


When More Is Less: Field Evidence on Unintended Consequences of Multitasking

Paulo B. Goes,^a Noyan Ilk,^b Mingfeng Lin,^c J. Leon Zhao^d

^a Eller College of Management, University of Arizona, Tucson, Arizona 85721; ^b Department of Business Analytics, Information Systems and Supply Chain, College of Business, Florida State University, Tallahassee, Florida 32306; ^c Department of Management Information Systems, Eller College of Management, University of Arizona, Tucson, Arizona 85721; ^d Department of Information Systems, College of Business, City University of Hong Kong, Kowloon Tong, Hong Kong

Contact: pgoes@eller.arizona.edu (PBG); nilk@business.fsu.edu,  <http://orcid.org/0000-0003-2152-288X> (NI); mingfeng@eller.arizona.edu (ML); jlzhao@cityu.edu.hk (JLZ)

Received: September 20, 2013

Revised: May 23, 2015; August 28, 2016

Accepted: January 15, 2017

Published Online in Articles in Advance: June 27, 2017

<https://doi.org/10.1287/mnsc.2017.2763>

Copyright: © 2017 INFORMS

Abstract. Online customer service chats provide new opportunities for firms to interact with their customers and have become increasingly popular in recent years for firms of all sizes. One reason for their popularity is the ability for customer service agents to multitask (i.e., interact with multiple customers at a time) thereby increasing the system “throughput” and agent productivity. Yet little is known about how multitasking impacts customer satisfaction—the ultimate goal of customer engagements. We address this question using a proprietary data set from an S&P 500 service firm that documents agent multitasking activities (unobservable to customers) in the form of server logs, customer service chat transcripts, and postservice customer surveys. We find that agent multitasking leads to longer in-service delays for customers and lower problem resolution rates. Both lead to lower customer satisfaction, although the impact varies for different customers. Our study is among the first to document the link between multitasking and customer satisfaction, and it has implications for the design of agent time allocation in contact centers and more broadly for how firms can best manage customer relations in new service channels enabled by information technology.

History: Accepted by Lorin Hitt, information systems.

Supplemental Material: The online appendix is available at <https://doi.org/10.1287/mnsc.2017.2763>.

Keywords: multitasking • customer satisfaction • service operations • information systems • IT policy and management

1. Introduction

Modern service organizations are under constant pressure to satisfy ever increasing service demands with limited operating budgets, of which human capital expenses constitute almost 70% (Human Capital Management Institute 2010). Not surprisingly, many service organizations adopt multitasking policies, under which employees handle multiple tasks simultaneously and switch from one task to another on a regular basis (Aral et al. 2012). This seems to be a natural and intuitive decision, when there is idle time during the execution of a particular task. Multitasking allows employees to better utilize their work time and to increase the number of “units” they can process in a given time frame, thus improving overall productivity. One of the most prominent applications of multitasking in the workplace can be observed in contact centers, especially in online live chats in which agents regularly shift their attention among different customers. Online customer service live chat has become a new channel for a firm to engage new customers and retain existing customers. It is being used not only in technology start-ups but also by many traditional service firms,

ranging from banks, car dealers, and insurance companies to hospitality, legal, and travel agencies (Reinsch et al. 2008). By putting more customers in touch with an agent at a given period of time, the benefit of multitasking in improving productivity seems obvious.

However, productivity of agents, especially in live chats, is much more than the “throughput” of how many customers an agent can handle. Customer satisfaction (Fornell et al. 1996), a critical factor in customer acquisition and retention, is a well-established aspect of any performance metrics for customer interactions that has rarely been addressed in the research and practice of multitasking. Firms will not be able to win new customers or retain existing customers if these customers are “processed” through queues but not satisfied. If customers are unable to resolve their issues satisfactorily via online visits, they may resort to more traditional means of customer contact, such as phone calls or in-person visits (Bavafa et al. 2017), thus rendering online chat irrelevant. Our goal is to fill this gap in the literature. Specifically, we address the following research question: *How does multitasking affect customer satisfaction in online service live chats?*

To investigate this problem, we utilize a comprehensive proprietary data set from the call center operations of an S&P 500 tax preparation services firm. This contact center receives requests from customers using a tax preparation software and provides tax filing and technical assistance through its human service agents, who communicate with customers via a live-chat interface. The environment is very dynamic with jobs arriving and leaving the system in a relatively short period of time. To increase the throughput of the online chat channel, the system automatically assigns multiple chat conversations to service agents. In other words, the arrival of customers in the queue and their assignment to agents are both exogenous; therefore, the multitasking level of agents is exogenously determined. This provides a uniquely ideal opportunity to study the impact of chat agent multitasking.

We focus on three important consequences of multitasking that have received relatively little attention in the empirical literature, especially in the chat center context. The first one is *problem resolution rates*. An inherent side effect of multitasking is its burden on the cognitive capacity of agents (Charron and Koehlin 2010). Specifically, task switching and interruptions due to multitasking may lead to difficulty with task reorientation and thus diminished performance (Adler and Benbunan-Fich 2012). We therefore hypothesize that multitasking will lead to lower problem resolution rates. The second consequence we study is *in-service delays* (i.e., response delays). When an agent chats with multiple customers, each customer will not receive as much attention as when they are the sole customer; therefore it will take longer for customers to hear back from agents, resulting in longer in-service delays. Furthermore, and even more important, service interruption and delays and lower problem resolution rates will lead to reduced customer satisfaction (Casado Díaz and Más Ruiz 2002, Tom and Lucey 1995). Given the critical importance of customer satisfaction for organizations (e.g., Fornell et al. 1996), the third and final outcome variable that we study in this paper is *customer satisfaction*.

In our empirical tests, we first examine problem resolution and response delay as a function of multitasking. We measure multitasking in two different ways, one focusing on the fact that an agent is *assigned to* multiple customers and the other focusing on the fact that an agent is *actively working with* multiple customers. These two complementary metrics allow us to open the “black box” of multitasking and understand how agents behave. We find that agents are less likely to resolve customer problems if they are actively engaged in two or more tasks concurrently. On the other hand, the assignment of a new customer alone is sufficient for the agents’ responsiveness to decrease, even if the agent is focusing on a single customer. We then turn

to the third outcome variable of interest and show that problem resolution and response delay further lead to negative impact of multitasking on customer satisfaction. These results are consistent after controlling for the mechanical effects of multitasking (such as delays solely caused by typing a longer text message). Our results are also robust to alternative specifications such as heterogeneity among agents, time-varying effects, and self-selection bias.

We further investigate whether such effects are uniform across varying agent workloads and different customer segments. We find that the marginal impact of having an additional customer is increasingly higher. Furthermore, we find that not all customers have the same degree of intolerance toward multitasking. To illustrate this, we use customer-specific variables observable at the time of their arrival and estimate a finite mixture model that accounts for heterogeneous customer responses to the negative impacts of multitasking. This suggests that contact centers can practically improve their customer-routing strategies—which customers should be assigned to multitasking agents, and which should not—to improve the system throughput while minimally reducing customer satisfaction.

Our study contributes to the literature and practice in several important ways. First, to our knowledge, our study is the first to address and quantify the relationship between multitasking and customer satisfaction, bridging these two fields that are each important in their own right. Second, we draw on data collected through multiple channels and multiple levels, including session-level and message-level data, as well as postservice customer follow-up survey data that are precisely linked to agents’ chat sessions. The extensiveness and comprehensiveness of the data is rare, to the best of our knowledge, in multitasking literature. Third, we go beyond the relationship between multitasking and satisfaction, and we address the heterogeneity of such impacts from both the agent’s and customer’s points of view. Our detailed findings can help managers develop and improve customer-routing strategies that can reduce the negative impact of multitasking, while retaining its productivity benefits.

2. Related Literature and Hypotheses

Literature on the relationship between workload, multitasking, and productivity has been long standing (Tan and Netessine 2014). For example, Aral et al. (2012) found that at low levels of multitasking (measured as the average number of projects simultaneously worked on), workers attain benefits from task complementarities and smoothing bursty work, leading to increased output. Cameron and Webster (2011) observed that busy workers may make themselves more accessible to their colleagues via multitasking. Reinsch et al. (2008)

argued that multitasking team members in a meeting will communicate more efficiently by initiating concurrent live-chat conversations among themselves. To our knowledge, however, there is no existing literature that directly links multitasking to customer satisfaction. To motivate our empirical analyses, we draw on several research streams to develop our hypotheses.

2.1. Multitasking and In-Service Delays

There has been a number of studies that considered the effects of multitasking during the process of a task (i.e., *in-service*). These studies can be broadly categorized into two groups. The first group approaches the problem from a queueing theory perspective and aims to develop analytical models for live-chat contact centers that enable multitasking. For instance, Campello et al. (2017) develop a stochastic model of how agents (case managers) process simultaneous jobs with multiple processing steps. Luo and Zhang (2013) model a contact center as a pool of many homogeneous servers, each operating under the processor-sharing protocol. They provide an asymptotic analysis and a fluid approximation of the system, which may help with staffing and admission control decisions. Tezcan and Zhang (2014) develop an analytical model for a similar system using linear programs (LPs). They propose a solution to the staffing LP and prove that this solution is asymptotically optimal. Furthermore, they characterize multitasking as a well-known queueing service discipline known as generalized processor sharing (Demers et al. 1989). Under this discipline, the server's work effort can be shared among the jobs in service using time-sharing mechanisms (Tan et al. 2005). A common approach to divide the work is to use round-robin scheduling in which time slots are assigned to each job in equal portions and in circular order. The server processes each job in turns and only for the amount of time defined by a job's allocated time slot. Evidently, the more jobs there are that share a single server, the more that these jobs wait idly between their turns.¹ It is important to note that such idle times are simply mechanical consequences of the multitasking activity (as a result of round-robin scheduling); the studies mentioned within this group do not consider additional task-switching-related factors.

On the other hand, the second body of literature emphasizes the switching costs associated with the multitasking activity and aims to specifically account for such costs. For example, Coviello et al. (2014) develop an economic production function that describes the slowdown in the output of a worker due to multitasking. On the basis of this function, they derive a law that can determine the output rate given the number of workers and the multitasking policy. Hall et al. (2015) define a

system in which machines are susceptible to work interruptions and task-switching costs. They include switching costs in their analyses and develop optimal algorithms for single-machine scheduling problems with multitasking.

Collectively, these two research streams emphasize the following two characteristics of the relationship between multitasking and in-service delays: (1) multitasking causes in-service delays between processing periods as a result of the nature of round-robin scheduling, and (2) these delays can be further inflated because of switching costs associated with shifting between different tasks. In our study, we expect to see similar effects in the form of service agent responsiveness. Agent responsiveness in a live-chat session can be defined as the agent's punctuality (i.e., promptness) in replying to individual customer messages. A typical chat conversation consists of multiple text messages being transmitted between the agent and the customer via a chat interface. These messages are generally short so that the other participant can respond quickly, thus creating a feeling of face-to-face conversation. According to an industry benchmark survey by TELUS International (2011), agent responses in a live-chat session should be provided within 30 seconds to generate a sense of real-timeliness. Otherwise, longer periods of idle time between messages could be considered in-service delays by customers.

Within the context of live-chat communication, we posit that multitasking will lead to increased in-service delays because of its processor-sharing nature as well as the switching costs associated with it. In particular, we expect the impact on in-service delays to grow nonlinearly with the increasing levels of multitasking activity.² We therefore hypothesize as follows.

Hypothesis 1 (H1). *Multitasking increases agent response delays in live-chat conversations. The marginal effect of more multitasking on response delays should be increasing.*

2.2. Multitasking and Task Resolution

There exists a rich body of research on the cognitive effects of multitasking including dual-task interference, increased cognitive load, and task accuracy (Cameron and Webster 2013). Psychology literature notes that loss of context associated with switching tasks can be a major delay factor during the resumption of the initial task (Czerwinski et al. 2004). Mark et al. (2008) argue that workers try to compensate for the loss of performance due to interruptions by working faster, which comes at a price of increased stress and frustration. Along with an increased multitasking workload, stress and frustration may result in cognitive overload, which leads to processing mistakes. In the end, service quality and worker performance may suffer (KC 2013).

There have been studies that investigated the validity of this argument in experimental settings. Using

brain imaging technologies, Charron and Koechlin (2010) show that frontal brain function is vulnerable to mistakes when an individual pursues two concurrent goals simultaneously. Adler and Benbunan-Fich (2012) conduct a controlled experiment to compare multitasking and non-multitasking conditions, and they find that increased levels of multitasking lead to a significant loss in accuracy. Bailey and Konstan (2006) demonstrate that task interruptions and continuous task switching increase the stress and anxiety levels of individuals and lead to up to twice the number of errors committed across the tasks.

