# On Customer (Dis)honesty in Priority Queues: The Role of Lying Aversion

Arturo Estrada Rodriguez, Rouba Ibrahim, Dongyuan Zhan

University College London, 1 Canada Square, London E14 5AB

arturo.rodriguez.18@ucl.ac.uk , rouba.ibrahim@ucl.ac.uk, d.zhan@ucl.ac.uk

Priority queues where people make deceitful claims to access service faster are everywhere. For example,

there have been multiple recent reports of people ineligibly "jumping" COVID-19 testing queues. Motivated

by the prevalence of such queue-jumping behaviour in practice, we study a two-priority $M/G/1$ queueing

system where priority levels are private information to customers. Customers make strategic claims about

their true priorities, and the manager must decide on the static scheduling policy in the system, based on

those claims. Customers are both delay sensitive and averse to lying i.e., they incur intrinsic psychological

costs when they are untruthful. The manager's objective is to minimize the total equilibrium expected delay

cost, which includes both under-prioritization and over-prioritization errors. We find that the equilibrium

proportion of customers who are untruthful is bounded away from 1, independently of waiting times in the

system. We also find that, in this system with endogenous customer claims, the optimal scheduling policy

deviates from the celebrated $c\mu$ rule. In particular, while the manager should always route high claims to

the high-priority class, it may be optimal to upgrade some low claims to the high-priority class as well, in

order to incentivize honesty in the system. To substantiate our theoretical analysis, we run controlled online

queueing experiments where we validate our theoretical insights with experimental data.

*Key words*: Scheduling, priority queues, strategic customers, lying aversion, experiments.

## 1. Introduction

Consider a priority queueing system where customers are delay sensitive. Assume that a customer's

true priority level is private information to her, i.e., it is not known by the service provider. Assume

that priority assignments are made according to customer claims about their own priorities, and that deceitful claims cannot be detected or punished. Then, we would expect, under the classical assumption of "homo economicus" (Becker 1968), that all customers would claim the highest priority level to access service faster. Fortunately, people are usually not purely selfish beings who only care about material outcomes. This is especially true when they are asked to report information that is private to them. In fact, people tend to be intrinsically averse towards misreporting private information due to moral or religious reasons, self-image concerns, or an unwillingness to deviate from socially-acceptable behaviour (Abeler et al. 2014, Gibson et al. 2013).

There exists a rich literature studying priority queueing systems where customer priority levels are assumed to be fully known and exogenously specified. In contrast, there is very little research on priority queues under imperfect information on customer priorities, despite the prevalence of such queues in practice. In this paper, we use both controlled experiments and a stylized queueing model to study priority queues where customer priority levels are unknown to the service provider, and where customers strategically convey private information about their own priority levels. On one hand, to save on waiting time, customers have a clear incentive to be untruthful, i.e., to claim that they should be given high priority. On the other hand, being untruthful is undesirable to people. When, why, and to what extent would people be untruthful in such settings? What are the system-level effects of this untruthful behaviour? In a setting where people can be untruthful with complete impunity, are their claims informative at all? How could the service provider mitigate the negative effects of that untruthful behaviour? We provide answers to those questions in this paper.

*Motivating examples.* Priority queues where people make deceitful claims in order to access service faster are everywhere. For example, during the COVID-19 pandemic in the United Kingdom, there have been reports of people lying about their employment status (BBC 2020b) or their symptoms (BBC 2020a) in order to receive a COVID-19 home-test faster. Interestingly, the online booking system relied on public honesty rather than employment checks (Weaver and Proctor 2020). There have also been recent reports, worldwide, of untruthful behavior relating to COVID-19 vaccination, with ineligible people jumping priority vaccine queues; e.g., in the United States

(NPR 2021a). Typically, clinics do not keep track of who is eligible and who is not among those vaccinated, and rely instead on an "honor system" for vaccination eligibility (NPR 2021a). Indeed, requiring additional diagnostic checks to verify claims would add barriers to access, and is therefore undesirable or logistically difficult to enforce (NPR 2021b). In another example, a chat-bot symptom checker application was recently commissioned by the National Health Service (NHS) to aid in booking general practitioner (GP) appointments for patients. This initiative was dropped, after initial pilot trials, because patients admitted that they would exaggerate their symptoms (in the text exchange with the chat-bot symptom-checker) in order to see the GP faster (Heather 2017). In telephone triage systems, which are widely used to triage patients with time-sensitive conditions, it is known that patients routinely exaggerate symptoms to get a doctor's appointment sooner (Kirton et al. 2020) or an ambulance faster (Jones 2020). Finally, with online shopping, people have been routinely able to exploit loopholes and jump online queues to shop for groceries, concert tickets, or retail products (Storm 2015, Isaiah 2020).

The aforementioned examples share several common features. First, in each example, there is asymmetric information between customers and the service provider: Only customers know their true priority levels (or have, at least, some partial information about it). Thus, there is a clear incentive for customers to misreport their own priority levels in order to begin service faster. Moreover, there is an opportunity to do so with complete impunity, as priority allocations are made solely based on customer claims, and without additional diagnostic checks. Additionally, in those settings it would be difficult to punish or fine untruthful behavior.

Second, queues in those settings are unobservable, i.e., a customer makes her priority claim without seeing other customers waiting. This is because the service requests themselves are made either online or through the phone. Moreover, those services typically involve a single (or rare) interaction with the service provider, so that there is a strong degree of anonymity for customers. Thus, reputation concerns do not play a (major) role, and community enforcement of desirable behavior, which may occur with physical queues, is not possible (Allon and Hanany 2012).

Third, market mechanisms, e.g., the possibility to purchase priority, which are standard tools to induce truthful, incentive-compatible, reporting are not applicable in those settings. For example, healthcare services provided by the NHS are publicly funded and cannot be individually priced.

Finally, while there is evidence of customer untruthfulness in each one of those examples, clearly it remains that not everyone is untruthful. Thus, we need to develop models which capture both materialistic incentives (the desire to save time) and the psychological drivers behind lying or telling the truth. Importantly, we also need to validate those models with experimental data. In this paper, we propose such a model, rely on that model to design better decision-making in the system, and validate our theoretical insights with experimental data.

*Contributions.* Even in simple one-on-one interactions involving monetary transactions, the extant literature has only begun to shed light on the underlying motives behind telling or concealing the truth (Abeler et al. 2019). To model aversion to untruthfulness in queues, we focus here on people's *intrinsic lying cost.* The existence of such an intrinsic cost has been experimentally validated with money as the incentive (Gibson et al. 2013, Özer et al. 2011). In particular, we theorize that people are reluctant to being untruthful because they incur an intrinsic cost depending on the size of their lie. This intrinsic aversion is especially relevant in the unobservable, anonymous, queueing settings that we are interested in, where customers can be untruthful with complete impunity.

*Operational insights.* To glean theoretical insights on system performance in queues with untruthful customers, we study a two-priority $M/G/1$ queueing system where the manager's information about customer priority levels, high ($H$) or low ($L$), is based solely on customer claims. Customers are both delay sensitive and averse to being untruthful. They observe the average waiting times in the two queues, and decide on their claimed priority levels: $H$ or $L$. We derive the optimal claim-based static scheduling policy in the system. The manager seeks to minimize the total expected delay cost in the system, which includes both under-prioritization (sending a true high-priority customer to the low class) and over-prioritization (sending a true low-priority customer to the high class) errors. Since customers know the manager's scheduling rule, their priority

claims are strategic. Waiting times are, themselves, affected by customer claims which, in turn, influence the manager's scheduling rule. Thus, a complex equilibrium emerges.

At equilibrium, we find that there is a proportion, bounded away from 1, of untruthful customers who wrongfully claim high priority. Moreover, this proportion is independent of the waiting times in the system. This insight is consistent with the extant literature: When money is at stake, rather than time, there is a solid body of experimental evidence which shows that the proportion of people who are untruthful is independent of monetary incentives (Abeler et al. 2019). For the manager's scheduling problem, we find that, if the distribution of the customer lying aversion cost is "sufficiently inelastic" (in a sense to be made more precise later), then the optimal routing policy that arises in equilibrium is to assign priorities according to customer claims. This provides theoretical validation to the "honor-based" prioritization rule which is commonly observed in practice, e.g., in our motivating examples above. However, if that distribution is "sufficiently elastic", then it is optimal to give some low-priority claims high priority, i.e., to upgrade these claims. Importantly, *this optimal routing policy deviates from the celebrated $c\mu$ rule*, which prioritizes customers in decreasing order of their $c\mu$ index, where $c$ denotes the delay cost and $\mu$ the service rate. This is noteworthy because the $c\mu$ rule has been shown to be optimal with uncertain, yet exogenous, customer priority claims. There, priority should be given according to the expected $c\mu$ values instead (Argon and Ziya 2009, Bren and Saghafian 2019). In our setting with uncertain, yet *endogenous*, claims, it turns out that, at equilibrium, customers with high claims have a strictly higher expected $c\mu$ value than customers with low claims. This means that, if the $c\mu$ rule were to hold, no customer with a low claim should be upgraded to the high priority class. However, in our model, such upgrading is optimal because customer claims are endogenous: The key insight here is that by giving some low-priority claims high priority, the manager incites more honesty in the system. Thus, customer claims become more informative, at equilibrium, which ultimately benefits the system as a whole.

*Experimental evidence.* To validate our theoretical insights, we run controlled experiments where people can misreport their private information in order to wait less in a virtual queue. Our experiments are specifically designed to tease out the behavioral implications of people's intrinsic aversion

to being untruthful. In particular, we rely on the experimental design in Fischbacher and Föllmi-Heusi (2013) to measure untruthfulness. In their original design, Fischbacher and Föllmi-Heusi (2013) use making money as the driver behind untruthfulness. Here, we adapt that design to our queueing setting where time is the incentive instead. We consider two virtual queues, representing different priority levels, with different waiting times. Participants roll a die at home and, based on the privately observed outcome of that die roll, make a claim about which queue they should wait in, short or long. The queues are unobservable: Participants do not see the dynamic state of the queue, and make their claims based solely on average-wait information. We run several experimental conditions, and we measure the proportions of people who claim that they should wait in the short queue. Consistent with our theoretical insights, our experimental data show that a significant proportion of participants refrains from misreporting, i.e., claims that they should wait in the longer queue. As in our model, we also observe that the proportion of participants who claim that they should wait in the longer queue does not decrease as the waiting time in queue increases. Finally, we find that offering participants the chance to wait in the short queue even if their die roll corresponds to the long queue (i.e., the possibility to be upgraded) leads to decreasing the probability of untruthfulness in the system. This provides experimental evidence to our main theoretical insight about the role of upgrading in increasing honesty in the system.

The remainder of the paper is structured as follows. In §2, we review the relevant literature. In §3, we describe our queueing model. In §4, we derive the optimal routing policy. In §5, we describe our experimental design, our hypotheses, and our results. Finally, in §6 we draw conclusions. We relegate additional experimental results and modelling extensions to the appendix.

## 2. Literature Review

There is a rich queueing-theoretic literature which studies scheduling decisions in priority queues where customer priority levels are perfectly known to the manager. In particular, the celebrated $c\mu$ rule has been shown to be optimal in various settings (Cox and Smith 1991, Van Mieghem 1995).

*Queues with uncertain priorities.* In practice, the true customer priority levels may not be perfectly known to the service provider. Van der Zee and Theil (1961) studies how exogenous misclassification errors affect optimal scheduling decisions in the system. Argon and Ziya (2009) assumes that the service provider receives a signal from each arriving customer, where the signal is the probability that the customer is of the high-priority class. Argon and Ziya (2009) shows that the Highest Signal First (HSF) scheduling policy, which is consistent with the $c\mu$ rule, yields the lowest long-run average waiting cost among all finite-class priority policies. Bren and Saghafian (2019) consider the case in which the class of an arriving customer is known but the service rate for each class has to be dynamically learnt through a data-driven optimization approach. They also derive an optimal scheduling policy that is similar to the HSF policy. Singh et al. (2021) assume that customer signals are the output of a data-driven classifier, and propose an integrated approach where the classifier and the prioritization policy are jointly optimized. In line with that stream of literature, we also investigate prioritization decisions when customer priority levels are not known to the service provider. However, unlike those papers, we focus on settings where customers have private information about their own priorities, and where they strategically manipulate the signals (claims) that they send to the service provider.

*Priority queues with strategic customers.* There is a rich queueing-economics literature which focuses on settings where customer information is private and where market mechanisms can be used to align customer and central planner incentives (Hassin and Haviv 2003). Similar to that literature, we focus on modelling queues with strategic customers. However, that literature typically takes the perspective of a firm which shares information with customers, or induces the desired outcomes from customers through incentive-compatible pricing. In contrast, we take here a customer-centric approach to information provision. Hu et al. (2021) takes a similar perspective and considers a service system where customers are delay sensitive and where they strategically manipulate the information that they share with the service provider. There, pricing controls are used by the service provider. In contrast, we focus on customers who are both delay sensitive and averse to untruthfulness, and we do not use pricing as a control.

*Behavioural queues.* Our work is broadly related to papers studying the behavioral foundations of queueing systems; see Allon and Kremer (2018) for background. Shunko et al. (2018) studies the behavioral impact of queue design on worker productivity in service systems that involve human servers. Buell (2021) identifies the negative effects of last-place aversion in queues. Armony et al. (2021) develop a game-theoretic model to assess the performance of pooling when behavioral servers choose their capacities strategically. Kim et al. (2020) study admission decisions in queues using behavioral models and controlled experiments. Wang and Zhou (2018) study how the queue configuration affects human servers' service time in a field experiment. Ülkü et al. (2020) investigate the relationship between waiting time and subsequent purchase decisions.

While there is ample empirical evidence illustrating the important role of social norms and preferences in queueing systems, Allon and Hanany (2012) is, to the best of our knowledge, the only work that investigates queueing intrusions with a formal mathematical model. In particular, Allon and Hanany (2012) show that the common observation of "queue jumping" can be part of social norms and can be explained on rational individual grounds. We depart from Allon and Hanany (2012) in three fundamental ways. First, Allon and Hanany focus on service systems where the manager is not involved in the way in which the queue is managed. Thus customers, through community enforcement, regulate the queue. In contrast, we focus our attention on unobservable priority queueing systems where the manager is the one who is in charge of controlling the queue, and customers do not see each other, and so are not able to prevent intrusions. Second, as Allon and Hanany focus on community enforcement, they study customer decision-making from the angle of social norms. In contrast, we study customer decision-making from the angle of social preferences by introducing the aversion to being untruthful into the customer decision-making process. Third, our methodological approach is different, as we conduct controlled experiments to shed light on the untruthful behavior of customers in queues.

*Intrinsic lying aversion.* The study of untruthful behavior is a complicated experimental endeavour. The paradigm in Fischbacher and Föllmi-Heusi (2013) is the most widely adopted in the

literature: It has been used in over 90 studies involving more than 44,000 subjects across 47 countries. Overall, these studies suggest that people exhibit an intrinsic lying aversion (Abeler et al. 2019). Özer et al. (2011) also assume that customers who deviate from truth-telling incur internal costs, however their setting on communicating forecasts in supply-chains is different from ours. In all studies adopting the FFH paradigm, *money* is used to incentivize participants to misreport their private information. In contrast, in priority queueing systems, people have an incentive to misreport private information (about their own priority levels) in order to shorten their *waiting times* in queue. It is unclear, a priori, whether people have similar cheating tendencies when saving time is the incentive, rather than making money. Indeed, it is well known that people do not treat time and money in an interchangeable manner (Leclerc et al. 1995).

## 3. Queueing Model

We model the service system as a single-server $M/G/1$ priority queuing system with two priority classes $k \in \{1, 2\}$. Customers who are assigned to priority class 1 receive non-preemptive priority over customers who are assigned to class 2. Customers arrive to the system according to a Poisson process with rate $\lambda$. An arriving customer has type $X$, where $X$ takes value $H$ with probability $p_H$, and value $L$ with probability $p_L = 1 - p_H$. The service times of customers with type $X = x$, where $x \in \{H, L\}$, are independent and identically distributed random variables with finite first and second moments $m_x$ and $n_x$, respectively. Customers are delay sensitive, and a customer of type $X = x$ has a per-time-unit waiting cost $c_x$. We define $\mu_x := 1/m_x$. Throughout this paper, we assume that $c_H \mu_H > c_L \mu_L$. This assumption implies that if customer priority levels were fully known, then serving high-priority customers first would minimize delay costs, as per the $c\mu$ rule. The traffic intensity is $\rho = \lambda(p_H m_H + p_L m_L)$. We assume that the system is stable, i.e., that $\rho < 1$.

We assume that customers have private information about their type $X$. Upon arrival, customers make a claim $Y = y \in \{H, L\}$ about their type. Because the system manager does not know the true customer types, it may be that $X \neq Y$, i.e., customers may be untruthful in their claims. In Appendix B, we consider a modelling extension where customers have uncertain information

about their own types, i.e., they hold some belief about their true type $X$. This setting may arise

in practice, e.g., in healthcare where patients may not know, with full certainty, the extent of their

illness. The main insights that we glean when customers are uncertain about their own types are

consistent with our insights in the main paper.

The system manager cannot observe individual customer types, but knows how customers behave

in aggregate. Also, customers know the system manager's claim-based prioritization policy. Wait-

ing times in the system, customer claims, and the system manager's prioritization policy, which

must be chosen optimally based on those claims, are all dependent on each other. Thus, a complex

equilibrium arises in the system, where customer claims must be consistent with the prioritiza-

tion policy and with the resulting waiting times. Here, we prove that there exists a unique Nash

equilibrium in the system, and we characterize that equilibrium.

## 3.1.    System Manager's Problem

The system manager defines a prioritization policy $\pi$ according to which customers will be served.

