# A NEW GENERATION OF VEHICLE ROUTING RESEARCH: ROBUST ALGORITHMS, ADDRESSING UNCERTAINTY

## DIMITRIS J. BERTSIMAS

*Massachusetts Institute of Technology, Cambridge, Massachusetts*

## DAVID SIMCHI-LEVI

*Northwestern University, Evanston, Illinois*

In recent years new insights and algorithms have been obtained for the classical, deterministic vehicle routing problem as well as for natural stochastic and dynamic variations of it. These new developments are based on theoretical analysis, combine probabilistic and combinatorial modeling, and lead to new algorithms that produce near-optimal solutions, and a deeper understanding of uncertainty issues in vehicle routing. In this paper, we survey these new developments with an emphasis on the insights gained and on the algorithms proposed.

In recent years, many service suppliers and distributors have recognized the importance of designing efficient distribution strategies to improve the level of customer's service and reduce freight transportation costs, which totaled more than 100 billion dollars in the United States in 1990. In a typical distribution system, vehicles (e.g., trucks or school buses) provide delivery, customer pick-up, or repair and maintenance services to customers that are geographically dispersed in a given area. In most applications (the distribution of soft drinks, beer, gasoline and pharmaceuticals, or the pick-up and delivery of students by school buses just to name a few), a common objective is to find a set of routes for the vehicles which satisfies a variety of constraints and so as to minimize the total fleet operating cost. The problem of minimizing total cost has traditionally been called the vehicle routing problem (VRP). In other applications (especially in repair and service contexts), it is important to minimize the total time the customers spend waiting to be served.

In the last decade, new insights and algorithms have been obtained for the classical deterministic vehicle routing problem as well as for natural stochastic and dynamic variations of it. These new developments are based on theoretical analysis, combine probabilistic and combinatorial modeling, and lead to new and effective algorithms and a deeper understanding of uncertainty issues in vehicle routing. In this paper, we survey these new developments with an emphasis on the insights gained and on the algorithms proposed.

The first set of results we survey (Sections 1–5) describes new near-optimal algorithms for VRPs. The complexity and wide applicability of the VRP have motivated researchers to develop heuristic algorithms (or simply heuristics) for its solution. Consequently, the problem has been analyzed extensively in contemporary journals; an excellent survey of the literature may be found in Fisher (1995). Traditionally, these heuristics are analyzed *empirically*, that is, the performance of a specific heuristic is evaluated on a set of *standard test problems*. As observed in Fisher a common limitation of this approach is the lack of robustness; a heuristic algorithm that works well on a set of standard test problems does not necessarily perform well on any particular application. The heuristic is then "patched up" to fix the troublesome cases, leading to an algorithm with growing complexity (see, Fisher). After considerable effort, a procedure is created that works well for the situation at hand (Fisher). Unfortunately, the resulting algorithm is usually extremely sensitive to changes in the data, and may perform poorly when transformed to other environments.

To overcome this difficulty, an in-depth analysis of some VRPs has been carried out that makes it possible to understand the underlying structure of these problems as a first step toward designing algorithms that can efficiently solve large-scale problems. This has led to the development of new algorithms that are more robust, i.e., algorithms that are independent of the specific environment and the variability in the data, because they are *designed* to handle a

general situation. In addition, this approach enables us to better understand models that integrate vehicle routing with other issues important to the firm, such as integrating inventory control and vehicle routing or design problems associated with distribution systems.

The second set of results we survey (Sections 6 and 7) describes new models and algorithms for VRPs in which uncertainty (in customer's demands, location or arrival time of a customer request for service) plays a major role. While the classical view of the VRP is static and deterministic, in many of the practical applications in which VRPs arise (e.g., distribution, inventory resupply, mobile repair), there are significant *stochastic* and *dynamic* components to the problem. Indeed, in many real-life logistics systems, demands arrive randomly in time, have a random size, and thus routing is a continuous process of collecting demands, forming tours, and dispatching vehicles.

To analyze problems of this type, new models have been proposed for VRPs under uncertainty. The analysis of these stochastic and dynamic VRPs provides structural insight into the effects on performance of traffic intensity, on-site service characteristics, the number, speed and capacity of vehicles employed, service region size, and the distribution of customer locations.

The discussion in this paper remains at an intuitive level; the interested reader is referred to Federgruen and Simchi-Levi (1995) and Powell, Jaillet and Odoni (1995) for proofs of some of the results presented below. The paper is organized as follows. In Sections 1 and 2 we review new algorithms for different versions of the capacitated VRP and analyze their performance from a worst and average perspective. In Sections 3 and 4 we extend some of the results to the VRP with capacity and time window constraints. Section 5 shows how the insights obtained from the previous analyses can be applied to general distribution systems. Section 6 presents a new approach based on a priori optimization for VRPs with stochastic but static customer demand. Section 7 presents a new approach based on combinatorial optimization and queueing theory for VRPs with stochastic and dynamic demand. The final section summarizes the most important points we are making.

## 1. VRP WITH EQUAL DEMANDS

Perhaps the simplest model for VRP that leads to important insights is defined as follows: A set of customers has to be served by a fleet of identical vehicles of fixed capacity $q$. The vehicles are initially located at a given depot. Associated with each customer is its demand which is the number of items that has to be delivered to that customer. The vehicle capacity specifies an upper bound on the number of items that can be delivered by a single vehicle. The objective is to find a set of routes for the vehicles of minimal total length. Each route begins at the depot, visits a subset of the customers, and returns to the depot without violating the capacity constraint.

If customers' demands are identical, they can be assumed, without loss of generality, to be equal to one. In that case, the capacity constraint states that the number of customers visited by a single vehicle cannot exceed $q$. This model is called the **equal demand** model.

The capacitated VRP with equal demands can hardly be considered a practical model. It has been, however, the subject of intensive analysis in the last five years. This is primarily due to two reasons: First, if the demands of the customers are not equal, but can be distributed among more than one vehicle, the problem can be reduced to the equal demand model. This is done by replacing each customer with demand $w$ by $w$ customers of unit demand, each one of them is located at the location of the original customer. Second, the insights that we get from the analysis of this model will be very useful in the analysis of models that integrate vehicle routing with other issues important to the firm.

Let $N$ denote the set of customers, $d_i$ the distance between node $i$ and the depot, $d_{max} \equiv \max_{i \in N} d_i$, is the distance to the furthest customer, and $d_{ij}$ the distance between customer $i$ and customer $j$. The distance matrix $\{d_{ij}\}$ is assumed to be symmetric and to satisfy the triangle inequality, i.e., $d_{ij} = d_{ji}$ for all $i, j$ and $d_{ij} \leq d_{ik} + d_{kj}$ for all $i, k, j$. We denote the optimal solution value of the capacitated VRP by $Z^*$ and that of any given heuristic **H** by $Z^H$.

In what follows, an $\alpha$-optimal traveling salesman tour plays an important role. An $\alpha$-optimal tour is a traveling salesman tour whose length is no more than $\alpha$ times the length of the optimal traveling salesman tour.

### 1.1. Worst-Case Analysis

A simple heuristic for the capacitated VRP, is the following tour partitioning heuristic suggested by Beasley (1983) and for which Altinkemer and Gavish (1990) provide an interesting worst-case analysis. In this heuristic, called the **optimal partitioning (OP)** heuristic, one constructs a traveling salesman tour through the customers and the depot. The tour is then *optimally* partitioned into segments, each containing at most $q$ customers, by formulating an appropriate shortest path problem.

This is done as follows: Given a traveling salesman tour through the customers and the depot, the points are numbered $x^{(0)}, x^{(1)}, \ldots, x^{(n)}$ in order of appearance on the tour, where $x^{(0)}$ is the depot. Let

$$C_{jk}$$

$$= \begin{cases} \text{the distance traveled by a vehicle that starts at} \\ \text{the depot, visits customers } x^{(j+1)}, x^{(j+2)}, \ldots, x^{(k)} \\ \text{in this order, and returns to the depot,} \\ +\infty, \end{cases}$$

$$\begin{aligned} &\text{if } k - j \leq q; \\ &\text{otherwise.} \end{aligned}$$

If we find the shortest path from $x^{(0)}$ to $x^{(n)}$ in the directed graph with distance cost $C_{jk}$, we will have chosen an optimal partition of the traveling salesman tour.

The performance of this heuristic clearly depends on the quality of the initial traveling salesman tour chosen in the first step of the algorithm. Hence, when the **OP** heuristic partitions an $\alpha$-optimal traveling salesman tour, it is denoted by **OP**($\alpha$). Altinkemer and Gavish (1990) proved the following result.

### Theorem 1

$$\frac{Z^{OP(\alpha)}}{Z^*} \leqslant 1 + \left(1 - \frac{1}{q}\right)\alpha .$$

For example, if Christofides' (1976b) algorithm ($\alpha = 1.5$) is used to obtain the initial traveling salesman tour in polynomial time, we have

$$\frac{Z^{OP(1.5)}}{Z^*} \leqslant \frac{5}{2} - \frac{3}{2q} .$$

The proof of the worst-case result for the **OP**($\alpha$) heuristic suggests that if we can improve the bound in Theorem 1 for $\alpha = 1$, then the bound can be improved for any $\alpha > 1$. However, the following theorem, proved by Li and Simchi-Levi (1990), says that this is impossible. That is, it shows that for $\alpha = 1$, $Z^{OP(1)}/Z^*$ tends to 2 when $q$ approaches infinity.

**Theorem 2.** *For any integer $q \geqslant 1$, there exists a problem instance with $Z^{OP(1)}/Z^*$ arbitrarily close to $2 - 2/(q + 1)$.*

We conclude that the worst-case performance of the optimal partitioning heuristic is quite disappointing; the cost of the solutions produced by **OP** can be quite far from the optimal cost. It is therefore natural to compare this performance with the average performance of **OP**. This is done in the next subsection.

### 1.2. Average-Case Analysis

For the purpose of characterizing the average performance of **OP**, we assume in the remainder of this section that the customers are points in the plane and that the distance between any pair of customers is given by the Euclidean distance. The next result, obtained by Haimovich and Rinnooy Kan (1985) fully characterizes the average performance of the **OP**($\alpha$) heuristic by comparing it to the best possible performance, i.e., to $Z^*$.

**Theorem 3.** *Let $x_k$, $k = 1, 2, \ldots, n$ be a sequence of independent random variables having a distribution $\mu$ with compact support in $\Re^2$. Let $d(y)$ be the Euclidean distance between the depot and $y \in \Re^2$ and let*

$$E(d) = \int_{\Re^2} d(y) \, d\mu(y) .$$

*Then, for any fixed $\alpha$, we have with probability one,*

$$\lim_{n \to \infty} \frac{Z^*}{n} = \lim_{n \to \infty} \frac{Z^{OP(\alpha)}}{n} = \frac{2}{q} E(d) .$$

The proof of Theorem 3 is based on constructing lower and upper bounds that converge asymptotically to the

same value. The fact that $Z^* \sim 2nE(d)/q$ is explained as follows Haimovich and Rinnooy Kan: Any solution for the capacitated VRP has two cost components; the first component is proportional to the total "radial" cost between the depot and the customers. The second component is proportional to the "circular" cost; the cost of traveling between customers. This cost is related to the cost of the optimal traveling salesman tour. It is well known (Beardwood, Halton and Hammersley 1959) that, for large $n$ the cost of the optimal traveling salesman tour grows like $\sqrt{n}$, while the total radial cost between the depot and the customers grows like $n$ because the number of vehicles used in any solution is at least $\lceil n/q \rceil$. Therefore, it is intuitive that when the number of customers is large enough the first cost component will dominate the optimal solution value.

## 2. VRP WITH UNEQUAL DEMANDS

In this section, we consider the more practical and usually more complicated, capacitated VRP with unequal demands. In this version, each customer $i$ has a demand $w_i$ and the capacity constraint states that the total amount delivered by a single vehicle cannot exceed the vehicle capacity $q$. As mentioned, if the demand of a customer can be split over several vehicles, the problem is reduced to the equal-weight case by treating a customer with demand $w$ as $w$ customers with unit size demand all at the same location. Therefore, the results of the previous section apply. Consequently, we consider the case where the demand of a customer may not be divided among vehicles. This constraint introduces bin-packing features into the routing problem and, consequently, requires different solution techniques. We refer to this version as the capacitated VRP with unsplit demands.

