

Holding Platforms Liable *

Xinyu Hua[†]
HKUST

Kathryn E. Spier[‡]
Harvard University

October 16, 2023

Abstract

Should platforms be held liable for the harms suffered by users? A two-sided platform enables interactions between firms and users. There are two types of firms: harmful and safe. The harmful firms impose larger costs on the users. If firms have deep pockets then platform liability is unwarranted. Holding the firms liable for user harms deters the harmful firms from joining the platform. If firms are judgment proof then platform liability plays an instrumental role in reducing social costs. With platform liability, the platform has an incentive to raise the interaction price to deter harmful firms and invest resources to detect and remove harmful firms from the platform. To prevent overinvestment in detection and removal, the residual liability assigned to the platform may be partial instead of full. The optimal level of platform liability depends on the impact on user participation, the intensity of platform competition, and whether users are involuntary bystanders or voluntary consumers.

* We would like to thank Gary Biglaiser, Lus Cabral, Jay Pil Choi, James Dana, Andrei Hagiu, Ørd Harstad, Ginger Jin, Yassine Lefouili, Hong Luo, Bentley MacLeod, Sarit Markovich, Haggai Porat, Urs Schweizer, Emil Temnyalov, Marshall Van Alstyne, Rory Van Loo, Julian Wright and seminar audiences at Boston University, Chinese University of Hong Kong, Columbia University, Fudan University, Georgetown University, Harvard University, the Kellogg School at Northwestern University, Asia Pacific Industrial Organization Conference (APIOC 2021), International Industrial Organization Conference (IIOC 2022), Economics of Platforms Seminar (TSE 2022), the 2022 Asia Meeting of the Econometric Society, American Law and Economics Association Conference (ALEA 2022), Society for Institutional & Organizational Economics Conference (SIOE 2022), JRC-TSE Workshop on Liability in the Digital Economy (2022), Society for the Advancement of Economic Theory Conference (SAET 2022). We also thank the support from the Hong Kong Research Grants Council (GRF Grant Number: 16500722).

[†]Hong Kong University of Science and Technology. xyhua@ust.hk.

[‡]Harvard Law School and NBER. kspier@law.harvard.edu.

1 Introduction

Online platforms are ubiquitous in the modern world. We connect with friends on Facebook, shop for products on Amazon, and search online for jobs, information, and entertainment. While the economic and social benefits created by platforms are undeniable, the costs and hazards for users are very real too. For example, platform users run the risk that their personal data and privacy will be compromised. Users of social networking sites and search engines may be misled by fraudulent advertisements and misinformation. Consumers who shop online run the risk of purchasing counterfeit, defective, or dangerous goods. Should internet platforms like Facebook and Amazon be liable for the harms suffered by users?

In the United States, platforms enjoy relatively broad immunity from lawsuits brought by users, although this immunity is being challenged in legislatures and the courts. Section 230 of the Communications Decency Act, enacted in 1996, shields platforms from liability for the digital content created by their participants.² Early proponents argued that the law was necessary to allow the internet to grow and flourish, but its application is controversial and many critics question the law's merits.³ In 2019, Facebook paid \$5 billion to settle charges that they failed to take adequate precautions to protect user data.⁴ The FTC has also been investigating how "platforms screen for misleading ads for scams and fraudulent and counterfeit products" and, "in 2022 alone, consumers reported losing more than \$12 billion to fraud that started on social media, more than any other contact method."⁵ Proposed federal legislation would hold platforms liable if they fail to protect users.⁶

Marketplace platforms have largely avoided responsibility for defective products and services sold by third-party vendors. In 2019 the Fourth Circuit held that Amazon.com is not a traditional seller and therefore not subject to strict tort liability.⁷ The following

¹See Buiten et al. (2020) for discussion of the European Commission's e-Commerce Directive. Hosting platforms in the EU may avoid liability for illegal content posted by users, assuming they are not aware of it, and are not responsible for monitoring the legality of the posted content.

²Section 230(c)(1) says that "No provider or user of an interactive computer service shall be treated as the publisher or speaker of any information provided by another information content provider." Proponents hoped Section 230 would address the "perverse incentives" created by

year, a California court found that Amazon could be held strictly liable for a defective laptop battery that was sold by third-party vendors but "Fulfilled by Amazon."⁸ Then, in 2021, Amazon was held strictly liable for harms caused by a defective hoverboard that was shipped directly to the consumer by an overseas third-party vendor. Although Amazon did not fulfill the hoverboard order, the court opined that Amazon was "instrumental" in its sale and that "Amazon is well situated to take cost-effective measures to minimize the social costs of accidents."⁹ In short, the law is far from settled.

Such settings include social and professional networking platforms such as Facebook and LinkedIn where the users enjoy same-side network benefits from sharing content with each other and the firms pay the platform to access user data or to engage in unintentional activities (e.g., advertising). Platform users may be harmed by the firms when their private data is breached or when they are exposed to harmful advertising or misinformation. Absent liability the harmful firms have no incentive to leave the platform, and the platform has an insufficient incentive to detect and remove them. Holding the firms and the platform jointly liable gets them to internalize the negative externalities on the user-bystanders.

If the firms have deep pockets, and must pay in full for the harms they cause, then platform liability is unwarranted. Holding just the firms liable achieves the first-best outcome. Platform liability is socially desirable when the firms are *judgment proof* and immune from liability.¹⁷ First, if the platform is held liable, the platform will raise the interaction price for the firms to reflect the platform's future liability costs. If the harmful firms are "marginal" (i.e., the harmful firms have a lower willingness to pay than the safe firms) then the higher interaction price deters the harmful firms from joining the platform. Second, if the harmful firms are "inframarginal" and undeterrable, the platform will invest resources to detect and remove the harmful firms from the platform.¹⁸ Interestingly, the optimal level of platform liability may be partial instead of full, as full liability could lead to excessive auditing by the platform.¹⁹

We then consider the more general setting with *heterogeneous* users where some join the platform and others do not. We show that platform liability has the added benefit of stimulating user participation. This happens for two reasons. First, users anticipate that the platform's auditing incentives are improved and that the platform is safer. Second, users view the larger damage award as a "rebate" for joining the platform. Because of the *user-participation effect* the optimal platform liability is higher than in the baseline model.

Next, we extend the baseline model to settings where *users are customers* of the firms, so interactions require the users' consent. Relevant settings include online marketplaces like eBay and Amazon where participants enjoy cross-side benefits from the sale of goods and services. As in the baseline model there are two types of sellers, harmful and safe. The harmful sellers have lower production costs but cause harms more frequently. The consumers are sophisticated and their willingness-to-pay reflects their rational expectations about product risks. The risk of harmful products depresses the price that consumers are willing to pay and, by extension, depresses the revenues that the platform can generate.

¹⁷Shavell (1986) provides the first rigorous treatment of the judgment proof problem, where injurers with limited assets tend to engage in risky activities too frequently and take too little care.

¹⁸If the firms are very judgment proof and can evade liability, then the harmful firms are inframarginal (i.e. the harmful firms have a strictly higher willingness to pay than the safe firms). If the firms are moderately judgment proof, then the harmful firms are "marginal."

¹⁹If the firms are very judgment proof, then the safe firms are marginal and the harmful firms get information rents. When choosing its audit intensity, the platform does not take into account the lost

If the harmful rms are marginal, then platform liability is unwarranted. Since consumers are willing to pay more for safer products, the platform has a private incentive to raise the interaction price to deter the harmful rms from joining the platform. If the harmful rms are inframarginal, however, then partial platform liability gives the platform an appropriate incentive to audit and remove the harmful rms.²⁰ Since the platform internalizes the average harm to consumers, the socially-optimal platform liability is lower than in

the basic platform (e.g., for social platforms).

platforms (e.g., for social platforms) are not as basic as the basic platform. The platform's liability is lower than in the basic platform because the platform internalizes the average harm to consumers, which gives it an incentive to audit and remove harmful rms. This is particularly true for social platforms, where the platform's liability is lower than in the basic platform because the platform internalizes the average harm to consumers, which gives it an incentive to audit and remove harmful rms.

platforms are not as basic as the basic platform. The platform's liability is lower than in the basic platform because the platform internalizes the average harm to consumers, which gives it an incentive to audit and remove harmful rms. This is particularly true for social platforms, where the platform's liability is lower than in the basic platform because the platform internalizes the average harm to consumers, which gives it an incentive to audit and remove harmful rms.

platform liability

platform liability is lower than in the basic platform because the platform internalizes the average harm to consumers, which gives it an incentive to audit and remove harmful rms. This is particularly true for social platforms, where the platform's liability is lower than in the basic platform because the platform internalizes the average harm to consumers, which gives it an incentive to audit and remove harmful rms.

extending liability to an injurer's lenders²⁵ and Dari Mattiacci and Parisi (2003) consider vicarious liability where liability is extended to the injurer's employer.²⁶ Arlen and MacLeod (2005a) show that holding managed care organizations liable for medical malpractice by their physicians can raise the physicians' incentives to take care. Our model, which has not been previously studied, investigates the design of platform liability when the platform can audit and remove harmful participants.²⁷

There is a vast literature on multi-sided platforms. The early studies (e.g., Caillaud and Jullien, 2003; Rochet and Tirole, 2003, 2006; Armstrong, 2006; and Weyl, 2010) have identified how cross-side externalities affect platform pricing schemes and users' participation incentives. The literature also examines the impact of seller competition²⁸ or the impact of platform competition on pricing.²⁹ Some recent studies pay attention to non-pricing strategies, including seller exclusion (Hagiu, 2009), information management (Julien and Pavan, 2019; Choi and Mukherjee, 2020), control right allocation (Hagiu and Wright, 2015, 2018), and platform governance (Teh, 2022).

