# DYNAMIC SYSTEM-OPTIMAL TRAFFIC ASSIGNMENT USING A STATE SPACE MODEL[1]

STÉPHANE LAFORTUNE and RAJA SENGUPTA
Department of Electrical Engineering and Computer Science

and

DAVID E. KAUFMAN and ROBERT L. SMITH
Department of Industrial and Operations Engineering, University of Michigan,
Ann Arbor, MI 48109, U.S.A.

**Abstract** — We propose a new mathematical formulation for the problem of optimal traffic assignment in dynamic networks with multiple origins and destinations. This problem is motivated by route guidance issues that arise in an Intelligent Vehicle–Highway Systems (IVHS) environment. We assume that the network is subject to known time-varying demands for travel between its origins and destinations during a given time horizon. The objective is to assign the vehicles to links over time so as to minimize the total travel time experienced by all the vehicles using the network. We model the traffic network over the time horizon as a discrete-time dynamical system. The system state at each time instant is defined in a way that, without loss of optimality, avoids complete microscopic detail by grouping vehicles into platoons irrespective of origin node and time of entry to network. Moreover, the formulation contains no explicit path enumeration. The state transition function can model link travel times by either impedance functions, link outflow functions, or by a combination of both. Two versions (with different boundary conditions) of the problem of optimal traffic assignment are studied in the context of this model. These optimization problems are optimal control problems for nonlinear discrete-time dynamical systems, and thus they are amenable to algorithmic solutions based on dynamic programming. The computational challenges associated with the exact solution of these problems are discussed and some heuristics are proposed.

## 1. INTRODUCTION

We consider the problem of dynamic traffic assignment in networks with multiple trip origins and destinations. Our approach is as follows. The traffic network, which may include both freeway corridors and surface streets, is modeled as a directed graph. The sets of origins and destinations are subsets of the set of vertices of the graph. The edges of the graph are links in the network. Some information is available about these links, in the form of impedance functions, which express link travel times in terms of the number of vehicles on the links, or link outflow functions, which constrain the departure or exit rate of vehicles from a link in terms of the number of vehicles on this link. We assume that the network is subject to known time-varying demands from vehicles for travel between their origins and destinations during a given finite time horizon (e.g. a period of a few hours). The objective is to assign the vehicles to links over time in order to minimize the total travel time experienced by all the vehicles using the network. Thus the resulting assignment will be system-optimal. We are dealing with a dynamic, as opposed to a static, problem because the demand is dynamic and optimal routes assigned to vehicles from their origins to their destinations depend on the entire set of demands over the whole time horizon considered.

Our motivation for studying this problem comes from route guidance issues that arise in an Intelligent Vehicle–Highway Systems (IVHS) environment (see Saxton, 1991). In IVHS, it is desired to perform anticipatory route guidance, i.e. to route the vehicles on the network on the basis of the future travel times they will experience on the links they will be traveling. However, these future travel times depend on the routing decisions made for other vehicles traveling on the network, and, thus, they have to be forecasted using a combination of historical and real-time information (see, e.g. Kaufman, *et al.*,

TR(B) 27:6-D

1990; Kaufman and Smith, 1993; and Wunderlich, 1990). Our objective is not to address directly this real-time forecasting/assignment problem where the demand is not known a priori, and the decisions may be affected by the occurrence of incidents and other unpredictable events. By making the simplifying assumption that the demand, although dynamic, is known beforehand, and by considering a fixed time horizon (with no incidents occurring), we obtain an optimization problem that may be solved completely, at least in principle. It is our thesis that the solution of the dynamic traffic assignment problem considered in this paper, as well as the properties of this solution, will provide considerable insight into the problem of real-time anticipatory route guidance in IVHS. This thesis is the primary motivation of the work that follows.

We model the traffic network over the time horizon as a discrete-time dynamical system. At each time instant, the system state consists of all platoons, each of which represents all vehicles on a certain link with the same destination and the same earliest possible time of departing the link. This state definition avoids complete microscopic detail by grouping vehicles into platoons irrespective of origin node and time of entry to network, yet it conforms to the requirement that the state should summarize all relevant past behavior so as to contain sufficient information for the determination of the future behavior of the system. Also, since this formulation is based on links, it contains no path enumeration. We provide a general form for the state transition function giving the possible states at time $t + 1$ as a function of the state at time $t$ and of the feasible assignment or routing decisions for platoons that exit a link or join the network in the time interval $(t, t + 1]$. Specific forms of the state transition function can model link travel times by either impedance functions or link outflow functions. Further, the two can be combined in a way that represents interaction between links due to recurring congestion and capacity reductions, in addition to single link impedances.