### 2.1. Worst-Case Analysis

In the worst-case analysis presented here, we assume that the numbers $w_1, w_2, \ldots, w_n$ and $q$ are rationals, and hence, without loss of generality, $q$ and $w_i$ are assumed to be integers.

The tour partitioning heuristic suggested for the equal demand case (i.e., **OP**($\alpha$)) can be trivially generalized for the unequal demand case. This is done by replacing the condition $k - j \leqslant q$ in the definition of the quantities $C_{jk}$ (see subsection 1.1) by $\sum_{i=j+1}^{k} w_i \leqslant q$. In that case, the heuristic is denoted **UOP** and if the initial tour is an $\alpha$-optimal traveling salesman tour the algorithm is called **UOP**($\alpha$). Altinkemer and Gavish (1987) prove the following result.

### Theorem 4

$$\frac{Z^{UOP(\alpha)}}{Z^*} \leqslant 2 + \left(1 - \frac{1}{q}\right)\alpha .$$

Observe the increase in the worst-case bound of the **UOP**($\alpha$) relative to that of **OP**($\alpha$). This increase is due to

the fact that when a tour partitioning heuristic assigns customers to vehicles, it can generate, in the worst-case, twice as many vehicles as in the optimal solution. By contrast, in the equal demand case, tour partitioning heuristics can always find a solution that uses the minimum number of vehicles. In view of this observation, it is not surprising that the bound of Theorem 4 cannot be reduced. Indeed, Li and Simchi-Levi show that when $\alpha = 1$, this bound is asymptotically tight as $q$ approaches infinity.

**Theorem 5.** *For any integer $q \geq 1$, there exists a problem instance with $Z^{UOP(1)}/Z^*$ arbitrarily close to $3 - 6/(q + 2)$.*

## 2.2. Average-Case Analysis

Results on the average performance of algorithms for the capacitated VRP with unsplit demands are closely related to results obtained for the bin-packing problem, a problem that has been analyzed extensively in the literature. An instance of the bin-packing problem is composed of the bin capacity (equal to 1) and a set of items each with a prespecified size no larger than 1. The problem is to find the smallest number of bins in which these items can be packed, subject to the constraint that the total size of items assigned to a single bin does not exceed 1.

In the probabilistic analysis of the capacitated VRP with unsplit demands we assume, without loss of generality, and in accordance to the convention in Coffman, Lueker and Rinnooy Kan (1988), that the vehicles' capacity $q$ equals 1, and the demand of each customer is no more than 1. Thus, vehicles and demands in the capacitated VRP correspond to bins and item sizes (respectively) in the bin-packing problem. Hence, for every routing instance there is a unique corresponding bin-packing instance.

### 2.2.1. Optimal Solution Value

Since the average performance of any heuristic has to be evaluated relative to the best possible performance, we start with the optimal solution value. Assume that the demands $w_1, w_2, \ldots, w_n$ are independent and identically distributed with distribution $\Phi$ defined on [0, 1]. In this section, we find the asymptotic optimal solution value for *any* distribution of the demands $\Phi$. This is done by showing that an asymptotically optimal algorithm for the bin-packing problem, with item sizes distributed like $\Phi$, can be used to solve the capacitated VRP with unsplit demands.

Given the demands $w_1, w_2, \ldots, w_n$, let $b^*$ be the number of bins used in the optimal solution for the corresponding bin-packing problem. As demonstrated in Rhee and Talagrand (1987), there exists a constant $\gamma > 0$ such that

$$\lim_{n \to \infty} \frac{b^*}{n} = \gamma \quad (a.s.).$$

The following theorem is proven in Bramel et al. (1992).

**Theorem 6.** *Let $x_k$, $k = 1, 2, \ldots, n$ be a sequence of independent random variables having a distribution $\mu$ with*

*compact support in $A \subset \mathfrak{R}^2$. Let $d(y)$ be the Euclidean distance between the depot and point $y \in \mathfrak{R}^2$ and let*

$$E(d) = \int_{\mathfrak{R}^2} d(y) \, d\mu(y).$$

*Let the demands $w_k$, $k = 1, 2, \ldots, n$ be a sequence of i.i.d. random variables having a distribution $\Phi$ with support on [0, 1] and assume that the demands and the locations of the customers are independent of each other. Then, almost surely,*

$$\lim_{n \to \infty} \frac{Z^*}{n} = 2\gamma E(d).$$

Intuitively, this implies that for large values of $n$, $n\gamma$ vehicles are needed to serve the customers and, on average, each one travels a distance of $2E(d)$. This is explained by the fact that the distribution $\Phi$ is independent of $n$ and thus with high probability the number of customers per vehicle is constant. Consequently, the total distance traveled is dominated by the radial distance to and from the depot.

Theorem 6 is proved by constructing lower and upper bounds that asymptotically converge to the same value; the asymptotic optimal solution value. The upper bound is based on the following feasible solution. It uses a special region partitioning scheme, in which the area where the customers are located is partitioned by means of a grid into many squares, also referred to as subregions. Customers within the same subregion are assigned to vehicles to minimize the number of vehicles used within each subregion. Every vehicle serves customers from only one square using the following routing strategy: The vehicle travels to the subregion where its customers are located, visits the customers in any order, and then returns to the depot. By choosing the grid such that the size of each square tends to zero at a rate slower than the rate in which $n$ increases, one can prove, under the assumptions of Theorem 6, that

$$\lim_{n \to \infty} \frac{Z^*}{n} \leq 2\gamma E(d) \quad (a.s.).$$

This upper bound is combined with a lower bound on the optimal solution value to prove Theorem 6.

### 2.2.2. Analysis of Classical Heuristics

The complexity and the economic importance of the capacitated VRP with unsplit demands have motivated the development of many heuristics for its solution; see, e.g., Christofides (1985) or Fisher (1995). Of special interest is the class of heuristics called (by Christofides, 1976a and 1985) two-phase methods. These heuristics are of two types: cluster first-route second, or route first-cluster second. In the first category, customers are clustered into groups and assigned to vehicles (phase I) and then efficient routes are designed for each cluster (phase II). In the second category, one constructs a traveling salesman tour through all the customers (phase I) and then partitions the tour into segments (phase II). One vehicle is assigned to

each segment and visits the customers according to their appearance on the traveling salesman tour.

Bienstock, Bramel and Simchi-Levi (1993) analyze the average performance of heuristics that belong to the latter class. To present their result, we need a precise definition of the class of the route first-cluster second method. Define this class as all those heuristics that first order the customers according to their locations and then partition this ordering to produce feasible clusters. These clusters consist of sets of customers that are consecutive in the initial order. Customers are then routed within their cluster depending on the specific heuristic.

Observe that this definition of the class of the route first-cluster second heuristics is more general than classical definitions. It is also clear that the **UOP**($\alpha$) heuristic described in subsection 2.1 belongs to this class of heuristics. The sweep algorithm suggested by Gillett and Miller (1974) can also be viewed as a route first-cluster second type of heuristic. In this algorithm, an arbitrary customer is selected as a starting customer. The other customers are ordered according to the angle between them, the depot and the starting customer. Customers are then assigned to vehicles following this initial ordering and efficient routes are designed for each vehicle.

Bienstock, Bramel and Simchi-Levi show that the performance of any heuristic in this class is strongly related to the performance of a nonoptimal bin-packing heuristic called next fit (**NF**). Thus, heuristics in this class can never be asymptotically optimal for the capacitated VRP with unsplit demands.

The next-fit bin-packing heuristic can be described in the following manner. Given a list of $n$ items where the size of item $i$ is $w_i$, start with item 1 and place it in bin 1. Suppose that we are packing item $j$; let bin $i$ be the highest indexed nonempty bin. If item $j$ fits in bin $i$, then place it there, else place it in a new bin indexed $i + 1$. Thus, the **NF** heuristic assigns items to bins according to the order they appear without using any knowledge of subsequent items in the list.

In their seminal work on the use of martingale inequalities for NP-complete problems, Rhee and Talagrand show that for any distribution of the item sizes, there exists a constant $\gamma^{NF} > 0$ such that $\lim_{n \to \infty} b^{NF}/n = \gamma^{NF}$ almost surely, where $b^{NF}$ is the number of bins produced by the **NF** packing and $\gamma^{NF}$ depends only on the distribution of the item sizes. This constant is used in the following theorem proved in Bienstock, Bramel and Simchi-Levi.

**Theorem 7.** *Let* **H** *be a generic route first-cluster second heuristic, that is,* **H** *is a heuristic that starts by ordering the customers in some manner depending only on their relative locations and not on their demands. Then, under the assumptions of Theorem 6 we have*

$$\lim_{n \to \infty} \frac{1}{n} Z^H \geq 2\gamma^{NF}E(d) \quad (a.s.).$$

We therefore conclude that the empirically well-studied route first-cluster second methods can never be asymptotically optimal for the capacitated VRP with unsplit demands except in some trivial cases, i.e., when $\gamma = \gamma^{NF}$. The next theorem completely characterizes the average performance of the **UOP**($\alpha$) heuristic by showing that it is the best possible heuristic in the route first-cluster second class.

**Theorem 8.** *Under the assumptions of Theorem 6, the* **UOP**($\alpha$) *heuristic is the best possible heuristic in the class of route first-cluster second, that is, for any fixed $\alpha \geq 1$ we have*

$$\lim_{n \to \infty} \frac{1}{n} Z^{UOP(\alpha)} = 2\gamma^{NF}E(d) \quad (a.s.).$$

In view of Theorems 6, 7, and 8 it is interesting to compare $\gamma^{NF}$ to $\gamma$ because the asymptotic error for any heuristic **H** in the class of route first-cluster second satisfies

$$\lim_{n \to \infty} Z^H/Z^* \geq \lim_{n \to \infty} Z^{UOP(\alpha)}/Z^* = \gamma^{NF}/\gamma.$$

This ratio was characterized by Karmarkar (1982) for the case when the item sizes are uniformly distributed on an interval $[0, a]$ for $0 < a \leq 1$. For instance, for $a$ satisfying $1/2 < a \leq 1$, we have

$$\gamma^{NF}/\gamma = \frac{2}{a} \left\{ \frac{1}{12a^3} (15a^3 - 9a^2 + 3a - 1) \right.$$
$$\left. + \sqrt{2} \left( \frac{1-a}{2a} \right) \tanh \left( \frac{1-a}{\sqrt{2}a} \right) \right\},$$

so that when the item sizes are uniform $[0, 1]$ the above ratio is 4/3 which implies that **UOP**($\alpha$) converge to a value which is 33.3% more than the optimal cost, which is a very disappointing performance for the best heuristic currently available in terms of worst-case behavior.

The sweep (**S**) algorithm also possesses the properties needed to apply the lower bound of Theorem 7 because this heuristic starts by choosing the order of the customers (according to the angle between them the depot and the starting customer) and then assigns them to vehicles following this order. In fact, Bienstock, Bramel and Simchi-Levi prove a much stronger result. To present the result, let $P_k$ be the long-term fraction of bins with $k$ items generated by the **NF** heuristic. As we shall see in Theorem 9, the expected distance traveled in the solution generated by the **sweep** depends on these values.

Without loss of generality, assume that the distribution of customer locations is indexed by polar coordinates (i.e., $\mu(r, \theta)$, $r \in [0, d_{max}]$, $\theta \in [0, 2\pi]$), where the starting angle ($\theta = 0$) is chosen arbitrarily. Let

$$\mu(\theta) = \int_0^{d_{max}} \mu(r, \theta)r \, dr.$$

Let $F(t|\theta)$ be the conditional probability of a point being at a distance no more than $t$ from the depot, given that the point is at angle $\theta$. We have the following theorem.