We restrict attention to static two-class priority policies, i.e., we assume that the system manager

will assign a priority level $k \in \{1, 2\}$ to customer claim $Y \in \{H, L\}$ with a certain probability. We

let $\Pi_Y$ denote the set of static two-class priority policies, based on the customer claims $Y$. In other

words, the manager must select routing probabilities to the two classes, 1 and 2, based on the

customer claims. The manager's problem is to minimize the total expected steady-state waiting

cost in the system, i.e.,

$$\underset{\pi \in \Pi_Y}{Min}\ C_\pi \equiv \underset{\pi \in \Pi_Y}{Min} \sum_{k \in \{1,2\}} \sum_{x \in \{H,L\}} \lambda_{k,x}^\pi c_x \mathbb{E}[W_k^\pi], \tag{1}$$

where $\lambda_{k,x}^\pi$ is the aggregate rate of arrivals of customers of true type $x$ who are assigned to priority

level $k$, under $\pi$, and $\mathbb{E}[W_k^\pi]$ is the resulting steady-state expected waiting time in priority level

$k$. That is, $\mathbb{E}[W_k^\pi]$ depends on $\lambda_{k,x}^\pi$ and on the composition of true $H$ and true $L$ customers in

priority level $k$. The system manager cannot observe the true type $x$ of each customer, but can

determine $\lambda_{k,x}^\pi$ and make inferences about the composition of true $H$ and true $L$ customers in each

priority level, based on the known (aggregate) customer behavior. Specifically, the arrival rate $\lambda_{k,x}^{\pi}$

is obtained by conditioning on customer claims as follows:

$$\lambda_{k,x}^{\pi} = \lambda \cdot \mathbb{P}(X=x) \cdot \sum_{y \in \{H,L\}} \mathbb{P}_{\pi}(\text{Claim } y \text{ is assigned to priority } k | X=x, Y=y)\mathbb{P}_{\pi}(Y=y|X=x),$$

where $\mathbb{P}_{\pi}(\text{Claim } y \text{ is assigned to priority } k | X=x, Y=y)$ is based on the system's manager's pri-

oritization policy $\pi$, and $\mathbb{P}_{\pi}(Y=y|X=x)$ captures a typical customer's claim $Y$, conditional on

her true type $X$. Since the system manager only observes the customer claims, and not their true

types, the assignment of claims to priority levels is independent of $X=x$, conditional on $Y=y$. In

Section 3.2, we explain how $\mathbb{P}_{\pi}(Y=y|X=x)$ arises at equilibrium in the customer choice problem.

For now, we treat these customer claim probabilities as given, and known to the system manager.

Using results on non-preemptive priority queues (see, e.g., Cobham 1954), the expected waiting

times in the two priority classes, $\mathbb{E}[W_1^{\pi}]$ and $\mathbb{E}[W_2^{\pi}]$, are given by:

$$\mathbb{E}[W_1^{\pi}] = \frac{\rho(n_H p_H + n_L p_L)}{2(p_H m_H(1 - \rho(1-\delta_H^{\pi})) + p_L m_L(1-\rho\delta_L^{\pi}))} \quad \text{and} \quad \mathbb{E}[W_2^{\pi}] = \frac{\mathbb{E}[W_1^{\pi}]}{1-\rho}, \tag{2}$$

where, depending on the routing policy $\pi$, $\delta_H^{\pi}$ is the steady-state fraction of misclassified $H$ cus-

tomers who are placed in priority level 2 (under-prioritization) and $\delta_L^{\pi}$ is the steady-state fraction

of misclassified $L$ customers who are placed in priority level 1 (over-prioritization). We can obtain

expressions for $\delta_H^{\pi}$ and $\delta_L^{\pi}$ by conditioning on customer claims:

$$\delta_H^{\pi} = \sum_{y \in \{H,L\}} \mathbb{P}_{\pi}(\text{Claim } y \text{ is assigned to priority } 2 | X=H, Y=y)\mathbb{P}_{\pi}(Y=y|X=H), \tag{3}$$

$$\delta_L^{\pi} = \sum_{y \in \{H,L\}} \mathbb{P}_{\pi}(\text{Claim } y \text{ is assigned to priority } 1 | X=L, Y=y)\mathbb{P}_{\pi}(Y=y|X=L). \tag{4}$$

As can be seen from (3) and (4), the misclassification errors, $\delta_H^{\pi}$ and $\delta_L^{\pi}$, depend on the customer

claims which, in turn, depend on the prioritization policy of the system manager. Thus, $\delta_H^{\pi}$ and $\delta_L^{\pi}$

arise in equilibrium, and the dependence between customer claims and the prioritization policy is

endogenously determined. In Section 4, we specify the optimal prioritization policy for the system

manager, and the optimal customer claim policy, which arise at equilibrium.

We close this section with Lemma 1, where we assume that $\delta_H^\pi$ and $\delta_L^\pi$ are exogenously given; in other words, we condition on some customer claim policy. Then, we study the monotonicity of the expected waiting cost as a function of those misclassification errors. Lemma 1 is similar to Lemma 7 in Singh et al. (2021) and Proposition 2 in Argon and Ziya (2009). We repeat this result here to draw parallels between our setting, where customers are strategic in their claims, and the models in those papers where customers do not have the ability to influence their priority classifications.

LEMMA 1. *The expected waiting cost in (1) is increasing in both $\delta_H^\pi$ and $\delta_L^\pi$. Moreover, it is more sensitive to $\delta_H^\pi$ than $\delta_L^\pi$, i.e., $\frac{\partial C_\pi}{\partial \delta_H^\pi} > \frac{\partial C_\pi}{\partial \delta_L^\pi} > 0$ if, and only if, $\delta_H^\pi + \delta_L^\pi < 1$. The additional sensitivity to $\delta_H^\pi$ over $\delta_L^\pi$, i.e., $\frac{\partial C_\pi}{\partial \delta_H^\pi} - \frac{\partial C_\pi}{\partial \delta_L^\pi}$, increases in $\rho$.*

Consistent with intuition, Lemma 1 shows that the waiting cost increases strictly in the misclassification errors $\delta_H^\pi$ and $\delta_L^\pi$. In Section 4, we show that the condition in Lemma 1 arises at equilibrium. In our setting, misclassification errors arise from the fact that customers make false claims about their types. This triggers allocation inefficiencies. Importantly, Lemma 1 demonstrates the asymmetric impact of under-prioritization errors, i.e., of routing a true $H$ customer to the lower priority level, particularly in congested systems. We will show that the system manager seeks to eliminate this adverse under-prioritization cost when selecting an optimal strategy.

## 3.2. Customer Problem

We now turn to specifying the way in which customers make claims. For a fixed prioritization policy $\pi$, define the function $g^\pi : \{H, L\} \to \{1, 2\}$, where $g^\pi(y)$ is the priority level assigned by the system manager to a customer who makes claim $y$. In our model, an arbitrary customer of true type $x$ decides on her claim $y$ in order to minimize her expected cost:

$$\underset{y \in \{H,L\}}{Min} \; c_x \mathbb{E}[W_{g^\pi(y)}^\pi] + \theta c_x (\mathbb{E}[W_x^\pi] - \mathbb{E}[W_y^\pi])^+. \tag{5}$$

We assume that customers are delay sensitive with delay cost $c_x$, depending on their true type $x$. Thus, there is an incentive for a customer to misreport her type if doing so leads to a shorter expected waiting time. Moreover, we assume that customers present *lying aversion*, modelled by

a *psychological cost* $\theta c_x(\mathbb{E}[W_x^\pi] - \mathbb{E}[W_y^\pi])^+$, which is incurred when the claim is untruthful. Here, $a^+ \equiv \max\{a, 0\}$ denotes the positive part function for a real number $a$. Indeed, there is a rich body of experimental data which strongly points to people's intrinsic preferences for truthfulness (Abeler et al. 2019) with heterogeneous lying aversion (Gibson et al. 2013). To capture heterogeneous preferences for truth-telling, we assume that each arriving customer is endowed with a random variate $\theta$, which represents their lying aversion, independently drawn from a distribution with support over some interval $[0, \bar{\theta}]$ where $\bar{\theta} > 1$. This is to ensure a high enough lying aversion so that not everyone consistently lies. We denote the probability density function of $\theta$ by $\phi$, and its cumulative distribution function by $\Phi$. We define the generalized failure rate of $\theta$ as $h(x) = x\phi(x)/(1 - \Phi(x))$; e.g., see Lariviere and Porteus (2001).

*Additional comments on the cost specification in (5).* We assume that a customer of true type $x$ who claims to be of type $y$ experiences a psychological cost $\theta c_x(\mathbb{E}[W_x^\pi] - \mathbb{E}[W_y^\pi])^+$. That is, our assumption is that down-reporting is not considered lying and that the psychological cost increases linearly in the material difference between the delay cost that corresponds to the customer's true priority level, $c_x\mathbb{E}[W_x^\pi]$, and the delay cost that corresponds to their claimed priority level, $c_x\mathbb{E}[W_y^\pi]$. This assumption is in line with experimental evidence in the Behavioral Economics literature, which suggests that people incur an intrinsic cost when they make an untruthful claim, and that this cost tends to increase in the "size" of the lie (Duch et al. 2021, Gneezy et al. 2018, Hilbig and Hessler 2013, Kajackaite and Gneezy 2017). We also recall that Özer et al. (2011) proposes a similar specification: There, a manufacturer incurs a psychological cost equal to $\beta|\hat{\xi} - \xi|$, where the parameter $\beta \geq 0$ captures the trustworthiness of the manufacturer, and $\xi$ and $\hat{\xi}$ represent the manufacturer's private forecast information, and the reported forecast information, respectively.

Our proposed expression for the "size of the lie", i.e., $c_x(\mathbb{E}[W_x^\pi] - \mathbb{E}[W_y^\pi])$, allows us to capture the intrinsic nature of lying aversion: This expression depends only on the customer's *true* priority level, and not on her *assigned* priority level. Unlike the true priority level, the assigned priority level depends on the system manager's prioritization policy. We emphasize that, because the manager

14

**Estrada, Ibrahim, and Zhan:** *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

may route customers inconsistently with their claims, the outcomes associated with lying or telling the truth are uncertain. It is far from straightforward to model lying costs in this case. For example, it is natural to ask: If a customer makes a deceitful claim, but is eventually assigned by the manager to her true priority level, should she incur a lying cost? Conversely, if a customer knows that she does not deserve to be given high priority, tells the truth, but the manager assigns her wrongly to the high-priority class, should she incur lying (guilt) costs? The literature is largely silent on lying behavior when there is payoff uncertainty; see Celse et al. (2019) for a notable recent exception. In our model, we assume that, because lying costs are intrinsic, people incur a cost when they lie, irrespective of the manager's scheduling policy. Ultimately, there is a need to provide experimental evidence substantiating our assumptions for the lying aversion cost function in (5): That is what we do in Section 5. There, we describe experimental data which supports the customer equilibrium behavior assuming the specification in (5). In so doing, we contribute to the nascent literature which models lying behavior under uncertain payoffs.

The customer decision model in (5) is appropriate to describe the customer's decision making in a setting where (i) the queue is unobservable to customers, and (ii) there are no reputation concerns since the service setting involves an anonymous, one-shot, interaction with the service provider. This is consistent with our motivating examples in the introduction. This said, it is important to note that a lying customer in our queueing setting does impose a negative externality on other customers in the queue. And, it is certainly conceivable that a customer may experience additional lying costs which are directly dependent on the size of this negative externality. Particularly, lying customers may incur higher lying costs because, in consequence of their lies, they are delaying high-priority customers who are truly deserving of fast service. Here, we focus solely on intrinsic lying costs which are independent of the size of that negative externality. Our modelling assumption is not unreasonable, for two main reasons. First, since queues are unobservable, such negative externality costs are more difficult to internalize. Indeed, this "out of sight, out of mind" phenomenon is well documented in Psychology (Schiano Lomoriello et al. 2018). Second, even if customers do

internalize costs about true high-priority customers who are delayed in consequence of their lies, then one can argue that the negative externality imposed is negligible in systems where there is enough capacity allocated to the high-priority class, so that high-priority waiting times are not too large. This is consistent with our motivating examples in the introduction. Since the negative externality imposed on high-priority customers, as a result of any one customer's lie, is indeed negligible, its associated costs could be assumed to be negligible as well. In general, the impact of negative externalities on lying behavior has been found to be ambiguous in the literature: While some studies show that such externalities have little effect (Fischbacher and Föllmi-Heusi 2013), others show that their effect depends on the stakes, in particular, if the stakes are high, then people's lying propensities are not affected by the externality (Gneezy and Kajackaite 2020). We relegate a study of the intricate effects of negative externalities to future work.

## 4. Optimal Scheduling Policy

In this section, we derive the optimal static prioritization policy for the system manager, and the optimal customer claim policy, which arise at equilibrium. The main result of this section is Proposition 1. For simplicity of exposition, we define the control variables $p_{1H}$ and $p_{1L}$ as follows:

$$p_{1H} := \mathbb{P}_\pi(k=1|Y=H) \text{ and } p_{1L} := \mathbb{P}_\pi(k=1|Y=L),$$

as the probabilities of giving high priority to a high and low claims, respectively. We let $p_{1H}^*$ and $p_{1L}^*$ denote corresponding optimal values, which arise at equilibrium. We recall that $h(x)$ is the generalized failure rate of $\theta$. The random variable $\theta$ is said to have an *increasing generalized failure rate* (IGFR) if $h(x)$ is weakly increasing in $x$ for $\Phi(x) < 1$; see Lariviere and Porteus (2001).

PROPOSITION 1. *If the lying aversion cost $\theta$ is IGFR, then there exists a unique Nash equilibrium which arises in the game between the customers and the manager, based on problems (1) and (5):*

*(a) All customers with true type $X = H$ make a truthful claim $Y = H$.*

*(b) A fraction $\Phi(1 - p_{1L}^*)$ of customers with true type $X = L$ are dishonest i.e., claim $Y = H$.*

16                   Estrada, Ibrahim, and Zhan: *On Customer (Dis)honesty in Priority Queues*

Article submitted to ; manuscript no. (Please, provide the manuscript number!)

(c) *The optimal routing policy is to assign high priority to all high-claim customers, i.e., $p_{1H}^* = 1$, and to assign high priority to low-claim customers with probability $p_{1L}^*$. We have that $p_{1L}^* = 0$ if, and only if, $h(1) \leq 1$, and $p_{1L}^* \in (0,1)$ as the unique solution to $h(1 - p_{1L}^*) = 1$, otherwise.*

The assumption that the lying aversion cost $\theta$ is IGFR in Proposition 1 is not restrictive because it captures most common distributions, e.g., uniform, exponential, normal, among many others.

Consistent with intuition, part (a) of Proposition 1 shows that customers of type $X = H$ always claim their true type. Part (b) shows that only a proportion (bounded away from 1) of $X = L$ type customers make a deceiving claim $Y = H$ in equilibrium. This observation is consistent with the extant literature where people consistently lie significantly less than they would have if they were only driven by materialistic incentives, e.g., making more money (Fischbacher and Föllmi-Heusi 2013). Proposition 1 shows that this lying proportion, $\Phi(1 - p_{1L}^*)$, is also independent of the waiting times in the system. This is also consistent with the extant literature: Abeler et al. (2019), and references therein, provide converging evidence that the average amount of lying does not significantly change as a function of monetary incentives. In Section 5, we validate those analytical results using data from controlled experiments.

## 4.1. Upgrading Low Claims

For the optimal routing policy, part (c) of Proposition 1 demonstrates that the system manager must prioritize customers according to their claims, under a condition which we interpret below. This is consistent with and provides theoretical support to the current "honor system" in place in the motivating examples of the introduction, where the manager routes customers according to their claims without additional checks. Interestingly, we find that it may also be optimal to *upgrade* low claims to the high-priority class, i.e., $p_{1L}^* > 0$, to further improve the system performance.

We note that, because all type $H$ customers claim truthfully, and the manager routes all high claims to the high-priority class, there is no under-prioritization (misclassifying true $H$ as $L$) in equilibrium, i.e., $\delta_H^{\pi^*} = 0$. Thus, allocation inefficiencies arise only from over-prioritization (misclassifying true $L$ as $H$) due to customer untruthfulness and because of upgraded customers, i.e.,

$\delta_L^{\pi^*} = \Phi(1 - p_{1L}^*) + p_{1L}^*(1 - \Phi(1 - p_{1L}^*))$, where $\Phi(1 - p_{1L}^*)$ is the proportion of true low class who claim high and are routed to the high class, and $p_{1L}^*(1 - \Phi(1 - p_{1L}^*))$ is the proportion of low claims who are upgraded to high priority. This result is consistent with Lemma 1, which demonstrates the asymmetry in cost between over-prioritizing and under-prioritizing. It is easy to see that upgrading low claims has an ambiguous effect on the over-prioritization error, $\delta_L^{\pi^*}$. On one hand, as the upgrading probability $p_{1L}$ increases, more low-type customers are given priority, which leads to an increase in over-prioritization and to additional delays in the high-priority class. On the other hand, increasing $p_{1L}$ also leads to a reduction in the the over-prioritization error by increasing the prevalence of truthfulness in the system. This is due to customers' lying aversion: An increase in the upgrading probability represents a higher incentive to report truthfully, since it entails a saving in the waiting cost, by being prioritized without incurring the cost of being untruthful. Thus, we observe that the lying probability, $\Phi(1 - p_{1L})$, decreases in $p_{1L}$.

Proposition 1 shows that an upgrading $p_{1L}^* \in (0, 1)$ such that $h(1 - p_{1L}^*) = 1$ balances this tension optimally, by minimizing the over-prioritization error. The generalized failure rate $h(\cdot)$ of the lying aversion cost $\theta$ has an intuitive economic interpretation. It represents the elasticity of the probability that a low-type customer claims truthfully, i.e., it measures how sensitive truthful behaviour is to the upgrading probability $p_{1L}$. Under the condition $h(1) \leq 1$, the probability that low-type customers will make a low claim is inelastic i.e., a change in $p_{1L}$ has little effect on truthful reporting, so that it would be better not to upgrade any low claims. For the case of a lying distribution such that $h(1) > 1$, truthful reporting is more sensitive to the upgrading control. In this case, the manager can resort to upgrading some low-claim customers to the high-priority class. The IGFR assumption on $\theta$ has an appealing implication. As the customer population is more honest, the aggregate truthful reporting behaviour becomes less elastic. The population is thus less responsive to the upgrading control (since, in any case, a significant proportion of customers already tells the truth). For further illustration, when $\theta$ has a uniform distribution over $[0, \bar{\theta}]$, the condition $h(1) \leq 1$ is equivalent to $\bar{\theta} \geq 2$. In other words, if customers exhibit strong enough lying

aversion, i.e., are honest enough, then it is optimal to route customers according to their claims. However, if $\bar{\theta} < 2$, i.e., if customers are sufficiently dishonest, then it is optimal to upgrade some low claims to the high class in order to incentivize honesty.

## 4.2. Incentivizing Honesty

Without customer claims, the system manager does not posses any information to be able to distinguish between customers. However, when customers make claims, and the system manager assigns priorities according to those claims, Proposition 1 shows that there is only partial misreporting of customer types, i.e., not all customers of true $L$ type claim to be of type $H$. Because of this partial dishonesty, we can easily show that customer claims are, themselves, informative at equilibrium, i.e., they give some indication of the true customer types. In particular, letting $\pi^*$ denote the optimal routing policy of the system manager, as described in Proposition 1, we can write:

$$\mathbb{P}_{\pi^*}(X = H | Y = H) = \frac{\mathbb{P}_{\pi^*}(X = H, Y = H)}{\mathbb{P}_{\pi^*}(Y = H)} = \frac{\mathbb{P}(X = H)}{\sum_{x \in \{H,L\}} \mathbb{P}_{\pi^*}(X = x, Y = H)} = \frac{p_H}{p_H + p_L \Phi(1 - p_{1L}^*)},$$

$$\mathbb{P}_{\pi^*}(X = L | Y = L) = 1.$$

Thus, $\mathbb{P}_{\pi^*}(X = H | Y = H) > p_H$ and $\mathbb{P}_{\pi^*}(X = L | Y = L) > p_L$, where we recall that $p_H$ and $p_L$ are the proportions of true $H$ and $L$ types.