There is a small but growing literature on platform liability. The policy papers by Buiten et al. (2020) and Lefouili and Madio (2022) discuss informally whether platforms should bear liability for harms caused by participants. A few working papers study copyright infringement and retail settings. De Chiara et al. (2021) examine hosting platforms' incentives to filter copyright-infringing materials. They focus on harms to copyright owners and do not consider platforms' pricing strategies. Jeon et al. (2022) examine how negligence-based liability changes platforms' incentives to remove IP-infringing products, which in turn affects brand owners' innovation incentives. Zenny (2023) considers the impact of platform liability on sellers' efforts to improve product safety, without discussing platforms' screening or auditing actions. Yasui (2022) discusses sellers' incentives to maintain reputation and platforms' ex-post efforts to discover and announce potential safety risks after consumers purchase products from sellers. Our paper considers a broad array of platform types and investigates the effects of liability on platform pricing, incentives to block bad actors, and social welfare.

Our paper is organized as follows. Section 2 presents the baseline model where users are homogeneous bystanders of the firms. Section 3 generalizes the baseline model by considering heterogeneous users with endogenous participation. Section 4 examines sev-

²⁵See also Boyer and Laont (1997) and Che and Spier (2008). Bebchuk and Fried (1996) argue informally for raising the priority of tort victims in bankruptcy above debt claims gives the debtholders an incentive to better monitor the borrower.

²⁶There are related legal studies. See Kraakman (1986) for a general taxonomy of gatekeeper enforcement strategies, Hamdani (2002) for liability on internet service providers, Hamdani (2003) on accountants and lawyers, and Van Loo (2020a) on big technology.

²⁷Our paper is also related to the studies comparing joint and several liability (JSL) to several liability (SL) for harms caused by multiple defendants (e.g., see Landes and Posner, 1980; Carvell et al., 2012). With JSL, the victim may recover full damages from a single deep-pocketed defendant. With SL, the victim's recovery from each defendant is limited by the defendant's share of responsibility.

²⁸See Nocke et al. (2007), Galeotti and Moraga-Gonzalez (2009), Hagiu (2009), Gomes (2014), Belleamme and Peitz (2019).

²⁹See Dukes and Gal-Or (2003), Hagiu (2006), Armstrong and Wright (2007), White and Weyl (2010), Karle et al. (2020), Tan and Zhou (2021).

eral extensions including a retail setting where the firms are sellers and the users are consumers and a setting with two competing platforms. Section 5 provides concluding thoughts. The proofs are in the appendix.

2 The Baseline Model

Consider a two-sided platform (P) with two kinds of participants, firms (S) and users (B). The platform is a monopolist and necessary for interactions between firms and users. Firms and users are small, have outside options of zero, and the mass of each is normalized to unity.

The platform provides two goods. First, the platform provides a quasi-public good that gives each user a private benefit $v > 0$. For simplicity, we first consider the special case where users are homogeneous and have the same valuation v . Section 3 generalizes the analysis to include heterogeneous users with different valuations. Second, the platform provides opportunities for the firms and the users to interact.

We assume that interactions between firms and users do not require the users' consent and so the users are effectively bystanders.³⁰ The benefits and costs of these interactions depend on the firms' type, $i \in \{H, L\}$, where α is the mass of type H and $1 - \alpha$ is the mass of type L in the firm population.³¹ The H -type firms have higher interaction benefits, $b_H > b_L$, but impose higher interaction losses on users, $d_H > d_L$ where $\alpha \in [0, 1]$ is the probability of harm and $d > 0$ is the level of harm per firm-user interaction.³² The firms privately observe their types.

We assume that the platform charges the firms a price ϕ per interaction and allows users to join the platform for free. This is broadly aligned with what we often observe in practice. Platforms such as Google and Facebook monetize quasi-public goods by selling online advertising to businesses and/or sharing user data and do not charge users for access. In theory, this pricing strategy can be very profitable for the platform in strategic environments with strong network effects.³³ Our assumption is also aligned with other papers in the platform literature.³⁴

The platform has the capability to detect and block the H -type firms. We will refer to the platform's efforts to detect the H -types as auditing. By virtue of their scale, data,

³⁰Section 4.1 extends the analysis to retail platforms where interactions require the users' consent.

³¹For simplicity, α is taken as exogenous. One may endogenize by allowing firms to invest resources to increase the likelihood being safe. Section 4.3 discusses an extension with firm moral hazard problems.

³²If $\alpha < \alpha^*$ then the H -types are marginal for all liability rules and auditing is unnecessary. The threshold α^* defined in (5) below is identically equal to zero, and all of our results apply.

³³Suppose that each user receives only if a large number of users join the platform. A user's decision to join depends on the price and their expectations about the number of other users. Following Harsanyi and Selten (1988), to avoid coordination failure, the platform should set a sufficiently low price (or even zero price) for the users. The appendix provides an illustrative example of the coordination game.

³⁴Armstrong (2006) shows that, with strong network effects, platforms have incentives to set negative prices. However, negative prices may be infeasible. Armstrong and Wright (2007) and Choi and Jeon (2021) justify non-negative prices on adverse selection and moral hazard grounds. Gans (2022) justifies this based on free disposal. See also Belleflamme and Peitz (2021).

and technological sophistication, platforms like Google may be in a good position to root out harmful platform participants.³⁵ Specifically, by spending effort $e \in [0, 1)$ per firm, the platform can detect H -type firms with probability

to join the platform.⁴³ A2 guarantees that the platform always gets non-negative profits

2.1 Motivating Examples

In our baseline model, bad actors on one side of a two-sided platform may harm users on the other side of the platform. In the following, we motivate the baseline model with three broad examples: fraudulent advertising, data misuse by technology partners, and the sale of harmful products. For each of these three settings, we will document the platform's financial incentives, the presence of bad actors, and the potential for user harm.

Advertisers. Many platforms rely on paid advertising as their main source of revenue. Although most online advertising is benign, fraudulent and misleading ads are a significant problem. (partners,)-29bb

ing app developers, who can use the data to improve their offerings and the experiences of platform users. When deciding whether to grant developers access to user data, platforms may consider the financial benefit (among other things). For example, in 2013, Facebook allegedly granted or denied access based on the developer spending at least \$250,000 on mobile advertising.⁵¹ In the Spring of 2023, Twitter, Reddit, and other platforms announced hefty charges for developers to access the platforms' API, leading some partners to reduce their data usage and others to terminate their contracts.⁵² There is substantial evidence that technology partners violate their platform developer agreements⁵³ and sell data to others.⁵⁴ In the wrong hands, platform data can be used for identity theft.⁵⁵ 8536 Td [(da

2.2 Equilibrium Analysis

In this subsection, we characterize the platform's pricing and auditing strategies, p and e , given the assignment of liability, w_s and w_p . A type- i firm will seek to join the platform when their expected profit per interaction is non-negative,

$$\pi_i - w_s - p \geq 0; \quad (4)$$

where π_i is the firm's interaction benefit, w_s is the firm's expected liability, and p is the price paid to the platform. Note that depending on the level of firm liability, w_s , the H -type may have higher or lower rents than the L -type. The rents of the two types are equal when

$$w_s = \bar{w} = \frac{H}{H+L} < d; \quad (5)$$

The threshold \bar{w} defined in (5) is critical for understanding the impact of platform liability on the interaction price and audit intensity. If the firms are sufficiently judgment-proof, $w_s < \bar{w}$

We now explore the platform's incentive to audit and remove the θ -type firms. The platform's aggregate profits are:

$$\pi(\theta) = (1$$

- (a) If $r_H(d, w) > r_H(w_s)$ then $0 < e < e^*$.
- (b) If $r_H(d, w) = r_H(w_s)$ then $0 < e = e^*$.
- (c) If $r_H(d, w) < r_H(w_s)$ then $0 < e < e^*$.

To summarize, when firm liability is below the threshold, $w_s \leq \bar{w}$, the H -type firms cannot be deterred from joining the platform by the interaction price p . The platform invests in auditing if and only if the joint platform-firm surplus is larger than the firms' information rent. Note that the platform's incentives to audit are stronger when w_p and w_s are larger. The platform's incentive to audit and remove the H -types is socially insufficient when the joint liability for the platform and firms is small (as in case 2(a)) but socially excessive if the joint liability is large (as in case 2(c)).