A limitation with these studies is the simplistic and artificial nature of tasks that are simulated within the lab environment. Our paper aims to contribute to this literature by studying multitasking in the context of complex, real-world tasks. Within the contact center environment, we postulate that service agents will be less effective in handling customer problems when they are engaged in multiple conversations at the same time. Accordingly, we hypothesize the following.

Hypothesis 2 (H2). *Multitasking reduces the agents' likelihood of resolving customer problems in live-chat conversations.*

2.3. Determinants of Customer Satisfaction

Customer satisfaction is a quintessential cornerstone for service organizations to achieve success and long-term profitability (Anderson et al. 1997). It is a complex construct with a variety of determinants including expectations, disconfirmation, performance, affect, and equity (Szymanski and Henard 2001). In the context of contact centers, we argue that in-service delays and task resolution are two of the major determinants of customer satisfaction.

To understand the impact of in-service delays on customer satisfaction, we must first distinguish different types of customer-centric waiting times. A common categorization is to define different types of waiting based on the point of time at which the wait is initiated. Customers can wait before the process (*preservice*), during the process (*in-service*), and after the process (*postservice*) (Dube-Rioux et al. 1989). Taylor (1994) defines *delay* as the preservice, postschedule wait, which occurs when a scheduled event does not begin on time. Casado Díaz and Más Ruiz (2002) broaden the concept of delay to in-service wait using a field study of delayed flights at an airport. In both studies, common results indicate that longer delays result in greater anger among customers, and anger leads to diminishing customer satisfaction with services. Similarly, Chebat and Filiatrault (1993) show that interrupted service (such as commonly seen in multitasking) has a negative effect on the perceived waiting duration, clients' mood, and the service quality. Furthermore, all of these studies emphasize the scheduled

aspect of delays. A delay is different than waiting in queue in the sense that it is a broken promise; customers are scheduled to receive services but are kept on hold during the process. Response delay in live-chat conversations is a similar form of this situation. Inevitably, customers would expect to receive timely responses to their messages once they have been admitted into service.

A consensus in the literature is that the outcome of a service encounter significantly and directly affects customer perceptions of the quality of service (Brady and Cronin 2001). In a seminal study, Zeithaml et al. (1996) find that customers who receive satisfactory resolution of their service problems have significantly higher loyalty and retention intentions, along with increased willingness to pay a premium for the service. They conclude that effective service resolution and recovery significantly improve all facets of customers' behavioral intentions. Later studies have solidified the argument that task resolution incompetency, along with longer in-service delays, would frustrate customers and lead to a decline in overall customer satisfaction. In this study, our goal is to link multitasking to customer satisfaction. We expect that such effects will be indirect—that is, they will be transmitted via task resolution and in-service delays—since customers cannot directly observe that the agent is multitasking. Accordingly, we hypothesize the following.

Hypothesis 3A (H3A). *Multitasking in live-chat conversations indirectly affects customer satisfaction via agent response delays.*

Hypothesis 3B (H3B). *Multitasking in live-chat conversations indirectly affects customer satisfaction via problem resolution.*

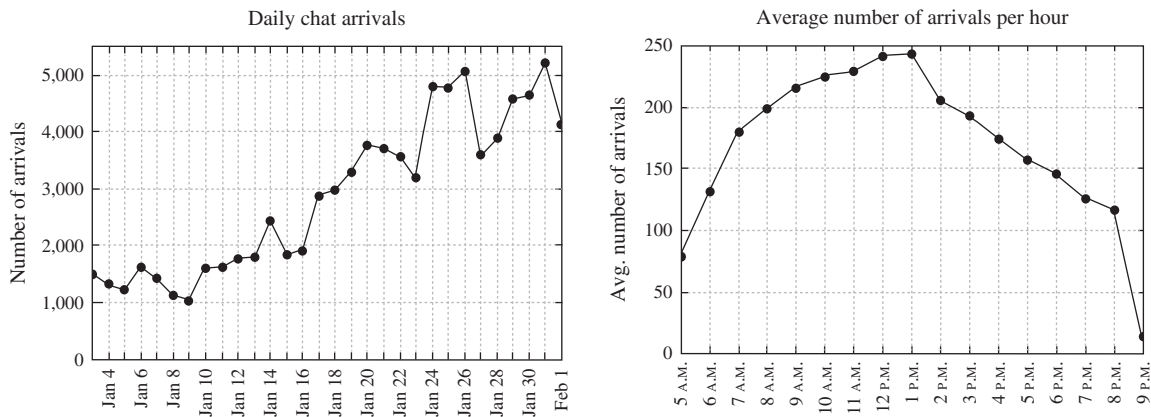
In addition to testing these hypotheses, we also develop and estimate a finite mixture model that takes possible customer heterogeneity into account while studying the impact of multitasking. This model has the flexibility to classify customers into segments based on customer-specific characteristics and simultaneously estimates the impact of multitasking on customer satisfaction for individual segments.

3. Empirical Context and Data

3.1. Contact Center

To estimate our models, we utilize multiple sources of data from a large-scale live-chat contact center between January 3, 2011 and February 2, 2011. The contact center is operated by a firm that develops tax preparation and filing support software and services. Many firm customers require customer support for technical problems as well as tax filing assistance. The firm categorizes incoming service requests into 20 distinct skill types and trains each service agent to respond to

Figure 1. Number of Chat Arrivals to the System



one or more types. In addition, skill types are classified into three broader domains with respect to the context of the conversation. Because of the intricate nature of U.S. income tax laws and the dire consequences of filing incorrect tax returns, the contact center receives a significant amount of traffic, especially during the tax filing season.

To help meet this demand, service agents multitask whenever there are customers waiting in the queue, with up to four parallel chat sessions. The queues are formed at the department level, and a particular queue is served by multiple agents. Customer-agent assignments are made automatically³ by the system on a first-come-first-served basis. Furthermore, multitasking assignments are determined by the system according to the length of the queue. The contact center operates from 5 A.M. to 10 P.M. daily. Agents are only active at certain time periods during the day, with four hours of work on average each day. Figure 1 depicts the daily and average hourly chat requests in the system.

3.2. Data

We access three data sources within the contact center and conduct an extensive data preparation and consolidation process to construct the measures and variables relevant to our analysis. The first data source is transaction server logs. Also called the metadata, this set of records includes information about each chat transaction’s arrival time to and departure time from the system, queue details, service-request-specific information such as skill key, information about agent assignments, and the amount of time spent in queue as well as in service. This data source is also used as one of the inputs to extract information about the agents’ multitasking activities.

The second data source is a repository that stores the actual chat conversations in the system in text form (i.e., transcripts) and initial problem descriptions provided by customers when they arrived in the queue. Chat transcripts consist of all the messages exchanged

between the agent and the customer. We analyze these transcripts to extract the time stamp of each message. For each customer message, we compute how long it took for the agent to respond to the customer by taking the difference between time stamps. We also check how many customers the agent was talking to during this response delay to obtain multitasking information.⁴ Specifically, we count all open sessions during the delay period as well as the effort of the agent (e.g., number of messages sent and number of words typed). Figure 2 illustrates this idea for a particular message in an illustrative conversation timeline. We also collect customer message lengths and customer response delays prior to agent messages to be used as control variables.

In addition, chat transcripts are accompanied by a text-based description of the problem and the reason for the contact request. Customers are required to enter this information before they are admitted into the system. Problem descriptions are important because they may reveal the initial mood of customers, which may influence their eventual satisfaction. We combine all problem description text with the first three messages of the customer and employ sentiment analysis (Pennebaker et al. 2001) on the aggregated text to extract the mood of customers at the beginning of the chat sessions.

Figure 2. (Color online) Illustrative Agent Responses in a Conversation Timeline

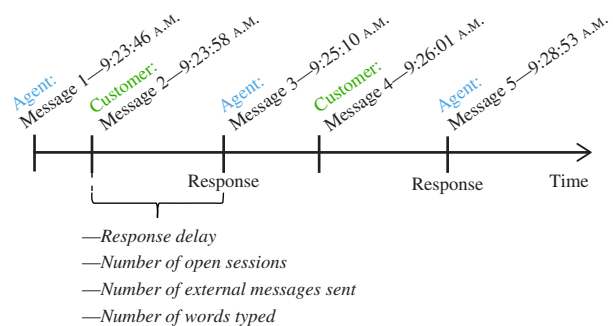


Table 1. Data Sources

	Transaction server logs	Chat transcripts	Survey response
<i>Chat Session ID</i>	✓	✓	✓
<i>System features</i>	✓		
<i>Session features</i>	✓		
<i>Problem description</i>		✓	
<i>Multitasking info.</i>	✓	✓	
<i>Customer satisfaction</i>			✓
<i>Problem resolution</i>			✓

The third data source contains the survey response information. Each customer completing a chat session receives an e-mail survey (sent one day after the chat session) from the company with questions about their chat experience, satisfaction, and problem resolution. During the period that our data covers, the overall survey response rate was 16%. A summary of the information contained in individual data sources is provided in Table 1.

3.3. Data Cleaning

We preprocess the consolidated data to remove the observations that have systematically different chat characteristics than the rest of the sample. First, we find and remove all chat sessions that were transferred between service agents. Second, we find and remove all agent-to-agent chat conversations.⁵ While the live-chat system permits agents to communicate with other agents, these conversations have different intent and dynamics than typical customer-agent interactions. Finally, we remove extreme observations from the consolidated data set. A rare but important situation is the system’s failure to properly terminate a chat conversation. In such cases, certain variables (e.g., service duration, response delay, customer slowness) can get artificially inflated. To mitigate this, we exclude observations for time-related variables whose values exceed the 99th percentile. If there exists an extreme observation at the message level, we remove the entire session related to this observation. Another issue is copying and pasting text while crafting responses. In our message-level analysis, we control for the delays associated with the agent’s typing behavior (i.e., mechanical efforts) using the number of words typed. However, on rare occasions, agents copy and paste long chunks of standard text into messages. For such cases, message lengths do not accurately reflect agents’ efforts. To mitigate this, we search for messages that were extremely long (over the 99th percentile) yet took less than the 50th percentile of the overall response time. Then, we replace the length of these messages with the mean message length in the data set. Table 2 shows the contact center summary statistics after data preprocessing and cleaning steps.

Table 2. Summary Statistics of the Contact Center (After Data Cleaning)

Total number of chat sessions	86,159
Total number of surveys returned	14,125
Total number of agents	1,209
Total number of distinct skills	20
Total number of agent responses	600,376
Avg. number of skills that an agent has	6
Avg. number of agent responses in a session	7
Avg. number of messages exchanged in a session	17
Avg. number of words typed per agent response	16

3.4. Variable Definitions and Measurement

We conduct our empirical analysis at two different levels: message and session. Message-level analysis is used to investigate the impact of multitasking on individual agent response delays (i.e., in-service delays). Session-level analysis is used for two purposes: (1) to investigate the impact of average session multitasking on problem resolution and (2) to investigate the indirect impact of average session multitasking on customer satisfaction. After data collection and processing, we construct the following variables.

3.4.1. Dependent Variables

Response Delay. To operationalize the response delay, we first define the concept *response*. A response is an agent message that follows the previous customer message. To be considered a response, the agent message should be directly preceded by the customer message in the chat conversation (see Figure 2). As a result, not all agent messages will be responses and not all customer messages will have a response. For each pair of customer message and agent response, we find the time that a customer waits to receive an agent response. This value gives us the response delay for the customer message. We account for back-to-back messages as well as canned agent messages or responses in our computations. For back-to-back customer messages, we define the response delay as the gap between the last customer message and the following agent message in our main models. As a robustness check, we also use the gap between the first customer message and the following agent response. For canned messages, we identify three types of messages in the data—greetings, goodbye messages, and survey reminders—and exclude them from our computations and analysis. We measure the response delays in seconds.

Problem Resolution. Problem resolution is an outcome indicator for the chat conversation that indicates whether or not the customer’s problem was successfully resolved by the agent at the end of the conversation. We collect problem resolution information from customer survey responses. This variable is coded in binary, with 1 indicating that the problem was resolved and 0 otherwise.