Overall, the optimization problem for traffic assignment becomes one of optimal control for a nonlinear integer-valued discrete-time dynamical system. The number of control variables (i.e. the number of routing decisions that have to be made) is bounded by the product of (a) the number of links, (b) the number of destinations, and (c) the time horizon considered. As formulated, the optimization problem is amenable to an algorithmic solution based on dynamic programming. Hence, Dijkstra's algorithm can be employed to determine the system-optimal routing decisions in the context of a forward dynamic programming search over the state space. Other search techniques (optimal or heuristic) could also be employed. In fact, our model has been developed in the hope of providing a compact yet detailed mathematical representation, to which methods in areas such as combinatorial optimization or mixed integer programming may be applied.

Our presentation is organized as follows. We compare our approach with related work on dynamic traffic assignment in Section 2. The dynamical system model is presented in Section 3, and two versions (with different boundary conditions) of the resulting optimal control problem are formulated in Section 4. This section also presents forward and backward dynamic programming recursions for solving the problem. The forward dynamic programming recursion is applied in Section 5 to a simple triangle network. Section 6 discusses important extensions of the model of Section 3 concerning background traffic, blocking controls, link outflow functions and the combination of impedance and link outflow functions. Issues of computational complexity and search heuristics are discussed in Section 7, while Section 8 concludes the paper.

## 2. COMPARISON WITH OTHER WORK

Merchant and Nemhauser (1978) formulated a mathematical program for system-optimal dynamic traffic assignment in a network with multiple origins and a single destination. They assumed that all links are uncapacitated and that in each time interval the number of vehicles departing a link is a function (nonlinear in general) only of the volume on that link. Therefore, congestion caused by blocking of one link by another congested link is not modeled. Carey (1987) reformulated this model as a convex program by assuming the link departure function to be a maximum outflow instead of an actual

outflow, hence constraining by nonlinear inequality instead of equality. Thus, vehicles may be held back to benefit the system-optimality criterion, but not due to congestion on other links since the model remains uncapacitated. Departure functions can be included in the dynamical system model presented in this paper (see Sections 6.3 and 6.4), and combination with impedance functions for link travel times will prevent unreasonably short link travel times that might otherwise occur with short time intervals $\Delta t$.

In the last few years, there has been renewed interest in the problem of optimal traffic assignment; see for instance the special issue of *Transportation Research: B*, edited by Gartner and Improta (1990). We now compare our approach with some of the recent work.

Papageorgiou *et al.* (1990) have studied the general requirements of dynamic models based on standard state space methods. Our approach is different from theirs in several respects. A major difference is that their model is macroscopic and continuous in nature, employing traffic volumes, traffic densities and rates of traffic volumes, whereas our model is microscopic in the sense that it tracks explicitly platoons of vehicles on links. Another difference is that the model of Papageorgiou *et al.* requires the specification of dynamical equations (that will be part of the complete state space model) to model the propagation of *composition rates* along links, where these rates represent the proportion of vehicles on a link flowing to a particular destination. Instead, we explicitly calculate and record in the state the exit time of a vehicle from its current link using a combination of impedance and link outflow functions.

Friesz *et al.* (1989), Wie *et al.* (1990), and Ran *et al.* (1993) have also approached dynamic traffic assignment from an optimal control perspective. In each case traffic flow is modeled as a real-valued continuous-time process. The model that we present is discrete time and integer valued. Furthermore, in the above works, the definition of dynamic user-optimality requires instantaneous user-optimal travel costs for all routes that are being used at each instant of time to be equal. This is different from the notion of optimal anticipatory routing that we are interested in, in the context of a system-optimal criterion (recall the discussion in Section 1). It should be noted that one of the features of the model of Boyce *et al.* (1991) is the consideration of exit flows as control variables; our use of blocking controls in Section 6.2 is conceptually similar to that.