Case 2: $w_s > \bar{w}$. Now suppose that firm liability is above the threshold, so the H -type firms are marginal. The platform's profit-maximizing strategy is to either charge $p_L = \frac{1}{2} w_s$ and deter the H -types from joining the platform or charge $p_H = \frac{1}{2} w_s < p_L$ and attract both types. Notably, if the platform chooses the latter strategy, then it will not invest in auditing, $e = 0$.⁶²

The platform will charge p_H and attract the H -types if and only if

$$(p_H - w_p) + (1 - \alpha)(p_H - \frac{1}{2} w_p) > (1 - \alpha)(p_L - \frac{1}{2} w_p):$$

Substituting the formulas for p_H and p_L and using the definition of \bar{w} in equation (5) this condition becomes:

$$(\frac{1}{2} w_s - w_p) > (1 - \alpha)(\frac{1}{2} w_s - \frac{1}{2} w_p): \tag{11}$$

The left-hand side is the joint platform-firm surplus of attracting the H -type firms on the platform: the fraction α of H -types

type in 470 platform 52 Tt 84Tt 3

on becomes:

2.3 Platform Liability

This subsection explores the social desirability and optimal design of platform liability for harm to user-bystanders, taking the level of firm liability w_s as fixed. We begin by presenting a benchmark where the platform is not liable for the harm $w_p = 0$.

Proposition 1. (*Firm-Only Liability.*) Suppose that the platform is not liable for harm to users, $w_p = 0$, and firm liability $w_s \in [0; d]$. There exists a unique threshold $w_s^* \in (0; \frac{H}{H-L})$, where $\frac{H}{H-L}$ is weakly increasing in α , such that:

1. If $w_s < w_s^*$ then the platform sets $p = p_L = p_L w_s$, attracts the L -type firms, and does not invest in auditing, $e = 0 < e^*$. The platform's auditing incentives are socially inefficient.
2. If $w_s \in (w_s^*; \frac{H}{H-L})$ then the platform sets $p = p_H = p_H w_s$, attracts the H -type firms, and does not invest in auditing, $e = 0 < e^*$. The platform's auditing incentives are socially inefficient.
3. If $w_s > \frac{H}{H-L}$ then the platform sets $p = p_L = p_L w_s$ and deters the H -type firms. The first-best outcome is achieved.

Should platforms be held liable for the harm suffered by users? Proposition 1 establishes that platform liability is unnecessary when the firms themselves are held sufficiently liable for harm to the users (case 3 in Proposition 1). In this case, the joint platform-firm surplus of including the H -types is low, so the platform has incentives to deter them by charging a high price. However, when the firms are more judgment proof and the platform faces no liability (cases 1 and 2 in Proposition 1), the private and social incentives diverge. The platform attracts the H -types and does not invest in costly auditing. In such cases, platform liability can be socially desirable, as shown in the next proposition.

Proposition 2. (*Optimal Platform Liability.*) Suppose firm liability $w_s \in [0; d]$. The socially optimal platform liability for harm to users is as follows:

1. If $w_s < w_s^*$ then $w_p = d - w_s - 1 \frac{L}{H} (w_s) \in (0; d - w_s]$ achieves the second-best outcome. The platform sets $p = p_L = p_L w_s$ and attracts the L -type firms. The platform's auditing incentives are socially efficient.
2. If $w_s \in (w_s^*; \frac{H}{H-L})$ then there exists a threshold $w_p^* > 0$

$f(v) > 0$ for $v \in [0, 1)$, with cumulative density $F(v)$.⁶⁴ As in the baseline model, the platform charges the firms price p per interaction and takes auditing effort e per firm. Note that there are economies of scale in (per-firm) auditing, so that both the private and the socially optimal incentives for auditing depend on the users' participation rate.⁶⁵ Users have the option to join the platform for free. As discussed in the baseline model, many platforms do not charge users in practice and this observation could emerge in equilibrium when there are strong same-side or cross-side network effects.⁶⁶

We assume that the users cannot directly observe the platform's audit intensity, or equivalently, the platform chooses its audit intensity after the users make their participation decisions.⁶⁷ Although the users do not observe the platform's auditing effort e when making their participation decisions, they observe the liability rule, W_s and W_p , and form correct beliefs about e in equilibrium.

In practice, the public does not directly observe platforms' enforcement efforts or technologies used in improving platform safety. In the words of former Facebook employee and whistleblower Frances Haugen, "Facebook became a \$1 trillion company by paying for its profits with our safety, including the safety of our children" and "almost no one outside of Facebook knows what happens inside Facebook."⁶⁸ The Digital Services Act in the European Union and the PACT Act recently proposed in the US contain many disclosure requirements,⁶⁹ which reflects lawmakers' concerns about the lack of transparency on platform safety and effort.⁷⁰

Consider the first-best outcome. Assumption A2 implies that it is socially efficient for

not join the platform and all the users join the platform.

Next, consider the second-best outcome. As in the baseline model, full deterrence of the H -types may not be possible. If the H -type users seek to join the platform, then costly auditing is necessary to reduce the social harm. In the second-best benchmark, social welfare is

$$S(e; \mathbf{w}) = \int_{\mathbf{w}} [v + (1 - e)(\theta_H - \theta_H d) + (1 - e)(\theta_L - \theta_L d)] f(v) dv - c(e); \quad (12)$$

where \mathbf{w} is the value of the marginal user,

$$\mathbf{w}(e; w) = (\theta_H - e)(\theta_H + \theta_L)(d - w); \quad (13)$$

Notice that $\mathbf{w}(e; w)$ is decreasing in e and w for all $d - w > 0$: higher levels of effort and liability stimulate user participation. Holding e constant, the users view w as a "rebate" for joining the platform. Therefore, the social planner would like to set $w = d$ (that is, $w_p = d - w_s$), so that all the users participate. Given full participation by the users, the socially efficient auditing effort is e^* , the same as in the baseline model.

We now characterize the equilibrium and the optimal platform liability. As in the baseline model, the L -type users are marginal if $w_s \leq \mathbf{w}$, while the H -types are marginal if $w_s > \mathbf{w}$. We consider each case in turn.

Case 1: $w_s \leq \mathbf{w}$. In this case, the L -type users are marginal and the platform charges $p^u = \theta_L - \theta_L w_s$. The platform's profit function may be written as:

$$\begin{aligned} \pi(e; \mathbf{w}) = & S(e; \mathbf{w}) + \int_{\mathbf{w}} (\theta_H - e)(\theta_H - \theta_L)(\mathbf{w} - w_s) \\ & + ((\theta_H - e)\theta_H + (\theta_H - e)\theta_L)(d - w) - \int_{\mathbf{w}} v f(v) dv; \quad (14) \end{aligned}$$

where \mathbf{w} is the marginal user defined in (13). Since the platform chooses its auditing effort e^u (if it is positive) satisfies⁷¹

$$\begin{aligned} \frac{\partial \pi(e^u; \mathbf{w})}{\partial e} = & \frac{dS(e^u; \mathbf{w})}{de} + \int_{\mathbf{w}} [(\theta_H - \theta_L)(\mathbf{w} - w_s) - \theta_H(d - w)] f(v) dv \\ & + \theta_H(d - w) \frac{\partial S(e^u; \mathbf{w})}{\partial \mathbf{w}} = 0; \quad (15) \end{aligned}$$

where

$$\frac{\partial S(e^u; \mathbf{w})}{\partial \mathbf{w}} = \frac{\partial \pi(e^u; \mathbf{w})}{\partial \mathbf{w}} - (\theta_H - e^u)(\theta_H - \theta_L)(\mathbf{w} - w_s) f(\mathbf{w}); \quad (16)$$

Equation (15) shows that the platform's auditing incentives diverge from the social planner's. The first line of equation (15) is familiar. As in the baseline model, when the platform increases e , the removed H -types lose their information rents, $(\theta_H - \theta_L)(\mathbf{w} - w_s)$

⁷¹See the proof of Proposition 3.

and the users' uncompensated loss is reduced by $H(d - w)$. If $w_p = w_p$ as defined in Proposition 2, these two effects offset each other. The last line of equation (15) identifies a new source of divergence: the platform's

2. If $w_s = d$ then $w_p^u = d$ w_s achieves the second-best outcome. The platform sets $p^u = \frac{L}{L + w_s}$ and chooses the efficient auditing effort $e^u = e^*$. All users participate.
3. If $w_s > d$ then $w_p^u = d$ w_s achieves the first-best outcome. The platform sets $p^u = \frac{L}{L + w_s}$ and deters the H-type firms. All users participate.

To summarize, as in the baseline model, if the firms have deep pockets and can be held fully liable ($w_s = d$), platform liability is unnecessary. However, if the firms are judgment proof, platform liability can motivate the platform to take more auditing effort or raise the interaction price, which removes or deters the harmful firms. Additionally, platform liability stimulates user participation. So, the optimal level of platform liability is weakly higher than in the baseline model. Note that, when the firms are very judgment proof (Case 1 in Proposition 3), the optimal platform liability leads to excessive auditing.