**Theorem 9.** *Under the assumptions of Theorem 6 we have*

$$\lim_{n \to \infty} \frac{1}{n} E(Z^S) = 2\gamma^{NF} \sum_{k=1}^{\infty} P_k \int_0^{2\pi} \int_0^{d_{max}}$$

$$\cdot (1 - F^k(t|\theta)) \, dt \, \mu(\theta) \, d\theta .$$

Intuitively, this is explained as follows: When the number of customers, $n$, is large enough, the sweep algorithm produces $n\gamma^{NF}$ vehicles, of which $nP_k\gamma^{NF}$ have exactly $k$ customers. Now consider all the customers having angle $\theta$ with $\theta_0 \le \theta \le \theta_0 + \epsilon$ for some $\theta_0 \ge 0$. Since the number of customers, $n$, is very large, the number of customers in this sector is also very large. Thus, a vehicle that serves a customer from that sector will, with high probability, serve no customers from other sectors. Focusing now on a single vehicle that serves customers set $S$ from this sector and choosing $\epsilon$ small enough, the total distance traveled by that vehicle can be made arbitrary close to twice the distance from the depot to the furthest customer in $S$. Consequently, if $S$ includes $k$ customers, the expected distance to the furthest customer is roughly

$$\int_0^{d_{max}} (1 - F^k(t|\theta_0)) \, dt .$$

We conclude that the main drawback of route first-cluster second methods is that a higher priority is given to the routing part of the problem than to the bin-packing part of it. These results show that in this case *the heuristic can never be asymptotically optimal*, except in some trivial cases (e.g., item sizes are uniformly distributed on $[a, b]$ for $a > 0.5$). This demonstrates that a heuristic has a potential of being asymptotically optimal only if the bin-packing component of the capacitated VRP is considered at the same time as the routing component. This is exactly the type of algorithm described in the next section.

### 2.2.3. A New Class of Heuristics

Bramel and Simchi-Levi (1995) used the insight obtained from the analysis of the asymptotic optimal solution value (see Theorem 6 and the discussion that follows it) to develop a new and highly effective class of heuristics for the capacitated VRP with unsplit demands. Specifically, this class of heuristics was motivated by the following observations.

A by-product of the proof of Theorem 6 is that the feasible solution used to find an upper bound on $Z^*$ is asymptotically optimal. In this upper bound the length of every tour that visits a set of customers $S$ consists of two parts. The first is the length of the *tour* that starts at the depot visits the subregion (where its customers are located) and then goes back to the depot. The second is the additional distance obtained by *inserting* all the customers in $S$ into this tour. It is clear, therefore, that if we can construct a heuristic that assigns customers to vehicles so as to minimize the sum of the length of all simple tours plus the total insertion costs of customers to each simple

tour, then the heuristic will have the right structure to be asymptotically optimal.

To construct such a heuristic we formulate the routing problem as a standard combinatorial problem commonly called (see, e.g., Pirkul 1987) the capacitated concentrator location problem. This problem can be described as follows: Given $m$ possible sites for concentrators of fixed capacity $Q$, we would like to locate concentrators at a subset of these $m$ sites and connect $n$ terminals, where terminal $i$ uses $w_i$ units of a concentrator's capacity, in such a way that each terminal is connected to *exactly one* concentrator, the concentrator capacity is not exceeded and the total cost is minimized. A site-dependent cost is incurred for locating each concentrator; that is, if a concentrator is located at site $j$, the *setup* cost is $v_j$ for $j = 1, 2, \ldots, m$. The cost of connecting terminal $i$ to concentrator $j$ is $c_{ij}$ (the *connection* cost) for $i = 1, 2, \ldots, n$ and $j = 1, 2, \ldots, m$.

The capacitated concentrator location problem can be formulated as the following integer linear program. Let

$$y_j = \begin{cases} 1, & \text{if a concentrator is located at site } j, \\ 0, & \text{otherwise,} \end{cases}$$

and let

$$x_{ij} = \begin{cases} 1, & \text{if terminal } i \text{ is connected to a concentrator at site } j, \\ 0, & \text{otherwise.} \end{cases}$$

### Problem P

Minimize $\sum_{i=1}^{n} \sum_{j=1}^{m} c_{ij} x_{ij} + \sum_{j=1}^{m} v_j y_j$

subject to

$$\sum_{j=1}^{m} x_{ij} = 1 \qquad \text{for all } i , \tag{1}$$

$$\sum_{i=1}^{n} w_i x_{ij} \le Q \qquad \text{for all } j , \tag{2}$$

$$x_{ij} \le y_j \qquad \text{for all } i, j , \tag{3}$$

$$x_{ij} \in \{0, 1\} \qquad \text{for all } i, j , \tag{4}$$

$$y_j \in \{0, 1\} \qquad \text{for all } j . \tag{5}$$

Constraints 1 ensure that each terminal is connected to exactly one concentrator, and constraints 2 ensure that the concentrator's capacity constraint is not violated. Constraints 3 guarantee that if a terminal is connected to site $j$, then a concentrator is located at that site. Constraints 4 and 5 ensure the integrality of the variables.

In formulating the capacitated VRP with unsplit demands as the capacitated concentrator location problem, every customer $x_j$ in the VRP is a possible concentrator site in the location problem. The length of the simple tour in the VRP that starts at the depot visits customer $x_j$ and then goes back to the depot is the setup cost in the location problem (i.e., $v_j = 2d_j$). Finally, the cost of inserting a customer into a simple tour in the VRP is the connection cost in the location problem (i.e., $c_{ij} = d_i + d_{ij} - d_j$).

Hence, finding a solution for the capacitated VRP is obtained by solving a facility location problem with the data as described before. The solution obtained for the capacitated concentrator location problem is transformed (in an obvious way) to a solution for the capacitated VRP.

The capacitated concentrator location problem, though NP-hard, can efficiently be solved by the familiar Lagrangian relaxation technique, as described in Pirkul (1987).

We can now describe an algorithm for the capacitated VRP, called the location-based heuristic (**LBH**), which is based on the insight obtained from the analysis of the asymptotic optimal solution (the details of the algorithm are presented in Bramel and Simchi-Levi, 1995).

### Algorithm LBH

*STEP 1.* Formulate the capacitated VRP as a capacitated concentrator location problem.

*STEP 2.* Solve the capacitated concentrator location problem.

*STEP 3.* Transform the solution obtained in Step 2 into a solution for the VRP.

This algorithm was tested empirically and analyzed analytically. For instance, Bramel and Simchi-Levi (1995) prove the following theorem.

**Theorem 10.** *Under the assumptions of Theorem 6, there is a version of the LBH which is asymptotically optimal, i.e.,*

$$\lim_{n \to \infty} \frac{1}{n} Z^{LBH} = 2\gamma E(d) \quad (a.s.).$$

Finally, we observe that the generalized assignment heuristic due to Fisher and Jaikumar (1981) can be viewed as a special case of the **LBH** in which the seed customers are first selected by a dispatcher. In the second step, customers are assigned to the seeds in an efficient way by solving a generalized assignment problem. The advantage of the **LBH** is that the selection of the seeds and the assignment of customers to seeds are done simultaneously, and not sequentially as in the generalized assignment heuristic. A by-product of the analysis, therefore, is that when the generalized assignment heuristic is carefully implemented (i.e., "good" seeds are selected), it is asymptotically optimal as well.

### 3. THE VRP WITH TIME WINDOW CONSTRAINTS

In many distribution systems each customer specifies, in addition to the load that has to be delivered to it, a period of time, called a *time window*, in which this delivery must occur. The objective is to find a set of routes for the vehicles, where each route begins and ends at the depot, serves a subset of the customers without violating the vehicle capacity and time window constraints, while minimizing the total length of the routes. We call this model the VRP with time windows.

Due to the wide applicability and the economic importance of the problem, variants of it have been extensively studied in the vehicle routing literature; for a review, see Solomon and Desrosiers (1988). Most of the work on the problem has focused on empirical analysis while very few papers have studied the problem from an analytical point of view in an attempt to characterize the theoretical behavior of heuristics and to use the insight obtained to construct effective algorithms for it. An exception is the recent work of Federgruen and van Ryzin (1992), Bramel and Simchi-Levi (1993) and Bramel, Li and Simchi-Levi (1994). We will describe the results of the second paper, which lends itself to a computationally attractive algorithm.

To present the results, let the quadruplet $(w_k, e_k, s_k, l_k)$ be the *parameters* of the $k$th customer which represents, respectively, the load that must be delivered, the earliest starting time for service, the time required to complete the service, called the *service time*, and the latest time service can end.

Surprisingly, the optimal solution of the VRP with time windows is directly related to the optimal solution of a *machine scheduling problem* defined as follows. Associated with each customer $k$ is a *job* whose parameters are the parameters of customer $k$. That is, the parameters of job $k$ are $(w_k, e_k, s_k, l_k)$, where $w_k$ is referred to as the *load* of job $k$, $e_k$ is referred to as the *release time* of job $k$, $s_k$ represents the processing time of job $k$, and $l_k$ is referred to as the *due date* of job $k$.

Consider the following machine scheduling problem defined by the parameters of the customers $N$ and an infinite sequence of *parallel* machines. Job (customer) $k$ becomes available for processing at time $e_k$ and must be finished processing by time $l_k$. The objective in this scheduling problem is to assign each job to a machine such that each machine has at most one job being processed on it at a given time; the processing time of each job starts no earlier than its release time and ends no later than its due date; and the total load of all jobs assigned to a machine is no more than 1, and the number of machines used is minimized.

Let $M^*$ be the minimum number of machines needed to schedule the set $N$ of jobs. The theory of subadditive processes (see Kingman 1976) implies that if $M_n^*$ is the minimum number of machines needed to schedule a set of $n$ jobs whose parameters are drawn independently from a distribution $\Phi$, then there exists a constant $\gamma > 0$ (depending only on $\Phi$) such that $\lim_{n \to \infty} M_n^*/n = \gamma$ (a.s.).

Bramel and Simchi-Levi (1993) show that asymptotically the VRP with time windows is no more difficult to solve than the corresponding scheduling problem. The main result is as follows.

**Theorem 11.** *Let $x_1, x_2, \dots, x_n$ and $E(d)$ be defined as in Theorem 6. Let the customer parameters $\{(w_k, e_k, s_k, l_k): k \in N\}$ be drawn independently from a distribution $\Phi$ with a*

*continuous density. Let $M_n^*$ be the minimum number of machines needed to feasibly schedule the $n$ jobs corresponding to these parameters, and let $\lim_{n \to \infty} M_n^*/n = \gamma$ (a.s.). Then*

$$\lim_{n \to \infty} \frac{1}{n} Z_n^* = 2\gamma E(d) \quad (a.s.) \,.$$

An important by-product of the analysis is the development of a new and highly efficient algorithm for the VRP with time windows. Computational evidence shows that the algorithm works very well on a set of standard test problems; see Bramel and Simchi-Levi (1993).

## 4. A COLUMN GENERATION TECHNIQUE

A classical method for solving the VRP with capacity and time window constraints, suggested by Balinski and Quandt (1964), is based on formulating the VRP as a set partitioning problem. The idea is as follows. Let the index set of all feasible routes be $\{1, 2, \ldots, R\}$ and let $c_r$ be the length of route $r$. Define

$$\alpha_{ir} = \begin{cases} 1, & \text{if customer } i \text{ is served in route } r, \\ 0, & \text{otherwise,} \end{cases}$$

for each $i = 1, 2, \ldots, n$ and $r = 1, 2, \ldots, R$. Finally, for $r = 1, 2, \ldots, R$, let

$$y_r = \begin{cases} 1, & \text{if route } r \text{ is in the optimal solution} \\ 0, & \text{otherwise.} \end{cases}$$

In the *set partitioning* formulation of the VRP, the objective is to select a minimum cost set of feasible routes such that each customer is included in some route. It is the following integer program.

## Problem S

Minimize $\sum_{r=1}^{R} y_r c_r$

subject to

$$\sum_{r=1}^{R} y_r \alpha_{ir} \geq 1, \quad \text{for all } i = 1, 2, \ldots, n \tag{6}$$

$$y_r \in \{0, 1\}, \quad \text{for all } r = 1, 2, \ldots, R \,.$$

Observe that we have written constraints 6 as inequality constraints. This is possible because the distances satisfy the triangle inequality and therefore each customer will be visited *exactly* once in the optimal solution.

