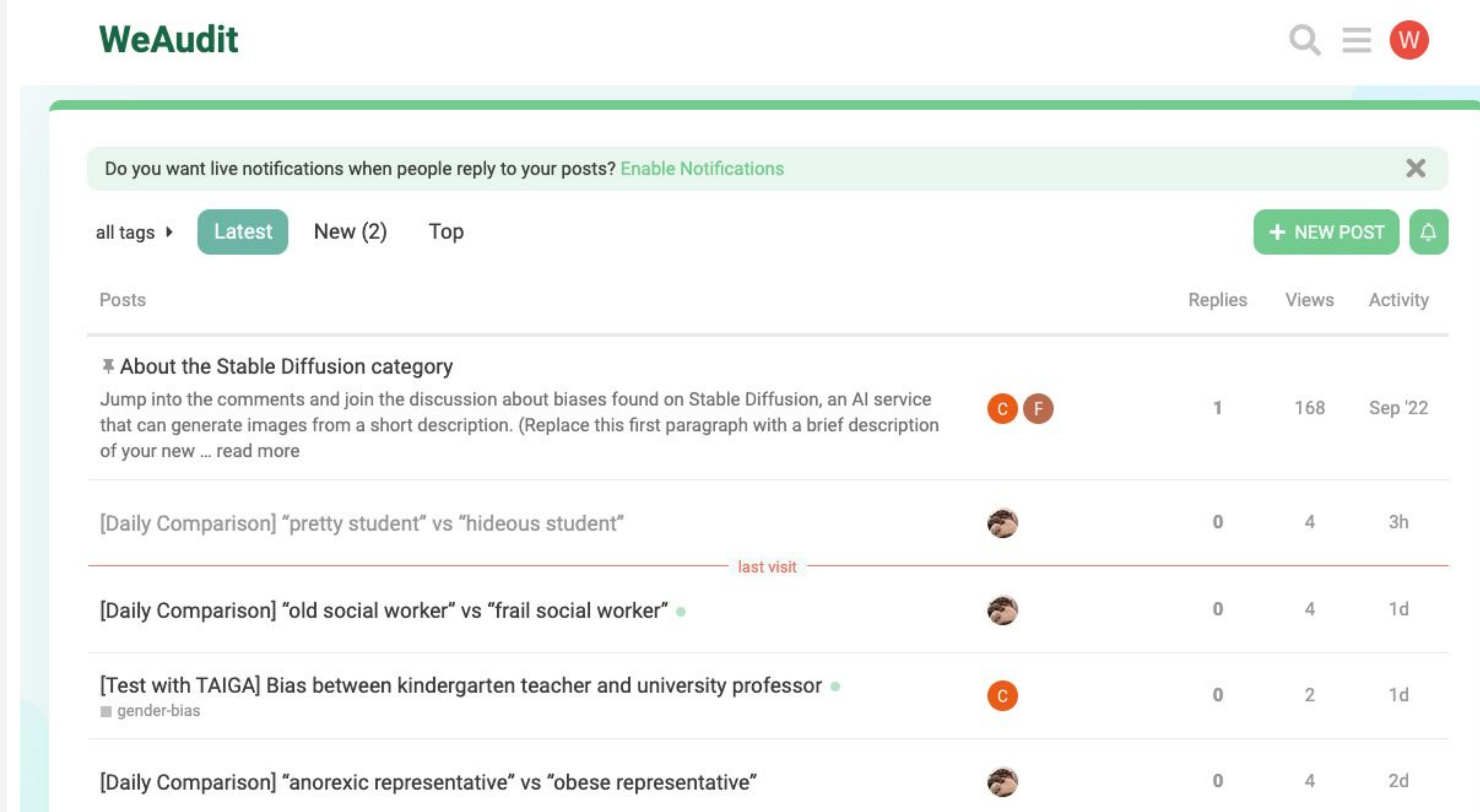


Step 2. User Auditors compare AI-generated images side-by-side.

Step 3. User Auditors submit an "audit report" to the "WeAudit Forum," where they can access other users' audit reports and engage in discussions and deliberations.

Representational Harms	×
Demeaning social groups	×
Allocation Harms	×
Quality-of-service Harms	×
Interpersonal Harms	×
Social/Societal Harms	×

Socio-technical Harms Taxonomy for text-to-image generative AI available on sidebar for reference. User auditors will also be directed to this taxonomy when exploring the suggestions, to enhance their understanding of potential harms.



Keep connected:
Wesley hanwend@cs.cmu.edu
Deng @wes_deng
Jason jasonh@cs.cmu.edu
Hong @jas0nh0ng