Remark

Lawmakers and commentators have historically expressed concern that the burden of liability might chill economic activity. These concerns were part of the rhetoric for platform immunity to liability in the early years. Section 230 of the Communications Decency Act was adopted to allow the internet to grow and flourish.⁷⁵ To be sure, defending against frivolous lawsuits can be costly and distract managers from the core business.⁷⁶ However, our analysis shows that platform liability can stimulate user participation, both directly and indirectly.⁷⁷

First, platform liability serves as a "rebate" to attract users. This effect is unique to the platform market. To see this, consider a non-platform market where a seller sells its product to consumers. Although products liability reduces consumers' uncompensated harm, it raises the seller's costs and leads the seller to raise the price of the product, which can neutralize the impacts on output.⁷⁸ By contrast, in a platform market with strong network effects, users have the option to join the platform for free. The platform does not adjust the price to fully reflect the users' uncompensated harm or the platform's liability costs and cross-side network benefits (i.e. revenue from the firms). Platform liability stimulates participation by reducing the "effective" price for users.

Second, platform liability can raise the audit intensity, which attracts users indirectly. For both platform and non-platform markets, when users cannot observe product safety or the platform's audit intensity, liability addresses the moral hazard problem and improves

⁷⁵Section 230 and has been called "the one line of federal code that has created more economic value in this country than any other." See <https://www.npr.org/sections/alltechconsidered/2018/03/21/591622450/section-230-a-key-legal-shield-for-facebook-google-is-about-to-change>.

⁷⁶Court errors and litigation costs are discussed in Section 4.3.

⁷⁷Some empirical studies observe a positive correlation between liability and innovation. Viscusi and Moore (1993) observe that when products liability is low or moderate, raising liability encouraged firms' investments in innovation. Galasso and Luo (2017) identify a positive correlation between liability and innovation.

⁷⁸If consumers have the same preference for product safety and can observe safety before purchase, then liability is irrelevant to output (Hamada, 1976).

safety. The increased safety reduces the joint costs for the platform and users (or the seller and consumers), thereby stimulating user participation. However, the moral hazard problem is not the only reason for the divergence between the platform's auditing incentive and the social incentive. As shown by equation (15), the divergence occurs also because the platform does not consider the benefit of auditing to the inframarginal users or the impact of increased participation on the firms' rents. Platform liability addresses these externalities and motivates the platform to raise audit intensity.

Remark Platform liability may be socially beneficial when users observe the platform's auditing effort e before making their participation decisions. In this setting, the platform's auditing incentives are stronger. Recall that in equation (15), when effort is not observable, the platform disregards the social benefit of increased participation (the last term). With observable effort and $w_s = 1$, the platform's effort (if it is positive) satisfies:

$$\frac{d(e^U; \tau)}{de} = \frac{dS(e^U; \tau)}{de} + \int_{\tau}^Z [(H - L)(\tau - w_s) - H(d - w)] f(v) dv + H(d - w) \frac{\partial S(e^U; \tau)}{\partial \tau} - \frac{\partial (e^U; \tau)}{\partial \tau} = 0: \quad (17)$$

where

$$\frac{\partial S(e^U; \tau)}{\partial \tau} \quad \frac{\partial (e^U; \tau)}{\partial \tau}$$

term in equation (17) drops out and we are left with $d(\dots) = de = dS(\dots) = de \int_{\nu}^R H(d$
 $w) f(\nu) d\nu = 0$. Private and social incentives diverge because the platform does not con-
sider the safety benefits that accrue to the participating users. Imposing full residual
liability on the platform,

the platform, but are sophisticated and form beliefs that are, in equilibrium, correct.⁸⁵ If the H -type firms seek to join the platform and the platform invests e in auditing, the conditional probability of harm per interaction is

$$E(j|e) = \frac{(1 - e) \theta_H + (1 - \theta_H) e}{(1 - e) \theta_H + (1 - \theta_H) e}; \quad (19)$$

which is a decreasing function of e .

rents from the L -type firms, $p^r = t^r (c_L w_s + c_L)$. Using (20) and $c_L = 0 - c_L$,

$$p^r = c_L - c_L w_s + \tau^r (d - w) \quad (21)$$

Comparing p^r to its counterpart p (see (6)) in the baseline model reveals an important difference: the interaction price paid by the firms (21) reflects the user-consumers' expected uncompensated harm, $\tau^r (d - w)$.

We now explore the platform's auditing incentives. Substituting p^r from (21), $S(e)$ from (2), and w from (5) into (7) gives the platform's profit function

$$\begin{aligned} \pi(e) = S(e) - v - (1 - e) (c_H - c_L) (w - w_s) \\ + [(1 - e) (c_H - c_L) + (1 - e) (c_L - c_H)] (d - w) \end{aligned} \quad (22)$$

The platform's profits $\pi(e)$ diverge from social welfare $S(e)$ for two reasons:

incremental social benefit of attracting the H

By contrast, when users are consumers, the retail price t^r paid by the users to the rms (and the price p^r paid by the rms to the platform) reflects the users' beliefs of the probability of harm. In Proposition 4, when the users are consumers, w_p^r satisfies

$$(w_H - w_L)(w_s) = (d - w_s - w_p^r) \quad (26)$$

Now the right-hand side reflects the users' uncompensated harm beyond their expectations. As in the baseline model, when rm liability w_s rises, both sides fall. However, the drop in the rms' rent on the left is *bigger* than the drop in the users' uncompensated harm (beyond their expectations) on the right. Holding w_p^r fixed, the platform would invest *too little* auditing. To restore the efficient incentives for auditing, platform liability should be raised. This is why platform liability and rm liability are complements in the retail platform extension.

Corollary 1. *Suppose $w_s > 0$. When the users are bystanders, the optimal platform liability decreases in w_s ; when the users are consumers, the optimal platform liability increases in w_s .*

Remark. The analysis above assumed that the platform removed discovered H -types from the platform. What would happen if the platform is required to disclose the audit results to the consumers, and the consumers decide for themselves whether to interact with the known H -types? Absent platform liability ($w_p = 0$), a rational consumer would decline to interact with a known H -type ex post.⁹⁴ Although ex post efficiency would be obtained without platform liability, the platform would have insufficient incentives to audit the sellers ex ante.⁹⁵ At the other extreme, with full platform liability ($w_p = d$), a rational consumer would interact with a known H -type.⁹⁶ That is, disclosure would not deter harmful interactions. These observations underscore the importance of granting retail platforms the discretion to remove bad actors rather than relying on disclosure alone.⁹⁷

4.2 Platform Competition

We now extend our baseline model (with user-bystanders) by considering two competing platforms, Platform 1 and Platform 2. Users are distributed symmetrically on a Hotelling

⁹⁴The joint surplus for a consumer and an H -type rm from their transaction is $\frac{1}{2}(w_H - w_L)$.

line with density $g(x) = g(1-x) > 0$ on $x \in [0; 1]$, Platform 1 is located at $x = 0$ while Platform 2 is located at $x = 1$. A user at location $x \in [0; 1]$ receives consumption value $v - x$ if they join Platform 1 but $v - (1 - x)$ if they join Platform 2, where $v > 0$ reflects the level of differentiation. Assume that v is sufficiently large such that the market is fully covered. The firms can join both platforms, while each user only joins one platform.⁸⁸ Thus, the platforms compete for users but not for firms.⁸⁹

In stage 1, the platforms set their prices simultaneously. The timing and the other assumptions are otherwise identical to the baseline model. Denote the platforms' prices and auditing efforts as p_j and e_j , $j = 1; 2$. We shall focus on the symmetric equilibrium where $p_1 = p_2$ and $e_1 = e_2$ and, accordingly, each platform serves half of the users. We will show that platform liability can still be socially beneficial in this competitive environment.

Case 1: $w_S = w$. In this case, the L -type firms are marginal and the platforms set $p_1 = p_2 = p_L = p_L w_S > 0$. Although the users do not observe the platforms' auditing efforts directly, they are sophisticated and form rational inferences in equilibrium. In

Now suppose $w_s \geq 2d$ (i.e.). If $w_p = d$, the users would be fully compensated for any harm and therefore each platform attracts half of the users. Each platform charges p_L if

$$\frac{1}{2}(1 - \alpha)(p_L - \alpha w_p) > \frac{1}{2}[(1 - \alpha)(p_H - \alpha w_p) + (1 - \beta)(p_H - \beta w_p)];$$

which holds given $p_H < p_L$ and $p_H - \alpha w_p = p_H - \alpha d < 0$. Hence, imposing full residual liability on the platforms gets the platforms to raise the interaction price and deter the H

3. If w_s is judgment proof, platform liability is unnecessary. The platform sets $L = w_s$ and deters the L -type firms.