Importantly, this optimal routing policy deviates from the celebrated $c\mu$ rule. In our game, the true priority level of an arriving customer (equivalently, the true index $c\mu$) is unknown to the system manager. Based on the available information, i.e., the customer claims, in equilibrium, the system manager can infer the "expected" $c\mu$ customer index, $I(y)$, for a customer who makes claim $Y = y$. Specifically, for a customer who claims $Y = y$, the system manager conditions on the true state of the customer, $X$, and calculates the expected $c\mu$ index, $I(y)$, as follows:

$$I(y) = c_H \mu_H \mathbb{P}_{\pi^*}(X = H | Y = y) + c_L \mu_L \mathbb{P}_{\pi^*}(X = L | Y = y) \quad \text{for} \quad y \in \{H, L\}.$$

Is it easy to see that $I(H) > I(L)$ if, and only if, $c_H \mu_H > c_L \mu_L$, which is our assumption throughout this paper. With exogenous customer classification errors, Argon and Ziya (2009) argues that customers should be prioritized in decreasing order of the expected $c\mu$ indices. Thus, if this were

to hold in our setting as well, then it would follow that arriving customers with higher indices $I(y)$ should always be served before customers with lower indices $I(y)$. In other words, no customer with a low claim should be upgraded to high priority. In our model, upgrading is optimal because customer claims are *endogenous*: By allowing for some low-claim customers to be in the high priority lane (i.e., by deviating from the $c\mu$ rule), the service provider incites customers to be more honest in their reporting, which ultimately benefits the system as a whole. In Section 5, we validate this key insight on the role of upgrading using data from controlled experiments.

To highlight the role of lying aversion in the system, we note that a standard economic model for customer reporting behaviour can be obtained by setting $\theta = 0$ in (5), i.e., by ignoring lying aversion. In this case, the equilibrium becomes trivial as all customers, independently of their true types, will make the same claim. This renders claims uninformative in equilibrium and, for whatever routing decision $(p_{1H}, p_{1L})$ of the system manager, the performance of the system is equivalent to a first-come-first-served (FCFS) discipline. This benchmark highlights the system-level effects that arise from considering lying aversion as a simple behavioural extension.

## 5. Experimental Evidence

In this section, we describe experimental evidence to support the theoretical insights of Section 4. In particular, our primary purpose here is to validate the customer equilibrium results in part (b) of Proposition 1. As such, we provide support to our modelling framework in Section 3, specifically to the assumed customer cost function in (5) which includes lying aversion.

We adopt the experimental paradigm in Fischbacher and Föllmi-Heusi (2013), which we hereafter refer to as the FFH paradigm, to investigate the extent to which people misreport their private information in order to wait less in a virtual queueing setting. In the FFH paradigm, participants privately observe the outcome of a random variable, e.g., the roll of a die. They are then asked to report that outcome and, subsequently, receive a payoff depending on their reported claims, e.g., the higher the claim, the higher the monetary reward. While dishonesty cannot be verified at the individual participant level (because the true die outcome is unknown to the experimenter),

one can make inferences about aggregate participant cheating behaviour, because the probability distribution of the random variable in the experiment is known to the experimenter. Importantly, participants can misreport their privately observed die outcomes with absolute impunity.

We conduct an online experiment that represents a one-shot individual decision-making situation where participants observe a given waiting time in each priority queue and a given upgrading probability and, based on this information, they make a decision to report honestly or to misreport their private information. Participants cannot observe the dynamic state of the queue, i.e., the number of participants in line, and the real-time behaviour of other participants. This is consistent with our motivating examples in the introduction. We take here the perspective of arriving customers, and focus on understanding how customers would make a claim, conditional on the observed waiting times in the two priority levels and on the scheduling policy, in particular the proportion of low claims that can be upgraded. That is, we do not study how the queueing system transitions into equilibrium, but rather we design the queueing setting to mimic the equilibrium and test whether participants' reporting behaviour is in line with our model-based results.

### 5.1.   Experimental Design

We begin by describing our design and procedure for the online queueing experiment.

*Pre-registration.* We set the target sample size for the experiment and our analysis plans a priori. We pre-registered our experiment, and the corresponding As Predicted document can be found at: `https://aspredicted.org/blind.php?x=7u6ar3`.

*Participants and exclusions.* A total of 2,373 participants (44.33% female, mean age $M_{age} = 37.97$, standard deviation SD = 11.70) were recruited on the Amazon Mechanical Turk (MTurk) platform. Participants with at least 0.95 HIT approval ratio (proportion of completed tasks) were recruited to take part in our queueing experiment. Participants were instructed that they must wait in a virtual queue, then answer a two-question survey in exchange for a 1 US dollar payment (the payment is independent of the wait time in queue). Participants were informed that the experiment could take up to 30 minutes. To ensure the independence of our observations, we exclude from our analysis the responses of 6 participants that presented the same Internet Protocol (IP) address.

In our experiment we have a completion rate of 95%. We exclude participants that did not complete the experiment. Moreover, while waiting in the virtual queue, and to ensure that participants do not engage in other activities while waiting, they asked to click a button that appears every 60 seconds in order to move ahead in the queue. The remaining waiting time for participants in a given queue stops from elapsing until they click that button. This mechanism is used to ensure that participants experience a real waiting cost. We recorded the time that participants took to click each button and excluded participants that, on average, took longer than 30 seconds to click those buttons once they appeared. From our original sample, 90% of participants took, on average, less than 30 seconds to click the buttons, and the mean and median average click time was 14 seconds and 4 seconds, respectively. After all those exclusions, we are left with a sample of 2,021 participants (45.47% female, mean age $M_{age} = 38.23$, standard deviation SD = 11.99). Our results are unchanged by the aforementioned exclusions; see Appendix D.

*Experimental procedure.* In our experimental design, participants are asked to privately roll a six-sided die and to record the outcome on a piece of paper. To avoid any fear of detection, participants are instructed to roll a die at home or, alternatively, to use Google's virtual die at the following link `https://www.google.com/search?q=dice+roller`. We do not observe the die rolls. Participants are randomly assigned to one of nine experimental conditions (see Table 1) which differ in the waiting times for the Short Queue and the Long Queue, and in the upgrading probability $p_{1L}$.

After participants roll the die and write down the outcome on a piece of paper, they are presented with the waiting times of two queues: Short Queue and Long Queue, according to their randomly-assigned experimental condition. They are instructed that if they report a number 5 then they will wait in the short queue, and if they report any number different from 5, then they will wait in the long queue with probability $1 - p_{1L}$ and in the short queue with probability $p_{1L}$.

Participants are asked, as part of the instructions for the experiment, some questions related to this rule. Participants are able to finish the instructions section only if those questions are correctly answered. Moreover, before they roll the die and report any number, they are placed in a practice

| Condition | Short Queue | Long Queue | Upgrading prob $p_{1L}$ | Sample |
|:---:|:---:|:---:|:---:|:---:|
| 1 | 2 min | 5 min | 0 | 226 |
| 2 | 2 min | 10 min | 0 | 220 |
| 3 | 2 min | 15 min | 0 | 217 |
| 4 | 2 min | 5 min | 0.5 | 222 |
| 5 | 2 min | 10 min | 0.5 | 227 |
| 6 | 2 min | 15 min | 0.5 | 222 |
| 7 | 2 min | 5 min | 0.9 | 220 |
| 8 | 2 min | 10 min | 0.9 | 233 |
| 9 | 2 min | 15 min | 0.9 | 234 |

**Table 1**      **Description of the experimental conditions.**

queue for 2 minutes. This ensures that participants understand what it feels like to wait in the virtual queue. Throughout the experiment, there is no mention of deception, honesty or any related concepts. After participants roll the die, report a number, and wait in queue, they answer two simple questions related to their experience in the queue, which concludes the experiment.

## 5.2. Derivation of Hypotheses

*Lying aversion.* For a routing policy with upgrading probability $p_{1L}$, part (b) of Proposition 1 shows that the proportion of people who misreport, i.e., claim to be of a high type when they are truly of a low type, is equal to $\Phi(1 - p_{1L})$, where $\Phi$ is the cumulative distribution function of the lying aversion cost $\theta$. This proportion is bounded away from 1, and arises because of the lying aversion cost in (5). This partial lying result leads to our first hypothesis.

H1. The proportion of participants who misreport their private information in order to wait in the short queue is bounded away from 1.

Hypothesis H1 is consistent with results in the extant literature, where money is used as the incentive to be untruthful, rather than time. Specifically, one central finding in Abeler et al. (2019)

is that, despite a clear monetary incentive to deviate from telling the truth, a significant proportion of participants are honest in their reporting. This finding also holds when other experimental designs are adopted as well (Ariely 2012). Our online doubly blind setting is ideal to test for the existence of intrinsic lying costs since it excludes the possibility of other confounding costs such as reputation costs, negative externality, perception by others, etc. Thus, we attribute the truthful behavior of participants to that intrinsic lying aversion, as in (5).

*Insensitivity to waiting times.* The expression for the lying proportion, $\Phi(1 - p_{1L})$ in part (b) of Proposition 1, also demonstrates that the proportion of low-priority customers who claim to be of high priority depends only on the underlying distribution of the lying aversion parameter, $\theta$, and not on the waiting times in the system. This leads to the formulation of our second hypothesis:

H2. Increasing the differences in waiting times between the two queues will not significantly increase the proportion of people who misreport their true type.

Hypothesis H2 is also consistent with results in the extant literature which uses money as the incentive, rather than time. Interestingly, Abeler et al. (2019) shows that the proportion of people who misreport their private information does not increase as the monetary incentive to do so increases; see chapter 1 in Ariely (2012) for similar evidence with a different experimental design. In particular, Abeler et al. (2019), and references therein, provide converging evidence that the average amount of lying does not change much if monetary stakes are increased from a few cents to 50 USD, a 500-fold increase! As discussed in Fischbacher and Föllmi-Heusi (2013), that can be explained as follows: The increased (and desirable) payoff associated with misreporting seems to be balanced by the psychological cost of that misreporting. In other words, while participants have a monetary incentive to lie when payoffs are increased, they also feel worse about lying when those stakes are higher. Indeed, it has been consistently observed that dishonesty costs tend to increase in the "size" of the lie (Duch et al. 2021, Gneezy et al. 2018), as per the model in (5). In Appendix E, we present an additional experiment, with shorter waits, which provides further support to our model with respect to Hypotheses H1 and H2.

*Lying with uncertain payoffs.* Little is known about how the psychological cost of misreporting is shaped when its consequences become uncertain. Celse et al. (2019) find that participants' dishonesty seems to be independent from payoff uncertainty. While Celse et al. (2019) explore the uncertainty in the payoff due to misreporting, we are interested in investigating the effect of uncertainty in the payoff due to telling the truth, i.e., the effect of the upgrading probability $p_{1L}$. If the finding in Celse et al. (2019) were to extend to our setting, then we would expect that changing the probability of being assigned to the short queue when claiming low priority, $p_{1L}$, would not affect the misreporting behaviour of participants. That is, the probability of reporting the number 5 should be the same across all experimental conditions. This leads to our third hypothesis:

H3. The proportion of participants who misreport their private information is not significantly affected by the probability of being upgraded to the short queue when reporting low priority.

We emphasize that hypothesis H3 lies in contrast with our results in part (b) of Proposition 1. There, we found that the proportion of people who misreport, $\Phi(1-p_{1L})$, decreases in the upgrading probability $p_{1L}$. Indeed, in our model, the probability $p_{1L}$ can be used as a tool to incentivize people to report truthfully, since it entails a saving in the waiting cost, by being potentially prioritized, without incurring the intrinsic psychological cost of being untruthful. Testing Hypotheses H1, H2, and H3 with experimental data in order to substantiate our analytical results lays out a solid, data-driven, micro-foundation for modelling customer dishonesty in priority queues.

## 5.3.   Empirical Strategy and Results

We now describe our experimental results. We relegate additional tables with detailed results to Appendix C. In Figure 1, we present the proportion of participants that reported the number 5 across experimental conditions; see also Table 2 in Appendix C. Recall that participants have the incentive to report the number 5 to reduce their waiting time in the queue.

We run two-sided exact binomial tests for the proportion of participants that reported the number 5 compared to the proportion that would have reported the number 5 under full honesty i.e., 1/6.
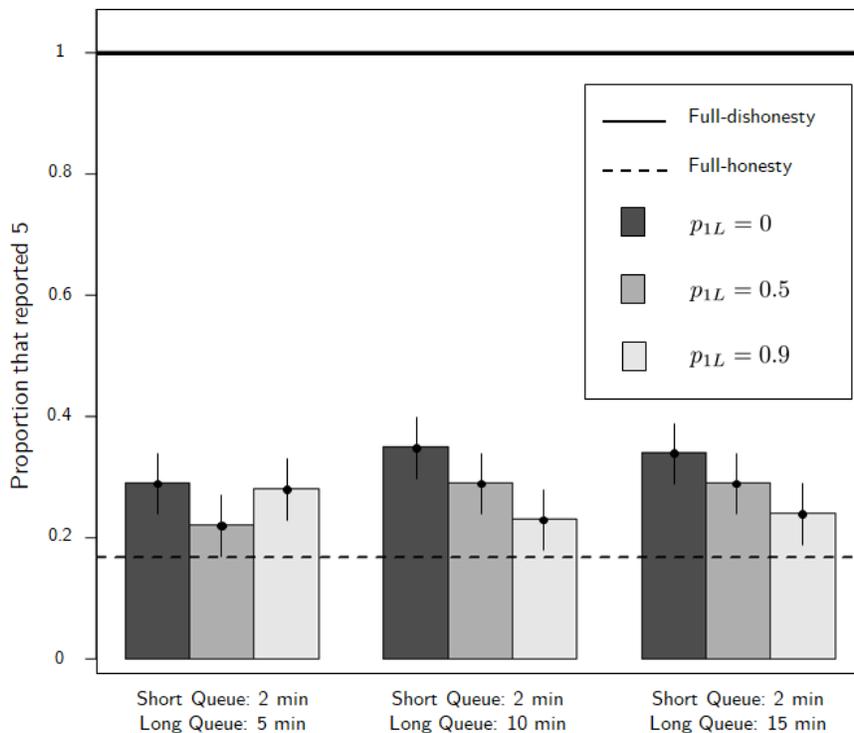
**Figure 1** **Proportions of participants that reported the number 5 across experimental conditions.**

We find that untruthfulness is significant under all conditions (p-values $< 0.05$). This can be readily seen in Figure 1, where the black horizontal dashed line represents the expected proportion of customers who would have reported the number 5 under the assumption of full honesty. Any value above the dashed line is based on untruthfulness, in expectation.

*Testing Hypothesis H1: The existence of lying aversion.* We find that participants misreport surprisingly little across all experimental conditions. It is clear from Figure 1 that the proportion of participants that reported the number 5 is far from 1, i.e., the full-dishonesty black horizontal solid line in the figure. Indeed, for all conditions, the two-sided exact binomial tests for the proportion of participants who reported the number 5 compared to the proportion that would have reported the number 5 under full dishonesty (i.e., 1) are significant (p-values $= 0$). In the absence of any possible punishments or reputation concerns, any reluctance to misreport can be attributed to an internal psychological cost due to lying (Fischbacher and Föllmi-Heusi 2013). This result strongly supports the existence of lying aversion, as per our model in (5).

26

Estrada, Ibrahim, and Zhan: *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

*Testing Hypothesis H2: Insensitivity to waiting times.* We run a Generalized Cochran-Mantel-Haenszel test (Agresti 2003) for the conditional association between the proportion of participants that reported the number 5 and the time incentive treatment, conditional on the levels of the upgrade treatment. We find that the conditional association is not significant ($M^2 = 1.36$, $df = 2$, p-value $= 0.51$). We also run a Chi-square test for the marginal association between the proportion of participants that reported the number 5 and the time incentive treatment. We find that the marginal association is not significant ($\chi^2 = 1.22$, $df = 2$, p-value $= 0.54$). These results indicate that a change in the time incentive does not affect significantly the misreporting behaviour of participants. Moreover, for each upgrade level (i.e., $p_{1L} = 0$, $p_{1L} = 0.5$, $p_{1L} = 0.9$), we run Jonckheere–Terpstra tests for an order, both ascending and descending. We find that the proportion of participants that claim the number 5 does not significantly increase monotonically (p-values $= 0.20, 0.11, 0.78$), nor decrease monotonically (p-values $= 0.79, 0.89, 0.22$) as the time incentive increases in each of the respective upgrading levels. Also, by aggregating the data (ignoring the upgrade treatment levels), both the ascending (p-value $= 0.24$) and descending (p-value $0.75$) trends are not significant in the Jonckheere–Terpstra test. Finally, we run a logistic regression (see Table 3 in Appendix C) and find that the coefficients for the different levels of the time incentive treatment are not significant in any of the model specifications. Overall, these results support the claim that a change in the time incentive does not affect participant untruthfulness behaviour, which supports hypothesis H2 and is in line with our equilibrium results in part (b) of Proposition 1.

*Testing Hypothesis H3: Lying with uncertain payoffs.* We run a Generalized Cochran-Mantel-Haenszel test for the conditional association between the proportion of participants that reported the number 5 and the upgrading treatment, conditional on the levels of the time incentive. We find that the conditional association is significant ($M^2 = 12.16$, $df = 2$, p-value $= 0.002$). We also run a Chi-square test for the marginal association between the proportion of participants that reported the number 5 and the upgrading treatment. We find that the marginal association is significant ($\chi^2 = 12.03$, $df = 2$, p-value $= 0.002$). These results indicate that the upgrading control affects

participant untruthfulness behaviour. For each time incentive level (i.e., $2-5$ min, $2-10$ min, $2-15$ min), we run Jonckheere–Terpstra tests for an order, both ascending and descending. We find that the proportion of participants who claim the number 5 does not significantly increase monotonically (p-values $= 0.61, 0.99, 0.97$) as the upgrading increases in any of the respective time incentive levels. Moreover, while that proportion does not significantly decrease monotonically for the low time incentive condition $2-5$ min (p-value $= 0.39$), it does significantly decrease monotonically as the upgrade increases for the higher wait-time conditions, i.e., $2-10$ min (p-value $= 0.009$) and $2-15$ min (p-value $= 0.03$). Also, by aggregating the data (ignoring the incentive treatment levels), the proportion of participants that claim 5 significantly decreases monotonically as the upgrade probability increases (p-value 0.005); this is also visually evident in Figure 1. Our results point to a potential interaction between the time incentive and the upgrading control; we describe statistical tests for this interaction in what follows.

