



Manufacturing & Service Operations Management

Publication details, including instructions for authors and subscription information:
<http://pubsonline.informs.org>

Customer Service Competition in Capacitated Systems

Joseph Hall, Evan Porteus,

To cite this article:

Joseph Hall, Evan Porteus, (2000) Customer Service Competition in Capacitated Systems. *Manufacturing & Service Operations Management* 2(2):144-165. <http://dx.doi.org/10.1287/msom.2.2.144.12353>

Full terms and conditions of use: <http://pubsonline.informs.org/page/terms-and-conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

© 2000 INFORMS

Please scroll down for article—it is on subsequent pages



INFORMS is the largest professional society in the world for professionals in the fields of operations research, management science, and analytics.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

Customer Service Competition in Capacitated Systems

Joseph Hall • Evan Porteus

Amos Tuck School of Business Administration, Dartmouth College, Hanover, New Hampshire 03755

Graduate School of Business, Stanford University, Stanford, California 94305

joseph.m.hall@dartmouth.edu • eporteus@stanford.edu

We investigate a simple dynamic model of firm behavior in which firms compete by investing in capacity that is used to provide a good or service to their customers. There is a fixed total market of customers whose demands for the good or service are random and who divide their patronage between the firms in each period. Periodically, the market shares of the two firms can change based on the realized level of customer service provided in the prior period. We assume that the expected level of customer service can be expressed as a function of the (per customer) capacity of the firms' service delivery systems, and that service declines as the capacity decreases. The firms differ in their customers' willingness to defect when confronted by service failure. The primary issue we address is the firms' capacity decisions in response to customer service concerns and competitive pressure. We provide conditions under which the firms' optimal (i.e., equilibrium) capacity levels in a period are proportional to the size of their respective customer bases in that period. Further, we develop expressions for the value of a firm's customers and the implicit cost of service failure. Results for both single-period and finite-horizon problems are investigated and applied to two examples: (1) competition between Internet service providers who operate systems that we approximate by simple loss-type queueing models, and (2) competition between make-to-stock producers who operate systems that we approximate by newsvendor inventory models. For both examples, solutions are derived and interpreted.

(Capacity Management; Game Theory; Customer Loyalty; Customer Service; Newsvendor Model)

1. Introduction

We develop a model of interfirm service competition for customers. Customers of Internet service providers (ISPs) may switch to a competing provider if they receive poor service. Retail customers for goods may begin shopping at a competing retailer if their usual choice of retail outlet is out of stock. This paper presents a dynamic model in which firms must decide each period on their capacity level for the coming period, and the level of service that they provide to their customers influences how many of those customers remain loyal and how many of them switch to competing firms. Our model yields simple measures of the value

of a firm's customers and the implicit costs associated with poor service. These costs, in turn, influence the firm's capacity decisions.

We apply the results of our general model to the examples cited above: call service competition between Internet service providers, and stock availability competition between suppliers. For the former we use a produce-to-order queueing model and for the latter we use a produce-to-stock newsvendor inventory model. In both cases, our models provide an endogenous measure of the present value of the future consequences of current service failures. In the case of the newsvendor model, we demonstrate that this measure takes the

simple form of an augmented stockout penalty. For both of the examples studied, we derive explicit solutions, present a numerical example, and discuss the managerial implications of our results.

The tension in our model is created by an interaction among customer service concerns, capacity choices, and competitive forces. We focus on the question of how customer service issues and knowledge of competitors can be incorporated into a firm's capacity-setting decisions. In general, there are four primary phenomena that characterize the interplay between customer service and competition:

- (1) a firm with poor customer service may lose customers to its competitors;
- (2) a firm with good customer service may capture some business from its competitors;
- (3) if firms exhibit poor customer service, customers may leave the market; and
- (4) if some firms exhibit good customer service, customers may enter the market.

This paper addresses phenomenon (1) and, indirectly, (2), but not (3) and (4). In particular, we assume that there are two firms and that there is a constant total number of customers who patronize either one or the other. That is, no existing customers leave the system and no new customers enter the system. Although a multitude of methods could be used to model Phenomena (1) and (2), we develop a simple, yet insightful, model. We apply an operations management perspective to the notion of customer service by focusing on those aspects of customer service that depend upon the customer load on the service delivery system, e.g., service denial, as with a busy signal or an item stockout.

The remainder of §1 discusses related literature and details our assumptions and other model background. Sections 2 and 3 contain analyses for the general one-period and multiperiod finite-horizon models, respectively. Sections 4 and 5 present example applications and numerical examples for ISP competition and product availability competition, respectively. Section 6 details an extension of the general model to more than two firms and §7 concludes.

Review of Existing Literature

The issues of customer retention and loyalty have received some attention in the management press in recent years. Reichheld and Sasser (1990) argue that service firms need to focus on zero customer defections,

analogous to the zero defects movement in manufacturing. They claim that "[Customer] defection rates are an accurate leading indicator of profit swings . . ." and argue for using defection rates as a "primary performance measure." The authors also focus on customer valuation as a means to estimate the cost of a customer defection. They state that "[m]ost [accounting] systems focus on current period costs and revenues and ignore expected cash flows over a customer's lifetime." The authors present the example of an enlightened Domino's Pizza franchisee who calculated that a customer was worth \$5000 over the 10-year life of his franchise contract. The role of capacity is (briefly) addressed in the example of a firm that ". . . is losing customers because of long lines [and] can estimate what percentage of defectors it would save by buying new cash registers. . . ." Reichheld (1996) examines many of these same ideas in more depth. In that work, the author develops a methodology to measure the net present value of a firm's customer base in light of customer defections. To support this analysis, the author defines a "customer loyalty coefficient" that measures ". . . how much economic force it takes to move these customers from one supplier to another." Note that this notion of loyalty is inherently firm-specific and not customer-specific. Rather than capturing individual customer behavior, it measures aggregate customer response at a supplier. Our model captures many of the issues raised by Reichheld and Sasser (1990) and Reichheld (1996). We focus on the problem of a firm's "best" capacity choice in light of potential customer defections to competing suppliers, where the suppliers may exhibit differing loyalty coefficients. Our solution procedure endogenously derives measures of customer value over multiple-period horizons and uses these measures as an input to the capacity decision in a natural way.

Jones and Sasser (1995) examine the link between customer satisfaction and customer loyalty. Their findings indicate that this link depends, in part, on the competitiveness of the industry. In this work, we do not deal directly with the notion of customer satisfaction; rather, we look at capacity-driven service failures, which implicitly or explicitly drive customer satisfaction, and link these failures directly to customer retention. One advantage of our approach is that service

failures and retention can often be measured directly and objectively, while customer satisfaction may be harder to define and measure.

Few behavioral studies have focused on the underlying drivers of customer defections. Keaveney (1995) presents one such example. This work is an exploratory study of customers' switching behavior across various service industries. Keaveney (1995) found that the response category "core service failures," defined as mistakes or other technical problems with the service, encompassed the most reasons cited by customers to explain their previous switching behavior (mentioned in 44% of the incidents). The second most cited category (34% of incidents) was "service encounter failures," which included problems with service employees' behaviors and attitudes. "Pricing" was the third most frequently mentioned category and "attraction by competing firms" was the sixth most frequent category, mentioned in 30% and 10% of the incidents, respectively. These results suggest that customers often behave reactively, i.e., service failures may prompt a change of patronage. The two most frequently mentioned reasons for switching provider—core service failures and service encounter failures—are both, in part, related to our model of service failures that are capacity driven. For "core service failures" the notion of capacity could be the system service rate or inventory level. "Service encounter failures" could also be fit into this framework if overloaded service employees are likely to express behaviors or attitudes that are interpreted by customers as "service failures." In this context, system capacity could be the number of service employees, or the level of periodic service training that employees receive.

A number of papers have examined the role of competition in make-to-order environments where queueing phenomena dominate delivery system behavior. Mendelson (1985) and Mendelson and Whang (1990) examine the behavior of a single firm, and Loch (1991) extends some elements of their work to a competitive environment. These models focus on the time that customers spend waiting for service and associate a cost with this waiting time. Customers possess a reservation value for their jobs that must exceed the sum of system admission price and expected delay cost if they are to demand service. Customers in these models

are treated as rational decision makers who are informed of the long-run system parameters, such as expected arrival and service rates, but not of current queue characteristics. Thus, the results presented in these models are static equilibria based upon expected system behavior. In contrast, our model is inherently dynamic. We assume that a firm finds itself with a certain market share at some point in time and wishes to determine what capacity it should provide over the next period. Kalai et al. (1992) and Li and Lee (1994) assume customers' supplier choice does depend on current queue characteristics; arriving customers enter a pool that is divided between two (or more) competing servers as each server completes its current service. Kalai et al. (1992) assume that the competitors are differentiated only by service rate; thus, customers choose the service provider that first becomes idle after they reach the head of the queue. Li and Lee (1994) extend this to an environment where customers possess a utility over price and quality of service in addition to delivery time. A key assumption of these models is that once customers arrive to the system they can switch between the two providers up until the point of entering service.

A number of papers have dealt with the notion of competition based upon product availability. Among these is the work of Li (1992). The author studies the optimal balance between make-to-order and make-to-stock, and examines a setting where firms compete based upon delivery time. The work of Li (1992) is similar to that of Kalai et al. (1992) and Li and Lee (1994) in that firms compete by racing to satisfy aggregate demand as it occurs—demand is not a priori assigned to one competitor or the other. The customer jockeying for service assumed in these three papers is realistic in some settings, while our assumption of customer loyalty is appropriate in other settings. Lippman and McCardle (1997) also make availability competition the subject of study. The authors present a natural competitive extension of the standard newsvendor model. In this model, potential customer demand for a perishable good is initially allocated among the firms in the market, demand is realized, and then unmet demand is reallocated among the firms with remaining inventory. Future demands are unaffected. In contrast, our model has no reallocation of unmet demands—

there is only the possibility of switching, which would begin to be felt in the next period. Thus, their model is inherently a single-period model (while ours is necessarily multiperiod) and consequently it does not readily allow study of longer-term customer service phenomena, such as lost goodwill costs. We note, though, that the model of Lippman and McCardle (1997) applies to a more general class of demand distributions than does ours.

Gans (1999a) constructs a model of individual consumer behavior in response to imperfect quality. The model presented is a myopic variant of a Bayesian choice rule; each contact between a consumer and a firm causes the consumer to update their prior beliefs about the quality of the firm's good or service. The focus of Gans (1999a) is developing measures of customer loyalty. Gans (1999b) builds upon these measures by presenting sensitivity analyses and models of quality provision in both monopolistic and competitive environments. The general setting of the models in Gans (1999a) and Gans (1999b) is similar in spirit to ours. The models differ in their emphasis, however. Our model focuses on firm behavior, while the work of Gans (1999a, 1999b) focuses on consumer response. Further, the equilibria presented in Gans (1999b) are based on long-run expected behavior, unlike the dynamic equilibria we develop.

McGahan and Ghemawat (1994) examine competition to retain customers in a context that is neither explicitly make-to-stock nor make-to-order, as their model is deterministic. They construct a two-stage game in which the players invest in capacity in the first stage to serve their existing customers. Some fraction of these customers are disaffected after the first stage, as a function of the capacity investment per customer. In the second stage, the firms compete on price, with the lower-price firm taking all customers disaffected after the first stage. They develop conditions such that the firm with the larger initial market share provides a higher level of service in the first stage of the game and charges a higher price in the second stage. The authors also empirically test the hypothesis that customer retention rates increase with firm size. In contrast to our model, McGahan and Ghemawat include price in their model of competition. The model we

present, however, is explicitly multistage and stochastic.

Our model stands in contrast to most of the above models in that we do not model customers as strictly rational in an economic sense. Rather (as Lippman and McCardle 1997 and the first stage game of McGahan and Ghemawat 1994), we model customers who may react to service failures by seeking out another service provider. However, unlike those papers, our formulation is explicitly multiperiod. This difference allows us to study the role of long-term competitive issues that result from repeated interactions amongst firms and customers.

Modeling Context

The following paragraphs outline our assumptions and model context. We later demonstrate that these assumptions are satisfied for the two examples we study. Our assumptions appear in italics and are numbered for reference. Explanations are included where appropriate. Throughout, i refers to an arbitrary firm and j refers to the other firm (in the case of two firms). That is, $j \neq i$. Other symbols are defined where first used and are summarized in the appendix.

