The Case for Designing Tech for Social Cohesion: The Limits of Content Moderation and Tech Regulation

By Lisa Schirch*

Apple CEO Steve Jobs described computers as "bicycles for the mind" that amplify human energy. But the metaphor of a bicycle suggests the internet is a road we can travel in any direction. In reality, digital platforms restrict what we can and cannot do.

Recognizing the power of computer infrastructure on human behavior, Stanford Psychology professor BJ Fogg taught a generation of Silicon Valley innovators how to design tech products to harness psychological insights in his course based on his book *Persuasive Technology*. Digital products are persuasive technologies; they engineer how humans communicate. The design of tech products may amplify some human behaviors, thoughts, and relationships and distort, obscure, or downgrade others. Small changes to algorithms and user interfaces on social media products can influence what people buy, whether they vote, who they vote for, etc.

Technology products also reflect the biases and perspectives of those designing affordances and algorithms.² Computer engineers embed their values into the affordances and algorithms that govern human interaction online. Harvard law professor Lawrence Lessig's book *Code and other Laws of Cyberspace* described the internet as a socio-technical institution; code is law. The architecture of technology products enables what people can and cannot do. Lessig warned that the digital architecture of the web could enable freedom and privacy, or the contrary; it could enable business and government to surveil and control.³ Lessig's point applies to polarization and social cohesion. Digital infrastructure can amplify polarization through the code. *And* digital infrastructure can persuade

^{*} Richard G. Stamann, Sr. Professor of the Practice of Peace Studies, University of Notre Dame.

^{1.} B.J. Fogg, *Persuasive Technology: Using Computers to Change What We Think and Do.* (Amsterdam: Morgan Kaufmann Publishers, 2003).

^{2.} See Ruha Benjamin, Race After Technology: Abolitionist Tools for the New Jim Code. Cambridge, Medford, MA: Polity, 2019; Cathy O'Neil. Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy. (New York: The Crown Publishing Group, 2016).

^{3.} Lawrence Lessig, Code and Other Laws of Cyberspace (New York: Basic Books, 1999).

people to build social cohesion.

Since the beginning of Silicon Valley's tech industry, there has been a thriving "tech for good" movement. Yet some of the tech products, particularly social media platforms, have become superhighways for disinformation, hate speech, and other forms of harmful content. The design of digital products shapes the direction we can pedal. Digital technology affordances and algorithms can amplify hate and disinformation online that spill over into real-world violence. Digital tools can also help us to build bridges online to improve social cohesion.

Tech platforms have engineered a new digital public sphere at a time when toxic polarization was already increasing globally.⁴ While not the origin of social and political division, there is wide agreement that harmful content on social media amplifies polarization. Toxic polarization refers to harmful levels of distrust and dysfunction in divided societies.

Divisive digital content influences the way political actors, traditional media, and the public frame issues even for people who do not use social media. The challenge of harmful content online is increasing. Political actors, cyber armies, and a growing for-profit disinformation industry amplify and incentivize individual producers of divisive digital propaganda aimed at polarizing societies with a "divide and conquer" strategy.

Since around 2017, some tech companies have built a Trust and Safety infrastructure with thousands of staff overseeing a global content moderation effort to remove, demote or disincentivize harmful content. But in 2022 the tech sector's relatively new "Trust and Safety" infrastructure laid off 120,000 tech workers and downsized human rights and content moderation teams due to reduced tech company stock prices, Elon Musk's Twitter acquisition, and other global factors. The politicization of content moderation is increasing. In the US, conservatives tend to criticize content moderation as censorship while liberals tend to view content moderation as a matter of life or death for threatened minority groups and democratic institutions targeted by online harmful content. Some repressive governments employ content moderation—sometimes shutting off internet access altogether—to subvert human rights and democracy activists.

This paper draws on nearly 60 interviews with staff at tech companies, critics of big tech, civil society groups impacted by tech-amplified social media, and new tech startups designing platforms to reduce polarization and support social cohesion.⁶ The interviews took place between 2021 and 2022,

^{4.} Thomas Carothers and Andrew O'Donohue, *Democracies Divided: The Global Challenge of Political Polarization* (Washington DC: The Brookings Institution. 2019).

^{5.} See Laurel Wamsley, "It's the End of the Boom Times in Tech, as Layoffs Keep Mounting," NPR, November 16, 2022.

^{6.} This research was commissioned by a working group formed to launch a Council on Tech and Social Cohesion, including the Center for Humane Technology, Search for Common Ground, the Toda Peace Institute, Braver Angels, More in Common, the University of Notre Dame, and the Alliance for

primarily in the US and Europe. The research aimed to identify and understand how tech companies were responding to harmful content online. Interviews revealed three distinct but complementary narratives or approaches to thinking about polarization and social cohesion in digital spaces.

The "User-Centered" Narrative describes harmful content online as generated by users, with social media products and search engines acting as a mirror of society. Several interviewees described the defeating feeling of playing "whack-a-mole" against the growing tide of individual and state-sponsored harmful digital content. This narrative points to the need for content moderation on user-generated content and increased digital media literacy to help the public navigate information and communication on the internet.

The "Tech Design Regulation" Narrative describes harmful content as amplified by tech product designs including the affordances and algorithms that are optimized for user engagement, advertising, and shareholder profit. Many social media companies optimize their product designs for user engagement to maximize their ad-based profits. Machine learning algorithms promote emotionally alarming, divisive, and attention-grabbing content, just as cars slow down driving past a car accident and news outlets use the "if it bleeds, it leads" principle. From this point of view, some tech products incentivize harmful content that drives toxic polarization. This narrative presses for government regulation to extend beyond privacy to regulating tech profit models, algorithms, affordances, and designs that amplify toxic content.

The "Social Cohesion by Design" Narrative describes tech products that amplify and scale social cohesion by designing affordances and algorithms optimized for these purposes. "Peacetech" engineers with training and expertise in social cohesion can design products that contribute to social cohesion. These digital products can support human agency to participate in civic action, bridge divided communities, and build trust between the public and institutions.

The first half of this article explores the complex relationship between toxic polarization and digital spaces and uses these three frames to help understand the role of digital spaces in toxic polarization. The second half of the paper focuses on examples and case studies of "social cohesion by

Peacebuilding. David Jay from the Center for Humane Technology and Althea Middleton-Detzner provided a list of and introduction to tech company staff that could be interviewed for this report. Funding to support the research came from two main sources. Search for Common Ground secured funding from KBF Canada to hire Althea Middleton-Detzner and the Toda Peace Institute supported research by Lisa Schirch, based at the University of Notre Dame. The two conducted most of the interviews together. Schirch wrote this report receiving important feedback from colleagues and interviewees.

design." The paper concludes with a call for governments and tech companies to move beyond content moderation to invest in technologies that will improve societies' ability to solve problems and prevent violence. The paper argues that governments can incentivize social cohesion by design. As a complement to content moderation and government regulation, designing tech to support social cohesion should be a primary strategy for addressing the crisis of toxic polarization.

I. UNDERSTANDING POLARIZATION AND SOCIAL COHESION

Polarization occurs when diverse identity groups in a society are divided along an axis into two sides.⁷ In general, polarization or differences of belief are not necessarily harmful and can be opportunities for positive social change. Polarization over the ethics of slavery, colonialism, and women's rights, for example, led to civil rights movements, policy proposals to improve equality, and social change.

In technical terms, there are different types of polarization. *Issue* polarization describes a normal situation where different groups hold different views but can listen to each other and solve problems that arise through democratic processes because of a shared sense of human dignity and trust. Issue polarization can be managed when there is social cohesion. A society with social cohesion addresses conflict or issue polarization as an opportunity for improving society. Conflict is a normal and important aspect of human relations that signals that there are issues needing attention. Groups of people with different experiences and interests often experience conflict.

Toxic polarization, also known as *affective polarization*, occurs when groups distrust and/or dehumanize others with us-vs-them narratives that view violence as necessary and justifiable against what they perceive as an existential threat.⁸ *Political polarization* refers to a society where political party affiliation becomes a defining element of identity, overshadowing how an individual may feel about an issue.⁹ In the United States, polarization is not just spilling over from elite political polarization; a growing number of people at the community level hold contempt for people of other political parties.¹⁰

A growing body of evidence suggests that political polarization exaggerates the actual policy differences between groups.¹¹ In other words,

^{7.} Shanto Iyengar et al., "The Origins and Consequences of Affect Polarization in the United States," *Annual Review of Political Science* 22 (2019): 129–146.

^{8.} Ibid.

^{9.} Peter T. Coleman, *The Way Out: How to Overcome Toxic Polarization* (New York: Columbia University Press, 2021).

^{10.} Daniel DellaPosta, "Pluralistic Collapse: The "Oil Spill" Model of Mass Opinion Polarization," *American Sociological Review*, 85, no. 3 (2020): 507–536.

^{11.} See Chris Bail, Break the Social Media Prism: How to Make Our Platforms Less Polarizing,

there is a perception gap. People think they disagree more than they actually do.¹²

Affective and political polarization can be toxic to society. Toxic polarization can reduce a society's ability to interact with each other and respond to complex problems like the climate crisis or the pandemic. As public mistrust of other social groups and public institutions decreases, so does an individual's belief that change is possible and that civic engagement is an effective route to change. ¹³ Societies with low levels of social cohesion have a weaker ability to solve problems together and have an increased likelihood of intergroup violence. ¹⁴

Social cohesion is the opposite of *toxic* polarization, as illustrated in Figure 1. The United Nations defines social cohesion as "the extent of trust in government and within society and the willingness to participate collectively toward a shared vision of sustainable peace and common development goals." Social cohesion is the glue that keeps a society together.

Toxic Polarization		Social Cohesion
Individuals feel isolation, humiliation, and frustration and behave as though they are not able to participate in decisions that affect them	Individual agency	Individuals feel a sense of safety and dignity, and behave with the skills and capacity to <i>influence</i> and <i>participate</i> in decisions that affect them
Individuals feel a sense of <i>exclusion</i> and behave with contempt and distrust toward other people	Horizontal cohesion within and between groups	Individuals feel a sense of belonging and inclusion and behave with empathy and trust toward other people

⁽Princeton, NJ: Princeton University Press, 2021).

^{12.} Daniel Yudkin, Stephen Hawkins and Tim Dixon, "The Perception Gap: How False Impressions are Pulling Americans Apart," *PsyArXiv* (September 14, 2019).