To test for such an interaction, we run logistic regressions (see Table 3 in Appendix C) and find that the interactions between the time incentive treatment and the upgrading treatment are not significant. Based on both the Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC), we find that the model specification that best fits the data, among the ones considered in Table 3, is Model 1 which includes an intercept, demographic covariates (insignificant in our data), and the upgrading probability covariate. Moreover, we conduct likelihood ratio tests between Model 1 and Model 2 ($LR = 1.43$, $df = 2$, p-value $= 0.49$), and Model 1 and Model 3 ($LR = 7.56$, $df = 6$, p-value 0.27) which show that the fit of Model 1 does not significantly improve by adding the time incentive predictor of Model 2 and the interaction terms between time incentive and upgrading of Model 3. Finally, the likelihood ratio test between Model 1 and the null model including only an intercept term ($LR = 13.86$, $df = 4$, p-value 0.007) shows that the fit of Model 1 is significantly better. These results show that Model 1 fits our experimental data the best. We also observe, in Table 3, that the coefficients for the upgrading covariates at both levels, $p_{1L} = 0.5$ and $p_{1L} = 0.9$, are significant and negative, and that the coefficient and marginal effects for

28

Estrada, Ibrahim, and Zhan: *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

$p_{1L} = 0.9$ are larger in absolute value than for $p_{1L} = 0.5$. This indicates that upgrading leads to a decrease in the lying probability, as per our results in part (b) of Proposition 1. Overall, we find sufficient statistical evidence in our data to reject hypothesis H3, i.e., to support the insight that untruthfulness decreases in the upgrading probability.

## 6. Conclusions

We studied the extent to which people misreport their private information in a two-class priority queueing system, where priorities are assigned according to customer claims, and where these claims are not verifiable nor punishable. Although standard economic theory predicts that all customers will claim that they should be given high priority, we showed, through the analysis of a theoretical queueing model and controlled experiments, that human behavior deviates far from this prediction. We theorized and experimentally validated the intrinsic costs that people incur when they lie.

Priority systems which rely on customer claims are prevalent in practice. Our results highlight the informative value of those claims, even in systems where customers can lie with total impunity. Thus, we provided theoretical and experimental evidence that customer claims should be sought by the system manager. This is particularly relevant to priority settings where market mechanisms, e.g., pricing the relevant services, are not applicable. In particular, we found theoretical evidence that supports the current "honor scheduling system" (routing customers according to their claims) which is in place in many real settings. Additionally, we found that, when the lying aversion cost is sufficiently elastic, a different prioritization rule, which deviates from the celebrated $c\mu$ rule, is optimal. Essentially, this rule prescribes upgrading some low-priority claims to the high-priority class in order to incentivize honesty in the system as a whole.

*Future research.* The study of untruthfulness in queues is rich and strongly motivated by practice. We hope that our work will spark more interest in this domain. Our experimental results are based on simple two-priority queueing settings. We advocate conducting more experimental work to further understand the role of intrinsic preferences for truth-telling in the reporting behaviour of people in more complicated queueing settings. Future work can, for example, investigate the extent

**Estrada, Ibrahim, and Zhan:** *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

29

to which people report untruthfully in a multi-priority queueing system, and unveil the relevant determinants in that case. While we foresee that some people will be honest and others will be untruthful, it is not clear, a priori, how customer claims will be distributed across the multiple priority queues: For example, will there be partial lying in this case?

In this paper, we focused on priority settings where the dynamic state of the queue is unobservable to customers, i.e., the length of the queue and the real-time behavior of other customers. We also assumed the absence of reputation concerns since the service setting involves an anonymous, one-shot, type of interaction with the service provider. These properties are common features shared by many priority queueing settings, but certainly not all systems. Future research could focus on understanding customer untruthfulness in priority queues where the dynamic state of the queue is observable, e.g., in emergency departments where patients may exaggerate their own symptoms at the triage stage to receive medical service faster. In such setting, where customers observe and are observed by others, reputation and negative externality concerns may play a major role.

# References

Abeler, J., A. Becker, and A. Falk (2014). Representative evidence on lying costs. *Journal of Public Economics 113*, 96–104.

Abeler, J., D. Nosenzo, and C. Raymond (2019). Preferences for truth-telling. *Econometrica 87*(4), 1115–1153.

Agresti, A. (2003). *Categorical data analysis*, Volume 482. John Wiley & Sons.

Allon, G. and E. Hanany (2012). Cutting in line: Social norms in queues. *Management Science 58*(3), 493–506.

Allon, G. and M. Kremer (2018). Behavioral foundations of queueing systems. *The Handbook of Behavioral Operations 9325*, 325–366.

Argon, N. T. and S. Ziya (2009). Priority assignment under imperfect information on customer type identities. *Manufacturing & Service Operations Management 11*(4), 674–693.

Ariely, D. (2012). *The (honest) truth about dishonesty: : how we lie to everyone–especially ourselves.* New York: Harper Collins Publishers.

30

**Estrada, Ibrahim, and Zhan:** *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

Armony, M., G. Roels, and H. Song (2021). Pooling queues with strategic servers: The effects of customer ownership. *Operations Research 69*(1), 13–29.

BBC (2020a). Coronavirus: People without symptoms 'misusing testing'. `https://www.bbc.co.uk/news/health-54088206`.

BBC (2020b). Matt hancock has urged coronavirus tests to be prioritised for key workers.jade says she lied about her employment status to get a test. `https://twitter.com/BBCBreakfast/status/1306109442295988224`.

Becker, G. S. (1968). Crime and punishment: An economic approach. In *The economic dimensions of crime*, pp. 13–68. Springer.

Bren, A. and S. Saghafian (2019). Data-driven percentile optimization for multiclass queueing systems with model ambiguity: Theory and application. *INFORMS Journal on Optimization 1*(4), 267–287.

Buell, R. W. (2021). Last-place aversion in queues. *Management Science 67*, 1430–1452.

Celse, J., S. Max, W. Steinel, I. Soraperra, and S. Shalvi (2019). Uncertain lies: How payoff uncertainty affects dishonesty. *Journal of Economic Psychology 71*, 117–125.

Cobham, A. (1954). Priority assignment in waiting line problems. *Journal of the Operations Research Society of America 2*(1), 70–76.

Cox, D. R. and W. Smith (1991). *Queues*, Volume 2. CRC Press.

Duch, R. M., D. Laroze, and A. Zakharov (2021). The moral cost of lying. Technical report, University of Oxford.

Fischbacher, U. and F. Föllmi-Heusi (2013). Lies in disguise—an experimental study on cheating. *Journal of the European Economic Association 11*(3), 525–547.

Gibson, R., C. Tanner, and A. F. Wagner (2013). Preferences for truthfulness: Heterogeneity among and within individuals. *American Economic Review 103*(1), 532–48.

Gneezy, U. and A. Kajackaite (2020). Externalities, stakes, and lying. *Journal of Economic Behavior & Organization 178*, 629–643.

Gneezy, U., A. Kajackaite, and J. Sobel (2018). Lying aversion and the size of the lie. *American Economic Review 108*(2), 419–53.

**Estrada, Ibrahim, and Zhan:** *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

31

Hassin, R. and M. Haviv (2003). *To queue or not to queue: Equilibrium behavior in queueing systems*, Volume 59. Springer Science & Business Media.

Heather, B. (2017). Babylon pilot was not pursued by ccgs. `https://www.hsj.co.uk/technology-and-innovation/babylon-pilot-was-not-pursued-by-ccgs/7021131.article`.

Hilbig, B. E. and C. M. Hessler (2013). What lies beneath: How the distance between truth and lie drives dishonesty. *Journal of Experimental Social Psychology 49*(2), 263–266.

Hu, M., R. Momot, and J. Wang (2021). Privacy management in service systems. Technical report, University of Toronto.

Isaiah, A. (2020). Skip ps5 direct queue: Youtuber alleges glitch on sony's website allows users to bypass the line. `https://www.techtimes.com/articles/254938/20201210/skip-ps5-direct-queue-bypass-using-glitch-sonys-website.htm`.

Jones, J. (2020). Take it from a paramedic: the ambulance service can't keep up with demand. `https://www.theguardian.com/commentisfree/2020/jan/29/paramedic-ambulance-service-call-triage`.

Kajackaite, A. and U. Gneezy (2017). Incentives and cheating. *Games and Economic Behavior 102*, 433–444.

Kim, S.-H., J. Tong, and C. Peden (2020). Admission control biases in hospital unit capacity management: How occupancy information hurdles and decision noise impact utilization. *Management Science 66*(11), 5151–5170.

Kirton, J. A., W. Thompson, M. Pearce, and J. M. Brown (2020). Ability of the wider dental team to triage patients with acute conditions: a qualitative study. *British dental journal 228*(2), 103–107.

Lariviere, M. A. and E. L. Porteus (2001). Selling to the newsvendor: An analysis of price-only contracts. *Manufacturing & service operations management 3*(4), 293–305.

Leclerc, F., B. H. Schmitt, and L. Dube (1995). Waiting time and decision making: Is time like money? *Journal of Consumer Research 22*(1), 110–119.

NPR (2021a). How people are jumping the covid-19 vaccine line. `https://www.npr.org/sections/health-shots/2021/02/09/965841419/does-loose-enforcement-of-vaccine-eligibility-rules-encourage-line-jumping`.

NPR (2021b). Is it ever ok to jump ahead in the vaccine line? `https://www.npr.org/sections/`
`health-shots/2021/02/06/964139633/is-it-ever-ok-to-jump-ahead-in-the-vaccine-line?`
`t=1623992132881`.

Özer, Ö., Y. Zheng, and K.-Y. Chen (2011). Trust in forecast information sharing. *Management Science 57*(6), 1111–1137.

Schiano Lomoriello, A., F. Meconi, I. Rinaldi, and P. Sessa (2018). Out of sight out of mind: Perceived physical distance between the observer and someone in pain shapes observer's neural empathic reactions. *Frontiers in psychology 9*, 1824.

Shunko, M., J. Niederhoff, and Y. Rosokha (2018). Humans are not machines: The behavioral impact of queueing design on service time. *Management Science 64*(1), 453–473.

Singh, S., I. Gurvich, and J. Van Mieghem (2021). Feature-based design of priority queues: Digital triage in healthcare. Technical report, Northwestern University.

Storm, D. (2015). 200 virtual line jumpers exploited 'technical backdoor' to score burning man tickets. `https://www.computerworld.com/article/2887148/` `200-virtual-line-jumpers-created-technical-backdoor-to-score-burning-man-tickets.` `html`.

Ülkü, S., C. Hydock, and S. Cui (2020). Making the wait worthwhile: Experiments on the effect of queueing on consumption. *Management Science 66*(3), 1149–1171.

Van der Zee, S. and H. Theil (1961). Priority assignment in waiting-line problems under conditions of misclassification. *Operations Research 9*(6), 875–885.

Van Mieghem, J. A. (1995). Dynamic scheduling with convex delay costs: The generalized c— mu rule. *The Annals of Applied Probability 5*(3), 809–833.

Wang, J. and Y.-P. Zhou (2018). Impact of queue configuration on service time: Evidence from a supermarket. *Management Science 64*(7), 3055–3075.

Weaver, M. and K. Proctor (2020). No checks to be made on essential-worker status for uk covid-19 tests. `https://www.theguardian.com/world/2020/apr/24/` `no-checks-to-be-made-on-essential-worker-status-for-uk-covid-19-tests`.

The appendices are organized as follows. In Appendix A, we present proofs for the technical results in the paper. In Appendix B, we present a modelling extension where customers are uncertain about their own types. In Appendix C, we present tables with detailed experimental results summarized in the paper. In Appendix D, we present the analysis of experimental results without the paper's exclusions. In Appendix E, we present the results of an additional experiment with shorter waits, in further support for hypotheses H1 and H2, i.e., where we keep the upgrading probability constant at 0.

## Appendix A:    Technical Proofs

We define $p_{LH} := \mathbb{P}(Y = L|X = H)$ and $p_{HL} := \mathbb{P}(Y = H|X = L)$, which capture the aggregate customer claiming behaviour. We let $p_{LH}^*$ and $p_{HL}^*$ denote the corresponding values at optimum. In what follows, for ease of notation, we often drop dependence on the manager's scheduling policy $\pi$.

### A.1.    Lemma 1

PROOF. Recall that, by assumption, we have that $c_H\mu_H > c_L\mu_L \iff c_Hm_L - c_Lm_H > 0$. We can write the expected waiting cost in the system as:

$$C_s(\delta_H, \delta_L) = g(\rho)f(\delta_H, \delta_L),$$

where

$$g(\rho) = \left(\frac{\rho^2}{1-\rho}\right)\left(\frac{p_Hn_H + p_Ln_L}{2(p_Hm_H + p_Lm_L)}\right),$$

$$f(\delta_H, \delta_L) = \frac{p_Hc_H(1-\rho(1-\delta_H)) + p_Lc_L(1-\rho\delta_L)}{p_Hm_H(1-\rho(1-\delta_H)) + p_Lm_L(1-\rho\delta_L)}.$$

After some algebraic manipulations we find:

$$\frac{\partial f}{\partial \delta_H} = \frac{\rho p_Hp_L(c_Hm_L - c_Lm_H)(1-\rho\delta_L)}{(p_Hm_H(1-\rho(1-\delta_H)) + p_Lm_L(1-\rho\delta_L))^2},$$

$$\frac{\partial f}{\partial \delta_L} = \frac{\rho p_Hp_L(c_Hm_L - c_Lm_H)(1-\rho(1-\delta_H))}{(p_Hm_H(1-\rho(1-\delta_H)) + p_Lm_L(1-\rho\delta_L))^2}.$$

Since $\rho \in (0,1)$ and $\delta_L, \delta_H \in [0,1]$, it follows that $\frac{\partial C_s}{\partial \delta_H} = g(\rho)\frac{\partial f}{\partial \delta_H} > 0$ and $\frac{\partial C_s}{\partial \delta_L} = g(\rho)\frac{\partial f}{\partial \delta_L} > 0$.

We define $\Delta f' \doteq \frac{\partial f}{\partial \delta_H} - \frac{\partial f}{\partial \delta_L}$. After some algebraic manipulations we find:

$$\Delta f' = \frac{p_Hp_L(c_Hm_L - c_Lm_H)\rho^2(1-\delta_H - \delta_L)}{(p_Hm_H(1-\rho(1-\delta_H)) + p_Lm_L(1-\rho\delta_L))^2}.$$

Since $\rho \in (0,1)$ we have that $\Delta C_s' = \frac{\partial C_s}{\partial \delta_H} - \frac{\partial C_s}{\partial \delta_L} = g(\rho)\Delta f' > 0$ if and only if $\delta_H + \delta_L < 1$.

Finally, taking the following derivatives:

$$\frac{\partial g(\rho)}{\partial \rho} = \frac{2-\rho}{(1-\rho)\rho}g(\rho),$$

$$\frac{\partial \Delta f'}{\partial \rho} = \frac{2 p_H p_L (m_H p_H + m_L p_L)(c_H m_L - c_L m_H)\rho(1 - \delta_H - \delta_L)}{(p_H m_H (1 - \rho(1 - \delta_H)) + p_L m_L (1 - \rho \delta_L))^3},$$

and recalling that $\Delta C_s' = \frac{\partial C_s}{\partial \delta_H} - \frac{\partial C_s}{\partial \delta_L} = g(\rho)\Delta f'$, we have that

$$\frac{\partial \Delta C_s'}{\partial \rho} = \frac{\partial g(\rho)}{\partial \rho}\Delta f' + g(\rho)\frac{\partial \Delta f'}{\partial \rho}.$$

Since $\rho \in (0,1)$, we have that $g(\rho) > 0$ and $\frac{\partial g(\rho)}{\partial \rho} > 0$. Moreover, notice that the sign of both $\frac{\partial \Delta f'}{\partial \rho}, \Delta f$ is

dictated by the sign of $1 - \delta_H - \delta_L$. It follows that $\frac{\partial \Delta C_s'}{\partial \rho} > 0$ if and only if $\frac{\partial \Delta f'}{\partial \rho}, \Delta f' > 0 \iff \delta_H + \delta_L < 1$. ∎

### A.2. Proposition 1

PROOF. There are 3 different cases to analyse: (1) when $p_{1H} < p_{1L}$, (2) when $p_{1H} = p_{1L}$, and (3) when $p_{1H} >$

$p_{1L}$. Recall that $p_{1H} = \mathbb{P}_\pi(k=1|Y=H)$ and $p_{1L} = \mathbb{P}_\pi(k=1|Y=L)$ are the manager's routing probabilities,

and $p_{LH} = \mathbb{P}(Y=L|X=H)$ and $p_{HL} = \mathbb{P}(Y=H|X=L)$ are the customer claim probabilities.

**Case 1: $p_{1H} < p_{1L}$.**

When an arbitrary customer type $x = H$ claims $y = H$ she has expected cost $c_H \mathbb{E}[W_1^\pi] p_{1H} + c_H \mathbb{E}[W_2^\pi](1 -$

$p_{1H})$. Conversely, when she claims $y = L$ she has expected cost $c_H \mathbb{E}[W_1^\pi] p_{1L} + c_H \mathbb{E}[W_2^\pi](1 - p_{1L})$. For this

case, it is easy to see that since $\mathbb{E}[W_2^\pi] \geq \mathbb{E}[W_1^\pi]$, independently on other customers' claims, the arbitrary

customer with type $x = H$ will always claim to be a low type, i.e., $p_{LH} = 1$.

When an arbitrary customer type $x = L$ claims $y = H$ she has expected cost $c_L \mathbb{E}[W_1^\pi] p_{1H} + c_L \mathbb{E}[W_2^\pi](1 -$

$p_{1H}) + \theta c_L(\mathbb{E}[W_2^\pi] - \mathbb{E}[W_1^\pi])$. Conversely, when she claims $y = L$ she has expected cost $c_L \mathbb{E}[W_1^\pi] p_{1L} +$

$c_L \mathbb{E}[W_2^\pi](1 - p_{1L})$. For this case, it is easy to see that since $\mathbb{E}[W_2^\pi] \geq \mathbb{E}[W_1^\pi]$, independently on other customers'

claims, the arbitrary customer with type $x = L$ will always claim to be a low type.