(A1). *There are only two firms that compete in the market.*

Therefore, either $(i, j) = (1, 2)$ or $(i, j) = (2, 1)$. Section 6 discusses an extension of our model to more than two competitors. We focus our attention on the two competitor case so that our results are not overshadowed by the complexity of the models.

(A2). *Firms seek to maximize the expected present value of their respective returns over a finite time horizon, consisting of periods $t = 1, 2, \dots, T$. The one-period discount factor α is strictly positive (and less than one).*

We do not allow the firms the choice of exiting from the market.

(A3). *Decisions are made once each period by firms and customers.*

At the start of each period, customers choose which firm they will patronize and, subsequently, firms choose their capacity levels. Assumption (A11) clarifies what each firm knows at the time they make their capacity decisions. We use "period" and "month" interchangeably; the period length is arbitrary.

(A4). Capacity is leased by each firm at a (strictly positive) unit monthly capacity cost of b and is available with zero leadtime.

Another interpretation is that the capacity market is so liquid that capacity can be bought and sold at the same unit price. In particular, total capacity costs for firm i in month t are $b\mu_{it}$, where μ_{it} represents the capacity selected by firm i in month t . Any positive capacity can be selected each month; we leave the study of capacity inflexibility to future work. For produce-to-order settings, one natural capacity measure is the rate at which orders can be serviced. For produce-to-stock settings, one natural capacity measure is the inventory level after replenishment.

(A5). The absolute size of the total market is constant over time and is normalized to equal one.

Thus, we use λ_{it} to represent both the absolute market size and the fractional market share for firm i in month t . As a consequence, $\lambda_{jt} = 1 - \lambda_{it}$ for all t . Within the context of the examples presented in §4 and §5, we will use λ_{it} to represent either mean arrival rate or mean number of units demanded, respectively.

(A6). Customers are continuously divisible and homogeneous.

Essentially, the total number of customers in the market is so large that the market share can be considered to be a continuous variable. Furthermore, the distribution of service demands is the same across all customers.

(A7). The normalized capacity of firm i in month t , y_{it} , is defined implicitly by $\mu_{it} \equiv y_{it}\lambda_{it}$. There is a single measure of customer service, $h_i(y_{it})$, where $\lambda_{it}h_i(y_{it})$ gives the expected number of firm i 's customers that experience service failures in month t when firm i has a normalized capacity of y_{it} . In addition, for positive real arguments, h_i is twice differentiable, decreasing, strictly convex, and takes on values between zero and one.

That is, customer service depends solely on the normalized capacity. In effect, this assumption rules out capacity pooling effects as service failures depend only on the normalized (per customer) capacity and not on the absolute level of market share or capacity level. In

some instances, it may be that as market share increases, less per customer capacity is required to provide the same level of service. This is the case in the simple example of a make-to-stock environment where individual customer demands are independent, identically distributed normal random variables; the same service level can be provided with proportionally less inventory per customer as total demand increases. Such a case is ruled out by this assumption. We define y_{it} implicitly to account for the possibility that $\lambda_{it} = 0$. Note that the inverse of the normalized capacity is the system utilization in a queueing setting in which all customer arrivals are served.

(A8). Firm i collects random revenues in month t with mean $p\lambda_{it}(1 - \beta_i h_i(y_{it}))$, where p is an exogenous (strictly positive) unit price and $\beta_i \in [0, 1]$.

That is, the mean return to a firm is the given unit price, reduced by a service failure factor, for each customer in its customer base. In particular, β_i can be interpreted as the fraction of service failures that result in no revenue in month t for firm i , either due to inability to collect revenue or a refund. We include β_i because industries (and firms within industries) may differ in the immediate revenue consequences of a service failure. For example, $\beta_i = 1$ when a customer who experiences a service failure will refuse to make a purchase or be unable to make a purchase (as in the case of an inventory stockout). In contrast, $\beta_i = 0$ when a customer who experiences a service failure will go through with a purchase, such as an ISP customer who has purchased the service in advance via a subscription fee.

(A9). The expected number of customers of firm i who switch to the other firm at the beginning of the next month is $\lambda_{it}\gamma_i h_i(y_{it})$, where $\gamma_i \in (0, 1]$.

This assumption casts customers in our model as experimenters rather than information gatherers and rational-choice decision makers—if, in the current month, customers are confronted with service failures, they may switch to the other provider in the next month, regardless of the prospects of obtaining better service after the switch. Another interpretation of this assumption is that firms view their customers only in the aggregate, focusing on their overall, rather than

individual, behavior. An alternative model formulation might assume that customers only switch firms if they expect to receive better service by doing so, based on the historical performance of both firms and expected switching of all other customers. Of course, such a model needs to assume that customers possess information about the selection of firms and about the other customers who are seeking service. In contrast, our assumption may be reasonable in environments where customers cannot easily obtain comparative performance information for the firms, such as when the firms are frequently upgrading or expanding their systems, or when the number of customers at each firm is fluctuating over time. The exogenous constant γ_i can be interpreted in a number of ways: an aggregate firm-specific measure of users' tradeoffs between the aggravation of service failure and the costs of switching firms, a measure of inherent differences in the services of the firms that make customers more or less loyal, or the fraction of service failures from which the firm does not gracefully recover. In the language of Reichheld (1996), γ_i is a "loyalty coefficient." Note that we make no assumption about the distribution of the number of customers who switch firms. See §7 for further discussion of this assumption.

(A10). *The end-of-horizon value of a firm can be expressed as a positive affine function of the size of its market share, i.e., as $V_i \lambda_{i,T+1} + W_i$, where $V_i > 0$ and $W_i \geq 0$ are exogenous constants.*

The parameter V_i represents the end-of-horizon value per customer at firm i ; i.e., it captures how the value of the firm increases in the size of its customer base. The parameter W_i represents the value of the business that is independent of the current market share.

(A11). *In addition to the basic problem parameters, firm i knows the following at the start of every month t : λ_{it} (and, therefore, λ_{jt}), β_{iv} , β_{jv} , γ_{iv} , γ_{jv} , h_{iv} , h_{jv} , V_{iv} , V_{jv} , and W_i .*

These informational assumptions are keys to our equilibrium analysis of the problem. We shall see that each firm can perfectly predict the capacity that the competing firm will select, and thus we need not assume that each firm is formally informed of the other firm's actions. However, our analysis requires that each firm

knows their market share prior to making their capacity decisions each month. This assumption is innocuous if the firms require customers to subscribe in advance of using the service. Our analysis also assumes that the firms can accurately estimate the size of the total customer base (which is normalized to unity in our analysis) and the parameters β_{iv} , β_{jv} , γ_{iv} , γ_{jv} , h_{iv} , and h_{jv} , all of which are assumed to be invariant over the problem horizon. Finally, the requirement that the firms know the end-of-horizon parameters can be interpreted as a requirement that they can value a business operation. We do not formally address the effects of uncertainty in these parameters.

Modeling Preliminaries

Henceforth, unless indicated otherwise, assume that (A1)–(A11) hold. We seek to characterize the equilibrium actions, and returns that result from those actions, for the firms in our idealized environment.

We search for a subgame perfect (Nash) equilibrium (see Kreps 1990) among the firms in the market. Note that the customers do not participate in the game that we define; aggregate customer behavior is assumed to follow, in expectation, an exogenous rule, as per (A9). A subgame perfect equilibrium gives, for each possible starting point for each period of the game, a Nash equilibrium for each firm. (See Kreps 1990, Kreps and Wilson 1982, Fudenberg and Tirole 1991], and Mas-Colell et al. 1995 for useful expositions of Nash equilibria and other appropriate concepts for dynamic games, such as sequential equilibria and Markov perfect equilibria.) Indeed, by working backwards in the usual way of dynamic programming, we shall find that there is a unique Nash equilibrium for each possible starting state in each period of the game. The strategies that result from pasting together these decisions across all states and periods comprise a subgame perfect equilibrium. We shall discuss the ramifications later.

As a consequence of our assumptions, the market share of firm i in month $t + 1$ consists of the market share of firm i in month t less the share lost to firm j plus the share gained from firm j . By (A9), once firm i has settled upon its normalized capacity at the beginning of month t , the expected fraction of firm i 's customers who will switch to firm j in month $t + 1$ is given by $\gamma_j h_{ji}(y_{it})$. Furthermore, the expected market share at firm i in month $t + 1$ can be expressed as:

$$E(\lambda_{i,t+1} | \lambda_{it}, \lambda_{jt}) = \lambda_{it} - \lambda_{it}\gamma_i h_i(y_{it}) + \lambda_{jt}\gamma_j h_j(y_{jt}). \quad (1)$$

By (A5), this expression can be rewritten as:

$$E(\lambda_{i,t+1} | \lambda_{it}) = \lambda_{it}(1 - \gamma_i h_i(y_{it}) - \gamma_j h_j(y_{jt})) + \gamma_j h_j(y_{jt}). \quad (2)$$

These expressions are used throughout our models, as they give the expected state a firm will find itself in next month as a function of the current states and actions of both firms. We will later demonstrate that the optimal value function is linear in the state, which allows us to focus our attention on the expected state transition. Expression (2) has an insightful interpretation. Firm, j , the competitor of firm i , has a month t market share of $\lambda_{jt} = 1 - \lambda_{it}$. Thus, the last term of (2) corresponds to the assumption that firm j has the entire market share in month t and indicates the fraction of that market that will switch to firm i at the beginning of month $t + 1$, which depends solely on what firm j does. To the extent that firm i has customers in month t , the first term of (2) indicates the effect their number will have on how many firm i has in month $t + 1$. That is, firm i will expect to lose the fraction $\gamma_i h_i(y_{it})$ due to its own service failures. Because firm i will also expect to gain a fraction of its competitor's customers, the more share firm i has in month t , the smaller is firm j 's share in month t , which reduces the share firm i will expect to gain from firm j in month $t + 1$. That is, $\gamma_j h_j(y_{jt})$ is the *direct* expected lost fraction of market share and $\gamma_j h_j(y_{jt})$ is the *indirect* expected lost fraction of market share.

2. The One-Period Model

We begin by investigating the form of the one-period problem that captures the immediate consequences of the capacity decision in one month. These consequences consist of the current period expected revenues and capacity expenses and the expected revenues to be received the following month as a result of the service provided to customers (of both providers). This formulation allows us to capture, in a simple problem, both the immediate and future direct consequences of a single decision. We formulate the problem for firm i in month T as:

$$\begin{aligned} \max_{\mu_{iT} \geq 0} \{ & p\lambda_{iT}(1 - \beta_i h_i(y_{iT})) - b\mu_{iT} \\ & + \alpha V_i E(\lambda_{i,T+1} | \lambda_{iT}) + \alpha W_i \}, \end{aligned} \quad (3)$$

where $E(\lambda_{i,T+1} | \lambda_{iT})$ is defined in Expression (2) and V_i and W_i are the end-of-horizon parameters for firm i (defined in (A10)). The first two terms represent the expected revenue to be collected during month T , the third term is the cost of capacity for month T and the fourth and fifth terms are the expected present value of the end-of-horizon values collected at the start of the next month, i.e., at the end of the problem horizon. Note that, using (2), we can rewrite (3) as a search for the optimal normalized capacity in month T :

$$\begin{aligned} \max_{y_{iT} \geq 0} \{ & \lambda_{iT}[p(1 - \beta_i h_i(y_{iT})) - b y_{iT} + \\ & \alpha V_i(1 - \gamma_i h_i(y_{iT}) - \gamma_j h_j(y_{jT}))] + \\ & \alpha V_j \gamma_j h_j(y_{jT}) + \alpha W_i \}. \end{aligned} \quad (4)$$

One ambiguity results from our focus on Problem (4) instead of Problem (3): if $\lambda_{iT} = 0$, the problem is independent of the decision variable, even though the optimal capacity choice, μ_{iT} , is clearly zero. This ambiguity is overcome by restating the search for the optimal normalized capacity as follows:

$$\begin{aligned} \max_{y_{iT} \geq 0} \{ & p(1 - \beta_i h_i(y_{iT})) - b y_{iT} + \\ & \alpha V_i(1 - \gamma_i h_i(y_{iT}) - \gamma_j h_j(y_{jT})), \end{aligned} \quad (5)$$

which assumes that the firm is interested in maximizing the contribution per customer. Henceforth in this section, the subscript T will be suppressed for clarity.

