

e - c o m p a n i o n

ONLY AVAILABLE IN ELECTRONIC FORM

Electronic Companion—“The Value of Partial Resource Pooling:
Should a Service Network Be Integrated or Product-Focused?”
by Barış Ata and Jan A. Van Mieghem, *Management Science*,
DOI 10.1287/mnsc.1080.0918.

A On-Line Appendix on Large Deviations Analysis

In this section, we give a quick review of large deviation analysis (Dembo & Zeitouni (1998) provides a standard reference); summarize relevant results from Avram et al. (2001); and specialize those to our setting to calculate various quantities of interest in closed form.

To justify the time scales implicitly assumed by the large deviations analysis in the context of our motivating application consider a hub or a sorting center. A typical FedEx hub¹⁸ has a sorting processing capacity in the order of tens of thousands of packages per hour, or several packages per second. With processing times of seconds, the appropriate time unit of measurement is also seconds. Given that FedEx quotes delays in the order of days (after subtracting the deterministic travel times from the quoted delay), it is also natural to measure the delays experienced in a hub in hours. The relevant performance criterion thus is whether queueing delays exceed several hours, i.e. several thousands of seconds, which is large compared to the time to process a package. Therefore, adopting a second as a time unit of measurement, the large deviations estimate is appropriate in our motivating application. Moreover, in our analysis the large deviations estimate is used only for class 2 which is precisely the class experiencing most congestion under the greedy scheduling rule.

A.1 Large Deviations Primer

Roughly speaking, most large deviations analysis can be divided into two major parts: proving a Large Deviations Principle (LDP), and solving the associated Variational Problem (VP), which in turn gives the rate function of the LDP. To be specific, a LDP for the steady-state distribution of Q amounts to the following approximation:

$$\mathbb{P}\left(\frac{Q}{u} \in A\right) \approx \exp\left\{-\left(\inf_{v \in A} I(v)\right) u\right\} \quad \text{for } u \text{ large,} \quad (33)$$

where A is a measurable subset of \mathbb{R}_+^2 and $I(v)$ is the rate function, defined as follows. Let ψ denote the Skorohod map, which essentially maps any sample path of the Brownian motion process into a queue length sample path, that is nonnegative, by exerting minimal amount of control at the boundaries of the positive quadrant. Then define the inner product $\langle \cdot, \cdot \rangle: \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ by

$$\langle v, w \rangle = v\Gamma^{-1}w \quad (34)$$

¹⁸For example, the processing capacity of the Dallas hub is 22,500 packages per hour; as of January 19, 2007 available at http://www.findarticles.com/p/articles/mi_m0EIN/is_2005_Sept_22/ai_n15625572.

with the associated norm $\|v\| = \sqrt{\langle v, v \rangle}$. Recall that Γ is the covariance matrix, so $\|\cdot\|$ is the natural norm to consider. Then for each $v \in \mathbb{R}_+^2$ the rate function is defined by the following variational problem

$$I(v) \equiv \inf_{T \geq 0} \inf_{x \in \mathcal{H}^2, \psi(x)(T)=v} \frac{1}{2} \int_0^T \|\dot{x}(t) - \theta\|^2 dt, \quad (35)$$

where \mathcal{H}^2 is the space of all absolutely continuous functions $x(\cdot) : [0, \infty) \rightarrow \mathbb{R}^2$ which have square integrable derivatives on bounded intervals and have $x(0) = 0$. Given $v \in \mathbb{R}_+^2$, if a path $x(\cdot)$ is such that $\psi \circ x(T) = v$ for some $T \geq 0$, and

$$\frac{1}{2} \int_0^T \|\dot{x}(t) - \theta\|^2 dt = I(v),$$

then, as in Avram et al. (2001), we will call x an optimal path for the VP (35) with optimal value $I(v)$. In addition, $\psi(x)$ is called an optimal reflected path. The VP (35) is solved explicitly in Avram et al. (2001), which also characterizes the corresponding optimal path explicitly.

A.2 Summary of Avram, Dai and Hasenbain (2001)

Avram et al. (2001) study a variational problem that arises in Large Deviations analysis of the steady-state distribution of SRBMs. For the two-dimensional case, the authors provide an analytical solution to the VP, which gives an appealing characterization of the optimal path (which in turn characterizes how rare events are most likely to occur) to a given point in the quadrant and also provides an explicit expression for the large deviations rate.

To be more specific, one constructs a “cone of boundary influence” that determines the nature of optimal paths in different regions of the quadrant. When a point $v \in \mathbb{R}_+^2$ is not in either of the cones, the optimal path is a direct, linear path. When v is contained in one of these cones, however, the optimal path first travels along a boundary, and then travels directly to v . Moreover, such a path leaves the boundary and enters the interior at a unique entrance angle which can be determined directly from the problem data.

To state the main theorem of Avram et al. (2001), we need to define the cone C_i associated with the face $F_i = \{x \in \mathbb{R}_+^2 : x_i = 0\}$ for $i = 1, 2$. For a face F_i , cone C_i defines a region of boundary influence on the solutions to the VP. The authors show that the boundary influence depends on two quantities which they term “exit velocity” and “entrance velocity,” which in turn define the “reflectivity” of a face.

Let a^i and \tilde{a}^i denote the exit and entrance velocities associated with face F_i , respectively, for $i = 1, 2$. The velocities for our specific Brownian model follow from formulas (3.2) and (3.4) of Avram et al. (2001):

$$a^1 = \frac{1}{\kappa^2 + 2\rho\kappa + 1} \begin{pmatrix} (1 - \kappa^2)\theta_1 - 2(\kappa + \rho)\theta_2 \\ -2\kappa(1 + \rho\kappa)\theta_1 - (1 - \kappa^2)\theta_2 \end{pmatrix}, \quad a^2 = \begin{pmatrix} -\theta_1 \\ \theta_2 - 2\rho\theta_1 \end{pmatrix},$$

$$\tilde{a}^1 = \frac{1}{\kappa^2 + 2\rho\kappa + 1} \begin{pmatrix} -(1 - \kappa^2)\theta_1 + 2(\kappa + \rho)\theta_2 \\ -2\theta_1(\rho + \kappa) + \theta_2(4\rho^2 + 4\rho\kappa + \kappa^2 - 1) \end{pmatrix}, \quad \tilde{a}^2 = \begin{pmatrix} \theta_1(4\rho^2 - 1) - 2\rho\theta_2 \\ 2\rho\theta_1 - \theta_2 \end{pmatrix},$$

where $\kappa = \sigma_1/\sigma_2$. Clearly, $\tilde{a}_i^i = -a_i^i$ for $i = 1, 2$, which indeed holds in general, cf. Avram et al. (2001). Moreover, the entrance velocity \tilde{a}^i is simply the ‘‘symmetry’’ of the exit velocity a^i around the face F_i with respect to the inner product introduced in (34).

A face F_i is called reflective if the i^{th} component of a^i is negative, that is $a_i^i < 0$. Equivalently, F_i is reflective if and only if $\tilde{a}_i^i > 0$. When F_i is reflective, C_i is defined to be the cone generated by e^i and \tilde{a}^i , where e^i is the directional vector on face F_i , normalized so that $\|e^i\| = 1$. Namely,

$$C_i = \{t_1 e^i + t_2 \tilde{a}^i : t_1, t_2 \geq 0\}, \quad i = 1, 2,$$

where

$$e^1 = \sqrt{1 - \rho^2} \begin{pmatrix} 0 \\ 1 \end{pmatrix} \quad \text{and} \quad e^2 = \sqrt{1 - \rho^2} \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

When F_i is not reflective the associated cone of boundary influence C_i is defined to be empty and the face F_i has no boundary influence on solutions to the VP for any $v \in \mathbb{R}_+^2$.

Avram et al. (2001) prove that the cone C_i identifies precisely the region in which the face F_i has boundary influence. Notice that any two cones C_1 and C_2 partition the state space \mathbb{R}_+^2 into three regions: $(\mathbb{R}_+^2 \cap C_1) \setminus C_2$, $(\mathbb{R}_+^2 \cap C_2) \setminus C_1$, and one of the following two regions $\mathbb{R}_+^2 \cap C_1 \cap C_2$ or $\mathbb{R}_+^2 \setminus (C_1 \cup C_2)$. Note that one of the latter two regions is always empty.

Before we can state the main result of Avram et al. (2001) some additional notation is needed. For $v \in \mathbb{R}_+^2$, let

$$\tilde{a}^0(v) = \frac{\|\theta\|}{\|v\|} v.$$

The following expressions determine the locally optimal value of the VP for various cases. For $v \in \mathbb{R}_+^2$, let

$$I^0(v) = \langle \tilde{a}^0(v) - \theta, v \rangle,$$

$$I^i(v) = \langle \tilde{a}^i(v) - \theta, v \rangle, \quad i = 1, 2.$$

Next, specializing Theorem 3.1 of Avram et al. (2001) to our setting we have the following theorem.

Theorem 4 (Avram et al.) *Consider the VP (35) with data (θ, Γ, R) . Also assume $\theta_1 < 0$ and $\theta_2 + \kappa\theta_1 < 0$. Then*

- (a) *If $v \notin C_1 \cup C_2$, then $I(v) = I^0(v)$;*
- (b) *If $v \in C_1 \setminus C_2$, then $I(v) = I^1(v)$;*
- (c) *If $v \in C_2 \setminus C_1$, then $I(v) = I^2(v)$;*
- (d) *If $v \in C_1 \cap C_2$, then $I(v) = \min \{I^1(v), I^2(v)\}$.*

Avram et al. (2001) also characterizes the optimal path associated with the VP in each of these cases either as a direct linear path from the origin to v , or as a piecewise linear path with only one break point, see Avram et al. (2001) for further discussion.

A.3 The large deviations rate function $I(v)$ for our Brownian model

In this section, we prove Proposition 2. Recall that we want to calculate the steady-state probability $\mathbb{P}(Q_2 > u)$. Clearly,

$$\mathbb{P}(Q_2 > u) = \mathbb{P}\left(\frac{Q_2}{u} > 1\right).$$

Therefore, letting $A = \{v \in \mathbb{R}_+^2 : v_1 \geq 0, v_2 \geq 1\}$, it follows from (33) that¹⁹

$$\mathbb{P}(Q_2 > u) = \exp\left\{-\left(\inf_{v \in A} I(v)\right) u\right\} \text{ for } u \text{ large.}$$

In other words, the rate r_q in (20) is exactly $\inf_{v \in A} I(v)$ where $A = \{v \in \mathbb{R}_+^2 : v_1 \geq 0, v_2 \geq 1\}$, which we now calculate. The scaling Lemma 5.3 of Avram et al. (2001) implies that

$$\inf_{v \in A} I(v) = \inf_{v \in \tilde{A}} I(v),$$

where $\tilde{A} = \{v \in \mathbb{R}_+^2 : v_1 \geq 0, v_2 = 1\}$. Theorem 3.1 and Lemma 6.1 of Avram et al. (2001) then yield

$$\inf_{v \in \tilde{A}} I(v) = \min \left\{ \inf_{v \in \tilde{A}} I^0(v), \inf_{v \in \tilde{A} \cap C_1} I^1(v), \inf_{v \in \tilde{A} \cap C_2} I^2(v) \right\}, \quad (36)$$

where $\inf_{\emptyset} I^i(v) = \infty$ for $i = 1, 2$ by convention. Therefore, next we compute each of the terms on the right hand side of (36). By straightforward calculations outlined in Avram et al. (2001) it

¹⁹Readers should keep in mind that the statement is an approximation, which becomes precise as u gets large; see Avram et al. (2001) for a precise statement.

follows that

$$\begin{aligned} I^0(v) &= \frac{1}{1-\rho^2} \left\{ \sqrt{\theta_1^2 - 2\rho\theta_1\theta_2 + \theta_2^2} \sqrt{(x-\rho)^2 + 1 - \rho^2} + x(\rho\theta_2 - \theta_1) + \rho\theta_1 - \theta_2 \right\}, \\ I^1(v) &= \frac{2}{(\kappa^2 + 2\rho\kappa + 1)} \{x[-\theta_1 + \theta_2(\kappa + 2\rho)] - (\theta_1\kappa + \theta_2)\}, \\ I^2(v) &= 2\{-\theta_1x + 2\rho\theta_1 - \theta_2\}, \end{aligned}$$

where $v = (1, x)'$. Recall that we focus on the case $\rho \leq 0$, in which case it is straightforward to show that

$$\inf_{v \in \tilde{A}} I^0(v) = I^0 \left(\left(\begin{array}{cc} 0 & 1 \end{array} \right)' \right) = \frac{1}{1-\rho^2} \left[\sqrt{\theta_1^2 - 2\rho\theta_1\theta_2 + \theta_2^2} + \rho\theta_1 - \theta_2 \right] = \sqrt{nr}r_1.$$

To calculate the other terms in (36) we need to characterize the cones C_1 and C_2 . Recall that the face F_1 has influence only if it is reflective, that is, $\tilde{a}_1^1 > 0$, in which case it follows that

$$\inf_{v \in \tilde{A} \cap C_1} I^1(v) = I^1 \left(\left(\begin{array}{cc} 0 & 1 \end{array} \right)' \right) = \frac{-2(\theta_1\kappa + \theta_2)}{\kappa^2 + 2\rho\kappa + 1} = \sqrt{nr}r_4.$$

Calculating $\inf_{v \in \tilde{A} \cap C_2} I^2(v)$ is more involved and requires considering different cases. Clearly, face F_2 has influence if and only if $\tilde{a}_2^2 > 0$, in which case we need to consider the following two cases, which are illustrated in Figure 6.