^{13.} Ethan Zuckerman, Mistrust: Why Losing Faith in Institutions Provides the Tools to Transform Them (New York: W.W. Norton and Co, 2021), 20.

^{14.} I. Olawole, A. Lichtenheld and R. Sheely, "Strengthening Social Cohesion for Violence Prevention: 10 Lessons for Policymakers and Practitioners," Mercy Corps (2022); A. Lichtenheld et al., "Understanding the Links Between Social Cohesion and Violence: Evidence from Niger," Mercy Corps. (2021).

^{15.} United Nations Development Program (UNDP), Strengthening Social Cohesion: Conceptual Framing and Programming Implications. (New York: UNDP, 2020).

People in society feel a sense of exclusion, contempt, and distrust toward leaders and institutions which are seen as corrupt and captured by elite interests

Vertical cohesion between institutions and the public People in society feel a sense of *inclusion, investment, and trust* in leaders and institutions which are seen as *transparent and accountable* to the public

Figure 1: The Polarization-Cohesion Spectrum

Social cohesion is both a *goal* and an *approach*. The UN uses the term social cohesion to describe *the goal* of its efforts in peacebuilding, dialogue, participatory governance, prevention of violent extremism, and bridgebuilding interventions. UN Peacebuilding initiatives¹⁶ have grown out of local peacebuilding¹⁷ and bridge-building efforts¹⁸ to coordinate diverse stakeholders and activities in support of social cohesion.

For more than four decades, the field of peacebuilding has been researching, experimenting, and practicing the science and art of facilitating dialogue, negotiation, and mediation to depolarize divided societies and address the root causes of conflict. ¹⁹ International organizations like the UN and World Bank invest large sums in peacebuilding, as they recognize its value in preventing violence, which negatively affects people, business interests, and the planet. The field of peacebuilding is an umbrella term that includes the concepts of conflict resolution, conflict management, and conflict transformation. Within the US, there are a wide range of movements and organizations whose work can be categorized as peacebuilding, including for example groups that aim to protect democracy, address social justice, or build bridges between groups.

The OECD uses the term social cohesion to *describe* a society that "works towards the well-being of all its members, fights exclusion and marginalization, creates a sense of belonging, promotes trust, and offers its members the opportunity of upward social mobility." The OECD defined social cohesion as characteristic of a society that values "the well-being of all its members, fights exclusion and marginalization, creates a sense of belonging, promotes trust, and offers its members the opportunity of upward social mobility."

Based on the work of Search for Common Ground, there are three

^{16.} United Nations, "Peacebuilding," https://www.un.org/peacebuilding.

^{17. &}quot;Peacebuilding," Wikipedia, https://en.wikipedia.org/wiki/Peacebuilding.

^{18. &}quot;Members," Bridge Alliance, https://www.bridgealliance.us/membercategories.

^{19.} Fletcher D. Cox and Timothy Sisk, *Peacebuilding in Divided Societies: Toward Social Cohesion* (Cham, Switzerland: Palgrave MacMillan, 2017).

^{20.} OECD, Perspectives on Global Development 2012: Social Cohesion in a Shifting World, (Paris: OECD Publishing, 2011).

^{21.} Mike Colledge and Chris Martyn, "Social Cohesion in the Pandemic Age," IPSOS (October 2020).

elements related to social cohesion.²²

- 1. **Individual Agency** exists when individuals feel a sense of safety, dignity, and capacity (skill) to *influence* and *participate* in decisions that affect their lives within society and with governing institutions. Individual agency requires an ability to communicate about difficult issues in a "healthy" way with communication skills that focus on problem-solving while recognizing the dignity of oneself and others.
- 2. **Horizontal Cohesion** exists when individuals feel a sense of *positive relationships, belonging, and trust* within and between identity groups based on politics, religion, ethnicity, class, education, region, or other shared identities. Horizontal cohesion requires skills for healthy expression of conflict and solving problems through inclusive, collaborative, non-violent processes in both bonding and bridging networks. It also includes efforts to improve horizontal cohesion through dialogue and research, building trust through working together in areas where there is common ground, and reality checking, as often people misperceive the intentions and beliefs of others. Horizontal cohesion is also called "horizontal social capital."

*Intra*communal cohesion, also known as "bonding social capital," refers to the quality of relationships within an identity group (e.g., relationships among Black Americans).

*Inter*communal cohesion, also known as "bridging social capital" refers to the quality of relationships between identity groups (e.g., between Black and white Americans).

3. **Vertical Cohesion** exists when individuals and groups in society feel a sense of *trust, transparency, accountability, and collaboration* with public institutions including government, as well as news media, academic institutions, and corporations. This is also called "vertical social capital." In an active democracy, citizens engage with governments. Civic engagement is an expression of vertical cohesion paired with individual agency. Vertical cohesion exists when public institutions recognize basic human rights and serve community members equitably. Public goods such as equal treatment under the law, safety, healthcare, and education are afforded to all.

Social cohesion enables a society to function in a way that addresses the needs of all members and to be resilient to shocks, stressors, and crises such as a pandemic or natural disaster. Social cohesion enables societies to work together to solve problems.²³ In the panoply of catastrophes facing humanity

^{22.} This schema synthesizes similar frameworks for social cohesion, and draws specifically from this report: Institutional Learning Team, "Building Social Cohesion in the Midst of Conflict: Identifying Challenges, Measuring Progress, and Maximizing Results," Search for Common Ground. (November 2020)

^{23. &}quot;Social Cohesion and the State: What Can the G20 Do to Improve Social Cohesion and Trigger Responsibility in Business and Politics?" Global Solutions: The World Policy Forum (2022), https://www.global-solutions-initiative.org/global-table/social-cohesion-through-business-and-

today, social cohesion enables societies to work together to solve problems including climate change, poverty, inequality, racism, and violence.²⁴ Cohesive societies are more likely to reduce income and unemployment disparity, are more likely to address problems collectively, and inspire a sense of belonging in people.²⁵

Societies with low levels of social cohesion have a weaker ability to solve problems together and have an increased likelihood of intergroup violence. During the Covid pandemic, countries with low levels of social cohesion suffered more deaths from Covid. A lack of social cohesion can mean that people did not feel a sense of agency to work for change, did not trust their neighbors to wear a mask or get a vaccine, and/or did not trust their government to give them accurate information about the pandemic and vaccine. Countries with higher levels of social cohesion had fewer deaths. Similarly, countries with high levels of social cohesion can make climate policies more acceptable to citizens.

A society with social cohesion approaches conflict as an opportunity to improve society. Conflict is a normal and important aspect of human relations that signals that there are issues needing attention. Groups of people with different experiences and interests often experience conflict. The goal of social cohesion is not to suppress conflict or to reduce differences between groups. Authoritarian governments tend to view conflict itself, such as citizens voicing a critique of government policy, as dangerous. The goal of social cohesion is to provide democratic processes and spaces for public deliberation and creative problem-solving to address conflicts between groups.

politics/.

24. Ibid.

^{25.} Danielle Baussan, "Social Cohesion: The Secret Weapon in the Fight for Equitable Climate Resilience," The Center for American Progress, (May 11, 2015) https://www.americanprogress.org/article/social-cohesion-the-secret-weapon-in-the-fight-for-equitable-climate-resilience/.

^{26.} Olawole, "Strengthening Social Cohesion" (2022); Lichtenheld et al., "Understanding the Links Between Social Cohesion and Violence: Evidence from Niger" (2021).

^{27.} Loring J. Thomas et al. "Geographical Patterns of Social Cohesion Drive Disparities in Early COVID Infection Hazard," *Proceedings of the National Academy of Sciences in the United States of America* (March 14, 2022).

^{28.} Adam Taylor, "Researchers Are Asking Why Some Countries Were Better Prepared for Covid. One Surprising Answer: Trust," *Washington Post*, February 1, 2022.

^{29.} Daniele Malerba, "The Effects of Social Protection and Social Cohesion on the Acceptability of Climate Change Mitigation Policies: What Do We (Not) Know in the Context of Low- and Middle-Income Countries?" *The European Journal of Development Research* 34 (May 6, 2022): 1358–1382.

Definitions

Toxic polarization occurs when people perceive other people as existential threats, distrust and dehumanize others with us-vs-them narratives and justify the use of violence against others. *Toxic polarization* includes three dimensions:

- Individual isolation and a loss of human agency to participate in civic life
- Divisions between groups into narratives of "us vs them" with emotional contempt for the "other"
- Lack of trust between the public and institutions in government and public-interest media

Social cohesion refers to the glue that keeps society together; it is the opposite of toxic polarization. Three dimensions of social cohesion include:

- Individual agency to participate in civic life
- Horizontal relationships within and between social groups
- Vertical relationships between public institutions and society

Bridge building and peacebuilding are types of **prosocial interventions** that support the goal of social cohesion in three ways. 1) Increasing individual agency to participate in civic life; 2) Bridging relationships between groups; and 3) Building public trust between society and governing institutions.

Technology or tech refers in this article to digital tools, with a particular focus on social media. **Affordances** are the features of a tech product that shape behaviors. The Like, Share, and Comment features of most social media products are examples of affordances. **Algorithms** are the computational settings of a tech product that determine what content users can see.

PeaceTech refers to technology that both supports the analysis of polarization and bridge building or peacebuilding interventions to support social cohesion.

II. SOCIAL MEDIA, HARMFUL CONTENT & POLARIZATION

Interviews with tech staff for this paper found that most reported high levels of concern about tech related harms such as polarization and noted that staff, in general, want to feel good about the company that employs them. Tech staff shared that there is a wide appetite for achieving company missions to "connect" people and build relationships. Yet even staff at companies who have hired tens of thousands of content moderators describe

an endless game of "whack-a-mole" to manage a "tsunami of harmful content" without adequate resources, particularly in the Global South where they lack staff who speak local languages.

Many of the tech insiders interviewed for this report questioned the link between technology and social cohesion. As noted earlier, research suggests polarization was increasing globally before the advent of digital technology.³⁰ There is_a robust literature on the impact of social media products on polarization.³¹ Research both supports and questions the link between technology and polarization.³² Some studies have found that polarization is growing more among groups with less internet usage.³³ But research surveys consistently find that social media platforms impact social cohesion by altering social networks and fragmenting public conversations on issues, rapidly spreading false information and the dysfunction of digital governance and norms.³⁴

A survey study by New York University's Stern Center for Business and Human Rights asserts that while big tech companies like Meta, Twitter, and Google were not the source or largest factor in rising U.S. political polarization, these products amplified "divisiveness" and its "corrosive consequences." According to the Pew Research Center, 64% of Americans believe social media is negatively affecting the US, and express concern about the misinformation and the hate and harassment they see on social media.³⁶

^{30.} See Carothers and O'Donohue, Democracies Divided.