THEOREM 1. *Given that the value to firm i of having market share λ at the end of the period is $V_i \lambda + W_i$ for each i , we have:*

(a) If

$$(p\beta_i + \alpha V_i \gamma_i) h'_i(0) \leq -b, \quad (6)$$

then the first-order condition

$$(p\beta_i + \alpha V_i \gamma_i) h'_i(y_i) = -b \quad (7)$$

has a unique solution y_i^0 ;

(b) Define

$$y_i^* = \begin{cases} y_i^0 & \text{if } (p\beta_i + \alpha V_i \gamma_i) h'_i(0) \leq -b \\ 0 & \text{otherwise;} \end{cases}$$

(c) Then, firm i will select normalized capacity y_i^* regardless

of what firm j does. That is, (y_1^*, y_2^*) forms a unique dominant pure strategy Nash equilibrium.

PROOF. For (a), since h_i is strictly convex by (A7) and p , α , γ_i , and V_i are strictly positive and β_i is positive by (A2), (A8), (A9), and (A10), it is easily verified that each firm's objective function, as defined by Expression (5), is strictly concave on the positive real line given its competitor's decision. Expression (6) is the condition that (7) has a positive solution (the derivative of firm i 's objective function is positive at zero), and thus y_i^0 is unique as claimed in (a). For (c), note that (7) contains no references to firm j , so y_i^* is chosen without regard for the action of firm j . Thus, the action vector (y_1^*, y_2^*) is, by definition, a unique dominant pure strategy Nash equilibrium. \square

One consequence of Part (c) of Theorem 1 (dominance of equilibrium) is that the Stackleberg (sequential-move) equilibrium and the simultaneous-move equilibrium for this problem are equivalent. Therefore, the equilibrium is unaffected if we allow one firm to act first and allow the other firm to observe that action before making its choice. Also notable in the solution is that the optimal normalized capacity for each firm is independent of their market share—equivalently, a firm's optimal capacity decision is a purely linear function of the size of the firm's current customer base (the state a firm finds itself in). Part (b) gives the optimal solution: if h_i is not sufficiently steep at zero, that is, (6) does not hold, then the first-order Condition (7) has no solution, as h_i is only defined for positive arguments. Otherwise (7) has a solution and it is unique. Theorem 2 below details another interesting characteristic of the solution in Theorem 1: If a firm's customers become more sensitive to service failures (either in terms of loss of current period revenue or switching behavior), that firm will provide better service. Other solution characteristics are presented in §3 following the multiperiod model. Proofs of all remaining results appear in the Appendix.

THEOREM 2. *For the single-period problem, y_i^* is an increasing function of γ_i and β_i .*

The simple one-period formulation investigated in this section only accounts for the consequences of a decision over a single month. We next address this shortcoming.

3. The Multiperiod Finite-Horizon Model

In the multiperiod model, by (A10), customers have an affine end-of-horizon value to the firm, which can be expressed as:

$$g_{i,T+1}(\lambda_{i,T+1} | \pi^i, \pi^j) = V_i \lambda_{i,T+1} + W_i,$$

where V_i and W_i are scalars that can differ for the two firms and are independent of the strategies chosen. We define π^i as the policy followed by firm i over months 1 through T , where $\pi_t^i(\lambda)$ is the action (i.e., normalized capacity decision) for firm i called for by strategy π^i when starting month t facing state λ . For the multiperiod problem, the expected returns to firm i can be computed recursively starting at the end of the problem horizon:

$$g_{it}(\lambda_{it} | \pi^i, \pi^j) = \lambda_{it}(p(1 - \beta_i h_i(y_{it})) - b y_{it}) + \alpha E(g_{i,t+1}(\lambda_{i,t+1} | \pi^i, \pi^j)), \quad (8)$$

where $y_{it} = \pi_t^i(\lambda_{it})$.

We next demonstrate that the Nash equilibria capacity decisions and returns to the firms are of a simple form. It is convenient for the discussion of our results to define y_{it}^0 , y_{it}^* , v_{it} , and w_{it} for each i , and $t = 1, 2, \dots, T$, as follows:

If

$$(p\beta_i + \alpha v_{i,t+1} \gamma_i) h_i'(0) \leq -b, \quad (9)$$

then let y_{it}^0 denote a solution to

$$(p\beta_i + \alpha v_{i,t+1} \gamma_i) h_i'(y_{it}^0) = -b; \quad (10)$$

$$y_{it}^* = \begin{cases} y_{it}^0 & \text{if } (p\beta_i + \alpha v_{i,t+1} \gamma_i) h_i'(0) \leq -b \\ 0 & \text{otherwise;} \end{cases} \quad (11)$$

$$v_{it} = p(1 - \beta_i h_i(y_{it}^*)) - b y_{it}^* + \quad (12)$$

$$\alpha v_{i,t+1}(1 - \gamma_i h_i(y_{it}^*) - \gamma_j h_j(y_{jt}^*)),$$

and

$$w_{it} = \alpha w_{i,t+1} + \alpha v_{i,t+1} \gamma_j h_j(y_{jt}^*). \quad (13)$$

Theorem 3 below demonstrates that when a solution to (10) is called for, it is unique, and that the suggestive notation is warranted.

THEOREM 3. *There exists a unique subgame perfect (Nash) equilibrium, (π^i, π^j) , and in each period t we have:*

$$g_{it}(\lambda_{it} | \pi^*, \pi^*) = v_{it}\lambda_{it} + w_{it},$$

where $y_{it}^* = \pi_t^*(\lambda_{it})$ and $y_{jt}^* = \pi_t^*(\lambda_{jt})$ (and thus v_{it} and w_{it}) are independent of the state λ_{it} .

The fact that each firm's normalized capacity decision is unaffected by the other firm's normalized capacity decision in the same period helps explain why our results are so simple. There are two reasons for this separability. The first already has been discussed: A service failure at one firm has no effect on the other firm until the following period. The second is that current high levels of normalized capacity do nothing to attract customers from the other firm. New customers arise only as a result of experiencing a service failure at the other firm. In short, each firm seeks to limit the customers it loses, but cannot influence the customers it gains. Fergani (1976) studies a related inventory model in a noncompetitive environment in which the firm can replenish its market size at a cost, such as through advertising. He gives conditions for a storable good and partial backlogging model in which a myopic policy is optimal.

Note that y_{it}^* , v_{it} , and w_{it} , as recursively defined by (9)–(13), are explicit functions of exogenous parameters. Thus, each firm purchases capacity in each month that is linear in the firm's market share in that period and the constants of proportionality, *i.e.*, the parameters y_{it}^* , can be calculated, *ex ante*, for the entire problem horizon. Because the equilibrium return function, $v_{it}\lambda + w_{it}$, is linear in the state variable λ , we need only concern ourselves with the mean of the distribution of the number of customers who switch providers. Our definitions of v_{it} and w_{it} allow us to express the end of horizon values defined in (A10) as $v_{i,T+1} = V_i$ and $w_{i,T+1} = W_i$.

The return function coefficients v_{it} and w_{it} have insightful interpretations. The coefficient w_{it} can be interpreted as the “fixed,” *i.e.*, independent of market share, value of being in business for firm i in month t under the equilibrium policy. This term results from the future switching of customers that is bound to occur—even if firm i should have no customers in the current month, it will have customers in the future, and these customers have an expected present value. Expression (13), which defines w_{it} recursively, supports this interpretation: The current “fixed” value of

being in business is the “fixed” value next period, appropriately discounted, plus the discounted value of the expected number of customers who switch from the firm's competitor at the start of the next period. In effect, w_{it} is keeping track of the value of all customers who are likely to switch to the firm in the future, starting with zero market share in period t , whereas v_{it} accounts for the value of increasing the current market share, as described next.

The coefficient v_{it} measures the value to the firm that is directly proportional to the current market share. Specifically, v_{it} is the expected value of the entire potential customer base (*i.e.*, market share of unity) in month t for firm i under the equilibrium policy. Multiplying the market share in month t , λ_{it} , by v_{it} gives the value of the customer base for firm i in month t , something akin to the net present value of the firm's customer base discussed by Reichheld (1996). Note that v_{it} is independent of the number of customers, and thus the marginal value of a customer is constant. Expression (12), the recursion for v_{it} , can be interpreted term by term. The first two terms represent the expected contribution in the current period, *i.e.*, the expected revenue less the costs under the optimal capacity decision; the last term gives the present value of the expected market share next period that will result from the current period's market share.

As noted above, the equilibrium presented in Theorem 3 has a very simple form. However, some of the underlying richness of our model is not apparent from the solution for a single period in isolation. In particular, we demonstrated that the equilibrium normalized capacity choice for a firm is independent of both the current period normalized capacity choice of the competing firm and the current market shares of both firms. However, the recursive nature of our model suggests that firms' choices are influenced by the anticipated future behavior of all firms in the market. A natural way to study these interactions is to suppose that one firm deviates from the equilibrium path, and examine the consequences. While a nonequilibrium choice is not credible, it is insightful to examine the effect that a future deviation from the path would have on current capacity decisions. Theorem 4 partially details this relationship.

THEOREM 4. *Holding all else constant, the equilibrium normalized capacity for a firm, y_{it}^**

- (a) *decreases as the firm's own next period normalized capacity decision, $y_{i,t+1}$, decreases below the equilibrium value, i.e., $y_{i,t+1} < y_{i,t+1}^*$*
- (b) *decreases as the firm's own next period normalized capacity decision, $y_{i,t+1}$, increases above the equilibrium value, i.e., $y_{i,t+1} > y_{i,t+1}^*$; and*
- (c) *increases as the competing firm's next period normalized capacity decision, $y_{j,t+1}$, increases.*

These interaction effects arise because the current period decision for firm i depends in part upon the expected value per customer in the next period, $v_{i,t+1}$, which in turn depends upon the decisions made in period $t + 1$ and all periods that follow. The result that deviations from the firm's own equilibrium normalized capacity decision in the next period have a negative effect upon the current normalized capacity decision is a consequence of $y_{i,t+1}^*$ maximizing $v_{i,t+1}$ over choices of $y_{i,t+1}$. Any other choice of $y_{i,t+1}$ leads to a lower value being assigned to the firm's customers in period $t + 1$. Thus, the best response for the firm in period t is to provide less normalized capacity than if $y_{i,t+1} = y_{i,t+1}^*$. Part (c) demonstrates that competing firms' normalized capacity levels in adjacent periods are positively related. If one firm believed the other firm was going to deviate from its equilibrium the following period by investing in more capacity, then the first firm would be induced to buy more capacity in the current period. This threatening deviation, while not credible, would drive the first firm to try to retain more of her customers this period because she will not get as many customers two periods later as a result of poor service provided by her competitor. These results demonstrate that while we arrive at a simple solution for a single period in isolation, a more complex structure exists within the multiperiod equilibrium.

Theorem 5 below presents some other characteristics of this equilibrium.

THEOREM 5. *If $\gamma_i < 1/2$ for each i , then under the equilibrium:*

- (a) *v_{it} and w_{it} are positive for all i and t ;*
- (b) *the optimal normalized capacity for a firm is:*
 - (b1) *increasing in the unit price p ,*
 - (b2) *increasing in $v_{i,t+1}$,*

- (b3) *increasing in the discount factor α , and*
- (b4) *decreasing in the capacity cost rate b ; and*
- (c) *the value per unit market share, v_{it} , is increasing (decreasing) in the period t if $V_i > (<) [p(1 - \beta_i h_i(y_{it}^*)) - by_{it}^*]/[1 - \alpha(1 - \gamma_i h_i(y_{it}^*) - \gamma_j h_j(y_{jt}^*))]$, for all i .*

Theorem 5(a) demonstrates that the firms in this market can be assured of a positive expected contribution in each period. Theorem 5(b) supports our intuition: Firms will provide better service when the margin increases, the marginal value of customers in the next period increases, and/or the effective cost of capital is smaller. Lastly, Theorem 5(c) provides conditions such that the value per unit market share is monotone over time and demonstrates that the direction of the relationship depends on the magnitude of the terminal value. Although we don't formally demonstrate it here, the terminal value that yields an equality in Part (c) must be the value of a unit of market share in the infinite horizon stationary problem. An intuitive explanation is that, for this special case, the value of a unit of market share would be constant for every finite-horizon problem, regardless of its length. Combining Parts (b2) and (c) of Theorem 5 provides conditions under which the level of service provided by each firm increases or decreases over time. For example, if both values of V_i are "large," customers can expect to receive better service as time passes. The requirements that $\gamma_i < \frac{1}{2}$ for each i play the role of sufficient stability criteria for the equilibrium.

