

with two switch-over levels for a class of $M/G/1$ queuing systems with variable arrival and service rate." *Stochastic Processes and Their Appl.*, vol. 6, pp. 213-222, 1978.

- [11] P. Varaiya, *Notes on Optimization*. New York: Van Nostrand Reinhold, 1972.



Zvi Rosberg was born in Germany on July 25, 1947. He received the B.Sc., M.A., and Ph.D. degrees from the Hebrew University, Jerusalem, Israel, in 1971, 1974, and 1978, respectively.

From 1972 to 1978 he was a Senior System Analyst in the General Computer Bureau of the Israeli government. From 1978 to 1979 he had a Research Fellowship at the Center of Operation Research and Econometric (CORE), Catholic University of Louvain, Louvain, Belgium. From 1979 to 1980 he was a Visiting Assistant Profes-

sor at the University of Illinois, Urbana. At present, he holds the position of Lecturer in the Faculty of Computer Science, Technion—Israel Institute of Technology, Haifa.

Pravin P. Varaiya (M'68-SM'78-F'80), for a photograph and biography, see this issue, p. 655.



Jean C. Walrand (S'71-M'74) was born in Belgium in 1951. He received the Ingenieur degree in electronics from the Université de Liege, Liege, Belgium, in 1974, and the Ph.D. degree in electrical engineering from the University of California, Berkeley, in 1979.

He was with the School of Electrical Engineering, Cornell University, Ithaca, NY, from 1979 to 1981. He is now with the Department of Electrical Engineering and Computer Science, University of California, Berkeley. His research interests

are in communication networks, stochastic control, and decentralized systems.

Distributed Dynamic Programming

DIMITRI P. BERTSEKAS, SENIOR MEMBER, IEEE

Abstract—We consider distributed algorithms for solving dynamic programming problems whereby several processors participate simultaneously in the computation while maintaining coordination by information exchange via communication links. A model of asynchronous distributed computation is developed which requires very weak assumptions on the ordering of computations, the timing of information exchange, the amount of local information needed at each computation node, and the initial conditions for the algorithm. The class of problems considered is very broad and includes shortest path problems, and finite and infinite horizon stochastic optimal control problems. When specialized to a shortest path problem the algorithm reduces to the algorithm originally implemented for routing of messages in the ARPANET.

RECENT advances in microcomputer technology have intensified interest in distributed computation schemes. Aside from modular expandability, other potential advantages of such schemes are a reduction in computation time for solving a given problem due to parallelism of computation, and elimination of the need to communicate problem data available at geographically dispersed data collection points to a computation center. The first advantage is of crucial importance in real-time applications where problem solution time can be an implementation bottleneck. The second advantage manifests itself for example in applications involving communication networks where there is a natural decentralization of problem data acquisition.

Manuscript received January 13, 1981; revised June 29, 1981. Paper recommended by W. E. Schmitendorf, Chairman of the Optimal Systems Committee. This work was supported in part by the National Science Foundation under Grant NSF/ECS 79-19880.

The author is with the Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, MA 02139.

The structure of dynamic programming naturally lends itself well to distributed computation since it involves calculations that to a great extent can be carried out in parallel. In fact it is trivial to devise simple schemes taking advantage of this structure whereby the calculation involved in each iteration of the standard form of the algorithm is simply shared by several processors. Such schemes require a certain degree of synchronization in that all processors must complete their assigned portion of the computation before a new iteration can begin. As a result complex protocols for algorithm initiation and processor synchronization may be necessary, and the speed of computation is limited to that of the slowest processor. These drawbacks motivate distributed algorithms whereby computation is performed asynchronously at various nodes and independently of the progress in other nodes. Their potential advantages are simpler implementation, faster convergence to a solution and, possibly, a reduction in information exchange between computation nodes.

This paper considers an asynchronous distributed algorithm for a broad class of dynamic programming problems. This class is described in Section II. The distributed computation model is described in Section III. It is shown in Section IV that the algorithm converges to the correct solution under very weak assumptions. For some classes of problems convergence in finite time is demonstrated. These include shortest path problems for which the distributed algorithm of this paper turns out to be essentially the same as the routing algorithm originally implemented in the ARPANET in 1969 [1]. To our knowledge there is no published proof of convergence of this algorithm.

II. PROBLEM FORMULATION

We use an abstract framework of dynamic programming, first introduced in [2], [3] which includes as special cases a number of specific problems of practical interest.