1. $\tilde{a}_1^2 = \theta_1(4\rho^2 - 1) - 2\rho\theta_2 \leq 0$. In this case $\mathbb{R}_+^2 \subset C_2$ and it follows that

$$\inf_{v \in \tilde{A} \cap C_2} I^2(v) = I^2 \left(\left(\begin{array}{cc} 0 & 1 \end{array} \right)' \right) = 2(2\theta_1\rho - \theta_2) = \sqrt{nr}r_2.$$

2. $\tilde{a}_1^2 = \theta_1(4\rho^2 - 1) - 2\rho\theta_2 > 0$. In this case $C_2 \subset \subset \mathbb{R}_+^2$ and it follows that $\inf_{v \in \tilde{A} \cap C_2} I^2(v) = I^2 \left(\left(\begin{array}{cc} x_0 & 1 \end{array} \right)' \right)$, where $x_0 = 2\rho - \theta_1/(2\theta_1\rho - \theta_2)$. Therefore, it follows that

$$\inf_{v \in \tilde{A} \cap C_2} I^2(v) = 2 \left[\frac{\theta_1^2}{2\theta_1\rho - \theta_2} - \theta_2 \right] = \sqrt{nr}r_3.$$

Finally, combining these derivations with (36) and Theorem 4 one can derive $\inf_{v \in A} I(v)$ explicitly as follows, depending on six cases.

Case I: $\tilde{a}_1^1 \leq 0$ so F_1 is *not* reflective. Then there are two subcases to consider:

Case I.a: $\tilde{a}_2^2 \leq 0$ so F_2 is *not* reflective. Then neither boundary has an effect on the rate, and it follows that $\inf_{v \in A} I(v) = I^0 \left(\left(\begin{array}{cc} 0 & 1 \end{array} \right)' \right) = \sqrt{nr}r_1$.

Case I.b: $\tilde{a}_2^2 > 0$ so F_2 is reflective. Then we need to consider the following two further subcases:

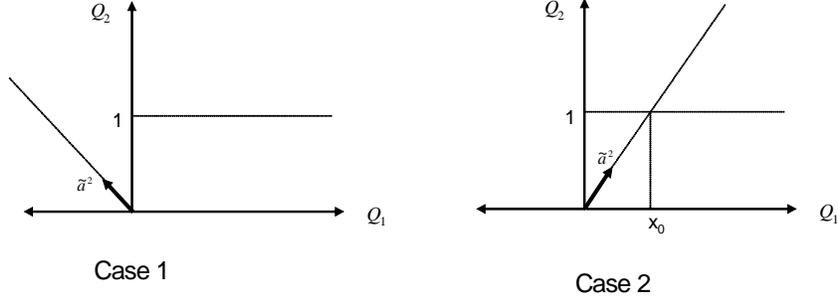


Figure 6 Depending on the sign of \tilde{a}_1^2 the cone of boundary influence C_2 may or may not subsume \mathbb{R}_+^2 . Case 1 illustrates $\mathbb{R}_+^2 \subset C_2$ and Case 2 illustrates $C_2 \subset \mathbb{R}_+^2$.

Case I.b.i: $\tilde{a}_1^2 \leq 0$. Then, $\mathbb{R}_+^2 \subset C_2$ and it follows that $\inf_{v \in A} I(v) = I^2 \left(\begin{pmatrix} 0 & 1 \end{pmatrix}' \right) = \sqrt{n}r_2$.

Case I.b.ii: $\tilde{a}_1^2 > 0$. Then, $C_2 \subset \subset \mathbb{R}_+^2$ and it follows that

$$\inf_{v \in A} I(v) = \min \left\{ I^0 \left(\begin{pmatrix} 0 & 1 \end{pmatrix}' \right), I^2 \left(\begin{pmatrix} x_0 & 1 \end{pmatrix}' \right) \right\},$$

where $x_0 = 2\rho - \theta_1 / (2\theta_1\rho - \theta_2)$. Therefore, it follows that $\inf_{v \in A} I(v) = \sqrt{n} \min \{r_1, r_3\}$.

Case II: $\tilde{a}_1^1 > 0$ so F_1 is reflective. Then we must consider the following two subcases:

Case II.a: $\tilde{a}_2^2 \leq 0$ so F_2 is *not* reflective. Then $\inf_{v \in A} I(v) = I^1 \left(\begin{pmatrix} 0 & 1 \end{pmatrix}' \right) = \sqrt{n}r_4$.

Case II.b: $\tilde{a}_2^2 > 0$ so F_2 is reflective. Then we need to consider the following two further subcases:

Case II.b.i: $\tilde{a}_1^2 = \theta_1(4\rho^2 - 1) - 2\rho\theta_2 \leq 0$. In this case $\mathbb{R}_+^2 \subset C_2$ and it follows that

$$\inf_{v \in A} I(v) = \min \left\{ I^1 \left(\begin{pmatrix} 0 & 1 \end{pmatrix}' \right), I^2 \left(\begin{pmatrix} 0 & 1 \end{pmatrix}' \right) \right\} = \sqrt{n} \min \{r_4, r_2\}.$$

Case II.b.ii: $\tilde{a}_1^2 = \theta_1(4\rho^2 - 1) - 2\rho\theta_2 > 0$. In this case, $C_2 \subset \subset \mathbb{R}_+^2$ and it follows that

$$\inf_{v \in A} I(v) = \min \left\{ I^1 \left(\begin{pmatrix} 0 & 1 \end{pmatrix}' \right), I^2 \left(\begin{pmatrix} x_0 & 1 \end{pmatrix}' \right) \right\}$$

where $x_0 = 2\rho - \theta_1 / (2\theta_1\rho - \theta_2)$. Therefore, it follows that $\inf_{v \in A} I(v) = \sqrt{n} \min \{r_4, r_3\}$.

Assuming $\rho \leq 0$ and $z > 0$ (that is, $\theta < 0$) it follows that $\tilde{a}_2^2 > 0$, and F_2 is always reflective, eliminating cases I.a and II.a above. Proposition 2 simply summarizes the remaining four cases of interest and it also expresses the conditions for each case in terms of primitives of the model described in Section 3.

B On-Line Appendix: Comparative Statics and the Discussion of Table 1

The impact of correlation on the integration value is already discussed in Section 5. Next, we discuss the comparative statics of the other parameters in Table 1.

1. *Guaranteed speed of service d* affects integration value through the ratio d_2/d_1 . An increase in d_2/d_1 increases value, the intuition behind which is as follows. As d_1 goes down the firm needs more flexible standardized excess capacity which increases the value of resource substitution. In contrast, as d_2 goes down the standardized excess capacity of resource 2 increases to provide the same reliability ε_2 in the dedicated network, whereas the excess capacity of the flexible resource remains the same. Therefore, the value of dynamic substitution decreases.
2. *Reliability ε* : Higher express service reliability (smaller ε_1) increases the value of integration. Indeed, such reliability increase again requires the firm to build more flexible standardized excess capacity thereby increasing the option value of resource substitution. However, the effect of ε_2 is more involved and will be discussed in the next section.
3. *Variance σ^2* affects the integration value through the ratio σ_1/σ_2 . An increase in σ_1/σ_2 increases value, and the intuition behind that is as follows. Keeping everything else constant, an increase in express class variability requires more (standardized) excess capacity of the flexible resource to provide the same quality of service, cf. (24), which increases the option value. In contrast, as the regular class becomes more variable the flexible resource capacity remains the same and the option value of dynamic substitution goes down because it becomes harder for the flexible resource to absorb larger fluctuations in regular demand. Moreover, the slow server capacity increases as σ_2 increases, decreasing the value of integration.
4. *Volume λ* : First consider the simplest comparative statics for volume by keeping variance constant (and thus ignoring the natural relationship (22) between volume and variance). In that case, an increase in volume of the express class λ_1 results in smaller (standardized) excess capacity of the flexible resource, cf. (24). Thus the option value decreases. In contrast, an increase in volume of the regular class λ_2 has a linear impact if the violation probability remains the same, which follows simply from the definition of value of integration V , cf. (28). Moreover, as the volume of the regular class increases we see a “statistical economies of

scale” effect: the standardized excess capacity z_2 of resource 2 decreases in λ_2 (which provides reliability of ε_2 in the dedicated network) while the standardized excess capacity z_1 of the flexible server remains the same so that the substitution frequency increases. Therefore, the relative impact of integration increases as λ_2 increases. More specifically, the probability of service failure decreases. Thus, the value of integration V increases superlinearly in λ_2 .

Although one gains some insight by considering isolated changes in volume or variance of each class, it is rare to have an isolated change in volume. Typically, when volume changes, the variance may change as captured by the arrival pooling parameter γ specified in (22). A more realistic comparative static on volume thus is to let λ change while keeping γ constant. Such combined volume-variance change in essence combines the effects of isolated changes in σ and λ . As the regular class volume-variance increases (keeping γ constant), the standardized excess capacity of resource 2 decreases, cf. (25), while the standardized excess capacity of the flexible resource remains the same. Therefore, the relative impact of integration value increases. This can also be interpreted as showing that the gains from statistical economies of scale dominate the reduction in the option value of flexible server’s excess capacity, which is due to the increase in variability of regular demand. However, the effect of a change in the express class volume-variance is more involved and will be discussed next.

Discussion of Subtle Comparative Statics. The two remaining comparative statics that impact the two types of firms differently are the effects of combined volume-variance of the express class and the reliability increase for the regular class. To understand these effects and their different consequences on the two types of firms, we found it helpful to break down how value is impacted by the three main drivers: arrival pooling, resource substitution, and the correlation effect as discussed in Section 5.

First consider an increase of the combined volume-variance of the express class; that is, λ_1 increases for a fixed γ_1 . For an express firm there is no strong substitution effect, while the dominant driver of value is negative correlation. As λ_1 changes, the standardized excess capacity z_1 of the express class changes to provide the same reliability in the dedicated network. The change in z_1 , of course, impacts how the express class queue behaves, which in turn affects the regular class queue through the correlation effect. Therefore, what matters in this case is the standardized relative excess capacity z_1 of the express class, which decreases as λ_1 increases, cf. (24). Therefore, the value of integration decreases for the express firm as λ_1 increases. In contrast, there is a strong

substitution effect for the regular firm. That is, the standardized excess capacity z_1 of the flexible server gets substituted for resource 2, and one unit of standardized excess capacity of the flexible server corresponds to σ_1/σ_2 units of standardized capacity of resource 2. Therefore, the standardized excess capacity of the flexible server measured in units of standardized capacity of resource 2, that is, the quantity $z_1\sigma_1/\sigma_2$, is the driver of value in this case. Moreover, $z_1 \sigma_1/\sigma_2$ increases in λ_1 (keeping γ_1 constant), resulting in higher value for the regular firm.