Corollary 1 below summarizes the results for the case where the two firms are identical (a symmetric equilibrium).

COROLLARY 1. *If $V_i = V$, $W_i = W$, $\beta_i = \beta$, $\gamma_i = \gamma$ and $h_i = h$, for each i , a symmetric equilibrium results: the two firms make identical decisions in each period and Recursions (12) and (13) can be simplified:*

$$v_t = p(1 - \beta h(y_t^*)) - by_t^* + \alpha v_{t+1} (1 - 2\gamma h(y_t^*)), \text{ and}$$

$$w_t = \alpha w_{t+1} + \alpha v_{t+1} \gamma h(y_t^*).$$

Note that for this symmetric case, if both firms follow the equilibrium policy, customers can expect to find the same service level at each. Our simple model of customer behavior assumes switching will still occur in this instance. Note also that this ongoing switching

(even when it appears senseless) can be critical—without it, a flat-fee provider ($\beta = 0$) would have little motivation to provide capacity. Put another way, even if customers do not act strategically, but rather somewhat whimsically, and if the providers believe this behavior will be observed, then there is a clearly defined equilibrium in which the consumers are far from exploited. In a sense, this demonstrates that consumers need not act strategically to avoid being exploited. Note also that if the users could collude and agree to set γ amongst themselves, then they could get even better service.

4. Example 1: Competition Between Internet Service Providers

Our basic model is potentially applicable to a broad range of service delivery systems that can be approximated by loss-type queueing models. Potential objects of study include systems where (1) the number of expected customer service failures can be expressed as a function of the normalized capacity of the system, and (2) the capacity decision is the service rate of a server. The former is characteristic of loss systems, in which arrivals are “lost” with a probability that depends on the number of customers already in the system. Examples include:

(1) competition between providers of systems that can be approximated as $M/M/1/\infty^1$ queues where a “service failure” results from a customer finding too long of a queue upon arrival and balking (refusing to queue) as a result;

(2) competition between providers of systems that can be approximated as $M/M/1/K$ queues where a “service failure” results from an arriving customer being blocked when there are K customers already in the system; and

(3) competition between providers of systems that can be approximated as $M/G/c/c$ queues where a “service failure” results from an arriving customer being blocked when all c servers are concurrently busy.

In this section, we focus on a version of Example (3):

We approximate the blocking probability of the systems operated by ISPs by the blocking probability of an $M/G/1/1$ system. We also assume that the two competitors charge a fixed and equal fee per month to each subscriber. Service failures are assumed to result from service denial: Customers of an ISP place calls² to the ISP’s system and are either admitted to the system and immediately proceed to be served, or they are denied admittance due to congestion.

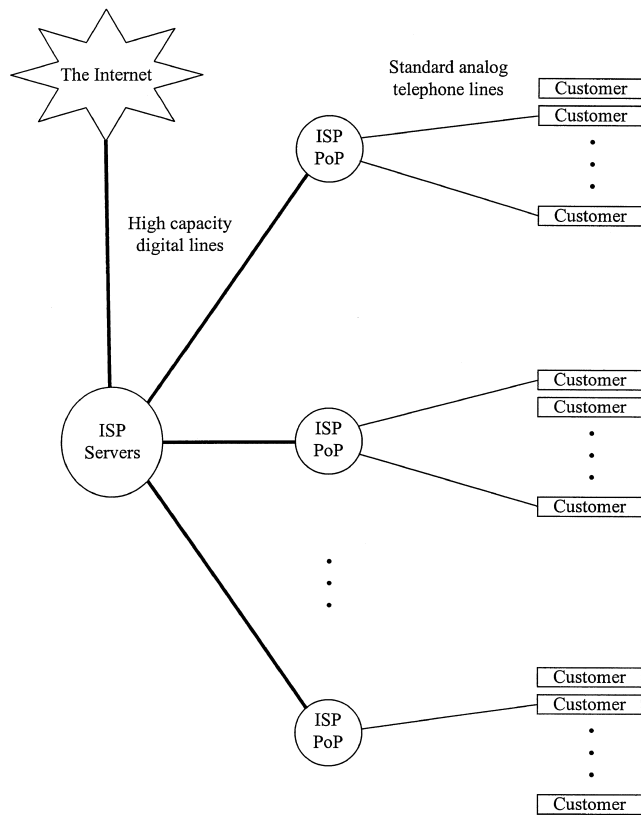
The blocking of customers’ calls to their ISPs that we discuss in this example is a significant issue for both ISPs and their customers. Inverse Network Technologies, Inc. collects and compiles data on ISP performance. For June, 1997, this organization reported an average call failure rate for 13 national ISPs of 12.6% during the hours of 6:00 p.m. to midnight. America Online, the dominant American ISP, was reported to have an evening call failure rate of 34.7% during this period (Inverse Network Technology 1997). These blocking rates have been discussed extensively in the popular press (see, e.g., Barrett et al. 1997, Kornblum 1997, and Tinnirello 1997). Pompili (1997) conducted an informal survey of ISP users during this same period. Users were asked whether they were considering switching from their current service provider to another service provider; the results indicated that anywhere from 57% (for America Online) to 13% (for MCI Internet) were considering such a move. Of course, these statistics reflect intent, not behavior. Wetzel (1997) details a survey of ISP customers that ranks service reliability as the most important factor cited by customers in choosing an ISP.

The growth in general popularity of the Internet in the 1990s has parallel the growth of the ISP industry. ISPs act as gateways to the Internet by providing dial-up access for both individual consumers and businesses, in addition to other related services. As of August, 1997, there were over 4,000 ISPs operating in the United States, with the geographic scope of individual ISPs ranging from municipal to international. The physical network of an ISP can be broken down into a few basic components (see Figure 1). First, an ISP has

¹Throughout we use this standard queueing systems taxonomy (see, e.g., Gross and Harris 1985). Thus, here we assume a Markovian (Poisson) call arrival process, a Markovian (exponential) service time distribution, one server, and an infinite arrival waiting capacity.

²Throughout this section, we refer to attempts to enter the system as “calls” and the subset of these calls that are successful in gaining admission as “visits.”

Figure 1 Simplified Schematic of an ISP Network



customers who connect to the ISP, using either an analog or digital modem, by way of the ISP's points-of-presence (PoPs), which are spread throughout the ISP's geographic service area. This customer-to-PoP connection is generally by way of a local telephone call, although this link can take other forms. The PoPs are linked to local telephone company (telco) switching centers and linked to the ISP's server computers by ISP-leased or ISP-owned high-capacity digital lines. The ISP's servers are then linked with the rest of the Internet by way of additional high-capacity digital lines.

Here, we focus on the dial-up connections that link the ISP's customers to the ISP's network. The concept of blocking discussed in this example is manifested in the busy signal that users experience when their ISP's local PoP is at full utilization, i.e., when all dial-up lines and modems are in use. The ISP's PoP capacity decision involves three costs: the charge to the ISP by the

local phone company for providing access to the local telephone network, the cost of modems, and the cost of the leased or owned links that connect together the telco, the PoP, and the ISP's servers.

The following additional assumptions support our analysis of this example. In this section, we assume that (A1)–(A15) hold.

(A12). *The overall customer arrival process at each firm is a Poisson process with rate equal to firm market share.*

We believe this assumption is reasonable. For instance, under certain regularity conditions, the limit of a superposition of many general (but rare event) interarrival time distributions, which here represent the interarrival time distributions of the individual customers' calls, yields a Poisson arrival process (Theorem 6.7 in Serfozo 1990).

(A13). *The system blocking probability can be approximated by that of an M/G/1/1 loss system.*

The firm market share, λ_{it} , is here also the (Poisson) arrival rate for firm i in month t . Such scaling entails no loss of generality. The capacity decision, μ_{it} , is here the chosen potential service rate (the usual mean rate at which customers can be served while the server is busy, also the inverse of the mean service time for a single customer). The M/G/1/1 queueing model is characterized by denial of queueing—any visit blocks all arriving calls from entering. The reality of a typical ISP's system is vastly more complex: An ISP can serve many customers simultaneously by sharing computing and communications capacity among users and by supplying a multitude of PoP dial-in telephone lines that allow users access to the system. Actual ISP capacity is a function of communications capacity, computing capacity, and the number of PoP dial-in lines. Without delving into the problem at that level, we have chosen to assume that each ISP has a single server that can only serve a single customer at a time. Capacity in this simplified system is therefore the expected number of visits per month if the facility is never idle. One can roughly interpret this quantity as the total number of customers that the ISP can serve simultaneously divided by the actual mean service time of a single customer. Of course, the blocking probability

that we obtain from this model will be an approximation at best. While we do not argue that our results have operational fidelity, we do believe they are insightful.

(A14). *Steady-state, infinite population results are valid.*

For large values of service rate and arrival rate, the system should settle down to steady-state behavior relatively quickly at the beginning of each month (see, e.g., Gross and Harris 1985 for a discussion of transient analysis of queueing models). We assume the customer groups that we study are large enough to result in negligible error under approximation by infinite size groups. Note that, by (A13), we have assumed the mean arrival rate at each firm is between zero and one arrivals per month. This was done strictly to simplify the exposition; the results scale up in an obvious way.

(A15). *Blocked calls are lost.*

This assumption is equivalent to saying that blocked calls are not resubmitted to the system. In reality, many blocked calls will be backlogged and resubmitted later, perhaps until the call is accepted by the system. We leave consideration of this effect to future work.

Consistent with this example, we make the following assignments:

$$\beta_i = 0, \text{ for all } i, \text{ and} \quad (14)$$

$$h_i(y_{it}) = \frac{1}{1 + y_{it}}, \text{ for all } i \text{ and } t. \quad (15)$$

Expression (14) is consistent with the common industry practice of charging a flat fee per month of service. Expression (15) is the blocking probability for an M/G/1/1 queueing system (see, e.g., Gross and Harris 1985). This form for $h_i(y_{it})$ is an interesting object of study, both because it fits our model of an ISP system and because it is one of the simplest functional forms that satisfies our requirements on $h_i(y_{it})$. It is easily verified that the functional form given by Expression (15) satisfies the requirements of (A7). Under these conditions, the equation defining y_{it}^* can be written as:

$$y_{it}^* = \max\left(\sqrt{\frac{\alpha v_{i,t+1} \gamma_i}{b}} - 1, 0\right).$$

The expressions for v_{it} and w_{it} are given by:

$$v_{it} = p - b y_{it}^* + \alpha v_{i,t+1} \left(1 - \frac{\gamma_i}{1 + \gamma_{it}^*} - \frac{\gamma_j}{1 + \gamma_{jt}^*}\right)$$

and

$$w_{it} = \alpha w_{i,t+1} + \frac{\alpha v_{i,t+1} \gamma_j}{1 + y_{jt}^*}.$$

Notable in this solution is that while y_{it}^* is superficially independent of p , which results from our assumption of flat-fee service ($\beta_i = 0$), it does depend on p through $v_{i,t+1}$.

We claim that one potential use of this model is to value customers. For example, this issue is relevant to the financial valuation of an ISP, i.e., what value should be assigned to the firm's subscribers? Within the context of our modeling assumptions, we can perform this valuation. What follows is a simple analysis of two competing ISPs.

For the sake of realism, we work directly with the raw market size, rather than the scaled version. The total market consists of 1,000 customers, each of whom places ten calls per month on average, yielding a total market of 10,000 (expected) call arrivals per month. We assume: $\alpha = 0.985$ (approximately equivalent to 20% annual discount rate); $b = \$0.10$ per visit per month; $p = \$2$ per call per month; $\gamma_1 = 0.2$; $\gamma_2 = 0.4$; $V_1 = V_2 = \$13$ per call; $W_1 = W_2 = \$400,000$. Thus, we (somewhat arbitrarily) assume that the cost to support one additional visit per month is \$0.10 and the monthly revenue per customer is \$20 (\$2 multiplied by 10 calls). We assume that Firm 1's customers are less sensitive to service failures than those of Firm 2, perhaps due to some auxiliary services provided by Firm 1 that make its customers more loyal. The model results for both firms at the start of a five year horizon are presented in Table 1.