Let S and C be two sets referred to as the *state space* and the *control space*, respectively. Elements of S and C are referred to as *states* and *controls* and are denoted by x and u , respectively. For each $x \in S$ we are given a subset $U(x) \subset C$ referred to as the *control constraint set at x* . Let F be the set of all extended real-valued functions $J: S \rightarrow [-\infty, \infty]$ on S . For any two functions $J_1, J_2 \in F$ we use the notation

$$J_1 \leq J_2 \quad \text{if } J_1(x) \leq J_2(x), \quad \forall x \in S, \quad (1a)$$

$$J_1 = J_2 \quad \text{if } J_1(x) = J_2(x), \quad \forall x \in S. \quad (1b)$$

Let $H: S \times C \times F \rightarrow [-\infty, \infty]$ be a mapping which is *monotone* in the sense that for all $x \in S$ and $u \in U(x)$ we have

$$H(x, u, J_1) \leq H(x, u, J_2), \quad \forall J_1, J_2 \in F \text{ with } J_1 \leq J_2. \quad (2)$$

Given a subset $\bar{F} \subset F$ the problem is to find a function $J^* \in \bar{F}$ such that

$$J^*(x) = \inf_{u \in U(x)} H(x, u, J^*), \quad \forall x \in S. \quad (3)$$

By considering the mapping $T: F \rightarrow F$ defined by

$$T(J)(x) = \inf_{u \in U(x)} H(x, u, J) \quad (4)$$

the problem is alternately stated as one of finding a fixed point of T within \bar{F} , i.e., a function $J^* \in \bar{F}$ such that

$$J^* = T(J^*). \quad (5)$$

We will assume throughout that T has a unique fixed point within \bar{F} .

We provide some examples that illustrate the broad scope of the problem formulation just given.

Example 1 (Shortest Path Problems): Let $(\mathcal{N}, \mathcal{L})$ be a directed graph where $\mathcal{N} = \{1, 2, \dots, n\}$ denotes the set of nodes and \mathcal{L} denotes the set of links. Let $N(i)$ denote the downstream neighbors of node i , i.e., the set of nodes j for which (i, j) is a link. Assume that each link (i, j) is assigned a positive scalar a_{ij} referred to as its length. Assume also that there is a directed path to node 1 from every other node. Then it is known (see [4, p. 67]) that the shortest path distances d_i^* to node 1 from all other nodes i solve uniquely the equations

$$d_i^* = \min_{j \in N(i)} \{a_{ij} + d_j^*\}, \quad \forall i \neq 1 \quad (6a)$$

$$d_1^* = 0. \quad (6b)$$

If we make the identifications

$$\begin{aligned} S &= C = \mathcal{N}, & U(x) &= N(x), \\ \bar{F} &= F, & J^*(x) &= d_x^* \end{aligned} \quad (7)$$

$$H(x, u, J) = \begin{cases} a_{xu} + J(u) & \text{if } x \neq 1 \\ 0 & \text{if } x = 1, \end{cases} \quad (8)$$

we find that the abstract problem (3) reduces to the shortest path problem.

Example 2 (Infinite Horizon Stochastic Optimal Control Problems): Let H be given by

$$H(x, u, J) = E\{g(x, u, w) + \alpha J[f(x, u, w)] | x, u\} \quad (9)$$

where we assume the following:

1) The parameter w takes values in a countable set W with given probability distribution $p(dw|x, u)$ depending on x and u , and $E\{\cdot | x, u\}$ denotes expected value with respect to this distribution.

2) The functions g and f map $S \times C \times W$ into $[-\infty, \infty]$ and S , respectively.

3) The scalar α is positive.

Because the set W is assumed countable the expected value in (9) is well defined for all $J \in F$ in terms of infinite summation provided we use the convention $+\infty - \infty = +\infty$ (see [3, p. 31]). It is possible to consider a more general probabilistic structure for W (see [3]) at the expense of complicating the presentation but this does not seem worthwhile in view of the computational purposes of the paper.

It is shown in [3] that with this definition of H the abstract problem (3) reduces under more specific assumptions to various types of standard stochastic optimal control problems. Thus if g is uniformly bounded above and below by some scalars and $0 < \alpha < 1$ the problem is equivalent to the standard infinite horizon discounted stochastic optimal control problem with bounded cost per stage (see [5, sect. 6.1–6.3]). Under these circumstances the mapping T of (4) has a unique fixed point J^* in the class of all bounded real-valued functions on S and J^* is the optimal value function of the corresponding stochastic optimal control problem.

If we assume that $0 \leq g(x, u, w)$ or $g(x, u, w) \leq 0$ for all $(x, u, w) \in S \times C \times W$, then we obtain stochastic optimal control problems of the type discussed extensively, for example, in [5, sect. 6.4–6.6, 7.1–7.4], and [3, ch. 5]. If J^* is the optimal value function for such a problem, then J^* is the unique fixed point of T over all functions $J \in F$ such that $0 \leq J \leq J^*$ if $0 \leq g(x, u, w)$ for all (x, u, w) , or $J^* \leq J \leq 0$ if $g(x, u, w) \leq 0$ for all (x, u, w) (see [5, p. 256]).