The manager correctly anticipates that high type customers will claim to be low type with probability

$p_{LH} = 1$ and low type customers will claim to be high type with probability $p_{HL} = 0$. Based on this, the

manager anticipates the following prioritization error probabilities:

$$\delta_H(p_{1H}, p_{1L}) = p_{LH}(1 - p_{1L}) + (1 - p_{LH})(1 - p_{1H}) = 1 - p_{1L},$$

$$\delta_L(p_{1H}, p_{1L}) = p_{HL} p_{1H} + (1 - p_{HL}) p_{1L} = p_{1L}.$$

It is easy to see that $\delta_H + \delta_L = 1$ for any routing policy such that $0 \leq p_{1H} < p_{1L} \leq 1$. The performance

of these routing policies is equivalent to the FCFS routing policy. Indeed, it is easy to see that whenever

$\delta_H + \delta_L = 1 \iff 1 - \delta_H = \delta_L$, the expected waiting cost in equation 1 is equal to the expected waiting cost

of a FCFS policy. Notice that if we relax our assumption that down-reporting does not carry a psychological cost, the overall performance of these routing policies will perform even worse than a FCFS routing policy (i.e., $\delta_H + \delta_L > 1$).

**Case 2:** $p_{1H} = p_{1L}$.

When an arbitrary customer type $x = H$ claims $y = H$ she has expected cost $c_H \mathbb{E}[W_1^\pi] p_{1H} + c_H \mathbb{E}[W_2^\pi](1 - p_{1H})$. Conversely, when she claims $y = L$ she has expected cost $c_H \mathbb{E}[W_1^\pi] p_{1L} + c_H \mathbb{E}[W_2^\pi](1 - p_{1L})$. Notice that whenever $p_{1H} = p_{1L} \iff 1 - p_{1H} = 1 - p_{1L}$, customers are indifferent in their claims and therefore they randomize with probability $p_{LH}$.

When an arbitrary customer type $x = L$ claims $y = H$ she has expected cost $c_L \mathbb{E}[W_1^\pi] p_{1H} + c_L \mathbb{E}[W_2^\pi](1 - p_{1H}) + \theta c_L (\mathbb{E}[W_2^\pi] - \mathbb{E}[W_1^\pi])$. Conversely, when she claims $y = L$ she has expected cost $c_L \mathbb{E}[W_1^\pi] p_{1L} + c_L \mathbb{E}[W_2^\pi](1 - p_{1L})$. For this case, it is easy to see that since $\mathbb{E}[W_2^\pi] \geq \mathbb{E}[W_1^\pi]$, independently on other customers' claims, the arbitrary customer with type $x = L$ will always claim to be a low type.

The manager correctly anticipates that high type customers will randomize with probability $p_{LH}$ and that low type customers will claim to be high type with probability $p_{HL} = 0$. Based on this, and given the fact that $p_{1H} = p_{1L} \iff 1 - p_{1H} = 1 - p_{1L}$, the manager anticipates the following prioritization error probabilities:

$$\delta_H(p_{1H}, p_{1L}) = p_{LH}(1 - p_{1L}) + (1 - p_{LH})(1 - p_{1H}) = 1 - p_{1L},$$

$$\delta_L(p_{1H}, p_{1L}) = p_{HL} p_{1H} + (1 - p_{HL}) p_{1L} = p_{1L}.$$

It is easy to see that $\delta_H + \delta_L = 1$ for any routing policy such that $p_{1H} = p_{1L}$. The performance of these routing policies is equivalent to the FCFS routing policy. Indeed, it is easy to see that whenever $\delta_H + \delta_L = 1 \iff 1 - \delta_H = \delta_L$, the expected waiting cost in equation 1 is equal to the expected waiting cost of a FCFS policy. Notice that this includes the special cases $p_{1H} = p_{1L} = 1$ and $p_{1H} = p_{1L} = 0$, which are indeed the FCFS routing policy.

**Case 3:** $p_{1H} > p_{1L}$.

When an arbitrary customer type $x = H$ claims $y = H$ she has expected cost $c_H \mathbb{E}[W_1^\pi] p_{1H} + c_H \mathbb{E}[W_2^\pi](1 - p_{1H})$. Conversely, when she claims $y = L$ she has expected cost $c_H \mathbb{E}[W_1^\pi] p_{1L} + c_H \mathbb{E}[W_2^\pi](1 - p_{1L})$. It is easy to see that since $\mathbb{E}[W_2^\pi] > \mathbb{E}[W_1^\pi]$, independently on other customers' claims, the arbitrary customer with type $x = H$ will always claim to be a high type.

When an arbitrary customer type $x = L$ claims $y = H$ she has expected cost $c_L \mathbb{E}[W_1^\pi] p_{1H} + c_L \mathbb{E}[W_2^\pi](1 - p_{1H}) + c_L \theta(\mathbb{E}[W_2^\pi] - \mathbb{E}[W_1^\pi])$. Conversely, when she claims $y = L$ she has expected cost $c_L \mathbb{E}[W_1^\pi] p_{1L} + c_L \mathbb{E}[W_2^\pi](1 - p_{1L})$. For this case, it is easy to see that since $\mathbb{E}[W_2^\pi] \geq \mathbb{E}[W_1^\pi]$, independently on other customers' claims, the arbitrary customer with type $x = L$ will claim to be a high type whenever:

$$c_L \mathbb{E}[W_1^\pi] p_{1H} + c_L \mathbb{E}[W_2^\pi](1 - p_{1H}) + c_L \theta(\mathbb{E}[W_2^\pi] - \mathbb{E}[W_1^\pi]) \leq c_L \mathbb{E}[W_1^\pi] p_{1L} + c_L \mathbb{E}[W_2^\pi](1 - p_{1L}) \iff \theta \leq p_{1H} - p_{1L}$$

and will claim to be low type whenever $\theta > p_{1H} - p_{1L}$. The manager correctly anticipates that all high type customers will claim to be high type (i.e., $p_{LH} = 0$) and that low type customers will claim to be high type with probability $p_{HL} = \mathbb{P}(\theta \leq p_{1H} - p_{1L}) = \Phi(p_{1H} - p_{1L})$. Based on this, the manager anticipates the following prioritization error probabilities:

$$\delta_H(p_{1H}, p_{1L}) = p_{LH}(1 - p_{1L}) + (1 - p_{LH})(1 - p_{1H}) = 1 - p_{1H},$$

$$\delta_L(p_{1H}, p_{1L}) = p_{HL} p_{1H} + (1 - p_{HL}) p_{1L} = p_{1L} + (p_{1H} - p_{1L})\Phi(p_{1H} - p_{1L}).$$

In this proof, for any function $q(\cdot)$, we use $q_z'(\cdot)$ to denote the first partial derivative of $q(\cdot)$ with respect to $z$. The manager defines the routing probabilities $p_{1H}, p_{1L} \in [0, 1]$ in order to minimize the expected waiting cost $C_s$ in the system which can be written as:

$$C_s = c \cdot f(p_{1H}, p_{1L}),$$

where

$$c = \left(\frac{\rho^2}{1 - \rho}\right)\left(\frac{p_H n_H + p_L n_L}{2(p_H m_H + p_L m_L)}\right) > 0,$$

$$f(p_{1H}, p_{1L}) = \frac{p_H c_H(1 - \rho(1 - \delta_H(p_{1H}, p_{1L}))) + p_L c_L(1 - \rho \delta_L(p_{1H}, p_{1L}))}{p_H m_H(1 - \rho(1 - \delta_H(p_{1H}, p_{1L}))) + p_L m_L(1 - \rho \delta_L(p_{1H}, p_{1L}))},$$

$$f_{p_{1H}}' = \frac{d(p_{1H}, p_{1L})\gamma(p_{1H}, p_{1L})}{g(p_{1H}, p_{1L})},$$

$$f_{p_{1L}}' = \frac{d(p_{1H}, p_{1L})(1 - \rho p_{1H})\kappa(p_{1H}, p_{1L})}{g(p_{1H}, p_{1L})},$$

$$d(p_{1H}, p_{1L}) = p_H p_L \rho(c_H m_L - c_L m_H)(1 - \Phi(p_{1H} - p_{1L})) > 0,$$

$$g(p_{1H}, p_{1L}) = (p_H m_H(1 - \rho(1 - \delta_H(p_{1H}, p_{1L}))) + p_L m_L(1 - \rho \delta_L(p_{1H}, p_{1L})))^2 > 0,$$

$$\gamma(p_{1H}, p_{1L}) = -(1 - \rho p_{1L}) + (1 - \rho p_{1H})h(p_{1H} - p_{1L}),$$

$$\gamma_{p_{1H}}' = -\rho h(p_{1H} - p_{1L}) + (1 - \rho p_{1H})h'(p_{1H} - p_{1L}),$$

$$\kappa(p_{1H}, p_{1L}) = 1 - h(p_{1H} - p_{1L}),$$

$$\kappa'_{p_{1L}} = h'(p_{1H} - p_{1L}),$$

$$h(p_{1H} - p_{1L}) = \frac{(p_{1H} - p_{1L})\phi(p_{1H} - p_{1L})}{1 - \Phi(p_{1H} - p_{1L})}.$$

We want to solve the following constraint minimization problem:

$$Min \ \ c \cdot f(p_{1H}, p_{1L})$$

$$s.t$$

$$p_{1L} \geq 0 \iff -p_{1L} \leq 0$$

$$p_{1L} \leq p_{1H} \iff p_{1L} - p_{1H} \leq 0$$

$$p_{1H} \leq 1 \iff p_{1H} - 1 \leq 0$$

We construct the lagrangian:

$$\mathcal{L} = cf(p_{1H}, p_{1L}) - \mu_1(p_{1L}) + \mu_2(p_{1L} - p_{1H}) + \mu_3(p_{1H} - 1).$$

We derive the KKT conditions:

*Stationarity*

$$\mathcal{L}'_{p_{1H}} = cf'_{p^*_{1H}} - \mu_2 + \mu_3 = 0,$$

$$\mathcal{L}'_{p_{1L}} = cf'_{p^*_{1L}} - \mu_1 + \mu_2 = 0.$$

*Complementary Slackness*

$$\mu_1(-p^*_{1L}) = 0,$$

$$\mu_2(p^*_{1L} - p^*_{1H}) = 0,$$

$$\mu_3(p^*_{1H} - 1) = 0.$$

*Dual Feasibility*

$$\mu_1, \mu_2, \mu_3 \geq 0.$$

*Primal Feasibility*

$$p^*_{1L} \geq 0,$$

38

**Estrada, Ibrahim, and Zhan:** *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

$$p_{1L}^* \leq p_{1H}^*,$$

$$p_{1H}^* \leq 1.$$

**Candidate 1:** $(p_{1H}^* = 1, p_{1L}^* = 0)$.

*Complementary Slackness*

$\mu_1 \neq 0 \Rightarrow p_{1L}^* = 0,$

$\mu_2 = 0,$

$\mu_3 \neq 0 \Rightarrow p_{1H}^* = 1.$

*Stationarity*

$\mu_3 = -cf'_{p_{1H}^*} = -\frac{cd(p_{1H}^*,p_{1L}^*)\gamma(p_{1H}^*,p_{1L}^*)}{g(p_{1H}^*,p_{1L}^*)} = -\frac{cd(1,0)\gamma(1,0)}{g(1,0)},$

$\mu_1 = cf'_{p_{1L}^*} = \frac{cd(p_{1H}^*,p_{1L}^*)(1-\rho p_{1H}^*)\kappa(p_{1H}^*,p_{1L}^*)}{g(p_{1H}^*,p_{1L}^*)} = \frac{cd(1,0)(1-\rho)\kappa(1,0)}{g(1,0)}.$

*Dual Feasibility*

$\mu_3 \geq 0 \iff \frac{cd(1,0)\gamma(1,0)}{g(1,0)} \leq 0 \iff \gamma(1,0) \leq 0 \iff h(1) \leq 1/(1-\rho),$

$\mu_1 \geq 0 \iff \frac{cd(1,0)(1-\rho)\kappa(1,0)}{g(1,0)} \geq 0 \iff \kappa(1,0) \geq 0 \iff h(1) \leq 1.$

*Primal Feasibility*

$p_{1L}^* \geq 0,$

$p_{1L}^* \leq p_{1H}^*,$

$p_{1H}^* \leq 1.$

**Candidate 2:** $(p_{1H}^* = 1, p_{1L}^* \in (0,1))$.

*Complementary Slackness*

$\mu_1 = 0,$

$\mu_2 = 0,$

$\mu_3 \neq 0 \Rightarrow p_{1H}^* = 1.$

*Stationarity*

$cf'_{p_{1L}^*} = 0 \iff \frac{cd(p_{1H}^*,p_{1L}^*)(1-\rho p_{1H}^*)\kappa(p_{1H}^*,p_{1L}^*)}{g(p_{1H}^*,p_{1L}^*)} = \frac{cd(1,p_{1L}^*)(1-\rho)\kappa(1,p_{1L}^*)}{g(1,p_{1L}^*)} = 0 \iff \kappa(1,p_{1L}^*) = 0 \iff h(1-p_{1L}^*) = 1,$

$\mu_3 = -cf'_{p_{1H}^*} = -\frac{cd(p_{1H}^*,p_{1L}^*)\gamma(p_{1H}^*,p_{1L}^*)}{g(p_{1H}^*,p_{1L}^*)} = -\frac{cd(1,p_{1L}^*)\gamma(1,p_{1L}^*)}{g(1,p_{1L}^*)}.$

*Dual Feasibility*

$\mu_3 \geq 0 \iff \frac{cd(1,p_{1L}^*)\gamma(1,p_{1L}^*)}{g(1,p_{1L}^*)} \leq 0 \iff \gamma(1,p_{1L}^*) \leq 0 \iff -\rho(1-p_{1L}^*) \leq 0.$

*Primal Feasibility*

Assume that $\Phi$ is IGFR (i.e., $h'(x) > 0$), then it follows that $h(x_1) \geq h(x_0)$ if and only if $x_1 \geq x_0$. We have that $h(0) = 0$. Since we assume that $h'(x) > 0$, then there is a unique $p_{1L}^* \in (0,1)$, whenever $h(1) > 1$ such that $h(1 - p_{1L}^*) = 1$.

$$p_{1L}^* \geq 0 \iff h(1) \geq 1,$$

$$p_{1L}^* \leq p_{1H}^*,$$

$$p_{1H}^* \leq 1.$$

**Candidate 3:** $(p_{1H}^* \in (p_{1L}^*, 1), p_{1L}^* = 0)$.

*Complementary Slackness*

$$\mu_1 \neq 0 \Rightarrow p_{1L}^* = 0,$$

$$\mu_2 = 0,$$

$$\mu_3 = 0,$$

*Stationarity*

$$cf'_{p_{1H}^*} = 0 \iff \frac{cd(p_{1H}^*, p_{1L}^*)\gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = 0 \iff \frac{cd(p_{1H}^*, 0)\gamma(p_{1H}^*, 0)}{g(p_{1H}^*, 0)} = 0 \iff \gamma(p_{1H}^*, 0) = 0 \iff h(p_{1H}^*) = \frac{1}{1 - \rho p_{1H}^*},$$

$$\mu_1 = cf'_{p_{1L}^*} = \frac{cd(p_{1H}^*, p_{1L}^*)(1 - \rho p_{1H}^*)\kappa(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = \frac{cd(p_{1H}^*, 0)(1 - \rho p_{1H}^*)\kappa(p_{1H}^*, 0)}{g(p_{1H}^*, 0)}.$$

*Dual Feasibility*

$$\mu_1 \geq 0 \iff \frac{cd(p_{1H}^*, 0)(1 - \rho p_{1H}^*)\kappa(p_{1H}^*, 0)}{g(p_{1H}^*, 0)} \geq 0 \iff \kappa(p_{1H}^*, 0) \geq 0 \iff h(p_{1H}^*) \leq 1 \text{ which is never the case since}$$

$h(p_{1H}^*) = \frac{1}{1 - \rho p_{1H}^*}$, $\rho \in (0,1)$, and $p_{1H}^* \in (0,1)$.

**Candidate 4:** $(p_{1H}^* \in (p_{1L}^*, 1), p_{1L}^* \in (0, p_{1H}^*))$.

*Complementary Slackness*

$$\mu_1 = 0,$$

$$\mu_2 = 0,$$

$$\mu_3 = 0.$$

*Stationarity*

$$cf'_{p_{1L}^*} = 0 \iff \frac{cd(p_{1H}^*, p_{1L}^*)(1 - \rho p_{1H}^*)\kappa(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = 0 \iff \kappa(p_{1H}^*, p_{1L}^*) = 0 \iff h(p_{1H}^* - p_{1L}^*) = 1,$$

$$cf'_{p_{1H}^*} = 0 \iff \frac{cd(p_{1H}^*, p_{1L}^*)\gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = 0 \iff \frac{cd(p_{1H}^*, p_{1L}^*)\gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = 0 \iff \gamma(p_{1H}^*, p_{1L}^*) = 0 \iff h(p_{1H}^* -$$

$p_{1L}^*) = \frac{1 - \rho p_{1L}^*}{1 - \rho p_{1H}^*} \neq 1.$

**Candidate 5:** $(p_{1H}^* = p_{1L}^* \in (0,1))$.

*Complementary Slackness*

$\mu_1 = 0,$

$\mu_2 \neq 0 \Rightarrow p_{1H}^* = p_{1L}^*,$

$\mu_3 = 0.$

*Stationarity*

$\mu_2 = c f_{p_{1H}^*}' = \frac{cd(p_{1H}^*, p_{1L}^*) \gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = -\frac{cd(p_{1H}^*, p_{1L}^*)(1 - \rho p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)}.$

*Dual Feasibility*

$\mu_2 \geq 0 \iff -\frac{cd(p_{1H}^*, p_{1L}^*)(1 - \rho p_{1H}^*)}{g(p_{1H}^*, p_{1L}^*)} \geq 0 \iff -(1 - \rho p_{1L}^*) \geq 0$ which is never the case.

**Candidate 6:** $(p_{1H}^* = p_{1L}^* = 0).$

*Complementary Slackness*

$\mu_1 \neq 0 \Rightarrow p_{1L}^* = 0,$

$\mu_2 \neq 0 \Rightarrow p_{1H}^* = p_{1L}^*,$

$\mu_3 = 0.$

*Stationarity*

$\mu_2 = c f_{p_{1H}^*}' = \frac{cd(p_{1H}^*, p_{1L}^*) \gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = -\frac{cd(0,0)}{g(0,0)}.$

*Dual Feasibility*

$\mu_2 \geq 0 \iff -\frac{cd(0,0)}{g(0,0)} \geq 0$ which is never the case.

**Candidate 7:** $(p_{1H}^* = p_{1L}^* = 1).$

*Complementary Slackness*

$\mu_1 = 0,$

$\mu_2 \neq 0 \Rightarrow p_{1H}^* = p_{1L}^*,$

$\mu_3 \neq 0 \Rightarrow p_{1H}^* = 1.$

*Stationarity*

$\mu_2 = -c f_{p_{1L}^*}' = -\frac{cd(p_{1H}^*, p_{1L}^*) \gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = -\frac{cd(1,1)(1 - \rho)}{g(1,1)}.$

*Dual Feasibility*

$\mu_2 \geq 0 \iff -\frac{cd(1,1)(1 - \rho)}{g(1,1)} = -(1 - \rho) \geq 0$ which is never the case.