^{31.} See for example, Paul M. Barrett, Justin Hendrix, and J. Grant Sims, "Fueling the Fire: How Social Media Intensifies Polarization," New York University Stern Center for Business and Human Rights (September 2021); Philipp Lorenz-Spreen, Lisa Oswald, Stephan Lewandowsky, Ralph Hertwig. "Digital Media and Democracy: A Systematic Review of Causal and Correlational Evidence Worldwide," SocArXiv 22 (Nov. 2021); Jay J. Van Bavel, Steve Rathje, Elizabeth Harris, Claire Robertson, and Anni Sternisko, "How Social Media Shapes Polarization," Trends in Cognitive Sciences 25, no. 11, (2021): 913–916; Almog Simchona, William J. Bradyc and Jay J. Van Bavel, "Troll and Divide: The Language of Online Polarization," PNAS Nexus 1, no. 1 (2022).

^{32.} See for example, Jonathan Stray, "Designing Recommender Systems to Depolarize," *First Monday* 27, no. 5 (May 2, 2022); Gideon Lewis-Kraus. "How Harmful Is Social Media?" *New Yorker*, June 3, 2022; Lydia Laurenson, "Polarisation and Peacebuilding Strategy on Digital Media Platforms," *Tokyo: Toda Peace Institute*, 2019, https://toda.org/assets/files/resources/policy-briefs/t-pb-44-laurenson-lydia-part-1 polarisation-and-peacebuilding-strategy.pdf.

^{33.} Levi Boxell, Matthew Gentzkow & Jesse M. Shapiro, "Greater Internet Use is Not Associated with Faster Growth in Political Polarization among U.S. Demographic Groups," *Proceedings of the National Academy of Sciences in the United States of America*, September 19, 2017.

^{34.} Sandra González-Bailón and Yphtach Lelkes, "Do Social Media Undermine Social Cohesion? A Critical Review," *Social Issues and Policy Review* 17 (2022) 1–26.

^{35.} Barrett, Hendrix, and Sims, "Fueling the Fire."

^{36.} Brooke Auxier, "64% of Americans Say Social Media have a Mostly Negative Effect on the Way Things are Going in the U.S. Today." *Pew Research Center* (October 15, 2020).

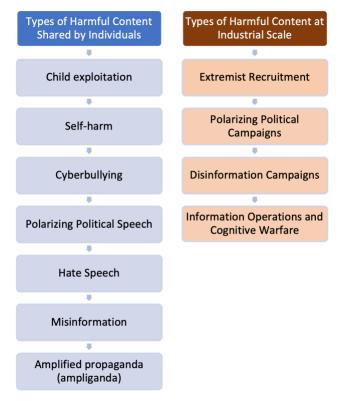


Figure 2: Typology of Harmful Content

Political actors are exploiting social media and search engines to spread false propaganda to divide citizens, aggravate existing social divisions, foment violence, and sway elections. False and deceptive information both online and offline synergize with hateful content, violent extremism, and repressive states to pit "us vs. them." Harmful content online contributes to "toxic polarization."

The problem of harmful content on these tech products started small. Early tech products like eBay and Flickr wrestled with "individual rule breakers" posting spam, fraud, and nudity. But an avalanche of other problems soon followed. The internet became a superhighway for child sexual abuse and exploitation. Some social media products became boxing arenas for verbal jousts and hateful commentary by average people. Individuals spreading harmful content and inadvertent rule breakers soon were joined by industrial-scale producers of harmful content. Figure 2 provides a typology of individual and industrial-scale harmful content.

Political actors from ISIS to Russia weaponize these affordances to operate mass influence operations. State-based cyber troops use

"computational propaganda" to wage "cognitive warfare³⁷ on both domestic and foreign publics. Tech savvy authoritarians maximize algorithmic rewards for outrage and division. By 2020, the University of Oxford Programme on Democracy and Technology warned of "industrialized disinformation" by over 80 countries with cyber armies spreading computational propaganda.³⁸ Cyber troops and a booming for-profit disinformation industry generate content to undermine public trust in democratic institutions and elections, discredit human rights activists, and widen preexisting divisions in society. Social media affordances enable ordinary people to amplify divisive propaganda by sharing false, deceptive, or polarizing information campaigns, also known as *ampliganda*.³⁹

Researchers around the world report social media playing a key role in further polarizing already divided societies, undermining public trust in democratic institutions, and increasing public support for autocrats. ⁴⁰ The impact of industrialized disinformation campaigns is what some call "the liar's dividend" or "epistemic insecurity" where the public senses chaos, feels confused, and views everything as questionable resulting in the collapse of truth. It is not uncommon to hear people refer to the weaponization of social media or refer to some tech products as "weapons of mass distraction" and "mass destruction."

Just like a small amount of toxins can pollute a river or lake, even a small amount of harmful content online can create toxic *information ecosystems* that enable autocratic political actors to undermine social cohesion and democracy. The Center for Humane Technology describes "polarization spills" on social media as unique. Unlike toxic oil spills, a polarization spill not only causes harm in dividing society. It also makes it difficult to govern. While an oil spill does not in itself make it more difficult to find regulatory solutions to prevent more oil spills, toxic polarization spills do make it more difficult for political actors to find regulatory solutions to digital amplification of polarization.⁴³

^{37.} Samuel C. Woolley and Philip N. Howard, *Computational Propaganda: Political Parties, Politicians, and Political Manipulation on Social Media* (New York: Oxford University Press, 2018). See also François du Cluzel, *Cognitive Warfare*, (Brussels: NATO Innovation Hub, 2020).

^{38.} Samantha Bradshaw, Hannah Bailey and Philip N. Howard, "Industrialized Disinformation: 2020 Global Inventory of Organized Social Media Manipulation," *Oxford, UK: Programme on Democracy & Technology* (2021), https://demtech.oii.ox.ac.uk/.

^{39.} Renée DiResta, "It's Not Misinformation. It's Amplified Propaganda," *The Atlantic*, October 9, 2021.

^{40.} Lisa Schirch, Ed., Social Media Impacts on Conflict and Democracy: The Techtonic Shift (Sydney: Routledge, 2021).

^{41.} Christina Nemr and William Gangware." Weapons of Mass Distraction: Foreign State-Sponsored Disinformation in the Digital Age" (Washington DC: Park Advisors, 2021).

^{42.} Sarah Jacobs Gamberini, "Social Media Weaponization: The Biohazard of Russian Disinformation Campaigns," *Joint Forces Quarterly* 99 (2020), https://wmdcenter.ndu.edu/Publications/Publication-View/Article/2422660/social-media-weaponization-the-biohazard-of-russian-disinformation-campaigns/.

^{43.} Center for Humane Technology, "Addressing the TikTok Threat," Your Undivided Attention

III. 3 APPROACHES TO REDUCING POLARIZING CONTENT ON DIGITAL SPACES

The research for this paper revealed three distinct but complementary approaches to thinking about polarization and social cohesion in digital spaces. The first approach blames users for generating harmful content. In this view tech products are neutral mirrors of society. Content moderation focuses on removing harmful user-generated content. The second approach blames tech companies for designing affordances and algorithms that amplify toxic content. This approach advocates government regulation of tech algorithms. The third approach focuses on designing new digital spaces with affordances and algorithms designed to support social cohesion. Table 1 below synthesizes these three approaches.

Table 1: Narratives on Harmful Digital Content

	Perception of the Challenge	Interventions
User- Centered Narrative	Tech insiders often frame the problem as user-generated content. In this view, technology is just a "mirror" of society.	Tech companies have built a "Trust and Safety" infrastructure to address how people use the internet to cause harm. The bulk of Trust and Safety initiatives focus on moderating content by developing data classifiers and using human moderators paired with machine learning and AI to remove harmful content. Tech insiders often refer to this as a "whack-a-mole" effort that cannot keep up with the scale of user-generated harmful content.
Tech Regulation Narrative	Tech critics often frame the problem as harmful tech products with profit models that incentivize affordances and algorithms that distort and amplify the worst aspects of human behavior.	Tech critics identify the need for government regulation of tech products not only in terms of data privacy and cybersecurity but also of tech profit models, product affordances, and algorithms.

ľ	PeaceTech Narrative	the problem as a lack of tech products that can scale social cohesion efforts to build individual agency, public trust, and bridge intergroup	moderation, tech regulation, and incentivizing prosocial tech product designs that amplify the best aspects of human behavior that improve social
		relationships.	

IV. THE USER-CENTERED NARRATIVE

Most of the tech company staff interviewed downplayed the responsibility of tech companies for harmful content or online polarization, asserting that technology is just a "mirror" reflecting what people already think. The logic of externalizing the problem of hateful content is part of a communication strategy for tech companies like Meta. For example, Facebook's Nick Clegg argued, "There is no editor dictating the frontpage headline millions will read on Facebook. Instead, there are billions of front pages, each personalized to our individual tastes and preferences, and each reflecting our unique network of friends, Pages, and Groups."⁴⁴ Some interviewees noted that journalists overstate the scale of toxic content. Facebook's Clegg is on record stating that the scale of harmful content online is relatively small, noting, "hate speech is viewed 7 or 8 times for every 10,000 views of content on Facebook."⁴⁵

Tech companies draw on a catalog of tech strategies to reduce harmful content. In response to widespread reports of escalating levels of digital toxicity, Silicon Valley's largest tech companies continue to invest in building a "Trust and Safety" infrastructure⁴⁶ that primarily uses content moderation to reduce digital harms that contribute to polarization. For example, Guy Rosen, Meta's VP of Integrity, rebutted critiques of Meta's role in toxic polarization with a list of Facebook's various strategies to reduce polarization.⁴⁷

Flooded with unsolicited advice from all corners of society, tech companies are open to ideas but ask for recommendations informed by what has already been tried. Tech insiders expressed frustration with outsiders offering ideas about how to fix tech without understanding the efforts

^{44.} Nick Clegg, "You and the Algorithm: It Takes Two to Tango," Medium, March 31, 2021.

