# Let's Solve the Social Dilemma

## A SAPIEN WHITEPAPER

sapien
humans first

# Contents

> **The Social Dilemma documentary ... traces the evolution of social media from its early days as a promising tool for social connection, into a multi-billion dollar corporate enterprise that makes money by hacking the brains and manipulating the behavior of users**

# Introduction

*The Social Dilemma* documentary was released on Netflix in September, 2020. It was viewed in 38 million homes within the first month of its release, raising public awareness of the threat posed by corporate-owned, profit-driven digital ecosystems.

The filmmakers make a compelling case that the Social Media Establishment is contributing to a host of serious problems, from digital addiction, depression, and suicide, to increased social polarization, the proliferation of fake news and conspiracy theories, vulnerability to political interference from foreign powers, and a general breakdown of shared notions of truth and objectivity.

Facebook receives a lot of attention in the film, but the problems discussed apply to all digital platforms that use algorithms to promote user behaviors that are optimized to benefit corporate owners.

Unfortunately, *The Social Dilemma* devotes the majority of its runtime to describing these problems, and far less time discussing possible solutions. And, as we will show, what solutions it does present are open to serious objections.

The question remains—*how do we solve The Social Dilemma?*

**We believe that any viable solution to the serious problems raised by *The Social Dilemma* will include key elements of the solution that Sapien is developing.**

To make this case we need to properly understand

1.  the nature of the problem posed by establishment social media platforms;

2.  the solution strategy that is proposed by the makers of *The Social Dilemma*;

3.  why this strategy is unlikely to be successful; and

4.  what a truly viable solution would look like.

Section 1 gives an overview of the core problems raised by modern social media as it is presented in *The Social Dilemma* film and as articulated by Tristan Harris and other members of the *Center for Humane Technology*. This section unpacks the business model and the technology that drives modern social media applications, and the causal narrative that connects this technology to the serious social problems discussed in the film.

Section 2 outlines the solution strategy that is presented in *The Social Dilemma* and elaborated in greater detail by members of the *Center for Humane Technology*. This strategy focuses on teaching media literacy skills to help individuals reclaim control of their attention, and on the broader goal of applying pressure on social media companies to reform the way they design their technology.

Section 3 raises skeptical objections to such a reform strategy, offering reasons to doubt that social media companies will make any substantive changes to a business model that has made them among the richest and most powerful industries on the planet.

Section 4 presents two conditions that, it is argued, any viable solution to *The Social Dilemma* must satisfy. These two conditions define a solution space that excludes any current member of the Social Media Establishment.

> "It's like playing chess against a computer that isn't restricted by human limitations, that can search through millions of possible moves to find the optimal next move. It's not a fair fight

# Understanding the Problem

*The Social Dilemma* documentary turns a spotlight on the negative impacts of social media on contemporary society. It traces the evolution of social media from its early days as a promising tool for social connection, into a multi-billion dollar corporate enterprise that makes money by hacking the brains and manipulating the behavior of users.[1]

The film is best viewed as a public awareness campaign sponsored by the *Center for Humane Technology*, a non-profit policy think tank founded by Tristan Harris and Aza Raskin that promotes ethical reform of social media technologies.[2] The messages and arguments presented in the film are largely based on the work of Harris and his colleagues at CHT.

The analysis described in *The Social Dilemma* identifies the root problem with the **business model** that creates profits for social media companies, that is driven by a persuasive technology stack that is algorithmically engineered to manipulate the behaviors of users.

Figure 1 gives a big-picture overview of the primary argument of the film.

The top row is a list of social problems that have grown in severity in recent years. The claim is that the social media business model operates in a way that exploits human vulnerabilities and limitations, creating pathological desires and behavioral habits that create risks for users and exacerbate higher level social problems. They are not arguing that social media is the SOLE cause of these problems, only that it has become an increasingly important contributing factor, and in some cases the primary driver.

**Figure 1:** The primary argument of The Social Dilemma.

## SOCIAL MEDIA COMPANIES SELL PREDICTIONS OF CONSUMER BEHAVIORS

To better understand this story we need to unpack the business model (Figure 2).

We start with users who consume and post to social media. Social media applications automatically track all aspects of user activity, down to every last click, every moment a user looks at the screen. This information is used to create detailed prediction models of a user's preferences and behavioral dispositions—what kinds of content a user is likely to engage with, what type of accounts a user is likely to follow and interact with, what ads a user is likely to click on, etc.

**Figure 2:** Here's a big-picture overview of the primary argument of the film.

These predictions are aggregated into statistical models that, together with a set of proprietary tools for using this data to create targeted personalized messaging, are sold to third parties that have an interest in delivering offers and messaging to users. This is the primary source of revenue for social media companies.

## SOCIAL MEDIA APPLICATIONS REINFORCE USER BEHAVIORS THAT OPTIMIZE FOR PROFIT

A key feature of this business model is that social media applications have been designed to modify the behaviours of users in order to optimize for commercially desirable outcomes, like engagement with the platform, growth of the user base, and receptivity of users to advertising. These optimizations are algorithmically driven, using machine learning to determine what interactions will be most effective for achieving these results.

These algorithmic mechanisms are hidden to the user, and can be difficult to explain to audiences unfamiliar with them. *The Social Dilemma*



**ENGAGEMENT ALGORITHM**   **GROWTH ALGORITHM**   **ADVERTISING ALGORITHM**

**Figure 3:** Social media algorithms personified in The Social Dilemma.

addresses this challenge with a narrative device where the algorithms are personified and shown making decisions in real time about what content and notifications to show the teenager who is using the app (Figure 3).

## SOCIAL MEDIA CREATES "LEADS" TO AUCTION TO ADVERTISERS

It is commonly said that social media sells user attention to advertisers. This doesn't quite capture what is new and concerning about this technology. It is more accurate to say that social media companies sell "qualified leads" to advertisers—they package users that have been identified as receptive to an offer, and present them for auction to the highest bidder.

But even this doesn't tell the full story; it doesn't take into account the power of social media to *shape and influence users*. The content we consume on social media doesn't simply reflect our pre-existing interests, it can also create interest where there was none before, or magnify interests that previously did not occupy much space in our thoughts. Beyond merely identifying receptive audiences, social media can turn a user into a receptive audience, creating qualified leads for advertisers.
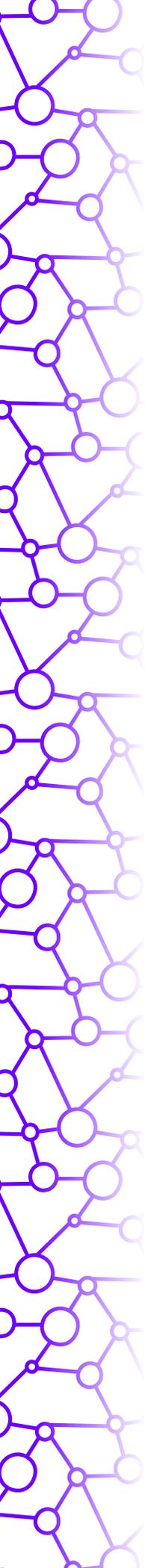
> It is commonly said that social media sells user attention to advertisers. This doesn't quite capture what is new and concerning about this technology. It is more accurate to say that social media companies sell "qualified leads" to advertisers—they package users that have been identified as receptive to an offer, and present them for auction to the highest bidder.
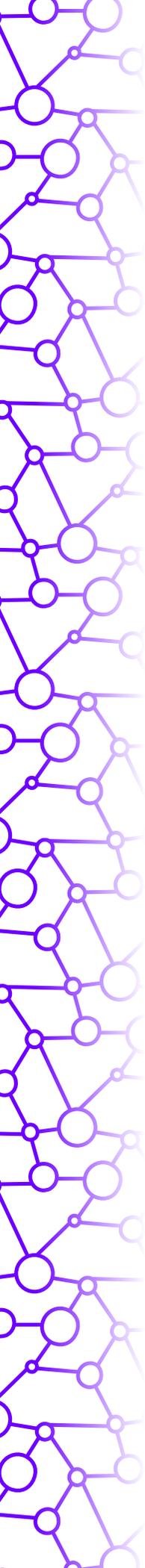
## SOCIAL MEDIA EXPLOITS HUMAN VULNERABILITIES