Several characteristics of this solution are worth noting. First, the firm with less sensitive customers (Firm 1) will provide less capacity per customer, leading to

Table 1 Results for ISP Example

	y_i	$h_i(y_i)$	$\gamma_i h_i(y_i)$	v_i	w_i
Firm 1	4.37	0.186	0.0372	\$14.60	\$479,000
Firm 2	6.14	0.140	0.0560	\$13.00	\$351,000

a higher call failure rate (18.6% vs. 14.0%) but fewer defections (3.7% vs. 5.6%). In the context of our model, the firm with a lower defection rate will enjoy a larger market share if the defection imbalance persists. Second, a customer (10 arrivals per month) is worth \$146 at Firm 1 and \$130 at Firm 2. Third, a majority of the value of each firm is independent of the size of its current market share—it's valuable for a firm to participate in this market in order to capture its competitor's disaffected customers. This results, in part, from our implicit modeling assumption that a firm captures its competitor's disaffected customers at no cost. Lastly, the firm that is better able to retain customers who have experienced service failure will earn higher returns than its competitor for a given market share level. In fact, Firm 2 will not equal or exceed Firm 1 in value unless it has 99.3% or more of the total market. The results of this analysis are supported by anecdotal evidence from the ISP industry. America Online is known to provide a user-friendly human interface via proprietary software and access to proprietary databases, in addition to standard Internet access (see, e.g., Pompili 1997). These characteristics would be likely make an America Online subscriber less sensitive to service failures than a subscriber to an alternate service. Data presented earlier in this section indicate that America Online provides service levels below the industry average while simultaneously holding the largest market share. According to one review of ISP services, "Though AOL enjoys the infamous reputation of serving up more busy signals than its competitors, it is equally famous for providing one of the simplest interfaces." (Pompili 1997).

A simple valuation of the two firms can also be performed. When the two firms possess equal market shares, their total present values are \$552,000 and \$416,000, for Firms 1 and 2, respectively. Here, the value of Firm 1 is approximately 33% greater than that of Firm 2 due to the lower sensitivity of Firm 1's customers to service failures. A simplistic valuation of the two firms might overlook this important driver of firm value. For example, one might be tempted to value the firms by evaluating an annuity based on the current period cash flow. Assuming the two firms have settled upon their optimal normalized capacities, the current period cash flows are \$7,815 and \$6,930 for Firms 1 and

2, respectively (again, assuming equal market shares). Annuities that pay this amount each period have a net present value of \$521,000 and \$462,000, respectively. By our model, Firm 1 is worth more than this simple analysis would suggest, while Firm 2 is worth less. The key to this difference is that the equal market shares are not sustainable in the long run, as the naïve annuity analysis assumes.

5. Example 2: Product Availability Competition

In this section, we demonstrate that our basic model has application to service delivery systems that are inventory-based. We focus on a firm's ability to fill customer orders from inventory, where inventory level is itself the periodic capacity decision faced by the firm. Our basic model is applicable to inventory-based systems to the extent that the system can be characterized by a service failure function of the form given in (A7). Examples include normally distributed demand where the standard deviation scales up linearly with the mean (covered below) and gamma distributed demand with constant shape parameter. We make a number of simplifying assumptions so that the underlying inventory model does not obscure the effects under study.

The following additional assumptions support our analysis of this example. In this section, we assume that (A1)–(A11) and (A16)–(A18) hold.

(A16). *No inventory is held between periods and unmet demand is not backlogged.*

This assumption allows us to reduce the problem state space to a single dimension, and is consistent with our aim of developing a simple model. The first part of this assumption will be met literally if the product is perishable with a lifetime of one period and can be met implicitly in other cases. For example, if the retailer is allowed to return unsold stock at the end of the period for full cost reimbursement, then this assumption is met (in this case, our formulation would need to be altered slightly to include the financial holding costs of the inventory). Similarly, the retailer could hold any

excess stock and reduce the purchase made in the following period by that amount. If leftover stock exceeded the amount desired in the next period, that excess amount would be returned for full cost reimbursement. Note that this last condition is superfluous in some instances. Theorem 5 presents conditions under which service levels increase overtime. Under these conditions, stock would only be returned if market share had declined, but market share could not have declined if stock was leftover, and thus there would never be excess stock to return. The assumption that unmet demand is not backlogged could be met in circumstances where one or more inferior substitutes for the desired good will always be available and one of these will be chosen by customers facing a stockout of the desired good.

(A17). *A firm's expected demand level is equal to its market share.*

Letting D_{it} be Firm i 's random demand in month t , we have $E(D_{it}) = \lambda_{it}$. This is merely a scaling convention.

(A18). *The random demand in each month for each firm is normally distributed with a coefficient of variation, c_v , that is constant across all potential market sizes (demand levels), the same for all firms, and small enough that the probability of negative demand is negligible.*

The standard deviation of the demand distribution thus scales proportionally with the mean, which is consistent with highly correlated demands across the individual consumers in a market. This rules out capacity pooling effects, consistent with (A7). We use $\Phi(\cdot)$, $\varphi(\cdot)$, and $I_N(\cdot)$ to represent the standard normal cumulative distribution, the standard normal density, and the standard normal loss function, respectively.

Our assumption of normal demand allows us to use relatively simple expressions in our analysis. By (A17), y_{it} is the ratio of capacity choice (here inventory stocking level) to expected monthly demand. Thus, $(y_{it} - 1)/c_v$ expresses the inventory stocking level in standard deviations from the mean demand. Using this expression, the expected unsatisfied demand for Firm i in month t , given an inventory stocking level of μ_{it} , can be written in terms of the standard normal loss function, $I_N(\cdot)$:

$$E(D_{it} - \mu_{it})^+ = c_v \lambda_{it} I_N \left(\frac{y_{it} - 1}{c_v} \right).$$

(See, for example, Porteus 1990 for details). In this model, unsatisfied demand is the root cause of service failure, and consistent with this we define:

$$\beta_i = 1, \text{ for all } i, \text{ and} \quad (16)$$

$$\lambda_{it} h_i(y_{it}) = c_v \lambda_{it} I_N \left(\frac{y_{it} - 1}{c_v} \right), \text{ for all } i \text{ and } t. \quad (17)$$

Expression (16) follows from (A16) because unfilled orders result in lost sales. Expression (17) equates the expected absolute number of service failures with the expected amount by which demand exceeds stock. It is easily verified that the functional form in (17) satisfies the requirements of (A7). The parameters γ_i here capture customer sensitivity to unfilled orders and serve to differentiate the two firms. By (16) and (17), the equations for y_{it}^* can be written:

if

$$\Phi \left(\frac{-1}{c_v} \right) \leq 1 - \frac{b}{p + \alpha v_{i,t+1} \gamma_i},$$

then

$$\Phi \left(\frac{y_{it}^0 - 1}{c_v} \right) = 1 - \frac{b}{p + \alpha v_{i,t+1} \gamma_i}, \quad (18)$$

$$y_{it}^* = \begin{cases} y_{it}^0 & \text{if } \Phi \left(\frac{-1}{c_v} \right) \leq 1 - \frac{b}{p + \alpha v_{i,t+1} \gamma_i} \\ 0 & \text{otherwise;} \end{cases}$$

The expressions for v_{it} and w_{it} are given by:

$$v_{it} = p \left(1 - c_v I_N \left(\frac{y_{it}^* - 1}{c_v} \right) \right) - b y_{it}^* + \alpha v_{i,t+1} \left(1 - \gamma_i c_v I_N \left(\frac{y_{it}^* - 1}{c_v} \right) - \gamma_j c_v I_N \left(\frac{y_{jt}^* - 1}{c_v} \right) \right) \quad (19)$$

and

$$w_{it} = \alpha w_{i,t+1} + \alpha v_{i,t+1} \gamma_j c_v I_N \left(\frac{y_{jt}^* - 1}{c_v} \right).$$

Note that (18) can be written in a standard critical fractile form:

$$y_{it}^0 = 1 + c_v \Phi^{-1} \left(\frac{p + \alpha v_{i,t+1} \gamma_i - b}{p + \alpha v_{i,t+1} \gamma_i} \right),$$

with unit underage cost $p + \alpha v_{i,t+1} \gamma_i - b$ and unit overage cost b . The underage cost consists of the immediate lost contribution, $p - b$, and a lost goodwill cost of $\alpha v_{i,t+1} \gamma_i$. The lost goodwill cost is the present value of an additional "customer" next month multiplied by the fraction of that customer that the firm expects to lose if that customer experiences a stockout, i.e., the expected cost (above the immediate lost contribution) of stocking one unit less than realized demand.

Our model is, to our knowledge, the first to endogenously evaluate the present value of the future consequences of current lost sales in a competitive environment. Instead of taking the order fill rate as exogenous (with an implied shortage cost), our model yields an equilibrium shortage cost and order fill rate based on simple cost parameters and a model of firm competition. That is, a lost sale results not only in an immediate loss of revenue, but in reduced future demand, due to the possibility of switching. Our results are reassuring in that the traditional approach of assigning a proportional shortage cost, as exemplified by Oral et al. (1972), is validated by our model. Indeed, (19) can be used to explicitly calculate the expected present value of a lost sale this period.

As with our previous example, we provide a short description of a numerical analysis. We assume that two retailers compete to supply a particular premium brand of bread to a base of customers. Retailer 1 has customers that are less sensitive to stockouts, perhaps because she is located in a more accessible location. Bread arrives from the bakery once per week, and any premium bread that remains unsold after one week is thrown away. In addition, if a retailer stocks out of the brand under study and a customer arrives with the intent to buy this brand, then there is an ample amount of a less preferred bread that the customer will purchase as a temporary substitute. However, the customer will consider such an occurrence to be a service failure that may provoke a switch in loyalty.

We assume a market size of 1,000 customers, each of whom demands one loaf of this premium bread per week on average, yielding a total expected demand of 1,000 loaves per week. We assume: $\alpha = 0.9965$ per week (approximately equivalent to 20% annual discount rate); $b = \$1.40$; $p = \$2$; $\gamma_1 = 0.25$; $\gamma_2 = 0.50$; V_1

$= V_2 = \$0$ per demand; $W_1 = W_2 = \$0$; and $c_v = 0.3$. Thus, we assume that the price the bakery charges is \$1.40 per loaf and the price charged to customers is \$2 per loaf. Here we assume the end-of-horizon values are zero. The model results for both firms at the start of a five year horizon are presented in Table 2, along with the results of a myopic newsvendor formulation based on the above costs.

We note a few characteristics of the solution. First, the retailer with less sensitive customers (Retailer 1) will provide less capacity per customer (0.23 vs. 0.50 standard deviations above the mean), leading to a higher expected fraction of customers facing service failures (8.69% vs. 6.05%) but fewer defections (2.17% vs. 3.03%). In the context of our model, the retailer with a lower defection rate will enjoy a larger market share if the defection imbalance persists. Second, a customer is worth \$5.87 at Retailer 1 and \$4.96 at Retailer 2. Third, the imputed lost goodwill costs of a stockout are significantly different for the two firms: \$1.46 and \$2.47 at Retailers 1 and 2, respectively. The myopic underage cost for both firms is \$0.60, which represents the immediate lost contribution on one sale. Retailer 2's larger imputed lost goodwill cost is due, in part, to Retailer 1's smaller defection rate—Retailer 2 must be more careful about losing customers as it doesn't enjoy such a (proportionally) large customer inflow, i.e., Retailer 2 is not as likely to get lost customers back quickly. Lastly, the retailer with less sensitive customers will expect to earn higher returns than its competitor for any current market share split.

6. Extension to More Than Two Competitors

In this section, we discuss a generalization of our basic model to an arbitrary number of competitors. Given n competitors in the market, the state of the system at

Table 2 Results for Inventory Example

	y_i	$h_i(y_i)$	$\gamma_i h_i(y_i)$	v_i	$\alpha v_i \gamma_i$	W_i
Retailer 1	1.07	0.0869	0.0217	\$5.87	\$1.46	\$29,400
Retailer 2	1.15	0.0605	0.0303	\$4.96	\$2.47	\$18,000
Myopic	0.843	0.214	—	—	—	—

any point in time can be described by a vector of dimension $n - 1$, assuming (as we do) that the total size of the market is known to all. The linear solution structure shown to hold for the two-competitor case will carry through for the n firm case given an additional assumption on the allocation of lost customers among the firms. This assumption is presented as (A19). Without loss of generality, we let the market share of firm n be determined.