Comparing Proposition 5 to Proposition 2 reveals how competition changes the socially-optimal level of platform liability. If the firms are very judgment proof, $w_s = 0$, then the socially-optimal level of platform liability is the same as for monopoly $w_p^C = w_p$. As before, platform liability encourages the platforms to detect and remove the L -type firms from the platforms. If the firms are modestly judgment proof, $w_s > 0$, then platform liability is socially beneficial when the platforms are sufficiently differentiated (large γ) but unnecessary when platform competition is fierce (small γ). By contrast, in the baseline model, platform liability was necessary to induce the platform to raise the interaction price to deter the bad actors. Here, when competition is fierce, the market mechanism gives the platforms the incentive to raise their interaction prices and deter the bad actors from participating.

Regulators across the globe have been focusing efforts on increasing competition and reducing market power in platform markets. For example, the Federal Trade Commission in the U.S. led a lawsuit against Facebook, asking the court to force it to sell WhatsApp and Instagram.¹⁰⁰ The Digital Services Act and Digital Markets Act in the European Union are geared towards establishing a level playing field (to foster innovation and competitiveness) and creating a safer digital space for users and others.¹⁰¹ Our analysis shows that policies that encourage platform competition should be complemented by changes in platform liability. When bad actors are judgment proof and undeterred, then platform liability plays an important role of encouraging platforms to invest efficiently to protect users from harm.

4.3 Other Extensions

Firm Moral Hazard. In our baseline model and main extensions, platforms played an instrumental role in solving the adverse selection problem by detecting and removing bad actors from the platforms. As discussed in Section 2.1, adverse selection is empirically relevant: Bad actors, masquerading as legitimate firms, post fraudulent advertisements, steal user data, and sell counterfeit products. Moral hazard is also empirically relevant: Otherwise legitimate app developers may sell user data to others and manufacturers may cut corners to lower costs and raise profit margins. When firms are judgment proof, platform liability can play an instrumental role in solving moral hazard problems, too.

Our baseline model can be easily adapted to reflect a moral hazard problem. Suppose all the firms are identical ex ante but may become either the L -type or H -type ex post. A firm can take (unobservable) care at cost $k > 0$, which reduces the probability of

¹⁰⁰See <https://www.reuters.com/technology/us-ftc-says-court-should-allow-antitrust-lawsuit-against-facebook-go-forward-2021-11-17/>

¹⁰¹See <https://digital-strategy.ec.europa.eu/en/policies/digital-services-act-package>. A report written by Cremer, et al. and published by the European Commission (2019) raised concerns about increased concentration in platform markets. See <https://ec.europa.eu/competition/publications/reports/kd0419345enn.pdf>

becoming an H -type. If the firms are very judgment proof ($w_s \rightarrow 0$), then the H -types earn information rents. It follows that ex ante the firms have no incentive to take care and, as in the baseline model, platform liability raises the platform's auditing incentives.¹⁰² When the firms are modestly judgment proof (w_s in a middle range), platform liability motivates the platform to raise the interaction price, which deters the H -type firms and, under certain conditions, motivates the firms to take ex ante effort.¹⁰³

Same-Side Harms. The previous analysis considered a setting with cross-side harms: Firms on one side of the platform harmed the users on the other side of the platform. In practice, some users on platforms may harm other users. For example, some influencers on TikTok create videos that draw attention but may induce children to engage in dangerous activities; celebrities' endorsement of cryptocurrency may persuade investors to buy risky tokens.¹⁰⁴ These influencers can monetize user attention by collaborating with brands or sharing advertising revenue with platforms.¹⁰⁵

Our model can be adapted to investigate such cases with same-side harms. Consider for example a social networking platform where most user-generated content is perfectly safe but some of it is socially harmful. Suppose further that the advertising revenue that the platform enjoys is proportional to the volume of shared content, both safe and harmful. If the users are judgment proof, and cannot be held accountable for the harmful content that they post, then holding the platform liable may make sense. Without platform liability, the platform has a financial incentive to facilitate the posting and sharing of all content, both safe and harmful; with platform liability, the platform has incentives to

that the firms are very judgment proof ($w_s < b$) so that the L -types are marginal. If the H -types do not join the platform, the platform would not take any auditing effort. But anticipating this, the H -types would deviate to join. In this case, there is no equilibrium where the H -types are fully deterred.¹⁰⁷ Therefore, platform liability can increase the platform's auditing incentives.

False Positives. Our analysis assumed that there were no "false positives." The auditing efforts of the platform did not erroneously remove the L -type firms. Several new insights emerge when the baseline model is extended to include false positives. First, the second-best auditing effort is lower than in our baseline model (since it is socially efficient for L -types to remain on the platform). Second, the platform has weaker incentives to invest in auditing than in the baseline model (since the platform loses revenue when it excludes the L -types). Third, the platform's incentives are even weaker relative to the

reduce their control of online activities, similar to the potential distortion caused by vicarious liability on organizations.¹¹⁴ Second, our model shows that platform liability may be socially desirable even if auditing is very costly *o mpl etel y ine ecta* detecting bad actors. Although platforms would not engage in auditing in this case, liability would force platforms to internalize the social harms and create an incentive for them to use the price mechanism to deter bad actors.

There is active debate over whether platforms may be treated as common carriers.¹¹⁵ Common carriers, including telephone companies, mail carriers, and transportation systems (e.g., railroads and airlines) have a duty to serve the general public and may not generally exclude users.¹¹⁶ Common carriers are, however, subject to regulations that ensure public safety and sometimes have discretion or even a duty to exclude parties that may cause harm to others. For example, under federal law, airlines must deny transport to passengers who refuse to be searched for weapons¹¹⁷ and airline pilots have "permissive removal" authority to deny service to passengers who appear nervous or potentially disruptive.¹¹⁸ Although the Digital Millennium Copyright Act limits liability for internet service providers (ISPs), it also requires ISPs to terminate the accounts of repeat infringers.¹¹⁹ In a lawsuit brought against Western Union, the court opined that the defendant was in fact *obligated* to discontinue service for illegal gambling communications.¹²⁰ Common carriers can be held liable if they fail to meet their duties¹²¹ and, in many jurisdictions, the standard of care exceeds "reasonable care."¹²²

This article advances the idea that liability can play an instrumental role making

platforms safer for users and for society more broadly. An open question is whether civil liability is the best mechanism to accomplish these goals, or whether regulation would prove more effective. Social media and other platforms share similarities to common carriers and public utilities and so, by analogy, one could in principle regulate them in similar ways. Platform liability arguably has substantial advantages over regulation. Specifically, given the complexity and diversity of platforms, it would be difficult (and perhaps inadvisable) for regulators to set uniform safety standards.¹²³ Moreover, given the rapidly changing market conditions, regulators would be chasing a moving target. Platforms, especially big tech platforms, have the relevant information to weigh the social costs and benefits. Liability would give platforms financial incentives to use their discretion for the greater good.

¹²³This view is shared by many platforms; eBay's 2022 Transparency Report states: "regulatory regimes or technology mandates that are 'one size fits all' can actually serve to limit the tools, resources and partnerships necessary to combat bad actors."

- [14] Che, Yeon-Koo, and Kathryn E. Spier, "Strategic Judgment Proofing," *RAND Journal of Economics* Vol. 39 (2008), pp. 926-948.
- [15] Chen, Yongmin and Xinyu Hua, "Ex ante Investment, Ex post Remedies, and Product Liability," *International Economic Review* Vol 53 (2012), pp. 845-866.
- [16] Chen, Yongmin and Xinyu Hua, "Competition, Product Safety, and Product Liability," *Journal of Law, Economics, & Organization* Vol. 33 (2017), pp. 237-267.
- [17] Choi, Albert, and Kathryn E. Spier, "Should Consumers Be Permitted to Waive Products Liability? Product Safety, Private Contracts, and Adverse Selection," *Journal of Law, Economics, & Organization* Vol. 30 (2014), pp. 734-766.
- [18] Choi, Jay Pil, and Arijit Mukherjee, "Optimal Certification Policy, Entry, and Investment in the Presence of Public Signals," *RAND Journal of Economics* Vol. 51 (2020), pp. 989-1013.
- [19] Choi, Jay Pil, and Doh-Shin Jeon, "A Leverage Theory of Tying in Two-sided Markets with Nonnegative Price Constraints," *American Economic Journal: Microeconomics*, Vol. 13 (2021), pp. 283-337.
- [20] Culotta, Aron, Ginger Zhe Jin, Yidan Sun, and Liad Wagman, "Safety Reviews on Airbnb: An Information Tale," (2022), working paper.
- [21] Daughety, Andrew F., and Jennifer F. Reinganum, "Product Safety: Liability, R&D, and Signaling," *American Economic Review*, Vol. 85 (1995), pp. 1187-1206.
- [22] Daughety, Andrew F. and Jennifer F. Reinganum, "Market, Torts, and Social Inefficiency," *RAND Journal of Economics* Vol. 37 (2006), pp. 300-323.
- [23] Daughety, Andrew F., and Jennifer F. Reinganum, "Communicating Quality: a Unified Model of Disclosure and Signaling," *RAND Journal of Economics* Vol. 39 (2008b), pp. 973-989.
- [24] Daughety, Andrew F., and Jennifer F. Reinganum, "Imperfect Competition and Quality Signaling," *RAND Journal of Economics* Vol. 39 (2008a), pp. 163-183.
- [25] Dari Mattiacci, Giuseppe, and Francesco Parisi, "The Cost of Delegated Control: Vicarious Liability, Secondary Liability and Mandatory Insurance," *International Review of Law and Economics*, Vol. 23 (2003), pp. 453-475.
- [26] De Chiara, Alessandro, Ester Manna, Antoni Rubi-Puig and Adrian Segura-Moreiras, "Efficient Copyright Filters for Online Hosting Platforms," (2021), working paper.
- [27] Dukes, Anthony, and Esther Gal-Or, "Negotiations and Exclusivity Contracts for Advertising," *Management Science*, Vol. 22 (2003), pp. 222-245.