Customer Satisfaction. Customer satisfaction is measured at the session level and reflects the customer's evaluation of the service encounter (the related chat session) (Tax et al. 1998). We collect customer satisfaction scores from two customer survey questions. The first question is a five-scale item that asks whether the customer's expectations were fulfilled at the end of the chat transaction. Customers choose from one of the following options: *fell way short of my expectations*, *fell somewhat short of my expectations*, *met my expectations*, *exceeded my expectations*, or *greatly exceeded my expectations*. The second question asks, based on the last chat transaction, what the likelihood is that the customer will promote the service to his or her friends (i.e., net promoter score). This question is measured using a 0–10 scale, with 10 being the highest. The company further “bins” this scale into three levels: 0–6 as *low*, 7 and 8 as *medium*, and 9 and 10 as *high*.

3.4.2. Independent Variables

Multitasking. Multitasking is described as the act of working on multiple jobs in parallel (Aral et al. 2012). In the live-chat context, we measure multitasking using two metrics. The first metric (henceforth MT_1) captures the substantial workload of the agent by counting the number of unique sessions in which the agent sent a message during a particular response delay. This metric focuses on the active attention of the agent and helps identify multitasking at a very granular level. However, it does not consider existing sessions in which the agent may be inactive. To consider the possibility that in-session but inactive sessions may still have an effect on the agent's cognitive capability, we define a second metric (henceforth MT_2) based on the number of all sessions open on the agent's service screen during a particular response delay.

For both metrics, we measure multitasking at the message level. A particular benefit of this approach is the ability to capture exogenous variations in multitasking values. As discussed earlier, service assignments are made automatically by the system. The live-chat environment is dynamic, and customers arrive and leave at random time points. As a result, even during one agent-customer conversation, the agent's level of multitasking values may exogenously change from time to time. We also include quadratic forms of the MT metrics in the message-level analysis to account for potential nonlinearity. For session-level models, we average the multitasking values of individual responses in the chat session: the multitasking rate of a chat session ($AvgMT$) is defined as the average of multitasking values of individual responses (MT) in that session.

3.4.3. Other Variables for Message-Level Analyses

Skill Heterogeneity. Contact centers commonly cross-train their staff in multiple skills to increase agent utilization (Tekin et al. 2009). It is expected that both the count and the heterogeneity of jobs that are being multitasked will have an influence on agent performance and service quality. To account for this issue, we observe heterogeneity among multitasked chat sessions. We use skill keys (problem types defined by the company) to distinguish different skill sets required in different sessions. For each agent response, we calculate the distribution of unique skill keys that are being worked on by that agent at that time using a Herfindahl Index. We construct two such heterogeneity indices to correspond to each of our multitasking metrics. For MT_1 , the heterogeneity measure is constructed as $SkillHetMT_1 = 1 - \sum_{a=1}^s (m_{ja}/m_j)^2$, where s is the total number of skill keys in the data set, m_{ja} is the number of sessions that have skill key index a and that were sent at least one message during response delay j , and m_j is the number of all sessions that were sent at least one message during response delay j . For MT_2 , the heterogeneity measure is constructed as $SkillHetMT_2 = 1 - \sum_{a=1}^s (o_{ja}/o_j)^2$, where s is the total number of skill keys in the data set, o_{ja} is the number of sessions that have skill key index a and that were open during response delay j , and o_j is the number of all open sessions during response delay j . To reflect the heterogeneity of skill sets, we subtract the summations from 1 in both measures, so that higher values correspond to higher heterogeneity of skills required in the sessions that the agent is simultaneously working on.

Multitasking Variation. It is possible that the variability of the multitasking during the session may lead to inconsistency in service delivery, which may have an additional effect on problem resolution and customer satisfaction. To consider this issue, we define a multitasking variation variable based on the standard deviation of multitasking within the chat session.

Typing Effort. The textual length of an agent's message is expected to affect response delays. Naturally, the longer the text that the agent types, the longer it would take to send the message. We consider such effects to be *mechanical* because of the nature of text-based communication. To disentangle the mechanical delays caused solely by typing a text message, we include a message length control in our analyses. We operationalize this variable by counting the total number of words typed by the agent in all the messages that he or she sent during a particular response delay.

Customer Message Length. The length of the customer message preceding the agent response may affect the delay of the response for two main reasons. First, as the customer's message gets lengthier, it would

take longer for the agent to read the entire message. Second, agents may be less enthusiastic about responding to customer messages that are wordy and long-winded. We therefore control for the customer's message length. This variable is defined as the length of the customer message preceding the response for each response in a given session.

Customer Slowness. Even though online chats are inherently asynchronous, they are still interactive; a smooth conversation will require the engagement of both the customer and the agent. If the customer is not responsive to agent queries, the agent may be less likely to respond promptly. To account for this issue, we compute the moving average of a customer's message delays (i.e., time delays between agent messages and corresponding customer messages) and include it in the message-level model.

Message Order. Not all agent responses are expected to take a similar amount of time. To account for the heterogeneity in response times due to the position of the response, we include an order variable in our multitasking versus response delay model. We bin the message order into three groups: the first two messages as the greeting phase, the last two messages as the feedback phase, and in-between messages as the solution phase.

Session Key. We use session fixed effects in the multitasking versus response delay model to analyze multitasking at the granularity level of a message. Session dummies help us control for any time-invariant characteristics and unobserved heterogeneity within individual chat sessions.

3.4.4. Other Variables for Session-Level Analyses

Queue Wait Time. Waiting time in queue refers to the amount of the time a customer waits after arriving to the system and before being admitted into service by an agent. It is widely considered to be a major determinant of service quality and customer satisfaction (Durrande-Moreau 1999). Furthermore, waiting time in queue is expected to correlate with multitasking, since company policies dictate that agents switch to multitasking when queues develop in the system. To disentangle these effects, we control for queue wait time in our models. Information regarding waiting times for each customer is extracted from transaction server logs. Waiting times are measured in seconds.

Session Duration (Service Time). Session duration is the total amount of time a customer spends in service, between the time the customer were admitted into service and the time he or she left the system. Similar to queue wait time, this information is obtained from transaction server logs and measured in seconds.

Agent Experience. For session-level models, the impact of multitasking may depend on agents' experience levels with the chat sessions. For example, an agent may become more skillful in solving a problem as he or she processes more problems of similar type. To account for this issue, we include an agent experience control in our models. We operationalize agent experience by counting the number of same skill key sessions that the agent has processed during the data collection month before each session.

Negative Emotion Prior. Customers arriving at a chat queue may already be in a particular emotional state. The positivity or negativity of that state will not only affect their interaction with the agent but also their ultimate evaluation of the entire customer service engagement. We therefore aim to control for such preexisting emotional states of customers in our analysis. The problem description that customers are prompted by the online chat system to enter is uniquely ideal for this purpose because it is written before customers are assigned to any particular agent. But since these texts are typically short and therefore not conducive to textual analyses, we combine them with the first three messages customers write after they enter the chat sessions. We then conduct sentiment analysis on those texts to extract the negative emotion of each customer at the beginning of a chat session using the standard text-mining software package LIWC (Pennebaker et al. 2001).

Survey Response Gap. The time difference between the date of a chat session and the survey response date for this chat session could be a factor affecting the customer satisfaction responses in the survey. For example, customers' scoring of an unpleasant transaction may be less harsh as more time passes. To account for this issue, we include a survey response gap variable in our model of customer satisfaction. The survey response gap is computed as the number of days between the chat date and the survey response date. Survey links are sent one day after the chat sessions.

Description Length. In our analyses, we control for the length of the initial customer message before he or she is connected to a chat agent.

Agent Key. During the data coverage period, there were 1,209 unique service agents who conducted chat sessions with customers. Naturally, these agents have difference experience levels as well as different behaviors. We include agent fixed effects in our models to account for unobserved heterogeneity across individual agents.

Skill Key. To control for differences among problem types, we include skill key dummies in our session-level models.

Table 3. Descriptive Statistics of Message-Level Variables

Variable	Obs.	Mean	Std. dev.	Min	Median	Max
(a) <i>RDelay</i>	600,376	56.47	60.04	1.00	38.00	597.00
(b) <i>MT</i> ₁	600,376	1.30	0.51	1.00	1.00	4.00
(c) <i>SkillHetMT</i> ₁	600,376	0.10	0.20	0.00	0.00	0.75
(d) <i>MT</i> ₂	600,376	1.91	0.66	1.00	2.00	4.00
(e) <i>SkillHetMT</i> ₂	600,376	0.28	0.26	0.00	0.44	0.75
(f) <i>Typing</i>	600,376	25.37	27.14	1.00	17.00	883.00
(g) $\log(\textit{Typing})$	600,376	2.76	1.06	0.00	2.83	6.78
(h) <i>CLength</i>	600,376	9.42	10.60	1.00	6.00	110.00
(i) $\log(\textit{CLength})$	600,376	1.91	0.94	0.00	1.95	4.71
(j) <i>CSlowness</i>	600,376	37.18	36.06	1	26	259
(k) $\log(\textit{CSlowness})$	600,376	3.17	1.05	0	3.26	5.56

Product Key. The contact center is operated by a company that develops tax preparation and filing support software. The company offers multiple software versions and products. These products vary based on their tax filing features and prices. To account for the inherent variation among different products, in our session-level models we include dummy variables for product keys.

Gender. Demographics variables such as gender are commonly used to determine membership in customer segmentation studies (Gupta and Chintagunta 1994). We use gender in our finite mixture model to classify a customer in terms of their tolerance of multitasking. To identify the gender, we extract each customer’s first name from chat transcripts using text processing techniques and then run the extracted name through a baby naming dictionary⁶ via an automated script.

Contact Date and Contact Time. To account for temporal variations, we include contact date and contact time interval dummies as controls.

Descriptive statistics for these variables are given in Tables 3–6. In total, we analyze 14,125 chat sessions and 600,376 individual response messages. We note that message-level and session-level models use their corresponding variables at their own level of granularity. Before estimating the models, we transform control variables into their natural logarithms to reduce skewness.⁷ As a robustness test, we also use these variables in the estimations without log transformations. Tables 3 and 4 provide the summary statistics and correlation results for message-level variables, whereas Tables 5 and 6 provide the same information for session-level variables.

4. Empirical Analyses and Results

4.1. Multitasking and Response Delay

4.1.1. Model Specification and Estimation. To investigate the relationship between multitasking and response delays, we first run ordinary least squares (OLS) estimations with session fixed effects using a

Table 4. Correlation Matrix for Message-Level Variables

	(a)	(b)	(c)	(d)	(e)	(f)	(g)	(h)	(i)	(j)	(k)
(a)	1.00										
(b)	0.39	1.00									
(c)	0.28	0.79	1.00								
(d)	0.14	0.43	0.37	1.00							
(e)	0.09	0.27	0.47	0.69	1.00						
(f)	0.45	0.62	0.46	0.26	0.16	1.00					
(g)	0.39	0.58	0.46	0.23	0.14	0.79	1.00				
(h)	0.05	0.01	0.02	0.04	0.04	0.03	0.05	1.00			
(i)	-0.01	-0.02	0.00	0.04	0.04	0.00	0.01	0.88	1.00		
(j)	0.05	0.02	0.03	0.05	0.04	0.02	0.01	0.34	0.33	1.00	
(k)	0.06	0.03	0.03	0.04	0.03	0.03	0.02	0.29	0.32	0.80	1.00

Note. See Table 3 for definitions of labels.

panel data set.⁸ The panel data set includes observations of each individual agent response within the entire chat session and for all chat sessions. Each observation is taken at the time of an agent response; hence, our unit of analysis is an individual message.

The dependent variable is the response delay (*RDelay*), and the key independent variable is the multitasking amount per response (*MT*), measured in two different forms. We include the quadratic form of the multitasking variable (*MT*²) to check for possible nonlinearity. We use skill heterogeneity across the processed sessions (*SkillHetMT*) in the estimation to account for possible cognitive delays due to switching between dissimilar tasks. Other control variables include agent’s typing effort (*Typing*), message length of the preceding customer message (*CLength*), the customer’s average message delay up to the particular response (*CSlowness*), and the dummy variables for the response order within the session (*Order*).