(A19). *The customers lost by each firm are distributed over the other firms in a manner that is independent of the system state. A customer lost by firm i becomes a customer of firm j next period with probability a_{ij} , where $\sum_{j \neq i} a_{ij} = 1$ for all i .*

By (A19), we can modify Expression (1) as follows:

$$E(\lambda_{i,t+1} | \lambda_{1,t}, \lambda_{2,t}, \dots, \lambda_{n-1,t}) = \lambda_{it} - \lambda_{it} \gamma_i h_i(y_{it}) + \sum_{\substack{j=1 \\ j \neq i}}^{n-1} a_{ji} \lambda_{jt} \gamma_j h_j(y_{jt}) + a_{ni} \left(1 - \sum_{j=1}^{n-1} \lambda_{jt} \right) \gamma_n h_n(y_{nt}).$$

All other assumptions take on their natural meaning for $n > 2$ firms. However, we note that the informational assumptions (A11) are no longer transparent; we will assume here that all competitors know (or can estimate with negligible error) the market shares and other relevant parameters of all competitors. We leave study of equilibrium under less complete informational assumptions to future work. It is necessary to define the following generalization of our v_{it} quantities: v_{ijt} = value to firm i of customers at firm j (valued relative to customers at Firm n) in period t . We require the additional subscript j because a firm may value the customers of its competitors differently due to differing values of a_{ji} and γ_j . For example, it is reasonable that customers who patronize a nearby competitor that has a history of losing its customers to a firm would be more valuable to that firm than those of a competitor who rarely loses its customers.

THEOREM 6. *If (A1)–(A11) and (A19) hold, then the equilibrium normalized capacity for firm i in any month t will be independent of the vector λ and the expected present value of the returns to firm i over periods t through T will be a linear function of the components of λ .*

We leave to future work the study of additional characteristics of this $n > 2$ firm model, including the development of an analog of our Theorem 5 for two firms. Recall that the sensitivity results of Theorem 5 rely upon the sufficient stability criteria $\gamma_i < \frac{1}{2}$ and $\gamma_j < \frac{1}{2}$. We note, however, that for special cases of the $n > 2$ firm model, such as identical firms, the analogous stability conditions are easily established.

7. Conclusion

We have investigated an explicitly dynamic model of firm behavior in which firms compete based upon customer service. Our assumptions allow us to formulate a model that implicitly measures the value of a customer over a multiperiod horizon, uses this value as a measure of the cost of a customer defection, and yields an optimal (equilibrium) capacity choice that reflects the possibility of defection in the event of service failure. We have provided conditions under which the equilibrium solution takes on a simple, intuitive form: Capacity levels scale directly and linearly in the number of customers being served. We have illustrated the use of our model in two contexts: competition between Internet service providers and inventory availability competition. For each example, we have detailed the functional form of the solution, provided a numerical example, and discussed the managerial implications.

The dynamic game that we construct is a game played by the firms in the market, not by the consumers. Our assumption (A9) details an exogenous “rule” of expected aggregate customer behavior which leads to the expected state transition given by Expressions (1) and (2). This assumption deserves further discussion. In a multiperiod problem, Expression (1) implies that customers may switch between the service providers multiple times during the problem horizon. Implicitly, we have assumed that customers are purely reactive in their switching behavior, with no memory of experiences prior to the current period and no information on the relative performance of the firms. Because this behavior is clearly unrepresentative of human behavior in general, we offer some arguments in support of our assumption. (1) We are able to obtain an explicit solution to the simple model that results. Such a solution provides a baseline that will aid in the

understanding of more realistic models as they are developed in the future. For example, it will be interesting to compare the level of service received under our model of consumer behavior with the level of service provided when consumers act more rationally. (2) The work of Keaveney (1995) supports the idea that service failures (rather than rational comparisons of firms) are a primary driver of switching behavior, although this work is silent on how customers transition from one supplier to another. (3) The form of our transition function is more palatable when a small fraction of the customers defects in any one period. In a stationary setting, the average tenure of a customer would be given by the mean of a geometric random variable, e.g., if 5% of customers defect each month, the average tenure is 20 months. In reality, it seems reasonable that a customer might switch back to a previous provider after such a period. (4) Section 6 demonstrates that our model can be extended to an arbitrary number of firms. With many firms in the market, intuition tells us that comparative information on the firms will not be as readily available or as simple to process, particularly in industries undergoing rapid change. It is also less likely that customers will switch to the same firm multiple times in quick succession. (5) This assumption may provide an approximately correct representation of an amalgam of more sophisticated behaviors. There may be so many different rational approaches used by consumers in practice that the mode we assume is as accurate as one based on a more sophisticated model of consumer choice. Of course, to our knowledge, the opposite may also be true. (6) The model form is amenable to aggregate empirical analysis.

The model we present is simple in many respects. Many potentially relevant issues are ignored in order to isolate some of the effects of customer service competition. Many of these open issues present opportunities for future work to either extend the results presented here or pursue alternate model formulations. Expanding upon our simple model of customer behavior presents one such opportunity. Among the other possibilities are: modeling different customer types or classes that exhibit different responses to service failure, modeling concurrent price and service competition, studying availability competition for general demand distributions over multiple periods, modeling

capacity inflexibility, and incorporating a firm's ability to directly attract competitors' customers and customers' ability to enter and exit the market.³

Appendix

Notation and Symbols

Exogenous parameters:

- α = Discount factor for one period, $0 < \alpha \leq 1$;
- b = Unit cost of service capacity per period;
- p = Price charged per unit market share per period;
- β_i = Fraction of service failures at firm i that result in no revenue in the current period;
- γ_i = Fraction of service failures at firm i that result in customer switching next period;
- T = Final decision period (period $T + 1$ represents "end of horizon" and is considered to occur at the end of period T);
- V_i = End of horizon value per customer;
- W_i = End of horizon value that is independent of the size of the customer base.

Decision variables:

- μ_{it} = Capacity selected by firm i in period t ;
- $\mu_{it} = y_{it} \lambda_{it}$, where y_{it} is the normalized capacity selected by firm i in period t .

Computed quantities:

- λ_{it} = Market share for firm i in period t ;
- $\lambda_{it} l_t(y_{it})$ = Expected number of customers who experience service failures at firm i in period t ;
- ν_{it} = Equilibrium value of the entire potential customer base for firm i in period t ;
- ν_{ijt} = Value to firm i of customers at Firm j (valued relative to customers at firm n) in period t (used in extension to $n > 2$ firms);
- w_{it} = Equilibrium value independent of size of customer base for firm i in period t ;
- $g_{it}(\lambda \mid \pi^i, \pi^j)$ = Expected present value to firm i over periods t through T given that it starts period t in state λ , and given that firms i and j follow strategies π^i and π^j over the time horizon, respectively (see Porteus 1982 for details on strategies for dynamic problems);
- $\pi_t^i(\lambda)$ = Action (i.e., normalized capacity decision) called for by strategy π^i for firm i when starting period t facing state λ ;

Proofs

³We are grateful for the comments and suggestions of an anonymous M&SOM senior editor and two anonymous reviewers. Valuable comments were also provided by Sunil Kumar, James Patell, J. Michael Harrison, and seminar participants at Stanford University. Financial support from the Stanford Integrated Manufacturing Association is gratefully acknowledged.

PROOF OF THEOREM 2. We first establish that (4) is supermodular with respect to (γ, γ) . Let C denote the set of allowed values (γ, y) , i.e.,

$$C = \{(\gamma, y) \mid \gamma \in (0, 1], y \in [0, \infty)\}.$$

Clearly C is lattice, as required by the theory. Let g_i represent our objective function for firm i :

$$g_i(\gamma_i, y_i) = p(1 - \beta_i h_i(y_i)) - b y_i + \alpha V_i(1 - \gamma_i h_i(y_i) - \gamma_j h_j(y_j)).$$

Supermodularity requires that $g_i(x^1 \wedge x^2) + g_i(x^1 \vee x^2) \geq g_i(x^1) + g_i(x^2)$. We focus on only the nontrivial cases: without loss of generality, assume $x^1 \wedge x^2 = (\gamma^2, y^1)$ and $x^1 \vee x^2 = (\gamma^1, y^2)$. Then,

$$\begin{aligned} & g(x^1 \wedge x^2) + g(x^1 \vee x^2) - g(x^1) - g(x^2) \\ &= -\alpha V_i(\gamma^2 h_i(y^1) + \gamma^1 h_i(y^2) - \gamma^1 h_i(y^1) - \gamma^2 h_i(y^2)) \\ &= -\alpha V_i((\gamma^2 - \gamma^1)(h_i(y^1) - h_i(y^2))) \geq 0. \end{aligned} \quad (1A)$$

Expression (1A) follows since h_i is decreasing by (A7). Thus g_i is supermodular in (γ_i, y_i) . Because the equilibrium solution, y_i^* , of Theorem 1 for firm i is the maximizer of $g_i(\gamma_i, y_i)$ over y_i , it follows from Topkis (1978) (see also Milgrom and Roberts (1992)) that y_i^* is increasing in γ_i . The proof is similar for the (β_i, y_i) case. \square

Lemma 1 below is an extension of Theorem 1 (solution to one-period problem) to encompass an arbitrary single period within a multiperiod problem. Lemma 1 is applied below to prove Theorem 3 (solution to multiperiod problem) by induction.

LEMMA 1. Suppose that $g_{i,t+1}(\lambda \mid \pi^i, \pi^j) = v_{i,t+1}\lambda + w_{i,t+1}$ for each i , where v_{it} and w_{it} are of indeterminate sign and independent of λ . Suppose also that the firms are interested in maximizing their value per customer, as in Problem (5).

- (a) There exists a unique Nash equilibrium in period t for each state λ . If firm i finds itself at the beginning of period t with a market share of λ_{it} , then the following hold:
- (b) (Attainment) The optimal normalized capacity to be provided by firm i , y_{it}^* , is given uniquely as follows:

$$(p\beta_i + \alpha v_{i,t+1}\gamma_i)h_i'(0) \leq -b, \quad (9)$$

then y_{it}^0 is the unique solution to

$$(p\beta_i + \alpha v_{i,t+1}\gamma_i)h_i'(y_{it}) = -b; \quad (10)$$

$$y_{it}^* = \begin{cases} y_{it}^0 & \text{if } (p\beta_i + \alpha v_{i,t+1}\gamma_i)h_i'(0) \leq -b \\ 0 & \text{otherwise;} \end{cases} \quad (11)$$

and is thus independent of both the state λ_{it} and the competitor's decision y_{jt} .

- (c) (Preservation) The resulting optimal expected return to firm i (over periods t through T) can be written as $v_{it}\lambda_{it} + w_{it}$, where v_{it} and w_{it} are indeterminate in sign, independent of λ_{it} and satisfy:

$$\begin{aligned} v_{it} &= p(1 - \beta_i h_i(y_{it}^*)) - b y_{it}^* + \\ &\quad \alpha v_{i,t+1}(1 - \gamma_i h_i(y_{it}^*) - \gamma_j h_j(y_{jt}^*)), \end{aligned} \quad (12)$$

and

$$w_{it} = \alpha w_{i,t+1} + \alpha v_{i,t+1}\gamma_j h_j(y_{jt}^*). \quad (13)$$

PROOF OF LEMMA 1. Lemma 1 does not follow directly from Theorem 1 since $v_{i,t+1}$ is indeterminate in sign and the maximization problem is thus, in general, not concave. However, the problem is strictly concave if $p\beta_i + \alpha v_{i,t+1}\gamma_i > 0$, since we have:

$$\frac{\partial^2 g_{it}}{\partial y_{it}^2} = -(p\beta_i + \alpha v_{i,t+1}\gamma_i)h_i''(y_{it}),$$

with $h_i(y_{it})$ strictly convex by (A7). Now, if the objective function is not strictly concave (i.e., $p\beta_i + \alpha v_{i,t+1}\gamma_i \leq 0$), then the objective function is also strictly decreasing in y_{it} since

$$\frac{\partial g_{it}}{\partial y_{it}} = -(p\beta_i + \alpha v_{i,t+1}\gamma_i)h_i'(y_{it}) - b,$$

and $h_i(y_{it})$ is decreasing by (A7). Hence, the unique optimal normalized capacity in this case is $y_{it}^* = 0$. If the objective function is strictly concave, application of Theorem 1 gives the optimal normalized capacity y_{it}^* . Because Condition (9) is sufficient to establish strict concavity of the problem, Part (b) is proved. Part (a) follows from the definition of a Nash equilibrium since Part (b) demonstrates that each firm's unique contribution maximizing action is independent of the action of the other firm—thus (y_{1t}^*, y_{2t}^*) forms the unique Nash equilibrium.