Example 3 (Finite Horizon Stochastic Optimal Control Problems): Let $S, C, U(x), W, p(dw|x, u), g$ and f be as in Example 2 and consider the set of equations

$$J_N(x_N) = 0, \quad x_N \in S \quad (10a)$$

$$\begin{aligned} J_k(x_k) &= \inf_{u_k \in U(x_k)} E\{g(x_k, u_k, w_k) \\ &\quad + J_{k+1}[f(x_k, u_k, w_k)] | x_k, u_k\}, \\ &\quad k = 0, 1, \dots, N-1, \quad x_k \in S, \end{aligned} \quad (10b)$$

where N is a positive integer. These are the usual dynamic programming equations associated with finite horizon sto-

chastic optimal control problems with zero terminal cost and stationary cost per stage and system function. It is possible to write these equations in the form (3) by defining a new state space consisting of an $(N + 1)$ -fold Cartesian product of S with itself, writing $J^* = (J_0, J_1, \dots, J_N)$, and appropriately defining H on the basis of (10). In fact this is a standard procedure for converting a finite horizon problem to an infinite horizon problem (see [5, p. 325]). This reformulation can also be trivially generalized to finite horizon problems involving a nonzero terminal cost and a nonstationary system and cost per stage.

III. A MODEL FOR DISTRIBUTED DYNAMIC PROGRAMMING

Our algorithm can be described in terms of a collection of n computation centers referred to as nodes and denoted $1, 2, \dots, n$. The state space S is partitioned into n disjoint sets denoted S_1, \dots, S_n . Each node i is assigned the responsibility of computing the values of the solution function J^* at all states x in the corresponding set S_i . A node j is said to be a neighbor of node i if $j \neq i$ and there exist a state $x_i \in S_i$ and two functions $J_1, J_2 \in \bar{F}$ such that

$$J_1(x) = J_2(x), \quad \forall x \notin S_j \tag{11a}$$

$$T(J_1)(x_i) \neq T(J_2)(x_i). \tag{11b}$$

The set of all neighbors of i is denoted $N(i)$. Intuitively j is not a neighbor of i if, for every $J \in \bar{F}$, the values of J on S_j do not influence the values of $T(J)$ on S_i . As a result, for any $J \in \bar{F}$, in order for node i to be able to compute $T(J)$ on S_i it is only necessary to know the values of J on the sets $S_j, j \in N(i)$, and, possibly, on the set S_i .

At each time instant, node i can be in one of three possible states—compute, transmit, or idle. In the compute state node i computes a new estimate of the values of the solution function J^* for all states $x \in S_i$. In the transmit state node i communicates the estimate obtained from the latest compute phase to one or more nodes m for which $i \in N(m)$. In the idle state node i does nothing related to the solution of the problem. Each node i also has one buffer per neighbor $j \in N(i)$ denoted B_{ij} where it stores the latest transmission from j , as well as a buffer B_{ii} where it stores its own estimate of values of the solution function J^* for all states $x \in S_i$. It is assumed that a node can receive a transmission from neighbors simultaneously with computing or transmitting.

We assume that computation and transmission for each node takes place in uninterrupted time intervals $[t_1, t_2]$ with $t_1 < t_2$, but do not exclude the possibility that a node may be simultaneously transmitting to more than one node nor do we assume that the transmission intervals to these nodes have the same origin and/or termination. We also make no assumptions on the length, timing, and sequencing of computation and transmission intervals other than the following.

Assumption 1: There exists a positive scalar P such that, for every node i , every time interval of length P contains at least one computation interval for i and at least one transmission interval from i to each node m with $i \in N(m)$.

The contents of each buffer B_{ij} where $j = i$ or $j \in N(i)$ at time t are denoted J_{ij}^t . Thus J_{ij}^t is, for every t , a function from S_j into $[-\infty, \infty]$ and may be viewed as the estimate by node i of the restriction of the solution function J^* on S_j available at time t . The rules according to which the functions J_{ij}^t are updated are as follows:

1) If $[t_1, t_2]$ is a transmission interval for node j to node i with $i \in N(j)$, the contents $J_{jj}^{t_1}$ of the buffer B_{jj} at time t_1 are transmitted and entered in the buffer B_{ij} at time t_2 , i.e.,

$$J_{ij}^{t_2} = J_{jj}^{t_1}. \tag{12}$$

2) If $[t_1, t_2]$ is a computation interval for node i , the contents of buffer B_{ii} at time t_2 are replaced by the restriction of the function $T(J_{ii}^{t_1})$ on S_i where, for all t , J_i^t is defined by

$$J_i^t(x) = \begin{cases} J_{ii}^t(x), & \text{if } x \in S_i, \\ J_{ij}^t(x), & \text{if } x \in S_j \text{ and } j \in N(i), \\ 0, & \text{otherwise.} \end{cases} \tag{13}$$

In other words we have

$$J_{ii}^{t_2}(x) = T(J_i^{t_1})(x) = \inf_{u \in U(x)} H(x, u, J_i^{t_1}), \quad \forall x \in S_i. \tag{14}$$

3) The contents of a buffer B_{ii} can change only at the end of a computation interval for node i . The contents of a buffer $B_{ij}, j \in N(i)$, can change only at the end of a transmission interval from j to i .