Next, consider the reliability increase of the regular class service. As the reliability of regular class service increases, that is, as ε_2 decreases, the standardized excess capacity z_2 of resource 2 increases, while that for the flexible resource remains the same. Naturally, one would expect the value of integration to go down, which is indeed the case for the express firm. Now consider a regular firm. Its strong substitution effect leads to partial resource pooling (meaning that both resource types almost behave as one resource pool) as discussed earlier that results in efficiencies in addition to those stemming from arrival pooling. With negative correlation, this partial resource pooling decreases the variability faced by resource 2. This indeed is reflected by the term $\varepsilon_2^{\sigma_2^2/\sigma_T^2}$ in (5), where $\sigma_2^2/\sigma_T^2 > 1$ for the regular firm. Thus, as ε_2 decreases, the service failure probability decreases at a faster rate, resulting in an increase in the value of integration for the regular firm. (Observe from Corollary 2 that the positive impact of ε_2 on a regular firm's value really is a second order effect.)

Readers may wonder why a change in the speed of service d_2 for regular service does not produce a similar effect as a change in ε_2 ? The reason is that the impact of variance reduction due to partial resource pooling is unaffected by the speed of service and depends only on demand characteristics and the reliability of regular class service, both of which are constant when considering a change in d_2 .

Our analysis also admits a similar study of comparative statics of the relative value of integration. For example, it is easy to see from Corollaries 1 and 2 that relative value increases superlinearly in regular class volume for an express firm while it increases sublinearly for a regular firm.

C On-Line Appendix: Proofs in Sections 4-7

Proof of Proposition 1: Given that the express class receives static priority service at server 1, its dynamics are unaffected by the regular class. Thus, Q_1 can be analyzed in isolation as a simple

one-dimensional standard reflected Brownian motion with stationary distribution

$$\mathbb{P}(Q_1 > x) = \exp\{2\theta_1 x\}, \quad (37)$$

cf. page 94 of Harrison (1985). Applying (18) gives our express violation probability:

$$\mathbb{P}(D_1 > d_1) = \exp\left\{2\theta_1 \frac{\lambda_1 d_1}{\sqrt{n}\sigma_1}\right\} = \exp\left\{2\sqrt{n} \frac{(\lambda_1 - \mu_1)}{\sigma_1} \frac{\lambda_1 d_1}{\sqrt{n}\sigma_1}\right\} = \exp\left\{\frac{2(\lambda_1 - \mu_1)\lambda_1}{\sigma_1^2} d_1\right\}. \quad \blacksquare$$

Proof of Proposition 4. We must determine the relevant case in Proposition 2. Given that $\lambda_1 \gg \lambda_2$, it follows from (22), (24), and (25) that $\left(1 - (\sigma_1/\sigma_2)^2\right) \frac{z_1}{z_2} < 0 < 2(\sigma_1/\sigma_2) + 2\rho$ and $\frac{z_1}{z_2}(4\rho^2 - 1) \approx 0 > 2\rho$. Therefore, we are in case 1 of Proposition 2, that is, $r = r_2$. Thus, Proposition 2 yields that

$$\begin{aligned} \mathbb{P}(D_2 > d_2) &= \exp\left\{-r_2 \frac{\lambda_2}{\sigma_2} d_2\right\} \\ &= \exp\left\{\frac{2\lambda_2(\lambda_2 - \mu_2^\varepsilon)}{\sigma_2^2} d_2\right\} \times \exp\left\{\frac{-4\rho(\lambda_1 - \mu_1^\varepsilon)}{\sigma_1} \frac{\lambda_2}{\sigma_2} d_2\right\}. \end{aligned} \quad (38)$$

On the other hand, by (23)-(24) it follows that

$$\varepsilon_i = \exp\left\{\frac{2(\lambda_i - \mu_i^\varepsilon)\lambda_i}{\sigma_i^2} d_i\right\}, \quad i = 1, 2. \quad (39)$$

Combining this with (38) gives the result. \blacksquare

Proof of Corollary 1. It follows from Proposition 4 that

$$\mathbb{P}(D_2 > d_2) = \varepsilon_2 \exp\left\{-2\rho \frac{\lambda_2 \sigma_1}{\lambda_1 \sigma_2} \frac{d_2}{d_1} \log(\varepsilon_1)\right\} = \varepsilon_2 \exp\left\{-2\rho \frac{d_2}{d_1} \log(\varepsilon_1) \frac{\lambda_2^{1-\gamma_2}}{\lambda_1^{1-\gamma_1}}\right\}.$$

Taylor's expansion of this combined with (28) gives (30). \blacksquare

Proof of Proposition 5. We must determine the relevant case in Proposition 2. Given that $\lambda_2 \gg \lambda_1$, it follows from (22), (24), and (25) that $\left(1 - (\sigma_1/\sigma_2)^2\right) \frac{z_1}{z_2} > 0 > 2(\sigma_1/\sigma_2) + 2\rho$. (Thus, if $\rho < 0$, we have that $\sigma_1 + 2\rho\sigma_2 < 0$ so that $\sigma_T^2 < \sigma_2^2$.) Therefore, we are either in case 3 or case 4 of Proposition 2 depending on the value of ρ . Thus, consider two cases:

Case (i): $\rho \leq -1/2$. Then $\frac{z_1}{z_2}(4\rho^2 - 1) \geq 0 > 2\rho$ and we are in case 3 of Proposition 2. That is, $r = \min(r_2, r_4)$. We show that $r_2 > r_4$:

$$r_2 = 2z_2 - 4\rho z_1 = \frac{\sigma_2}{\lambda_2} \frac{\log\left(\frac{1}{\varepsilon_2}\right)}{d_2} - 2\rho \frac{\sigma_1}{\lambda_1} \frac{\log\left(\frac{1}{\varepsilon_1}\right)}{d_1}.$$

For λ_2/λ_1 large, substituting $\sigma_i = \lambda_i^{\gamma_i}$ for $i = 1, 2$ gives

$$r_2 = \frac{\log\left(\frac{1}{\varepsilon_2}\right)}{d_2} \frac{1}{\lambda_2^{1-\gamma_2}} - 2\rho \frac{\log\left(\frac{1}{\varepsilon_1}\right)}{d_1} \frac{1}{\lambda_1^{1-\gamma_1}} \approx -2\rho \frac{1}{d_1 \lambda_1^{1-\gamma_1}} \log\left(\frac{1}{\varepsilon_1}\right).$$

Similarly,

$$\begin{aligned} r_4 &= \frac{2(\mu_1 - \lambda_1 + \mu_2 - \lambda_2)}{\sigma_1^2 + 2\rho\sigma_1\sigma_2 + \sigma_2^2} \sigma_2 = \frac{-\frac{\sigma_1^2 \log(\varepsilon_1)}{\lambda_1 \frac{d_1}{d_1}} - \frac{\sigma_2^2 \log(\varepsilon_2)}{\lambda_2 \frac{d_2}{d_2}}}{\sigma_1^2 + 2\rho\sigma_1\sigma_2 + \sigma_2^2} \sigma_2 \\ &= \frac{-\frac{\log(\varepsilon_1)}{d_1} \lambda_1^{2\gamma_1-1} - \frac{\log(\varepsilon_2)}{d_2} \lambda_2^{2\gamma_2-1}}{\lambda_1^{2\gamma_1} + 2\rho\lambda_1^{\gamma_1}\lambda_2^{\gamma_2} + \lambda_2^{2\gamma_2}} \lambda_2^{\gamma_2} \approx -\frac{\log(\varepsilon_2)}{d_2} \frac{1}{\lambda_2^{1-\gamma_2}}. \end{aligned}$$

Given that λ_2/λ_1 is large, $r_4 \ll r_2$, so that $r = r_4$.

Case (ii): $\rho > -1/2$. Then $\frac{z_1}{z_2}(4\rho^2 - 1) < 2\rho < 0$ and we are in case 4 of Proposition 2. That is, $r = \min(r_3, r_4)$. We show that $r_3 > r_4$:

$$r_3 = \frac{2z_1^2}{z_2 - 2\rho z_1} + 2z_2 = \frac{2\left(\frac{\sigma_1}{2d_1\lambda_1} \log\left(\frac{1}{\varepsilon_1}\right)\right)^2}{\frac{\sigma_2}{2d_2\lambda_2} \log\left(\frac{1}{\varepsilon_2}\right) - \rho \frac{\sigma_1}{d_1\lambda_1} \log\left(\frac{1}{\varepsilon_1}\right)} + \frac{\sigma_2}{d_2\lambda_2} \log\left(\frac{1}{\varepsilon_2}\right).$$

For λ_2/λ_1 large, substituting $\sigma_i = \lambda_i^{\gamma_i}$ for $i = 1, 2$ gives

$$r_3 = \frac{\left(\frac{1}{2d_1} \log\left(\frac{1}{\varepsilon_1}\right)\right)^2 \frac{1}{\lambda_1^{2-2\gamma_1}}}{\frac{\log(\varepsilon_2)}{d_2} \frac{1}{\lambda_2^{1-\gamma_2}} - 2\rho \frac{\log\left(\frac{1}{\varepsilon_1}\right)}{d_1} \frac{1}{\lambda_1^{1-\gamma_1}}} - \frac{1}{d_2\lambda_2^{1-\gamma_2}} \log\left(\frac{1}{\varepsilon_2}\right) \approx -\frac{\log\left(\frac{1}{\varepsilon_1}\right)}{2\rho d_1} \frac{1}{\lambda_1^{1-\gamma_1}},$$

while r_4 remains as in case (i). Given that λ_2/λ_1 is large, $r_4 \ll r_3$, so that again $r = r_4$.

Thus, in either case, we have that

$$\mathbb{P}(D_2 > d_2) = \exp\left\{-r_4 \frac{\lambda_2}{\sigma_2} d_2\right\} \quad (40)$$

$$= \exp\left\{\frac{2[(\lambda_1 - \mu_1) + (\lambda_2 - \mu_2)]\lambda_2 d_2}{\sigma_1^2 + 2\rho\sigma_1\sigma_2 + \sigma_2^2}\right\}. \quad (41)$$

Combining this with Proposition 3 gives the result. ■

Proof of Corollary 2. It follows from Proposition 5 that $\mathbb{P}(D_2 > d_2) = \varepsilon_2 e^\delta$, where

$$\delta = \left(\frac{\sigma_2^2}{\sigma_T^2} - 1\right) \log(\varepsilon_2) + \frac{\sigma_1^2}{\sigma_T^2} \frac{\lambda_2}{\lambda_1} \frac{d_2}{d_1} \log(\varepsilon_1).$$

Taylor's expansion yields

$$\mathbb{P}(D_2 > d_2) = \varepsilon_2 [1 + \delta + O(\delta^2)]. \quad (42)$$

Because $\lambda_2 \gg \lambda_1$ and $\sigma_i = \lambda_i^{\gamma_i}$ for $i = 1, 2$, it follows that $\delta \leq 0$ and $|\delta| \ll 1$. To approximate δ first observe that

$$\frac{\sigma_2^2}{\sigma_T^2} = 1 - 2\rho \frac{\sigma_1}{\sigma_2} + O\left(\frac{\sigma_1^2}{\sigma_2^2}\right) = 1 - 2\rho \frac{\lambda_1^{\gamma_1}}{\lambda_2^{\gamma_2}} + O\left(\frac{\lambda_1^{2\gamma_1}}{\lambda_2^{2\gamma_2}}\right). \quad (43)$$

Using this it follows that

$$\frac{\sigma_1^2 \lambda_2}{\sigma_T^2 \lambda_1} = \frac{\sigma_2^2 \sigma_1^2 \lambda_2}{\sigma_T^2 \sigma_2^2 \lambda_1} = \left[1 - 2\rho \frac{\sigma_1}{\sigma_2} + O\left(\frac{\sigma_1^2}{\sigma_2^2}\right) \right] \frac{\sigma_1^2 \lambda_2}{\sigma_2^2 \lambda_1} = \frac{\lambda_1^{2\gamma_1-1}}{\lambda_2^{2\gamma_2-1}} + O\left(\frac{\lambda_1^{3\gamma_1-1}}{\lambda_2^{3\gamma_2-1}}\right).$$

Combining this with (43) we have that

$$\delta = \frac{\lambda_1^{2\gamma_1-1}}{\lambda_2^{2\gamma_2-1}} \frac{d_2}{d_1} \log(\varepsilon_1) - 2\rho \frac{\lambda_1^{\gamma_1}}{\lambda_2^{\gamma_2}} \log(\varepsilon_2) + O\left(\frac{\lambda_1^{3\gamma_1-1}}{\lambda_2^{3\gamma_2-1}}\right), \quad (44)$$

which combined with (42) and (28) gives the Taylor expansion for a regular firm's integration value.