This formulation was first used successfully by Cullen, Jarvis and Ratliff (1981) to design heuristic methods for the VRP. More recently, Desrochers, Desrosiers and Solomon (1992) have used it in conjunction with other methods to generate optimal or near-optimal solutions to the VRP.

The set of all feasible routes is extremely large and one cannot expect to generate it completely. Even if this set is given, it is not clear how to solve the set partitioning problem because it is a large-scale integer program. To overcome the first difficulty, Desrochers, Desrosiers and Solomon use the celebrated column generation technique,

which makes it possible to solve the linear programming relaxation of problem S without having to enumerate all the routes. This is done by enumerating a portion of all possible routes, and solving the resulting linear programming relaxation with this partial route set. The solution to the linear program is then used to determine if there are any routes not included which can reduce the solution value. This is the *column generation* step. Using the values of the optimal dual variables (with respect to the partial route set), we generate a new route and resolve the linear programming relaxation of the set partitioning problem. This is continued until one can show that an optimal solution to the linear program is found; one that is optimal for the complete route set. Finally, to get an integer solution to the set partitioning problem, the linear program is combined in a branch-and-bound routine.

It is well known that a branch-and-bound strategy works well only if the lower bound used in the bounding step is very tight. Fortunately, many researchers have reported that the linear programming relaxation of the set partitioning problem provides a solution close to the optimal integer solution; see, e.g., Desrochers, Desrosiers and Solomon. That is, the solution to the linear programming relaxation of S provides a strong lower bound on the solution to the VRP. Recently Bramel and Simchi-Levi (1993b) demonstrate why this is true in general. They prove the following theorem.

**Theorem 12.** *Let the customer locations be independently and identically distributed according to a distribution $\mu$ with compact support in $\mathbb{R}^2$. Let the customer parameters be independently and identically distributed like $\Phi$. Let $Z^{LP}$ be the value of the optimal fractional solution to S, and let $Z^*$ be the value of the optimal integer solution to S; that is, the value of the optimal solution to the VRP. Then,*

$$\lim_{n \to \infty} \frac{1}{n} Z^{LP} = \lim_{n \to \infty} \frac{1}{n} Z^* \quad (a.s.) \,.$$

Observe that an alternative formulation of problem S is obtained when constraints 6 are replaced by equality constraints, i.e.,

$$\sum_{r=1}^{R} y_r \alpha_{ir} = 1, \quad \text{for all } i = 1, 2, \ldots, n \,.$$

We call this problem SE. Since any feasible solution to the linear programming relaxation of problem SE is feasible for the linear programming relaxation of S, the cost of the optimal solution to the linear programming relaxation of SE is no smaller than the cost of the optimal solution to the linear programming relaxation of problem S. Consequently, Theorem 12 also holds when $Z^{LP}$ represents the optimal solution value of the linear programming relaxation of problem SE.

The theorem thus implies that the optimal solution value of the linear programming relaxation of problem S (and problem SE) tends to the optimal solution of the vehicle routing problem as the number of customers tends

to infinity. This is important since, as shown by Bramel and Simchi-Levi (1993), other classical formulations of the VRP can lead to diverging linear and integer solution values.

# 5. APPLICATIONS TO DISTRIBUTION SYSTEMS

The results obtained from the analysis of the capacitated VRP have been used to analyze, develop, and implement efficient algorithms for models that integrate vehicle routing with other issues important to a firm. This includes models that integrate inventory considerations with transportation costs as well as design issues associated with distribution systems. For a discussion on applications to inventory-routing models we refer the reader to Anily and Federgruen (1990, 1993), Gallego and Simchi-Levi (1990), Federgruen and Simchi-Levi (1992), Chan, Federgruen and Simchi-Levi (1993), and Bramel and Simchi-Levi (1995). Next we describe the use of the results in the context of system design.

## 5.1. Applications to Systems Design

The results described in subsection 1.2 characterize the average behavior of the $OP(\alpha)$ as well as the asymptotic optimal solution value of the single-depot capacitated VRP with equal demands. These results enable us (see Simchi-Levi 1992), to develop analytical models to assist the design and control of distribution systems.

As an example, consider a company that delivers consumer goods to a number of stores located in an area of size $A$. The company has decided to open a number of warehouses in the region and has carried out a market survey to estimate the number of potential customers, denoted by $n$, and the probability distribution of their demands. At this preliminary stage of the analysis the company assumes that all potential customers have the same probability distribution with $\bar{w}$ being the expected customer's demand.

Based on the information available from the survey, the company wants to determine the number and locations of warehouses; how to allocate customers to depots and what should be the routing strategies to minimize total system cost. This cost includes a cost associated with the average distance traveled by all vehicles plus a fixed setup cost, denoted by $c$, for each established depot.

The insight obtained from the analysis of the capacitated VRP can be used to propose a three-stage hierarchical approach in which decisions about the number of centers and their locations (first stage), customers allocations (second stage), and routing strategies (third stage) are combined to reduce total system cost.

As we have seen, the total radial cost between the depot and the customers dominates the objective function. This cost is related to the cost of the $K$-median problem. Thus, the $K$-median problem provides an insight to our model. For instance, when the demand of a customer can be split over several vehicles, the total distance traveled in a distribution system with $K$ centers optimally located in the area is asymptotically

$$\frac{2\bar{w}n}{q} \beta \sqrt{A/K},$$

where $\beta = 0.377196\ldots$ and $q$ is the vehicle capacity. This is true as long as $1/K_n = o(n)$ and $K_n = O(n/\log n)$. Furthermore, this asymptotic value is achieved by placing centers in a regular hexagonal pattern throughout the area and each service center serves all the customers inside its hexagon.

It follows that, for large enough $n$, the best number of warehouses minimizes the following function (recall that $c$ is a fixed setup cost for establishing a depot)

$$TC(K) = \frac{2\bar{w}n}{q} \beta \sqrt{A/K} + cK.$$

Let $\delta = (\beta (\bar{w}n/qc) \sqrt{A})^{2/3}$, then the best number of stations is $\lceil \delta \rceil$ or $\lfloor \delta \rfloor$, whichever yields the best $TC(K)$. Furthermore, the analysis also shows where to locate the centers and how to allocate customers to each center, thus providing answers to the strategic and tactical problems.

What should be the routing strategy used on a daily basis? Note that in this model, we assume that every working day the servers have exact information on the customers that need service and their actual demands. Hence, every working day, each center faces an instance of the single depot capacitated VRP, for which efficient heuristics exist. For example, by using the $OP(\alpha)$ heuristic on a daily basis the minimal total cost ($TC(K)$) can actually be achieved when the number of customers $n$ is large enough.

We can now summarize the hierarchical design: Choose the number of stations as the one that minimizes $TC(K)$, locate the facilities at the center of hexagonal patterns (strategic decision), each customer will be served by its closest service station (tactical decision), and use the $OP(\alpha)$ heuristic on a daily basis to find efficient routes for servers (operational control).

The main difficulty in using a hierarchical approach to design is in estimating the difference between the total system cost associated with the hierarchical design and the minimal total cost (obtained by integrating the three levels of decisions). In theory, the approach provides a design in which the total system cost approaches the minimal total cost, as the number of customers tends to infinity. In practice, however, the rate of convergence of the hierarchical design cost to the optimal cost may be quite slow. To estimate this rate of convergence Simchi-Levi (1992) performs a series of numerical experiments for problems of moderate size. For instance, for three distribution centers, located according to the hierarchical design, the relative error between the hierarchical design cost and a lower bound on the optimal solution value decreases from 35% to 16% as the number of customers increases from 100 to 500.

# 6. STOCHASTIC AND STATIC VEHICLE ROUTING

In this section, we consider a natural variation of the classical VRP in which demand at each location is unknown at the time when the tour is designed, but is assumed to

follow a known probability distribution. This situation arises in practice whenever a company (e.g., UPS), on any given day, is faced with the problem of deliveries/collections to/from a set of customers, each of which has a random demand.

Vehicle routing problems with random demands have received limited attention in the literature. Stewart and Golden (1983), Dror and Trudeau (1986), Dror, Laporte and Trudeau (1989), and Laporte and Louveaux (1990) use stochastic programming techniques to solve optimally small sized problems. Compared with this technique, the approach discussed in this section is completely different.

An obvious strategy to the vehicle routing with random demands is to redesign the routes when the demand becomes known. There are, however, several difficulties with such an approach: 1) computing resources might not be available, 2) even if resources are available it might be very time consuming to redesign the routes, 3) redesigning the routes might create confusion to drivers, and finally, 4) regularity and personalization of service by having the same vehicle and driver visit a particular customer every day is not guaranteed if one redesigns the routes. To overcome these difficulties, the strategy of designing *an a priori route* among all potential customers has been proposed in Jaillet (1988) for the traveling salesman problem (see also Jaillet and Odoni 1988) and Bertsimas (1992) for the VRP as an alternative to the strategy of redesigning the routes. The idea is to find an a priori solution to the combinatorial problems and update this solution when the demand is realized; see Bertsimas, Jaillet and Odoni (1990).

This strategy is described as follows: Determine a fixed a priori sequence among all potential customers. Depending on when information about customer's demand becomes available update routes as follows.

### Fixed a Priori Strategy A

Under this strategy the vehicle visits all the customers in the same fixed order as under the a priori sequence, but serves only customers requiring service that day. The total expected distance traveled corresponds to the fixed length of the a priori sequence plus the expected value of the additional distance that must be covered whenever the demand on the sequence exceeds vehicle capacity.

### Adaptive a Priori Strategy B

This strategy is defined similarly to the fixed strategy A with the sole difference that routes are adaptively updated by skipping customers with no demand on a particular instance.

To illustrate the difference between the two strategies consider the following example. If the a priori sequence is (0, 1, 2, 3, 4, 5, 6, 0), the depot is node 0, the vehicle has capacity 3, and the demand of the customers is $D_1 = 0$, $D_2 = 2, D_3 = 1, D_4 = 0, D_5 = 2, D_6 = 0$, then under strategy A the resulting routes are (0, 1, 2, 3, 0), (0, 4, 5, 6, 0), while under strategy B the resulting routes are (0, 2, 3, 0), (0, 5,

0). Note that at node 3 the capacity is reached and the vehicle is forced to return to the depot.

There is an important difference in the philosophy of the two updating strategies. Strategy A models situations in which the demand (if any) of a particular customer becomes known only when the customer is visited. The vehicle is then forced to return to the depot when its capacity is reached. Under strategy B, however, the actual demand is known either before the tour starts (customers call, or the operator calls them, or in the case of package deliveries the addresses are known) or becomes gradually known (the driver calls the depot for the next visit), so that savings can occur by skipping customer locations with zero demand. Notice that depending on how information about customer demand becomes known redesigning the routes might or might not be a possibility.

An important question that arises is how to choose the a priori sequence. One possible solution is to choose the a priori sequence as the one with *minimal expected total length*. This value corresponds to the expected total length of the fixed set of routes plus the expected value of the extra distance that might be required by a particular realization of the demand. The extra distance is due to the fact that demand on the route may occasionally exceed the capacity of the vehicle and force it to go back to the depot before continuing on its route. The problem of selecting the a priori sequence of minimum expected length is called the stochastic VRP (SVRP). We now give some examples in which VRPs with stochastic demand arise.

### Application Areas

In a strategic planning scenario, consider a delivery and collection company which has decided to begin service in a particular area. The company has carried out a market survey and identified a number $n$ of potential major customers who, during any collection/distribution period, have a significant probability of requiring a visit. The company wishes to estimate the resources necessary to serve these customers. At this stage of planning, the company can only assign probability distributions for the demand of all potential customers. To address the planning problem the company will wish to estimate approximately the expected amount of travel that will be necessary on a typical day to serve the subset of the $n$ customers that will require a visit.

In a routing context, Lambert, Laporte and Louveaux (1993) consider the problem in which a central bank has to collect cash on a daily basis from several but not all of its branches. The capacity $q$ of the vehicle used may not correspond to any physical constraint but to an upper bound on the amount of cash that a vehicle might carry for safety reasons. The supply of cash at each particular branch is stochastic and has a distribution which may be different among branches. The bank faces a similar problem, when it wishes to deliver cash to different automatic teller machines.