Table 5. Descriptive Statistics of Session-Level Variables

Variable	Obs.	Mean (freq.)	Std. dev.	Min	Median	Max
(l) <i>AvgMT</i> ₁	14,125	1.32	0.32	1.00	1.25	4.00
(m) <i>AvgSkillHetMT</i> ₁	14,125	0.10	0.13	0.00	0.06	0.67
(n) <i>sdMT</i> ₁	14,125	0.36	0.28	0.00	0.45	2.12
(o) <i>AvgMT</i> ₂	14,125	1.94	0.60	1.00	2.00	4.00
(p) <i>AvgSkillHetMT</i> ₂	14,125	0.30	0.24	0.00	0.36	0.75
(q) <i>sdMT</i> ₂	14,125	0.19	0.25	0.00	0.00	1.41
(r) $\log(\textit{QueueWait})$	14,125	3.77	2.57	0.00	4.16	7.88
(s) $\log(\textit{Duration})$	14,125	6.91	0.73	3.78	6.90	8.82
(t) $\log(\textit{AExperience})$	14,125	2.88	1.66	0.00	2.83	6.94
(u) $\log(\textit{NegEmotion})$	14,125	0.28	0.46	0.00	0.00	2.48
(v) $\log(\textit{DescLength})$	14,125	3.93	0.63	0.69	3.99	6.45
(w) $\log(\textit{TotalTyping})$	14,125	5.01	0.81	0.00	5.13	7.30
(x) $\log(\textit{AvgCLength})$	14,125	2.30	0.55	0.69	2.32	4.30
(y) $\log(\textit{AvgCSlowness})$	14,125	3.62	0.49	0.00	3.62	5.55
(z) $\log(\textit{SurveyGap})$	14,125	0.95	0.44	0.69	0.69	5.06
(aa) <i>Resolution</i>	14,125	0.63	0.48	0.00	1.00	1.00
(bb) <i>Satisfaction</i> ₁	14,021	2.99	1.40	1.00	3.00	5.00
(cc) <i>Satisfaction</i> ₂	14,069	2.24	0.88	1.00	3.00	3.00

Table 6. Correlation Matrix for Session-Level Variables

	(l)	(m)	(n)	(o)	(p)	(q)	(r)	(s)	(t)	(u)	(v)	(w)	(x)	(y)	(z)
(l)	1.00														
(m)	0.78	1.00													
(n)	0.66	0.51	1.00												
(o)	0.62	0.53	0.58	1.00											
(p)	0.39	0.67	0.38	0.70	1.00										
(q)	0.05	0.03	0.22	-0.02	-0.03	1.00									
(r)	0.30	0.25	0.30	0.46	0.33	-0.08	1.00								
(s)	0.06	0.05	0.22	0.02	0.03	0.27	0.09	1.00							
(t)	0.05	-0.14	0.06	0.05	-0.23	-0.03	0.05	-0.11	1.00						
(u)	0.01	0.03	0.03	0.05	0.07	0.02	0.02	0.07	-0.06	1.00					
(v)	-0.01	0.04	-0.03	0.00	0.06	-0.01	-0.11	0.01	-0.11	0.27	1.00				
(w)	0.36	0.27	0.49	0.18	0.10	0.26	0.16	0.60	0.02	0.02	-0.10	1.00			
(x)	0.03	0.06	0.06	0.06	0.09	0.04	0.04	0.16	-0.06	0.16	0.47	0.05	1.00		
(y)	0.10	0.11	0.10	0.11	0.11	0.08	0.04	0.29	-0.05	0.08	0.25	0.05	0.41	1.00	
(z)	0.00	0.00	0.02	-0.01	0.00	0.02	0.02	0.03	0.01	-0.01	0.00	0.03	-0.01	0.03	1.00

Note. See Table 5 for definitions of labels.

An econometric concern in our analysis is unobserved heterogeneity. Session fixed effects help us control for time-invariant characteristics of chat sessions such as unique agent characteristics and problem variation as well as unobserved heterogeneity for time-of-day and date effects. A second concern in our analysis is simultaneity—that is, multitasking may be caused by delays between messages. However, this is an unlikely scenario in our context because it is the chat system policy that dictates admission control decisions. It is not up to the agents to decide whether they have more than one customer at a time and who that customer is. And since customer arrival and departure times are also largely random, the multitasking variable is not only exogenous but also can vary exogenously from one message to another, even within the same chat conversation. These features allow us to identify the effect of multitasking, measured in the two different ways that we described, on service performance.

4.1.2. Results. Table 7 reports the first estimation results of various specifications. Columns (1) and (2) employ two different multitasking variables independently. Subsequent specifications use both multitasking variables in the same estimation. Columns (4)–(7) use one-way (agent and date levels) and two-way clustered standard errors as well as bootstrapped errors to check the robustness of results.

The results from Table 7 provide support for H1: Multitasking indeed has a positive and statistically significant effect on response delays. The results are consistent across multiple specifications for response delays. The direction and statistical significance remain consistent for both forms of *MT* variables. However, the degree of impact is different for different definitions of multitasking. Merely having a parallel session in progress (even if the agent is idle in that session) increases response delays by three seconds per message.⁹ On the other hand, it takes 34 seconds longer

for an agent to respond to a particular customer's message if this agent was actively multitasking with a new parallel session. This amounts to a slowdown of 60% in agent responses. From the customer's perspective, an interpretation of this finding is that each customer spends four extra minutes¹⁰ in service when the agent is handling one additional customer in parallel.

We are also interested in the marginal effects of increasing multitasking levels on response delays. The coefficients on the multitasking squared terms (for both *MT* measures) are positive and significant, implying a convex relationship. This means that multitasking an additional customer has a greater negative impact on response delays for agents who concurrently process more jobs as opposed to agents who concurrently process fewer jobs. This finding provides support for the second argument of H1.

Before we conclude this section, we turn to the findings on the message order. Our base case for binned order variables is $Bin = 1$ (i.e., the beginning phase of the chat transaction). We find that it takes about 5.6 seconds longer to respond to message in the solution phase than in the base case. On the other hand, compared with the base case again, messages in the closure phase take about 7.3 seconds quicker in response time. This finding highlights the inherent response delay differences between different stages of the conversation.

4.2. Multitasking and Problem Resolution

4.2.1. Model Specification and Estimation. To test the hypothesis that multitasking decreases service agents' ability to successfully resolve problems, we estimate a logit model on the session-level data. The outcome variable is 1 if a customer's problem is resolved at the end of that particular chat session and 0 otherwise. We collect problem resolution information from customer survey responses. Our main independent variable of

Table 7. Response Delay Analysis (Message Level)

Variable	Response delay (<i>RDelay</i>)						
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
MT_1	32.172*** (0.435)		33.874*** (0.441)	33.874*** (1.344)	33.874*** (1.820)	33.874*** (2.151)	33.874*** (0.491)
MT_1^2	9.879*** (0.410)		7.389*** (0.408)	7.389*** (1.032)	7.389*** (0.788)	7.389*** (1.166)	7.389*** (0.425)
$SkillHetMT_1$	-25.889*** (0.806)		-29.955*** (0.827)	-29.955*** (2.498)	-29.955*** (2.955)	-29.955*** (3.633)	-29.955*** (0.791)
MT_2		14.835*** (0.294)	3.121*** (0.317)	3.121*** (0.557)	3.121*** (0.500)	3.121*** (0.630)	3.121*** (0.317)
MT_2^2		8.714*** (0.211)	5.568*** (0.275)	5.568*** (0.462)	5.568*** (0.358)	5.568*** (0.469)	5.568*** (0.301)
$SkillHetMT_2$		0.061 (0.754)	11.639*** (0.801)	11.639*** (1.079)	11.639*** (1.185)	11.639*** (1.290)	11.639*** (0.756)
$\log(Typing)$	13.910*** (0.086)	21.071*** (0.073)	13.817*** (0.086)	13.817*** (0.301)	13.817*** (0.142)	13.817*** (0.303)	13.817*** (0.099)
$\log(CLength)$	-2.171*** (0.088)	-2.748*** (0.084)	-2.177*** (0.088)	-2.177*** (0.198)	-2.177*** (0.167)	-2.177*** (0.231)	-2.177*** (0.090)
$\log(CSlowness)$	1.689*** (0.073)	1.929*** (0.076)	1.681*** (0.072)	1.681*** (0.102)	1.681*** (0.071)	1.681*** (0.091)	1.681*** (0.078)
$Order(Bin = 2)$	5.476*** (0.162)	5.928*** (0.172)	5.645*** (0.162)	5.645*** (0.372)	5.645*** (0.280)	5.645*** (0.413)	5.645*** (0.162)
$Order(Bin = 3)$	-7.749*** (0.201)	-9.113*** (0.211)	-7.304*** (0.201)	-7.304*** (0.482)	-7.304*** (0.480)	-7.304*** (0.621)	-7.304*** (0.213)
<i>Intercept</i>	15.381*** (0.372)	-7.685*** (0.420)	10.845*** (0.451)	10.845*** (1.229)	10.845*** (0.603)	10.845*** (1.212)	10.845*** (0.425)
Session f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Standard errors	Robust	Robust	Robust	Clustered (agent)	Clustered (date)	Clustered (agent and date)	Bootstrapped
No. of obs.	600,376	600,376	600,376	600,376	600,376	600,376	600,376
R-squared	0.415	0.384	0.416	0.416	0.416	0.416	0.416

Notes. *MT* values are centered. Columns (1) and (2) use the two different multitasking/associated variables individually. The remaining models use all the multitasking variables together. Robust standard errors are shown in parentheses. Standard errors are clustered at the agent level in column (4), at the date level in column (5), and at both agent and date levels in column (6). Bootstrapped standard errors are used in column (7) (100-sample bootstrapping). For two-way clustering, we use the code from <http://acct.wharton.upenn.edu/~dtayl/code.htm> (accessed May 15, 2016), used in Gow et al. (2010). f.e., fixed effects.

* $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$.

interest is the average multitasking level across the session (*AvgMT*). We include the average skill heterogeneity across the session (*AvgSkillHetMT*) to account for possible cognitive effects caused by working on dissimilar tasks simultaneously. We also include the standard deviation of multitasking across the session (*sdMT*) to account for the variability of multitasking possibly affecting the problem resolution. The control variables include waiting time in queue (*QueueWait*), session duration (*Duration*), the experience level of the agent with the particular skill for the session (*AExperience*), negative emotions of the customer at the beginning of the session (*NegEmotion*), the brevity of the description of the problem at the beginning of the session (*DescLength*), the total number of words that the agent typed (*TotalTyping*), average customer message length (*AvgCLength*), response delay (*AvgCSlowness*), and the number of days between the contact date and the survey

response date (*SurveyGap*). These values are computed using text processing techniques on chat transcripts and server logs. Similar to the message-level model, we use two measures of *MT* in our estimations. We include agent fixed effects in the model to control for unobserved agent-specific heterogeneity. We also add skill key, product key, and contact date and contact time interval fixed effects to the model to control for other forms of heterogeneity. Skill key and product key fixed effects aim to address inherent differences in problem complexity and product differences. Contact date and contact time interval fixed effects are used to control for time-specific trends and differences.

Another empirical issue to consider in the analysis is self-selection bias. Customer satisfaction surveys are especially prone to self-selection bias because of possible systematic motivational differences between those who choose to respond to surveys and those

Table 8. Problem Resolution Analysis (Session Level)

Variable	Resolution (Survey)				
	(1)	(2)	(3)	(4)	(5)
<i>AvgMT</i> ₁	−0.747*** (0.120)		−0.724*** (0.148)	−0.724*** (0.147)	−0.724*** (0.179)
<i>AvgSkillHetMT</i> ₁	0.341 (0.277)		0.285 (0.373)	0.285 (0.367)	0.285 (0.409)
<i>sdMT</i> ₁	0.098 (0.105)		0.138 (0.110)	0.138 (0.109)	0.138 (0.120)
<i>AvgMT</i> ₂		−0.244*** (0.064)	−0.053 (0.079)	−0.053 (0.080)	−0.053 (0.054)
<i>AvgSkillHetMT</i> ₂		0.176 (0.139)	0.036 (0.189)	0.036 (0.199)	0.036 (0.200)
<i>sdMT</i> ₂		−0.100 (0.089)	−0.120 (0.091)	−0.120 (0.088)	−0.120 (0.074)
log(<i>QueueWait</i>)	−0.028* (0.015)	−0.029* (0.016)	−0.028* (0.016)	−0.028* (0.016)	−0.028 (0.018)
log(<i>Duration</i>)	−0.992*** (0.320)	−0.947*** (0.320)	−0.988*** (0.320)	−0.988*** (0.330)	−0.988** (0.396)
log(<i>AExperience</i>)	−0.052 (0.034)	−0.043 (0.034)	−0.050 (0.034)	−0.050 (0.033)	−0.050 (0.033)
log(<i>NegEmotion</i>)	−0.317*** (0.082)	−0.318*** (0.082)	−0.318*** (0.082)	−0.318*** (0.082)	−0.318*** (0.09)
log(<i>DescLength</i>)	−0.394 (0.275)	−0.425 (0.274)	−0.398 (0.275)	−0.398 (0.280)	−0.398 (0.322)
log(<i>TotalTyping</i>)	0.468*** (0.041)	0.383*** (0.036)	0.471*** (0.042)	0.471*** (0.048)	0.471*** (0.042)
log(<i>AvgCLength</i>)	−0.364*** (0.048)	−0.350*** (0.048)	−0.365*** (0.048)	−0.365*** (0.048)	−0.365*** (0.039)
log(<i>AvgCSlowness</i>)	0.105** (0.049)	0.075 (0.049)	0.108** (0.049)	0.108** (0.055)	0.108*** (0.035)
log(<i>SurveyGap</i>)	0.159*** (0.048)	0.163*** (0.048)	0.160*** (0.048)	0.160*** (0.051)	0.160*** (0.052)
λ	−3.730* (2.053)	−3.867* (2.048)	−3.751* (2.052)	−3.751* (2.101)	−3.751 (2.457)
Agent f.e.	Yes	Yes	Yes	Yes	Yes
Skill key f.e.	Yes	Yes	Yes	Yes	Yes
Product key f.e.	Yes	Yes	Yes	Yes	Yes
Contact date f.e.	Yes	Yes	Yes	Yes	Yes
Contact time f.e.	Yes	Yes	Yes	Yes	Yes
Standard errors	Robust	Robust	Robust	Clustered (agent)	Clustered (date)
No. of obs.	13,507	13,507	13,507	13,507	13,507

Notes. Columns (1) and (2) use the two different multitasking/associated variables individually. The remaining models use all the multitasking variables together. Robust standard errors are shown in parentheses. Standard errors are clustered at the agent level in column (4) and at the date level in column (5). f.e., fixed effects.