Note that by definition of the neighbor set $N(i)$, the value $T(J_i^t)(x)$ for $x \in S_i$ does not depend on the values of J_i^t at states $x \in S_m$ with $m \neq i$, and $m \notin N(i)$. We have assigned arbitrarily the default value zero to these states in (13). Our objective is to show that for all $i = 1, \dots, n$

$$\lim_{t \rightarrow \infty} J_{ij}^t(x) = J^*(x), \quad \forall x \in S_j, j = i \text{ or } j \in N(i).$$

It is clear that an assumption such as Assumption 1 is necessary in order for such a result to hold. Since iteration (14) is of the dynamic programming type it is also clear that some restrictions must be placed on the mapping H that guarantee convergence of the algorithm under the usual circumstances where the algorithm is carried out in a centralized synchronous manner (i.e., when there is only one computation node). The following assumption places somewhat indirect restrictions on H but simplifies the convergence analysis.

Assumption 2: There exist two functions \underline{J} and \bar{J} in F such that the set of all functions $J \in F$ with $\underline{J} \leq J \leq \bar{J}$ belongs to \bar{F} and furthermore

$$J \geq T(\bar{J}), \quad T(\underline{J}) \geq \underline{J} \tag{15}$$

$$\lim_{k \rightarrow \infty} T^k(\bar{J})(x) = J^*(x), \quad \forall x \in S \quad (16a)$$

$$\lim_{k \rightarrow \infty} T^k(\underline{J})(x) = J^*(x), \quad \forall x \in S \quad (16b)$$

where T^k denotes composition of the mapping T with itself k times.

Note that in view of the monotonicity of H [cf. (2)] and the fact $J^* = T(J^*)$, Assumption 2 implies

$$\begin{aligned} \bar{J} &\geq T(\bar{J}) \geq T^2(\bar{J}) \geq \dots \geq J^* \geq \dots \\ &\geq T^2(\underline{J}) \geq T(\underline{J}) \geq \underline{J}. \end{aligned}$$

Furthermore if $J \in F$ satisfies $\underline{J} \leq J \leq \bar{J}$ then

$$\lim_{k \rightarrow \infty} T^k(J)(x) = J^*(x) \quad \text{for all } x \in S.$$

Assumption 2 leaves open the question of how to find suitable functions \underline{J} and \bar{J} . On the other hand for most problems of interest the choices are clear. In particular we have the following:

1) For shortest path problems (Example 1) it is straightforward to verify that the choices

$$\underline{J}(i) = 0, \quad \forall i = 1, \dots, n \quad (17a)$$

$$\bar{J}(i) = \begin{cases} \infty & \text{if } i \neq 1 \\ 0 & \text{if } i = 1 \end{cases} \quad (17b)$$

satisfy Assumption 2.

2) For finite horizon stochastic optimal control problems (Example 3) for which the function g in (10) is uniformly bounded below it can be easily verified that the functions $\bar{J} = (\bar{J}_0, \bar{J}_1, \dots, \bar{J}_N)$ and $\underline{J} = (\underline{J}_0, \underline{J}_1, \dots, \underline{J}_N)$ where for all k and x

$$\bar{J}_k(x) = \infty, \quad \underline{J}_k(x) = -\infty \quad (18)$$

satisfy Assumption 2.

3) For discounted infinite horizon stochastic optimal control problems with bounded cost per stage (Example 2 with $\alpha \in (0, 1)$ and g uniformly bounded above and below) it is easily shown that every pair of functions \bar{J}, \underline{J} of the form

$$\bar{J}(x) = \bar{\beta}, \quad \underline{J}(x) = \underline{\beta}, \quad \forall x \in S$$

where

$$\begin{aligned} \frac{1}{1-\alpha} \sup_{x \in S} \inf_{u \in U(x)} E\{g(x, u, w)\} &\leq \bar{\beta} < \infty \\ -\infty < \underline{\beta} &\leq \frac{1}{1-\alpha} \inf_{x \in S} \inf_{u \in U(x)} E\{g(x, u, w)\} \end{aligned}$$

satisfy Assumption 2.

4) For infinite horizon stochastic optimal control problems with nonpositive cost per stage (Example 2 with $g \leq 0$) it can be shown that the functions \bar{J}, \underline{J} with

$$\underline{J}(x) = J^*(x), \quad \bar{J}(x) = 0, \quad \forall x \in S$$

satisfy Assumption 2 ([5, pp. 261, 298]). If the cost per stage is nonnegative (Example 2 with $g \geq 0$), then, under a

mild assumption (which is satisfied in particular if $U(x)$ is a finite set for each x), it can be shown that the choices \bar{J}, \underline{J} with

$$\underline{J}(x) = 0, \quad \bar{J}(x) = J^*(x), \quad \forall x \in S$$

satisfy Assumption 2 ([5, pp. 263, 298]). The choice of \underline{J} and \bar{J} can be further sharpened and simplified under more specific assumptions on problem structure but we will not pursue this matter further.

Our convergence result will be shown under the assumption that the contents J_{ij}^0 of the buffers B_{ij} at the initial time $t = 0$ satisfy

$$\underline{J}(x) \leq J_{ij}^0(x) \leq \bar{J}(x), \quad \forall x \in S_j. \quad (19)$$

The broad range of initial conditions allowed by (19) eliminates the need to reset the contents of the buffers in an environment where it is necessary to execute periodically the algorithm with slightly different problem data as for example in routing algorithms for communication networks. This is particularly true for cases 1)–3) above where condition (19) implies that the initial buffer contents can be essentially arbitrary.