^{45.} Ibid. Without full access to internal research, it is difficult to challenge these numbers. Yet there is wide skepticism that the problem is small given the wide perception of the vast scale of false, deceptive, and hateful content on social media. A meta-analysis of research on the scale of mis/disinformation on social media related to the COVID-19 pandemic found that up to one third of Covid-related content was false or deceptive.

^{46.} See, e.g., the Trust and Safety Professionals Association, https://www.tspa.org.

^{47.} Guy Rosen, "Investments to Fight Polarization." Meta, May 27, 2020.

already underway and the reductive nature of proffered solutions. They argue that some attempts to fix tech harms have reinforced the problem or created new ones. Interviewees noted that building classifiers to identify harmful content is complex and difficult; reducing tech harms goes well beyond simply adding a button or tweaking product designs. Reducing tech harms goes well beyond simply adding a button or tweaking product designs. There is no one "silver bullet" to reduce tech harms.

As of 2022, tech companies are taking a variety of steps to reduce digital harm. *Guidelines* strategies refer to how people can use the tech product. *User Interface* strategies determine how products present content. *Moderation* strategies determine what content is available. *Algorithm-based* strategies determine how tech products rank and recommend content to users. *Policies and partnership* strategies refer to the ways companies engage with outside groups and events, such as civil society or elections. *Company infrastructure* strategies refer to how tech companies organize their internal teams to prevent or respond to harm.

A. Incentives for Addressing Toxic Polarization on Tech Products

Tech companies have incentives and disincentives for responding to online polarization. Media reports and public pressure to remove harmful content are powerful incentives for tech companies to act. Yet significant challenges inhibit corporate action, including the complexity of the task and the scale and pace of toxic content.

Incentives include staff desire to achieve their tech company mission to "connect" people and grow the user base of people who want a safe place to communicate. Some identify a broader commitment to social responsibility to prevent harms. Several interviewees noted that a tech company that brands itself as strengthening community but then is charged with enabling genocide or undermining democracy has a serious problem. A tech company that faces widespread charges of harming society is failing its mission, which will make it more difficult to retain and attract good staff. Interviewees noted that people want to feel good about the company that employs them, that their efforts are contributing toward a positive corporate mission.

Within tech companies, interviewees noted that there is a "huge appetite" for achieving company missions that align with the public good, and great concern about tech-related harms. Some also noted that reports of tech harms have reduced the number of applicants applying to big tech companies, and drove a brain drain away from big tech as some staff left after not seeing enough effort or will to implement needed changes. Other interviewees noted that recent media reports from whistleblowers leaking internal documents have generated distrust, leading to more secrecy and restriction of information and data for researchers. Tech companies may

also respond to digital harms as a way of managing not only reputational risks from media attention, but also digital harms or public boycotts that might spur investors to withdraw support. Tech companies are also trying to prevent further government regulation or sanctions for harmful content.

Yet significant challenges inhibit corporate action. Many companies simply lack the staff necessary to manage the scale of industrial-scale disinformation and hate speech in the global town squares they have created. Escalating amounts of harmful content has created a sense of futility that moderation is a Sisyphean game of "whack-a-mole."

B. Analyzing Nuance at Scale

Content moderation is difficult, as machine learning algorithms need to be taught what is considered harmful. But classifying disinformation, hate speech, and other forms of harmful content requires analysis and debate on what views are protected as free speech.

A main challenge of moderation is to find a way to analyze nuance at scale. Facebook has over 3 billion users, creating an unimaginable amount of content requiring classification systems in dozens of different languages in contexts that change rapidly. Metaphors for hate speech may evolve quickly as companies censor one term, and users create new terms for the same hateful content. People rapidly innovate new ways of dehumanizing and demonizing others without using explicit hate speech, or even mentioning the group in question. In Myanmar, for example, people on social media praised the qualities of the Buddhist Burmese with the purposes of excluding and erasing Muslim groups.

C. The Politicization of Content Moderation

Tech companies face dilemmas to define the limits of free speech online, and the social norms for digital spaces.⁴⁸ On the left, human rights and democracy activists argue that tech companies do not moderate enough. On the right, conservative activists argue that tech companies violate free speech by removing posts deemed hateful, false, or deceptive. Content moderation, as a strategy for addressing harm, is a highly contentious process.

Tech company efforts to avoid partisan decisions on content moderation have proved unavoidable. Some tech staff assert they are committed to free speech, and thus minimize content moderation. Some use the term "social engineering" to the deliberate psychological manipulation of users through content moderation. Conservative critics of companies like Facebook and Google note that efforts to reduce harms are a form of social engineering.

^{48.} Valerie C. Brannon, "Free Speech and the Regulation of Social Media Content," *Congressional Research Service* (March 27, 2019).

For example, one content moderation program redirects user search queries for white supremacy content to organizations such as Life After Hate, founded and run by former white supremacists who are working to prevent the spread of white supremacy. Some groups view this as a form of censorship.⁴⁹

D. Profit Model Considerations

Several interviewees noted they were never in a room where anyone spoke about how a product or algorithm change aimed at reducing harm might reduce profits. Several insiders asserted they never directly observed tension between profits over safety or public goods like social cohesion. Many interviewees insisted that harmful content does not benefit the company's profit model and that harmful content is bad for business. As an example of this argument, Facebook's Nick Clegg stated in a recent article,

[It's] not in Facebook's interest—financially or reputationally—to continually turn up the temperature and push users towards ever more extreme content. The company's long-term growth will be best served if people continue to use its products for years to come. If it prioritized keeping you online an extra 10 or 20 minutes, but in doing so made you less likely to return in the future, it would be self-defeating. And bear in mind, the vast majority of Facebook's revenue comes from advertising. Advertisers don't want their brands and products displayed next to extreme or hateful content—a point that many made explicitly last summer during a high-profile boycott by a number of household-name brands.⁵⁰

Yet, other interviewees insisted the profit model of user engagement underlies all company decisions about designs and algorithms. Other interviewees noted that while profits might not be discussed during a crisis, the overarching push for growth, user engagement, and profits remain as a central framework for employees seeking to climb the ranks. Other interviewees noted the ad-based profit models are an unacknowledged obstacle to the bigger changes that might reduce harm and increase benefits. One interviewee noted that over the long term, some people are going to leave tech products that generate anger, recrimination, and conflict, and will gravitate towards tech products that create empathy, connection, belonging, dignity, and a sense of inclusion. One interviewee in a tech startup noted that "[i]f you build a system to give people justice, transparency, and a place where they feel heard, and they feel fairly treated, they will come back, and they will reward you with more money."

While tech company spokespeople like Clegg have challenged the claim

^{49.} Bronwyn Howell, "Consequences of the Christchurch Call: Social Engineering by Internet Platforms?" *American Enterprise Institute*, September 23, 2019.

^{50.} Nick Clegg, "You and the Algorithm: It Takes Two to Tango," *Medium*, March 31, 2021.

that tech company profit models incentivize polarizing content, other observers noted that the boycott Clegg references had little visible impact on Facebook. More than a thousand of the 9 million companies that advertise on Facebook joined the Stop Hate for Profit boycott of Facebook, including large advertisers. The boycott did result in a short-term decrease in company profits.⁵¹ While the boycott harmed Facebook's reputation, boycotts against social media companies have not met a threshold to cause shareholder harm to the company. To date, user boycotts and advertiser boycotts have had little impact on profits.

E. The Limits of Content Moderation

Tech companies are investing far more in efforts to reduce digital harm rather than promote prosocial content. By the end of 2022, an increasing number of tech insiders and analysts expressed dismay at the limits of content moderation.⁵² Moderating user-generated content is expensive, slow, and requires a vast global infrastructure because of the inability of AI automation to identify content to remove.

Interviewees noted that there are studies indicating frustration and counterintuitive impacts of content moderation. Harvard Kennedy School found that improving the amount of truthful information had a more powerful effect than removing misinformation.⁵³ Correcting people on Twitter leads to more toxic and less accurate future retweets. Researchers found causal evidence on Twitter that the experience of being corrected increases the partisan slant and language toxicity of a user's subsequent retweets and had no significant effect on the user's primary tweets. Researchers inferred that those individuals felt defensive after being publicly corrected by another user, which shifted their attention away from accuracy concerns. The researchers note this presents an important challenge for social correction approaches.⁵⁴

To date, there has been relatively little effort to look beyond content moderation to design technology that contributes to social cohesion. Some interviewees noted that it is natural that a company would start from the place where they are getting the most criticism by removing "bad stuff" from showing up. A negative experience can be more impactful than a

^{51.} Tiffany Hsu and Eleanor Lutz. "More Than 1,000 Companies Boycotted Facebook. Did It Work?" New York Times, August 1, 2020).

^{52.} See for example, Ravi Iyer, "Content Moderation is Dead," *The Psychology of Technology Newsletter*, October 7, 2022;, Valerie C. Brannon and Whitney K. Novak, "Online Content Moderation and Government Coercion," *Congressional Research Service* (May 13, 2022).

^{53.} Alberto Acerbi, Sacha Altay and Hugo Mercier, "Research note: Fighting misinformation or fighting for information?" *Harvard Misinformation Review*, January 12, 2022.

^{54.} M. Mosleh, C. Martel, D. Eckles and D. Rand, "Perverse Downstream Consequences of Debunking: Being Corrected by Another User for Posting False Political News Increases Subsequent Sharing of Low Quality, Partisan, and Toxic Content in a Twitter Field Experiment," *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (May 2021): 1–13.

positive one for users.

V. THE TECH REGULATION NARRATIVE

Tech critics hold a different view. Most large social media products are not a simple reflection of society.⁵⁵ Rather than blaming users for harmful content online, tech critics point to how the affordances and algorithms of tech products incentivize and reward harmful content. Instead of focusing on the "symptom" of harmful content, tech critics argue that the focus should be on the system incentivizing harmful content. The tech regulation narrative suggests that regulation should go beyond privacy and antitrust issues to address tech profit models and the affordances and algorithms baked into the design of some social media products which create perverse incentives for toxic polarization.

