

Relatório sobre os Desafios da Moderação de Conteúdo em Plataformas Digitais

Análise Crítica da Moderação de Conteúdo sob a ótica de Justiça, Transparência, Impacto Social e Governança.

Introdução

A moderação de conteúdo é o processo pelo qual as plataformas online monitoram e gerenciam o conteúdo criado pelo usuário para garantir a conformidade com suas políticas e diretrizes. Em 2025, essa tarefa evoluiu de uma simples remoção de spam para uma complexa operação global que se situa na intersecção da tecnologia de IA, direitos humanos, segurança pública e ética. Este relatório analisa os desafios e as responsabilidades inerentes a essa atividade através de quatro pilares fundamentais.

1. Viés e Justiça

A busca pela justiça na moderação é minada por vieses intrínsecos, tanto humanos quanto algorítmicos, que podem levar à aplicação desigual das regras e à supressão desproporcional de vozes.

Viés Algorítmico: Os sistemas de IA, responsáveis pela triagem inicial de bilhões de posts diariamente, são treinados com dados históricos que refletem preconceitos sociais existentes. Consequentemente, esses sistemas podem identificar incorretamente como "tóxicas" ou "ofensivas" as linguagens e dialetos de grupos minoritários ou marginalizados, enquanto falham em detectar formas mais sutis de discurso de ódio de grupos dominantes, ocasionado por mecanismos distintos como viés de representação, viés de medição, viés de agregação ou viés de avaliação. A proteção contra a desinformação, por exemplo, pode acabar por silenciar ativistas e fontes de notícias independentes se o algoritmo for treinado para favorecer narrativas estatais ou de fontes tradicionais.

Viés Humano e Cultural: A moderação humana não é imune a preconceitos. Moderadores, muitas vezes trabalhando sob intensa pressão e com diretrizes ambíguas, aplicam suas próprias lentes culturais e políticas ao julgar o conteúdo. Uma política de remoção de conteúdo "politicamente sensível", por exemplo, pode ser interpretada de maneiras drasticamente diferentes por moderadores em São Paulo, Jacarta ou Dublin. Isso resulta em inconsistência e na percepção de que a moderação é arbitrária e tendenciosa, erodindo a confiança do usuário na justiça do sistema.

2. Transparência e Explicabilidade

A falta de clareza sobre "por que" uma decisão de moderação foi tomada é uma das maiores fontes de frustração para os usuários e um obstáculo à legitimidade das plataformas.

Transparência: As plataformas devem ir além de publicar suas diretrizes. A transparência efetiva envolve a publicação de relatórios detalhados e periódicos que quantifiquem o volume de conteúdo removido, as categorias de violação, o número de apelações e a taxa de reversão dessas decisões. Podemos separar esse ponto em três modelos diferentes a transparência técnica, transparência operacional e a transparência de propósito. A clareza sobre como as políticas são criadas e atualizadas, incluindo o envolvimento de especialistas externos, é fundamental para que o público possa escrutinar e entender o processo, mesmo que haja limitações no processo.

Explicabilidade: Quando um conteúdo é removido, a notificação genérica de "violação dos padrões da comunidade" é insuficiente. O usuário tem o direito de saber qual regra específica foi violada, idealmente com a indicação do trecho ou elemento do conteúdo que motivou a ação. Considerando que a automação exige a necessidade de um ser humano na participação direta de resoluções no meio digital, não devemos nos alienar com a obrigação dos prestadores do serviço ou produto em sua responsabilidade ética de sua utilização. Além disso, o processo de apelação deve ser claro, acessível e oferecer uma reavaliação genuína da decisão inicial, idealmente por uma equipe diferente da que tomou a decisão original.

3. Impacto Social e Direitos

As decisões de moderação não ocorrem no vácuo; elas têm impactos profundos e diretos na sociedade, influenciando o debate público, a segurança e os direitos fundamentais.

Liberdade de Expressão vs. Segurança: Este é o dilema central. Uma moderação excessivamente rigorosa pode criar um sentimento de censura, onde usuários evitam discutir tópicos controversos por medo de sanções. Por outro lado, uma moderação frouxa permite a proliferação de discurso de ódio, assédio, incitação à violência e desinformação, que podem causar danos reais, minar processos democráticos e colocar

em risco a saúde pública. No entanto, é válido salientar que a constituição brasileira apesar de em seu artigo 5º, inciso IV, permitir a capacidade de se expressar livremente deixa claro que tal ato deve estar vinculado ao ator, ou seja, qualquer dano que o ato realizar, terá uma pessoa atrelada que será devidamente responsabilizada. Dessa forma a segurança será garantida permitindo reparar qualquer dano.

Pressão Governamental e Censura: Em todo o mundo, governos utilizam suas leis locais para pressionar plataformas a remover conteúdo crítico ou dissidente. A moderação de conteúdo torna-se, assim, um campo de batalha para os direitos humanos, onde as empresas devem decidir entre cumprir uma ordem legal local que pode violar padrões internacionais de liberdade de expressão ou resistir e arriscar o bloqueio de seus serviços no país.

4. Responsabilidade e Governança

A imensa influência das plataformas exige modelos robustos de governança e mecanismos claros de responsabilização para garantir que seu poder seja exercido de forma ética e responsável.

Responsabilidade Corporativa: As empresas de tecnologia devem assumir a responsabilidade pelo ecossistema de informação que criam. Isso inclui investir pesadamente em equipes de moderação bem treinadas e com apoio psicológico adequado, em tecnologia de detecção de vieses e em pesquisa e desenvolvimento de políticas mais justas e contextuais.

Governança e Supervisão Externa: Modelos como o Conselho de Supervisão da Meta representam um passo em direção à supervisão independente, permitindo que um corpo externo de especialistas revise decisões complexas e estabeleça precedentes. No entanto, a soberania final ainda reside na plataforma.

O Papel da Regulação: Leis como o *Digital Services Act (DSA)* da União Europeia e as discussões em torno de projetos como o PL 2630, juntamente com a atual PL 2628/22, no Brasil indicam um movimento global em direção à correção. O objetivo não é que o governo dite o que pode ser dito, mas que estabeleça obrigações de transparência, gestão de risco e o devido processo para as plataformas, tornando-as legalmente responsáveis por seus sistemas, e não por cada conteúdo individual.

Conclusão

A moderação de conteúdo em 2025 é uma tarefa de governança global com implicações diretas na justiça social, nos direitos humanos e na estabilidade democrática. Não há soluções fáceis. O caminho a seguir exige uma abordagem multifacetada, combinando aprimoramento tecnológico com supervisão humana, transparência radical, processos de apelação justos e um modelo de governança que equilibre a responsabilidade corporativa com uma regulação inteligente e protetora dos direitos fundamentais.