- [28] Epple, Dennis, and Artur Raviv, "Product Safety: Liability Rules, Market Structure, and Imperfect Information," *American Economic Review*, Vol. 68 (1978), pp. 80-95.
- [29] Farooqi, Shehroze, Maaz Musa, Zubair Shaq, and Fareed Za ar, "Canary-Trap: Detecting Data Misuse by Third-party Apps on Online Social Networks," arXiv:2006.15794v1 [cs.CY], (2020), <https://arxiv.org/pdf/2006.15794.pdf>.
- [30] Fu, Qiang, Jie Gong, and Ivan Png, "Law, Social Responsibility, and Outsourcing," *International Journal of Industrial Organization* (2018), pp. 114-146.
- [31] Galasso, Alberto, and Hong Luo, "Tort Reformation in the Age of Information," *Journal of Law and Economics*, Vol. 57 (2014), pp. 1-32.

- [58] Simon, Marilyn J., \Imperfect Information, Costly Litigation, and Product Quality," *Bel I Journal of Economics*, Vol. 12 (1981), pp. 171-184.
- [59] Sitaraman, Ganesh, \Deplatforming," *Yale Law Journal* forthcoming.
- [60] Shavell, Steven, \The Judgment Proof Problem," *International Review of Law and Economics* , Vol. 6 (1986), pp. 45-58.
- [61] Spence, A. Michael, \Consumer Misperceptions, Product Failure, and Producer Liability," *Review of Economic Studies* , Vol. 44 (1977), pp. 561-572.
- [62] Spence, A. Michael, \Monopoly, Quality, and Reputation," *Bel I Journal of Economics* , Vol. 6 (1975), pp. 417-429.
- [63] Tadelis, Steven, \Reputation and Feedback Systems in Online Platform Markets," *Annual Review of Economics*, Vol. 8 (2016), pp. 321-340.
- [64] Tan, 1-4]

Appendix A

An Example of the Coordination Game. This example illustrates the idea that, given the same-side network effects, the platform finds it optimal to set a sufficiently small price (or even zero price) for the users.

Suppose that there are two potential users, 1 and 2, who independently choose whether to join the platform or not. Each user receives a private benefit v_i , if and only if both users join the platform. In addition, when joining the platform, a user incurs costs c_i , $i=1,2$, which can include entry costs, opportunity costs, and the expected harm caused by the firms on the platform. The platform charges the same membership fee,

Now we prove the remaining results in the lemma. Using the definition of $r_H(w_s)$ in the lemma, (8) implies $e > 0$ if and only if $(r_H(w_s) - r_L(w_s)) < 0$. This gives the condition for cases 1 and 2. Totally differentiating (10), and using the fact the social welfare function is concave, gives $de = dw_s = S''(e) > 0$ and $de = dw_p = S''(e) > 0$. When $e > 0$ (an interior solution), increasing the level of liability for either the firm or the platform increases the platform's auditing effort. Equation (10) implies $e > e^*$ if and only if $r_H(w_s) - r_L(w_s) > 0$. This gives the condition for subcases 2(a), 2(b) and 2(c).

Proof of Proposition 1. Note that $w < d < \frac{1}{L}$ by Assumption A1. Suppose $w_p = 0$ and $w_s = w$. From Lemma 1, a necessary and sufficient condition for $e = 0$ is (8) or

$$r_H(w_s) - r_L(w_s) > (r_H(w) - r_L(w)) \left(\frac{w}{w_s} \right).$$

Substituting for w from (5),

$$r_H(w_s) - r_L(w_s) > (r_H(w) - r_L(w)) \left(\frac{w}{w_s} \right);$$

which is equivalent to $w_s < \frac{1}{L}$. Since $w_s < \frac{1}{L}$

Proof of Proposition 3. We start by showing that the platform does not charge the users (i.e. $m = 0$) if $\lambda_L (\lambda_H + (1 - \lambda_L)d)$ is sufficiently large. To see this, first consider the scenario where the L -type firms are marginal ($w_S = \bar{w}$). Given the belief e and damage award $w = w_S + w_p$, a user will participate when

$$v - m + [(\lambda_H + (1 - \lambda_L)d)](d - w) > 0$$

The platform's equilibrium price charge to the firms is the same as in the baseline model (see Lemma 1). Thus, the platform's profits are

$$[1 - F(m + (\lambda_H + (1 - \lambda_L)d)(d - w))][\lambda(e) + m] - c(e)$$

where $1 - F(\cdot)$ is the users' participation rate and

$$\lambda(e) = (1 - e) (\lambda_L \lambda_L w_S + \lambda_H w_p) + (1 - \lambda_L) (\lambda_L \lambda_L w) \quad (30)$$

When $e = 0$, $w_S = 0$ and $w_p = d$, $\lambda(e)$ achieves the lowest value

$$\lambda_L (\lambda_H + (1 - \lambda_L)d)$$

which is positive by Assumption A2. Taking differentiation of the profit function with respect to m , we have

$$[1 - F(\cdot)] - f(\cdot)[\lambda(e) + m]$$

which is negative if $\lambda(e)$ is sufficiently large. Hence, if $\lambda_L (\lambda_H + (1 - \lambda_L)d)$ is sufficiently large, the platform would set $m = 0$.

Next, consider the scenario where the H -type firms are marginal ($w_S > \bar{w}$). If the platform accommodates all the H -type firms, a user will participate when

$$v - m + [\lambda_H + (1 - \lambda_L)d](d - w) > 0$$

If the platform deters all the H -types firms by charging a larger price, a user will participate when

$$v - m + (1 - \lambda_L)(d - w) > 0$$

Similar to the earlier analysis, we can show that, if $\lambda_L \lambda_L d$ is sufficiently large, the platform would set $m = 0$.

In the remaining analysis, we maintain the assumption that $\lambda_L (\lambda_H + (1 - \lambda_L)d)$ is sufficiently large, which also implies $\lambda_L \lambda_L d$ is sufficiently large, such that the platform does not charge the users.

Now we prove condition (15), which highlights the potential divergence between the private and social incentives for auditing. Given w , (12) implies

$$\frac{dS(e; \bar{w})}{de} = \frac{\partial S(e; \bar{w})}{\partial e} - \frac{\partial S(e; \bar{w})}{\partial \bar{w}} \lambda_H(d - w)$$

Using (14), if the equilibrium auditing effort is positive, then e^U satisfies

$$\begin{aligned} \frac{\partial S(e^U; \nu)}{\partial e} &= \frac{\partial S(e^U; \nu)}{\partial e} + \int_{\nu}^Z (H - L)(\nu - w_s) - H(d - w) f(v) dv \\ &= \frac{dS(e^U; \nu)}{de} + \int_{\nu}^Z (H - L)(\nu - w_s) - H(d - w) f(v) dv + H(d - w) \frac{\partial S(e^U; \nu)}{\partial \nu} \\ &= 0: \end{aligned}$$

Next, we show that, if $w_s < \nu$, then $w_p^U > w_p$. Totally differentiating (12) with respect to w_p gives

$$\frac{dS(e^U; \nu)}{dw_p} = \frac{\partial S(e^U; \nu)}{\partial e} \frac{de}{dw_p} + \frac{\partial S(e^U; \nu)}{\partial \nu} \frac{d\nu}{dw_p} - H(d - w) \frac{dw}{dw_p} + \frac{\partial S(e^U; \nu)}{\partial \nu} \frac{d\nu}{dw_p}; \quad (31)$$

where $\frac{\partial S(e^U; \nu)}{\partial \nu} < 0$ and $\frac{d\nu}{dw_p} < 0$. Similar to the analysis in the baseline model, we can show that, given

Proof of Proposition 4. We prove two claims respectively for $w_s \leq w$ and $w_s > w$.

Claim 1 : Suppose $w_s \leq w$. The platform sets $p^r = p_L = p_L(w_s, d, w)$ and attracts the H -type firms where $e^r = E(je^r)$ are the equilibrium posterior beliefs. Let $E(je)$, $e^0 = E(j0)$, and $r_H(w_s) = (p_H - c_L)(w - w_s)$.