Again, it may be objected that these sorts of persuasion and marketing tactics are really nothing new. Part of marketing a business is working to create demand for your product. But there is another important step in *The Social Dilemma's* argument that needs to be considered.

The tech industry insiders interviewed in the documentary spend considerable time describing how the persuasive technologies built into social media are designed to exploit our psychological vulnerabilities.

The most direct route to behavior modification is by triggering our basic human needs for security, approval and belonging. Many features of social media applications have been designed for just this purpose. Notifications, infinite scroll, like and share buttons, were all purposefully designed tools for inducing feedback that the user can control, acting on our brain's reward centers like food pellets for a rat's feeding machine, or slot machine payouts for a gambler.

The content we consume provides another important set of behavior-modifying stimuli. Social media platforms use algorithms to serve up the content that a user sees in their timeline and in their recommendations lists. This content is optimized to incentivize users to stay on and engage with the platform.

It's important to understand that the algorithms don't care whether the content a user sees is true or false, supported or unsupported, uplifting or depressing, positive or negative. All they track is how users *respond* to the content. Did a user click on it? If so, how long did they engage with the content? Did they comment on it? Did they share it? Who did they share it with? Did they watch an ad that was adjacent to the content?

On the basis of a user's history of engagement on the platform, the algorithm serves up content and recommendations that it predicts will optimize for engagement, growth, and ad revenue.
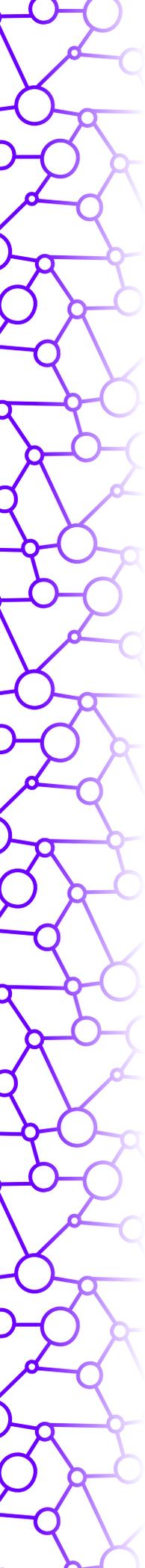
This is the logic that drives many worrisome features of social media, like the "extremist rabbit-hole" effect. Users start out searching for a piece of information and get sucked into a recommendation sequence that serves up ever more extremist content, for no other reason than because such content reliably provokes a response and captures attention.

### IT'S NOT A FAIR FIGHT, AND WE ARE OVERWHELMED

As tools for marketing persuasion and behavior modification, modern social media platforms really have brought something new to the table. Their features are designed to generate behavioral feedback loops that are habit-forming. The content they present to users is selected for its ability to draw and hold attention and stimulate engagement on the platform, in ways that are designed to maximize growth and profits for the platform. This business model has given rise to a new vocabulary to describe it—"surveillance capitalism", the "attention economy", etc.[3]

There is one more concern that we need to address, because it's central to the larger story of the risks that social media poses to human welfare, and to the concerns that motivate Tristan Harris and his colleagues at the *Center for Humane Technology*.

The behavior modification strategies that are coded into social media platforms are partially the product of intentional engineering by human designers, but increasingly they are the product of machine learning programs searching for inputs (how users are prompted to
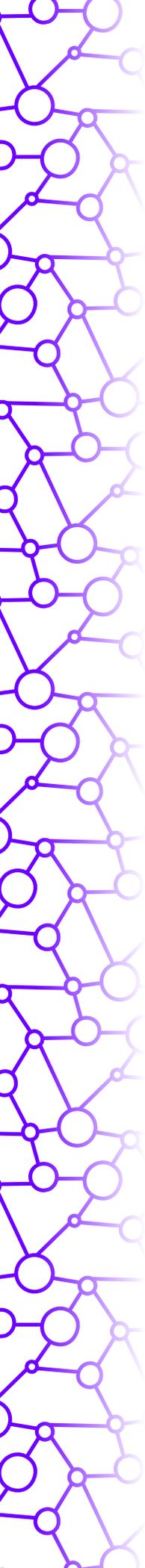
engage on the platform) that generate optimal outputs (increases in engagement, growth and ad revenue). In other words, we have super-computers running machine learning programs that are constantly looking for ways of nudging users toward certain behaviors, with no restrictions on the way those outcomes are achieved.

The result has been what Tristan Harris calls "a race to the bottom of the brain stem".[4] The algorithms have discovered what we already know, that human behavior is most easily influenced by triggering responses related to our primal needs for security, approval and be-longing. But the algorithms are able to identify very specific patterns of stimuli, that no human being could come up with, that are per-sonalized for each user, to produce the desired effects on the user's behavior. It's like playing chess against a computer that isn't restricted by human limitations, that can search through millions of possible moves to find the optimal next move. It's not a fair fight.[5]

Put another way, social media algorithms are discovering strategies for achieving optimal results by effectively *hacking the human brain*. Algorithmically-driven optimization overwhelms our normal psycho-logical functioning by exploiting our brain's natural vulnerabilities.

A complete list of such vulnerabilities would be very long. We have limited memory, attention and ability to process information. Our mo-tivational system is tied to the activity of neurotransmitters like dopa-mine and serotonin that can be manipulated through rewards. As so-cial animals we have basic needs for social approval, acceptance and validation that are important drivers of human behavior. We often rely on shortcuts to make rapid decisions under conditions of uncer-tainty. We are naturally disposed to hostility toward out-groups that appear threatening to us. We rely on our group identities to determine which authorities and narratives to trust.

What social media machine learning algorithms have *taught them-selves* is how to use the features and content available on social media platforms to nudge human behavior by exploiting such vulnerabilities.