We can see that only candidates 1 and 2 comply with the KKT conditions. In particular, we notice that candidate 1 (i.e., $p_{1H}^* = 1, p_{1L}^* = 0$) holds whenever $h(1) \leq 1$ whereas candidate 2 (i.e., $p_{1H}^* = 1, p_{1L}^* \in (0,1)$)

such that $h(1 - p_{1L}^*) = 1)$ holds whenever $h(1) > 1$. Since these two candidates are mutually exclusive, we conclude that the KKT necessary conditions are also sufficient.

Finally, notice that in this case, we have that all high type customers claim their type and that a proportion equal to $\Phi(1 - p_{1L}^*) < 1$ of low type customers claim deceptively. When the manager sets the routing prioritization policy $p_{1H}^* = 1, p_{1L}^* \in [0,1)$ the prioritization error probabilities are $\delta_H = 0$ and $\delta_L = p_{1L}^* + (1 - p_{1L}^*)\Phi(1 - p_{1L}^*) < 1$. It is easy to see that $\delta_H + \delta_L < 1$, which implies that the routing policy performs better than a FCFS discipline. Indeed, recall that for a FCFS policy (i.e., $p_{1H} = 1, p_{1L} = 1 \Rightarrow \delta_H = 0, \delta_L = 1$) we have that $\delta_H + \delta_L = 1$. From Lemma 1, it is easy to see that since the routing policy yields $\delta_H = 0, \delta_L < 1$ the associated expected waiting cost is strictly lower than that of a FCFS policy. Finally, since the routing policies in cases 1 (i.e., $p_{1H} < p_{1L}$) and 2 (i.e., $p_{1H} = p_{1L}$) perform the same as a FCFS policy, we conclude that the routing policy in this case 3 performs better, and thus arises in equilibrium. ∎

## Appendix B:   Customers with Uncertain Type Information

So far, we have assumed that customers have perfect private information about their own types. In this section, we extend this notion and consider the case where customers have only partial, imperfect, information about their true types. Based on this private imperfect information, they make inferences about their true underlying types. In what follows, we restrict attention to the case where $\theta$ is uniformly distributed over $[0, \bar{\theta}]$. We assume that customers draw their private information, $S = s \in \{H, L\}$, from a joint distribution $\mathbb{P}(X = x, S = s)$ which is given as follows:

$$\mathbb{P}(X = x, S = s) = \begin{cases} p_H s_H & \text{if } (x,s) = (H,H), \\ p_H(1 - s_H) & \text{if } (x,s) = (H,L), \\ p_L(1 - s_L) & \text{if } (x,s) = (L,H), \\ p_L s_L & \text{if } (x,s) = (L,L), \end{cases}$$

where $s_H = \mathbb{P}(S = H | X = H)$ and $s_L = \mathbb{P}(S = L | X = L)$. Consider an arriving customer with private information $s = H$. Based on this private information, this customer computes her conditional probability of being of $H$ type as follows:

$$\mathbb{P}(X = H | S = H) = \frac{p_H s_H}{p_H s_H + p_L(1 - s_L)}.$$

42

**Estrada, Ibrahim, and Zhan:** *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

Similarly an arriving customer with private information $S = L$ computes her conditional probability of being of $L$ type as follows:

$$\mathbb{P}(X = L | S = L) = \frac{p_L s_L}{p_L s_L + p_H (1 - s_H)}.$$

We note that whenever $s_H = s_L = 1$, we have that $\mathbb{P}(X = H | S = H) = \mathbb{P}(X = L | S = L) = 1$. In this case, the private information of a customer is a perfect indicator of her type. On the other hand, whenever $s_H = s_L = 1/2$, we have that $\mathbb{P}(X = H | S = H) = p_H$ and $\mathbb{P}(X = L | S = L) = p_L$. In this case, the private information of a customer does not carry additional information about her type. Based on this, we consider that $s_H = \mathbb{P}(S = H | X = H) \geq 1/2$ and $s_L = \mathbb{P}(S = L | X = L) \geq 1/2$ in what follows.

## B.1. Customer Problem

In our model, an arbitrary customer does not know her true type. Instead, she has some private signal $s$, and must decide on her claim $y$ in order to minimize her expected cost:

$$\underset{y \in \{H, L\}}{Min} \; \mathbb{E}\left[ c_X \mathbb{E}[W^{\pi}_{g^{\pi}(y)}] + \theta c_X \left(\mathbb{E}[W^{\pi}_X] - \mathbb{E}[W^{\pi}_y]\right)^+ \cdot \mathbf{1}_{y \neq s} \right], \tag{6}$$

where $\mathbf{1}(\cdot)$ is the indicator function. In (6), we assume that customers incur a lying cost whenever they misreport their private information $s$ i.e., they make a claim $y \neq s$. In particular, since customers are uncertain about their true types $X$, in so doing we capture the fact that misreporting the private information that a customer has is necessary to trigger her psychological cost. For instance, consider a customer with private information $s = H$. This customer does not know with certainty her true type $X$. Even if she knows that there is some chance to be of a type $X = L$, she will not incur a psychological cost when reporting to be of $y = H$ type, since she is truthfully revealing her private information $s = H$.

## B.2. Optimal Scheduling Policy

PROPOSITION 2. *There exists a unique Nash equilibrium which arises in the game between the customers and the manager, based on problems (1) and (6). At equilibrium:*

(a) *A fraction $= s_H + (1 - s_H)(1 - p^*_{1L})(1 + u)/\bar{\theta}$ of customers with true type $x = H$ make a claim $y = H$. The remaining type $H$ customers claim $y = L$.*

(b) *A fraction $1 - s_L + s_L(1 - p^*_{1L})(1 + u)/\bar{\theta}$ of customers with true type $x = L$ claim $y = H$. The remaining type $L$ customers claim $y = L$.*

(c) *The optimal routing policy is to assign high priority to all customers that claim to be high type (i.e.,*

$p_{1H}^* = 1$) *and to randomly assign with probability* $p_{1L}^*$ *high priority to customers that claim to be low type.*

*We have that* $p_{1L}^* = 0$ *if, and only if,* $\bar{\theta} > 2(1+u)$ *and* $p_{1L}^* = \frac{2(1+u)-\bar{\theta}}{2(1+u)}$ *otherwise, where* $u = \frac{p_H(1-s_H)c_H}{p_L s_L c_L}$.

PROOF. There are 3 qualitatively different cases to analyse: (1) $p_{1H} < p_{1L}$, (2) $p_{1H} = p_{1L}$, and (3) $p_{1H} > p_{1L}$.

We recall that

$$p_{LH} = \mathbb{P}[Y = L | X = H] = \mathbb{P}[Y = L | S = H, X = H]\mathbb{P}[S = H | X = H] + \mathbb{P}[Y = L | S = L, X = H]\mathbb{P}[S = L | X = H]$$

and

$$p_{HL} = \mathbb{P}[Y = H | X = L] = \mathbb{P}[Y = H | S = H, X = L]\mathbb{P}[S = H | X = L] + \mathbb{P}[Y = H | S = L, X = L]\mathbb{P}[S = L | X = L].$$

**Case 1:** $p_{1H} < p_{1L}$**.**

When an arbitrary customer with private information $s = H$ claims $y = H$ she has expected cost $\mathbb{P}(X = H | S = H)(c_H(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L | S = H)(c_L(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi]))$. Conversely, when she claims $y = L$, she has expected cost $\mathbb{P}(X = H | S = H)(c_H(\mathbb{E}[W_1^\pi]p_{1L} + (1 - p_{1L})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L | S = H)(c_L(p_{1L}\mathbb{E}[W_1^\pi] + (1 - p_{1L})\mathbb{E}[W_2^\pi]))$. Since $\mathbb{E}[W_2^\pi] \geq \mathbb{E}[W_1^\pi]$, independently of other customers' claims, the arbitrary customer with private information $s = H$ will always claim to be a low type.

When an arbitrary customer with private information $s = L$ claims $y = H$ she has expected cost $\mathbb{P}(X = H | S = L)(c_H(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L | S = L)(c_L(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi]) + \theta c_L(\mathbb{E}[W_2^\pi] - \mathbb{E}[W_1^\pi]))$. Conversely, when she claims $y = L$ she has expected cost $\mathbb{P}(X = H | S = L)(c_H(\mathbb{E}[W_1^\pi]p_{1L} + (1 - p_{1L})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L | S = L)(c_L(p_{1L}\mathbb{E}[W_1^\pi] + (1 - p_{1L})\mathbb{E}[W_2^\pi]))$. Since $\mathbb{E}[W_2^\pi] \geq \mathbb{E}[W_1^\pi]$, independently on other customers' claims, the arbitrary customer with private information $s = L$ will always claim to be a low type.

The manager correctly anticipates that customers of type $H$ will claim to be low type with probability $p_{LH} = 1$ and customers with a low type will claim to be high type with probability $p_{HL} = 0$. Based on this, the manager anticipates the following prioritization error probabilities:

$$\delta_H(p_{1H}, p_{1L}) = p_{LH}(1 - p_{1L}) + (1 - p_{LH})(1 - p_{1H}) = 1 - p_{1L},$$

$$\delta_L(p_{1H}, p_{1L}) = p_{HL}p_{1H} + (1 - p_{HL})p_{1L} = p_{1L}.$$

It is easy to see that $\delta_H + \delta_L = 1$ for any routing policy such that $0 \leq p_{1H} < p_{1L} \leq 1$. The performance of this routing policy is equivalent to the FCFS routing policy. Indeed, it is easy to see that whenever $\delta_H + \delta_L = 1 \iff 1 - \delta_H = \delta_L$, the expected waiting cost in equation 1 is equal to the expected waiting cost of a FCFS policy.

**Case 2:** $p_{1H} = p_{1L}$.

When an arbitrary customer with private information $s = H$ claims $y = H$ she has expected cost $\mathbb{P}(X = H|S = H)(c_H(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = H)(c_L(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi]))$. Conversely, when she claims $y = L$, she has expected cost $\mathbb{P}(X = H|S = H)(c_H(\mathbb{E}[W_1^\pi]p_{1L} + (1 - p_{1L})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = H)(c_L(p_{1L}\mathbb{E}[W_1^\pi] + (1 - p_{1L})\mathbb{E}[W_2^\pi]))$. Notice that whenever $p_{1H} = p_{1L} \iff 1 - p_{1H} = 1 - p_{1L}$, customers are indifferent in their claims and therefore they randomize with some probability $p$.

When an arbitrary customer with private information $s = L$ claims $y = H$ she has expected cost $\mathbb{P}(X = H|S = L)(c_H(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = L)(c_L(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi]) + \theta c_L(\mathbb{E}[W_2^\pi] - \mathbb{E}[W_1^\pi]))$. Conversely, when she claims $y = L$ she has expected cost $\mathbb{P}(X = H|S = L)(c_H(\mathbb{E}[W_1^\pi]p_{1L} + (1 - p_{1L})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = L)(c_L(p_{1L}\mathbb{E}[W_1^\pi] + (1 - p_{1L})\mathbb{E}[W_2^\pi]))$. Since $\mathbb{E}[W_2^\pi] \geq \mathbb{E}[W_1^\pi]$, independently on other customers' claims, the arbitrary customer with private information $s = L$ will always claim to be a low type.

The manager correctly anticipates that high type customers will claim to be low type with probability $p_{LH} = 1 - s_H(1 - p)$ and that low type customers will claim to be high type with probability $p_{HL} = (1 - s_L)p$. Based on this, and given the fact that $p_{1H} = p_{1L} \iff 1 - p_{1H} = 1 - p_{1L}$, the manager anticipates the following prioritization error probabilities:

$$\delta_H(p_{1H}, p_{1L}) = p_{LH}(1 - p_{1L}) + (1 - p_{LH})(1 - p_{1H}) = 1 - p_{1L},$$

$$\delta_L(p_{1H}, p_{1L}) = p_{HL}p_{1H} + (1 - p_{HL})p_{1L} = p_{1L}.$$

It is easy to see that $\delta_H + \delta_L = 1$ for any routing policy such that $p_{1H} = p_{1L}$. The performance of this routing policy is equivalent to the one of the FCFS routing policy. Indeed, it is easy to see that whenever $\delta_H + \delta_L = 1 \iff 1 - \delta_H = \delta_L$, the expected waiting cost in equation 1 is equal to the expected waiting cost of a FCFS policy. Notice that this includes the special cases $p_{1H} = p_{1L} = 1$ and $p_{1H} = p_{1L} = 0$, which indeed correspond to the FCFS routing policy.

**Case 3:** $p_{1H} > p_{1L}$.

When an arbitrary customer with private information $s = H$ claims $y = H$ she has expected cost $\mathbb{P}(X = H|S = H)(c_H(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = H)(c_L(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi]))$. Conversely, when she claims $y = L$ she has expected cost $\mathbb{P}(X = H|S = H)(c_H(\mathbb{E}[W_1^\pi]p_{1L} + (1 - p_{1L})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = H)(c_L(p_{1L}\mathbb{E}[W_1^\pi] + (1 - p_{1L})\mathbb{E}[W_2^\pi]))$. Since $\mathbb{E}[W_2^\pi] \geq \mathbb{E}[W_1^\pi]$, independently on other customers' claims, the arbitrary customer with private information $s = H$ will always claim to be a high type.

When an arbitrary customer with private information $s = L$ claims $y = H$ she has expected cost $\mathbb{P}(X = H|S = L)(c_H(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = L)(c_L(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi]) + \theta c_L(\mathbb{E}[W_2^\pi] - \mathbb{E}[W_1^\pi]))$. Conversely, when she claims $y = L$ she has expected cost $\mathbb{P}(X = H|S = L)(c_H(\mathbb{E}[W_1^\pi]p_{1L} + (1 - p_{1L})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = L)(c_L(p_{1L}\mathbb{E}[W_1^\pi] + (1 - p_{1L})\mathbb{E}[W_2^\pi]))$. Since $\mathbb{E}[W_2^\pi] \geq \mathbb{E}[W_1^\pi]$, independently on other customers' claims, the arbitrary customer with private information $s = L$ will claim to be a high type whenever: $\mathbb{P}(X = H|S = L)(c_H(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = L)(c_L(p_{1H}\mathbb{E}[W_1^\pi] + (1 - p_{1H})\mathbb{E}[W_2^\pi]) + \theta c_L(\mathbb{E}[W_2^\pi] - \mathbb{E}[W_1^\pi])) \leq \mathbb{P}(X = H|S = L)(c_H(\mathbb{E}[W_1^\pi]p_{1L} + (1 - p_{1L})\mathbb{E}[W_2^\pi])) + \mathbb{P}(X = L|S = L)(c_L(p_{1L}\mathbb{E}[W_1^\pi] + (1 - p_{1L})\mathbb{E}[W_2^\pi]))$

$$\Longleftrightarrow$$

$$\theta \leq p_{1H} - p_{1L} + \frac{\mathbb{P}(X = H|S = L)c_H}{\mathbb{P}(X = L|S = L)c_L}(p_{1H} - p_{1L})$$

$$\Longleftrightarrow$$

$$\theta \leq (p_{1H} - p_{1L})(1 + \frac{p_H(1 - s_H)c_H}{p_L s_L c_L}) = (p_{1H} - p_{1L})(1 + u),$$

where we recall that $u = \frac{p_H(1-s_H)c_H}{p_L s_L c_L}$. That is, a customer with a low signal will claim to be high type with probability $\mathbb{P}(\theta \leq (p_{1H} - p_{1L})(1 + u)) = (p_{1H} - p_{1L})(1 + u)/\bar{\theta}$, where the constant $u = \frac{p_H(1-s_H)c_H}{p_L s_L c_L}$ captures the effect of uncertainty on customers' aggregate claiming behaviour. Notice that as the uncertainty of types increases (i.e., as $s_H$ and $s_L$ decrease), the proportion of customers that misreport their private information increases. Naturally, if the upper bound $\bar{\theta}$ is not sufficiently high, all customers with information low signal will cheat. We are interested in analysing the role of uncertainty in cheating behaviour. We assume that $\bar{\theta}$ is sufficiently large such that there is always some truthful reporting (potentially very small), i.e., $\bar{\theta} > 1 + u$.