■

Proof of Theorem 1. For notational simplicity assume $\lambda_2^e = \lambda_1^r = 1$ and $\lambda_1^e = \lambda_2^r = N \gg 1$.

Then it follows from Corollary 1 that

$$V^e \approx 2\rho p_2 \varepsilon_2 \log(\varepsilon_1) \frac{d_2}{d_1} \frac{1}{N^{1-\gamma_1}} \ll 1.$$

Similarly, it follows from Corollary 2 that

$$V^r \approx -p_2 \log(\varepsilon_1) \frac{d_2}{d_1} N^{2-2\gamma_2} + 2\rho p_2 \log(\varepsilon_2) N^{1-\gamma_2} \gg 1.$$

Therefore, $V^r \gg V^e$. On the other hand, it also follows from Corollary 1 that

$$\frac{V^e}{C_D^e} \approx \frac{2\rho p_2 \varepsilon_2 \log(\varepsilon_1) \frac{d_2}{d_1} \frac{1}{N^{1-\gamma_1}}}{N p_1 \varepsilon_1 + p_2 \varepsilon_2} \approx 2\rho \frac{p_2 \varepsilon_2}{p_1 \varepsilon_1} \frac{d_2}{d_1} \log(\varepsilon_1) \frac{1}{N^{2-\gamma_1}}.$$

Similarly, it follows from Corollary 2 that

$$\begin{aligned} \frac{V^r}{C_D^r} &= \frac{-p_2 \log(\varepsilon_1) \frac{d_2}{d_1} N^{2-2\gamma_2} + 2\rho p_2 \log(\varepsilon_2) N^{1-\gamma_2}}{p_1 \varepsilon_1 + N p_2 \varepsilon_2} \\ &\approx -\frac{d_2 \log(\varepsilon_1)}{d_1} \frac{1}{\varepsilon_2} \frac{1}{N^{2\gamma_2-1}} + 2\rho \frac{\log(\varepsilon_2)}{\varepsilon_2} \frac{1}{N^{\gamma_2}} \approx -\frac{d_2 \log(\varepsilon_1)}{d_1} \frac{1}{\varepsilon_2} \frac{1}{N^{2\gamma_2-1}}. \end{aligned}$$

Since $2\gamma_2 - 1 < 2 - \gamma_1$ for all $\gamma_1, \gamma_2 \in [1/2, 1)$ we have that $\frac{1}{N^{2\gamma_2-1}} \gg \frac{1}{N^{2-\gamma_1}}$, which clearly implies

$$\frac{V^r}{C_D^r} \gg \frac{V^e}{C_D^e}. \quad \blacksquare$$

Proof of Proposition 6. It is clear that (in the integrated network) as one reduces the slow server capacity (or, equivalently the standardized excess capacity z_2) regular service failure probability increases. We want to reduce z_2 to the critical value, say z_2^* , which corresponds to the service failure probability ε_2 . One can use Proposition 2 to calculate service failure probabilities for different values of z_2 . Clearly, as in proof of Proposition 4, case 1 of Proposition 2 applies before any capacity reduction. Assuming we are in case 1 of Proposition 2 throughout the interval (z_2^*, z_2^d) as we reduce z_2 from z_2^d to z_2^* , where $z_2^d = (\mu_2^d - \lambda_2)/\sigma_2$, we have that

$$\mathbb{P}(D_2 > d_2) = \exp\left\{-r_2 \frac{\lambda_2}{\sigma_2} d_2\right\} = \varepsilon_2,$$

where r_2 is evaluated at $z_2 = z_2^*$. Substituting $r_2 = 2z_2 - 4z_1\rho$ and combining that with (7) and (22), and solving for the slow-server capacity μ_2 gives the result.

To conclude the proof, we must justify the assumption that for all $z_2 \in (z_2^*, z_2^d)$ case 1 of Proposition 2 applies, which can be proved along the lines of the proof of Proposition 4, since $\lambda_1 \gg \lambda_2$ and

$$\frac{z_2^*}{z_2} = 1 + \rho \frac{2d_2 \log(\varepsilon_1)}{d_1 \log(\varepsilon_2)} \frac{\lambda_2^{1-\gamma_2}}{\lambda_1^{1-\gamma_1}} \approx 1. \quad \blacksquare$$

Proof of Proposition 7. As in the proof of Proposition 6, it is clear that (in the integrated network) as one reduces the slow-server capacity (or, equivalently, the standardized excess capacity z_2) regular service failure probability increases. We want to reduce z_2 to the critical value of z_2^* , which corresponds to the service failure probability ε_2 . One can use Proposition 2 to calculate service failure probabilities for different values of z_2 . As in Proposition 5, $r = r_4$ before any capacity reduction (we are in either case 3 or case 4 of Proposition 2 depending on the value of ρ). Assuming that $r = r_4$ throughout the interval (z_2^*, z_2^d) as we reduce z_2 from z_2^d to z_2^* , where $z_2^d = (\mu_2^d - \lambda_2)/\sigma_2$, we have that

$$\mathbb{P}(D_2 > d_2) = \exp \left\{ -r_4 \frac{\lambda_2}{\sigma_2} d_2 \right\} = \varepsilon_2,$$

where r_4 is evaluated at $z_2 = z_2^*$. Substituting $r_4 = \frac{2(\sigma_1 z_1 + \sigma_2 z_2)}{\sigma_1^2 + 2\rho\sigma_1\sigma_2 + \sigma_2^2} \sigma_2$ and combining that with (7) and (22), and solving for the slow-server capacity μ_2 gives the result.

To conclude the proof, we must justify the assumption that for all $z_2 \in (z_2^*, z_2^d)$ we have that $r = r_4$, which can be proved along the lines of the proof of Proposition 5, since $\lambda_2 \gg \lambda_1$ and

$$\frac{z_2^*}{z_2} = 1 - \frac{d_2 \log(\varepsilon_2)}{d_1 \log(\varepsilon_1)} \frac{\lambda_1^{2\gamma_1-1}}{\lambda_2^{2\gamma_2-1}} - \left[\frac{\rho \lambda_1^{\gamma_1}}{\lambda_2^{1-\gamma_2} d_2} + \frac{\lambda_1^{2\gamma_1}}{2\lambda_2 d_2} \right] \log \left(\frac{1}{\varepsilon_2} \right) \approx 1. \quad \blacksquare$$

Proof of Theorem 2. For notational simplicity assume $\lambda_2^e = \lambda_1^r = 1$ and $\lambda_1^e = \lambda_2^r = N \gg 1$. Then it follows Proposition 6 that $V^e = c_2 \Delta\mu_2 = \frac{c_2 \rho \log(\varepsilon_1)}{d_1 N^{1-\gamma_1}}$. Similarly, it follows from Proposition 7 that

$$V^r = c_2 \Delta\mu_2 = \frac{c_2}{2d_1} \log \left(\frac{1}{\varepsilon_1} \right) - c_2 \left[\frac{\rho}{N^{1-\gamma_2} d_2} + \frac{1}{2N d_2} \right] \log \left(\frac{1}{\varepsilon_2} \right).$$

Therefore, $V^r \gg V^e$. On the other hand,

$$\frac{V^e}{C_D^e} = \frac{\frac{c_2 \rho \log(\varepsilon_1)}{d_1 N^{1-\gamma_1}}}{N p_1 \varepsilon_1 + p_2 \varepsilon_2} \approx \frac{c_2 \rho \log(\varepsilon_1)}{d_1 p_1 \varepsilon_1} \frac{1}{N^{2-\gamma_1}} \quad \text{and} \quad \frac{V^r}{C_D^r} = \frac{V^r}{p_1 \varepsilon_1 + N p_2 \varepsilon_2} \approx \frac{c_2 \log(\frac{1}{\varepsilon_1})}{2 p_2 \varepsilon_2 d_2} \frac{1}{N}.$$

Therefore, $\frac{V^r}{C_D^r} \gg \frac{V^e}{C_D^e}$. \blacksquare

D On-Line Appendix: A Simulation Study

In this appendix, we provide a simulation study and a sensitivity analysis as a sanity check of our analytic approximation. We shall see that the simulation study agrees with the recommendation stemming from our analytical study of the integration value. In what follows, we first establish a realistic base case for our simulation studies, specifying the demand data to be used. We then report the simulation results followed by a detailed sensitivity analysis.

First, we describe the demand data we used for the express firm (e.g. FedEx). We assume that the total arrival rate²⁰ is 12,000 requests per hour corresponding to 200 requests per minute, which is indeed the right order of magnitude for a FedEx hub as mentioned in Section 4. Then assuming that 80% of the requests are express, we use the demand data for the express firm shown in Table 3. (We first report results assuming demand is bivariate normally distributed; at the end of this section we discuss this and also report results for the bivariate lognormal distribution). For the regular firm (e.g. UPS) we use the symmetric demand data which is shown in Table 4. We assume a correlation of $\rho = -0.4$. (In what follows, we will present the results for a spectrum of negative correlation values.) Moreover, as discussed in Section 4 given that FedEx and UPS quote delays in the order of days (after subtracting the deterministic travel times from the quoted delays), it is natural to measure the delays experienced in a hub in hours. Thus, the relevant performance criterion is whether queuing delays at a hub exceed several hours. For simplicity, we set $d_1 = d_2 = 1$ hour = 60 minutes in our simulation study for both types of firms. Finally, we need to set capacities for each firm to establish our base case. Taking $\varepsilon_1 = 0.992\%$ and $\varepsilon_2 = 2.683\%$, the corresponding service rates of the two servers (in the dedicated case, cf. Proposition 3) are $(\mu_1, \mu_2) = (162.60, 41.1)$ for the express firm, resulting in the utilization of 98.401% for the flexible server and 97.324% for the regular server. Similarly, we have $(\mu_1, \mu_2) = (41.40, 162.04)$ for the regular firm corresponding to $\varepsilon_1 = 0.005\%$ and $\varepsilon_2 = 18.57\%$, resulting in the utilization of 96.614% for the flexible server and 98.741% for the regular server. Although we were not able to obtain exact utilization numbers, Wright (2006) of UPS confirmed that the percentage utilization of their trucks and the airplanes are both "in the high nineties with the utilization of trucks being higher than that of the airplanes", which is consistent with our utilization numbers.

In this base case, we consider the question of whether a regular firm derives more value from in-

²⁰Unfortunately, no data on covariances is publicly available; nor were we able to get such data through our industry contacts.

	Class 1	Class 2
Arrival rate per minute	160	40
Standard deviation	104.03	38.2

Table 3 Hub-level demand characteristics for the express firm used for the simulation study.

	Class 1	Class 2
Arrival rate per minute	40	160
Standard deviation	26	152.8

Table 4 Hub-level demand characteristics for the regular firm used for the simulation study.

tegration, and test the predictive power of our analytic approximations. To this end, we report both the integration value calculated by our analytic formulas as well as that estimated by simulation. As a preliminary, we first report the service violation probability for each class for the dedicated network. We present both our analytic approximations and the estimates by simulation, which are summarized in Tables 5 and 6. Consider the express firm. The analytic approximations give the regular service violation probability for the integrated network as $\mathbb{P}^I(D_2 > d_2) = 2.683\%$, while the simulation estimate is $2.417\% \pm 0.003\%$, which are indeed close. Next, consider the regular firm. The analytic approximations give the regular service violation probability for the integrated system as $\mathbb{P}^I(D_2 > d_2) = 18.56\%$, while the simulation estimate is $18.633\% \pm 0.168\%$, which are also close to each other.

Having reported various probability estimates, we can now compare the value of integration for the two firms. We take $p_2 = 1$ for simplicity; and the value of p_1 is irrelevant for this comparison. First, consider the express firm. The analytic approximations give 0.986 as the integration value, while the simulation estimate of that is 0.927 ± 0.002 ; and the two differ by $5.971 \pm 0.337\%$. As for the regular firm, the analytic approximations gives 4.176 as the value of integration, while the simulation estimate is 4.199 ± 0.007 , and the two differ by $0.535 \pm 0.286\%$. For an alternative and arguably more insightful way to look at the same data, consider the value of integration suggested

	$\mathbb{P}(D_1 > d_1)$	$\mathbb{P}(D_2 > d_2)$
Analytic approximation	0.992%	2.683%
Simulation estimate	$0.746\% \pm .004\%$	$2.417\% \pm .003\%$

Table 5 Service violation probabilities in the dedicated network for the express firm.