- (c) Recall the recursion for the value function:

$$\begin{aligned} g_{it}(\lambda_{it} \mid \pi^i, \pi^j) &= \lambda_{it}(p(1 - \beta_i h_i(y_{it})) - b y_{it}) + \\ &\quad \alpha E(g_{i,t+1}(\lambda_{i,t+1} \mid \pi^i, \pi^j)). \end{aligned} \quad (8)$$

Given that $g_{i,t+1}(\lambda \mid \pi^i, \pi^j) = v_{i,t+1}\lambda + w_{i,t+1}$, we can write (8) as:

$$\begin{aligned} g_{it}(\lambda_{it} \mid \pi^i, \pi^j) &= \lambda_{it}(p(1 - \beta_i h_i(y_{it})) - b y_{it}) + \\ &\quad \alpha v_{i,t+1} E(\lambda_{i,t+1}) + \alpha w_{i,t+1} \\ &= \lambda_{it}(p(1 - \beta_i h_i(y_{it})) - b y_{it}) + \\ &\quad \alpha v_{i,t+1} (\lambda_{it}(1 - \gamma_i h_i(y_{it}) - \gamma_j h_j(y_{jt})) + \gamma_j h_j(y_{jt})) + \alpha w_{i,t+1} \end{aligned} \quad (3A)$$

$$\begin{aligned} &= \lambda_{it}(p(1 - \beta_i h_i(y_{it})) - b y_{it}) + \\ &\quad \alpha v_{i,t+1}(1 - \gamma_i h_i(y_{it}) - \gamma_j h_j(y_{jt})) + \\ &\quad \alpha v_{i,t+1}\gamma_j h_j(y_{jt}) + \alpha w_{i,t+1}. \end{aligned} \quad (4A)$$

Expression (3A) follows using (2), and (4A) follows by combining terms. Given that y_{it}^* and y_{jt}^* are independent of λ_{it} , Expression (4A) is of the desired form. \square

PROOF OF THEOREM 3. By Definition 9.B.1 of Mas-Colell et al. (1995), the subproblem consisting of the final τ decision periods forms a subgame of the overall game, as all information sets are singletons. We refer to this subgame as the τ -subgame. Given V_i, V_j, W_i and W_j , by Lemma 1, (y_{1T}^*, y_{2T}^*) is the unique (subgame perfect) Nash equilibrium in the 1-subgame. Now, for the 2-subgame we apply Lemma 1 to construct a reduced form of the 2-subgame in

which the 1-subgame is replaced by the terminal values that result from play of the 1-subgame's unique subgame perfect strategy. By Lemma 1, (y_{1T-1}^*, y_{2T-2}^*) is the unique subgame perfect equilibrium of this reduced-form game. By Proposition 9.B.3 of Mas-Colell et al. (1995), $((y_{1T-1}^*, y_{2T-2}^*), (y_{1T}^*, y_{2T}^*))$ is the unique subgame perfect equilibrium of the 2-subgame. Let the induction hypothesis be that $((y_{1T-\tau+1}^*, y_{2T-\tau+1}^*), \dots, (y_{1T}^*, y_{2T}^*))$ is the unique subgame perfect equilibrium for the τ -subgame. We now construct a reduced form of the $\tau + 1$ -subgame in which the τ -subgame is replaced by the terminal values that result from play of the τ -subgame's unique subgame perfect strategy. By Lemma 1, $(y_{1T-\tau}^*, y_{2T-\tau}^*)$ is the unique subgame perfect equilibrium of this reduced-form game. By Proposition 9.B.3, $((y_{1T-\tau}^*, y_{2T-\tau}^*), (y_{1T-\tau+1}^*, y_{2T-\tau+1}^*), \dots, (y_{1T}^*, y_{2T}^*))$ is the unique subgame perfect equilibrium of the $(\tau + 1)$ -subgame. Thus, by induction, the equilibrium $((y_{11}^*, y_{21}^*), \dots, (y_{1T}^*, y_{2T}^*))$ is the unique subgame perfect equilibrium. The remaining elements of Theorem 3 follow directly from Lemma 1. \square

PROOF OF THEOREM 4. We present the proof for the case in which there is sufficient differentiability. See the proof of Theorem 2 above for an example of how a complete proof would proceed.
(a)–(b) Combining the first-order condition for y_{it}^* given in (10) with the recursion for $\nu_{i,t+1}$ given in (12), we have:

$$h'_i(y_{it}) = -b/[p\beta_i + \alpha\gamma_i(p(1 - \beta_i h_i(y_{i,t+1})) - by_{i,t+1} + \alpha\nu_{i,t+2}[1 - \gamma_i h_i(y_{i,t+1}) - \gamma_j h_j(y_{j,t+1})])].$$

Implicitly differentiating with respect to $y_{i,t+1}$ yields

$$\frac{\partial y_{it}}{\partial y_{i,t+1}} = \alpha b \gamma_i (- (p\beta_i + \alpha\nu_{i,t+2}\gamma_i) h'_i(y_{i,t+1}) - b) / [h''_i(y_{it})[p\beta_i + \alpha\gamma_i(p(1 - \beta_i h_i(y_{i,t+1})) - by_{i,t+1} + \alpha\nu_{i,t+2}[1 - \gamma_i h_i(y_{i,t+1}) - \gamma_j h_j(y_{j,t+1})])^2].$$

The sign of this expression depends entirely on

$$\frac{\partial g_{i,t+1}}{\partial y_{i,t+1}} = - (p\beta_i + \alpha\nu_{i,t+2}\gamma_i) h'_i(y_{i,t+1}) - b,$$

which is zero for the case $y_{i,t+1} = y_{i,t+1}^*$, positive if $y_{i,t+1}^* < y_{i,t+1}$, and negative if $y_{i,t+1}^* > y_{i,t+1}$.

(c) Similar to (a) and (b), we can establish

$$\frac{\partial y_{it}}{\partial y_{j,t+1}} = \alpha b \gamma_i (- \alpha\nu_{i,t+2}\gamma_j h'_j(y_{j,t+1})) / [h''_i(y_{it})[p\beta_i + \alpha\gamma_i(p(1 - \beta_i h_i(y_{i,t+1})) - by_{i,t+1} + \alpha\nu_{i,t+2}[1 - \gamma_i h_i(y_{i,t+1}) - \gamma_j h_j(y_{j,t+1})])^2],$$

which is unambiguously positive for all choices of $y_{j,t+1}$. \square

Lemma 2 establishes that under our sufficient stability criteria ($\gamma_i < \frac{1}{2}$ for all i) the value per unit market share is increasing in the unit price of the good or service. This result is applied in the proof of Theorem 5, Part (b1) to show that the optimal normalized capacity is increasing in the price p .

LEMMA 2. Given that $\gamma_i < \frac{1}{2}$ for all i , ν_{it} is increasing in p for all i and t .

PROOF OF LEMMA 2. Let $\nu_{it}(p)$ and $y_{it}^*(p)$ denote ν_{it} and y_{it}^* as functions of p . Suppose $p^1 \geq p^2$. In period T , because $y_{iT}^*(p^1)$ is at least as good as $y_{iT}^*(p^2)$ when the price is p^1 ,

$$\begin{aligned} \nu_{iT}(p^1) - \nu_{iT}(p^2) &\geq p^1(1 - \beta_i h_i(y_{iT}^*(p^2))) - by_{iT}^*(p^2) + \\ &\quad \alpha V_i(1 - \gamma_i h_i(y_{iT}^*(p^2)) - \gamma_j h_j(y_{jT}^*(p^2))) - p^2(1 - \beta_i h_i(y_{iT}^*(p^2))) + \\ &\quad by_{iT}^*(p^2) - \alpha V_i(1 - \gamma_i h_i(y_{iT}^*(p^2)) - \gamma_j h_j(y_{jT}^*(p^2))) \\ &= (p^1 - p^2)(1 - \beta_i h_i(y_{iT}^*(p^2))) \geq 0. \end{aligned}$$

This proves the base case. Now, assume $\nu_{i,t+1}(p^1) - \nu_{i,t+1}(p^2) \geq 0$. As above,

$$\begin{aligned} \nu_{it}(p^1) - \nu_{it}(p^2) &\geq p^1(1 - \beta_i h_i(y_{it}^*(p^2))) - by_{it}^*(p^2) + \\ &\quad \alpha\nu_{i,t+1}(p^1)(1 - \gamma_i h_i(y_{it}^*(p^2)) - \gamma_j h_j(y_{jt}^*(p^2))) - \\ &\quad p^2(1 - \beta_i h_i(y_{it}^*(p^2))) + by_{it}^*(p^2) - \\ &\quad \alpha\nu_{i,t+1}(p^2)(1 - \gamma_i h_i(y_{it}^*(p^2)) - \gamma_j h_j(y_{jt}^*(p^2))) \\ &= (p^1 - p^2)(1 - \beta_i h_i(y_{it}^*(p^2))) + \alpha(\nu_{i,t+1}(p^1) - \nu_{i,t+1}(p^2)) \cdot \\ &\quad (1 - \gamma_i h_i(y_{it}^*(p^2)) - \gamma_j h_j(y_{jt}^*(p^2))) \geq 0. \end{aligned} \quad (5A)$$

The inequality in (5A) follows from the induction hypothesis and the assumption that $\gamma_i < \frac{1}{2}$ for all i . Thus, by induction, the ν_{it} 's are increasing in p . \square

PROOF OF THEOREM 5.

(a) Recall

$$\nu_{it} = p(1 - \beta_i h_i(y_{it}^*)) - by_{it}^* + \alpha\nu_{i,t+1}(1 - \gamma_i h_i(y_{it}^*) - \gamma_j h_j(y_{jt}^*)), \quad (12)$$

$$w_{it} = \alpha w_{i,t+1} + \alpha\nu_{i,t+1}\gamma_j h_j(y_{jt}^*). \quad (13)$$

This proof proceeds by backwards induction on t . In the final decision period, $V_i > 0$ by (A10), and $\gamma_i < \frac{1}{2}$ for all i by assumption, thus the last term of (12) is positive for $t = T$ for all choices of y_{it} . A feasible solution to the problem is $y_{it} = 0$, for which the remaining terms of (12) are positive. As this is a maximization problem, we have constructed the bound $\nu_{iT} \geq 0$. Given $\nu_{i,t+1} \geq 0$, the same argument can be applied in period t , thus by backwards induction on t , $\nu_{it} \geq 0$ for all t . Now, given $W_i \geq 0$ by (A10), a similar induction argument can be constructed to conclude that $w_{it} \geq 0$ for all t .

(b1) Given $g_{i,t+1}(\lambda | \pi^t, \pi^t) = \nu_{i,t+1}\lambda + w_{i,t+1}$, we want to establish that (8) is supermodular with respect to (p, y) . Let C denote the lattice set of allowed values (p, y) , i.e.,

$$C = \{(p, y) | p \in (0, \infty], y \in [0, \infty)\}.$$

Supermodularity requires that $g_{it}(x^1 \wedge x^2) + g_{it}(x^1 \vee x^2) \geq g_{it}(x^1) + g_{it}(x^2)$. We focus on only the nontrivial cases: without loss of generality, assume $x^1 \wedge x^2 = (p^2, y^1)$ and $x^1 \vee x^2 = (p^1, y^2)$. In period t we have:

$$\begin{aligned} g_{it}(x^1 \wedge x^2) + g_{it}(x^1 \vee x^2) - g_{it}(x^1) - g_{it}(x^2) &= \lambda_{iT}(h_i(y^1) - \\ &\quad h_i(y^2))(\beta_i(p^1 - p^2) + \alpha\gamma_i(\nu_{i,t+1}(p^1) - \nu_{i,t+1}(p^2))) \geq 0. \end{aligned} \quad (6A)$$

The inequality in (6A) follows from Lemma 2. Thus g is super-modular in all periods t and y_{iT} is increasing in p by the same argument given for Theorem 2.