* $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$.

who do not. As stated earlier, 16% of customers in our entire data set completed the customer satisfaction surveys. While this response rate is on par with the industry average of 10%–15% for external surveys,¹¹ it could still indicate a selectivity bias. We address this potential bias using the two-step Heckman correction method (Heckman 1979). First, we estimate the likelihood of a customer responding to the survey using probit regression. The probit model uses session duration, queue wait time, agent experience, prior negative

emotion, and description length as predictor variables and constructs the inverse Mills ratio (λ); λ accounts for the fact that self-reported resolution information is only observed when customers respond to the surveys. Then, we include λ in the estimation of our main model.

4.2.2. Results. Table 8 reports the results for multitasking versus problem resolution analysis for various specifications. The first two columns provide the estimations for MT_1 and MT_2 individually, whereas col-

umn (3) uses both multitasking variables in the same estimation. Columns (4) and (5) provide the results with clustered standard errors. We note that because of the inclusion of agent fixed effects and some of the agents' indifference in terms of the dependent variable, 618 observations were removed from the logit model estimation.

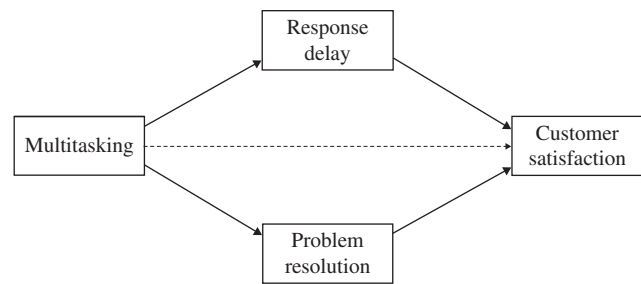
Three primary results emerge based on the estimation of this model. First, among all the specifications, the coefficient of $AvgMT_1$ is statistically significant and negative, indicating that working on a parallel session is negatively associated with problem resolution, providing support for H2. The coefficient of $AvgMT_1$ (columns (3)–(5)) indicates that if an agent actively communicates with this additional customer throughout the existing session, then the estimated odds of resolving the existing problem decreases by approximately 51%. This result is robust to unobserved heterogeneity and self-selection bias, and it is consistent across different standard error structures. Second, we do not observe the same degree of impact for the $AvgMT_2$, since its coefficient (when estimated together with $AvgMT_1$) is substantially smaller and insignificant. This finding supports the argument that there is a difference between two different definitions of multitasking: (1) being assigned an additional task and (2) actively working on an additional task. The problem resolution capability of agents diminishes when the agents actively switch between different tasks throughout their work periods. Considering the fact that multitasking policies inherently promote task-switching behavior,¹² this situation appears to create an interesting conundrum for agents—whether to follow the policy for better utilization of their time or ignore the policy in return for more successfully resolving their assigned tasks. Finally, our finding about the decline in problem resolution as a result of multitasking not only raises concerns about the quality of service under multitasking regimes but also challenges the argument that favors multitasking for its presumed productivity benefits. For many service organizations with repeat customers, such presumed productivity gains may cease to materialize in the long run, because customers whose problems have not been resolved may return later (Mehrotra et al. 2012), either to the same customer service channel or to a different one (e.g., phone or in-person).

4.3. Multitasking and Customer Satisfaction

4.3.1. Model Specification and Estimation. To study the relationship between multitasking and customer satisfaction, we conduct a path analysis for the model shown in Figure 3.

This analysis aims to identify the possible direct effect of multitasking on customer satisfaction, as well as possible indirect effects via response delay and problem resolution paths. For this purpose, we fit three OLS

Figure 3. Indirect Relationship Between Multitasking and Customer Satisfaction



regression models to the session-level data and estimate the path coefficients using the results from the models. In these models, our main variable of interest is the average multitasking level across the session ($AvgMT$), represented in two forms. Our dependent variables are average response delay within the session ($AvgRDelay$), problem resolution ($Resolution$), and customer satisfaction ($Satisfaction$). We measure *Satisfaction* using two items. The first item is a survey question that asks whether the expectations of the customer were fulfilled at the end of the chat session. The second item is a survey question that asks, based on the last chat transaction, the likelihood that the customer will promote the service to friends and relatives.

The three regression models in the path analysis are all at the session level; therefore we use the same procedure discussed in Section 4.2.1 to address econometric concerns. Specifically, we include control variables and agent, skill key, product key, contact date, and contact time interval fixed effects in the model to control for different forms of heterogeneity. In addition, we address the self-selection bias using the Heckman correction method.

4.3.2. Results. As discussed earlier, we use two measures for customer satisfaction, which are expectation fulfillment and net promoter scores. Tables 9 and 10 present the results for analyses with each measure, respectively.

Panel A of these tables shows the regression results for the direct paths from multitasking to response delay and problem resolution, as well as the direct effect on customer satisfaction. $AvgMT_1$ is significant in the first two regressions, and the direction of the coefficient is consistent with our analyses at the message level (see Section 4.1) and the logit model (see Section 4.2). On the other hand, $AvgMT_1$ is not a significant predictor of customer satisfaction when response delay and problem resolution are controlled for. The results are consistent whether $AvgMT_1$ is used independently or jointly with $AvgMT_2$. This finding points to the indirect effects of multitasking on customer satisfaction, which are calculated in panel B of the tables. In these panels, we break down our estimations for

Table 9. Effects of Multitasking on Expectation Fulfillment via Average Response Delay and Resolution

Panel A: Regression analyses									
Variable	Avg. response delay (<i>AvgRDelay</i>)			Resolution (<i>Resolution</i>)			Expectation fulfillment (<i>Satisfaction₁</i>)		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
<i>AvgMT₁</i>	63.807*** (1.390)		49.788*** (1.748)	-0.150*** (0.024)		-0.053* (0.031)	-0.064 (0.053)		0.007 (0.066)
<i>AvgSkillHetMT₁</i>	-30.365*** (3.169)		-51.799*** (4.127)	0.064 (0.054)		0.072 (0.073)	-0.060 (0.112)		-0.119 (0.152)
<i>sdMT₁</i>	-13.149*** (1.187)		-23.093*** (1.176)	0.020 (0.020)		0.111*** (0.021)	0.054 (0.042)		0.155*** (0.044)
<i>AvgMT₂</i>		16.031*** (0.748)	-1.408* (0.754)		-0.047*** (0.012)	-0.018 (0.015)		-0.030 (0.026)	-0.051 (0.032)
<i>AvgSkillHetMT₂</i>		-4.475*** (1.725)	17.769*** (2.059)		0.034 (0.027)	-0.008 (0.036)		0.013 (0.056)	0.041 (0.076)
<i>sdMT₂</i>		2.357** (1.093)	1.953** (0.988)		-0.020 (0.017)	-0.002 (0.017)		-0.081** (0.035)	-0.078** (0.036)
log(<i>QueueWait</i>)	-0.025 (0.151)	-0.035 (0.168)	-0.130 (0.149)	-0.005* (0.003)	-0.006* (0.003)	-0.005* (0.003)	-0.020*** (0.006)	-0.021*** (0.006)	-0.021*** (0.006)
log(<i>Duration</i>)	10.782*** (0.458)	5.795*** (0.492)	10.817*** (0.365)	-0.208*** (0.059)	-0.196*** (0.059)	-0.167*** (0.059)	-0.091 (0.124)	-0.088 (0.124)	-0.035 (0.124)
log(<i>AExperience</i>)	-1.993*** (0.375)	-2.321*** (0.409)	-1.773*** (0.362)	-0.010 (0.006)	-0.008 (0.006)	-0.009 (0.006)	-0.021 (0.013)	-0.020 (0.013)	-0.019 (0.013)
log(<i>NegEmotion</i>)	-0.307 (0.528)	-0.725 (0.577)	-0.228 (0.510)	-0.067*** (0.015)	-0.066*** (0.015)	-0.070*** (0.015)	-0.003 (0.032)	-0.004 (0.032)	-0.007 (0.032)
log(<i>DescLength</i>)	2.357*** (0.468)	3.383*** (0.511)	2.403*** (0.451)	-0.088* (0.051)	-0.091* (0.051)	-0.106** (0.051)	-0.079 (0.106)	-0.088 (0.106)	-0.105 (0.106)
log(<i>TotalTyping</i>)	3.134*** (0.448)	3.687*** (0.438)	18.797*** (0.596)	0.094*** (0.008)	0.077*** (0.007)	-0.044*** (0.011)	0.106*** (0.016)	0.108*** (0.014)	0.021 (0.023)
log(<i>AvgCLength</i>)	1.959*** (0.524)	0.725 (0.571)	1.420*** (0.505)	-0.072*** (0.009)	-0.069*** (0.009)	-0.065*** (0.009)	-0.096*** (0.019)	-0.094*** (0.019)	-0.089*** (0.019)
log(<i>AvgCSlowness</i>)	2.335*** (0.553)	4.872*** (0.600)	1.814*** (0.526)	0.023** (0.009)	0.016* (0.009)	0.004 (0.009)	-0.001 (0.020)	-0.001 (0.019)	-0.025 (0.019)
log(<i>SurveyGap</i>)				0.030*** (0.009)	0.030*** (0.009)	0.031*** (0.009)	0.007 (0.018)	0.007 (0.018)	0.008 (0.018)
λ				-0.810** (0.378)	-0.817** (0.378)	-0.865** (0.380)	-0.718 (0.794)	-0.776 (0.793)	-0.840 (0.795)
<i>AvgRDelay</i>							-0.002*** (0.000)	-0.002*** (0.000)	-0.002*** (0.000)
<i>Resolution</i>							2.003*** (0.018)	2.003*** (0.018)	2.014*** (0.018)
Agent f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Skill key f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Product key f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Contact date f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Contact time f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
No. of obs.	14,021	14,021	14,021	14,021	14,021	14,021	14,021	14,021	14,021

the effect of multitasking on customer satisfaction into (1) indirect effect via the response delay path, (2) indirect effect via the resolution path, (3) direct effect, and (4) total effect. We compute the indirect effects on the outcome variables by taking the product of coefficients from the corresponding columns of the tables. We find that the total effects of *AvgMT₁* on expectation fulfillment and net promoter scores are significant and negative. Yet these total effects are primarily driven by indi-

rect effects through response delays and problem resolution paths. This finding provides support for H3A and H3B for the *MT₁* variable. In other words, while customers cannot observe the efforts of agents directly, they still feel the side effects of multitasking. When a service agent actively engages in communication with multiple customers, the service he or she provides is less likely to fulfill these customers' expectations because of increasing in-service delays and decreas-

Downloaded from informs.org by [132.239.1.231] on 28 June 2017, at 02:18. For personal use only, all rights reserved.