	$\mathbb{P}(D_1 > d_1)$	$\mathbb{P}(D_2 > d_2)$
Analytic approximation	0.0048%	18.56%
Simulation estimate	$0.0075\% \pm .0009\%$	$18.633\% \pm 0.168\%$

Table 6 Service violation probabilities in the dedicated network for the regular firm.

by the analytic approximations. The integration value is 0.986 for the express firm and 4.176 for the regular firm, suggesting that the value of integration is higher for the regular firm. On the other hand, the simulation estimate of the integration value is 0.927 ± 0.002 for the express firm, and that for the regular firm is 4.199 ± 0.007 , also suggesting that the value of integration is higher for the regular firm in-line with the insight derived from our analytic approximations.

To validate our approximations further, we next consider various values of correlation. For each firm we first report the value of integration which is estimated both by our analytic approximations as well as by simulation; see Tables 7 and 8. Table 7 presents the results for an express firm, while Table 8 presents them for a regular firm.

Comparing Tables 7 and 8 we reach the following conclusions. First, for every correlation value both the simulation results and the analytic approximations suggest that integration value is higher (by an order of magnitude) for the regular firm. That is, the simulation study validates our main result derived from the analytic approximations. Moreover, for the regular firm the integration value calculated by analytic approximations is very accurate for all correlation values. For the express firm, the accuracy is good yet decreases as the correlation approaches zero. This may be explained by our analytic results which show that the value of integration is a small second order effect for FedEx except if correlation is very negative. Therefore, for other correlation values it is harder to pick up the very small effects. This problem does not exist for a regular firm (for which the value is first-order). In other words, the poorer accuracy results from estimating a second order effect. Nonetheless, the two estimates are sufficiently close that the predictions of simulation agrees with those of our analytical results.

We next report the regular service violation probability in the integrated network, and the tightness factor F (or the % improvement in the regular service violation probability through integration), cf. (28). Tables 9 and 10 present the regular service violation probability in an integrated network for express and regular firms, respectively. Tables 11 and 12 present the tightness factors for the express and regular firms. We see from Tables 9 and 10 that for the regular firm, even the probability estimates of our analytic approximations are fairly close to the estimates by

Correlation ρ	V (approximation)	V (simulation)	$\ V(\text{app})-V(\text{sim})\ /V(\text{sim})$
-.05	.558	.801 \pm .001	30.337% \pm 0.174%
-.10	.639	.821 \pm .001	22.168% \pm 0.190%
-.15	.716	.839 \pm .001	14.660% \pm 0.203%
-.20	.788	.858 \pm .002	8.159% \pm 0.428%
-.25	.853	.878 \pm .001	2.847% \pm 0.221%
-.30	.910	.897 \pm .002	1.449% \pm 0.227%
-.35	.954	.913 \pm .002	4.491% \pm 0.229%
-.40	.986	.927 \pm .002	6.365% \pm 0.230%
-.45	1.010	.944 \pm .002	6.992% \pm 0.227%
-.50	1.027	.957 \pm .002	7.315% \pm 0.225%
-.55	1.039	.968 \pm .002	7.335% \pm 0.222%
-.60	1.048	.976 \pm .002	7.377% \pm 0.220%
-.65	1.055	.984 \pm .003	7.215% \pm 0.328%
-.70	1.060	.992 \pm .003	6.855% \pm 0.324%
-.75	1.063	.998 \pm .003	6.513% \pm 0.321%
-.80	1.066	1.009 \pm .003	5.649% \pm 0.315%

Table 7 The estimates of value of integration (by analytic approximations and by simulation) for the express firm as a function of correlation.

Correlation ρ	V (approximation)	V (simulation)	$\ V(\text{app})-V(\text{sim})\ /V(\text{sim})$
-0.050	20.005	19.804 ± 0.164	$1.015\% \pm 0.844\%$
-0.100	20.459	20.228 ± 0.228	$1.142\% \pm 1.153\%$
-0.150	20.907	20.566 ± 0.172	$1.658\% \pm 0.857\%$
-0.200	21.347	20.912 ± 0.159	$2.080\% \pm 0.782\%$
-0.250	21.781	21.228 ± 0.208	$2.605\% \pm 1.015\%$
-0.300	22.206	21.662 ± 0.262	$2.511\% \pm 1.255\%$
-0.350	22.623	21.887 ± 0.201	$3.363\% \pm 0.958\%$
-0.400	23.032	22.217 ± 0.222	$3.668\% \pm 1.046\%$
-0.450	23.432	22.691 ± 0.308	$3.266\% \pm 1.421\%$
-0.500	23.822	22.998 ± 0.282	$3.583\% \pm 1.286\%$
-0.550	24.202	23.300 ± 0.216	$3.871\% \pm 0.972\%$
-0.600	24.572	23.593 ± 0.274	$4.150\% \pm 1.224\%$
-0.650	24.930	23.843 ± 0.197	$4.559\% \pm 0.871\%$
-0.700	25.278	24.246 ± 0.180	$4.256\% \pm 0.780\%$
-0.750	25.614	24.456 ± 0.158	$4.735\% \pm 0.681\%$
-0.800	25.938	24.689 ± 0.166	$5.059\% \pm 0.711\%$

Table 8 The estimates of value of integration (by analytic approximations and by simulation) for the regular firm as a function of correlation.

simulation. In contrast, for the express firm, it seems that the probability estimates from analytical approximations are off by a factor of approximately 2 to 3. While this may seem discouraging at first, it leads to the following insight: As long as the two estimates (the one by analytic approximations and the one by simulation) correspond to reductions of the same order of magnitude in the regular service violation probability, they will result in similar predictions of integration value. To see this, consider Table 11, which shows the tightness factor as a function of correlation for an express firm. As can be seen from Table 11, our analytic approximations give rise to tightness factors which are close to the simulation results; Table 12 displays similar (indeed, stronger) results for the regular firm. Given that the tightness factor is the determinant of integration value (cf. (28)), this is indeed very encouraging. Thus, we conclude that our analytic approximations can predict the value of integration fairly accurately even when the probability estimates for the express firm may be off by a factor of 2 to 3. Since our primary focus is the value of integration, we believe that the use of the approximate formulas we derived based on large deviations approximations is appropriate.

Finally, we test the validity of our analytic approximations as the server capacities (hence the server utilizations and the promised reliabilities ε) change. The insights gained so far carry over to this case as well. Tables 13, 14, 15, and 16 summarize the impact of these changes on various probability estimates, while Tables 17, 18, 19, and 20 report the impact of the changes on the estimates of integration value V and the tightness factor F .

Lognormal: All the simulation results above sampled bivariate normal random variables and used an ad-hoc method to correct for negative demand samples.²¹ To eliminate any errors from that procedure, we re-ran the entire simulation using the bivariate lognormal distribution. We report the corresponding results in Tables 21 to 26 below (and continuing at the end of the appendix). While the results are slightly better (worse) for the regular (fedex) firm, the key insights remain: the value of integration is much larger for a regular firm and our analytic results provide fairly good approximations.

²¹Any negative demand sample in a period was truncated to zero and kept in a “petty cash” account to adjust future demand. For example, say sample demand is -5 in period i and 100 in period $i + 1$ and petty cash was 0 in period i . Then the simulated demand was 0 and petty cash -5 in period i , and $100 - 5 = 95$ and petty cash = 0 in period $i + 1$. The petty cash ensures that the first moment remains unchanged by truncation.

Correlation ρ	$\mathbb{P}^I(D_2 > d_2)$ (approximation)	$\mathbb{P}^I(D_2 > d_2)$ (simulation)	$\ \mathbb{P}^I(\text{app}) - \mathbb{P}^I(\text{sim})\ / \mathbb{P}^I(\text{sim})$
-.05	1.287%	.348% \pm .002%	269.828% \pm 2.138%
-.10	1.086%	.298% \pm .001%	264.430% \pm 1.227%
-.15	.893%	.254% \pm .001%	251.575% \pm 1.390%
-.20	.713%	.213% \pm .000%	234.742% \pm 0.000%
-.25	.550%	.178% \pm .000%	208.989% \pm 0.000%
-.30	.408%	.146% \pm .000%	179.452% \pm 0.000%
-.35	.298%	.120% \pm .000%	148.333% \pm 0.000%
-.40	.218%	.099% \pm .000%	120.202% \pm 0.000%
-.45	.159%	.079% \pm .000%	101.266% \pm 0.000%
-.50	.116%	.062% \pm .000%	87.097% \pm 0.000%
-.55	.085%	.046% \pm .000%	84.783% \pm 0.000%
-.60	.062%	.034% \pm .000%	82.353% \pm 0.000%
-.65	.045%	.026% \pm .000%	73.077% \pm 0.000%
-.70	.033%	.019% \pm .000%	73.684% \pm 0.000%
-.75	.024%	.015% \pm .000%	60.000% \pm 0.000%
-.80	.018%	.012% \pm .000%	50.000% \pm 0.000%

Table 9 The estimates of regular service violation probability (by analytic approximations and by simulation) in an integrated network for the express firm as a function of correlation.

Correlation ρ	$\mathbb{P}^I(D_2 > d_2)$ (approximation)	$\mathbb{P}^I(D_2 > d_2)$ (simulation)	$\ \mathbb{P}^I(\text{app}) - \mathbb{P}^I(\text{sim})\ / \mathbb{P}^I(\text{sim})$
-0.050	6.064%	5.577% \pm 0.070%	8.732% \pm 1.382%
-0.100	5.781%	5.322% \pm 0.058%	8.625% \pm 1.197%
-0.150	5.501%	5.109% \pm 0.057%	7.673% \pm 1.215%
-0.200	5.225%	4.882% \pm 0.070%	7.026% \pm 1.557%
-0.250	4.955%	4.636% \pm 0.058%	6.881% \pm 1.354%
-0.300	4.689%	4.408% \pm 0.071%	6.375% \pm 1.741%
-0.350	4.428%	4.215% \pm 0.059%	5.053% \pm 1.491%
-0.400	4.173%	3.992% \pm 0.059%	4.534% \pm 1.568%
-0.450	3.923%	3.770% \pm 0.074%	4.058% \pm 2.083%
-0.500	3.679%	3.548% \pm 0.071%	3.692% \pm 2.117%
-0.550	3.441%	3.362% \pm 0.068%	2.350% \pm 2.113%
-0.600	3.210%	3.169% \pm 0.062%	1.294% \pm 2.021%
-0.650	2.986%	2.976% \pm 0.047%	0.336% \pm 1.610%
-0.700	2.769%	2.775% \pm 0.043%	0.216% \pm 1.138%
-0.750	2.559%	2.607% \pm 0.038%	1.841% \pm 2.862%
-0.800	2.356%	2.430% \pm 0.035%	3.045% \pm 2.794%

Table 10 The estimates of regular service violation probability (by analytic approximations and by simulation) in an integrated network for the regular firm as a function of correlation.

Correlation ρ	F (approximation)	F (simulation)	$\ F(\text{app})-F(\text{sim})\ /F(\text{sim})$
-.05	.520	.852 \pm .000	38.967% \pm 0.000%
-.10	.595	.873 \pm .000	31.844% \pm 0.000%
-.15	.667	.892 \pm .000	25.224% \pm 0.000%
-.20	.734	.910 \pm .000	19.341% \pm 0.000%
-.25	.795	.925 \pm .000	14.054% \pm 0.000%
-.30	.848	.939 \pm .000	9.691% \pm 0.000%
-.35	.889	.950 \pm .000	6.421% \pm 0.000%
-.40	.919	.959 \pm .000	4.171% \pm 0.000%
-.45	.941	.967 \pm .000	2.689% \pm 0.000%
-.50	.957	.975 \pm .000	1.846% \pm 0.000%
-.55	.968	.981 \pm .000	1.325% \pm 0.000%
-.60	.977	.986 \pm .000	0.913% \pm 0.000%
-.65	.983	.989 \pm .000	0.607% \pm 0.000%
-.70	.988	.992 \pm .000	0.403% \pm 0.000%
-.75	.991	.994 \pm .000	0.302% \pm 0.000%
-.80	.993	.995 \pm .000	0.201% \pm 0.000%

Table 11 The estimates of the tightness factor F (by analytic approximations and by simulation) in an integrated network for the express firm as a function of correlation.