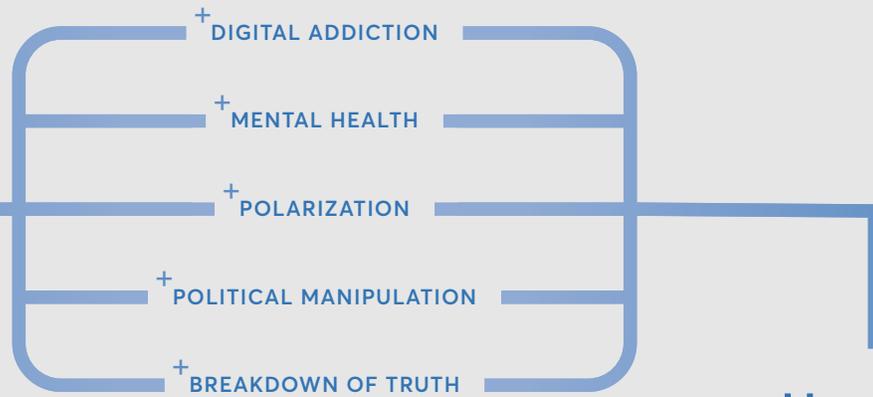## SOCIAL MEDIA DOWNGRADES OUR HIGHEST HUMAN CAPACITIES

The further claim that Tristan Harris and his colleagues make is that these "hacks" don't simply exploit our vulnerabilities—they *overwhelm* our capacity to resist them, and thereby *erode* our higher capacities for reflective problem-solving, rational decision making and collective action toward common goals.

The argument for this claim isn't fully articulated in *The Social Dilemma*, it's presented more as an assertion. But it's an argument that Tristan Harris has been making for several years. Figure 4 gives a picture of the reasoning.

The social media business model influences users in ways that exploit their cognitive vulnerabilities and result in various kinds of pathology and dysfunction. Information overload results in shortening of attention spans. Social media use triggers dopamine hits that lead to addictive habits, which lead to social isolation. Our need for social validation is overloaded by a constant stream of likes and shares that induces a form of social narcissism that makes us vulnerable to depression and self-abuse. Our vulnerability to outrage is fed by a constant stream of outrage-inducing content that drives polarization and pathological tribalism. And so on.

Each of these causal chains is a hypothesis that is open to testing from various scientific disciplines. Some are better supported than

**Figure 4:** Social media dysfunction leads to "human downgrading"

others, but the *Center for Humane Technology* has collected a large body of research that is relevant to these claims.[6]

The collective impact of this suite of social pathologies is what Harris calls "human downgrading".[7] His claim is that over time, social media use is contributing to systemic social problems that are effectively downgrading our ability to think critically, solve complex problems, make smart decisions, and cooperate on shared projects that serve the common good.

### HUMAN DOWNGRADING IS AN EXISTENTIAL RISK

*The Social Dilemma* presents human downgrading as an *existential risk*, meaning that it poses a threat to the continuing existence of human life as we know it. This may sound like hyperbole, but the key point is that the *primary* concern of the filmmakers isn't with polarization or addiction or mental health issues *per se*, but rather the
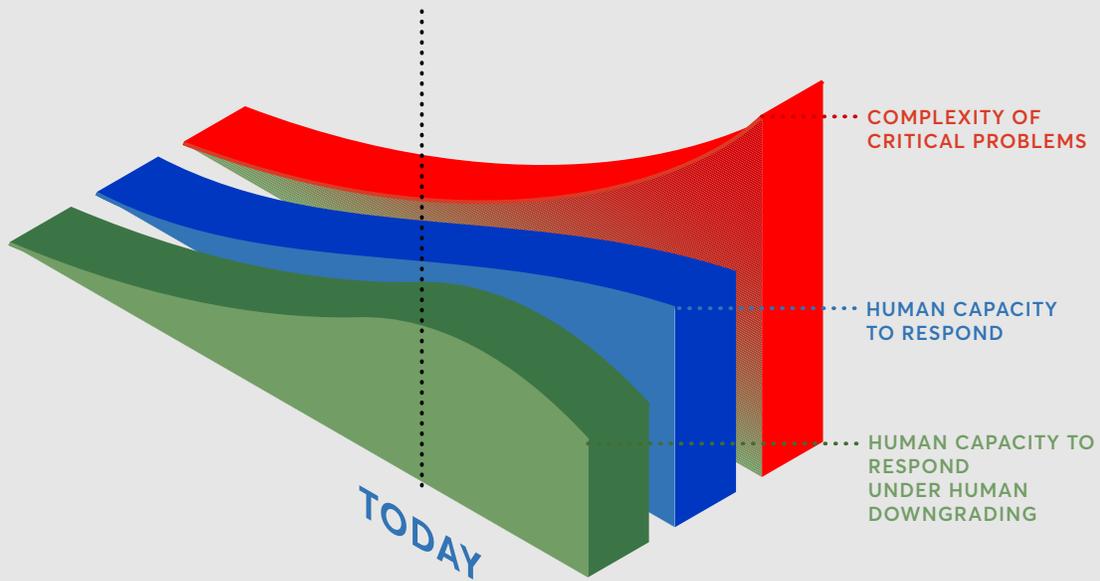
**Figure 5:** Human downgrading reduces our capacity to solve critical problems.

collective impact of all these pathologies on our ability as a species to solve critical problems that may threaten our existence.

Figures 5 and 6 are reproduced from a presentation delivered by Tristan Harris in 2019.[8] They describe a scenario in which accelerating technology presents increasingly complex problems that require coordinated collective action to solve. Managing climate change is the classic example, but one can also consider the threat of global pandemics, or the risks of runaway artificial intelligence, or any number of potential threats posed by disruptive technologies with unpredictable impacts.

The dotted blue line describes our current capacity to respond to such problems. The red line indicates the increasing scale and complexity of the problems we can expect to face in the future. Already, the concern is that such problems will outstrip our capacity to solve them. However, the introduction of modern social media, and the attendant "human downgrading" to which it is contributing, only make the problem worse.

**Figure 6:** The gap in human problem solving capacity that needs to be overcome.

The green line indicates our *eroding capacity* to respond to critical problems under human downgrading. The gap between this line and where we'll need to be to adequately respond to future problems only grows as time progresses.

This is presented as a cautionary narrative, and it begs a fuller argument to ground it. But it is the root of the urgency in the messaging of *The Social Dilemma*. This is the mission of the *Center for Humane Technology*, to raise awareness of these risks and to propose new models for the ethical design of digital technologies that could help to avoid these outcomes.

> Changes need to be made that not only halt the process of human downgrading, but ultimately reverse the trend—our technology must actually promote human upgrading

# The Social Dilemma Solution Strategy

Now that we understand how *The Social Dilemma* views the problem posed by modern social media, we can look at the solution strategy that the film proposes.

As noted earlier, the film doesn't spend much time talking about solutions, but the general strategy is clear. It has two parts (Figure 7).