IV. CONVERGENCE ANALYSIS

Our main result is the following proposition.

Proposition 1: Let Assumptions 1 and 2 hold and assume that for all $i = 1, \dots, n$

$$\underline{J}(x) \leq J_{ij}^0(x) \leq \bar{J}(x), \quad \forall x \in S_j, \quad j = i \text{ or } j \in N(i). \quad (20)$$

Then for all $i = 1, \dots, n$

$$\lim_{t \rightarrow \infty} J_{ij}^t(x) = J^*(x), \quad \forall x \in S_j, \quad j = i \text{ or } j \in N(i).$$

Proof: Since the contents of a buffer can change only at the end of a computation or transmission interval at a node we can analyze convergence in terms of a discrete time process. We focus attention at the sequence of times $\{t_k\}$ with $0 < t_1 < t_2 < \dots$ where each t_k is the end of a computation interval for one or more nodes.

Let us consider for all $t \geq 0$

$$\bar{J}_{ij}^t: S_j \rightarrow [-\infty, \infty], \quad \underline{J}_{ij}^t: S_j \rightarrow [-\infty, \infty],$$

where for each $x \in S_j$, the value $\bar{J}_{ij}^t(x)[\underline{J}_{ij}^t(x)]$ represents the contents of buffer B_{ij} at time t if the algorithm were executed with the same timing and order of computation and transmission intervals but with initial condition $\bar{J}(x)[\underline{J}(x)]$ instead of $J_{ij}^0(x)$ for each buffer B_{ij} and $x \in S_j$. The monotonicity assumption (2) and the definition of the algorithm [cf. (13), (14)] clearly imply that for all t

$$\begin{aligned} \underline{J}_{ij}^t(x) \leq J_{ij}^t(x) \leq \bar{J}_{ij}^t(x), \quad \forall x \in S_j, \quad i = 1, \dots, n, \\ j = i \text{ or } j \in N(i). \end{aligned} \quad (21)$$

It will thus suffice to show that

$$\lim_{t \rightarrow \infty} \underline{J}_{ij}^t(x) = J^*(x), \quad \forall x \in S_j, \quad i=1, \dots, n, \quad j \in N(i) \quad (22)$$

$$\lim_{t \rightarrow \infty} \bar{J}_{ij}^t(x) = J^*(x), \quad \forall x \in S_j, \quad i=1, \dots, n, \quad j \in N(i). \quad (23)$$

In view of the fact $\bar{J} \geq T(\bar{J})$ we have clearly

$$\bar{J}(x) = \bar{J}_{ij}^0(x) \geq \bar{J}_{ij}^1(x), \quad \forall x \in S_j, \quad i=1, \dots, n, \quad j \in N(i) \quad (24)$$

with potential strict inequality only for nodes j for which t_1 was the end of a computation interval. For $t \in [t_1, t_2)$ the content of B_{ij} is either \bar{J}_{ij}^0 or \bar{J}_{ij}^1 so from (24) we must have

$$\begin{aligned} \bar{J}_{ij}^t(x) &\geq \bar{J}_{ij}^1(x), & \forall x \in S_j, \quad i=1, \dots, n, \\ & & j \in N(i), \quad t \in [t_1, t_2) \\ \bar{J}_{ij}^t(x) &\geq \bar{J}_{ij}^0(x), & \forall x \in S_j, \quad i=1, \dots, n, \\ & & j \in N(i), \quad t \in [t_1, t_2). \end{aligned} \quad (25)$$

The last relation can also be written as

$$\bar{J}_{jm}^t(x) \geq \bar{J}_{jm}^t(x), \quad \forall x \in S_m, \quad j=1, \dots, n, \quad m \in N(j), \quad t \in [t_1, t_2). \quad (26)$$

In view of the monotonicity of H it follows from (26) and (14) that

$$\bar{J}_{jj}^t(x) \geq \bar{J}_{jj}^2(x), \quad \forall x \in S_j \quad (27)$$

with potential strict inequality only for nodes j for which t_2 was the end of a computation interval. Combining (25) and (27) we obtain

$$\bar{J}_{ij}^t(x) \geq \bar{J}_{ij}^2(x), \quad \forall x \in S_j, \quad i=1, \dots, n, \quad j \in N(i)$$

with potential strict inequality only for nodes j for which either t_1 or t_2 was the end of a computation interval. The preceding argument can be repeated to show that for all $k, i=1, \dots, n$, and $j \in N(i)$ we have

$$\bar{J}_{ij}^t(x) \geq \bar{J}_{ij}^t(x) \geq \bar{J}_{ij}^{t-k}(x) \geq \bar{J}_{ij}^{t-k-1}(x), \quad \forall x \in S_j, \quad t \in [t_k, t_{k+1}). \quad (28)$$