Correlation ρ	F (approximation)	F (simulation)	$\ F(\text{app})-F(\text{sim})\ /F(\text{sim})$
-0.050	0.673	0.689 ± 0.003	$2.322\% \pm 0.851\%$
-0.100	0.689	0.704 ± 0.004	$2.131\% \pm 1.112\%$
-0.150	0.704	0.716 ± 0.003	$1.676\% \pm 0.824\%$
-0.200	0.719	0.728 ± 0.004	$1.236\% \pm 1.085\%$
-0.250	0.733	0.741 ± 0.004	$1.080\% \pm 1.068\%$
-0.300	0.747	0.754 ± 0.005	$0.928\% \pm 1.314\%$
-0.350	0.762	0.764 ± 0.004	$0.262\% \pm 0.518\%$
-0.400	0.775	0.777 ± 0.004	$0.257\% \pm 0.510\%$
-0.450	0.789	0.790 ± 0.005	$0.127\% \pm 0.383\%$
-0.500	0.802	0.802 ± 0.005	$0.000\% \pm 0.627\%$
-0.550	0.815	0.812 ± 0.004	$0.369\% \pm 0.497\%$
-0.600	0.827	0.823 ± 0.004	$0.486\% \pm 0.491\%$
-0.650	0.839	0.833 ± 0.003	$0.720\% \pm 0.364\%$
-0.700	0.851	0.845 ± 0.003	$0.710\% \pm 0.359\%$
-0.750	0.862	0.854 ± 0.002	$0.937\% \pm 0.237\%$
-0.800	0.873	0.864 ± 0.002	$1.042\% \pm 0.234\%$

Table 12 The estimates of the tightness factor F (by analytic approximations and by simulation) in an integrated network for the regular firm as a function of correlation.

μ_1	161.40	161.80	162.60	163.80	165.00
Server 1's Utilization	99.133%	98.888%	98.401%	97.680%	96.970%
$\mathbb{P}(D_1 > d_1)$ (approximation)	8.343%	4.103%	.992%	.118%	.014%
$\mathbb{P}(D_1 > d_1)$ (simulation)	$6.972 \pm .016\%$	$3.317 \pm .008\%$	$.746 \pm .004\%$	$.085 \pm .000$	$.007 \pm .000\%$
$\mathbb{P}^I(D_2 > d_2)$ (approximation)	.694%	.471%	.218%	.068%	.021%
$\mathbb{P}^I(D_2 > d_2)$ (simulation)	$.371 \pm .000\%$	$.225 \pm .000\%$	$.099 \pm .000\%$	$.024 \pm .000\%$	$.007 \pm .000\%$

Table 13 The impact of changing server 1 processing rate μ_1 on the delay probability estimates for an express firm (keeping the server 2 processing rate constant at $\mu_2 = 41.1$).

μ_2	40.7	40.9	41.1	41.8	42.5
Server 2's Utilization	98.280%	97.800%	97.324%	95.694%	94.118%
$\mathbb{P}(D_2 > d_2)$ (approximation)	10.000%	5.180%	2.683%	.268%	.027%
$\mathbb{P}(D_2 > d_2)$ (simulation)	$9.890 \pm .006\%$	$4.963 \pm .003\%$	$2.417 \pm .003\%$	$.208 \pm .001\%$	$.013 \pm .000\%$
$\mathbb{P}^I(D_2 > d_2)$ (approximation)	.811%	.420%	.218%	.022%	.002%
$\mathbb{P}^I(D_2 > d_2)$ (simulation)	$.328 \pm .001\%$	$.181 \pm .000\%$	$.099 \pm .000\%$	$.009 \pm .000\%$	$.001 \pm .000\%$

Table 14 The impact of changing server 2 processing rate μ_2 on the delay probability estimates for an express firm (keeping the server 1 processing rate constant at $\mu_1 = 162.6$).

μ_1	40.75	40.95	41.40	42.00	42.50
Server 1's Utilization	98.160%	97.680%	96.618%	95.238%	94.118%
$\mathbb{P}(D_1 > d_1)$ (approximation)	8.484%	4.394%	.993%	.139%	.027%
$\mathbb{P}(D_1 > d_1)$ (simulation)	$7.300 \pm .007\%$	$3.668 \pm .007\%$	$.744 \pm .002\%$	$.080 \pm .000\%$	$.010 \pm .000\%$
$\mathbb{P}^I(D_2 > d_2)$ (approximation)	.278%	.182%	.070%	.020%	.007%
$\mathbb{P}^I(D_2 > d_2)$ (simulation)	$.329 \pm .001\%$	$.195 \pm .001\%$	$.074 \pm .001\%$	$.017 \pm .000\%$	$.003 \pm .000\%$

Table 15 The impact of changing server 1 processing rate μ_1 on the delay probability estimates for a regular firm (keeping the server 2 processing rate constant at $\mu_2 = 162.04$).

μ_2	161.30	161.70	162.04	163.00	163.30
Server 2's Utilization	99.194%	98.949%	98.741%	98.160%	97.979%
$\mathbb{P}(D_2 > d_2)$ (approx.)	9.962%	4.900%	2.680%	.488%	.287%
$\mathbb{P}(D_2 > d_2)$ (simulation)	$10.245 \pm .010\%$	$5.014 \pm .004\%$	$2.698 \pm .003\%$	$.461 \pm .002\%$	$.271 \pm .002\%$
$\mathbb{P}^I(D_2 > d_2)$ (approx.)	.335%	.144%	.070%	.009%	.005%
$\mathbb{P}^I(D_2 > d_2)$ (simulation)	$.333 \pm .001\%$	$.142 \pm .001\%$	$.074 \pm .001\%$	$.008 \pm .000\%$	$.002 \pm .000\%$

Table 16 The impact of changing server 2 processing rate μ_2 on the delay probability estimates for a regular firm (keeping the server 1 processing rate constant at $\mu_1 = 41.4$).

μ_1	161.40	161.80	162.60	163.80	165
V (approximation)	.796	.885	.986	1.046	1.065
V (simulation)	.818 \pm .002	.876 \pm .002	.927 \pm .002	.957 \pm .002	.965 \pm .002
% difference	2.825 \pm .383	.922 \pm .367	5.971 \pm .337	8.471 \pm .320	9.360 \pm .315
F (approximation)	.741	.824	.919	.975	.992
F (simulation)	.846 \pm .000	.906 \pm .000	.959 \pm .000	.990 \pm .000	.998 \pm .000
% difference	14.150 \pm .059	9.990 \pm .042	4.385 \pm .021	1.610 \pm .007	.622 \pm .002

Table 17 The impact of changing server 1 processing rate μ_1 on the estimates of integration value V and the tightness factor F for an express firm (keeping the server 2 processing rate constant at $\mu_2 = 41.1$).

μ_2	40.70	40.90	41.10	41.80	42.50
V (approximation)	3.676	1.904	.986	.099	.010
V (simulation)	3.825 \pm .004	1.913 \pm .002	.927 \pm .002	.080 \pm .001	.005 \pm .000
% difference	4.053 \pm .180	.470 \pm .147	5.971 \pm .337	19.205 \pm 1.257	50.646 \pm .375
F (approximation)	91.892%	91.892%	91.892%	91.892%	91.892%
F (simulation)	96.679 \pm .004%	96.354 \pm .001%	95.921 \pm .003%	95.715 \pm .017%	95.287 \pm .048%
% difference	5.209 \pm .024%	4.855 \pm .006%	4.385 \pm .021%	4.160 \pm .110%	3.694 \pm .307%

Table 18 The impact of changing server 2 processing rate μ_2 on the estimates of integration value V and the tightness factor F for an express firm (keeping the server 1 processing rate constant at $\mu_1 = 162.6$).

μ_1	40.75	40.95	41.40	42.00	42.50
V (approximation)	3.844	3.997	4.176	4.257	4.278
V (simulation)	3.790 \pm .007	4.005 \pm .007	4.199 \pm .007	4.290 \pm .007	4.312 \pm .007
% difference	1.406 \pm .319	.200 \pm .298	.535 \pm .286	.771 \pm .288	.809 \pm .292
F (approximation)	.896	.932	.974	.993	.997
F (simulation)	.878 \pm .000	.928 \pm .000	.973 \pm .000	.994 \pm .000	.999 \pm .000
% difference	2.035 \pm .103	.439 \pm .080	.106 \pm .045	.129 \pm .021	.166 \pm .004

Table 19 The impact of changing server 1 processing rate μ_1 on the estimates of integration value V and the tightness factor F for a regular firm (keeping the server 2 processing rate constant at $\mu_2 = 162.04$).

μ_2	161.30	161.70	162.04	163.00	163.30
V (approximation)	15.404	7.609	4.176	.766	.451
V (simulation)	$15.858 \pm .025$	$7.794 \pm .009$	$4.199 \pm .007$	$.726 \pm .005$	$.430 \pm .004$
% difference	$2.948 \pm .297$	$2.437 \pm .205$	$.535 \pm .286$	5.230 ± 1.150	4.606 ± 1.628
F (approximation)	.966	.971	.974	.981	.983
F (simulation)	$.967 \pm .000$	$.972 \pm .000$	$.973 \pm .000$	$.984 \pm .000$	$.991 \pm .000$
% difference	$.111 \pm .042$	$.109 \pm .044$	$.106 \pm .045$	$.275 \pm .059$	$.836 \pm .027$

Table 20 The impact of changing server 2 processing rate μ_2 on the estimates of integration value V and the tightness factor F for a regular firm (keeping the server 1 processing rate constant at $\mu_1 = 41.4$).

Correlation ρ	V (approximation)	V (simulation)	$\ V(\text{app})-V(\text{sim})\ /V(\text{sim})$
-0.050	0.558	0.974 ± 0.015	$42.710\% \pm 1.765\%$
-0.100	0.639	1.000 ± 0.015	$36.100\% \pm 1.917\%$
-0.150	0.716	1.024 ± 0.011	$30.078\% \pm 1.502\%$
-0.200	0.788	1.036 ± 0.018	$23.938\% \pm 2.644\%$
-0.250	0.853	1.067 ± 0.016	$20.056\% \pm 2.398\%$
-0.300	0.910	1.084 ± 0.015	$16.052\% \pm 2.324\%$
-0.350	0.954	1.098 ± 0.014	$13.115\% \pm 2.216\%$
-0.400	0.986	1.116 ± 0.014	$11.649\% \pm 2.217\%$
-0.450	1.010	1.124 ± 0.011	$10.142\% \pm 1.759\%$
-0.500	1.027	1.127 ± 0.015	$8.873\% \pm 2.426\%$
-0.550	1.039	1.140 ± 0.017	$8.860\% \pm 2.719\%$
-0.600	1.048	1.147 ± 0.018	$8.631\% \pm 2.868\%$

Table 21 [Like TABLE 7 but with lognormal demand]The estimates of value of integration (by analytic approximations and by simulation) for the express firm as a function of correlation.

E OnlineAppendix: Proofs in Section 8

Proof of Proposition 8. The first order conditions for the capacity optimization problem (32) gives

$$\lambda_i p_i \frac{\partial}{\partial \mu_i} \mathbb{P}^{ded}(D_i > d_i) + c_i = 0. \quad (45)$$

By Proposition 1, we have that

$$\frac{\partial}{\partial \mu_i} \mathbb{P}^{ded}(D_i > d_i) = -2 \frac{\lambda_i}{\sigma_i^2} d_i \mathbb{P}^{ded}(D_i > d_i).$$

Substituting this into (45) gives

$$\mathbb{P}^{ded}(D_i > d_i) = \frac{c_i}{2p_i d_i} \frac{\sigma_i}{\lambda_i},$$

from which we have that

$$z_i^{ded} = \frac{1}{2d_i} \frac{\sigma_i}{\lambda_i} \log \left(\frac{2p_i d_i}{c_i} \left(\frac{\lambda_i}{\sigma_i} \right)^2 \right). \quad \blacksquare$$

The proofs of Propositions 9 and 10 assume that the cost minimized in the capacity optimization problem (32) is (strictly) convex, which is easy to check for the dedicated network. Under this assumption, the first order conditions are necessary and sufficient. However, since there are eight different cases for calculating the regular service failure probability (as can be seen in Appendix A.3) depending on problem parameters, in particular, the optimal capacity choices, solving the first order conditions can be tricky. To be specific, it requires identifying which one of the eight different regimes is the relevant one under the optimal capacity decisions, which in turn depends on the particular regime. Thus, solving the first order conditions analytically does not seem tractable in general. Luckily, the analysis is simpler for the canonical firm (eg. express and regular), although it still involves guessing the parameter regime.