**PART 1: EMPOWER USERS TO PROTECT THEMSELVES FROM THE HARMFUL EFFECTS OF SOCIAL MEDIA.**

If users are going to use social media, they should be empowered to reassess and reboot their relationship to technology, to reclaim control over their digital lives. Helpful strategies include:[9]

- learn more about how the social media business model works and the problems it promotes
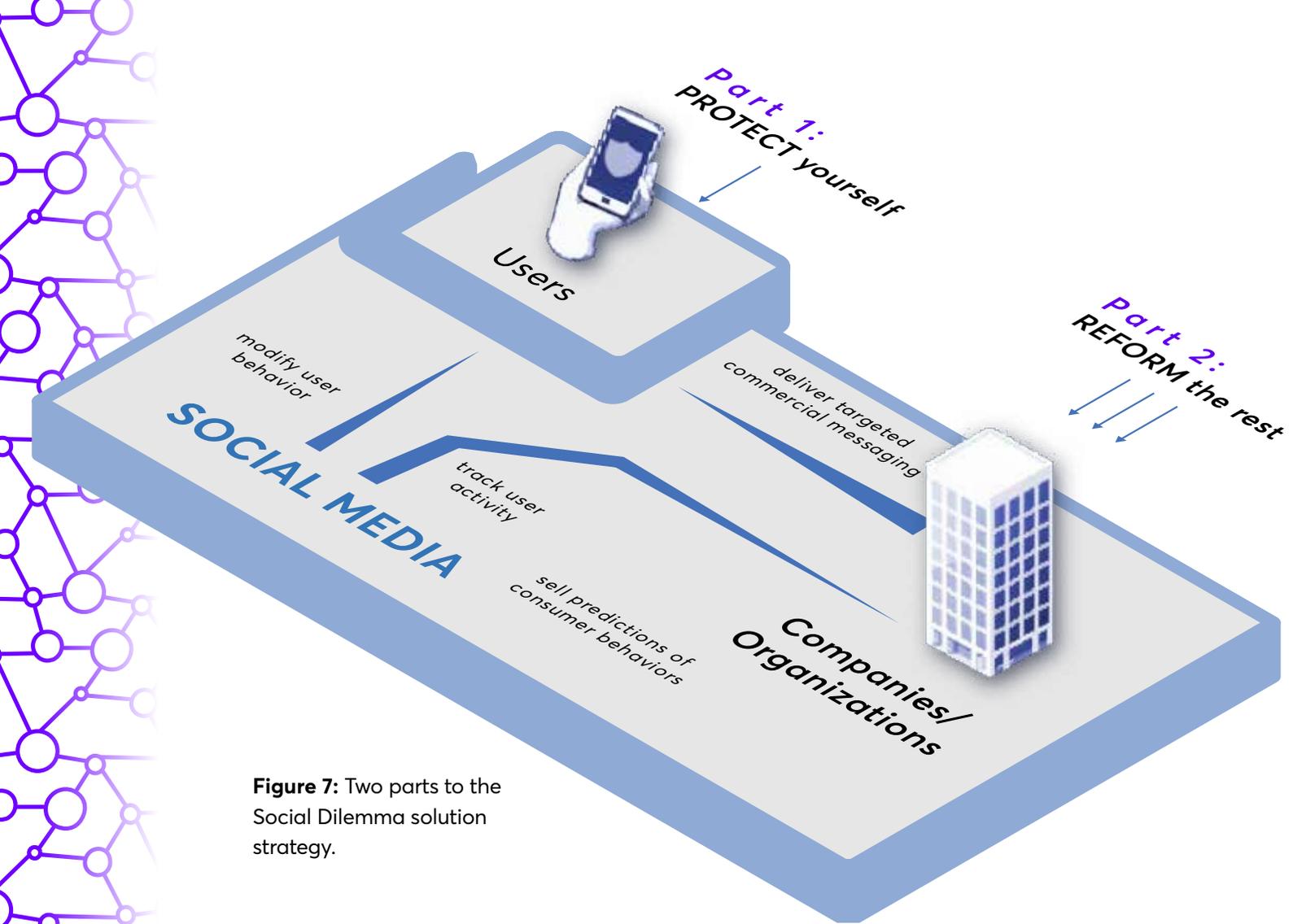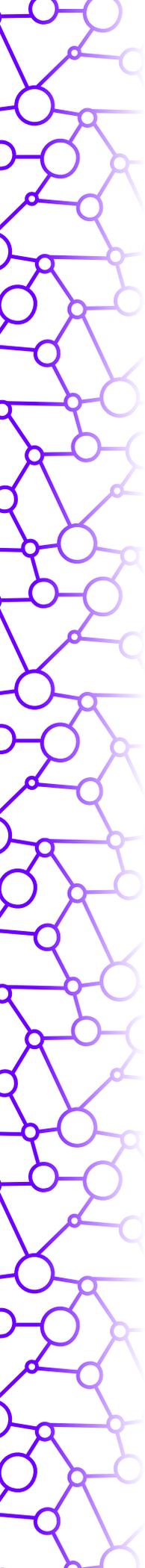
- turn off notifications

**Figure 7:** Two parts to the Social Dilemma solution strategy.

- use alternatives to "toxic" apps

- track screen time

- monitor and set limits on cell phone and social media usage

- organize social groups where all parties are committed to helping one another reclaim their attention and moderate their use of social media

**PART 2: REFORM THE SOCIAL MEDIA ESTABLISHMENT.**

A strategy that only teaches social media awareness and self-defense to individual users won't solve the larger problems regarding the negative impacts of social media on society. Even if a million users commit to healthier practices, there are billions of human beings continuing to use social media as they always have. To address this problem we need to reform the way that social media companies design the technology that powers their platforms.

**HUMANE DESIGN IMPLIES RADICAL REFORM OF THE BUSINESS MODEL**

The central concept of the *Center for Humane Technology*'s reform strategy is "**humane design**". This term has been used for several years to describe the general goal of designing technologies with genuine human needs and interests in mind, to develop technologies that promote human flourishing rather than sacrifice it for the sake of profits or other shortsighted motives.[10] The *Center for Humane Technologies* is named for this principle.

Precisely what humane design entails for social media is an open-ended question, but from the perspective of the CHT, the criteria for success is clear. Changes need to be made that not only halt the process of human downgrading, but ultimately *reverse* the trend— our technology must actually promote human *upgrading*. This is the only outcome that will solve the existential risk problem we are facing.

It's clear that a "humans first" approach to social media design requires a radical shift in the social media business model. Tristan Harris describes the goal in stark terms: "We have to end Attention & Surveillance Capitalism".[11]

## EXTERNAL PRESSURE AND INTERNAL PRESSURE

How do we change the way the entire Social Media Establishment does business? How do we persuade the richest and most powerful tech companies in the world to reform their practices?

*The Social Dilemma* documentary describes a two-front strategy (Figure 8):

1. apply **external pressure** from the public, the media, shareholders, and government regulatory bodies

