



Institute for a Better Internet

Trusted Online Content Moderation

Mike Matthys
mikem@betterinternet.us

Trusted Online Content Moderation

- **Current social media content moderation is not trusted**
 - People on both sides view content moderation as being biased and broken
- **Proposed Solution**
 - Define generalized content moderation guard rails for each platform's independent rules
 - Certify platforms that transparently publish content moderation rules
 - Outsource appeals of user/platform content disputes to trusted non-govt 3rd party
- **Benefits**
 - Reduce legal costs, PR issues, and costs of handling appeals of disputes for the companies
 - Provide trusted non-govt solution that embraces both online safety and neutrality



Guard Rails of Content Moderation

- **Safety**

- Protect users from content that is criminal or imminently harmful

- **Neutrality**

- Don't favor one side on an issue – except to protect against imminent harm

- **Transparency**

- Publish clear & detailed content moderation rules as well as enforcement actions

- **Accountability**

- Appeals outsourced to independent non-govt entity with power to ensure platform compliance



Ensuring Online Safety

Criminal and imminently harmful content should be blocked

- **Criminal content**

- Exhorting violence, criminal/terrorist planning, child porn, trafficking, extortion

- **Cyber-security threats & viruses**

- **Imminently harmful content**

- Hateful/bullying speech targeting specific non-public persons, interference in election voting, threats of violence, doxxing

- **Disinformation from designated terrorist orgs & overseas governments**



Ensuring Neutrality

- **Avoid taking sides on controversial issues**
 - Avoid role of true/false arbiter - except to protect against imminent harm to a person(s)
 - True/false test inevitably breaks neutrality and perception of neutrality “Who decides?”
- **Ensure platform algorithms amplify or depress content on neutral basis**
 - Protect against viewpoint amplification/SPAM from bots, troll farms & inauthentic user networks
- **Ensure enforcement is consistent with published rules & applied fairly to all users**
 - Violations of safety rules enforced regardless of user popularity, viewpoint or political affiliation
- **Government officials can offer info, but inherently cannot act as neutral fact-checkers**



Neutrality is not Uncontrolled Free Speech

- **Content moderation innovations within the Guard Rails are encouraged**
 - Content moderation can contribute to Quality of Service & Features
- **Some categories of content moderation allowed, even if not imminently harmful**
 - Adult content
 - Requirement to authenticate users
 - Spam and click-bait for purpose of commercial gain
- **Smaller platforms and user sub-groups may be publicly non-neutral**
 - Platforms smaller than threshold – must publish their non-neutral standards
 - Examples - Reddit subgroups, family-friendly groups, religious groups, etc
- **Key is for any content moderation to be Safe, Neutral & Transparent**



Ensuring Transparency

- **Publish content moderation standards & related enforcement actions**
 - Users & content creators should easily know the specific rules/standards they violated
- **For any enforcement action, provide clear explanation & justification to user**
 - Specific content rules broken, explanation of review/adjudication process, steps to appeal
 - Involvement and communications between platform and 3rd party fact-checkers
- **Identify and publish qualifications of 3rd party fact-checkers**
 - Publish their names, qualifications, affiliations, funding, and history of previous decisions
- **Report all communications to/from government employees/contractors within 24 hours**



Ensuring Accountability

- **Platforms benefit when user appeals outsourced to trusted 3rd party**
 - 3rd party entity can ensure perception of neutrality/fairness between users and platforms
- **FINRA as a model for a non-govt trusted entity with enforcement powers**
 - FINRA, funded by financial industry, outsources appeals of investor/advisor disputes
 - Governance designed for neutrality/fairness; no hire/fire influence from govt or industry
 - Access to 8000 arbitration judges for handling scale volume of online appealed disputes
 - FINRA has power to levy material fines, suspend/ban financial advisors to ensure compliance
- **Proposed OMRA entity would operate similar to FINRA – for online platforms**
 - Platforms would need to accept OMRA judgements for legitimacy
 - Industry-wide buy-in ensures benefits to platform companies (PR, govt anti-trust, legal)
 - FINRA = Financial Industry Regulatory Authority OMRA = Online Media Regulatory Authority



True/False Standard Inevitably Breaks Trust

- **Who Decides ? Is it true/false or just a disagreement ?**
 - Whichever choice is made on controversial public topics, a large share of users will claim bias
- **Misinformation/Disinformation concept is disconnected from online safety**
 - Misinformation should only be blocked to prevent imminent harm/danger
- **Outsourcing to third-party fact-checkers does not solve the problem**
 - Fact-checkers, especially govt or media-based fact-checkers, have their own biases, politics & funders
- **Harm of censorship is worse than the harm of falsehoods**
 - Progress and fairness happen when all voices can be heard
 - If the government (or monopolies on govt's behalf) acts as arbiter of truth, this leads to authoritarianism



A Better Standard: Imminently Harmful or Not

- **Harmful/Not Harmful is a test that supports both Safety and Neutrality**
- **Imminently Harmful/Not Harmful is easier to adjudicate than True/False**
 - Avoid issues of “ Who Decides? ” Avoid disagreements among subject-matter “experts”
 - Arbitration judges capable of assessing harm without needing subject expertise
- **Scalable across millions of contentious topics, disagreeable opinions, etc**
 - Avoid controversial disagreements between users and company-selected media fact-checkers
- **Should “misinformation” be moderated ? Only if it leads to imminent harm**



IBI's Proposal:

New entity to outsource appeals & certify platforms' content rules

- **Tweak Section 230 to adjust legal incentives for content moderation**
 - Replace “Objectionable” with “Imminently Harmful”
- **Create a non-govt entity (modeled after FINRA) for 4 functions:**
 1. Outsource appeals of selected user/platform disputes
 2. Certification of each platform for operating within the Guard Rails
 3. Review data reports periodically received from companies
 4. Review/publish 24-hour reports of govt communications with companies
- **Non-govt entity is funded and supported by the industry**
 - Social media companies “Opt In” if they want Section 230 legal benefits
 - Realistic financial enforcement (including penalties) accepted by companies



In Section 230: “Imminently Harmful” Rather Than “Objectionable”

- “Imminently Harmful” can be defined & adjudicated by arbitration judges
- The term “Otherwise Objectionable” is vague & not defined anywhere
- Original intent of Sec230 was to prevent porn & dangerous content



FINRA as a Model for a Non-Govt Regulatory Entity

- **FINRA handles investor/advisor disputes for a portion of the financial industry**
 - \$800M budget, 3000 employees, 8000 arbitration judges
 - FINRA has handed out fines to companies and suspended 100s of financial advisors/brokers/etc
 - Started from NASDAQ/NYSE with governance to avoid regulatory capture by the industry
- **Proposed OMRA (Online Media Regulatory Authority) will operate like FINRA**
 - As FINRA is loosely overseen by SEC, OMRA would be loosely overseen by the FCC or FTC
 - Power to levy fines large enough to motivate the companies
- **FINRA outsources investor appeals, OMRA outsources online user appeals**
 - Companies can build user trust and off-load the costs of PR controversies and legal issues
- **OMRA governance carefully designed to ensure all stake-holders fairly represented**
 - Social media platforms (large & small), content creators, news publishers, right and left, user representatives
 - Balanced governance: No regulatory capture, no political appointees, and no govt funding



OMRA Will Reduce Litigation Costs

- **User EULAs modified to drive users to online arbitration with OMRA**
 - Users will have simple online mechanism to appeal their disputes to OMRA
- **OMRA structure will include two-tiers of appeals**
 - Accepted user disputes are resolved via online video arbitration with 1 arbitration judge
 - Accepted 2nd tier appeal is escalated to an online video arbitration with a panel of judges
- **OMRA precedents will be organized and searchable online**
 - User dispute categories, resolutions, opinions, enforcement actions, penalties if any
 - Platforms can point users to similar disputes previously resolved by OMRA
- **Clear published content moderation rules will reduce user disputes**
 - When users understand the specific rules violated, they're less likely to dispute the enforcement



Low Cost of an OMRA Model of Regulation

- **OMRA can be industry-funded on similar scale as FINRA**
 - Top 10 companies are profitable with annual revenues exceeding \$500B
 - Meta already committed \$130M for its hand-picked Oversight Board
- **Minimal reporting requirements and regulatory processes**
 - Data reports on content moderation categories, enforcement actions, fact-checkers, user disputes, etc
 - Reporting requirements limited to larger online media companies
- **Cost efficient use of online video rather than F2F meetings for arbitrations**
 - FINRA successfully uses online video meetings for thousands of investor-vs-industry arbitrations
- **Some business/legal costs for dispute appeals can be outsourced to OMRA**



Dealing with Massive Scale

- **“Certified” online platforms will have fewer & simpler appeals**
 - Accepted appeals would adjudicate based on whether platform followed its own published rules
 - Transparency and clear enforcement explanations to users will reduce appeals
- **Companies will handle initial content disputes - OMRA will receive only the appeals**
 - OMRA will rely on its published precedents – and may or may-not accept each appeal based on merits
 - Refundable nominal fee (e.g. \$100)to discourage frivolous appeals by users
- **OMRA penalties incentivize platforms to improve algorithms & transparency**
 - Improve content moderation algorithms and transparency, neutrality rating of fact-checkers
 - Platform moderation algorithms & published standards will provide the heavy-lifting to handle scale
- **OMRA can be industry-funded at similar scale as FINRA (with far fewer employees)**
 - Cost efficient use of online video rather than F2F meetings for arbitrations



Monopoly Service Providers Cannot Pick Their Users

- **Key is equal treatment of all users (those who follow the rules)**
 - Common Carrier definition: monopoly public service provider (train, toll road, utility, telecom)
 - Monopoly platforms can avoid Common Carrier legal status if they are legislatively mandated to treat all users equally and can no longer pick & choose their users/customers
- **Vast majority of users click Terms of Service without reading them**
- **Private Rights of monopoly service providers are superseded by public speech rights of users and content publishers/communicators**
- **Largest online platforms qualify as monopolies**
 - Google 90% monopoly of search and online video (YouTube)
 - Facebook 85% of social media (Facebook, Instagram, WhatsApp)
 - Amazon 80% of retail books
 - Apple & Google 90% of online mobile app store



Govt Cannot Legally Influence Content Moderation

- **Monopoly platforms cannot censor on behalf of government**
 - Govt officials can provide info/expertise, but cannot serve as “fact checkers”
 - Govt efforts to determine true/false is a constitutional threat per US Supreme Court
- **Government suffers from partisan political incentives toward bias**
 - Govt agency leaders are appointed and influenced by partisan politicians
- **Government content standards is a recipe for Authoritarian States**
 - Who Decides ? Do standards change upon 4 year election cycles ?
 - Supreme Court: Govt efforts to decide true/false are more harmful than falsehoods



Guard Rails & OMRA Also Apply to Online Payments

- **Online payment & site monetization are key components of free speech**
- **Examples of near-monopoly online payment platforms covered include:**
 - PayPal and Venmo
 - Credit Card brands/networks
 - GoFundMe and other donation sites
 - Ad-based hosting platforms
 - Mobile app payment gateways