Let k_1 be the first integer for which $2P \leq t_{k_1}$ where P is as in Assumption 1. Then each node must have completed at least one computation phase in the interval $[0, P]$ and at least one transmission phase to all nodes in the interval $[P, 2P]$. It is easily seen that this together with (28), the monotonicity of H , and the definition of the algorithm implies that for all $t \in [t_{k_1}, t_{k_1+1})$

$$T(\bar{J})(x) \geq \bar{J}_{ij}^t(x) \geq \bar{J}_{ij}^{t-k_1}(x), \quad \forall x \in S_j, \quad i=1, \dots, n, \quad j \in N(i).$$

This argument can be repeated and shows that if $m(k)$ is the largest integer m such that $2mP \leq t_k$ then for all $t \in [t_k, t_{k+1})$

$$T^{m(k)}(\bar{J})(x) \geq \bar{J}_{ij}^t(x) \geq \bar{J}_{ij}^{t-k}(x), \quad \forall x \in S_j, \quad i=1, \dots, n, \quad j \in N(i). \quad (29)$$

Similarly we obtain for all $t \in [t_k, t_{k+1})$

$$\underline{J}_{ij}^t(x) \geq \underline{J}_{ij}^t(x) \geq T^{m(k)}(\underline{J})(x), \quad \forall x \in S_j, \quad i=1, \dots, n, \quad j \in N(i). \quad (30)$$

By combining (21), (29), and (30) and using Assumption 2 we obtain (22) and (23), and the proof of the proposition is complete. Q.E.D.

An alternative proof of Proposition 1 may be found in [8]. Note that (21), (29), and (30) provide useful estimates of rate of convergence. In fact by using these relations it is possible to show that in some special cases convergence is attained in finite time.

Proposition 2: In the cases of Examples 1 and 3 with the choices of \bar{J}, \underline{J} given by (17) and (18), respectively and J_{ij}^0 satisfying (20), there exists a time $\bar{t} > 0$ such that for all $i=1, \dots, n$ and $t \geq \bar{t}$ there holds

$$J_{ij}^t(x) = J^*(x), \quad \forall x \in S_j, \quad j=i \text{ or } j \in N(i).$$

Proof: For Example 3 it is easily seen that there holds

$$T^k(\bar{J})(x) = T^k(\underline{J})(x) = J^*(x), \quad \forall x \in S, \quad k \geq N+1.$$

The proof follows from (21), (29), and (30). For Example 1 it is easily seen that

$$T^k(\bar{J})(i) = J^*(i), \quad \forall i=2, \dots, n, \quad k \geq n.$$

Also for each i , $T^k(\underline{J})(i)$ represents the length of a path starting from i with k links, and each link has positive length. Therefore there exists a \bar{k} such that $T^{\bar{k}}(\underline{J})(i)$ represents length of a path from i to node 1, for otherwise the paths corresponding to $T^k(\underline{J})(i)$ would cycle indefinitely without reaching node 1 and we would have $T^k(\underline{J})(i) \rightarrow \infty$. Since $T^k(\underline{J})(i) \leq J^*(i)$ and $J^*(i)$ is the shortest distance from i to 1 we obtain

$$T^k(\underline{J})(i) = J^*(x), \quad \forall i=2, \dots, n, \quad k \geq \bar{k}.$$

The result again follows from (21), (29), and (30). Q.E.D.

It is possible to construct examples showing that in the case of the shortest path problem the number of iterations needed for finite convergence of $T^k(\underline{J})$ depends on the link lengths in a manner which makes the overall algorithm nonpolynomial.

In many problems of interest the main objective of the algorithm is to obtain a minimizing control law μ^* , i.e., a function $\mu^*: S \rightarrow C$ with $\mu^*(x) \in U(x)$ for all $x \in S$ such that

$$H[x, \mu^*(x), J^*] = \min_{u \in U(x)} H(x, u, J^*), \quad \forall x \in S. \quad (31)$$

It is thus of interest to investigate the question of whether control laws $\mu^t: S \rightarrow C$ satisfying

$$\mu^t(x) \in U(x), \quad \forall x \in S \quad (32)$$

and

$$H[x, \mu^t(x), J_i^t] = \min_{u \in U(x)} H(x, u, J_i^t), \quad \forall x \in S_i, i = 1, \dots, n \quad (33)$$

where J_i^t is given for all t by (13), converge in some sense to a control law μ^* satisfying (31). The following proposition shows that convergence is attained in finite time if the sets $U(x)$ are finite and H has a continuity property which is satisfied for most problems of practical interest. A related convergence result can be shown assuming the sets $U(x)$ are compact (cf. [5, Proposition 5.11]).