2. apply **internal pressure** from within social media companies themselves, from the engineers, managers and executives that are actually responsible for developing the technology.
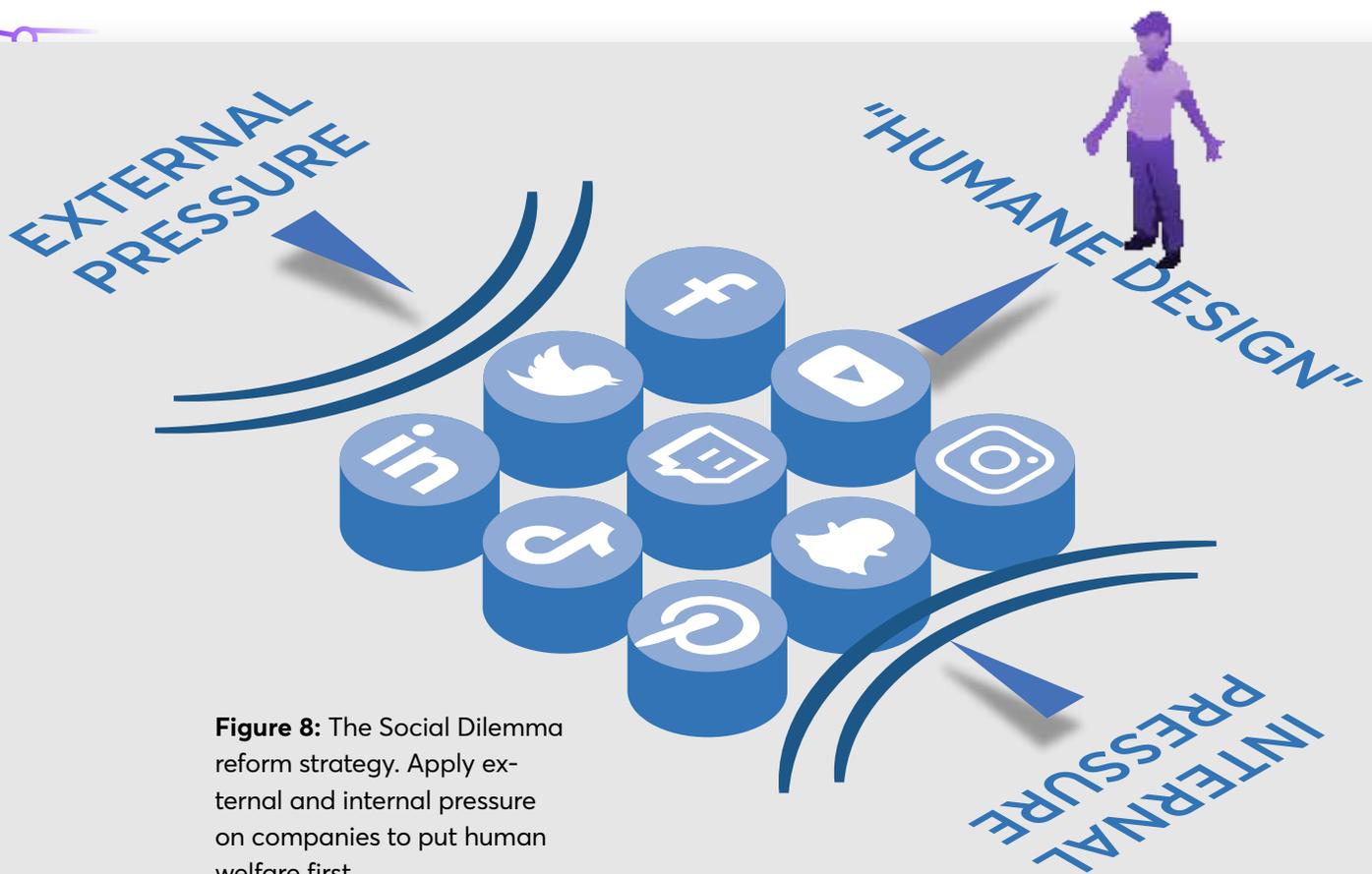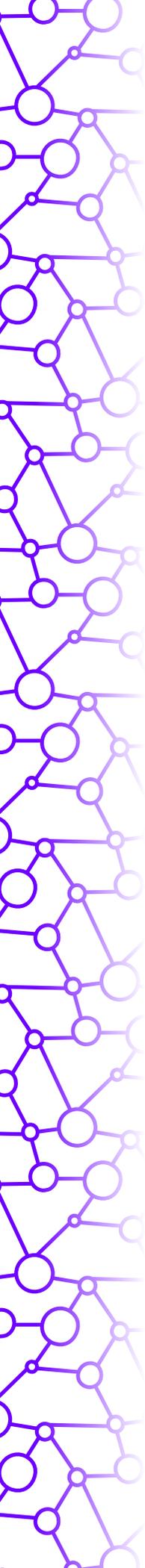


**Figure 8:** The Social Dilemma reform strategy. Apply external and internal pressure on companies to put human welfare first.

*External* pressure strategies include **mass awareness and education campaigns** that motivate users to reconsider their social media usage, like *The Social Dilemma* documentary, and Tristan Harris's *Your Undivided Attention* podcast. It also includes **educating policymakers** about social risks and policy proposals that protect society and reward humane technologies, to drive and support government regulatory pressures on the tech industry.
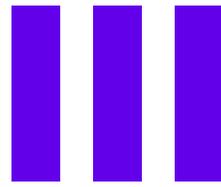
*Internal* pressure strategies include **working directly with technologists** by delivering educational presentations and practical workshops on humane design principles, and by **supporting crucial conversations** within tech companies.

The hope is that pressure from these two fronts will motivate real change in the business model that currently drives the operation of social media companies.

One reason for optimism around this reform strategy is that, unlike a problem like climate change whose origins are diffused across the whole planet, the power to dramatically change the direction of social media technology lies in the hands of no more than 1000 people, an ironic advantage of the centralization of power and influence in the hands of a few dominant corporations. If those industry leaders can be persuaded to change course, the whole industry will shift.

> **Big business does not have an inspiring track record when it comes to accepting responsibility for public health or ethical concerns involving the use of its products.**

# III Challenges Facing the Social Dilemma Solution Strategy

*The Social Dilemma* documentary, along with a vigorous promotional effort on the part of Tristan Harris, has been very successful in raising awareness within the public of these problems with social media. More people than ever are examining their social media usage and asking critical questions about the social media industry.

And there is no doubt that we need education in "social media literacy" that empower users to reclaim control over their attention and develop more intentional habits that minimize the potential harms.

However, there are reasons to be skeptical that a reform strategy of the kind described above will be adequate to the task.
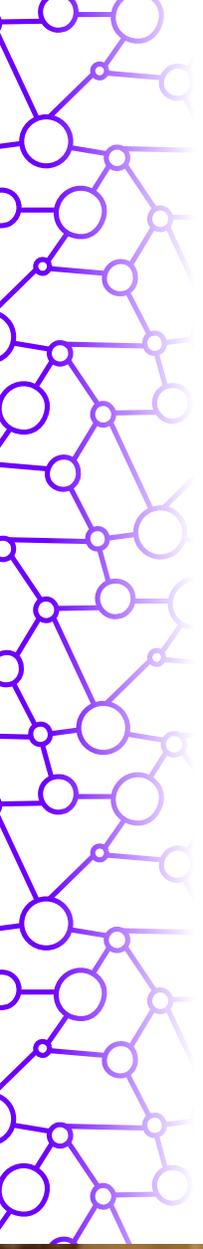
**Figure 9:** Tobacco industry leaders deny that nicotine is addictive.

## OBJECTION 1: FOR BIG TECH IT IS MUCH EASIER TO *APPEAR* VIRTUOUS THAN TO *BE* VIRTUOUS

Big business does not have an inspiring track record when it comes to accepting responsibility for public health or ethical concerns involving the use of its products. Historically the most common pattern is, whenever possible, to push back on such claims and defend the safety of the product and the reputation of the company.