In a distribution context, the delivery of packages from a post office has important stochastic components, where the

probability that a certain building requires a visit can be found from historical data and the capacity $q$ corresponds to the physical constraint that a truck can carry only a fixed weight or volume. Other examples reported in the literature include a "hot meals" delivery system (Bartholdi et al. 1983) and routing of forklifts in a cargo terminal or in a warehouse.

## 6.1. Problem Definition

Given a complete network, let the nodes be $\{0, 1, \ldots, n\}$, where node 0 denotes the depot and the set $V = \{1, 2, \ldots, n\}$ denotes the set of customer locations. The distances $d(i, j)$ are assumed to be symmetric and to satisfy the triangle inequality: $d(i, j) \leq d(i, k) + d(k, j)$. Let the capacity of the vehicle be $q$ and let $D_i, i = 1, \ldots, n$ be the random variable that describes the demand of customer $i$. We assume that the probability distribution of $D_i$ is discrete and is known. Let $p_i(k) = Pr\{D_i = k\}, i = 1, \ldots, n$ and $k = 0, 1, \ldots, K$. We further assume that $K \leq q$, i.e., no single location has a demand exceeding the capacity $q$. We further assume that the demands are independent.

There are $(K + 1)^n$ possible realizations of the demand and therefore $(K + 1)^n$ possible instances of the problem. Given demands $D_1, D_2, \ldots, D_n$, denote by $L^*_{VRP}(D_1, D_2, \ldots, D_n)$ the length of total distance traveled in the optimal solution to the associated VRP. Note that since the demand is stochastic this is a random variable. We call the expectation of this random variable the expected length of *the re-optimization strategy*, since we redesign (re-optimize) the routes at every problem instance. This expected length is thus given by

$$E[Z^*_{REOPT}]$$
$$= \sum_{i_1, \ldots, i_n} p_1(i_1) \ldots p_n(i_n) L^*_{VRP}(i_1, \ldots, i_n), \quad (7)$$

where the summation is over all demand instances for the nodes. Clearly, the exact estimation of $E[Z^*_{REOPT}]$ is a computationally intractable problem, since it involves $(K + 1)^n$ terms, each of which involves the exact solution of a VRP. So, in a strategic planning scenario, in which a company needs to have an estimate of the expected travel cost, the expected length of the re-optimization strategy is not a realistic alternative computationally.

Let us now consider the two proposed a priori strategies A and B. Given an a priori sequence $\tau$ let $L^i_\tau(i_1, \ldots, i_n)$ be the length of the a priori sequence $\tau$ which will result under strategy $i = a, b$ if the demand pattern is $i_1, \ldots, i_n$. We denote with

$$E[L^a_\tau] = \sum_{i_1, \ldots, i_n} p_1(i_1) \ldots p_n(i_n) L^a_\tau(i_1, \ldots, i_n), \quad (8)$$

the expected length of the a priori sequence $\tau$ under strategy A and

$$E[L^b_\tau] = \sum_{i_1, \ldots, i_n} p_1(i_1) \ldots p_n(i_n) L^b_\tau(i_1, \ldots, i_n), \quad (9)$$

the expected length of the a priori sequence $\tau$ under strategy B.

The VRP with stochastic demand (SVRP) is then defined as the following optimization problem.

## Problem SVRP

$$E[Z^*_i] = \min_\tau E[L^i_\tau], \quad \text{for } i = a, b.$$

The similarity of (7) with (8) and (9) might lead one to think that the evaluation of (8) and (9) is also an intractable problem. We next show, however, that there is an efficient algorithm for the evaluation of (8) and (9).

Bertsimas proposes an $O(K^2 n^2)$ algorithm to compute the expected sequences under both strategies A and B, which implies that once an a priori sequence is selected its performance can be evaluated efficiently.

We next address the central questions of analytically comparing the performance of the re-optimization strategy and the fixed (A) and adaptive (B) a priori strategies. To achieve this goal we propose heuristics for the SVRP and analyze their performance from both a worst case and average point of view.

## 6.2. Worst-Case Analysis

The following heuristic proposed in Bertsimas is a modification of the tour partitioning heuristic for the deterministic VRP.

## Cyclic Heuristic

1. Given an initial sequence $\tau \triangleq \tau_1 = (0, 1, 2, \ldots, n, 0)$, consider the sequences $\tau_i = (0, i, \ldots, n, 1, \ldots, i - 1, 0), i = 2, \ldots, n$.
2. Compute $E[L^a_{\tau_i}]$ for all $i = 1, \ldots, n$.
3. The sequence with the minimum expected length among $E[L^a_{\tau_i}], i = 1, \ldots, n$ is the proposed solution $\tau_H$ to the SVRP under the fixed strategy A.

When the initial sequence is an $\alpha$-optimal traveling salesman tour the heuristic is denoted $CH(\alpha)$ and its value is $E[Z^{CH(\alpha)}_a]$. Let $E[Z^*_a]$ be the expected length of the optimal sequence under the fixed strategy A. The following theorem is proved in Bertsimas.

**Theorem 13.** *Assume that the demands of the customers are identically distributed. If the initial sequence given to the cyclic heuristic is an $\alpha$-optimal traveling salesman tour, then under the triangle inequality*

$$\frac{E[Z^{CH(\alpha)}_a]}{E[Z^*_a]} \leq 1 + \alpha + \frac{1}{n}\left(\frac{q}{E[D]} - 1\right) = 1 + \alpha + O\left(\frac{1}{n}\right).$$

Moreover, if $Pr\{D = 0\} = 0$, i.e., all customers have some demand, then one can strengthen the previous bound to $E[Z^{CH(\alpha)}_a]/E[Z^*_{REOPT}] \leq 1 + \alpha + O(1/n)$, which means that the cyclic heuristic is within a factor of $1 + \alpha$ of the re-optimization strategy.

The space filling curve heuristic has been introduced and analyzed by Platzman and Bartholdi (1989) for the

Euclidean traveling salesman problem. Bertsimas and Howell (1992) analyze it for the adaptive strategy B for Euclidean problems. The heuristic can be described as follows.

## Space-Filling Curve Heuristic

1. Given the $n$ coordinates $(x_i, y_i)$ of the points in the plane compute the number $f(x_i, y_i)$ for each point. The function $f: R^2 \to R$ is called the Sierpinski curve (for details on the computation of $f(x, y)$ see Platzman and Bartholdi).
2. Sort the numbers $f(x_i, y_i)$ and visit the corresponding initial points $(x_i, y_i)$ in that order, producing a tour $\tau_{SF}$, which is the proposed a priori sequence for the SVRP under strategy B.

Let $E[Z_b^{SF}]$ be the expected length of the heuristic for the SVRP under strategy B. Let $E[Z_b^*]$ be the expected length of the optimal sequence under strategy B. Bertsimas and Howell show that:

**Theorem 14.** *For the Euclidean instances with arbitrary demand distributions*

$$\frac{E[Z_b^{SF}]}{E[Z_b^*]} \leqslant \frac{E[Z_b^{SF}]}{E[Z_{REOPT}^*]} = O(\log n).$$

Bertsimas and Grigni (1989) show that there exists an example in which the logarithmic bound is achievable. Note that previous theorems bound the degree of suboptimality of a priori optimization compared with the re-optimization strategy in the worst case. Gendreau, Laporte and Seguin (1993) propose an exact algorithm for the **SVRP** and solve to optimality problems with up to 50 nodes. We next illustrate that on average the results can be improved sharply.

### 6.3. Average-Case Analysis

Let $X_1, X_2, \ldots$ be an infinite sequence of independent, identically distributed random points in the unit square and assume that the depot is at $(0, 0)$. Let $E[r]$ be the expected distance from the origin and let $X^{(n)}$ denote the first $n$ points of the sequence. Let $E[Z_{REOPT}^*(X^{(n)})]$ be the expected length of the re-optimization strategy and $E[Z_a^*(X^{(n)})]$, $E[Z_b^*(X^{(n)})]$ be the expected lengths of the two a priori strategies $a$, $b$, respectively. Let $E[Z_a^{CH(NN)}(X^{(n)})]$ be the expected length of the cyclic heuristic if the initial sequence given to the cyclic heuristic is the nearest neighbor tour.

Let $D$ be the random variable that describes the demand of each customer. Bertsimas (1992) for the capacitated case (case 1 in Theorem 15) and Jaillet (1988) for the uncapacitated case (case 2 in Theorem 15) prove:

**Theorem 15.** *The asymptotic behavior of the three updating strategies is:*

1. *If $q$ is a constant (does not depend on $n$), then with probability 1*

$$\lim_{n\to\infty} \frac{E[Z_{REOPT}^*(X^{(n)})]}{n} = \lim_{n\to\infty} \frac{E[Z_a^{CH(NN)}(X^{(n)})]}{n}$$

$$= \lim_{n\to\infty} \frac{E[Z_a^*(X^{(n)})]}{n} = \lim_{n\to\infty} \frac{E[Z_b^*(X^{(n)})]}{n}$$

$$= \frac{2E[r]E[D]}{q}.$$

2. *If $\lim_{n\to\infty} q/\sqrt{n} = \infty$, and $p = Pr\{D > 0\}$ then almost surely*

$$\lim_{n\to\infty} \frac{E[Z_{REOPT}^*(X^{(n)})]}{\sqrt{n}} = \beta_{TSP}\sqrt{p},$$

$$\lim_{n\to\infty} \frac{E[L_{\tau_a}^a(X^{(n)})]}{\sqrt{n}} = \beta_{TSP},$$

$$\lim_{n\to\infty} \frac{E[Z_b^*(X^{(n)})]}{\sqrt{n}} = \beta(p),$$

*where $\beta_{TSP}\sqrt{p} \leqslant \beta(p) \leqslant \min[0.92\sqrt{p}, \beta_{TSP}]$, with $\beta_{TSP}$ the constant appearing in the celebrated Beardwood, Halton and Hammersley paper.*

Notice that in the capacitated case the cyclic heuristic initialized by the nearest-neighbor tour is asymptotically optimal and equivalent to the re-optimization strategy, which is particularly important since nearest-neighbor tours are widely used in practice. In the uncapacitated case the fixed strategy A is within a factor of $1/\sqrt{p}$, while strategy B is within a small constant factor. From the re-optimization strategy in terms of performance, similar asymptotic theorems can be proved for the case that the distribution of customer locations has a continuous part with density $f(x)$.

### 6.4. Reflections

The previous results attempted to offer analytical evidence that the performance of a priori strategies for the SVRP are close to the strategy of re-optimization when the number of customers tends to infinity. This is true from worst case as well as average case points of view.

From a practical standpoint, we believe that the a priori strategies provide an alternative to the strategy of re-optimization, and they can be useful in the absence of intense computational power. Bertsimas and Howell (1992) for the uncapacitated and Bertsimas, Chervi and Peterson (1995) for the capacitated case report extensive computational results which support the results of the analytical investigations.

## 7. STOCHASTIC AND DYNAMIC VEHICLE ROUTING IN EUCLIDEAN SERVICE REGIONS

As we have seen in the previous sections the classical view of VRP is static and deterministic. The previous section represented a partial departure from this paradigm, addressing the VRP with stochastic but static demands. Yet, in many of the practical applications in which VRPs arise (e.g., distribution, inventory resupply, mobile repair), there

are significant stochastic but also *dynamic* components to the problem. Indeed, in many real logistics systems, demands arrive randomly in time and thus routing is a continuous process of collecting demands, forming tours, and dispatching vehicles.

As a canonical example of a logistics application with strong probabilistic and dynamic components, consider the following utility repair problem: A utility firm (electric, gas, water and sewer, highway, etc.) is responsible for maintaining a large, geographically dispersed facilities network. The network is subject to failures (major and minor) which occur randomly both in time and space (location). The firm operates a fleet of repair vehicles which are dispatched from a depot to respond to failures. Routing decisions are made based on a real-time log of current failures and perhaps some characterization of the future failure process. Vehicle crews spend a random amount of time servicing each failure before they are free to move on to the next failure. The firm would like to operate its fleet in a way that minimizes the average downtime due to failures.

