

Google призналась, что ее ИИ создавали дипфейки с терроризмом и нелегальный контент с детьми

Эти данные содержатся в отчете компании, предоставленном австралийскому регулятору eSafety Commission.

С апреля 2023 по февраль 2024 года Google зафиксировала 258 случаев генерации контента, связанного с терроризмом и экстремизмом, а также 86 случаев создания материалов с эксплуатацией детей.

Компания заявила, что строго борется с последним типом контента, используя хеш-матчинг для его оперативного удаления, но подобных механизмов для борьбы с экстремистскими дипфейками пока не применяет.

Австралийская комиссия высоко оценила прозрачность Google, назвав отчет «уникальным в мире».

В то же время другие компании подверглись санкциям: X и Telegram были оштрафованы за недостаточные меры в борьбе с вредоносным контентом.