Consider the tobacco industry's history of resistance to charges that smoking is addictive and causes cancer. It is now well established that the main tobacco companies were aware of the addictive qualities of nicotine as early as the 1950s, and had been manipulating nicotine levels in cigarettes for the purpose of influencing their addictive properties. Yet in 1994 the top executives of the seven largest American tobacco companies testified in Congress (and on national television) that they did not believe that cigarettes were addictive[12] [13] (Figure 9).
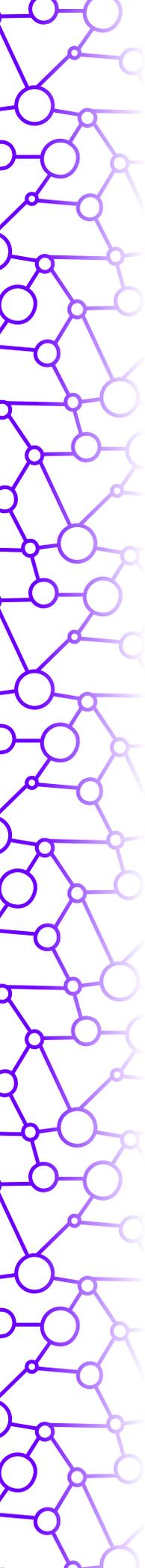
They also insisted that, while cigarettes *may* cause lung cancer, heart disease and other health problems, "the evidence is not conclusive".

Similar stories of corporate resistance can be found in almost any large industry. The sugar and soft drink industries, for example, have a history of pro-sugar propaganda that challenges claims that sugar consumption can be addictive or that it is correlated with increases in heart disease and obesity, even while internal evidence was accumulating that this was the case.[14]

The tech industry is no exception. In July 2020, industry leaders Mark Zuckerberg of Facebook, Tim Cook of Apple, Sundar Pichai of Google and Jeff Bezos of Amazon were called to testify before a Congressional antitrust committee. The running motif of the meeting was a carefully managed defense against all claims that these companies engage in any form of exploitative, unethical, illegal or harmful practice[15] (Figure 10).

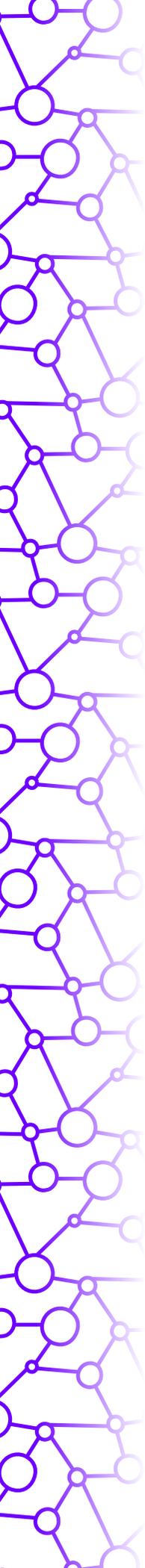**Figure 10:** Big tech leaders deflect criticisms of their practices.

Another common pattern in corporate responses to external criticism is to deny any wrongdoing that would admit legal or ethical culpability, while also creating internal teams or divisions within the company that are dedicated to investigating and developing recommendations for improvement on issues raised by critics. Facebook and Google both have dedicated AI Ethics teams, for example, though the activities and directives of these teams are not transparent.[16] [17] [18]

A cynic might view such teams as part of a broader public relations/reputation management strategy that allows tech companies to signal their responsiveness to public concerns and their commitment to making progress, without making any significant changes to how they do business. Tristan Harris was himself employed as a design ethicist for Google, but left because he felt powerless to create change from within the organization.

But even if we grant the sincerity of the people working internally to address ethics issues, the fact remains that the incentive structure of large corporate enterprises is strongly oriented toward protecting the financial interests of shareholders and maintaining a growth trajectory within the industry. These incentives increasingly come to dominate corporate decision making as companies grow in scale. The sheer size of the social media tech giants is itself a reason to expect that the default response to demands for change will be resistance.

## OBJECTION 2: REAL CHANGE IS UNLIKELY WITHOUT A VIABLE ALTERNATIVE BUSINESS MODEL

It is unrealistic to expect an entire industry to abandon a business model that is responsible for its growth and continuing dominance in the world.

The Social Media Establishment runs on an extractive attention/surveillance economy business model. Tristan Harris admits that we can't realistically expect Facebook and Google to change in any significant way if it requires them sacrificing their own business interests for the sake of the public good.

A much more conducive environment for reform is one where all parties can point to *at least one viable alternative business model* that shows how social platforms can be both financially successful and humanely designed.

Such an alternative would provide a genuine option for social media users who are awakened to the problems discussed in *The Social Dilemma*, and a model for a different business strategy that technology companies could explore.

In the absence of such an alternative, there is little reason to expect external pressures alone to generate significant change within the Social Media Establishment.

The general point is reluctantly acknowledged by Harris, and the need for alternatives is recognized. In one presentation Harris shows the graphic redrawn in Figure 11, adding a third category of pressure, "aspirational pressure", as an important driver of humane technological change.[19] Aspirational pressure is described as *pressure arising from the presence of alternative models* that exemplify the ideals of humane technology, that can inspire creative exploration for entrepreneurs and provide real options for users.

*The Social Dilemma* documentary says almost nothing about alternative business models, nor does it discuss any concrete examples of

social platform alternatives that operate on different business models. The *Center for Humane Technology* website also has very little to say on the topic; visitors to the site might easily get the impression that no such alternatives exist. That is unfortunate, because alternative models do exist.