A. Optimized for User-Engagement and Profit

The Center for Human Technology insists that technology is not neutral.⁵⁶ Tech product design features foster "a race to the bottom of the brainstem" and an "attention economy" that rewards divisive content, resulting in "polarization spills." Tristan Harris of the Center for Humane Technology describes Twitter as a "gladiator stadium, like a Roman Coliseum, where people are being told that they need to debate free speech and ideas in a marketplace of ideas with balls and chains and arrows and swords."⁵⁷ According to a Pew Survey, a minority of highly active users post the majority of tweets, and nearly half of Twitter users in the US are silent observers of extreme and violent users.⁵⁸

In 2021, Facebook whistleblower Francis Haugen revealed damning internal reports documenting staff concerns that the company was driving polarization in countries around the world. Countless researchers and journalists from outlets such as the *Wall Street Journal* and the *New York Times* were documenting the evidence.⁵⁹ In 2020, the *Wall Street Journal* published an article claiming Facebook was ignoring and undermining efforts to address polarization. The article suggested Facebook's internal research found that its algorithms were increasing polarization by exploiting "the human brain's attraction to divisiveness." The article cited a 2018 slide

^{55.} Dean Eckles. "Algorithmic Transparency and Assessing Effects of Algorithmic Ranking." Testimony before the Senate Subcommittee on Communications, Media, and Broadband. (December 9, 2021), https://www.commerce.senate.gov/services/files/62102355-DC26-4909-BF90-8FB068145F18.

^{56.} The Center for Humane Technology, "The Myth of Neutrality," March 31, 2022.

^{57.} Tristan Harris, "Humane Technology on 60 Minutes," *Your Undivided Attention Podcast*, November 10, 2022; Tristan Harris and Aza Raskin, "Elon, Twitter and the Gladiator Arena," *Your Undivided Attention Podcast*, October 27, 2022.

^{58.} Meltem Odabas, "5 Facts about Twitter 'Lurkers'" Pew Research Center. (March 16, 2022).

^{59.} John D. McKinnon and Ryan Tracy, "Facebook Whistleblower's Testimony Builds Momentum for Tougher Tech Laws," *The Wall Street Journal*, October 5, 2021.

from an internal presentation that noted that "if left unchecked [Facebook algorithms optimized for profit would offer] more and more divisive content to gain user attention and increase time on the platform." The article stated that Facebook researcher and sociologist Monica Lee gave a presentation in 2016 that detailed how Facebook was fueling extremism. The 2016 slides state that extremist content that is "racist, conspiracy-minded and pro-Russian" is found in a third of all large German Facebook groups, and "64% of all extremist group joins are due to our recommendation tools."

Tech critics point to engagement-based profit models that incentivize and optimize for polarizing and extremist content that keeps users engaged with emotional content. Harvard Business School Professor Shoshana Zuboff refers to the social media profit model as *surveillance capitalism*. Tech companies capture more private user information and attention to ads the longer a user uses a product. User-engagement metrics translate to company profit as ad companies pay more to access more users.⁶¹

The Center for Humane Technology calls this the "race to the bottom of the brainstem." The user-engagement profit model driving social media tech companies translates into "design choices that will create a more addicted, distracted, outraged, polarized, validation seeking, and narcissistic society" while leaving people vulnerable to political actors waging psychological influence campaigns.⁶² Harris states, "Twitter's business model of engagement is about making sure that every post, every moment of anger, every moment of controversy is as maximally visible and interactive with as many other people as possible."⁶³ Just as CNN found that its profits increased when it offers round the clock crisis coverage, social media companies profit more when their "trauma inflating" algorithms amplify anger and injustice.⁶⁴

The film *The Social Dilemma* portrays three challenges social media poses for society. First, there is a *mental health dilemma* that relates to internet addiction, depression, anxiety, and a loss of an ability to have agency, or the ability to make decisions. Second, there is a *discrimination dilemma* that relates to the subjective biases and prejudices in algorithms that amplify oppressive dynamics. Third, there is a *democracy dilemma* that relates to the role of some tech products in undermining public trust in democratic institutions, public interest journalism, and elections.

^{60.} Jeff Horwitz and Deepa Seetharaman, "Facebook Shut Efforts to Become Less Polarizing," *The Wall Street Journal*, May 27, 2020.

^{61.} Shoshana Zuboff, *The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power* (New York: Public Affairs, 2019).

^{62.} Tristan Harris, "Humane Technology on 60 Minutes," Your Undivided Attention Podcast, November 10, 2022.

^{63.} Ibid.

^{64.} Tristan Harris and Aza Raskin, "Can Psychedelic Therapy Reset Our Social Media Brains? with Rick Doblin," *Your Undivided Attention Podcast*, December 15, 2022.

Tech design affordances may reduce individual agency, a marker of social cohesion. *Engagement-driven affordances* such as "likes" and "shares" fuel social comparisons and foster greater use or even addiction-like obsessions, deteriorating mental health and well-being, and keeping users scrolling rather than taking actions to improve the quality of life for themselves and their communities. *Engagement-driven algorithms* rank content to promote and recommend divisive content designed to keep people on the tech product longer. *Engagement-driven data collection* on user location, ideas, behaviors, beliefs, networks, and identities creates a databank of information that governments and political actors can use to surveil the public. This surveillance may fuel distrust between the public and institutions with access to their data. Institutions may use this surveillance to repress certain identity groups or civil society groups advocating for human rights or democracy.

B. Disrupting Public Interest Journalism

In addition to tech affordances and algorithms that amplify polarizing content, there are also a wider set of digital impacts on social cohesion.⁶⁵ Digital advertising is diverting money away from public interest media where it helped fund local news and investigative journalism.⁶⁶ The decline in the availability and quality of legacy media (newspapers, radio, TV) enables disinformation online and offline to spread.

Media fragmentation leads to users encountering similar conspiracies, partisan, or false information in a hybrid online/offline media ecosystem that reinforces political divides. Public surveys document a decline in public trust in journalism.⁶⁷ The growth of partisan media and a growing disinformation industry seem to contribute to epistemic insecurity where the public is unsure who to believe or what is true.⁶⁸

Cumulatively, the design of digital spaces and their optimization for user engagement and profit comes at the expense of social cohesion.

C. The Limits of Tech Regulation

Perhaps unwittingly, social media products are a de facto digital public sphere: a space for discussion of issues that affect people's lives. But the engineers and tech innovators who created most social media products had no training in designing public spheres. Few companies consulted social

 $^{\,}$ 65. Lisa Schirch, Digital Space and Peace Processes. (Geneva, Switzerland: Interpeace, Fondation Hirondelle, ICT4Peace. May 2022).

^{66.} Derek Wilding, Peter Fray, Sacha Molitorisz, and Elaine McKewon, *The Impact of Digital Platforms on News and Journalistic Content* (Sydney: University of Technology NSW, 2019).

^{67.} Katherine Fink, "The Biggest Challenge Facing Journalism: A Lack of Trust," *Journalism* 20, no. 1 (2019).

^{68.} Emanuel Adler and Alena Drieschova, "The Epistemological Challenge of Truth Subversion to the Liberal International Order," *International Organization*, 75 no. 2 (2021): 359–386.

scientists. Most tech companies lack staff with appropriate backgrounds to anticipate and respond to governing and addressing toxic polarization. Tech companies also lack the political legitimacy to do the policing of these new town squares, particularly for moderating "political entrepreneurs" who use political messages to instill fear and division within potential voters, and for industrial-level harmful digital content created by cyber armies and disinformation industries.⁶⁹

A Congressional Research Service report offers a summary of the Townsquare doctrine, a legal theory that says that certain types of technically private spaces still have certain types of protections for freedom of expression. In other words, when a technology product grows and becomes widely used, it has public responsibilities. To Most governments are not yet prepared to keep up with tech innovations that create new public spaces that demand new types of guidelines and regulations. Tech companies face dilemmas to define the limits of free speech on their products and the social norms for these spaces. While Harvard Law's Lessig implored humanity to understand that "code is law," he also writes about the inability for government laws to adequately regulate code. Digital spaces have become resistant to regulation.

To users and government regulators, a company can tout its product as a neutral communication platform where anyone can communicate. To advertisers and investors, a company can tout its product as an "advertising" or "marketing" platform where "users" and their private information and attention are the product being sold. It might take researchers 20 years to determine exactly how much technology companies are responsible for harming human agency, polarizing communities, and undermining trust in democratic institutions. But there is a precedent for not waiting for the absolute scientific consensus on tech impacts on polarization when the stakes are so high. Policymakers have taken action to restrict potentially harmful medicines and toxins even before science proves harm.

Government regulation of tech companies has focused on privacy and cybersecurity concerns, not the affordances and algorithms that amplify toxic polarization. Tech products optimized for user engagement, advertising, and profit incentivize the spread of false and hateful posts.

^{69.} Jennifer McCoy and Murat Somer, "Toward a Theory of Pernicious Polarization and How It Harms Democracies: Comparative Evidence and Possible Remedies," *The Annals of the American Academy of Political and Social Science* 681 no. 1, (December 20, 2018): 234–271.

^{70.} Valerie C. Brannon. "Free Speech and the Regulation of Social Media Content." *Congressional Research Service* (March 27, 2019).

^{71.} *Id*.

^{72.} Lessig, Code and Other Laws of Cyberspace.

^{73.} Tarleton Gillespie, "The Politics of 'Platforms," New Media and Society 12 no. 3 (2010): 47–364

^{74.} Jonathan Haidt, "Yes, Social Media Really Is Undermining Democracy: Despite What Meta Has to Say," *The Atlantic*, July 28, 2022.

Regulating algorithms can curb tech platform's prioritizing profit over people. Yet the speed of the movement for government regulation of technology platforms harmful impacts on society is nowhere close to catching up to the impacts of digital amplification of disinformation and hate speech on intergroup relations, the escalation of threats to electoral integrity, or the decline in public trust in institutions.

VI. THE SOCIAL COHESION BY DESIGN NARRATIVE

A third way of approaching technology companies' roles in responding to or preventing harmful content goes a step further. In *addition* to content moderation and tech regulation, tech companies can design tech products with affordances and algorithms that support social cohesion.

Like governments, technology companies have a tremendous amount of power to steer human behavior. Governments contribute to social engineering by providing public schools, enforcing a criminal justice system, and building roads and bridges. These activities encourage people to behave in "prosocial" ways that encourage humanizing and expressing concern for others. Societies encourage social cohesion when they use benevolent manipulation to incentivize and structure prosocial behavior.

At the 2022 Trust and Safety Research Conference at Stanford University, former Twitter VP of trust and safety Del Harvey urged tech companies to look beyond content moderation toward designing for health. Harvey, oft described as Silicon Valley's "chief sanitation officer" for her role in removing harmful digital content, 75 explained that it is not enough for tech companies to remove harmful content. Tech companies could learn from public health to move beyond reaction to prevention. Yet when asked if tech companies had an appetite to "design for health," other panelists with Harvey indicated they did not observe such an interest.