1. If $(p_H - c_H)d + (p_H - c_L)(d - w) > r_H(w_s)$ then the platform does not audit, $e^r = 0 < e^0$.
2. If $(p_H - c_H)d + (p_H - c_L)(d - w) < r_H(w_s)$ then $e^r > 0$. The platform's auditing effort decreases in firm liability $dw_s < 0$ and increases in platform liability $de^r = dw_p > 0$.
 - (a) If $(p_H - c_H)(d - w) > r_H(w_s)$ then $0 < e^r < e^0$.
 - (b) If $(p_H - c_H)(d - w) = r_H(w_s)$ then $0 < e^r = e^0$.
 - (c) If $(p_H - c_H)(d - w) < r_H(w_s)$ then $0 < e^r < e^0$.

Proof of Claim 1 : Since $w_s \leq w$, it is not possible for the platform to deter the H -types without deterring the L -types, too. If the L -type is willing to participate, then the H -type also prefers to participate.

To begin, we construct values $e^r; p^r; t^r$ that maximize the platform's profits subject to the platform's incentive compatibility constraint and the participation constraints of the consumers and the L -type firms (as the L -type firm is marginal). Then, we will verify that these values are an equilibrium of the game.

$$\max_{e; p; t} \pi(e; p) = (1 - e)(p - c_H w_p) + (1 - e)(p - c_L w_p) - c(e) \quad (36)$$

subject to

$$e = \arg \max_{e \geq 0} \pi(e; p) \quad (37)$$

$$0 \leq t \leq E(je)(d - w_s - w_p) \quad (38)$$

$$t \leq (p - c_L)w_s + c_L - p \leq 0 \quad (39)$$

(37) is the platform's incentive compatibility constraint, (38) is the consumer's participation constraint, and (39) is the L -type firm's participation constraint.¹²⁴

The L -type's participation constraint (39) must bind. To see this, consider two cases. First, suppose that neither (38) nor (39) binds. Then the platform would increase the price p which would increase the platform's profits in (36) and maintain the consumer's participation constraint (38). Second, suppose that (38) binds while (39) does not. Again, the platform would increase the price p marginally. The direct effect of increasing p is

¹²⁴The H -type's participation constraint is satisfied if (39) holds, and is therefore not included in the program.

that the platform's profits in (36) increase. Since $\frac{\partial^2 \pi(e; p)}{\partial e \partial p} = \frac{\partial \pi(e; p)}{\partial p} < 0$; increasing p also (weakly) decreases the platform's effort in (37), which in turn raises $E(je)$ and, since (38) binds, reduces t . However, since t is not in (36), the platform's profits still increase.

Since the L -type's constraint (39) binds, $p = t - (w_s + c_L)$ and we can rewrite the optimand (36) as a function of e and t :

$$(1 - e)(t - (w_s + c_L) - w_p) + (1 - e)(t - (w_s + c_L) - w_p) - c(e): \quad (40)$$

Next, we show that the consumer's participation constraint (38) binds. Suppose not. Then, the platform would increase t and its profits would rise. Since both participation constraints (38) and (39) bind, we have

$$p = 0 - E(je)(w_s + w_p) - (w_s + c_L): \quad (41)$$

Since $w_s = 0 - c_L$ and $w = w_s + w_p$ the solution to the platform's optimization problem is:

$$e^r = \arg \max_{e \geq 0} \pi(e; p^r) \quad (42)$$

$$t^r = 0 - E(je^r)(w) \quad (43)$$

$$p^r = w_s - E(je^r)(w): \quad (44)$$

We now verify that the values $e^r; p^r; t^r$ defined in (42), (43), and (44) are an equilibrium of the game. Suppose that the platform charges p^r in (44), and that the firms and consumers believe that the probability of harm is $\tau^r = E(je^r)$ where e^r is defined in (42). The consumers are (just) willing to pay t^r in (43) and the L -type firms are (just) willing to pay p^r in (44). If the consumers and the firms all participate, the platform exerts effort e^r in (42). Therefore the equilibrium beliefs $\tau^r = E(je^r)$ are consistent.

Next, we verify that Assumption A2 guarantees that the platform's profits are positive.

We now show that the algebraic condition in case 1 is necessary and sufficient for a corner solution, $e^r = 0$. We first show the condition is necessary. If $e^r = 0$ then $E(j^0) = 0$. Since the consumer's participation constraint (38) binds we have $e^r = 0$; since the L -type firm's participation constraint (39) binds we have $p^r = \frac{L}{L} w_s^0(d-w)$. Finally, for $e^r = 0$ to satisfy the platform's IC constraint (37) we need $(e; p) = 0$ or equivalently $p^r - H W_p = 0$. Substituting p^r , this condition becomes

$$\frac{L}{L} w_s^0(d-w) - H W_p = 0 \quad (45)$$

Adding and subtracting terms this becomes

$$(H - H d) - (H - L) \frac{L}{L} w_s^0(d-w) + H W_p + H W + (H - 0)(d-w) = 0 \quad (46)$$

and rearranging this expression gives

$$(H - H d) + (H - 0)(d-w) - (H - L) \frac{L}{L} w_s^0(d-w) = 0 \quad (47)$$

The right-hand side is $r_H(w_s)$. This confirms that the condition in case 1 is necessary.

Next, we show that the condition in case 1 is sufficient. Suppose the condition holds and $e^r > 0$. Since $E(j^{e^r}) < 0$, $t^r > 0$ and $p^r > \frac{L}{L} w_s^0(d-w)$. Assumption A2 implies $p^r - H W_p > 0$, so the platform does not audit, $e^r = 0$.

Now consider case 2. The condition implies $p^r - H W_p < 0$ so the platform is losing money from each H -type transaction. The equilibrium effort $e^r > 0$ and consumers' equilibrium beliefs $p^r = E(j^{e^r})$ satisfy equation (23). The platform charges $p^r = \frac{L}{L} w_s^r(d-w)$ and consumers believe that the platform will exert effort e^r to audit,

the rms charge the consumers $t^r = 0$ ($d - w$). The platform's price extracts the marginal H -type firm's surplus, that is, $p^r = t^r - (c_H w_s + c_H)$ or

$$p^r = c_H - c_H w_s - (d - w) \quad (50)$$

The platform's profits are

$$\begin{aligned} p^r - w_p &= (1 - \alpha_L)(d - w_p) + (\alpha_H - \alpha_H d) + (1 - \alpha_H)[\alpha_H - \alpha_L - (\alpha_H - \alpha_L)w_s] \\ &= (1 - \alpha_L)(d - w_p) + (\alpha_H - \alpha_H d) + (1 - \alpha_H)(\alpha_H - \alpha_L)(w - w_s) \\ &< (1 - \alpha_L)(d - w_p) \end{aligned}$$

where the inequality follows from Assumption A1 and $w_s > w$: Therefore, if $w_s > w$, the platform charges $p^r = \alpha_L - \alpha_L(d - w_p)$ and deters the H -types.

We now proceed to proof Proposition 4. Suppose $w_s = w$, so the L -type is marginal. From Claim 1, we have $e^r = e$ if and only if

$$(\alpha_H - \alpha_L)(w - w_s) - (\alpha_H - \alpha_H)(d - w) = 0 \quad (51)$$

Substituting that $w = w_p + w_s$ and isolating w_p on the left-hand side establishes the result. Suppose $w_s > w$. The results follow from Claim 2.

Online Appendix B

This appendix contains the analysis of four additional extensions: heterogeneous users with observable effort, firm moral hazard, false positives and litigation costs.

B1. Heterogeneous Users with Observable Effort

Section 3 shows that platform liability can be socially desired when heterogeneous users make participation decisions but do not observe the platform's auditing effort. Now we consider the setting where the platform can commit to its auditing effort before the users make participation decisions.

If $w_s > \bar{w}$ then the analysis is the same as case 2 in Section 3. As shown in Section 3, if $w_s > \bar{w}$, the platform would not take any auditing effort, and imposing full residual liability on the platform implements the first-best outcome. The following analysis examines case 1 where $w_s < \bar{w}$.

When auditing effort is observable, equation (14) implies that the platform's effort (if it is positive) satisfies

$$\frac{d(e^U; \bar{w})}{de} = \frac{dS(e^U; \bar{w})}{de} + \int_{\bar{w}}^Z [(H - L)(\bar{w} - w_s) - H(d - w)] f(v) dv$$

$$H(d - w) [(1 - e^U)(H - L)(\bar{w} - w_s)] f(\bar{w}) = 0 \quad (52)$$

where $\bar{w} = \bar{w}(e; w)$.

When $w_s = \bar{w}$, $\frac{d(e^U; \bar{w})}{de} = \frac{dS(e^U; \bar{w})}{de}$ if and only if $w_p^U = d - w_s$. Therefore, imposing full residual liability on the platform implements the second-best outcome: the platform chooses $e^U = e$ and all the users join the platform.