Table 9. (Continued)

Path	Panel B: Path analysis			
	AvgMT ₁	AvgMT ₂	Both MT variables together	
			AvgMT ₁	AvgMT ₂
Indirect effect via <i>Response Delay</i>	-0.097*** (0.019)	-0.027*** (0.005)	-0.083*** (0.016)	0.002 (0.001)
Indirect effect via <i>Resolution</i>	-0.300*** (0.047)	-0.095*** (0.024)	-0.107* (0.062)	-0.037 (0.031)
Total indirect eff.	-0.396*** (0.051)	-0.122*** (0.025)	-0.190*** (0.064)	-0.035 (0.031)
Direct effect	-0.064 (0.053)	-0.030 (0.026)	0.007 (0.066)	-0.051 (0.032)
Total effect	-0.460*** (0.068)	-0.152*** (0.035)	-0.183** (0.089)	-0.086* (0.044)

Note. f.e., fixed effects.
 * $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$.

ing problem resolution. Furthermore, these customers would be less likely to promote the service to their friends and relatives.

5. Robustness and Additional Analyses

5.1. Robustness

5.1.1. Matching Analysis at the Session Level. As discussed earlier, a useful characteristic of our context is the exogenous nature of multitasking assignments. In this contact center setting, customers are served on a first-come-first-served basis, and multitasking decisions are made by a contact distribution system based on the length of the queues. First, to validate that this admissions policy was indeed implemented in our data set, we search the data for records that were exceptions to the policy guidelines. In particular, we search for cases where a customer was *jumped over*—that is, where one customer was admitted later than another customer with the same problem type who came after him or her. We find an extremely small number of cases (less than 0.1% of all records) that contradict the policy. This finding supports the argument regarding the exogenous nature of customer admissions. Next, to further test the robustness of the session-level results obtained from our setup, we employ a propensity score matching (PSM) approach. This approach helps overcome the potential concerns regarding the queue lengths affected by individual agent characteristics such as speed or capability.

To conduct the PSM analysis, we first estimate a logit model to identify whether a customer would be assigned to a multitasking agent during his or her session (i.e., treated or not). Since our multitasking variables are continuous by definition, we use different thresholds (between 1.3 and 1.8, with 0.1 increments) to define what constitutes multitasking treatment in a session. For the treatment model, we use *queue wait time*,

agent experience, *agent key*, *problem skill key*, *contact date*, and *contact time interval* as predictor variables. For each observation that was actually treated, we find a non-treated observation that is highly similar in terms of its treatment model score, using nearest neighbor matching. Tables 11 and 12 provide the results from propensity score matching for each of the multitasking variables (MT_1 and MT_2). From Table 11, we can see that both multitasking definitions have a significant and positive effect on average response delays, regardless of the threshold value. This result is consistent with the regression estimations in panel A in Tables 9 and 10 (the first three columns), and it also complements the findings in the message-level analysis in Table 7. From Table 12, we can see that only $AvgMT_1$ has a consistent negative and significant (for certain thresholds) effect on problem resolution. Again, this finding is consistent with our logit model in Table 8 and the regression models in Tables 9 and 10 (columns (4)–(6)).

One common concern for propensity score matching is the existence of unobservables that may be driving the difference between the treatment (i.e., multitasking) and comparison (i.e., nonmultitasking) groups. To investigate how sensitive our matching estimates are to the possible existence of an unobserved confounding variable, we conduct sensitivity analysis (Rosenbaum 2002). We find that to attribute the effect to an unobserved confounder, this unobserved confounder would need to produce, on average, a threefold increase in the odds of multitasking. This is much higher than typical thresholds used in prior literature (e.g., Sun and Zhu 2013, p. 2327), suggesting that unobservable information is unlikely to drive our findings.

5.1.2. Matching Analysis at the Message Level. A different concern in our analyses is that at the message level, a service agent’s decision on whom to respond to (in the case of handling multiple simultaneous ses-

Table 10. Effects of Multitasking on Net Promoter Score via Average Response Delay and Resolution

Panel A: Regression analyses									
Variable	Avg. response delay (<i>AvgRDelay</i>)			Resolution (<i>Resolution</i>)			Promotion (<i>Satisfaction₂</i>)		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
<i>AvgMT₁</i>	63.679*** (1.390)		70.291*** (1.708)	-0.151*** (0.024)		-0.147*** (0.029)	-0.077* (0.040)		-0.072 (0.048)
<i>AvgSkillHetMT₁</i>	-30.706*** (3.171)		-52.979*** (4.267)	0.069 (0.054)		0.056 (0.072)	-0.100 (0.085)		0.069 (0.115)
<i>sdMT₁</i>	-12.982*** (1.186)		-15.439*** (1.249)	0.022 (0.020)		0.029 (0.021)	0.063* (0.032)		0.078** (0.034)
<i>AvgMT₂</i>		16.040*** (0.784)	-1.969** (0.890)		-0.048*** (0.012)	-0.009 (0.015)		-0.015 (0.019)	-0.013 (0.024)
<i>AvgSkillHetMT₂</i>		-4.480** (1.724)	16.866*** (2.129)		0.035 (0.027)	-0.008 (0.036)		0.047 (0.042)	0.021 (0.057)
<i>sdMT₂</i>		2.403** (1.092)	4.465*** (1.021)		-0.019 (0.017)	-0.023 (0.017)		-0.048** (0.027)	-0.062** (0.027)
<i>log(QueueWait)</i>	-0.027 (0.151)	-0.026 (0.168)	-0.096 (0.153)	-0.005* (0.003)	-0.006* (0.003)	-0.006* (0.003)	-0.012** (0.005)	-0.012*** (0.005)	-0.012*** (0.005)
<i>log(Duration)</i>	10.789*** (0.457)	5.834*** (0.492)	10.459*** (0.461)	-0.204*** (0.059)	-0.192*** (0.059)	-0.203*** (0.059)	-0.132 (0.094)	-0.130 (0.094)	-0.129 (0.094)
<i>log(AExperience)</i>	-1.984*** (0.374)	-2.323*** (0.409)	-1.864*** (0.374)	-0.010 (0.006)	-0.009 (0.006)	-0.010 (0.006)	-0.025** (0.010)	-0.025** (0.010)	-0.025** (0.010)
<i>log(NegEmotion)</i>	-0.359 (0.528)	-0.736 (0.576)	-0.393 (0.527)	-0.066*** (0.015)	-0.065*** (0.015)	-0.066*** (0.015)	-0.064*** (0.024)	-0.064*** (0.024)	-0.064*** (0.024)
<i>log(DescLength)</i>	2.407*** (0.468)	3.437*** (0.510)	2.485*** (0.467)	-0.085* (0.051)	-0.088* (0.051)	-0.085* (0.051)	-0.054 (0.080)	-0.058 (0.080)	-0.055 (0.080)
<i>log(TotalTyping)</i>	3.190*** (0.448)	3.598*** (0.438)	3.071*** (0.451)	0.094*** (0.008)	0.077*** (0.007)	-0.094*** (0.008)	0.024** (0.012)	0.030*** (0.011)	0.026** (0.012)
<i>log(AvgCLength)</i>	1.892*** (0.524)	0.659 (0.571)	1.420*** (0.505)	-0.072*** (0.009)	-0.069*** (0.009)	-0.072*** (0.009)	-0.077*** (0.014)	-0.076*** (0.014)	-0.077*** (0.014)
<i>log(AvgCSlowness)</i>	2.335*** (0.553)	4.841*** (0.600)	1.880*** (0.522)	0.022** (0.009)	0.016* (0.009)	0.022** (0.009)	-0.016 (0.015)	-0.017 (0.015)	-0.017 (0.015)
<i>log(SurveyGap)</i>				0.031*** (0.009)	0.031*** (0.009)	0.031*** (0.009)	0.004 (0.014)	0.005 (0.014)	0.005 (0.014)
λ				-0.785** (0.378)	-0.793** (0.378)	-0.841** (0.339)	-0.730 (0.601)	-0.760 (0.601)	-0.738 (0.601)
<i>AvgRDelay</i>							-0.002*** (0.000)	-0.001*** (0.000)	-0.001*** (0.000)
<i>Resolution</i>							0.998*** (0.013)	0.998*** (0.013)	0.998*** (0.013)
Agent f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Skill key f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Product key f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Contact date f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Contact time f.e.	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
No. of obs.	14,069	14,069	14,069	14,069	14,069	14,069	14,069	14,069	14,069

sions) is based on the agent's choice, and therefore MT_1 may be endogenous. First, we note that this is an unlikely event in our context; the response preferences are mainly driven by an external factor that is exogenous to agents. The reason for this exogeneity lies within the nature of multitasking policies. The main premise of multitasking is to improve agent productivity by switching to a different customer whenever there is downtime with the other customer(s). Productivity

benefits is the major reason that multitasking policies are implemented by firms, and therefore, multitasking behavior is primarily driven by the downtimes caused by customers rather than by arbitrary agent decisions. In other words, service agents are more likely to follow the policy for its intended purpose, rather than making arbitrary decisions on whom to serve by themselves.

Nevertheless, to further alleviate this concern, we perform matching analysis at the message level. For

Table 10. (Continued)

Path	Panel B: Path analysis			
	<i>AvgMT</i> ₁	<i>AvgMT</i> ₂	Both MTs together	
			<i>AvgMT</i> ₁	<i>AvgMT</i> ₂
Indirect effect via <i>Response Delay</i>	−0.024* (0.014)	−0.008** (0.003)	−0.026* (0.016)	0.001 (0.001)
Indirect effect via <i>Resolution</i>	−0.151*** (0.024)	−0.048*** (0.012)	−0.147*** (0.029)	−0.009 (0.015)
Total indirect eff.	−0.175*** (0.028)	−0.056*** (0.013)	−0.173*** (0.033)	−0.008 (0.015)
Direct effect	−0.077* (0.040)	−0.015 (0.019)	−0.072 (0.050)	−0.013 (0.024)
Total effect	−0.252*** (0.044)	−0.071*** (0.023)	−0.245*** (0.054)	−0.021 (0.028)

Note. f.e., fixed effects.
 p* < 0.1; *p* < 0.05; ****p* < 0.01.

this analysis, we apply one-to-one matching and match each message in which the agent had performed multitasking (i.e., sent at least one message to a nonfocal customer) to another message in a different session in which the agent had strictly worked with a single customer. For the matching procedure, we search for similar characteristics including the same agent, same skill,

same product type, same order bin, and similar agent and customer message lengths.¹³ If there are multiple matches for a particular message, we randomly pick one match and discard the rest. Next, we obtain the response delays for each message in the two groups (multitasked versus non-multitasked) and run a paired *t*-test to see whether the response delay means across

Table 11. Propensity Score Matching Results for Average Response Delay at the Session Level

Threshold to be considered as treated	Multitasking variable	ATET	95% conf. interval (low)	95% conf. interval (high)
1.3	<i>AvgMT</i> ₁	21.871*** (0.330)	21.223	22.518
1.4	<i>AvgMT</i> ₁	24.062*** (0.350)	23.375	24.749
1.5	<i>AvgMT</i> ₁	27.985*** (0.427)	27.149	28.821
1.6	<i>AvgMT</i> ₁	29.788*** (0.451)	28.904	30.672
1.7	<i>AvgMT</i> ₁	33.222*** (0.582)	32.080	34.363
1.8	<i>AvgMT</i> ₁	36.486*** (0.613)	35.192	38.116
1.3	<i>AvgMT</i> ₂	14.774*** (0.795)	13.215	16.333
1.4	<i>AvgMT</i> ₂	13.214*** (0.644)	11.952	14.475
1.5	<i>AvgMT</i> ₂	11.468*** (0.671)	10.152	12.783
1.6	<i>AvgMT</i> ₂	11.099*** (0.662)	9.801	12.397
1.7	<i>AvgMT</i> ₂	9.996*** (0.518)	8.980	11.011
1.8	<i>AvgMT</i> ₂	8.564*** (0.474)	7.635	9.493

p* < 0.1; *p* < 0.05; ****p* < 0.01.