Several interviewees for this report noted that there had been some internal experiments to incentivize positive content to build social cohesion. For example, Facebook created a "Common Ground" team mandated to improve intergroup relationships. But this did not last long. Citing concerns from conservatives in the U.S., Facebook hired Republican leader Joel Kaplan in 2011 to be policy chief and to vet proposed changes. Kaplan expressed concern that changes to encourage better communication skills were "paternalistic," calling the vetting process "Eat Your Veggies." According to the *Wall Street Journal* and tech staff interviewed for this report, Kaplan approved some changes but blocked other proposals because they lacked "rigor and responsibility" related to effectiveness and might have led to unintended consequences. Facebook disbanded the Common Ground Initiative citing political bias, social engineering, and cognitive

manipulation around the end of 2018. The Central Integrity Team also disbanded, though other dispersed integrity teams continued.⁷⁶

Elsewhere in the world, a range of new pro-social technologies or "peacetech" aims to decrease polarization, improve social cohesion, and advance computational democracy. These strategies offer a compelling alternative paradigm for thinking about Trust and Safety and the dilemmas of content moderation.

eBay's dispute resolution product designer Colin Rule believes that computer code can act as a mediator and foster positive social interaction. A tech product can provide the structure to coach users on what they can say to increase the chance of a positive encounters. In this case, Rule asserts that tech product designs are a form of "benevolent manipulation." As system designers, tech companies can provide the "walls" to structure positive behavior and enhance social cohesion online. Other tech products could learn from eBay's example of coaching users to communicate more effectively.

Rule estimates that 90% of individual bad behavior is from someone with a first offense. A warning and "the first mistake is free" approach offers an education opportunity to reinforce community guidelines. Rule notes that tech products could send a message describing why a piece of content was harmful. For example, a prompt could tell users "you said a hurtful thing on the forums, and people were upset about it, so your content got flagged. Please watch this short video to learn more about healthy conflict and effective communication." If someone makes another offense, their access to the product or forum can then be reduced. They might, for example, only be allowed to post 30 times a month. Upon next offense, they would be limited to 15 posts a month. And then on the fourth offense, they might be "de-platformed" or lose all posting ability for 6 months. Building on this example, other tech companies could use harmful content as an opportunity to coach users in effective communication. While not universally welcome, this approach might fuel less anger and rebellion than immediate censorship.

Peacetech is an umbrella term referring to forms of technology that improve social cohesion. Peacetech enables pro-social content. It is part of a broader field of public interest technology that uses technology to advance public interest, generate public benefits, and promote public good. Peacetech may include other public interest technologies such as civictech which informs citizens on public interest issues and services, connects people with others, and facilitates communication with their government; and govtech which helps governments to facilitate communication with citizens to improve public services and public engagement.

VII. THE HISTORY AND SCOPE OF PEACETECH

The idea to design technology to support social cohesion dates back to the 1980s.⁷⁷ The term "peacetech" emerged from multiple places in the early 2000s.⁷⁸ The Swiss think tank <u>ICT4Peace</u> began research on peacetech in 2003.⁷⁹ In 2004, the US Institute of Peace in Washington, DC initiated what is now known as the PeaceTech Lab.⁸⁰ In the mid-2000s, local tech innovators in Sri Lanka and Kenya designed tech products to support early warning of violence and citizen journalism. In 2007, the tech company <u>Ushahidi</u> began using tech for the prevention of election violence. In the same year, Stanford psychology professor B.J. Fogg began teaching courses and researching ways technology could be used to support peace, which he called "peace technology."⁸¹ Building on this research, Stanford Peace Innovation Lab continues to create models for peacetech.

In 2020, the UN Secretary-General released a Roadmap for Digital Cooperation, detailing a robust "digital transformation" agenda supporting the innovation of tech products that support the UN's Department of Peacebuilding and Political Affairs. ⁸² The United Nations is also investing in a suite of technology tools to support the UN Department of Political and Peacebuilding Affairs (DPPA). Countless NGOs are also exploring peacetech and digital peacebuilding. The NGO Build Up works with partners around the world to support civil society in learning how to use peacetech and authored *Search for Common Ground*, a Digital Peacebuilding Guide, which provides insight into how to choose what type of technology to use to support social cohesion. ⁸³ The NGO swisspeace conducts research and an online course on digital peacebuilding. The Alliance for Peacebuilding hosts a community of practice on "Digital Peacebuilding" that offers monthly meetings to learn about new types of peacetech.

There are now centers around the world devoted to peacetech, including

^{77.} For a longer history, see Lisa Schirch, "25 Spheres of Digital Peacebuilding and PeaceTech," (Tokyo: Toda Peace Institute, 2020).

^{78.} See for example, Yiannis Laouris, "Information Technology in the Service of Peacebuilding: The Case of Cyprus," *World Futures* 60, no. 1 (December 2003): 67–79; Helena Puig Larrauri and Anne Kahl. "Technology for Peacebuilding," *Stability: International Journal of Security and Development*, 2 no. 3, (2013); Ioannis Tellidis & Stephanie Kappler. "Information and Communication Technologies in Peacebuilding: Implications, Opportunities and Challenges." *Cooperation and Conflict* 51, no. 1, (March 2016); Pamina Firchow, Charles Martin-Shields, Atalia Omer, and Roger Mac Ginty. "PeaceTech: The Liminal Spaces of Digital Technology in Peacebuilding," *International Studies Perspectives* 18, no. 1 (February 2017); Lisa Schirch, "Social Media Impacts on Social & Political Goods: A Peacebuilding Perspective," Toda Peace Institute. Policy Brief #38. (April 2019).

^{79.} ICT4Peace, "History," https://ict4peace.org/about-us/history/.

^{80.} PeaceTech Lab, https://www.peacetechlab.org/.

^{81.} B.J. Fogg, *Peace Technology: Why a Class about Facebook Apps?* Scribd, (2007) https://document.pub/document/dr-bj-fogg-facebook-peace-technology.html?page=1.

⁸² UN Secretary General. Roadmap for Digital Cooperation, (New York: United Nations, 2020).

^{83.} Search for Common Ground, Build Up, ConnexUs, "Digital Peacebuilders Guide," https://howtobuildup.stonly.com/kb/guide/en/digital-peacebuilders-guide-X49wcx4IFi/Steps/1469015.

the University of Waterloo's Grebel Peace Incubator, the University of Bristol's Interdisciplinary PeaceTech Group, and the University of Notre Dame's PeaceTech and Polarization Lab. In Florence, Italy, the European University Institute held the first Global Peacetech Conference in November 2022.84

VIII. FUNCTIONS OF PRO-SOCIAL PEACETECH

Prosocial technology contributes to social cohesion in four broad ways. First new tech platforms help to analyze digital harms and polarization. Second, technology products can improve human agency to participate in civic issues affecting their lives. Third, technology products can support intra-group and inter-group communication and joint problem-solving. Fourth, tech platforms can improve public trust and inclusion in governance.

A. Tech for Analyzing Digital Harms and Polarization

Understanding the dynamics of polarization is an essential element of planning effective social cohesion programs. Every context has a unique information ecosystem and a unique set of conflict dynamics. For many decades, conflict analysis and context assessment tools have been essential to developing effective peacebuilding and development programs. Society Growing polarization and state-sponsored disinformation campaigns highlight the need to add an analysis of information ecosystems and how digital spaces and their interaction with offline spaces drive conflict. So

Strategic planning on the use of digital tools to support social cohesion begins with first analyzing information ecosystems.⁸⁷ The United Nations is investing in a suite of technology tools to support social media analysis.⁸⁸ Sparrow is a social media analysis tool created by and for the UN Department of Political and Peacebuilding Affairs (DPPA) for analyzing Twitter to identify trending topics, hashtags, and key influencers.

Another example is called Phoenix. The peacebuilding NGO Build Up

^{84.} See Kalypso Nicolaidis and Michele Giovanardi. Global PeaceTech: Unlocking the Better Angels of Our Techne, EUI RSC, 2022/66, Global Governance Programme-481, Europe in the World, https://cadmus.eui.eu/handle/1814/74985.

^{85.} Lisa Schirch, Conflict Assessment and Peacebuilding Planning: Toward a Participatory Approach to Human Security (Boulder, Colorado: Lynne Rienner Press, 2013).

^{86.} Fondation Hirondelle, Demos, Harvard Humanitarian Initiative and ICREDES, "Influencers and Influencing for Better Accountability in the DRC" (July 2019).

^{87.} Branka Panic, Data for Peacebuilding and Prevention Ecosystem Mapping: The State of Play and the Path to Creating a Community of Practice (New York: NYU Center on International Cooperation, 2020).

^{88.} See for example, United Nations, *Digital Technologies and Mediation in Armed Conflict*. Helsinki: Department of Political and Peacebuilding Affairs; Centre for Humanitarian Dialogue, 2019; Global Pulse, *E-Analytics Guide: Using Data and New Technology for Peacemaking, Preventive Diplomacy and Peacebuilding* (New York: United Nations, 2019).

and the technology company DataValuePeople partnered to create Phoenix, an open-source, non-commercial, customizable process and tool to support peacebuilders and mediators who want to work ethically with social media data to inform programming. Local communities first develop contextually grounded problem statements that address peacebuilding objectives. The groups then use Phoenix to create a data pipeline to add social media sources, along with labeling and visualization tools. Phoenix offers new ways to understand the drivers of conflict and the opportunities for peace.⁸⁹

B. Tech to Support Individual Agency

Some tech products support individual agency so that people have the capacities and belief that they can participate in civic action to work on issues that affect their lives. These platforms can help people feel that they have a voice by providing tools for them to share their identity, experiences, beliefs, and passions. Some platforms offer affordances such as hashtags to enable isolated individuals to find each other to form larger movements, such as with the hashtags #MeToo and #BlackLivesMatter.⁹⁰

Other tech products help individuals to reality check their perceptions to help individuals recognize they there is more common ground between people than commonly assumed. For example, digital quizzes such as The Perception Gap, developed by the bridge-building organization More in Common, provide individuals with an opportunity to reflect and test whether their perceptions of other groups match reality. People in different countries could take the quiz and find out how realistic their view was of people on the other end of the political spectrum. This serves an important role in "reality testing" and challenging people's presumptions about other groups. Helping individuals realize they do not accurately understand their political opponents might prompt them to be curious to learn more so that they may correct their perceptions and understanding.