There are other closely related problems to this canonical example that arise in practice. For instance, consider a firm that delivers a product from a central depot to customers based on orders that arrive in real time. Orders are entered into a log and delivery vehicles are dispatched with the objective of minimizing some combination of the delivery cost and the average wait for delivery (service level). Such an order process is likely to be found in firms that serve a large population of customers (or potential customers) each of whom orders relatively infrequently (e.g., home heating oil distributors, mail order firms, etc.).

Additional important examples are found in finished goods distribution and freight consolidation. Consider, for example, an automobile manufacturer. Cars are produced at an assembly plant and put into finished goods inventories (parking lots) to await distribution by a fleet of car-hauling trucks. Each car is designated for a particular dealer. Conceptually, the inventory can be thought of as a "log" of locations that must be visited by the delivery vehicles. New entries to this log are made every time a new vehicle is added to the inventory, and entries are deleted when automobiles are delivered to their designated dealers. For a fixed production rate, lowering the waiting time in this case is, by Little's theorem, equivalent to reducing the inventory of finished goods either on the lot or in transit (i.e., the "pipeline" inventory). The manufacturer would like to minimize the sum of its delivery and holding costs and also provide quick delivery times to its dealers.

Similar distribution problems are found in freight consolidation (e.g., less-than-truckload (LTL) shipping) and parcel post systems. Here distribution centers receive partial loads designated for specific locations in a service region. These partial loads are queued (stored in a distribution terminal) and eventually consolidated into full truckloads for delivery. Whereas travel cost is an important consideration in these systems, lowering the wait for delivery is also a concern, both for improving the delivery time and, as in the previous example, for reducing inventory costs (terminal space, insurance costs, etc.).

For these applications, static, deterministic vehicle routing models do not capture the essential tradeoffs needed to understand and effectively operate these systems. Indeed, we see that minimizing travel distance—the classical objective—is often less important than minimizing delivery times. This is true for several reasons:

1. the firm may compete primarily on the basis of service level rather than delivery cost;
2. inventory costs may dominate delivery costs, in which case the firm may want to quickly off-load its stock of finished goods;
3. the firm may be employing a just-in-time (JIT) inventory and production policy and would like, as a matter of operating policy, to minimize pipeline inventory and speed deliveries;
4. many of the operating costs (labor, depreciation, terminal costs) may be fixed and thus minimizing distance has a marginal effect on reducing costs; or
5. waiting time *is* the objective, as is the case in the utility repair application just mentioned.

It is unfortunate that classical models and techniques have little to say about vehicle routing when stochastic and dynamic elements are included. This is due in large part to the inherent difficulties of combining vehicle routing and congestion models. In particular, including a time element along with randomness usually destroys the combinatorial structure required for classical vehicle routing methods. Similarly, the strong dependencies present in travel times usually violate the assumptions required to apply traditional queueing models. Indeed, Psaraftis (1988) points out that although congestion (queueing) and vehicle routing theory are both very rich subjects, little work has been done to combine them. Recently, a *theory of vehicle routing under congestion* has been developed in Bertsimas and van Ryzin (1990, 1992, 1993) and Bertsimas and Xu (1992). In this section, we review these developments.

We remark that the scope of our review is strictly limited to dynamic vehicle routing problems, which is a subset of a broader family of models called *dynamic transportation models*. This broader area includes (we provide just one reference per area):

- dynamic fleet management (Powell 1986);
- dynamic traffic assignment (Friesz et al. 1989);
- dynamic air traffic control (Vranas, Bertsimas and Odoni 1993);
- dynamic shortest path problems (Psaraftis and Tsitsiklis 1993).

## 7.1. Problem Definition

The problem we investigate, which we call the dynamic traveling repairman problem (DTRP) is a Euclidean

model of a dynamic VRP. Demands for service arrive according to a renewal process with intensity $\lambda$ and finite variance to a connected, bounded Euclidean service region $\mathcal{A}$ of area $A$. Upon arrival, demands assume an independent and identically distributed (i.i.d.) location in $\mathcal{A}$ according to a continuous density $f(x)$ defined over $\mathcal{A}$. Demands are serviced by $m$ identical vehicles that travel at constant velocity $v$. At each location, vehicles spend some time $s$ in on-site service, which is a random variable with finite first and second moments denoted by $\bar{s}$ and $\overline{s^2}$, respectively. Successive service times are independent and identically distributed.

Initially, we shall assume that there are no capacity constraints on the vehicles. Later, we consider the case where there is an upper bound, $q$, on the number of demands that can be served before a vehicle must return to a designated depot location.

A policy for routing the vehicles is called stable if the number of unserved demands in the system is bounded almost surely for all times $t$. Let $\mathcal{M}$ denote the set of stable policies. If a policy is stable, $\rho \equiv \lambda \bar{s}/m$ is the fraction of time vehicles spend in on-site service. We write $T_\mu$ to indicate the system time of a particular policy $\mu \in \mathcal{M}$. The DTRP is then defined as the following optimization problem.

### Problem DTRP

$$\min_{\mu \in \mathcal{M}} T_\mu .$$

We let $T^*$ denote the optimal value in this minimization. The model can be also extended to handle a mixed objective that includes costs on both waiting time and the average distance traveled per demand served. This extension is discussed next.

### 7.2. The Uncapacitated Vehicle DTRP

In this subsection, we assume that the vehicles are uncapacitated. The principal question we address is how the optimal system time behaves as a function of the parameters of the problem ($\rho$, $m$, $f(x)$, etc.), how the geometry of the system affects the performance, when the system is stable, and what are optimal or near-optimal policies that are easily implementable.

To address the dependence of the optimal system time on the demand distribution we distinguish two cases. Under the class of *spatially unbiased policies*, the system time of a random demand does not depend on its location, while if we do not impose this restriction, the class of policies is called *spatially biased*. The motivation for considering this distinction is that it offers insights in the tradeoff of performance and "fairness." In the spatially unbiased case, the average delay experienced by customers is the same regardless of their location, while in the spatially biased case customer location might affect the experienced delay. In the spatially unbiased case we create "fairness" by distributing a higher average delay equally among customers, while in the spatially biased case we

achieve better overall performance (smaller delay) at the expense of creating inequalities among customer delays.

In light traffic ($\rho \to 0$) the optimal policy is as follows.

### Stochastic Median Policy

Solve the $m$ median problem and locate the $m$-vehicles in the medians $x_1^*, \ldots, x_m^*$ of $\mathcal{A}$. Serve customers in a first-come, first-serve (FCFS) order, returning to the corresponding median after each service is completed.

Bertsimas and van Ryzin (1990) show that the optimal expected system time, $T^*$, in this case satisfies:

**Theorem 16.** *In light traffic, the m median policy is optimal and satisfies*

$$T^* \to \frac{E[\|X - x^*\|]}{v} + \bar{s} \quad as \ \rho \to 0 .$$

The principal insight from Theorem 16 is that in light traffic, the problem is essentially a location problem. In heavy traffic ($\rho \to 1$), however, the behavior radically changes. For the case of $m$ uncapacitated vehicles with general distribution $f(x)$ the following bounds are derived in Bertsimas and van Ryzin (1990, 1992).

**Theorem 17.** *The system is stable if and only if $\rho < 1$. Under the class of spatially unbiased policies the optimal system time satisfies*

$$\gamma^2 \frac{\lambda[\int_{\mathcal{A}} f^{1/2}(x) \, dx]^2}{m^2 v^2 (1-\rho)^2} \leq T^* \leq \frac{\beta_{TSP}^2}{2} \frac{\lambda[\int_{\mathcal{A}} f^{1/2}(x) \, dx]^2}{m^2 v^2 (1-\rho)^2}$$
$$+ O\left(\frac{1}{(1-\rho)^{3/2}}\right) , \tag{10}$$

*where $\gamma = 2/3\sqrt{\pi}$ and $\beta_{TSP}$ is the constant appearing in the theorem of Beardwood, Halton and Hammersley.*

*Under the class of spatially biased policies the optimal system time satisfies*

$$\gamma^2 \frac{\lambda[\int_{\mathcal{A}} f^{2/3}(x) \, dx]^3}{m^2 v^2 (1-\rho)^2} \leq T^* \leq \frac{\beta_{TSP}^2}{2} \frac{\lambda[\int_{\mathcal{A}} f^{2/3}(x) \, dx]^3}{m^2 v^2 (1-\rho)^2}$$
$$+ O\left(\frac{1}{(1-\rho)^{3/2}}\right) . \tag{11}$$

The previous results reveal several interesting features of the problem.

1. The stability condition is independent of the geometry of the system and, surprisingly, of the speed $v(v > 0)$.
2. Compared with traditional queues in which the system time behaves like $\Theta(1/1 - \rho)$, the system time in the **DTRP** grows much faster, like $\Theta(1/(1 - \rho)^2)$. The reason is the distributed character of the system that gives rise to spatial queues.
3. The effect of having $m$ vehicles in the system time is nonlinear. Doubling the number vehicles (but keeping the same traffic intensity) decreases the system time by a factor of 4.
4. The theorem characterizes how the demand distribution influences the system time. The system time under

the class of spatially unbiased policies is higher than the system time under the class of spatially biased policies. Note, however, that in the case of uniform distribution ($f(x) = 1/A$) the optimal policy will necessarily be spatially unbiased and the behavior will be $\Theta(\lambda A/m^2 v^2(1 - \rho)^2)$. In other words, in the case of uniform distribution there is no advantage of using spatially biased policies.

The lower bounds in Theorem 17 are established for all stable policies using arguments from geometrical probability, queueing theory, and optimization. The upper bounds are established from analyzing specific policies. Bertsimas and van Ryzin (1990, 1992) propose several stable policies $\mu$ that have the same asymptotic behavior as the bounds of Theorem 17, i.e., for the uniform case ($f(x) = 1/A$), for example:

$$T_\mu \sim \gamma_\mu^2 \frac{\lambda A}{v^2 m^2 (1 - \rho)^2} \quad \text{as } \rho \to 1 \,,$$

where the constant $\gamma_\mu$ depends only on the policy $\mu$. Using a variety of methods, Bertsimas and van Ryzin (1990) analyze several heuristic policies, among which two are particularly interesting because of their simplicity and ease of implementation: Visit the next customer that is given from the space filling curve heuristic of Platzman and Bartholdi, and serve the nearest neighbor. The system time under these two policies behaves for the uniform case like $\gamma_\mu^2 \lambda A/v^2 m^2 (1 - \rho)^2$, even for a moderate value of $\rho$. The constants $\gamma_\mu$ for these policies are estimated from simulation. Though not competitive with the best traveling salesman policies that will be discussed, these heuristics are simple to implement and of low complexity, which may make them particularly attractive in practice.

The provably best policy found until now, which achieves the upper bound in (10), is one based on batching arrivals into sets of a fixed size, forming optimal traveling salesman tours on these sets, and then serving these tours FCFS as in a $GI/G/1$ queue and is described as follows.

## The Unbiased (U) Traveling Salesman Policy

Let $k$ be a fixed positive integer. From a central point in the interior of $\mathcal{A}$, subdivide the service region into $k$ wedges $\mathcal{A}_1, \mathcal{A}_2, \ldots, \mathcal{A}_k$ such that $\int_{\mathcal{A}_i} f(x) \, dx = 1/k$, $i = 1, 2, \ldots, k$. (One could do this by "sweeping" the region from the depot using an arbitrary starting ray until $\int_{\mathcal{A}_1} f(x) \, dx = 1/k$, continuing the sweep until $\int_{\mathcal{A}_2} f(x) \, dx = 1/k$, etc.) Within each subregion, form sets of size $n/k$ ($n$ is a parameter to be determined). As sets are formed, deposit them in a queue and service them FCFS with the first available vehicle by forming a traveling salesman tour on the set and following it in an arbitrary directions. Optimize over $n$.

In the previous policy we needed to select a priori a parameter $n$. An adaptive variant of the previous policy with exactly the same behavior was proposed in Bertsimas and Xu.