The manager correctly anticipates that customers with true type $H$ will claim to be high type with probability $1 - p_{LH} = \mathbb{P}(Y = H|X = H) = s_H + (1 - s_H)(p_{1H} - p_{1L})(1 + u)/\bar{\theta}$ and L type customers will claim

to be high type with probability $p_{HL} = \mathbb{P}(Y = H | X = L) = (1 - s_L) + s_L(p_{1H} - p_{1L})(1 + u)/\bar{\theta}$. Based on this, the manager anticipates the following prioritization error probabilities:

$$\delta_H(p_{1H}, p_{1L}) = p_{LH}(1 - p_{1L}) + (1 - p_{LH})(1 - p_{1H}),$$

$$\delta_L(p_{1H}, p_{1L}) = p_{HL}p_{1H} + (1 - p_{HL})p_{1L}.$$

In this proof, for any function $q(\cdot)$, we use $q'_z(\cdot)$ to denote the first partial derivative of $q(\cdot)$ with respect to $z$. The manager defines the routing probabilities $p_{1H}, p_{1L} \in [0, 1]$ in order to minimize the expected waiting cost $C_s$ in the system which can be written as:

$$C_s = c_1 f(p_{1H}, p_{1L}),$$

where,

$$c_1 = \left(\frac{\rho^2}{1 - \rho}\right)\left(\frac{p_H n_H + p_L n_L}{2(p_H m_H + p_L m_L)}\right) > 0,$$

$$f(p_{1H}, p_{1L}) = \frac{p_H c_H(1 - \rho(1 - \delta_H(p_{1H}, p_{1L}))) + p_L c_L(1 - \rho\delta_L(p_{1H}, p_{1L}))}{p_H m_H(1 - \rho(1 - \delta_H(p_{1H}, p_{1L}))) + p_L m_L(1 - \rho\delta_L(p_{1H}, p_{1L}))},$$

$$f'_{p_{1H}} = \frac{c_2 \gamma(p_{1H}, p_{1L})}{g(p_{1H}, p_{1L})},$$

$$f'_{p_{1L}} = \frac{c_2(1 - \rho p_{1H})\kappa(p_{1H}, p_{1L})}{g(p_{1H}, p_{1L})},$$

$$c_2 = p_H p_L \rho(c_H m_L - c_L m_H)(s_H + s_L - 1) > 0,$$

$$g(p_{1H}, p_{1L}) = (p_H m_H(1 - \rho(1 - \delta_H)) + p_L m_L(1 - \rho\delta_L))^2 > 0,$$

$$\gamma(p_{1H}, p_{1L}) = -(1 - \rho p_{1L}) + (p_{1H} - p_{1L})((1 + u)/\bar{\theta})(2 - \rho(p_{1H} + p_{1L})),$$

$$\gamma'_{p_{1H}} = 2((1 + u)/\bar{\theta})(1 - \rho p_{1H}) > 0,$$

$$\gamma'_{p_{1L}} = \rho - 2((1 + u)/\bar{\theta})(1 - \rho p_{1L}),$$

$$\kappa(p_{1H}, p_{1L}) = 1 - 2((1 + u)/\bar{\theta})(p_{1H} - p_{1L}),$$

$$\kappa'_{p_{1L}} = 2(1 + u)/\bar{\theta},$$

$$\kappa'_{p_{1H}} = -2(1 + u)/\bar{\theta}.$$

We want to solve the following constraint minimization problem:

$$Min \ \ c_1 f(p_{1H}, p_{1L})$$

$$s.t$$

$$p_{1L} \geq 0 \iff -p_{1L} \leq 0$$

$$p_{1L} \leq p_{1H} \iff p_{1L} - p_{1H} \leq 0$$

$$p_{1H} \leq 1 \iff p_{1H} - 1 \leq 0$$

We construct the lagrangian:

$$\mathcal{L} = c_1 f(p_{1H}, p_{1L}) - \mu_1(p_{1L}) + \mu_2(p_{1L} - p_{1H}) + \mu_3(p_{1H} - 1).$$

We derive the KKT conditions:

*Stationarity*

$$\mathcal{L}'_{p_{1H}} = c_1 f'_{p^*_{1H}} - \mu_2 + \mu_3 = 0,$$

$$\mathcal{L}'_{p_{1L}} = c_1 f'_{p^*_{1L}} - \mu_1 + \mu_2 = 0.$$

*Complementary Slackness*

$$\mu_1(-p^*_{1L}) = 0,$$

$$\mu_2(p^*_{1L} - p^*_{1H}) = 0,$$

$$\mu_3(p^*_{1H} - 1) = 0.$$

*Dual Feasibility*

$$\mu_1, \mu_2, \mu_3 \geq 0.$$

*Primal Feasibility*

$$p^*_{1L} \geq 0,$$

$$p^*_{1L} \leq p^*_{1H},$$

$$p^*_{1H} \leq 1.$$

**Candidate 1:** $(p^*_{1H} = 1, p^*_{1L} = 0)$.

*Complementary Slackness*

$\mu_1 \neq 0 \Rightarrow p^*_{1L} = 0,$

$\mu_2 = 0,$

$\mu_3 \neq 0 \Rightarrow p_{1H}^* = 1.$

*Stationarity*

$\mu_3 = -c_1 f'_{p_{1H}^*} = -\frac{c_1 c_2 \gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = -\frac{c_1 c_2 \gamma(1,0)}{g(1,0)},$

$\mu_1 = c_1 f'_{p_{1L}^*} = \frac{c_1 c_2 (1 - \rho p_{1H}^*) \kappa(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = \frac{c_1 c_2 (1 - \rho) \kappa(1,0)}{g(1,0)}.$

*Dual Feasibility*

$\mu_3 \geq 0 \iff \frac{c_1 c_2 \gamma(1,0)}{g(1,0)} \leq 0 \iff \gamma(1,0) \leq 0 \iff \bar{\theta} \geq (2 - \rho)(1 + u),$

$\mu_1 \geq 0 \iff \frac{c_1 c_2 (1-\rho) \kappa(1,0)}{g(1,0)} \geq 0 \iff \kappa(1,0) \geq 0 \iff \bar{\theta} \geq 2(1 + u).$

*Primal Feasibility*

$p_{1L}^* \geq 0,$

$p_{1L}^* \leq p_{1H}^*,$

$p_{1H}^* \leq 1.$

**Candidate 2:** $(p_{1H}^* = 1, p_{1L}^* \in (0, 1)).$

*Complementary Slackness*

$\mu_1 = 0,$

$\mu_2 = 0,$

$\mu_3 \neq 0 \Rightarrow p_{1H}^* = 1.$

*Stationarity*

$c_1 f'_{p_{1L}^*} = 0 \iff \frac{c_1 c_2 (1 - \rho p_{1H}^*) \kappa(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = \frac{c_1 c_2 (1 - \rho) \kappa(1, p_{1L}^*)}{g(1, p_{1L}^*)} = 0 \iff \kappa(1, p_{1L}^*) = 0 \iff p_{1L}^* = \frac{2(1 + u) - \bar{\theta}}{2(1 + u)},$

$\mu_3 = -c_1 f'_{p_{1H}^*} = -\frac{c_1 c_2 \gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = -\frac{c_1 c_2 \gamma(1, p_{1L}^*)}{g(1, p_{1L}^*)}.$

*Dual Feasibility*

$\mu_3 \geq 0 \iff \frac{c_1 c_2 \gamma(1, p_{1L}^*)}{g(1, p_{1L}^*)} \leq 0 \iff \gamma(1, p_{1L}^*) \leq 0 \iff -\rho \bar{\theta}/4(1 + u) \leq 0.$

*Primal Feasibility*

$p_{1L}^* \geq 0 \iff \bar{\theta} \leq 2(1 + u),$

$p_{1L}^* \leq p_{1H}^*,$

$p_{1H}^* \leq 1.$

**Candidate 3:** $(p_{1H}^* \in (p_{1L}^*, 1), p_{1L}^* = 0).$

*Complementary Slackness*

$\mu_1 \neq 0 \Rightarrow p_{1L}^* = 0,$

$$\mu_2 = 0,$$

$$\mu_3 = 0.$$

*Stationarity*

$$c_1 f'_{p^*_{1H}} = 0 \iff \frac{c_1 c_2 \gamma(p^*_{1H}, p^*_{1L})}{g(p^*_{1H}, p^*_{1L})} = 0 \iff \frac{c_1 c_2 \gamma(p^*_{1H}, 0)}{g(p^*_{1H}, 0)} = 0 \iff \gamma(p^*_{1H}, 0) = 0 \iff p^*_{1H} = \frac{((1+u)/\bar{\theta}) - \sqrt{((1+u)/\bar{\theta})^2 - ((1+u)/\bar{\theta})\rho}}{((1+u)/\bar{\theta})\rho},$$

$$\mu_1 = c_1 f'_{p^*_{1L}} = \frac{c_1 c_2 (1 - \rho p^*_{1H}) \kappa(p^*_{1H}, p^*_{1L})}{g(p^*_{1H}, p^*_{1L})} = \frac{c_1 c_2 (1 - \rho p^*_{1H}) \kappa(p^*_{1H}, 0)}{g(p^*_{1H}, 0)}.$$

*Dual Feasibility*

$$\mu_1 \geq 0 \iff \frac{c_1 c_2 (1 - \rho p^*_{1H}) \kappa(p^*_{1H}, 0)}{g(p^*_{1H}, 0)} \geq 0 \iff \kappa(p^*_{1H}, 0) \geq 0 \iff \rho^2/4 \leq 0 \text{ which is never the case.}$$

*Primal Feasibility*

$$p^*_{1H} < 1 \iff \rho < (2((1+u)/\bar{\theta}) - 1)/((1+u)/\bar{\theta}),$$

$$p^*_{1H} \in \mathbb{R} \iff \rho < ((1+u)/\bar{\theta}).$$

**Candidate 4:** $(p^*_{1H} \in (p^*_{1L}, 1), p^*_{1L} \in (0, p^*_{1H}))$.

*Complementary Slackness*

$$\mu_1 = 0,$$

$$\mu_2 = 0,$$

$$\mu_3 = 0.$$

*Stationarity*

$$c_1 f'_{p^*_{1L}} = 0 \iff \frac{c_1 c_2 (1 - \rho p^*_{1H}) \kappa(p^*_{1H}, p^*_{1L})}{g(p^*_{1H}, p^*_{1L})} = 0 \iff \kappa(p^*_{1H}, p^*_{1L}) = 0 \iff p^*_{1L} = \frac{2((1+u)/\bar{\theta}) p^*_{1H} - 1}{2((1+u)/\bar{\theta})},$$

$$c_1 f'_{p^*_{1H}} = 0 \iff \frac{c_1 c_2 \gamma(p^*_{1H}, p^*_{1L})}{g(p^*_{1H}, p^*_{1L})} = 0 \iff \frac{c_1 c_2 \gamma(p^*_{1H}, p^*_{1L})}{g(p^*_{1H}, p^*_{1L})} = 0 \iff \gamma(p^*_{1H}, p^*_{1L}) = -\rho\bar{\theta}/4(1+u) \neq 0.$$

**Candidate 5:** $(p^*_{1H} = p^*_{1L} \in (0, 1))$.

*Complementary Slackness*

$$\mu_1 = 0,$$

$$\mu_2 \neq 0 \Rightarrow p^*_{1H} = p^*_{1L},$$

$$\mu_3 = 0.$$

*Stationarity*

$$\mu_2 = c_1 f'_{p^*_{1H}} = \frac{c_1 c_2 \gamma(p^*_{1H}, p^*_{1L})}{g(p^*_{1H}, p^*_{1L})} = -\frac{c_1 c_2 (1 - \rho p^*_{1H})}{g(p^*_{1H}, p^*_{1L})}.$$

*Dual Feasibility*

$$\mu_2 \geq 0 \iff -\frac{c_1 c_2 (1 - \rho p^*_{1H})}{g(p^*_{1H}, p^*_{1L})} = -(1 - \rho p^*_{1H}) \geq 0 \text{ which is never the case.}$$

**Candidate 6:** $(p_{1H}^* = p_{1L}^* = 0)$.

*Complementary Slackness*

$\mu_1 \neq 0 \Rightarrow p_{1L}^* = 0$,

$\mu_2 \neq 0 \Rightarrow p_{1H}^* = p_{1L}^*$,

$\mu_3 = 0$.

*Stationarity*

$\mu_2 = c_1 f'_{p_{1H}^*} = \frac{c_1 c_2 \gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = -\frac{c_1 c_2}{g(0,0)}$.

*Dual Feasibility*

$\mu_2 \geq 0 \iff -\frac{c_1 c_2}{g(0,0)} \geq 0$ which is never the case.

**Candidate 7:** $(p_{1H}^* = p_{1L}^* = 1)$.

*Complementary Slackness*

$\mu_1 = 0$,

$\mu_2 \neq 0 \Rightarrow p_{1H}^* = p_{1L}^*$,

$\mu_3 \neq 0 \Rightarrow p_{1H}^* = 1$.

*Stationarity*

$\mu_2 = -c_1 f'_{p_{1L}^*} = -\frac{c_1 c_2 \gamma(p_{1H}^*, p_{1L}^*)}{g(p_{1H}^*, p_{1L}^*)} = -\frac{c_1 c_2 (1-\rho)}{g(1,1)}$.

*Dual Feasibility*

$\mu_2 \geq 0 \iff -\frac{c_1 c_2 \bar{\theta}(1-\rho)}{g(1,1)} = -(1-\rho) \geq 0$ which is never the case.

We can see that only candidates 1 and 2 comply with the KKT conditions. In particular, we notice that candidate 1 (i.e., $p_{1H}^* = 1, p_{1L}^* = 0$) holds whenever $\bar{\theta} \geq 2(1+u)$ whereas candidate 2 (i.e., $p_{1H}^* = 1, p_{1L}^* = \frac{2(1+u)-\bar{\theta}}{2(1+u)}$) holds whenever $\bar{\theta} \leq 2(1+u)$. Since these two candidates are mutually exclusive, we conclude that the KKT necessary conditions are also sufficient.

Finally, notice that in this case, when the manager sets the routing prioritization policy $p_{1H}^* = 1, p_{1L}^* = 0$ the prioritization error probabilities are $\delta_H = (1-s_H)(1 - (1+u)/\bar{\theta})$ and $\delta_L = 1 - s_L(1 - (1+u)/\bar{\theta})$. When the manager sets the routing prioritization policy $p_{1H}^* = 1, p_{1L}^* = \frac{2(1+u)-\bar{\theta}}{2(1+u)}$, the prioritization error probabilities are $\delta_H = (1-s_H)\bar{\theta}/4(1+u)$ and $\delta_L = 1 - s_L\bar{\theta}/4(1+u)$. In both routing policies, it is easy to see that $\delta_H + \delta_L < 1$, since $s_H + s_L > 1$ and $(1+u)/\bar{\theta} < 1$, which implies that the routing policy performs better than a FCFS policy. Finally, since the routing policies in cases 1 (i.e., $p_{1H} < p_{1L}$) and 2 (i.e., $p_{1H} = p_{1L}$) perform the

same as a FCFS policy. (To see why, one can easily compare the objective values in this case.) We conclude

that the routing policy in this case 3 performs better, and thus arises in equilibrium.  ■

## Appendix C:   Details of Experimental Results

| Condition | $\mathbb{P}(\text{claim} = \text{short queue})$ |
|:---:|:---:|
| 1 | $0.29 \pm 0.06$ |
| 2 | $0.35 \pm 0.06$ |
| 3 | $0.34 \pm 0.06$ |
| 4 | $0.22 \pm 0.05$ |
| 5 | $0.29 \pm 0.06$ |
| 6 | $0.29 \pm 0.06$ |
| 7 | $0.28 \pm 0.06$ |
| 8 | $0.23 \pm 0.05$ |
| 9 | $0.24 \pm 0.05$ |

**Table 2      Proportions and half widths of 95% confidence intervals of participants who reported the die roll 5**

**across experimental conditions.**

## Appendix D:   Experimental Analysis without Exclusions

In this section, we present the analysis of experimental results without the exclusions mentioned in the main

paper. Without exclusions, the sample size in each experimental condition is given in Table 4.

In Table 5 and Figure 2, we present the proportions of participants who reported the die roll 5 across

experimental conditions. Recall that participants have the incentive to report the number 5 to shorten their

waiting time.

We run two-sided exact binomial tests for the proportion of participants that reported the number 5

compared to the proportion that would have reported the number 5 under full honesty (i.e., 1/6). We find

that untruthfulness is significant under all conditions (p-values $< 0.009$). This can be readily seen in Figure

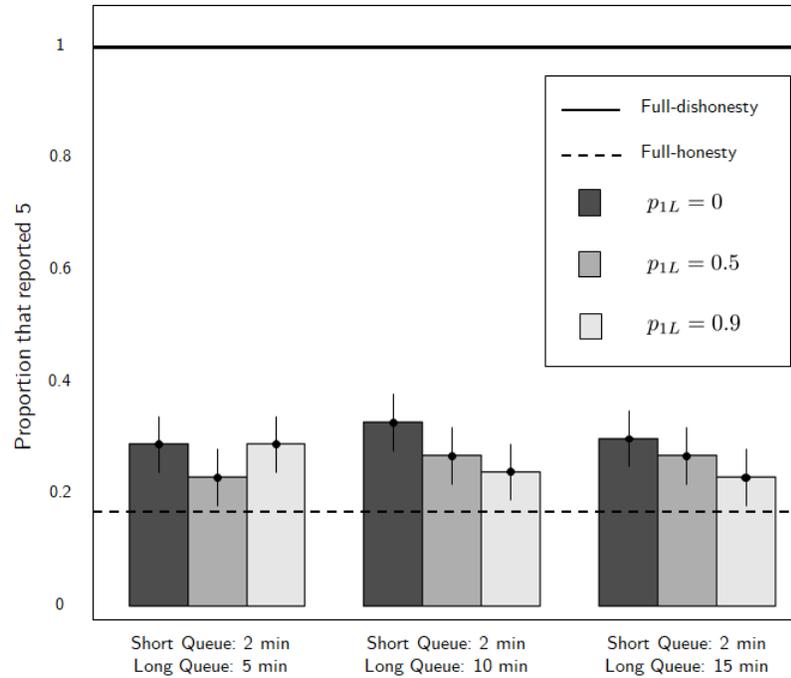2 where the black horizontal dashed line represents the expected proportion of customers who would have

**Figure 2**    **Proportions of participants that reported the number 5 across experimental conditions.**

reported the number 5 under the assumption of full honesty. Any value above the dashed line is based on untruthfulness (in expectation).

*Testing Hypothesis H1: The existence of lying aversion.* We find that participants misreport surprisingly little across all experimental conditions. It is clear from both Table 5 and Figure 2 that the proportions of participants who reported the number 5 is far from 1, i.e., the full-dishonesty black horizontal solid line in Figure 1. Indeed, for all conditions, the two-sided exact binomial tests for the proportion of participants that reported the number 5 compared to the proportion that would have reported the number 5 under full dishonesty (i.e., 1) are significant (p-values $= 0$). In the absence of any possible punishments or reputation concerns, any reluctance to misreport can be attributed to the psychological cost incurred due to internalized social preferences, i.e., the intrinsic lying cost in (5); see Fischbacher and Föllmi-Heusi (2013).

*Testing Hypothesis H2: The effect of the time incentive.* We run a Generalized Cochran-Mantel-Haenszel test for the conditional association between the proportion of participants who reported the number 5 and the time incentive treatment, conditional on the levels of the upgrade treatment, i.e., the probability $p_{1L}$. We find that the conditional association is not significant ($M^2 = 0.256$, $df = 2$, p-value $= 0.88$). We

also run a Chi-square test for the marginal association between the proportion of participants who reported the number 5 and the time incentive treatment. We find that the marginal association is not significant ($\chi^2 = 0.247$, $df = 2$, p-value = 0.88). These results indicate that a change in the time incentive does not affect significantly the misreporting behaviour of participants. Moreover, for each upgrade level (i.e., $p_{1L} = 0$, $p_{1L} = 0.5$, $p_{1L} = 0.9$), we run Jonckheere–Terpstra tests for an order, both ascending and descending. We find that the proportion of participants that claim the number 5 does not significantly increase monotonically (p-values $= 0.24, 0.21, 0.84$), nor decrease monotonically (p-values $= 0.62, 0.79, 0.15$) as the time incentive increases in each of the respective upgrading levels. Also, by aggregating the data (ignoring the upgrade treatment levels), both the ascending (p-value $= 0.48$) and descending (p-value 0.51) trends are not significant in the Jonckheere–Terpstra test. Finally, we run a logistic regression (see Table 6) and find that the coefficients for the different levels of the time incentive treatment are not significant in any of the model specifications. Overall, these results support the claim that a change in the time incentive does not affect participant untruthfulness behaviour. This is in line with our model equilibrium predictions.