There are a variety of tech products to help people learn effective communication skills and to model how to have a healthy conversation on difficult issues or conflicts. For example, Games for Peace uses Minecraft games between Israeli and Palestinian youth. ⁹² In addition, individual influencers on TikTok are offering conflict resolution advice using hashtags such as #resolveconflict. Another example comes from Karin Tamerius, founder of Smart Politics, who created an "Angry Uncle" Chatbot to help coach people in effective communication skills for having political

^{89.} See Build Up, https://howtobuildup.org/programs/digital-conflict/phoenix/.

^{90.} Sarah J. Jackson, Moya Bailey, and Brooke Foucault Welles, #HashtagActivism: Networks of Race and Gender Justice, (Cambridge, MA: MIT Press, 2020).

^{91.} See The Perception Gap, https://perceptiongap.us/.

^{92.} See Games for Peace, https://www.gamesforpeace.org/.

conversations at holiday dinners.⁹³ The Canadian-based Suzuki Foundation created a climate conversation coach bot called CliMate. Other organizations offer cooperative video games between groups in conflict.⁹⁴

Based on his experience building eBay's Online Dispute Resolution (ODR) system, Colin Rule helped to set up the mechanism for eBay users (sellers and buyers) who had disputes. Rule found in his dispute resolution work with eBay is that if you have a dispute between a buyer and a seller, it is not helpful to give an open text box to the buyer. Tech product designers do not need to let "everyone" talk to "everyone." This gives people too much ability to generate more anger and havoc for themselves and others. They may engage in threats and insults because that is the way they think they can get a sense of fairness. They are angry and frustrated and want the other side to know that. Instead of giving the complainant an open textbox where they vent that anger, products can instead structure more constructive communication by giving them a forum where they can make selections. What kind of problem do you have? What kind of solution do you want?95 Users leave with a positive sense of resolution and empowerment, a key element of social cohesion. eBay has resolved millions of disputes through this system. eBay bots coach complainants to rephrase and reframe their messaging to take out insults. The seller has an incentive for that buyer to be happy because the buyer is unhappy, and they leave them negative feedback that is going to impact their ability to sell on the site.⁹⁶

Similarly, some have suggested that popups, a box, symbol, or window that appears when you begin writing on a computer, might offer users feedback on their tone. On Twitter, such a concept could include informing users with a popup stating, "I see you might be headed for an uncivil conversation?" ⁹⁷

C. Tech to Support Horizontal Cohesion

As described above, there are two forms of horizontal cohesion. Intragroup cohesion is known as "bonding" social capital. Inter-group cohesion is known as "bridging" social capital.⁹⁸

There are several examples of new tech startup companies that focus on

^{93.} Karin Tamerius, "How to Have a Conversation With Your Angry Uncle Over the Holidays," *New York Times*, November 18, 2018.

^{94.} David Suzuki Foundation. "How and Why to Have Climate Change Conversations," https://davidsuzuki.org/what-you-can-do/how-and-why-to-have-climate-change-conversations/.

^{95.} Amy J. Schmitz and Colin Rule. "Lessons Learned on eBay," The New Handshake: Online Dispute Resolution and the Future of Consumer Protection, *American Bar Association Section on Dispute Resolution*. (2017): 33–46.

^{96.} Interview with Colin Rule, February 18, 2022.

^{97.} Molly Wood, "Twitter Hires Social Scientists to Help Figure Out Our Conversation Problem," Marketplace, September 25, 2019.

^{98.} Robert D. Putnam, *Bowling Alone: The Collapse and Revival of American Community* (New York: Simon & Schuster, 2020).

intra-group bonding, particularly for individuals meeting online for social or work purposes. Gatheround is a video conferencing tech product that describes itself as "a team bonding and community engagement platform for people-focused organizations seeking to build relationships and strengthen teams in an era of disconnection and distraction." Co-founder Lisa Conn, formerly director of the Common Ground Initiative at Facebook, differentiated Gatheround from Zoom because it is "designed for how humans connect."99 As in real life, individuals on Gatheround do not see themselves, just the other participants to whom they are talking. The video focuses on people at a "nose-biting" distance to encourage participants to be kind to each other. Gatheround offers conversation prompts for people to share their experiences, so they feel more heard and seen. There are no affordances to mute or turn off the camera, making it impossible for people to check out of the conversation. There are no backgrounds so people see where you are sitting. Gatheround has a share/facilitation feature where, like a talking stick in a dialogue, a question is asked, and people form a line with equal time to speak. Conn describes how this disrupts existing power dynamics and structures to provide more equity. 100

A second example of tech-supported intra-group cohesion is Marco Polo, a social media product focused on well-being and happiness in a closed social network. Marco Polo offers a video chat or video voicemail with a front-facing camera. Marco Polo emerged from a sense that people had turned to lower-quality text-based communication and had stopped calling each other and having conversations. Text-based products may increase the quantity of relationships at the cost of the quality of relationships. Cofounder Vlada Bortnik describes the "thoughtful, human-centered design" as focusing on the quality of the connection. As a person records a video chat, they look into the camera, like looking in the mirror. There are no filters or glamor, but rather an encouragement by design to be authentic and intimate. This may make it more likely to present positive body language and less likely to spew hate at someone. The home screen in Marco Polo is chronological. There are no counts of friends, likes, or emojis, as the product does not want to have "vanity metrics." The experience in Marco aims to be enriching and nourishing, to increase happiness, and attempt to curb loneliness. Marco Polo does not want to increase anxiety by urging users to be competitive. Instead, the design aims to be a tool for humans intrinsically motivated to use rather than be manipulative of a person's time. Marco Polo does not sell user data and does not advertise on Google because it wants to protect the privacy of users. The staff at Marco Polo assert that because people on Marco Polo connect more intentionally with their closest friends, the threat of encountering harmful content is lower.

^{99.} Interview with Lisa Conn, December 18, 2021.

^{100.} Ibid.

Marco Polo does allow people to block someone in their network. But because the video chat is asynchronous, a user cannot talk over somebody. Marco Polo aims to be a product that encourages people to listen to each other.¹⁰¹

Other tech-support platforms aim to improve inter-group cohesion to build "bridging" social capital between people that belong to different social groups. Intergroup relations can improve in a variety of ways. These tech products offer affordances for people to explore differences as well as common ground. Some products seek to create safe or "brave" spaces for dialogue across the lines of conflict. The examples here illustrate that technology can scale empathy and understanding between groups, as well as increase a group's capacity for solving problems.

For example, the tech designers at Soliya set out to pair technology with the power of dialogue in 2003, before the rise of social media. Soliya is a Virtual Exchange product to foster "high impact inter- and cross-cultural education facilitated through digital technology." Soliya hosts dialogues between 15,000 young people per year in small, diverse groups to share their perspectives on identity and current events. Soliya has held a special focus on intercultural dialogue between young adults in the West and the Arab and Muslim World. Soliya is unique in part because participants' videos show up in a circle, surrounding a prompt for the dialogue that asks a question and participants focused on a topic. ¹⁰²

Another approach to building social cohesion is to invite people to engage in a conversation not to prove who is right but rather to see who can change another person's views. This exercise requires users to listen to other points of view and then try to build a bridge between worldviews. Scottish teenager Kal Turnbull founded Reddit's "ChangeMyView" community which invites people to "post an opinion you accept may be flawed, to understand other perspectives on the issue and to encourage users to enter with a mindset for conversation, not debate." Users rewarded compelling arguments with a delta symbol (Δ) to indicate when someone changed their mind.

Turnbull extended the subreddit community by creating a new website. "Change a View helps internet commenters see eye-to-eye, where the forum breaks us out of our online filter bubbles, and where we relearn how to talk to each other online." Users note the digital space feels like an "oasis" while journalists call it "our best hope for civil discourse." Change a View uses Jigsaw's comment-ranking engine, called Perspective API, which

^{101.} Interview with Vlada Bortnik, March 2, 2022.

^{102.} Marta Guarda, "Giving Voice and Face to Other Cultures: The Soliya Connect Program and the Development of Intercultural Communicative Competence," *Carte d'Occasione* 5 (2013), 111–131.

^{103.} See https://www.reddit.com/r/changemyview/.

^{104.} Arielle Pardes, "Change My View' Reddit Community Launches Its Own Website," Wired, April 6, 2019.

scores comments, demotes harmful content, and eases moderator loads. Change a View provides a template for how to improve difficult conversations online.

Researchers examine the affordances of "ChangeMyView" that enable effective communication on the platform, namely the "game" elements and its social norms. Gamification is a method of turning an activity into a game to increase motivation. Gamification provides enjoyment and social approval through competition, with the award of a delta sign, which accumulates into "delta scores." Participants told researchers that the incentive of earning a delta encourages them "to be civil to one another." Users observed that the people who were able to change the views of others were "polite in their posts." Users also noted the role of moderation of trolls and people who were rude or not open-minded. Jigsaw is experimenting with using the ChangeMyView comment ranking engine to detoxify online conversation. It could be adopted by news agencies for comment sections or by major social media platforms.

The peacebuilding organization Build Up launched "The Commons" as an intervention to depolarize political conversations on Twitter and Facebook in the USA. Paid facilitators initiated thousands of conversations across some of the most polarized individuals and polarizing topics. The goal was to help people engaged in polarized conversations to have more positive conversations, to increase interest in promoting civility, and to change how they engage with people on social media. Facilitators found polarized conversations by curating a list of top hashtags and content creators at the center of US political conversations. Bots would then identify who was open to a conversation, and humans would engage with them. 107

There are a wide variety of other new tech startups aiming to improve intergroup dialogue. Some are exploring how to use virtual reality to foster intergroup dialogue and empathy. The group HackthePlanet offers a variety of VR programs to build intergroup understanding. ¹⁰⁸ A platform called Kazm bills itself as a "conversation engine." Kazm describes itself as a social platform built for community, as opposed to other platforms built for an audience. Kazm's affordances offer support and tools for dialogue facilitators and community administrators. ¹⁰⁹ Kazm offers bridge-building and dialogue groups like Living Room Conversations affordances including a way to have members connect via video dialogue, join events, access content, and comment on discussion boards. Kazm also offers video

^{105.} Shagun Jhaver, P. Vora, and A. Bruckman. *Designing for Civil Conversations: Lessons Learned from ChangeMyView, GVU Center Technical Reports* (2017).