To build intuition, note that combining the analysis of the dedicated network with the arrival pooling effect, we see that z_1/z_2 is small for an express firm while it is large for a regular firm. Although it is not obvious that this will be the case in the integrated network apriori, we will show that it is indeed the case. In what follows, we will solve the first order conditions only for the canonical firms (eg. express and regular firms), and believe that deriving a solution in general is not analytically tractable. In the case of canonical firm, we initially guess that the solution to the first order conditions have the following properties: $z > 0$ and z_1/z_2 is small for an express firm while it is large for a regular firm. Then we will show that the solution we derive satisfies these assumptions.

Proof of Proposition 9. First observe that $z^{int} > 0$ and z_1^{int}/z_2^{int} is small since $\lambda_1 \gg \lambda_2$ for an express firm. Since σ_1/σ_2 is also large for an express firm, it is easy to see that we are in case 1 of Proposition 1 (in a neighborhood of z^{int}). That is,

$$\mathbb{P}^{int}(D_2 > d_2) = \exp \left\{ -r \frac{\lambda_2}{\sigma_2} d_2 \right\},$$

where $r = 2z_2^{int} - 4\rho z_1^{int}$. Then note that the first order conditions in the capacity optimization problem are as follows:

$$\lambda_1 p_1 \frac{\partial}{\partial \mu_1} \mathbb{P}^{int}(D_1 > d_1) + c_1 + \lambda_2 p_2 \frac{\partial}{\partial \mu_1} \mathbb{P}^{int}(D_2 > d_2) = 0, \quad (46)$$

$$\lambda_2 p_2 \frac{\partial}{\partial \mu_2} \mathbb{P}^{int}(D_2 > d_2) = 0. \quad (47)$$

Recall that $\mathbb{P}^{int}(D_1 > d_1) = \exp \left\{ -2 \frac{\lambda_1}{\sigma_1} z_1 d_1 \right\}$. Then it is easy to see that

$$\frac{\partial}{\partial \mu_1} \mathbb{P}^{int}(D_1 > d_1) = -\frac{2\lambda_1 p_1}{\sigma_1^2} \mathbb{P}^{int}(D_1 > d_1). \quad (48)$$

Similarly,

$$\frac{\partial}{\partial \mu_2} \mathbb{P}^{int}(D_2 > d_2) = \frac{\partial}{\partial r} \mathbb{P}^{int}(D_2 > d_2) \frac{\partial r}{\partial z_2} \frac{\partial z_2}{\partial \mu_2} = -\frac{2\lambda_2 p_2}{\sigma_2^2} \mathbb{P}^{int}(D_2 > d_2), \quad (49)$$

$$\frac{\partial}{\partial \mu_1} \mathbb{P}^{int}(D_2 > d_2) = \frac{\partial}{\partial r} \mathbb{P}^{int}(D_2 > d_2) \frac{\partial r}{\partial z_1} \frac{\partial z_1}{\partial \mu_1} = \frac{4\rho\lambda_2 d_2}{\sigma_1 \sigma_2} \mathbb{P}^{int}(D_2 > d_2). \quad (50)$$

We will show that the first order conditions are satisfied, but before checking the first order conditions (46)-(47), note that

$$\mathbb{P}^{int}(D_1 > d_1) = \exp \left\{ -2 \frac{\lambda_1}{\sigma_1} z_1^{int} d_1 \right\} = \frac{1}{2p_1 d_1} \left(\frac{\sigma_1}{\lambda_1} \right)^2 \left[c_1 + 2\rho c_2 \frac{\sigma_2}{\sigma_1} \right]. \quad (51)$$

Similarly,

$$\mathbb{P}^{int}(D_2 > d_2) = \exp \left\{ -r \frac{\lambda_2}{\sigma_2} d_2 \right\} = \frac{c_2}{2p_2 d_2} \left(\frac{\sigma_2}{\lambda_2} \right)^2. \quad (52)$$

Then combining (46)-(52) it is straightforward to show that the first order conditions (46)-(47) are satisfied by z^{int} . ■

Proof of Proposition 10. First, observe that $z^{int} > 0$ and z_1^{int}/z_2^{int} is large since $\lambda_1 \ll \lambda_2$ for a regular firm. Since σ_1/σ_2 is also small for a regular firm, it is easy to see from Proposition 1 that

$$\mathbb{P}^{int}(D_2 > d_2) = \exp \left\{ -r \frac{\lambda_2}{\sigma_2} d_2 \right\}$$

in a neighborhood of z^{int} where

$$r = r_4 = 2 \frac{\sigma_1 \sigma_2}{\sigma_T^2} z_1 + 2 \frac{\sigma_2^2}{\sigma_T^2} z_2.$$

Recall that the first order conditions for the capacity optimization problem are as follows:

$$\begin{aligned}\lambda_1 p_1 \frac{\partial}{\partial \mu_1} \mathbb{P}^{int}(D_1 > d_1) + c_1 + \lambda_2 p_2 \frac{\partial}{\partial \mu_1} \mathbb{P}^{int}(D_2 > d_2) &= 0, \\ \lambda_2 p_2 \frac{\partial}{\partial \mu_2} \mathbb{P}^{int}(D_2 > d_2) &= 0.\end{aligned}$$

Also recall that

$$\mathbb{P}^{int}(D_1 > d_1) = \exp \left\{ -2 \frac{\lambda_1}{\sigma_1} z_1 d_1 \right\}.$$

Then it is easy to see that

$$\frac{\partial}{\partial \mu_1} \mathbb{P}^{int}(D_1 > d_1) = -\frac{2\lambda_1 p_1}{\sigma_1^2} \mathbb{P}^{int}(D_1 > d_1). \quad (53)$$

Similarly,

$$\frac{\partial}{\partial \mu_2} \mathbb{P}^{int}(D_2 > d_2) = \frac{\partial}{\partial r} \mathbb{P}^{int}(D_2 > d_2) \frac{\partial r}{\partial z_2} \frac{\partial z_2}{\partial \mu_2} = -\frac{2d_2 \lambda_2}{\sigma_T^2} \mathbb{P}^{int}(D_2 > d_2), \quad (54)$$

$$\frac{\partial}{\partial \mu_1} \mathbb{P}^{int}(D_2 > d_2) = \frac{\partial}{\partial r} \mathbb{P}^{int}(D_2 > d_2) \frac{\partial r}{\partial z_1} \frac{\partial z_1}{\partial \mu_1} = -\frac{2d_2 \lambda_2}{\sigma_T^2} \mathbb{P}^{int}(D_2 > d_2). \quad (55)$$

We will show that the first order conditions are satisfied, but before doing so note that

$$\mathbb{P}^{int}(D_1 > d_1) = \exp \left\{ -2 \frac{\lambda_1}{\sigma_1} z_1^{int} d_1 \right\} = \frac{c_1 - c_2}{2p_1 d_1} \left(\frac{\sigma_1}{\lambda_1} \right)^2, \quad (56)$$

$$\mathbb{P}^{int}(D_2 > d_2) = \exp \left\{ -r \frac{\lambda_2}{\sigma_2} d_2 \right\} = \exp \left\{ -\frac{\lambda_2}{\sigma_2} d_2 \left[2 \frac{\sigma_1 \sigma_2}{\sigma_T^2} z_1^{int} + 2 \frac{\sigma_2^2}{\sigma_T^2} z_2^{int} \right] \right\} = \frac{c_2}{2p_2 d_2} \left(\frac{\sigma_T}{\lambda_2} \right)^2 \quad (57)$$

Then combining (53)-(57), it is straightforward to check that the first order conditions are satisfied by z^{int} . ■

Proof of Theorem 3. We will attach a superscript e (r) to denote the express (regular) firm to avoid confusion as needed. For a meaningful comparison of the express and regular firms, let $\lambda_1^e + \lambda_2^e = \lambda_1^r + \lambda_2^r$ be large with $\lambda_2^e = \lambda_1^r = O(1)$. Let C denote the cost rate for an optimally designed dedicated network and C^e denote the corresponding cost rate for an integrated express firm. It follows from Proposition 8 that

$$C = \sum_{i=1}^2 c_i \lambda_i + \sum_{i=1}^2 \frac{c_i}{2d_i} \frac{\sigma_i^2}{\lambda_i} \left[1 + \log \left(\frac{2p_i d_i}{c_i} \left(\frac{\lambda_i}{\sigma_i} \right)^2 \right) \right].$$

Similarly, it follows from Proposition 9 that

$$\begin{aligned}C^e &= \sum_{i=1}^2 \left(c_i \lambda_i + \frac{c_i}{2d_i} \frac{\sigma_i^2}{\lambda_i} \right) + \rho \frac{c_2}{d_1} \frac{\sigma_1}{\lambda_1} \sigma_2 + \frac{c_1}{2d_1} \frac{\sigma_1^2}{\lambda_1} \log \left(\frac{2p_1 d_1}{c_1 + 2\rho c_2 \frac{\sigma_2}{\sigma_1}} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) \\ &\quad + \rho \frac{c_2}{d_1} \frac{\sigma_1}{\lambda_1} \sigma_2 \log \left(\frac{2p_1 d_1}{c_1 + 2\rho c_2 \frac{\sigma_2}{\sigma_1}} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) + \frac{c_2}{2d_2} \frac{\sigma_2^2}{\lambda_2} \log \left(\frac{2p_2 d_2}{c_2} \left(\frac{\lambda_2}{\sigma_2} \right)^2 \right).\end{aligned}$$

The value of integration for the express firm $V^e = C - C^e$ is then given by

$$V^e = \frac{c_1}{2d_1} \frac{\sigma_1^2}{\lambda_1} \log \left(1 + 2\rho \frac{c_2 \sigma_2}{c_1 \sigma_1} \right) - \rho \frac{c_2 \sigma_1}{d_1 \lambda_1} \sigma_2 - \rho \frac{c_2 \sigma_1}{d_1 \lambda_1} \sigma_2 \log \left(\frac{2p_1 d_1}{c_1 + 2\rho c_2 \frac{\sigma_2}{\sigma_1}} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right).$$

Recall that $\log(1+x) \approx x - x^2/2$ for small x . Then since σ_2/σ_1 is small for an express firm, we conclude that

$$V^e \approx \frac{c_1}{2d_1} \frac{\sigma_1^2}{\lambda_1} \left[2\rho \frac{c_1 \sigma_2}{c_2 \sigma_1} - 2\rho^2 \left(\frac{c_2 \sigma_2}{c_1 \sigma_1} \right)^2 - \rho \frac{c_2 \sigma_1}{d_1 \lambda_1} \sigma_2 - \rho \frac{c_2 \sigma_1}{d_1 \lambda_1} \sigma_2 \log \left(\frac{2p_1 d_1}{c_1 + 2\rho c_2 \frac{\sigma_2}{\sigma_1}} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) \right].$$

That is,

$$V^e \approx -\rho \frac{c_2 \sigma_1}{d_1 \lambda_1} \sigma_2 \log \left(\frac{2p_1 d_1}{c_1 + 2\rho c_2 \frac{\sigma_2}{\sigma_1}} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) - \rho^2 \frac{c_2^2}{d_1 c_1} \frac{\sigma_2^2}{\lambda_1}.$$

Notice that $V^e = O(\log(\lambda_1^e)/(\lambda_1^e)^{1-\gamma_1})$ which is small, i.e. $o(1)$, since λ_1^e is large.

