

Trusted Online Content Moderation

Mike Matthys mikem@4betterinternet.com

Trusted Online Content Moderation

- Search & social media platforms are not widely trusted for safety & neutrality
 - Both political parties doubt efforts for safe & neutral content moderation; Surveys show users lack trust in platforms
 - Opaque enforcement actions with unclear explanations to affected users will continue to drive user distrust in platform
- Proposed Solution: Mandate transparency on all aspects of content moderation
 - Report on each enforcement action: user, specific rule broken, specific content category, use of named fact-checker, govt request
 - Published enforcement action reports can be sorted, tabulated, analyzed to measure online safety and viewpoint neutrality
 - Any communications with govt or entities/contractors funded by govt (except specific law enforcement /natl security actions)
- Glare of publicity will help ensure both online safety and viewpoint neutrality
- Congress action needed to ensure solutions survive future White House admins



Why Focus on Transparency?

- Published platform comparisons will improve content moderation performance
 - 3rd party reports will enable peer-to-peer comparison of platforms like reports to SEC on financial performance
 - Observers, critics, media and academics can study and measure content moderation performance & results
- Ensure content moderation requests by govt & govt-funded entities are appropriate / visible
 - · Specific national security and law enforcement actions would be excepted
- Expanded transparency is workable for companies to implement
 - Large companies already gather, report and publish subsets of relevant generalized information and statistics
 - See downloadable CSV files at www.transparencyreport.google.com or www.transparencyreport.google.com
- Broad support for transparency; limited support for Section 230 changes in divided Congress
 - Sec230 change requires difficult agreement on standards for online safety, viewpoint neutrality, role of govt, liabilities



Ensuring Full Transparency

- Publish reports on all enforcement actions taken
 - Specific content categories affected, type of action (including demonetize or de-amplify), content rules broken, repeat offenders
 - Disputes and appeals of dispute resolutions, role of specific 3rd-party fact checkers in enforcement actions
- For any enforcement action, provide clear explanation & justification to user
 - Specific content and specific rules broken, explanation of review/adjudication process, steps to appeal
- Publish criteria and ranking of websites as authoritative sources for search results
- Report relevant company communications to/from govt & govt-funded entities
 - Whistleblowers protected and rewarded if companies are purposely evading public transparency
 - Actual national security and law enforcement actions are excepted
- Fact-checkers
 - Identify and publish affiliations, funding & online history of fact-checkers



Transparency ≠ Divulging Tech Trade Secrets & IP

- Majority of enterprise value is their proprietary algorithms and software
 - Output/results are regulated in other industries rather than "how it works" trade secrets
 - Academic researchers cannot be prevented from eventually joining/assisting competitor companies
- Monitoring and regulating content moderation does not require knowledge of proprietary algorithms and software
 - Online safety and consistency in enforcement actions can be tested and monitored
- Echo chambers & features of their proprietary algorithms become less important if enforcement of online safety is transparent and can be monitored



IBI's Proposal:

Mandate transparency, with financial penalties to ensure accountability

- Require all platforms to publish mandated transparency data and reports quarterly
 - All content rules, user dispute resolutions, whistleblower protections, fact-checkers
 - All enforcement actions including specific content category that triggered action and specific type of enforcement action
 - All website ranking categories for use as authoritative sources for search results
 - All communications with govt or govt-funded entities (except natl security & law enforcement actions)
- For each enforcement action, provide clear explanation/justification to user
 - Specific content and specific rule(s) broken, explanation of appeals process
- Define criminal/financial penalties to ensure accountability for transparency
 - Federal criminal actions can be brought by federal DOJ or FTC
 - Audits/discovery triggered by civil or regulatory enforcement actions
 - Protection and reward for whistleblowers (employees or contractors of the online companies)
 - Financial penalties TBD based on intent, repeat of offenses, and harm caused by missing or incorrect enforcement actions
- A transparency-only solution for both online safety and viewpoint neutrality



Other Industries Provide Safety & Performance Statistics

- Automobiles are measured for safety, mileage, performance, resale value, etc.
- Airlines are measured for on-time arrivals
- Banks are measured for innumerable financial and fairness measures
- Heavy industry measured for employee safety and environmental protections

When statistics are measured and reports are published – safety and performance improves