A solution to *The Social Dilemma* will require a viable example of an alternative social platform business model, to provide the aspirational pressure that a social media reform movement needs to succeed.
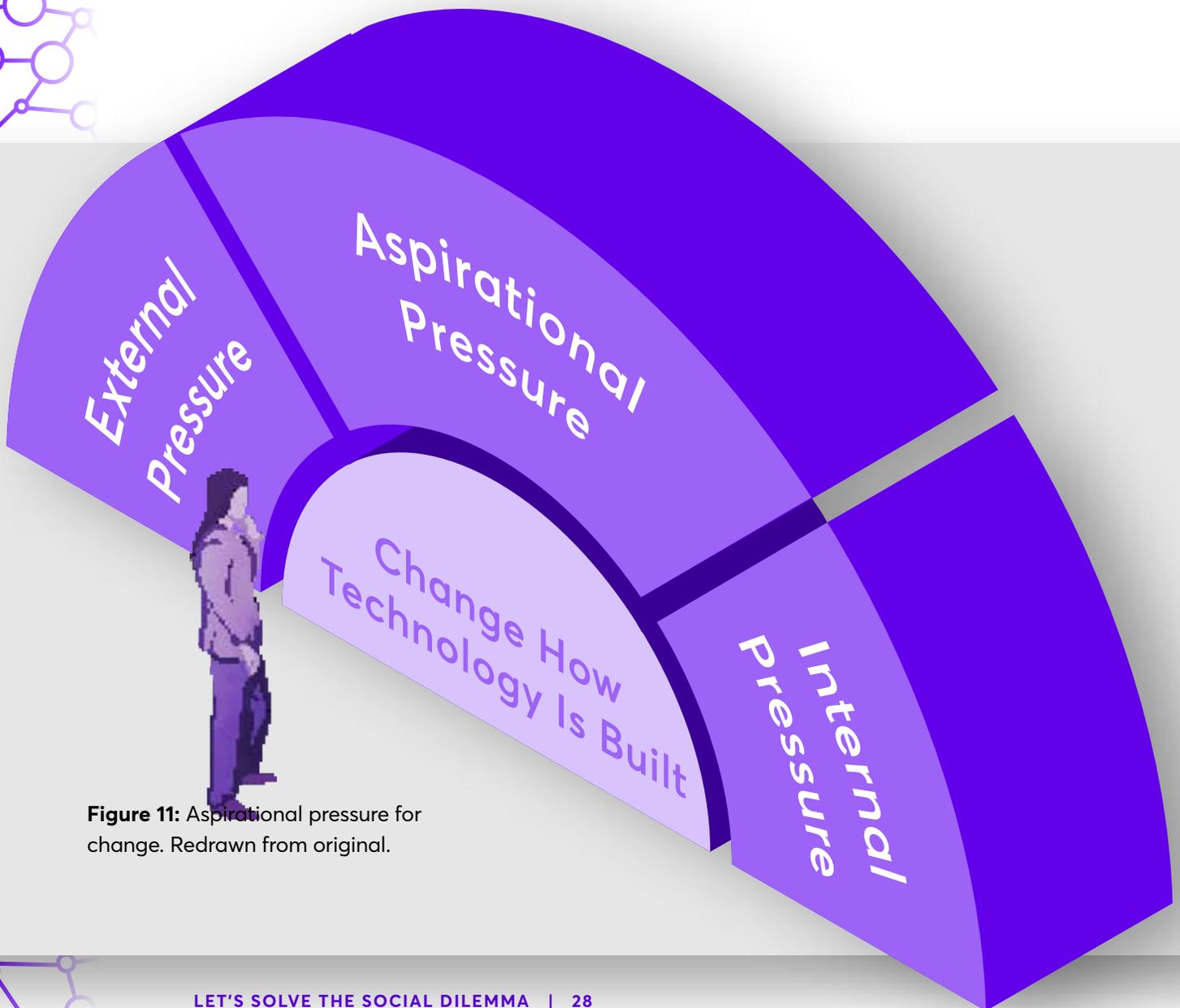
**External Pressure**

**Aspirational Pressure**

**Change How Technology Is Built**

**Internal Pressure**

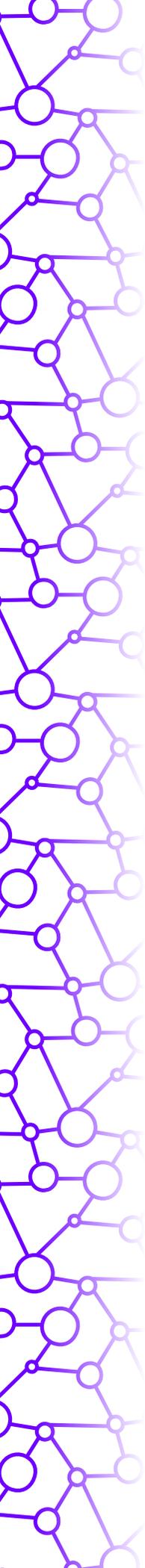**Figure 11:** Aspirational pressure for change. Redrawn from original.

> **A social platform must prioritize the agency, sovereignty and wellbeing of human beings over commercial interests**

# IV What a Viable Solution Would Look Like

A solution to *The Social Dilemma* cannot require users to abandon the many positive features of digital social networks that provide real value to billions of people. Nor can it require social platforms to be run as free public utilities. A viable solution must be allowed to earn a profit that scales as the volume of user activity increases.

But the business model must be radically different from the current one that has powered the Social Media Establishment for the past ten years, that extracts value from user-generated content, and manipulates and sells users as targets of commercial advertising to third parties. *The Social Dilemma* is awakening the public to a dark truth, that the social platforms we use every day are inexorably reshaping society in ways that undermine human sovereignty and wellbeing, and our collective capacity to solve critical problems.

We propose the following two conditions that any viable solution must satisfy.

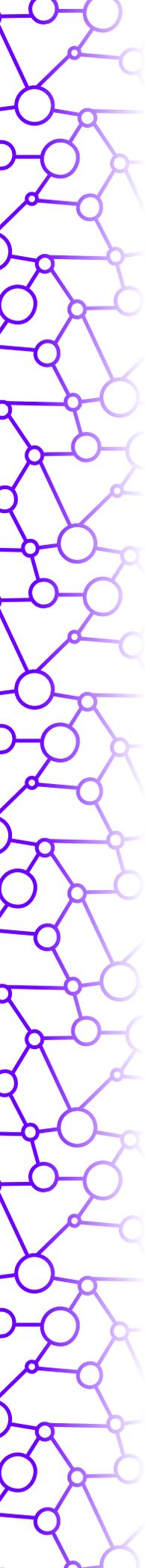## 1. THE *HUMANS FIRST* CONDITION

A social platform must prioritize the agency, sovereignty and wellbeing of human beings over commercial interests.

## 2. THE *VALUE ALIGNMENT* CONDITION

The business model that powers a social platform must function in a way that maintains alignment between the interests of individual users, the broader social communities in which they participate, and the interests of the business entity that develops and manages the platform.

The first condition requires that the platform is anchored in principles of humane design. The second condition requires that the business model reinforces and rewards these principles via a mechanism that ensures that the pursuit of profitability will also promote the freedom, autonomy and genuine well-being of platform communities and their members.

Every Social Media Establishment platform violates both of these conditions right out of the gate.
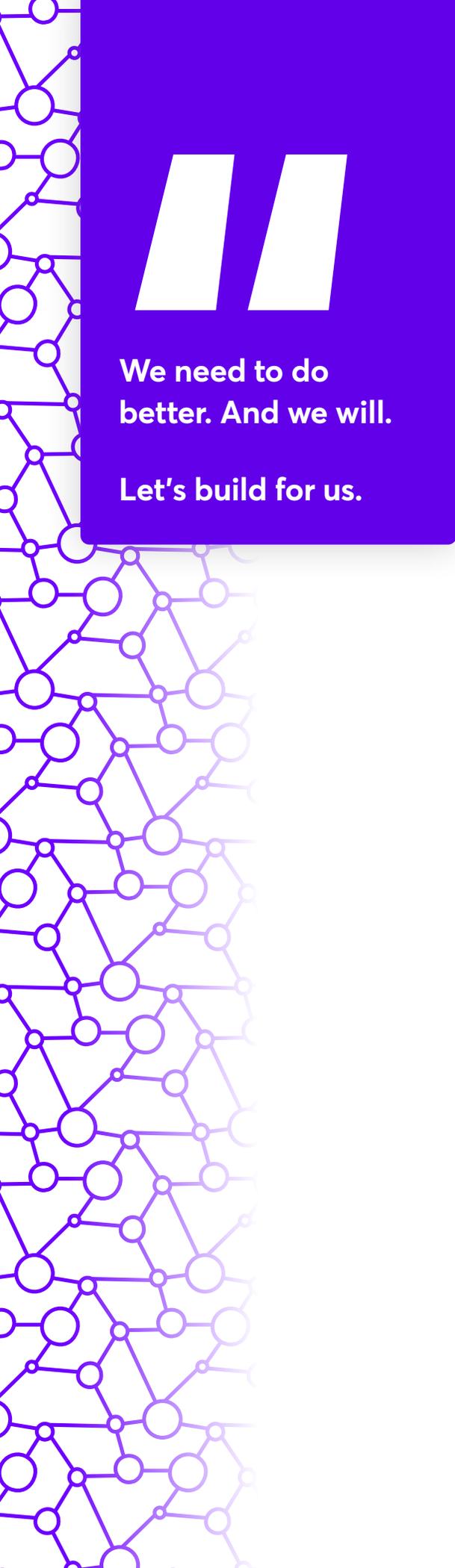
## A SPACE OF VIABLE SOLUTIONS

It is important to note that the *Humans First* condition and the *Value Alignment* condition do not pick out a unique solution. One can imagine a set of solutions that satisfy both conditions. It is a challenge for entrepreneurs, technologists and designers to come up with innovative, practically implementable solutions that meet these ideals.