The analysis for the regular firm is similar. Letting C^r denote the cost rate for an integrated regular firm, it follows from Proposition 10 that

$$\begin{aligned} C^r &= \frac{c_1 - c_2}{2d_1} \frac{\sigma_1^2}{\lambda_1} + c_1 \left(\lambda_1 + \frac{1}{2d_1} \frac{\sigma_1^2}{\lambda_1} \log \left(\frac{2p_1 d_1}{c_1 - c_2} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) \right) \\ &\quad + c_2 \left\{ \lambda_2 + \frac{\sigma_T^2}{2\lambda_2 d_2} \log \left(\frac{2p_2 d_2}{c_2} \frac{\lambda_2^2}{\sigma_T^2} \right) - \frac{\sigma_1^2}{2d_1 \lambda_1} \log \left(\frac{2p_1 d_1}{c_1 - c_2} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) \right\}. \end{aligned}$$

Then the value of integration $V^r = C - C^r$ for the regular firm is given by

$$\begin{aligned} V^r &= \frac{c_2}{2d_1} \frac{\sigma_1^2}{\lambda_1} + \frac{c_1}{2d_1} \frac{\sigma_1^2}{\lambda_1} \log \left(1 - \frac{c_2}{c_1} \right) + \frac{c_2}{2d_2 \lambda_2} (\sigma_2^2 - \sigma_T^2) + \frac{c_2}{2d_1} \frac{\sigma_1^2}{\lambda_1^2} \log \left(\frac{2p_1 d_1}{c_1 - c_2} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) \\ &\quad + \frac{\lambda_2 c_2}{2d_2} \left[\frac{\sigma_2^2}{\lambda_2^2} \log \left(\frac{2p_2 d_2}{c_2} \frac{\lambda_2^2}{\sigma_2^2} \right) - \frac{\sigma_T^2}{\lambda_2^2} \log \left(\frac{2p_2 d_2}{c_2} \frac{\lambda_2^2}{\sigma_T^2} \right) \right]. \end{aligned}$$

We assume $c_2/c_1 \ll 1$, which is a reasonable assumption in most settings including the FedEx-UPS example. Then

$$\log \left(1 - \frac{c_2}{c_1} \right) \approx -\frac{c_2}{c_1} - \frac{1}{2} \left(\frac{c_2}{c_1} \right)^2.$$

Thus,

$$\begin{aligned} V^r &\approx \frac{c_2}{2d_1} \frac{\sigma_1^2}{\lambda_1} \left[\log \left(\frac{2p_1 d_1}{c_1 - c_2} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) - \frac{1}{2} \frac{c_2}{c_1} \right] + \frac{c_2}{2\lambda_2 d_2} (\sigma_2^2 - \sigma_T^2) \\ &\quad + \frac{\lambda_2 c_2}{2d_2} \left[\frac{\sigma_2^2}{\lambda_2^2} \log \left(\frac{2p_2 d_2}{c_2} \frac{\lambda_2^2}{\sigma_2^2} \right) - \frac{\sigma_T^2}{\lambda_2^2} \log \left(\frac{2p_2 d_2}{c_2} \frac{\lambda_2^2}{\sigma_T^2} \right) \right]. \end{aligned}$$

Letting $g(x) = x \log(\alpha/x)$ with $\alpha = 2p_2 d_2/c_2$, we have that

$$\frac{\sigma_2^2}{\lambda_2^2} \log \left(\frac{2p_2 d_2}{c_2} \frac{\lambda_2^2}{\sigma_2^2} \right) - \frac{\sigma_T^2}{\lambda_2^2} \log \left(\frac{2p_2 d_2}{c_2} \frac{\lambda_2^2}{\sigma_T^2} \right) = g(x_1) - g(x_2),$$

where $x_1 = (\sigma_2/\lambda_2)^2 > x_2 = (\sigma_2/\lambda_2)^2$. Then by a simple Taylor's expansion

$$g(x_1) - g(x_2) \approx g'(x_1)(x_1 - x_2) = \left[\log \left(\frac{2p_2 d_2}{c_2} \left(\frac{\lambda_2}{\sigma_2} \right)^2 \right) - 1 \right] \frac{\sigma_2^2 - \sigma_T^2}{\lambda_2^2}.$$

Then

$$\begin{aligned} V^r \approx & \frac{c_2}{2d_1} \frac{\sigma_1^2}{\lambda_1} \left[\log \left(\frac{2p_1 d_1}{c_1 - c_2} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) - \frac{c_2}{2c_1} \right] - \rho \frac{c_2}{d_2} \frac{\sigma_2}{\lambda_2} \sigma_1 - \frac{c_2 \sigma_1^2}{4d_2 \lambda_2} \\ & - \rho \frac{c_2}{2d_2} \frac{\sigma_2}{\lambda_2} \left[\log \left(\frac{2p_2 d_2}{c_2} \left(\frac{\lambda_2}{\sigma_2} \right)^2 \right) - 1 \right] - \frac{c_2}{2d_2} \frac{\sigma_1^2}{\lambda_2} \left[\log \left(\frac{2p_2 d_2}{c_2} \left(\frac{\lambda_2}{\sigma_2} \right)^2 \right) - 1 \right]. \end{aligned}$$

Since all but the last term on the right are small, we write

$$V^r \approx \frac{c_2}{2d_1} \frac{\sigma_1^2}{\lambda_1} \left[\log \left(\frac{2p_1 d_1}{c_1 - c_2} \left(\frac{\lambda_1}{\sigma_1} \right)^2 \right) - \frac{c_2}{2c_1} \right],$$

and thus, we conclude that $V^r = O(1)$ since $\lambda_1^r = O(1)$. Therefore,

$$V^r \gg V^e,$$

proving that the value of integration is higher for the regular firm. The assertion regarding the relative value of integration follows along similar lines. ■

Correlation ρ	V (approximation)	V (simulation)	$\ V(\text{app})-V(\text{sim})\ /V(\text{sim})$
-0.050	20.005	19.820 ± 0.292	$0.933\% \pm 1.509\%$
-0.100	20.459	20.131 ± 0.289	$1.629\% \pm 1.480\%$
-0.150	20.907	20.387 ± 0.249	$2.551\% \pm 1.268\%$
-0.200	21.347	20.691 ± 0.236	$3.170\% \pm 1.190\%$
-0.250	21.781	21.190 ± 0.275	$2.789\% \pm 1.352\%$
-0.300	22.206	21.507 ± 0.301	$3.250\% \pm 1.466\%$
-0.350	22.623	21.889 ± 0.261	$3.353\% \pm 1.247\%$
-0.400	23.032	22.239 ± 0.192	$3.566\% \pm 0.902\%$
-0.450	23.432	22.478 ± 0.170	$4.244\% \pm 0.794\%$
-0.500	23.822	22.885 ± 0.116	$4.094\% \pm 0.530\%$
-0.550	24.202	23.171 ± 0.100	$4.450\% \pm 0.453\%$
-0.600	24.572	23.496 ± 0.151	$4.580\% \pm 0.676\%$

Table 22 [Like TABLE 8 but with lognormal demand]The estimates of value of integration (by analytic approximations and by simulation) for the regular firm as a function of correlation.

Correlation ρ	$\mathbb{P}^I(D_2 > d_2)$ (approximation)	$\mathbb{P}^I(D_2 > d_2)$ (simulation)	$\ \mathbb{P}^I(\text{app}) - \mathbb{P}^I(\text{sim})\ / \mathbb{P}^I(\text{sim})$
-0.050	1.287%	0.484% \pm 0.010%	165.909% \pm 5.610%
-0.100	1.086%	0.424% \pm 0.010%	156.132% \pm 6.187%
-0.150	0.893%	0.361% \pm 0.007%	147.368% \pm 4.891%
-0.200	0.713%	0.308% \pm 0.009%	131.494% \pm 6.968%
-0.250	0.550%	0.257% \pm 0.008%	114.008% \pm 6.876%
-0.300	0.407%	0.216% \pm 0.009%	88.426% \pm 8.192%
-0.350	0.297%	0.179% \pm 0.007%	65.922% \pm 6.753%
-0.400	0.217%	0.151% \pm 0.006%	43.709% \pm 5.947%
-0.450	0.159%	0.124% \pm 0.005%	28.226% \pm 5.388%
-0.500	0.116%	0.105% \pm 0.004%	10.476% \pm 4.375%
-0.550	0.085%	0.089% \pm 0.002%	4.494% \pm 4.295%
-0.600	0.062%	0.076% \pm 0.002%	18.421% \pm 4.297%

Table 23 [Like TABLE 9 but with lognormal demand]The estimates of regular service violation probability (by analytic approximations and by simulation) in an integrated network for the express firm as a function of correlation.

Correlation ρ	$\mathbb{P}^I(D_2 > d_2)$ (approximation)	$\mathbb{P}^I(D_2 > d_2)$ (simulation)	$\ \mathbb{P}^I(\text{app}) - \mathbb{P}^I(\text{sim})\ / \mathbb{P}^I(\text{sim})$
-0.050	6.064%	6.246% \pm 0.092%	2.914% \pm 2.861%
-0.100	5.781%	6.009% \pm 0.083%	3.794% \pm 2.658%
-0.150	5.501%	5.801% \pm 0.086%	5.172% \pm 2.812%
-0.200	5.225%	5.567% \pm 0.088%	6.143% \pm 2.968%
-0.250	4.955%	5.304% \pm 0.074%	6.580% \pm 2.607%
-0.300	4.689%	5.091% \pm 0.081%	7.896% \pm 2.932%
-0.350	4.428%	4.859% \pm 0.065%	8.870% \pm 2.439%
-0.400	4.173%	4.650% \pm 0.042%	10.258% \pm 1.621%
-0.450	3.923%	4.454% \pm 0.039%	11.922% \pm 1.543%
-0.500	3.679%	4.234% \pm 0.035%	13.108% \pm 1.437%
-0.550	3.441%	4.052% \pm 0.051%	15.079% \pm 2.138%
-0.600	3.210%	3.871% \pm 0.050%	17.076% \pm 2.143%

Table 24 [Like TABLE 10 but with lognormal demand]The estimates of regular service violation probability (by analytic approximations and by simulation) in an integrated network for the regular firm as a function of correlation.

Correlation ρ	F (approximation)	F (simulation)	$\ F(\text{app})-F(\text{sim})\ /F(\text{sim})$
-0.050	0.520	0.834 ± 0.003	$37.650\% \pm 0.449\%$
-0.100	0.595	0.855 ± 0.002	$30.409\% \pm 0.326\%$
-0.150	0.667	0.876 ± 0.002	$23.858\% \pm 0.348\%$
-0.200	0.734	0.894 ± 0.003	$17.897\% \pm 0.551\%$
-0.250	0.795	0.912 ± 0.003	$12.829\% \pm 0.573\%$
-0.300	0.848	0.926 ± 0.003	$8.423\% \pm 0.593\%$
-0.350	0.889	0.939 ± 0.003	$5.325\% \pm 0.605\%$
-0.400	0.919	0.949 ± 0.002	$3.161\% \pm 0.408\%$
-0.450	0.941	0.958 ± 0.002	$1.775\% \pm 0.410\%$
-0.500	0.957	0.964 ± 0.001	$0.726\% \pm 0.206\%$
-0.550	0.968	0.970 ± 0.001	$0.206\% \pm 0.206\%$
-0.600	0.977	0.974 ± 0.001	$0.308\% \pm 0.103\%$

Table 25 [Like TABLE 11 but with lognormal demand]The estimates of the tightness factor F (by analytic approximations and by simulation) in an integrated network for the express firm as a function of correlation.

Correlation ρ	F (approximation)	F (simulation)	$\ F(\text{app})-F(\text{sim})\ /F(\text{sim})$
-0.050	0.673	0.665 ± 0.006	$1.203\% \pm 0.921\%$
-0.100	0.689	0.677 ± 0.005	$1.773\% \pm 0.757\%$
-0.150	0.704	0.687 ± 0.005	$2.475\% \pm 0.751\%$
-0.200	0.719	0.699 ± 0.005	$2.861\% \pm 0.741\%$
-0.250	0.733	0.714 ± 0.005	$2.661\% \pm 0.724\%$
-0.300	0.747	0.725 ± 0.005	$3.034\% \pm 0.716\%$
-0.350	0.762	0.738 ± 0.004	$3.252\% \pm 0.563\%$
-0.400	0.775	0.749 ± 0.003	$3.471\% \pm 0.416\%$
-0.450	0.789	0.759 ± 0.003	$3.953\% \pm 0.413\%$
-0.500	0.802	0.772 ± 0.002	$3.886\% \pm 0.270\%$
-0.550	0.815	0.781 ± 0.003	$4.353\% \pm 0.402\%$
-0.600	0.827	0.791 ± 0.003	$4.551\% \pm 0.398\%$

Table 26 [Like TABLE 12 but with lognormal demand]The estimates of the tightness factor F (by analytic approximations and by simulation) in an integrated network for the regular firm as a function of correlation.