Proposition 3: Let the assumptions of Proposition 1 hold. Assume also that for every $x \in S$, $u \in U(x)$ and sequence $\{J^k\} \subset \bar{F}$ for which $\lim_{k \rightarrow \infty} J^k(x) = J^*(x)$ for all $x \in S$ we have

$$\lim_{k \rightarrow \infty} H(x, u, J^k) = H(x, u, J^*). \quad (34)$$

Then for each state $x \in S$ for which $U(x)$ is a finite set there exists $\bar{t}_x > 0$ such that for all $t \geq \bar{t}_x$ if $\mu^t(x)$ satisfies (32), (33) then

$$H[x, \mu^t(x), J^*] = \min_{u \in U(x)} H(x, u, J^*).$$

Proof: Assume the contrary, i.e., that there exists a state $x \in S_i$ for which $U(x)$ is finite, and an increasing sequence $\{t_k\}$ with $t_k \rightarrow \infty$ such that

$$H[x, \mu^{t_k}(x), J_i^{t_k}] = \min_{u \in U(x)} H(x, u, J_i^{t_k}), \quad \forall k = 1, 2, \dots \quad (35)$$

and

$$H[x, \mu^{t_k}(x), J^*] > \min_{u \in U(x)} H(x, u, J^*) = J^*(x), \quad \forall k = 1, 2, \dots \quad (36)$$

Since $U(x)$ is finite, there exists a $\bar{u} \in U(x)$ such that $\mu^{t_k}(x) = \bar{u}$ for an infinite subset of indexes K . From Proposition 1 we have that $J_i^{t_k}$ converges pointwise to J^* on the set $S_i \cup_{j \in N(i)} S_j$. Using the definition of the neighbor set $N(i)$, (34) and (36), it follows that

$$\lim_{k \rightarrow \infty} H(x, \bar{u}, J_i^{t_k}) = H(x, \bar{u}, J^*) > J^*(x).$$

On the other hand from (35) and Proposition 1 we have

$$\lim_{\substack{k \rightarrow \infty \\ k \in K}} H(x, \bar{u}, J_i^{t_k}) = J^*(x)$$

which contradicts the previous relation. Q.E.D.

V. DISCUSSION AND CONCLUSIONS

The analysis of this paper shows that natural distributed dynamic programming schemes converge to the correct solution under very weak assumptions on the problem structure, and the timing and ordering of computation and

internode communication. The restrictions on the initial conditions are also very weak. This means that, for problems that are being solved continuously in real time, it is not necessary to reset the initial conditions and resynchronize the algorithm each time the problem data changes. As a result the potential for tracking slow variations in optimal control laws is improved, and algorithmic implementation is greatly simplified.

The crucial assumption in the analysis of this paper is the monotonicity property of the mapping H [cf. (2)]. Indeed this property is largely responsible for most of the basic results in dynamic programming (see [3]). On the other hand the mapping H of many dynamic programming models possesses a contraction property which is sufficient to guarantee the validity of the distributed algorithm of this paper even in the absence of the monotonicity assumption (2). To be more specific let \bar{F} be the set of all uniformly bounded real valued functions on S equipped with the sup-norm

$$\|J\| = \sup_{x \in S} |J(x)|, \quad \forall J \in \bar{F}. \quad (37)$$

Assume that, in place of the monotonicity assumption (2), H has the following properties:

$$T(J) \in \bar{F}, \quad \forall J \in \bar{F} \quad (38)$$

$$\|T(J) - T(J')\| \leq \rho \|J - J'\|, \quad \forall J, J' \in \bar{F} \quad (39)$$

where ρ is a scalar with $0 < \rho < 1$. Then T has a unique fixed point J^* in \bar{F} and it is possible to show that the conclusion of Proposition 1 holds provided Assumption 1 is in effect and the initial buffer contents J_{ij}^0 are uniformly bounded functions on the corresponding sets S_j . It is not necessary to assume Assumption 2 for this result. The proof is very similar to the one of Proposition 1 and utilizes the contraction assumption (39) to show that the sequences

$$\sup_{i,j} \sup_{x \in S_j} \{|J_{ij}^t(x) - J^*(x)|\}$$

decrease monotonically to zero as $t \rightarrow \infty$. Note that since the value of H need not depend on u , this result shows the validity of our algorithm applied to an arbitrary fixed point problem of the form $J = T(J)$ for which the mapping $T: \bar{F} \rightarrow \bar{F}$ satisfies (38) and (39). The result just stated can also be proved by applying a general convergence theorem for distributed fixed point algorithms given in [8].

The use of the sup-norm (37) is essential for the validity of the result described above. Indeed for the important class of Markovian decision problems involving a finite state space and minimization of average cost per stage (see Howard [6]), a distributed asynchronous version of the usual dynamic programming algorithm due to White [7] (see [5, sect. 8.2]) may fail to converge to the correct solution. This is illustrated in the following example constructed by J. Tsitsiklis. In this example the basic mapping H does not satisfy the monotonicity assumption (2), and the corresponding mapping T is not a contraction mapping with respect to the sup-norm (37). It is a contraction mapping with respect to a different norm.

Example (J. Tsitsiklis): Let the state space consist of two states $S = \{0, 1\}$, and the control space consist of a single control $C = \{0\}$. Consider the Markov chain with state space S for which at each stage if the state is 0 a cost g_0 is incurred and a transition to state 1 occurs with probability $p_0 \in (0, 1)$. If the state is 1 a cost g_1 is incurred and a transition to state 0 occurs with probability $p_1 \in (0, 1)$. Consider the mapping $T: R^2 \rightarrow R^2$ defined by

$$T(J)(0) = g_0 + p_0[J(1) - J(0)] \quad (40)$$

$$T(J)(1) = g_1 + p_1[J(1) - J(0)]. \quad (41)$$

Because there is only one control available at each stage the definition of T does not involve a minimization as in (4). It is clear however that T arises from a mapping H of the form considered in this paper except that this mapping does not satisfy the monotonicity condition (2).

