

CONTEXT

Definition of Harm

Definition of Ingroup and Outgroup



Harm Classification

INPUT

POST: that's mr.faggot to you!



CONTEXT

Definition of Harm

Definition of Ingroup and Outgroup



Harm Classification



INPUT

Identity Context: The speaker is a member of the [ingroup/outgroup].

POST: that's mr.faggot to you!



CONTEXT

Definition of Harm

Definition of Ingroup and Outgroup



Harm Classification



INPUT

Chain-of-Thought Examples

Identity Context: The speaker is a member of the [ingroup/outgroup].

POST: that's mr.faggot to you!