*Testing Hypothesis H3: The effect of the upgrading control.* We run a Generalized Cochran-Mantel-Haenszel test for the conditional association between the proportion of participants who reported the number 5 and the upgrading treatment, conditional on the levels of the time incentive. We find that the conditional association is significant ($M^2 = 7.38$, $df = 2$, p-value = 0.02). We also run a Chi-square test for the marginal association between the proportion of participants who reported the number 5 and the upgrading treatment. We find that the marginal association is significant ($\chi^2 = 7.38$, $df = 2$, p-value = 0.02). These results indicate that the upgrading control affects participant untruthfulness behaviour. For each time incentive level (i.e., $2 - 5$ min, $2 - 10$ min, $2 - 15$ min), we run Jonckheere–Terpstra tests for an order, both ascending and descending. We find that the proportion of participants that claim the number 5 does not significantly increase monotonically (p-values $= 0.53, 0.95, 0.92$) as the upgrading increases in any of the respective time-incentive levels. For the case of a decreasing trend we find (p-values $= 0.46, 0.03, 0.08$). Also, by aggregating the data (ignoring the incentive treatment levels), the proportion that claim 5 significantly decreases monotonically (p-value 0.03). We run logistic regressions (see Table 6) and find that the interactions between the time incentive treatment and the upgrading treatment are not significant. Based on both the Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC), we find that the model specification that best fits the data among those considered in Table 6 is that of Model 1 (i.e., a model that only considers the

upgrade predictor). Moreover, we conduct likelihood ratio tests between Model 1 and Model 2 ($LR = 0.26$, $df = 2$, p-value $= 0.87$), and Model 1 and Model 3 ($LR = 4.89$, $df = 6$, p-value $0.55$) which show that the fit of Model 1 does not significantly improves by adding the time incentive predictor and the interaction term between time incentive and upgrading. Finally, the likelihood ratio test between Model 1 and the null model which only includes an in intercept term ($LR = 8.51$, $df = 4$, p-value $0.07$) shows that Model 1 is significant at $p - value < 0.1$ and the likelihood ratio test between model 1 and a reduced model (i.e., with age and gender as the only predictors) ($LR = 7.19$, $df = 2$, p-value $0.03$) shows that the upgrading predictor indeed improve model fit. Since Model 1 is the best fit among the considered models, we can observe that the coefficients for the upgrading treatment at both levels $p_{1L} = 0.5$ and $p_{1L} = 0.9$ are significant and of negative sign. Moreover, the absolute value of the coefficient (and marginal effect) for $p_{1L} = 0.9$ is larger than that of $p_{1L} = 0.5$. This indicates that the probability that a participant claims the number 5 decreases as the upgrading increases. Overall, these results reject hypothesis H3.

| | Model 1 | ME | Model 2 | ME | Model 3 | ME |
|---|---|---|---|---|---|---|
| (Intercept) | -0.71 *** | | -0.80 *** | | -0.88 *** | |
| | (0.19) | | (0.20) | | (0.23) | |
| Age | -0.00 | -0.00 | -0.00 | -0.00 | -0.00 | -0.00 |
| | (0.00) | (0.00) | (0.00) | (0.00) | (0.00) | (0.00) |
| GenderMale | 0.13 | 0.03 | 0.13 | 0.03 | 0.14 | 0.03 |
| | (0.10) | (0.02) | (0.10) | (0.02) | (0.10) | (0.02) |
| $p_{1L} = 0.5$ | -0.30 * | -0.06 * | -0.30 * | -0.06 * | -0.37 | -0.06 * |
| | (0.12) | (0.02) | (0.12) | (0.02) | (0.22) | (0.02) |
| $p_{1L} = 0.9$ | -0.40 ** | -0.08 *** | -0.40 *** | -0.08 *** | -0.06 | -0.08 *** |
| | (0.12) | (0.02) | (0.12) | (0.02) | (0.21) | (0.02) |
| (2-10 min) | | | 0.13 | 0.03 | 0.30 | 0.03 |
| | | | (0.12) | (0.02) | (0.20) | (0.02) |
| (2-15 min) | | | 0.12 | 0.02 | 0.22 | 0.02 |
| | | | (0.12) | (0.02) | (0.21) | (0.02) |
| (2-10 min): $p_{1L} = 0.5$ | | | | | 0.06 | |
| | | | | | (0.30) | |
| (2-15 min): $p_{1L} = 0.5$ | | | | | 0.14 | |
| | | | | | (0.30) | |
| (2-10 min): $p_{1L} = 0.9$ | | | | | -0.57 | |
| | | | | | (0.30) | |
| (2-15 min): $p_{1L} = 0.9$ | | | | | -0.44 | |
| | | | | | (0.30) | |
| N | 2021 | 0 | 2021 | 0 | 2021 | 0 |
| AIC | 2389.29 | 2389.29 | 2391.85 | 2391.85 | 2393.73 | 2393.73 |
| BIC | 2417.35 | 2417.35 | 2431.13 | 2431.13 | 2455.45 | 2455.45 |
| Pseudo R2 | 0.01 | | 0.01 | | 0.02 | |

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$.

**Table 3    Summary of Logistic Regression results. The models are logistic specifications with the fraction of**

**participants who claim 5 as the dependent variable. Columns with ME present the Marginal Effects associated to**

**the respective Models 1, 2 and 3.**

56

**Estrada, Ibrahim, and Zhan:** *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

| Condition | Short Queue | Long Queue | Upgrading prob $p_{1L}$ | Sample |
|:---:|:---:|:---:|:---:|:---:|
| 1 | 2 min | 5 min | 0 | 266 |
| 2 | 2 min | 10 min | 0 | 265 |
| 3 | 2 min | 15 min | 0 | 266 |
| 4 | 2 min | 5 min | 0.5 | 256 |
| 5 | 2 min | 10 min | 0.5 | 261 |
| 6 | 2 min | 15 min | 0.5 | 269 |
| 7 | 2 min | 5 min | 0.9 | 259 |
| 8 | 2 min | 10 min | 0.9 | 266 |
| 9 | 2 min | 15 min | 0.9 | 265 |

**Table 4**     Description of the experimental conditions.

| Condition | $\mathbb{P}(\text{claim} = \text{short queue})$ |
|:---:|:---:|
| 1 | $0.29 \pm 0.05$ |
| 2 | $0.33 \pm 0.06$ |
| 3 | $0.30 \pm 0.06$ |
| 4 | $0.23 \pm 0.05$ |
| 5 | $0.27 \pm 0.05$ |
| 6 | $0.27 \pm 0.05$ |
| 7 | $0.29 \pm 0.06$ |
| 8 | $0.24 \pm 0.05$ |
| 9 | $0.23 \pm 0.05$ |

**Table 5**     Proportions and half widths of 95% confidence intervals of participants who reported the die roll 5

across experimental conditions.

| | Model 1 | ME | Model 2 | ME | Model 3 | ME |
|---|---|---|---|---|---|---|
| (Intercept) | -0.76 *** | | -0.78 *** | | -0.85 *** | |
| | (0.18) | | (0.19) | | (0.21) | |
| Age | -0.00 | -0.00 | -0.00 | -0.00 | -0.00 | -0.00 |
| | (0.00) | (0.00) | (0.00) | (0.00) | (0.00) | (0.00) |
| GenderMale | 0.08 | 0.02 | 0.08 | 0.02 | 0.08 | 0.02 |
| | (0.09) | (0.02) | (0.09) | (0.02) | (0.09) | (0.02) |
| $(p_{1L} = 0.5)$ | -0.24 * | -0.05 * | -0.24 * | -0.05 * | -0.30 | -0.05 * |
| | (0.11) | (0.02) | (0.11) | (0.02) | (0.20) | (0.02) |
| $(p_{1L} = 0.9)$ | -0.27 * | -0.06 * | -0.27 * | -0.06 * | -0.01 | -0.06 * |
| | (0.11) | (0.02) | (0.11) | (0.02) | (0.19) | (0.02) |
| (2-10 min) | | | 0.05 | 0.01 | 0.19 | 0.01 |
| | | | (0.11) | (0.02) | (0.19) | (0.02) |
| (2-15 min) | | | 0.01 | 0.00 | 0.08 | 0.00 |
| | | | (0.11) | (0.02) | (0.19) | (0.02) |
| (2-10 min):$(p_{1L} = 0.5)$ | | | | | 0.03 | |
| | | | | | (0.28) | |
| (2-15 min):$(p_{1L} = 0.5)$ | | | | | 0.14 | |
| | | | | | (0.28) | |
| (2-10 min):$(p_{1L} = 0.9)$ | | | | | -0.44 | |
| | | | | | (0.27) | |
| (2-15 min):$(p_{1L} = 0.9)$ | | | | | -0.35 | |
| | | | | | (0.28) | |
| N | 2373 | 0 | 2373 | 0 | 2373 | 0 |
| AIC | 2782.06 | 2782.06 | 2785.80 | 2785.80 | 2789.16 | 2789.16 |
| BIC | 2810.92 | 2810.92 | 2826.20 | 2826.20 | 2852.65 | 2852.65 |
| Pseudo R2 | 0.01 | | 0.01 | | 0.01 | |

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$.

**Table 6** **Summary of Logistic Regression results. The models are logistic specifications with the fraction of participants who claim 5 as the dependent variable. Columns with ME present the Marginal Effects associated to the respective Models 1, 2 and 3.**

58

Estrada, Ibrahim, and Zhan: *On Customer (Dis)honesty in Priority Queues*
Article submitted to ; manuscript no. (Please, provide the manuscript number!)

## Appendix E:   Additional Experimental Results

We present an additional experiment which further supports the results presented in the main paper with respect to hypotheses H1 and H2. The design and procedure of the experiment is identical to the one presented in §5. In this experiment, the focus is on the time-incentive treatment alone. Particularly, we consider lower levels of time incentive than those in the main paper.

*Pre-registration.* We set the target sample size for the experiment and our analysis plans a priori. We pre-registered our experiment, and the corresponding As Predicted document can be found at: `https://aspredicted.org/blind.php?x=965nu4`.

*Participants and exclusions.* A total of 1,082 participants (41.68% female, mean age $M_{age} = 36.79$, standard deviation SD = 11.43) were recruited on the Amazon Mechanical Turk (MTurk) platform. Participants with at least 0.95 HIT approval ratio (proportion of completed tasks) were recruited to take part in our queueing experiment. Participants were instructed that they must wait in a virtual queue, then answer a two-question survey in exchange for a 1 US dollar payment. Participants were informed that the experiment would take up to 15 minutes. To ensure the independence of our observations, we exclude from our analysis the responses of 17 participants that presented the same Internet Protocol (IP) address.

In our experiment we have a completion rate of 95%. We exclude participants that did not complete the experiment. Moreover, while waiting in the virtual queue, participants are asked to click a button that appears every 30 seconds in order to move ahead in the queue. The remaining waiting time for participants in a given queue stops from elapsing until they click that button. This mechanism is used to ensure the continued attention of participants throughout their waiting in the virtual queue i.e., to ensure that they experience a real waiting cost. We recorded the time participants took to click each button and excluded those that presented and average time between clicks greater than 60 seconds. That is, we excluded participants that on average took longer than 30 seconds to click the buttons once they appeared. From our original sample, 97% of participants took on average less than 30 seconds to click the buttons, and the mean and median average click time was 6.16 seconds and 2.57 seconds respectively. After exclusions, we have a final sample of 976 participants (42.93% female, mean age $M_{age} = 37.16$, standard deviation SD = 11.57). Our results are unchanged by the aforementioned exclusions. For transparency, we include the analysis of results without exclusions as well.

*Experimental procedure.* The experimental procedure is identical to the one in §5. In this experiment, participants are randomly assigned to one of 4 experimental conditions (see Table 7) which differ in the waiting times for the short queue and the long queue. Participants who report the die roll 5 wait in the short queue, and those who report any other number wait in the long queue. Participants have the incentive to report the number 5 to reduce their waiting time.

| Condition | Short queue | Long queue | Sample size (with exclusions) | Sample size (without exclusions) |
|---|---|---|---|---|
| Baseline | 2 min | 2 min | 245 | 267 |
| 2 | 2 min | 3 min | 245 | 271 |
| 3 | 2 min | 4 min | 240 | 269 |
| 4 | 2 min | 5 min | 246 | 275 |

**Table 7    Description of the experimental conditions.**

### E.1.    Empirical Strategy and Results

In Table 8 and Figure 3 (with exclusions), we present the proportions of participants who reported the number 5 across experimental conditions.

| Condition | $\mathbb{P}(\text{claim} = \text{short queue})$ (with exclusions) | $\mathbb{P}(\text{claim} = \text{short queue})$ (without exclusions) |
|---|---|---|
| Baseline | $0.26 \pm 0.06$ | $0.27 \pm 0.05$ |
| 2 | $0.32 \pm 0.06$ | $0.32 \pm 0.05$ |
| 3 | $0.37 \pm 0.06$ | $0.36 \pm 0.06$ |
| 4 | $0.38 \pm 0.05$ | $0.37 \pm 0.06$ |

**Table 8    Proportions and half widths of $95\%$ confidence intervals of participants who reported the number 5 across experimental conditions.**
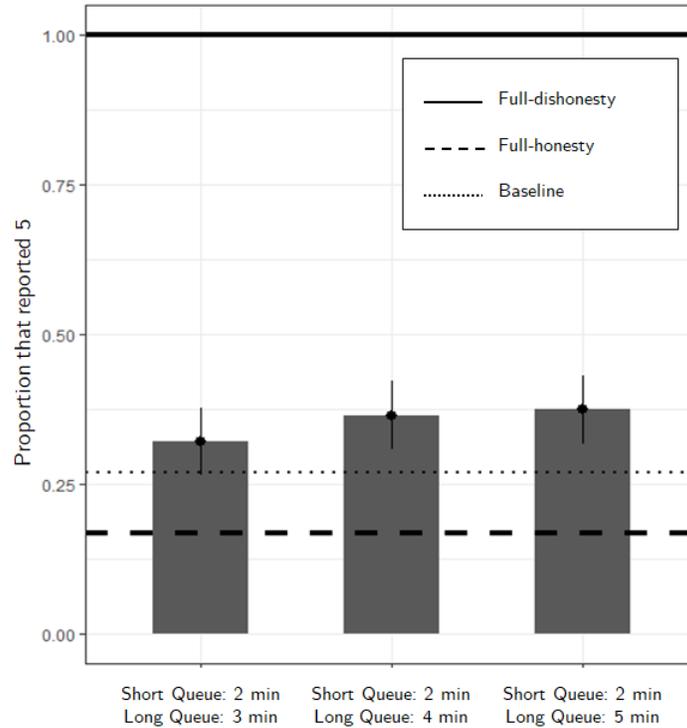
**Figure 3    Proportion of participants that reported the number 5 across experimental conditions.**

We run two-sided exact binomial tests for the proportion of participants who reported the number 5 compared to the proportion that would have reported the number 5 under full honesty (i.e., 1/6). We find that untruthfulness is significant under all conditions (p-values $< 0.001$).

*Testing Hypothesis H1: The existence of lying aversion.* We find that participants misreport surprisingly little across all experimental conditions. It is clear from Figure 3 that the proportion of participants that reported the number 5 is far from 1, i.e., the full-dishonesty black horizontal solid line in the figure. Indeed, for all conditions, the two-sided exact binomial tests for the proportion of participants who reported the number 5 compared to the proportion that would have reported the number 5 under full dishonesty (i.e., 1) are significant (p-values $= 0$). This result strongly supports the existence of lying aversion.

*Testing Hypothesis H2: Insensitivity to waiting times.* Notice that in this experiment we run a baseline condition (i.e., 2-2 min). The reason for this is that when there is no time incentive, our model predicts that people will randomize their claims. This boundary case represents a difference, at a qualitative level, in the reporting behaviour of people. We explore whether the reporting behaviour in this boundary case differs from the one where a time incentive is included. In other words, we explore whether the case of no waiting time

leads to a different claiming behavior than the case with waiting time. We will also check that increasing the

magnitude of the time incentive, once it is present, does not lead to a significantly different claiming behavior.

We run a Chi-square test for the association between the proportion of participants who reported the number

5 and the time incentive treatment plus the baseline (i.e., baseline, 2-3 min, 2-4 min, 2-5 min). We find that

the association is significant: $\chi^2 = 9.28$, $df = 3$, p-value = 0.02, with exclusions ($\chi^2 = 8.32$, $df = 3$, p-value

= 0.04, without exclusions). We then run a Chi-square test for the association between the proportion of

participants who reported the number 5 and the time-incentive treatment without the baseline (i.e., 2-3 min,

2-4 min, 2-5 min). We find that the association is not significant: $\chi^2 = 1.84$, $df = 2$, p-value = 0.39, with

exclusions, ($\chi^2 = 1.92$, $df = 2$, p-value = 0.38, without exclusions). These results indicate that *the presence*

of a time incentive affects significantly the misreporting behaviour of participants. Moreover, in the presence

of a time incentive, these results indicate that *a change* in the time incentive does not affect significantly

the misreporting behaviour of participants. Moreover, we run Jonckheere–Terpstra tests for an order in the

presence of a time incentive, both ascending and descending. We find that the proportion of participants that

claim the number 5 does not significantly increase monotonically (p-values = 0.14, with exclusions), (p-values

= 0.14, without exclusions), nor decrease monotonically (p-values = 0.85, with exclusions), (p-values = 0.85,

without exclusions) as the time incentive increases (i.e., 2-3 min, 2-4 min, 2-5 min). Finally, we run a logistic

regression (see Table 9) and find that the coefficients for the different levels of the time-incentive treatment

are not significant. Overall, these results support the claim that in the presence of a time incentive, a change

in such time incentive does not affect participant untruthfulness behaviour, which is in line with our model.

| | Model 1 | ME | Model 2 | ME |
|---|---|---|---|---|
| (Intercept) | -0.29 | | -0.28 | |
| | (0.31) | | (0.29) | |
| Age | -0.01 | -0.00 | -0.01 | -0.00 |
| | (0.01) | (0.00) | (0.01) | (0.00) |
| GenderMale | -0.07 | -0.02 | -0.06 | -0.01 |
| | (0.16) | (0.04) | (0.15) | (0.03) |
| (2-4 min) | 0.18 | 0.04 | 0.18 | 0.04 |
| | (0.19) | (0.04) | (0.18) | (0.04) |
| (2-5 min) | 0.23 | 0.05 | 0.22 | 0.05 |
| | (0.19) | (0.04) | (0.18) | (0.04) |
| N | 731 | 0 | 815 | 0 |
| AIC | 956.97 | 956.97 | 1063.50 | 956.97 |
| BIC | 979.94 | 979.94 | 1087.02 | 979.94 |
| Pseudo R2 | 0.01 | | 0.01 | |

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$.

**Table 9** **Summary of Logistic Regression results. The models are logistic specifications with the fraction of participants who claim 5 as the dependent variable. Columns with ME present the Marginal Effects associated to the respective Models 1 (with exclusions), 2 (without exclusions).**