## The Unbiased (U) Adaptive Traveling Salesman Policy

Subdivide the service region into $m$ smaller regions with the same expected demand. Each vehicle operates independently in the smaller region according to the following policy: Let $k$ be a fixed positive integer. Subdivide the smaller service region into $k$ subregions $\mathcal{A}_1, \mathcal{A}_2, \ldots, \mathcal{A}_k$ such that $\int_{\mathcal{A}_i} f(x) \, dx = 1/km$, $i = 1, 2, \ldots, k$. Design a tour among the subregions. Follow this tour to move from one subregion to another. Within each subregion, the demands that are present at the time the server enters the subregion are serviced using a traveling salesman tour. After the server finishes the current traveling salesman tour, it will move to the next subregion in the a priori tour and the process is repeated.

The unbiased traveling salesman policy is a parametric policy, because it uses parameters $k, n$ that depend on the data of the problem. In contrast, the policy just described is adaptive, in the sense that it does not use any parameters and adapts its behavior even when the data of the problem are changing (the vehicles serve all the customers that they find in each subregion). As already mentioned, the system time under both traveling salesman policies (parametric and adaptive) is the same and is given in the upper bound of (10). For the case of biased policies the best policy found to date that achieves the upper bound of (11) is as follows.

## The Biased (B) Traveling Salesman Policy

Approximate the arbitrary density $f(x)$ with a piecewise constant density. Let $\mathcal{A}_1, \mathcal{A}_2, \ldots, \mathcal{A}_J$ be a partition of $\mathcal{A}$ such that $f(x) = \mu_j$ for all $x \in \mathcal{A}_j, j = 1, 2, \ldots, J$. Let $A_j$ denote the area of $\mathcal{A}_j$. For a given positive integer $k$, partition each subset $\mathcal{A}_j$ further into $k_j = \mu_j^{2/3} A_j k$ regions of area $A_j/k_j = (\mu_j^{2/3} k)^{-1}$ ($k$ is a scale factor that will be chosen arbitrarily large; hence, we assume an integer $k_j$ can be found such that $k_j/k$ is sufficiently close to $\mu_j^{2/3} A_j$). Within each of these subregions, form demands into sets of size $n/k$ as they arrive. As sets are formed, deposit them in a queue and service them FCFS with the first available vehicle as follows:

1. form a traveling salesman tour on the set;
2. connect the tour to the depot through an arbitrary point in the tour; and
3. follow the resulting tour in an arbitrary direction servicing demands as they are encountered. Optimize over $n$.

These results imply that compared with the lower bounds in Theorem 17 the traveling salesman policies are guaranteed to be within about 80% of the optimal policy in heavy traffic. We conjecture, however, that these policies are indeed optimal in heavy traffic, i.e., the factor of 1.8 from optimal is due to slack in the lower bound of Theorem 17. Preliminary work in Papastavrou and Chandru (1992) and Bertsimas and Xu (1992) support the conjecture.

The traveling salesman policies can be extended to a mixed objective involving both waiting time and travel cost. By increasing the size $n$ of the sets that are formed, travel distance per demand can be reduced at the expense of increasing the mean system time. Indeed, one can show that to minimize system time, we essentially maximize the amount of travel per demand served. Thus, travel cost and system time are conflicting objectives that can be balanced by sizing routes in an appropriate way. One could place costs on *both* the travel time per demand and the system time and select a set size that minimizes their sum. The result would be to form larger sets which improve travel efficiency at the expense of increased system time.

## 7.3. Capacitated Dynamic Vehicle Routing

The scenario is the same as before except that the region $\mathcal{A}$ is now serviced by a homogeneous fleet of $m$ vehicles operating out of a set $\mathcal{D}$ of $|\mathcal{D}| = m$ depots, where each vehicle is restricted to visiting at most $q$ customers before returning to its respective depot. (The depot locations need not be distinct.) Let $\bar{r}$ be the expected distance of a random demand from the closest of the $m$ depots.

To address the question of how the performance of the optimal policy depends on the system's parameters Bertsimas and van Ryzin (1993) prove the following theorem.

**Theorem 18.** *The system is stable if and only if $\rho + 2\lambda\bar{r}/mvq < 1$. The optimal system time satisfies:*

$$\frac{\gamma^2}{9} \frac{\lambda\left(1 + \frac{1}{q}\right)^2 G}{m^2 v^2 \left(1 - \rho - \frac{2\lambda\bar{r}}{mvq}\right)^2} \leq T^*$$

$$\leq \frac{\beta^2}{2} \frac{\lambda\left(1 + \frac{1}{q}\right)^2 G}{m^2 v^2 \left(1 - \rho - \frac{2\lambda\bar{r}}{mvq}\right)^2}$$

$$+ O\left(\frac{1}{\left(1 - \rho - \frac{2\lambda\bar{r}}{mvq}\right)^{3/2}}\right), \tag{12}$$

*where $G = [\int_{\mathcal{A}} f^{1/2}(x) \, dx]^2$, $[\int_{\mathcal{A}} f^{2/3}(x) \, dx]^3$ for the spatially unbiased and spatially biased cases, respectively, and $\gamma$ is the same numerical constant from the uncapacitated bound (10).*

These capacitated **DTRP** results provide some intuitively satisfying insights. Unlike the uncapacitated case the stability condition $\rho + 2\lambda\bar{r}/mvq < 1$ depends on the geometry of the system through $\bar{r}$ and on the speed $v$; however, for $q \to \infty$ the dependence vanishes. The second term in this stability condition has the interpretation of a *radial collection cost* in the sense of Haimovich and Rinnooy Kan. That is, $2\bar{r}/v$ is essentially the average time required to reach a set of $q$ customers from the nearest depot (the radial cost). Dividing by $q$ gives the average radial travel time per customer, and, hence, multiplying by $\lambda$ we obtain the fraction of time the server spends in radial travel. The previous condition says that as long as this fraction plus the fraction

of time spent on-site is less than one, the system will be stable. Furthermore, the waiting time grows like the inverse square of the stability difference, $1 - \rho - 2\lambda\bar{r}/vq$, just as it does in the uncapacitated case. Note that the average radial distance $\bar{r}$ plays a crucial role in the system's behavior in this case. Indeed, Bertsimas and van Ryzin (1993) show that if one has the option of locating the depots anywhere within $\mathcal{A}$, then minimizing $\bar{r}$ (i.e., locating the depot at the medians) is always optimal in heavy traffic.

For $q$ finite, Bertsimas and van Ryzin (1993) construct policies, $\mu$, for which

$$T_\mu \sim \gamma_\mu^2 \frac{\lambda A \left(1 - \frac{1}{q}\right)^2}{m^2 v^2 \left(1 - \rho - \frac{2\lambda\bar{r}}{mqv}\right)^2}$$

as $\rho + 2\lambda\bar{r}/mqv \to 1$, and therefore have a constant factor guarantee from the optimal policy.

The best of these is a policy based on modifying the traveling salesman policies using the tour partitioning heuristic of Haimovich and Rinnooy Kan. We consider first the case that all depots coincide.

### The Dynamic Tour Partitioning Policy

For some fixed integer $k \geq 1$, divide $\mathcal{A}$ into $k$ subregions $\mathcal{A}_1, \mathcal{A}_2, \ldots, \mathcal{A}_k$, such that $\int_{\mathcal{A}_i} f(x) \, dx = 1/k$ $i = 1, 2, \ldots, k$ using radial cuts centered at the depot. Within each region, collect demands into sets $N_1, N_2, \ldots$ of size $n/k$ as they arrive and construct optimal tours on these sets. Starting at a randomly selected point in $N_1$, split the tour into $l = \lceil n/q \rceil$ segments of $q$ demands each (except, perhaps, for the last segment). Connect the end points of the segments to the depot to form $l$ tours of at most $q$ demands each. As sets are formed deposit them in a queue. Service the queue FCFS with the first available vehicle by following the collection of tours. Optimize over $n$.

This policy achieves the upper bound in (12). Bertsimas and van Ryzin (1993) generalize the results of Theorem 18 when there are exactly $p$ vehicles per depot, i.e., $m = kp$ under the following symmetry assumption. Suppose these $k$ depots induce Voronoi cells that are identical in shape and size. Given a collection of points $X$, a Voronoi cell around a point $O$ is the set of all points which are closer to $O$ than to any other point in $X$; see Preparata and Shamos (1985). Then, if one applies a $p$-vehicle policy in each cell, the resulting system time will be within a constant factor of the lower bound in heavy traffic. This is due to the fact that each cell has an arrival rate of $\lambda/k$ and serves an area of size $A/k$, each of which has the same mean radial distance $\bar{r}$. Therefore, since each region operates with $p$ vehicles we have

$$T \sim \frac{\beta^2 (\lambda/k)(A/k)}{2p^2 v^2 \left(1 - \rho - \frac{2(\lambda/k)\bar{r}}{vqp}\right)^2}$$

$$= \frac{\lambda\beta^2 A}{2m^2 v^2 \left(1 - \rho - \frac{2\lambda\bar{r}}{vqm}\right)^2} \quad \text{as } \rho + \frac{2\lambda\bar{r}}{qvm} \to 1$$

and, hence, the policy has a constant factor performance guarantee.

If $k$ is large and the depots are located at the $k$ median locations, then Haimovich and Magnanti (1988) show that the Voronoi cells approach a uniform, hexagonal partition of $A$ (i.e., a honeycomb pattern). Since this simultaneously produces uniform Voronoi cells and minimizes $\bar{r}$, one can show that assigning $p$ vehicles to each of the $k$ medians is again provably good.

## 7.4. Reflections

The analysis of dynamic VRPs yields simple expressions for the system time that provide structural insight into the effects of traffic intensity, on-site service characteristics, the number, speed, and capacity of vehicles employed, service region size, and the distribution of customer locations.

A reoccurring finding in the analysis is that static vehicle routing methods when properly adapted can yield optimal or near-optimal policies for dynamic routing problems *with uniform, stationary demands*. This is an encouraging result on several levels. On a theoretical level, it suggests that there is, indeed, a connection between static and dynamic problems at least for models with uniform stationary demands; that is, the DTRP has geometrical characteristics that are intimately related to the corresponding characteristics for static VRPs. On a practical level, the results imply that most of the exact algorithms, heuristics, and insights which have been developed over years of investigation of static VRPs can, in fact, form the basis for effective policies in dynamic environments.

## 8. SUMMARY

We attempt in this final section to present a number of important observations.

1. We believe that there are three primary benefits for analytical analysis of combinatorial problems in general and VRPs in particular: a) analytical analysis fosters new insights into the algorithmic structure required to solve large-sized problems (see the LBH); b) it makes it possible to analyze the performance of classical heuristics (see the analysis of the class route first-cluster second); and c) it leads to a better understanding of models that integrate vehicle routing with other issues important to the firm, such as inventory control. As pointed by one of the referees, we should qualify this statement with the observation that, in practice, issues such as variation in travel times, and crew scheduling complicate VRPs significantly. We hope, however, that this deeper understanding will have an impact on vehicle routing practice.

2. Asymptotic analysis leads to interesting, qualitative insights on the structure of the asymptotic optimal solution of both static and dynamic VRPs. These insights can be used to develop new algorithms that generate solutions with a structure similar to the asymptotic structure. Moreover, preliminary computational results indicate that even for moderate problem sizes, the limiting behavior is indeed present; see Bramel and Simchi-Levi (1993) for the performance of the algorithm developed for the VRP with capacity and time window constraints, Anily and Federgruen (1990) for the performance of algorithms developed for inventory-routing problems, and Bertsimas and van Ryzin (1993) for the dynamic vehicle routing problem.

3. A priori optimization in VRPs is an attractive policy in the absence of intensive computational power.

4. In dealing with stochasticity in the VRP, new structural insights are gained by considering the VRP under congestion. These insights can guide the construction of practical algorithms for VRPs in a dynamic and stochastic environment.

## ACKNOWLEDGMENT

## REFERENCES

ALTINKEMER, K., AND B. GAVISH. 1987. Heuristics for Unequal Weight Delivery Problems with a Fixed Error Guarantee. *O. R. Letts.* **6**, 149–158.