^{106.} Pardes, "Change My View' Reddit Community Launches Its Own Website."

^{107.} Anooj Bandari, "The Commons: Where Are We at in 2021? *Medium*, September 27, 2021, https://howtobuildup.medium.com/the-commons-where-are-we-at-in-2021-e58ae31fc196.

^{108.} See Hack the Planet, https://www.hack-the-planet.io/.

^{109.} See Kazm, https://about.kazm.com/.

coaching to guide the conversation, prompts to focus a dialogue, videos, polls, and word clouds. Noting its role in social cohesion, Kazm representatives state that there are no ads or trolls on Kazm.¹¹⁰

D. Tech to Support Vertical Social Cohesion and Public Trust

Other new tech startups aim to improve vertical cohesion by enabling citizens to participate in governance through "civtech" which enables citizens to engage in collective problem-solving on policy topics, and "govtech" by enabling governments to design more inclusive processes for consulting with citizens on public issues.

These tech products recognize that social cohesion does not require preventing the expression of tension or conflict or making people be superficially "nice" to each other. Social cohesion also does not require people to form personal relationships or have direct contact.

For example, Ushahidi is a crowdsourcing and mapping tool that enables citizens to report potential violence to governments on simple mobile phone applications. In 2007, Ushahidi helped to prevent election violence in Kenya. Local people reported where tensions were rising in the streets. This information was shared with local civil society mediation and peace teams as well as police. Since then, Ushahidi has grown significantly and now the platform is used to enable citizen reporting and coordination about civil society and governments to respond to public issues. It provided real-time information to defuse electoral-related violence in the streets. ¹¹¹ Ushahidi has been used in Haiti and Nepal to coordinate relief efforts, monitor and report on corruption in Indonesia, help address sexual violence in Egypt, and map police violence in Portland, Oregon.

Other platforms enable inclusion and participation in decision-making by making it easier for people to participate and by creating incentives to identify common ground or consensus. The 2014 tech start-up Remesh began with the mission to create a technology that would, in the words of founder Andrew Konya, "represent the will of the people and amplify their collective voice." Conflict mediators, civil society groups, or governments can use Remesh to dialogue with and poll the public. Remesh software can extract key themes and draw insights from a dynamic and open-ended "conversation" with up to 1,000 people. 112 The UN used Remesh in Libya to gather stakeholder opinions on a proposed interim government. In Yemen, the UN used Remesh to listen to public perceptions of a cease-fire and opinions on the prospects for a peace process. The UN is now

^{110.} See Alliance for Peacebuilding, "Scaling Facilitated Dialogue with Conversation Engines," YouTube (November 18, 2021), https://www.youtube.com/watch?v=qwkhc-8jZaI.

^{111.} Juliana Rotich. "Ushahidi: Empowering Citizens through Crowdsourcing and Digital Data Collection." Field Action Science Reports. no. 16. (2017), 36–38.

^{112.} Interview with Andrew Konya, March 20, 2022.

considering using Remesh for peace support in Sudan, Mali, Afghanistan, and Iraq.¹¹³

Inspired by insights from social cohesion efforts in nonviolent communication and attempts at collective decision-making in the Occupy Movement, Colin Megill designed the tech platform Pol.is to improve computational democracy. Experiments in Taiwan and the UK illustrate that Pol.is can help a divided public find areas of common ground and develop policy solutions on polarized public issues. In Taiwan, the government has used Pol.is dozens of times on different issues resulting in government action 80% of the time.¹¹⁴

The Pol.is platform is optimized for consensus building, finding common ground, and fostering citizen engagement. Polis provides "a real-time system for gathering, analyzing, and understanding what large groups of people think in their own words, enabled by advanced statistics and machine learning." Pol.is enables "collective intelligence" and fosters mutual "listening at scale" through digital citizen assemblies that use tools to support "computational democracy." The platform gathers both qualitative data and quantitative data. Unlike other platforms, on Polis users do not reply to each other's posts. Rather users submit an idea (one at a time) that others can up-vote or down-vote. This affordance enables users to reward ideas that address the interests of most people and generate new and better solutions. The lack of a "reply" affordance prevents trolling and abuse, and thus removes the pain and heat from discussions. Pol.is operates on opensource code allowing anyone to use the platform to host public dialogues seeking to find consensus.¹¹⁵ Pol.is seems to incentivize the development of creative options that meet the interests of diverse stakeholders and enables "thinking outside the box" to envision positive future coexistence.

Few of the tech start-ups designing new products to support social cohesion are reaching the scale necessary to address toxic polarization. While these examples offer insight into design affordances that may support social cohesion, to date, the most popular tech platforms like Twitter and Instagram continue to offer polarizing affordances that enable social comparisons and polarizing algorithms that prioritize user engagement. These newer platforms face a challenge in drawing people away from networks where their existing friends are posting content. More research is needed to find out what prevents or encourages users to explore platforms that emphasize prosocial design.

Big tech companies with the scale to shift societies away from

^{113.} Jordan Bilich, Michael Varga, Daanish Masood and Andrew Konya, "Faster Peace via Inclusivity: An Efficient Paradigm to Understand Populations in Conflict Zones," AI for Social Good Workshop at NeurIPS, Vancouver, Canada (2019).

^{114.} Josh Smith, Toby O'Brien and Harry Carr, "Polis and the Political Process," *Demos*, August 3, 2020.

^{115.} Interview with Colin Megill, February 20, 2022.

polarization and toward social cohesion will need to learn from and adapt the design affordances and algorithms from smaller startup tech companies. For example, Twitter recently drew inspiration from Pol.is to create incentives for individual agency and participation in negotiating the validity or truthfulness of digital posts. Pol.is engineers optimized the platform to contribute to social cohesion and to put guardrails on the platform to limit harmful content. Learning from Pol.is' affordances and algorithms, Twitter staff developed a program called Community Notes (formerly Birdwatch) to empower Twitter users to add helpful notes to Tweets that might be misleading. Wired Magazine calls Twitter's experiment "one of the most exciting content moderation innovations ever to come out of not just Twitter, but any major platform."116

Aviv Ovadya and Jonathan Stray have been writing about the potential of big tech companies to adopt the types of bridging ranking systems found in platforms like Pol.is and Remesh.¹¹⁷ There are more opportunities for big tech companies to test the use of affordances and algorithms in divided communities.

IX. CONCLUSION

Digital town squares are increasingly important for information sharing and deliberation. But disinformation and other harmful content plague digital spaces and are amplifying toxic polarization. Toxic polarization prevents society from solving pressing problems and can contribute to violence. Toxic polarization online requires a multi-stakeholder response. This paper explored three complementary approaches for responding to harmful digital content.

The user-centered narrative assumes users alone are to blame for harmful content. To date, tech companies have focused primarily on content moderation to remove or weaken the impact of user-generated content. Content moderation is important, but it is not keeping pace with the scale of harmful digital content and toxic polarization. Even staff at companies who have hired tens of thousands of content moderators expressed dismay at the task of managing a "tsunami of harmful content" without adequate resources, particularly in the Global South where they lack staff who speak local languages. Given recent tech layoffs in Trust and Safety teams, this approach to reducing harmful content seems unlikely to reduce toxic polarization in the near future.

^{116.} Carl Miller, "Elon Musk Embraces Twitter's Radical Fact-Checking Experiment," WIRED Magazine, November 28, 2022.

^{117.} Jonathan Stray, "Designing Recommender Systems to Depolarize," First Monday. 27, No. 5-2 (July 11, 2021). https://arxiv.org/abs/2107.04953. Aviv Ovadya, "Bridging-Based Ranking: How Platform Recommendation Systems Might Reduce Division and Strengthen Democracy," Belfer Center, Harvard Kennedy School (2022).

While this paper did not address digital media literacy as a strategy to reduce user-generated harmful content,¹¹⁸ there is a movement afoot to improve digital communication norms and strengthen public immunity to harmful disinformation by inoculating people to "prebunk" conspiracy theories and other false and deceptive content online.¹¹⁹ A mass digital media literacy public education effort will take time and has only begun in a few countries such as Finland.¹²⁰

The *tech regulation narrative* asserts that the design of technology platforms is not neutral. Design affordances and algorithms can amplify toxic polarization or help to build social cohesion. While tech regulation is important, digital spaces are resistant to regulation and digital polarization spills are undermining policy solutions. While governments focus on issues like privacy and cybersecurity, the challenge of regulating algorithms and design affordances on tech platforms will likely be slower and more challenging. Governments will need to create incentives for tech companies to reduce harmful content amplified by their algorithms and design features, either by changing their profit model and/or paying taxes on their polarization spills to help fund social cohesion efforts.¹²¹

The *pro-social design narrative* asserts that tech algorithms and affordances can amplify social cohesion. Designing technology to support social cohesion is an alternative and a complement to these other approaches. Pro-social tech platforms already exist, and we can learn from these tech design affordances and algorithms that support social cohesion. Computer engineers with training in social cohesion designed some of the technology platforms described in this article. Others started as initiatives of the UN or NGOs in partnership with tech startups to create products that would support bridge-building and peacebuilding work. Scaling social cohesion requires partnerships between practitioners and tech platforms to design better platforms and improve how people use tech in democratic processes. Big tech companies, new tech startups along with private and public funders can invest in building new tech platforms aimed to improve social cohesion.

Addressing the tsunami of false and hateful information online requires this type of innovation—designing and scaling tech products that support

^{118.} Media Literacy: A Definition and More, Center for Media Literacy, Center for Media Literacy, http://www.medialit.org/media-literacy-definition-and-more.

^{119.} Stephan Lewandowsky and Sander van der Linden, "Countering Misinformation and Fake News through Inoculation and Prebunking." *Null* 32, no. 2 (2021): 348–384, https://www.tandfonline.com/doi/full/10.1080/10463283.2021.1876983.

^{120.} Jenny Gross, "How Finland Is Teaching a Generation to Spot Misinformation," New York Times, January 10, 2023.

^{121.} Helena Puig, "Societal Divides as a Taxable Negative Externality of Digital Platforms," Ashoka Tech and Humanity (2023), https://www.next-now.org/sites/default/files/2023-03/Societal%20Divides%20as%20a%20taxable%20negative%20externality%20of%20digital%20platforms_0.pdf.

social cohesion. Content moderation, tech regulation, digital media literacy, and designing tech for social cohesion can be complementary. Together, they offer a way forward to address the *system* and not just the *symptom* of harmful content online.