Table 12. Propensity Score Matching Results for Problem Resolution at the Session Level

Threshold to be considered as treated	Multitasking variable	ATET	95% conf. interval (low)	95% conf. interval (high)
1.3	<i>AvgMT</i> ₁	−0.013 (0.009)	−0.032	0.005
1.4	<i>AvgMT</i> ₁	−0.016 (0.010)	−0.036	0.004
1.5	<i>AvgMT</i> ₁	−0.044*** (0.013)	−0.069	−0.020
1.6	<i>AvgMT</i> ₁	−0.058*** (0.009)	−0.077	−0.039
1.7	<i>AvgMT</i> ₁	−0.085*** (0.003)	−0.091	−0.078
1.8	<i>AvgMT</i> ₁	−0.092*** (0.013)	−0.117	−0.067
1.3	<i>AvgMT</i> ₂	0.016 (0.016)	−0.015	0.047
1.4	<i>AvgMT</i> ₂	0.007 (0.012)	−0.016	0.030
1.5	<i>AvgMT</i> ₂	−0.014 (0.018)	−0.049	0.021
1.6	<i>AvgMT</i> ₂	−0.002 (0.014)	−0.030	0.026
1.7	<i>AvgMT</i> ₂	0.002 (0.013)	−0.024	0.028
1.8	<i>AvgMT</i> ₂	0.004 (0.011)	−0.018	0.026

p* < 0.1; *p* < 0.05; ****p* < 0.01.

Downloaded from informs.org by [132.239.1.231] on 28 June 2017, at 02:18. For personal use only, all rights reserved.

Table 13. Paired *t*-Test Results After Matching at the Message Level

Variable	Mean	N	Std. err.	95% conf. int. (low)	95% conf. int. (high)	<i>t</i>	Sig.
<i>MT</i>	69.026	34,382	0.283	68.471	69.582		
<i>No MT</i>	52.876	34,382	0.274	52.339	53.412		
<i>Difference</i>	16.150	34,382	0.371	15.423	16.878	43.503	0.000

the two groups differ. Table 13 provides the results from this analysis.

We find that paired *t*-test results are fully consistent with our earlier findings—that is, there is a statistically significant difference ($t = 43.503, p = 0.000$) in response delays between multitasked ($\mu = 69.026$) and non-multitasked ($\mu = 52.876$) groups.

We run two further robustness tests and find that the results are still robust: (1) rather than using a random pick, we average the response delay from all matches, and (2) rather than using just a *t*-test, we use a multivariate regression that incorporates agent, skill key, product key, contact date, contact time, and order bin fixed effects. For brevity, we do not report those results here.

5.1.3. Nonlogged Control Variables. Our main models employ natural log transformations of control variables to reduce the skewness of data. To ensure that this process does not bias the findings, we estimate the response delay and problem resolution models using the original scales of these variables. To conserve space, we provide the full results in Appendix A, available online. The results are highly consistent with those provided in Section 4.

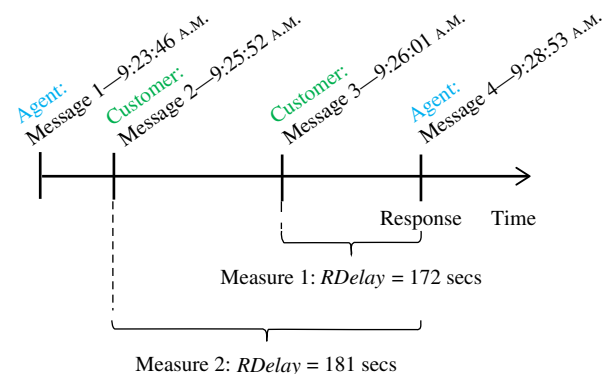
5.1.4. Alternative Response Delay Measure for Back-to-Back Messages. An empirical challenge in measuring response delay is the handling of back-to-back messages, since some agent responses were preceded by multiple consecutive customer messages. We previously computed response delays using the final message in customers' messages. This idea is illustrated as point 1 in Figure 4.

This approach is expected to lead to a more conservative estimate of the impact of multitasking; this is why we use it in the main models. To ensure the robustness of our findings, we measure response delays using the first message in the message sequence (illustrated as point 2 in Figure 4) as a robustness test. We also recompute the multitasking measurement values accordingly. To conserve space, we provide the full results in Appendix B, available online. These results are consistent with the original analysis and support the finding that multitasking has a positive and statistically significant effect on response delays.

5.2. Acting on Our Findings: A Segmentation Model for a Customer-Specific Routing Policy

Our analyses so far focus on the cause-effect relationships between multitasking and the dependent variables of interest. These findings can, first and foremost, inform managerial decisions in terms of whether to multitask agents and to what extent agents should be multitasking. Yet given budgetary constraints, companies may not always be able to avoid multitasking, so the natural question is, *can we selectively let agents multitask in some situations and not multitask in others?* Can our analyses be extended to provide design guidelines for the customer assignment systems to preserve customer satisfaction while increasing the throughput? Because a full treatment of this question is beyond the scope of this paper, here we focus on one angle: customer heterogeneity in terms of the customer's reactions to—or, rather, tolerance of—agent multitasking. If we can use some information about customers to classify them into sensitive and nonsensitive clusters, then the system can assign less sensitive customers to multitasking agents. Doing so would help us leverage the productivity benefits of multitasking (in terms of more efficient use of agents' times) while reducing its negative impact on customer satisfaction.

We propose a finite mixture approach that takes such possible customer heterogeneity into account. Finite mixture approaches have been commonly used in the marketing domain for customer segmentation (Wedel and Kamakura 2000) and are being increasingly adopted in information systems research (Bapna et al. 2011). Following common practice, we assume

Figure 4. (Color online) Two Approaches to Handle Back-to-Back Messages

that customers in our data set can be classified into S discrete segments (i.e., classes). Customers within the same segment are assumed to be homogeneous in terms of their satisfaction behavior based on session-specific variables. However, sensitivities toward these variables may differ across different segments. We define the probability structure for this model as follows:

Let $P_i(k | s)$ denote the probability that at the end of chat session i ; the customer will have the satisfaction score of k conditional on the customer belonging to segment s . Using a multinomial logistic regression model, this probability can be represented as

$$P_i(k | s) = \frac{\exp(\mathbf{B}_s \mathbf{X}_{ik})}{\sum_{l=1}^K \exp(\mathbf{B}_s \mathbf{X}_{il})}, \quad (1)$$

where \mathbf{X}_{ik} is the column vector of explanatory variables that can help explain the satisfaction behavior and \mathbf{B}_s is the row vector of coefficients associated with \mathbf{X}_{ik} . Furthermore, let P_{is} denote the probability that a customer will belong to segment s ; P_{is} depends on a vector of customer-specific (i.e., concomitant) variables \mathbf{D}_i , which can be represented as

$$P_{is} = \frac{\exp(\gamma_s \mathbf{D}_i)}{\sum_{s=1}^S \exp(\gamma_s \mathbf{D}_i)}, \quad (2)$$

where γ_s ($s = 1, 2, \dots, S$) is the vector of coefficients that represents the effects of customer-specific variables on the probability of segment membership. Combining these two equations, the probability of an arbitrary session i having satisfaction level k is

$$\begin{aligned} P_i(k) &= \sum_{s=1}^S P_{is} \times P_i(k | s) \\ &= \sum_{s=1}^S \frac{\exp(\gamma_s \mathbf{D}_i)}{\sum_{s=1}^S \exp(\gamma_s \mathbf{D}_i)} \times \frac{\exp(\mathbf{B}_s \mathbf{X}_{ik})}{\sum_{l=1}^K \exp(\mathbf{B}_s \mathbf{X}_{il})}. \end{aligned} \quad (3)$$

Estimation of the unknown coefficients \mathbf{B}_s and γ_s can be conducted using the maximum likelihood method. Once these estimates are obtained, we can assign each customer to a segment by finding the segment for which the customer has the highest posterior probability.

To demonstrate the application of this model, we apply it within the context of multitasking—customer satisfaction analysis. We assume that there are two segments of customers who have different levels of sensitivity toward in-service delay.¹⁴ At the time of their arrival to the system, we have limited information about the customers. Specifically, we can use the following variables to classify the customers into one of the two segments: gender (*Gender*), problem domain (*Skill domain*), product type (*Product key*), and hourly time interval of the contact (*Contact time*). The values of

these variables are recurrent and easy to obtain, making them feasible for use in the segmentation process. A customer satisfaction estimation can then be conducted for each segment.

Table 14 presents the finite mixture estimation results with $Satisfaction_1$ as the dependent variable and $AvgMT_1$ as one of the predictor variables. We provide results for both one-class and two-class models to compare different fits. Four different model choice criteria (Akaike information criterion (AIC), Bayesian information criterion (BIC), consistent Akaike information criterion (CAIC), and R -squared) collectively suggest that a two-class model is a better fit than a one-class model for this data set. For the two-class model, the p -value for the Wald statistics of $AvgMT_1$ is less than 0.001, indicating that the overall effect of multitasking is highly significant. The coefficient of this effect is negative and statistically significant for both segments, which is consistent with our previous findings. On the other hand, at the individual-segment level, we find the coefficient of $AvgMT_1$ to be approximately two times larger than that of segment 2, indicating a much stronger negative impact of multitasking for customers that are members of segment 1. A separate Wald statistic shows that the between-segment difference in terms of $AvgMT_1$ indeed exists and is significant ($p < 0.1$). In terms of identifying the segments, we find the coefficients of gender (*Gender*), problem domain (*SkillDomain*), and product type (*ProductID*) to be significant, meaning that these variables can be used for customer segment prediction.

In summary, these results suggest that multitasking has a negative overall effect on customer satisfaction; however, not all customers are affected to the same degree. By using information from the time of customer arrival to the system, contact centers can segment customers into two groups with different multitasking sensitivity levels. This segmentation scheme can be further integrated into a routing policy that can make decisions on which customers should be multitasked first when multiple customers are waiting for service.

6. Conclusion

In this research, we studied the negative consequences of multitasking in service organizations. We used a diverse set of operational data (including contact transcripts, server logs, and customer survey responses) from an online chat contact center to investigate the effects of multitasking on problem resolution rates, response delays, and customer satisfaction. Our results are based on two complementary metrics of multitasking: (1) whether an agent was nominally assigned an additional customer and (2) whether he or she was actively working with that additional customer. We also expanded our analyses to account for

Table 14. Two-Segment Finite Mixture Model

	Expectation fulfillment (<i>Satisfaction₁</i>)			
	(1)		(2)	
	One-class model		Two-class model	
AIC (LL)	43,524.57		42,554.21	
BIC (LL)	44,000.16		43,796.87	
CAIC (LL)	44,063.16		43,961.87	
R-squared	0.076		0.233	
	Segment	Level	Segments	
	1		1	2
Predictors				
<i>AvgMT₁</i>	−0.268*** (0.039)		−0.439*** (0.082)	−0.215*** (0.071)
<i>AvgSkillHetMT₁</i>	0.029 (0.087)		0.107 (0.184)	0.060 (0.163)
<i>sdMT₁</i>	0.049 (0.033)		0.037 (0.073)	0.054 (0.057)
log(<i>QueueWait</i>)	−0.013*** (0.004)		−0.017* (0.008)	−0.005 (0.007)
log(<i>Duration</i>)	−0.129*** (0.012)		−0.151*** (0.036)	−0.141*** (0.024)
log(<i>AExperience</i>)	0.025*** (0.006)		0.030** (0.012)	0.035*** (0.011)
log(<i>NegEmotion</i>)	−0.024* (0.015)		0.066 (0.040)	−0.088*** (0.027)
log(<i>DescLength</i>)	0.027** (0.013)		0.153*** (0.033)	−0.061** (0.029)
log(<i>TotalTyping</i>)	0.230*** (0.012)		0.349*** (0.035)	0.181*** (0.023)
log(<i>AvgCLength</i>)	−0.153*** (0.015)		−0.280*** (0.035)	−0.065** (0.030)
log(<i>AvgCSlowness</i>)	0.015 (0.015)		0.013 (0.035)	0.039 (0.029)
log(<i>SurveyGap</i>)	0.042*** (0.014)		0.034 (0.029)	0.059** (0.028)
Skill key f.e.	Yes		Yes	Yes
Contact date f.e.	Yes		Yes	Yes
Model for segments				
<i>Gender</i>	N/A		0.132*** (0.039)	−0.132*** (0.039)
<i>Skill domain</i>	N/A	Getting started	−0.274** (0.109)	0.274* (0.109)
	N/A	Shop/buy	0.851*** (0.163)	−0.851*** (0.163)
	N/A	Working on file	−0.576*** (0.136)	0.576*** (0.136)
<i>Product key</i>	N/A	19 levels	Yes	Yes
<i>Contact time</i>	N/A	18 levels	Yes	Yes
Relative segment size	1.00		0.51	0.49

Note. f.e., fixed effects; LL, log likelihood.

* $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$.

customer heterogeneity and identified tailored routing strategies based on customer segmentation. Prior research has mainly followed analytical and experimental approaches to examine the impact of multitasking on productivity (i.e., throughput) related measures. To our knowledge, this is one of the very few studies that investigates the multitasking phenomenon

from a service quality perspective using observational data.