## FROM IDEAL CONCEPT TO REAL-WORLD CONCEPTION

As presented here, these two conditions represent high-level goals that require articulation and elaboration to specify how they would be implemented in a real platform. The concept of "autonomy", for example, can be analyzed in different ways, and one can imagine different *conceptions* of autonomy being implemented in different platforms and business models. Similarly, one can imagine different conceptions of what counts as "alignment of values" across different domains.

The key point is that, while these principles don't pick out a unique solution, the space of solutions excludes anything remotely resembling the extractive business model that currently drives Facebook, YouTube, Instagram, Twitter, Reddit, etc.

> **We need to do better. And we will.**
>
> **Let's build for us.**

# Conclusion

We have argued that a solution to The Social Dilemma must point to at least one concrete alternative to the Social Media Establishment that employs a business model that (1) puts humans first and (2) maintains alignment between the interests of members, communities and the platform itself. Such alternatives provide the "aspirational pressure" that is needed to support a successful reform movement and provide real options for social media users.

We believe that a decentralized framework that implements a flexible, tokenized value economy is an ideal architecture for developing a platform that fulfills all these goals, but we have not argued that it is the only framework that can do the job.

What is clear is that we cannot stand idle as the Social Media Establishment continues to tear apart the fabric of our society without resistance. The original architects were wrong to pair exponential technology with a business model that exploits and amplifies the worst aspects of human nature; we are
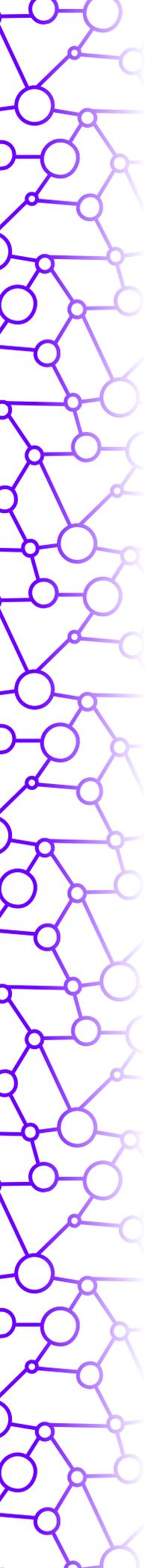
now accelerating towards a dark future that monetizes discord and discontent at great external cost. The precipice draws nearer.

Thankfully, collective human ingenuity can be channeled to fix what's broken. We can learn from our mistakes, course correct, and design our social systems to align with our best interests. We can build a foundation for social media that promotes healthy conversations, preserves digital sovereignty, and maximizes human flourishing.
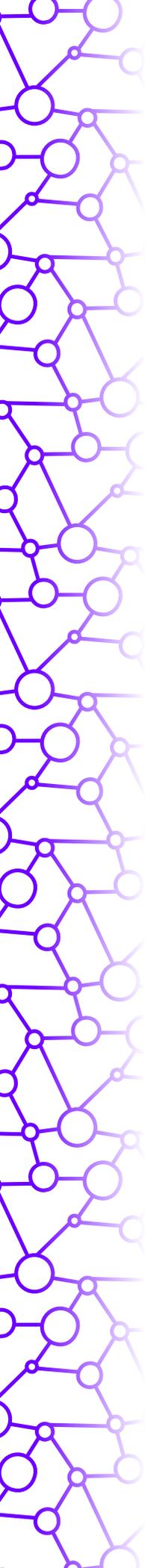
We need to do better. And we will.

Let's build for us.

# Endnotes

*1*      "The Social Dilemma." https://www.thesocialdilemma.com/.

*2*      "Center for Humane Technology." https://www.humanetech.com/.

*3*      "The Age of Surveillance Capitalism - The Guardian." 20 Jan. 2019, https://www.theguardian.com/technology/2019/jan/20/shoshana-zuboff-age-of-surveillance-capitalism-google-facebook.

*4*      "Former Google insider Tristan Harris's dire warning on ...." 17 Jan. 2020, https://www.theaustralian.com.au/weekend-australian-magazine/former-google-insider-tristan-harriss-dire-warning/news-story/1aaac97bb409012bdaad34fe6df7eb87.

*5*      "Opinion | Our Brains Are No Match for Our Technology - The ...." 5 Dec. 2019, https://www.nytimes.com/2019/12/05/opinion/digital-technology-brain.html.

*6*      "Ledger of Harms." https://ledger.humanetech.com/.

*7*      "Tristan Harris: Tech Is 'Downgrading Humans.' It's ... - Wired." 23 Apr. 2019, https://www.wired.com/story/tristan-harris-tech-is-downgrading-humans-time-to-fight-back/.

*8*      "A Path to Humane Technology - Tristan Harris" Nov 14, 2019, https://youtu.be/-oFcGfQ8bWM

*9*      "Take Control - *Center for Humane Technology*." https://www.humanetech.com/take-control.

*10*     "Resources | Humane by Design." https://humanebydesign.com/resources.

11    "A Path to Humane Technology - Tristan Harris" Nov 14, 2019, https://youtu.be/-oFcGfQ8bWM?t=850 [timestamp: 14:10]

12    "Tobacco Chiefs Say Cigarettes Aren't Addictive - The New ...." 15 Apr. 1994, https://www.nytimes.com/1994/04/15/us/tobacco-chiefs-say-cigarettes-aren-t-addictive.html.

13    "1994 - Tobacco Company CEOs Testify Before Congress" https://youtu.be/e_ZDQKq2F08

14    "Hacking the American Mind | Robert Lustig Website." https://robertlustig.com/hacking/.

15    "The 5 biggest little lies tech CEOs told Congress — and us ...." 29 Jul. 2020, https://www.washingtonpost.com/technology/2020/07/29/big-tech-ceo-hearing-lies/.

16    "Facebook forms ethics team to prevent bias in AI software." 3 May. 2018, https://www.cnbc.com/2018/05/03/facebook-ethics-team-prevents-bias-in-ai-software.html.

17    "Facebook announces award recipients of the Ethics in AI ...." https://research.fb.com/blog/2020/06/facebook-announces-award-recipients-of-the-ethics-in-ai-research-initiative-for-the-asia-pacific/.

18    "Our Principles – Google AI." https://ai.google/principles/.

19    "A Path to Humane Technology - Tristan Harris" Nov 14, 2019, https://youtu.be/-oFcGfQ8bWM?t=933 [timestamp: 15:33]