Now by applying a well-known result (e.g., [5, p. 345]) we have that T has a unique fixed point $J^* = (J^*(0), J^*(1))$, and $J^*(0)$ is the average gain of the process. Furthermore, the standard dynamic programming algorithm which consists of the successive generation of $T(J), T^2(J), \dots$, starting from an arbitrary initial $J \in R^2$ converges to J^* . Indeed T is an affine mapping involving the matrix

$$\begin{bmatrix} -p_0 & p_0 \\ -p_1 & p_1 \end{bmatrix}.$$

It can be easily seen that the eigenvalues of this matrix lie strictly within the unit circle and as a result T is a contraction mapping with respect to some norm on R^2 . However, T is not a contraction mapping with respect to the sup-norm.

Now consider the distributed algorithm of Section III with two computation nodes 0 and 1 corresponding to the two states. Consider a sequence of events whereby node 0 does many iterations before transmitting at time t_1 the final value $J^{t_1}(0)$ to node 1, while in the meantime node 1 is idle. Then node 1 does many iterations before transmitting at time t_2 the final value $J^{t_2}(1)$ while node 0 is idle, and the process is repeated. If $J^0(1)$ is the estimate of $J^*(1)$ available at nodes 0 and 1 at time 0, we have using (40) and (41)

$$J^{t_1}(0) \approx \frac{1}{1+p_0} [g_0 + p_0 J^0(1)]$$

$$J^{t_2}(1) \approx \frac{1}{1-p_1} [g_1 - p_1 J^{t_1}(0)].$$

By eliminating $J^{t_1}(0)$ in the relations above, we obtain

$$J^{t_2}(1) \approx \frac{1}{1-p_1} \left(g_1 - \frac{p_1}{1+p_0} g_0 \right) - \frac{p_0 p_1}{(1+p_0)(1-p_1)} J^0(1).$$

Thus the estimate of $J^*(1)$ is updated approximately according to the equation

$$J(1) \leftarrow \frac{1}{1-p_1} \left(g_1 - \frac{p_1}{1+p_0} g_0 \right) - \frac{p_0 p_1}{(1+p_0)(1-p_1)} J(1)$$

and it follows that, if p_1 is sufficiently close to unity, then $J(1)$ will oscillate between positive and negative numbers of increasingly large magnitude. This shows that the natural distributed version of the dynamic programming algorithm for average cost Markovian decision problems is not guaranteed to converge to the correct solution when the sequencing of computation and internode communication is arbitrary.

REFERENCES

- [1] J. McQuillan, G. Falk, and I. Richer, "A review of the development and performance of the ARPANET routing algorithm," *IEEE Trans. Commun.*, vol. COM-26, pp. 1802-1811, 1978.
- [2] D. P. Bertsekas, "Monotone mappings with application in dynamic programming," *SIAM J. Contr. Optimiz.*, vol. 15, pp. 438-464, 1977.
- [3] D. P. Bertsekas and S. E. Shreve, *Stochastic Optimal Control: The Discrete-Time Case*. New York: Academic, 1978.
- [4] E. L. Lawler, *Combinatorial Optimization: Networks and Matroids*. New York: Holt, Rinehart, and Winston, 1976.
- [5] D. P. Bertsekas, *Dynamic Programming and Stochastic Control*. New York: Academic, 1976.
- [6] R. Howard, *Dynamic Programming and Markov Processes*. Cambridge, MA: M.I.T., 1960.
- [7] D. J. White, "Dynamic programming, Markov chains, and the method of successive approximations," *J. Math. Anal. Appl.*, vol. 6, pp. 373-376, 1963.
- [8] D. P. Bertsekas, "Distributed computation of fixed points," Massachusetts Inst. Technol., Cambridge, LIDS Rep. P-1135, Aug. 1981.



Dimitri P. Bertsekas (S'70-SM'77) received the Diploma of Mechanical and Electrical Engineering from the National Technical University of Athens, Athens, Greece, in 1965, the M.S.E.E. degree from George Washington University, Washington, DC, in 1969, and the Ph.D. degree in electrical engineering from Massachusetts Institute of Technology, Cambridge, in 1971.

He has served on the faculty of Stanford University and the University of Illinois. He is currently Professor of Electrical Engineering and Computer Science at M.I.T. His research interests are in optimization theory and algorithmic aspects of communication networks. He is the author of *Dynamic Programming and Stochastic Control* (New York: Academic, 1976) and *Constrained Optimization and Lagrange Multiplier Methods* (New York: Academic, 1982) and coauthor of *Stochastic Optimal Control: The Discrete-Time Case* (New York: Academic, 1978). He is also an Editor of the *SIAM Journal on Control and Optimization*, and the *Journal of Optimization Theory and Applications*.