ALTINKEMER, K., AND B. GAVISH. 1990. Heuristics for Equal Weight Delivery Problems With Constant Error Guarantees. *Trans. Sci.* **24**, 294–297.

ANILY, S., AND A. FEDERGRUEN. 1990. One Warehouse Multiple Retailer Systems With Vehicle Routing Costs. *Mgmt. Sci.* **36**, 92–114.

ANILY, S., AND A. FEDERGRUEN. 1993. Two-Echelon Distribution Systems With Vehicle Routing Costs and Central Inventories. *Opns. Res.* **41**, 37–47.

BALINSKI, M. L., AND R. E. QUANDT. 1964. On an Integer Program for a Delivery Problem. *Opns. Res.* **12**, 300–304.

BARTHOLDI, J. J., L. K. PLATZMAN, R. LEE COLLINS AND W. H. WARDEN. 1983. A Minimal Technology Routing System for Meals on Wheels. *Interfaces* **13**, 1–8.

BEASLEY, J. 1983. Route First-Cluster Second Method for Vehicle Routing. *Omega* **11**, 403–408.

BEARDWOOD, J., J. HALTON AND J. HAMMERSLEY. 1959. The Shortest Path Through Many Points. *Proc. Cambridge Phil. Soc.* **55**, 299–327.

BERTSIMAS, D. 1992. A Vehicle Routing Problem With Stochastic Demand. *Opns. Res.* **40**, 574–585.

BERTSIMAS, D., AND G. VAN RYZIN. 1990. A Stochastic and Dynamic Vehicle Routing Problem in the Euclidean Plane. *Opns. Res.* **39**, 601–615.

BERTSIMAS, D., AND G. VAN RYZIN. 1992. Stochastic and Dynamic Vehicle Routing With General Demand and Interarrival Time Distributions. *Adv. Appl. Prob.* **25**, 947–978.

BERTSIMAS, D., AND G. VAN RYZIN. 1993. Stochastic and Dynamic Vehicle Routing in the Euclidean Plane With Multiple Capacitated Vehicles. *Opns. Res.* **41**, 60–76.

BERTSIMAS, D., AND H. XU. 1992. Optimal Adaptive Policies for Stochastic and Dynamic Vehicle Routing Problems. ORC Working Paper, MIT (submitted for publication).

BERTSIMAS, D., AND L. HOWELL. 1992. Further Results on the Probabilistic Traveling Salesman Problem. *Eur. J. Opnl. Res.* (to appear).

BERTSIMAS, D., AND M. GRIGNI. 1989. Worst Case Examples for the Spacefilling Curve Heuristic for the Euclidean Traveling Salesman Problem. *O. R. Letts.* **8**, 241–244.

BERTSIMAS, D., P. CHERVI AND M. PETERSON. 1991. Computational Approaches to Stochastic Vehicle Routing Problems. *Trans. Sci.* **29**, 342–352.

BERTSIMAS, D., P. JAILLET AND A. ODONI. 1990. A Priori Optimization. *Opns. Res.* **38**, 1019–1033.

BIENSTOCK, D., J. BRAMEL AND D. SIMCHI-LEVI. 1993. A Probabilistic Analysis of Tour Partitioning Heuristics for the Capacitated Vehicle Routing Problem With Unsplit Demands. *Math O. R.* **18**, 786–802.

BRAMEL, J., AND D. SIMCHI-LEVI. 1993a. Probabilistic Analysis and Practical Algorithms for the Vehicle Routing Problem With Time Windows. *Opns. Res.* (to appear).

BRAMEL, J., AND D. SIMCHI-LEVI. 1993b. On the Effectiveness of the Set Partitioning Formulation for the Vehicle Routing Problem. *Opns. Res.* (to appear).

BRAMEL, J., AND D. SIMCHI-LEVI. 1995. A Location Based Heuristic for General Routing Problems. *Opns. Res.* **43**, 649–660.

BRAMEL, J., C. L. LI AND D. SIMCHI-LEVI. 1994. Probabilistic Analysis of the Vehicle Routing Problem With Time Windows. *Am. J. Math. and Mgmt. Sci.* **13**, 267–322.

BRAMEL, J., E. G. COFFMAN JR., P. SHOR AND D. SIMCHI-LEVI. 1992. Probabilistic Analysis of Algorithms for the Capacitated Vehicle Routing Problem With Unsplit Demands. *Opns. Res.* **40**, 1095–1106.

CHAN, A., A. FEDERGRUEN AND D. SIMCHI-LEVI. 1993. Probabilistic Analyses and Practical Algorithms for Inventory-Routing Models. *Opns. Res.* (to appear).

CHRISTOFIDES, N. 1976a. The Vehicle Routing Problem. *R.A.I.R.O. Recherche Operationelle* **10**, 55–76.

CHRISTOFIDES, N. 1976b. Worst-Case Analysis of a New Heuristic for the Travelling Salesman Problem. Report 388, Graduate School of Industrial Administration, Carnegie-Mellon University, Pittsburgh.

CHRISTOFIDES, N. 1985. Vehicle Routing. In *The Traveling Salesman Problem*, E. L. Lawler, J. K. Lenstra, A. H. G. Rinooy Kan and D. B. Shmoys (eds.). John Wiley, New York, 431–448.

COFFMAN, E. G., JR., G. S. LUEKER AND A. H. G. RINNOOY KAN. 1988. Asymptotic Methods in the Probabilistic Analysis of Sequencing and Packing Heuristics. *Mgmt. Sci.* **34**, 266–290.

CULLEN, F., J. JARVIS AND D. RATLIFF. 1981. Set Partitioning Based Heuristics for Interactive Routing. *Networks* **11**, 125–144.

DESROCHERS, M., J. DESROSIERS AND M. SOLOMON. 1992. A New Optimization Algorithm for the Vehicle Routing Problem with Time Windows. *Opns. Res.* **40**, 342–354.

DROR, M., AND P. TRUDEAU. 1986. Stochastic Vehicle Routing With Modified Saving Algorithm. *Eur. J. Opnl. Res.* **23**, 228–235.

DROR, M., G. LAPORTE AND P. TRUDEAU. 1989. Vehicle Routing with Stochastic Demands: Properties and Solution Frameworks. *Trans. Sci.* **23**, 166–176.

FEDERGRUEN, A., AND D. SIMCHI-LEVI. 1995. Analytical Analysis of Vehicle Routing and Inventory-Routing Problems. In *Handbooks in Operations Research and Management Science*, the volume on Network Routing, M. Ball, T. Magnanti, C. Monma and G. Nemhauser (eds.) 297–374.

FEDERGRUEN, A., AND G. RYZIN. 1992. Probabilistic Analysis of a Generalized Bin Packing Problem With Applications to Vehicle Routing and Scheduling Problems. Working Paper, Graduate School of Business, Columbia University, New York.

FISHER, M. L. 1995. Vehicle Routing. In *Handbooks in Operations Research and Management Science*, the volume on Network Routing, M. Ball, T. Magnanti, C. Monma and G. Nemhauser (eds.) 1–33.

FISHER, M. L., AND R. JAIKUMAR. 1981. A Generalized Assignment Heuristic for Vehicle Routing. *Networks* **11**, 109–124.

FRIESZ, T. L., J. LUQUE, R. TOBIN AND B. WIE. 1989. Dynamic Network Traffic Assignment Considered as a Continuous Time Optimal Control Problem. *Opns. Res.* **37**, 893–901.

GALLEGO, G., AND D. SIMCHI-LEVI. 1990. On the Effectiveness of Direct Shipping Strategy for the One Warehouse Multi-Retailer R-Systems. *Mgmt. Sci.* **36**, 240–243.

GENDREAU, M., G. LAPORTE AND R. SEGUIN. 1993. The Vehicle Routing Problem With Stochastic Customers and Demands. University of Montreal, Canada.

GILLETT, B. E., AND L. R. MILLER. 1974. A Heuristic Algorithm for the Vehicle Dispatch Problem. *Opns. Res.* **22**, 340–349.

HAIMOVICH, M., AND A. H. G. RINNOOY KAN. 1985. Bounds and Heuristics for Capacitated Routing Problems. *Math O. R.* **10**, 527–542.

HAIMOVICH, M., AND T. L. MAGNANTI. 1988. Extremum Properties of Hexagonal Partitioning and the Uniform Distribution in Euclidean Location. *SIAM J. Disc. Math.* **1**, 50–64.

JAILLET, P. 1988. A priori Solution of a Traveling Salesman Problem in Which a Random Subset of the Customers Are Visited. *Opns. Res.* **6**, 929–936.

JAILLET, P., AND A. ODONI. 1988. The Probabilistic Vehicle Routing Problem. In *Vehicle Routing: Methods and Studies*, B. L. Golden and A. A. Assad (eds.). North-Holland, Amsterdam.

KARMARKAR, N. 1982. Probabilistic Analysis of Some Bin-Packing Algorithms. Proc. 23rd Annual Symposium. *Foundations of Computer Science*, 107–111.

KINGMAN, J. F. C. 1976. Subadditive Processes. In *Lecture Notes in Math. 539*, Springer-Verlag, Berlin, 168–222.

LAMBERT, V., G. LAPORTE AND F. V. LOUVEAUX. 1993. Designing Collection Routes Through Bank Branches. *Comput. in Opns. Res.* (to appear).

LAPORTE, G., AND F. V. LOUVEAUX. 1990. Formulations and Bounds for the Stochastic Capacitated Vehicle Routing Problem With Uncertain Supplies. In *Economic Decision-Making: Games, Econometrics and Optimization*. J. Gabzewicz, J. F. Richard and L. Wolsey (eds.). North-Holland, Amsterdam.

LI, C. L., AND D. SIMCHI-LEVI. 1990. Analysis of Heuristics for the Multi-Depot Capacitated Vehicle Routing Problems. *ORSA J. Comput.* **2**, 64–73.

PAPASTAVROU, J., AND V. CHANDRU. 1992. On the Dynamic Traveling Repairman Problem. Technical Report, Purdue University, West Lafayette, Ind.

PIRKUL, H. 1987. Efficient Algorithms for the Capacitated Concentrator Location Problem. *Comp. in Opns. Res.* **14,** 197–208.

PLATZMAN, L. K., AND J. J. BARTHOLDI, III. 1989. Spacefilling Curves and the Planar Traveling Salesman Problem. *J. Assoc. Comp. Mach.* **36,** 719–737.

POWELL, W. 1986. A Stochastic Model for the Dynamic Vehicle Allocation Problem. *Trans. Sci.* **20,** 117–129.

POWELL, W., P. JAILLET AND A. ODONI. 1995. Stochastic and Dynamic Networks and Routing. In *Handbooks in Operations Research and Management Science*, the volume on Network Routing, M. Ball, T. Magnanti, C. Monma and G. Nemhauser (eds.).

PSARAFTIS, H. 1988. Dynamic Vehicle Routing Problems. In *Vehicle Routing: Methods and Studies*, B. Golden and A. Assad (eds.). North-Holland, Amsterdam.

PSARAFTIS, H., AND J. TSITSIKLIS. 1993. Dynamic Shortest Paths in Acyclic Networks With Markovian Arc Costs. *Opns. Res.* **41,** 91–101.

PREPARATA, F., AND M. SHAMOS. 1985. *Computational Geometry*. Springer-Verlag, New York.

RHEE, W. T., AND M. TALAGRAND. 1987. Martingale Inequalities and NP-Complete Problems. *Math O. R.* **12,** 177–181.

SIMCHI-LEVI, D. 1992. Hierarchical Planning for Probabilistic Distribution Systems in Euclidean Spaces. *Mgmt. Sci.* **38,** 198–211.

SOLOMON, M. M., AND J. DESROSIERS. 1988. Time Window Constrained Routing and Scheduling Problems. *Trans. Sci.* **22,** 1–13.

STEWART, W., AND B. GOLDEN. 1983. Stochastic Vehicle Routing: A Comprehensive Approach. *Eur. J. Opnl. Res.* **14,** 371–385.

VRANAS, P., D. BERTSIMAS AND A. ODONI. 1993. Dynamic Ground Holding Policies for a Network of Airports. *Trans. Sci.* (to appear).