When $w_s < \bar{w}$, the last term on the right-hand side of equation (52) is negative. Moreover, if $w_p < w_p^U$, where $w_p^U \in (0; d - w_s)$ is the optimal platform liability in Proposition 2 of the baseline model, then the second term on the right-hand side of equation (52) is non-positive. Therefore, $\frac{dS(e^U; \bar{w})}{de} > 0$; that is, the platform's auditing incentive is socially insufficient. The social planner chooses w_p to maximize social welfare:

$$\frac{dS(e^U; \bar{w})}{dw_p} = \frac{dS(e^U; \bar{w})}{de} \frac{de^U}{dw_p} + \frac{\partial S(e^U; \bar{w})}{\partial \bar{w}} \frac{\partial \bar{w}}{\partial w_p}; \quad (53)$$

where $\frac{\partial \bar{w}}{\partial w_p} = (1 - e^U)$

1. If $w_s < w$, then $w_p^u > w_p$ as long as $\frac{de^u}{dw_p} > 0$. The platform sets $p^u = L - L w_s$.
The second-best outcome is not achieved.

2. If

may become either the L -type or H -type ex post. If a firm takes (unobservable) care with cost $k > 0$, the probability of becoming an H -type is b . If the firm does not take care, the probability of being an H -type rises to $b^0 > b$: The platform commits to its price p before the firms decide to take care or not. The firms privately learn their realized types and decide whether to join the platform.

For simplicity, we maintain the following assumption

$$k < (b^0 - b)(d_L - d) + (b - b^0)(d_H - d): \quad (54)$$

Assumption (54) leads to several implications.

First, since $(b - b^0)(d_H - d) < 0$, $k < (b^0 - b)(d_L - d)$. If the H -types never join the platform, it is socially efficient for the (ex ante identical) firms to invest k .

Second, Assumption (54) implies

$$k < (b^0 - b)[(d_L - d) - (d_H - d)] = (b^0 - b)(d_L - d):$$

Even if both types join the platform, it is efficient for the firms to invest k .

Finally, Assumption (54) implies

$$(b - b^0)(d_H - d) + (1 - b^0)(d_L - d) > (1 - b)(d_L - d):$$

that is, social welfare is larger if all the firms invest k and join the platform than if no firm invests and only the L -types join the platform.

In the first-best benchmark, all the firms invest k ex ante and only the L -types join the platform. Given k , there exists $w^k \in (0, d)$ such that, if and only if $w_s > w^k$,

$$k < (b - b^0)(d_H - d)(w_s^k):$$

The profit difference,

$$\pi^0 - \pi^L = (b_L - b)(w_L - w_S - w_p) - k(1 - b_L);$$

decreases in w_p . That is, the platform has stronger incentives to charge p_0 if w_p is lower. When $k > \frac{(b_L - b)^2}{(1 - b_L)}(w_L - w_S)$, then the platform never charges p_0 , so platform liability is unnecessary. When $k < \frac{(b_L - b)^2}{(1 - b_L)}(w_L - w_S)$, then $\pi^0 - \pi^L > 0$ if $w_p = 0$ but may become negative if w_p is large, so it is optimal to set $w_p = 0$.

Case 2.2: $w_S > w^k; \frac{H}{H}$: Given $w_S < \frac{H}{H}$, the H -types may have incentives to join the platform. Moreover, given $w_S > w^k$, we have $k < (b_H - b)(w_S - w)$, which implies $p_0 > p_H = \frac{H}{H} w_S > 0$. If the platform charges p_L , the firms would not invest k and the platform's profit is

$$\pi^L = (1 - b_L)(w_L - w_S - w_p);$$

If the platform charges p_H , the L -types receive information rent $(b_H - b_L)(w_S - w)$. Since $k < (b_H - b)(w_S - w)$, the firms would invest k and always join the platform. Then the platform's profit is

$$\pi^H = (b_H - b)(w_S - w) + (1 - b_L)(w_L - w_S - w_p);$$

If the platform charges

$w_s \geq (w^k; \bar{w})$, only under a non-empty set of $w_p > 0$, the platform charges p_0 and the first-best outcome is achieved.²⁷ That is, if $w_s \geq (w^k; \bar{w})$, platform liability is socially desired.

If $w_s = \bar{w}$, $\pi^H = 0$ and $\pi^L = 0$ only under $w_p = 0$, so it is optimal to set $w_p = 0$. If $w_s \geq (\bar{w}; \frac{H}{H})$, the platform never charges p_0 . Since it is efficient for all the firms to invest k and the profit difference $\pi^H - \pi^L$ decreases in w_p , it is optimal to set $w_p = 0$, under which the platform charges p_H and the firms invest k .

Case 2.3: $w_s \geq (w; w^k)$. Given $w_s < w^k$, we have $k > (b^H - b^L)(w_s - w)$, which implies $p_0 < p_H$. If the platform charges p_L , the firms would not invest k and the platform's profit is

$$\pi^L = (1 - b^L)(\pi^L - w_s - w_p):$$

If the platform charges p_H , the L -types receive information rent $(b^H - b^L)(w_s - w)$. Since

B3. False Positives (Type-I Errors)

Now we extend the baseline model by considering false positives. Suppose that the auditing effort of the platform may erroneously remove the L -type firms with probability e , where $e < 1$. The first-best benchmark is the same as in the baseline model. For the second-best benchmark, suppose that the H -type firms seek to join the platform. Social welfare is:

$$S(e) = v + (1 - e)($$

B4. Litigation Costs

We extend the baseline model by considering litigation costs. When a user gets harmed by a firm and files a lawsuit, the litigation costs are $Z_p; Z_s; Z_b$, respectively for the platform, the firm, and the user. Denote $Z = Z_p + Z_s + Z_b$. Assume that $Z_b < w_s + w_p$ and $\lambda_L < \lambda_L d + Z > 0$.¹²⁸ So, litigation is credible and it is efficient to have interactions between the L -type firms and users. If the H -type firms seek to join the platform, social welfare is

$$S(e) = v + (1 - e)(\lambda_H - \lambda_H(d + z)) + (1 - e)(\lambda_L - \lambda_L(d + z)) - c(e):$$

The socially optimal auditing effort $\bar{e} > 0$ satisfies

$$-(\lambda_H - \lambda_H(d + z)) - c'(e) = 0:$$

The two types of firms have the same rent when:

$$w_s + z_s = w = \frac{\lambda_H}{\lambda_H - \lambda_L} \frac{\lambda_L}{\lambda_L}: \quad (58)$$

Case 1: $w_s + z_s = w$. The platform sets $p^z = \lambda_L - \lambda_L(w_s + z_s)$ to extract the L -type firms' rent. The platform chooses $e > 0$ if and only if $p^z - \lambda_H(w_p + z_p) < 0$, which can be rewritten as

$$\lambda_H - \lambda_H(w + z_p + z_s) - (\lambda_H - \lambda_L)(w - w_s - z_s) < 0:$$

The platform's profits can be written as

$$\pi(e) = S(e) - (1 - e)(\lambda_H - \lambda_L)(w - w_s - z_s) + [(1 - e)\lambda_H + (1 - e)\lambda_L](d + z_b - w) - v:$$

Denote the equilibrium auditing effort as e^z . If $e^z > 0$, the first-order condition is

$$\pi'(e^z) = S'(e^z) - (\lambda_H - \lambda_L)(w - w_s - z_s) - \lambda_H(d + z_b - w) = 0: \quad (59)$$

The users' uncompensated loss caused by the types, $\lambda_H(d + z_b - w)$, increases in z_b ; and the firms' information rent, $(\lambda_H - \lambda_L)(w - w_s - z_s)$, decreases in z_s . Therefore, as compared to the baseline model, the platform's auditing incentives are even weaker relative to the social incentives. We can show the following results.

Lemma 3. Suppose $w_s + z_s = w$. The platform sets $p^z = \lambda_L - \lambda_L(w_s + z_s)$ and attracts the H -types. Let $r_H^z(w_s) = (\lambda_H - \lambda_L)(w - w_s - z_s)$ denote the H -types' information rents.

1. If $\lambda_H - \lambda_H(w + z_p + z_s) > r_H^z(w_s)$ then the platform does not audit, $e^z = 0 < \bar{e}$.
2. If $\lambda_H - \lambda_H(w + z_p + z_s) < r_H^z(w_s)$ then $e^z > 0$.

(a) If $\lambda_H(d + z_b - w) > r_H^z(w_s)$ then $0 < e^z < \bar{e}$.

¹²⁸We also assume that z is lower than the benefit of improved platform incentives.

(b) If $r_H^z(d + z_b - w) = r_H^z(w_s)$ then $0 < e^z = \bar{e}$.

(c) If $r_H^z(d + z_b - w) < r_H^z(w_s)$ then $0 < \bar{e} < e^z$.

Case 2: $w_s + z_s > w$. The platform's profit-maximizing strategy is to either charge $p = r_L^z(w_s + z_s)$ and deter the H -types from joining the platform or charge $p = r_H^z(w_s + z_s)$ and attract both types. The platform will charge $p = r_H^z(w_s + z_s)$ and attract the H