(b2) Proof is similar to (b1).

(b3) Proof is similar to (b1).

(b4) Proof is similar to (b1).

(c) We proceed by backwards induction. We first assume that

$$V_i > \frac{p(1 - \beta_i h_i(y_{iT}^*)) - b y_{iT}^*}{1 - \alpha(1 - \gamma_i h_i(y_{iT}^*) - \gamma_j h_j(y_{iT}^*))}, \text{ for all } i. \quad (7A)$$

By (12), in period T , we have

$$v_{iT} = p(1 - \beta_i h_i(y_{iT}^*)) - b y_{iT}^* + \alpha V_i(1 - \gamma_i h_i(y_{iT}^*) - \gamma_j h_j(y_{iT}^*)).$$

We can write out the difference:

$$V_i - v_{iT} = V_i(1 - \alpha(1 - \gamma_i h_i(y_{iT}^*) - \gamma_j h_j(y_{iT}^*))) - p(1 - \beta_i h_i(y_{iT}^*)) + b y_{iT}^*,$$

and conclude by (7A) that $V_i - v_{iT} > 0$ for all i . Now, assume $v_{i,t+2} - v_{i,t+1} > 0$ for all i . By Theorem 5(b2), $y_{i,t+1} \geq y_{it}$ and thus $h_i(y_{i,t+1}) \leq h_i(y_{it})$ for all i . We now form the difference $v_{i,t+1} - v_{it}$:

$$\begin{aligned} v_{i,t+1} - v_{it} &= p\beta_i(h_i(y_{it}^*) - h_i(y_{i,t+1}^*)) + b(y_{it}^* - y_{i,t+1}^*) + \\ &\quad \alpha v_{i,t+2}(1 - \gamma_i h_i(y_{i,t+1}^*) - \gamma_j h_j(y_{i,t+1}^*)) - \\ &\quad \alpha v_{i,t+1}(1 - \gamma_i h_i(y_{it}^*) - \gamma_j h_j(y_{it}^*)) \end{aligned}$$

Now form an inequality using the induction hypothesis:

$$\begin{aligned} v_{i,t+1} - v_{it} &> p\beta_i(h_i(y_{it}^*) - h_i(y_{i,t+1}^*)) - b(y_{i,t+1}^* - y_{it}^*) \\ &\quad + \alpha v_{i,t+2}(\gamma_i(h_i(y_{it}^*) - h_i(y_{i,t+1}^*)) + \gamma_j(h_j(y_{it}^*) - h_j(y_{i,t+1}^*))). \end{aligned}$$

The terms involving Firm j can now be dropped as their total is positive. If $y_{i,t+1} = y_{it}$, the result $v_{i,t+1} - v_{it} > 0$ follows immediately. Otherwise, divide through by $y_{i,t+1}^* - y_{it}^*$:

$$\frac{v_{i,t+1} - v_{it}}{y_{i,t+1}^* - y_{it}^*} > -(p\beta_i + \alpha v_{i,t+2}\gamma_i) \frac{(h_i(y_{i,t+1}^*) - h_i(y_{it}^*))}{y_{i,t+1}^* - y_{it}^*} - b.$$

Now, by the intermediate value theorem and the conditions on h_i from (A7), we have

$$\frac{v_{i,t+1} - v_{it}}{y_{i,t+1}^* - y_{it}^*} > -(p\beta_i + \alpha v_{i,t+2}\gamma_i) h_i'(y_{i,t+1}^*) - b. \quad (8A)$$

Either $y_{i,t+1}^*$ has been chosen such that the right-hand side of (8A) is zero, in which case the result $v_{i,t+1} - v_{it} > 0$ follows immediately, or $y_{i,t+1}^* = 0$. However, if $y_{i,t+1}^* = 0$, then $y_{it}^* = 0$ by assumption and the result follows as detailed above. Thus, we can conclude that $v_{i,t+1} - v_{it} > 0$, and the result holds by backwards induction. The argument for the condition

$$V_i < \frac{p(1 - \beta_i h_i(y_{iT}^*)) - b y_{iT}^*}{1 - \alpha(1 - \gamma_i h_i(y_{iT}^*) - \gamma_j h_j(y_{iT}^*))}, \text{ for all } i,$$

is analogous. \square

Lemma 3, below, is the n firm analog of Lemma 1, characterizing

the solution to a single period problem in isolation. This serves as the basis for the proof of Theorem 6.

LEMMA 3. Suppose that $g_{i,t+1}(\lambda \setminus \pi^1, \pi^2, \dots, \pi^n) = v_{i,t+1}\lambda_1 + \dots + v_{i,n-1,t+1}\lambda_{n-1} + w_{i,t+1}$ for each i , where λ_k is the k^{th} component of the $n - 1$ dimensional vector λ and v_{ijt} and w_{it} are of indeterminate sign and independent of λ .

(a) There exists a unique dominant pure strategy Nash equilibrium in period t for each state λ . If firm i finds itself at the beginning of period t with a market share of λ_{it} , then the following hold:

(b) (Attainment) The optimal normalized capacity to be provided by firm i , y_{it}^* , is the unique solution to:

if

$$\left(p\beta_i + \alpha v_{i,t+1}\gamma_i - \sum_{j=1, j \neq i}^{n-1} \alpha v_{ij,t+1} a_{ij}\gamma_j \right) h_i'(0) \leq -b,$$

then y_{it}^0 is the unique solution to

$$\left(p\beta_i + \alpha v_{i,t+1}\gamma_i - \sum_{j=1, j \neq i}^{n-1} \alpha v_{ij,t+1} a_{ij}\gamma_j \right) h_i'(y_{it}) = -b;$$

$$y_{it}^* = \begin{cases} y_{it}^0 & \text{if } \left(p\beta_i + \alpha v_{i,t+1}\gamma_i - \sum_{j=1, j \neq i}^{n-1} \alpha v_{ij,t+1} a_{ij}\gamma_j \right) h_i'(0) \leq -b \\ 0 & \text{otherwise;} \end{cases}$$

and is thus independent of both the state λ and all competitors' decisions.

(c) (Preservation) The resulting optimal expected return to firm i (over periods t through T) can be written as $\sum_{j=1}^{n-1} v_{ijt}\lambda_{jt} + w_{it}$, where the v_{ijt} and w_{it} are indeterminate in sign, independent of λ and satisfy:

$$v_{ijt} = \begin{cases} p(1 - \beta_i h_i(y_{it}^*)) - b y_{it}^* + \alpha v_{i,t+1}(1 - \gamma_i h_i(y_{it}^*) - a_{ni}\gamma_n h_n(y_{it}^*)) + \\ \quad \sum_{k=1, k \neq i}^{n-1} \alpha v_{ik,t+1} (a_{ik}\gamma_k h_k(y_{it}^*) - a_{nk}\gamma_n h_n(y_{it}^*)) \text{ for } j = i \\ \alpha v_{ij,t+1}(1 - \gamma_i h_i(y_{it}^*) - a_{ni}\gamma_n h_n(y_{it}^*)) + \\ \quad \sum_{k=1, k \neq j}^{n-1} \alpha v_{ik,t+1} (a_{jk}\gamma_j h_j(y_{it}^*) - a_{nk}\gamma_n h_n(y_{it}^*)) \text{ for } j \neq i, \end{cases}$$

and

$$w_{it} = \alpha w_{i,t+1} + \sum_{j=1}^{n-1} \alpha v_{ij,t+1} a_{nj}\gamma_n h_n(y_{it}^*).$$

PROOF OF LEMMA 3. Follows by analogy to Lemma 1. \square

PROOF OF THEOREM 6. Follows by analogy to Theorem 3 using Lemma 3. \square

References

- Barrett, A., P. Eng, K. Rebello. 1997. For \$19.95 a month, unlimited headaches at AOL. *Business Week* (January 27) 35.
- Fergani, Y. 1976. A Market Oriented Stochastic Inventory Model. Ph.D. Dissertation, Stanford University, Stanford, CA.
- Fudenberg, D., J. Tirole. 1991. *Game Theory*. MIT Press, Cambridge, MA.
- Gans, N. 1999a. Customer learning and loyalty when quality is un-

- certain. Working Paper, OPIM Department, The Wharton School, University of Pennsylvania.
- . 1999b. Customer loyalty and supply strategies for quality competition. Working Paper, OPIM Department, The Wharton School, University of Pennsylvania, Philadelphia, PA.
- Gross, D., C. M. Harris. 1985. *Fundamentals of Queueing Theory*. John Wiley and Sons, New York.
- Inverse Network Technology. 1997. Inverse Network Technology announces new results of Internet service provider (ISP) performance. Press Release, July 22.
- Jones, T. O., E. W. Sasser, Jr. 1995. Why satisfied customers defect. *Harvard Bus. Rev.* (November-December) 88–99.
- Kalai, E., M. I. Kamien, M. Rubinovitch. 1992. Optimal service speeds in a competitive environment. *Management Sci.* **38** (8) 1154–1163.
- Keaveney, S. 1995. Customer switching behavior in service industries: An exploratory study. *J. Marketing* **59**, (April) 71–82.
- Kornblum, J. 1997. AOL agrees to refunds. *WWW.News.Com* January 29.
- Kreps, D. M. 1990. *A Course in Microeconomic Theory*. Princeton University Press, Princeton, NJ.
- , R. Wilson. 1982. Sequential equilibria. *Econometrica* **50** (4) 863–894.
- Li, L. 1992. The role of inventory in delivery time competition. *Management Sci.* **38** (2) 182–197.
- , Y. S. Lee. 1994. Pricing and delivery time performance in a competitive environment. *Management Sci.* **40** (5) 663–646.
- Lippman, S. A., K. F. McCardle. 1997. The competitive newsboy. *Oper. Res.* **45** (1) 54–65.
- Loch, C. 1991. Pricing in markets sensitive to delay. Ph.D. Dissertation, Stanford University, Stanford, CA.
- Mas-Collel, A., M. D. Whinston, J. R. Green. 1995. *Microeconomic Theory*. Oxford University Press, New York.
- McGahan, A. M., P. Ghemawat. 1994. Competition to retain customers. *Marketing Sci.* **13** (2) 165–176.
- Mendelson, H. 1985. Pricing computer services: Queueing effects. *Comm. ACM* **28** (3) 312–321.
- , S. Whang. 1990. Optimal incentive-compatible priority pricing for the M/M/1 queue. *Oper. Res.* **38** (5) 870–883.
- Milgrom, P., J. Roberts. 1992. *Economics, Organization, and Management*. Prentice-Hall, Newark, NJ.
- Oral, M., M. Salvador, A. Reisman, B. Dean. 1972. On the evaluation of shortage costs for inventory control of finished goods. *Management Sci.* **18** (6) B344–B351.
- Pompili, T. 1997. Choosing your ISP. *PC Magazine* **16** (15) 216–217.
- Porteus, E. L. 1982. Conditions for characterizing the structure of optimal strategies in infinite-horizon dynamic programs. *J. Optim. Theory and Appl.* **36** (3) 419–432.
- . 1990. Stochastic inventory theory. D. P. Heyman, M. J. Sobel, eds. *Handbooks in OR & MS, Vol. 2*. Elsevier Science Publishers, New York.
- Reichheld, F. F. 1996. *The Loyalty Effect*. Harvard Business School Press, Boston, MA. 64–67.
- , E. W. Sasser, Jr. 1990. Zero defections: Quality comes to services. *Harvard Bus. Rev.* (September-October) 105–111.
- Serfozo, R. F. 1990. Point processes. D. P. Heyman, M. J. Sobel, eds. *Handbooks in OR & MS, Vol. 2*. Elsevier Science Publishers, New York.
- Tinnirello, P. C. 1997. ISPs: Intermittent service providers. *PC Week* **14** (49) 89.
- Topkis, D. 1978. Minimizing a submodular function on a lattice. *Oper. Res.* **26** (2) 305–321.
- Wetzel, R. 1997. Customers rate ISP services. *PC Week* November 10, 105–124.

The consulting Senior Editor for this manuscript was Lawrence Wein. This manuscript was received November 19, 1998, and was with the authors 6 1/2 months for 2 revisions. The average review cycle time was 56 days.