Our findings have several implications for research and practice. Previous research has shown that under certain conditions it is possible to increase productivity via multitasking (Aral et al. 2012, O'Leary et al. 2011). We found that increased productivity may come at the

expense of diminished quality of service and customer satisfaction. A managerial implication of this finding is that focusing solely on the productivity measures may not be the best benchmark for firms. Service organizations need to take a holistic view of their service goals, from both quality-of-service and productivity perspectives.

An interesting aspect of our context is that in live-chat communication, customers cannot directly observe the task-processing behavior of agents, and they may be unaware of the on-going multitasking activity. Yet, inadvertently, they still feel the side effects of multitasking. First, customers experience longer in-service delays in the presence of multitasking. This result holds even if the agent focuses on only one customer during service. However, the effect is much greater when the agent routinely switches between multiple tasks or when the number of multitasked jobs is increased (i.e., nonlinear impact). Second, customers are less likely to get their problems resolved if their agents are communicating with different customers in parallel. We found that the same result does not hold when a service agent is assigned a new task, but he or she does not actively work on this task. This is an interesting finding, which brings a new perspective to the literature about the impact of discretionary (i.e., voluntary) task switching (Madjar and Shalley 2008, Payne et al. 2007). Our results further suggest that the presumed productivity gains of multitasking may cease to materialize in the long run because customers who have received unsatisfactory service may return to the system and create future demands. Overall, we demonstrate that multitasking has an indirect, but significant, negative impact on service quality and customer satisfaction through in-service delay and problem resolution paths.

Finally, our study contributes to the extensive literature on the well-known control problem in contact centers known as the call routing problem (Aksin et al. 2007) from a unique angle. Existing work in this area has tackled this problem mostly from a theoretical perspective using stylized analytical models. We showed that explanatory and predictive (i.e., segmentation) models built on observational data can be simultaneously used to develop simple but actionable routing strategies. We believe that such routing strategies are widely applicable and can have significant practical implications for contact centers.

Acknowledgments

The authors thank the department editor, associate editor, and the three anonymous reviewers for their careful reading of the manuscript and their many thoughtful suggestions, which greatly improved the paper.

Endnotes

¹To illustrate, consider a hypothetical example in which there is a single server with a 100-second total service capacity. If there were two jobs that share this service capacity, each time slot would be 50 seconds, and a given job would need to wait $100 - 50 = 50$ seconds between its processing turns. On the other hand, if there were four jobs that share this service capacity, each time slot would be $100/4 = 25$ seconds, and a given job would need to wait $100 - 25 = 75$ seconds between its processing turns.

²This nonlinearity is assumed to be a result of additional slowdown due to cognitive load as well as the diminishing marginal benefits of bursty work at increased levels of multitasking.

³One exception is transferred cases (i.e., agents can manually transfer a customer to another agent). We exclude all transferred cases from the data set in our analysis since they account for only a very small portion of the data.

⁴Multitasking information is not stored in the databases of the contact center. Therefore, we develop procedures to measure multitasking levels from the existing data sources.

⁵Our data cleaning is performed after multitasking computations. We use the entire system data (including agent-to-agent sessions) for multitasking computations because otherwise we may miss certain sessions that were multitasked but preemptively removed from the data set. Once multitasking values are computed, we remove agent-to-agent sessions from the data.

⁶See, for example, <http://www.gpeters.com/names/baby-names.php> (accessed May 15, 2016).

⁷An exception is dummy variables such as gender.

⁸Section 5 reports results of robustness tests for the analyses reported here such as matching.

⁹Being idle in a chat session is a rare and unlikely situation for agents. Therefore, this value is expected to be a conservative estimate for the true impact of multitasking.

¹⁰This value is for the average session in the data set (which has seven agent responses) and assumes that the agent is equally distributing his or her workload among two parallel sessions.

¹¹See, for example, <https://www.surveygizmo.com/survey-blog/survey-response-rates/> (accessed May 15, 2016).

¹²This is because the presumed productivity benefits of multitasking can only be realized when agents switch between different tasks.

¹³This matching procedure is, in a sense, stricter than propensity score matching, because we require matches to be identical on many dimensions. For robustness, we also apply the typical propensity score matching procedure. The average treatment effect on the treated (ATET) is consistent with our main findings for both MT_1 and MT_2 . Furthermore, the propensity matching method allows us to test the sensitivity of the results to unobservables. We found that for the unobserved confounder to be driving our results, this unobserved factor needs to produce a 4.8 times (for MT_1) or 1.8 times (for MT_2) increase in the odds of multitasking, which is again much higher than the thresholds reported in the literature (e.g., Sun and Zhu 2013).

¹⁴While there may be different values for S that would fit the data better, we pick $S = 2$ since our main goal is to demonstrate the applicability of finite mixture models for deriving actionable strategies. A two-segment model could guide the contact center to implement a simple two-class routing policy—that is, use different routing rules for customers who are more or less sensitive to multitasking. Models with more segments will also require more customer covariates.

References

Adler R, Benbunan-Fich R (2012) Juggling on a high wire: Multitasking effects on performance. *Internat. J. Human-Comput. Stud.* 70(2):156–168.

- Aksin Z, Armory M, Mehrotra V (2007) The modern call center: A multi-disciplinary perspective on operations management research. *Productions Oper. Management* 16(6):665–688.
- Anderson EW, Fornell C, Rust RT (1997) Customer satisfaction, productivity, and profitability: Differences between goods and services. *Marketing Sci.* 16(2):129–145.
- Aral S, Brynjolfsson E, Van Alstyne M (2012) Information, technology and information worker productivity. *Inform. Systems Res.* 23(3):849–867.
- Bailey BP, Konstan JA (2006) On the need for attention-aware systems: Measuring effects of interruption on task performance, error rate, and affective state. *Comput. Human Behav.* 22(4):685–708.
- Bapna R, Goes P, Wei KK, Zhang Z (2011) A finite mixture logit model to segment and predict electronic payments system adoption. *Inform. Systems Res.* 22(1):118–133.
- Bavafa H, Hitt LM, Terwiesch C (2017) The impact of e-visits on visit frequencies and patient health: Evidence from primary care. Working paper, University of Wisconsin–Madison, Madison. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2363705.
- Brady MK, Cronin JJ (2001) Some new thoughts on conceptualizing perceived service quality: A hierarchical approach. *J. Marketing* 65(3):34–49.
- Cameron A-F, Webster J (2011) Relational outcomes of multicommuting: integrating incivility and social exchange perspectives. *Organ. Sci.* 22(3):754–771.
- Cameron A-F, Webster J (2013) Multicommuting: Juggling multiple conversations in the workplace. *Inform. Systems Res.* 24(2):352–371.
- Campello F, Ingolfsson A, Shumsky RA (2017) Queueing models of case managers. *Management Sci.* 63(3):882–900.
- Casado Diaz AB, Más Ruiz FJ (2002) The consumer's reaction to delays in service. *Internat. J. Service Indust. Management* 13(2):118–140.
- Charron S, Koechlin E (2010) Divided representation of concurrent goals in the human frontal lobes. *Science* 328(5976):360–363.
- Chebat J-C, Filiatrault P (1993) The impact of waiting in line on consumers. *Internat. J. Bank Marketing* 11(2):35–40.
- Coviello D, Ichino A, Persico N (2014) Time allocation and task juggling. *Amer. Econom. Rev.* 104(2):609–623.
- Czerwinski M, Horvitz E, Wilhite S (2004) A diary study of task switching and interruptions. *Proc. ACM Conf. Human Factors Comput. Systems (CHI '04)* (ACM, New York), 175–182.
- Demers A, Keshav S, Shenkar S (1989) Analysis and simulation of a fair queueing algorithm. *ACM SIGCOMM Comput. Comm. Rev.* 19(4):1–12.
- Dube-Rioux L, Schmitt BH, Leclerc F (1989) Consumers' reactions to waiting: When delays affect the perception of service quality. Srull TK, eds. *Advances in Consumer Research* Vol. 16 (Association for Consumer Research, Provo, UT), 59–63.
- Durrande-Moreau A (1999) Waiting for service: Ten years of empirical research. *Internat. J. Service Indust. Management* 10(2):171–194.
- Fornell C, Johnson MD, Anderson EW, Cha J, Bryant BE (1996) The American Customer Satisfaction Index: Nature, purpose, and findings. *J. Marketing* 60(4):7–18.
- Gow ID, Ormazabal G, Taylor DJ (2010) Correcting for cross-sectional and time-series dependence in accounting research. *Accounting Rev.* 85(2):483–512.
- Gupta S, Chintagunta PK (1994) On using demographic variables to determine segment membership in logit mixture model. *J. Marketing Res.* 31(1):128–136.
- Hall NG, Leung JY-T, Li C-L (2015) The effects of multitasking on operations scheduling. *Production Oper. Management* 24(8):1248–1265.
- Heckman JJ (1979) Sample selection bias as a specification error. *Econometrica* 47(1):153–161.
- Human Capital Management Institute (2010) Managing an organization's biggest cost: The workforce. White paper, Human Capital Management Institute, Marina Del Rey, CA.
- KC DS (2013) Does multitasking improve performance? evidence from the emergency department. *Manufacturing Service Oper. Management* 16(2):168–183.
- Luo J, Zhang J (2013) Staffing and control of instant messaging contact centers. *Oper. Res.* 61(2):328–343.
- Madjar N, Shalley CE (2008) Multiple tasks' and multiple goals' effect on creativity: Forced incubation or just a distraction? *J. Management* 34(4):786–805.
- Mark G, Gudith D, Klocke U (2008) The cost of interrupted work: more speed and stress. *Proc. ACM Conf. Human Factors Comput. Systems (CHI '08)* (ACM, New York), 107–110.
- Mehrotra V, Ross K, Ryder F, Zhou Y-P (2012) Routing to manage resolution and waiting time in call centers with heterogeneous servers. *Manufacturing Service Oper. Management* 14(1):66–81.
- O'Leary MB, Mortensen M, Woolley AW (2011) Multiple team membership: A theoretical model of its effects on productivity and learning for individuals and teams. *Acad. Management Rev.* 36(3):461–478.
- Payne SJ, Duggan GB, Neth H (2007) Discretionary task interleaving: heuristics for time allocation in cognitive foraging. *J. Experiment. Psych.* 136(3):370–388.
- Pennebaker JW, Francis RJ, Booth ME (2001) *Linguistic Inquiry and Word Count (LIWC)* (Lawrence Erlbaum Associates, Mahwah, NJ).
- Reinsch NL, Turner JW, Tinsley CH (2008) Multicommuting: A practice whose time has come? *Acad. Management Rev.* 33(2):391–403.
- Rosenbaum PR (2002) *Observational Studies* (Springer, New York).
- Sun M, Zhu F (2013) Ad revenue and content commercialization: Evidence from blogs. *Management Sci.* 59(10):2314–2331.
- Szymanski DM, Henard DH (2001) Customer satisfaction: A meta-analysis of the empirical evidence. *J. Acad. Marketing Sci.* 29(1):16–35.
- Tan TF, Netessine S (2014) When does the devil make work? An empirical study of the impact of workload on worker productivity. *Management Sci.* 60(6):1574–1593.
- Tan Y, Mookerjee VS, Moinszadeh K (2005) Optimal processing policies for an e-commerce web server. *INFORMS J. Comput.* 17(1):99–110.
- Tax SS, Brown SW, Chandrashekar M (1998) Customer evaluations of service complaint experiences: Implications for relationship. *Jour. Marketing* 62(2):60–76.
- Taylor S (1994) Waiting for service: The relationship between delays and evaluations of service. *J. Marketing* 58(2):56–69.
- Tekin E, Hopp WJ, Van Oyen MP (2009) Pooling strategies for call center agent cross-training. *IIE Trans.* 41(6):546–561.
- TELUS International (2011) Best practices: Online chat sales. TELUS International, Vancouver, Canada.
- Tezcan T, Zhang J (2014) Routing and staffing in customer service chat systems with impatient customers. *Oper. Res.* 62(4):943–956.
- Tom G, Lucey S (1995) Waiting time delays and customer satisfaction in supermarkets. *J. Services Marketing* 9(5):20–29.
- Wedel M, Kamakura WA (2000) *Market Segmentation: Conceptual and Methodological Foundations*, 2nd ed. (Kluwer Academic Publishers, Boston).
- Zeithaml VA, Berry LL, Parasuraman A (1996) The behavioral consequences of service quality. *J. Marketing* 60(2):31–46.