Finally, Janson (1991a, 1991b) has addressed the dynamic user-equilibrium traffic assignment problem for networks with multiple origins and destinations and known time-varying travel demands. He has developed a mathematical program for dynamic user-equilibrium assignment and proposed an algorithm for the solution of this problem. The algorithm is based on a two-step decomposition of the problem that is then solved iteratively. His simulations have shown that the iterative procedure exhibits good near-convergence properties. Our approach differs from Janson's work in at least two respects. First, we are considering a system-optimal objective rather than a user-equilibrium one. Second, our discretization of time is fundamentally different. In Janson (1991a, 1991b), the discrete time interval $\Delta t$ must be chosen large enough so that vehicles completely traverse any link in one time interval. Otherwise, the links need to be broken into smaller ones resulting in an increase of the dimensionality of the mathematical program. In our formulation, $\Delta t$ must be chosen small enough so that vehicles will not traverse more than one link in one time interval. The number of decision variables in our optimization problem increases linearly in the number of time intervals.

## 3. THE STATE SPACE MODEL

We present in this section a dynamical system model of a traffic system having multiple origins and destinations. (The discussion in this section can also be found in Lafortune *et al.* (1991).) Each link in the network may be described by a first-order difference equation, and the network as a whole is represented by another first-order difference equation that is an aggregation of the link equations. Thus our dynamical equation possesses an attractive and simple modular structure.

### 3.1. Definition of the main variables

The geographical network is viewed as a finite-vertex directed graph in which every edge is associated with one and only one ordered pair of vertices. Thus any origin or destination must necessarily be a vertex and any highway or arterial link connecting two points with no intervening point of interest is an edge on the directed graph. The structure is formalized as follows:

1. $|X|$ = number of elements of a finite set $X$.
2. $V = \{1, \ldots, |V|\}$ = set of vertices or nodes $v$ of a network.
3. $E = \{1, \ldots, |E|\}$ = set of directed edges or links $\ell$ of a network.
4. $D$ = set of destination nodes $d$. This is a subset of $V$.

The set of origin nodes does not have to be explicitly referred to in the dynamical equations that follow. Therefore, without loss of generality, we assume that this set is equal to $V$.

Because the formulation is discrete time, we define:

1. $t$ = discrete time index. Thus $t \in \mathbb{N}$.
2. $T_\ell$ = maximum possible number of sampling instants spent by a vehicle on the link $\ell \in E$.

Thus, it is assumed that there exists an upper bound on the time that a vehicle may spend on a given link. We consider such an assumption tenable if a link $\ell$ is never loaded in excess of a defined capacity $c_\ell$; the capacity $c_\ell$ does not represent the blockage capacity, but rather the maximum capacity that is deemed acceptable on link $\ell$. This restriction will be introduced when the region of admissible behavior of the model is defined. Moreover, it may be noted that $T_\ell$ is independent of the time at which a vehicle joins the link $\ell$, i.e. of prevailing traffic conditions at the time of entry. The variable $T_\ell$ is a function of highway or arterial capacity, which is also time invariant. (In Section 6, this model will be generalized to allow vehicles to spend more than $T_\ell$ units of time on link $\ell$.)

The state of the network at any given sampling time must reflect the location of each vehicle on the network at that time. Accordingly we define a state variable that groups vehicles on a link according to their destinations and exit times from the link. Thus for all $\ell \in E$, $d \in D$ and $1 \leq k \leq T_\ell$ define

$x_k^{d\ell}(t)$ = number of vehicles on link $\ell$ and traveling to destination $d$, present on the link in the time interval $[t, t + k)$.

Note that the traffic represented by $x_k^{d\ell}(t)$ is assumed to be on its next link at the time instant $t + k$. Also, this variable is integer-valued since our objective is to develop a model that is microscopic in nature.

The state of the network is affected by vehicles entering or leaving links, which suggests the use of input and output variables. In accordance with this intuition we define for all $\ell \in E$, $d \in D$ and $1 \leq k \leq T_\ell$:

$u^{d\ell}(t)$ = number of vehicles traveling to destination $d$ that are on link $\ell$ at time $t + 1$ but are not on it at time $t$ (i.e. these vehicles entered $\ell$ during $(t, t + 1]$).
$y^{d\ell}(t)$ = number of vehicles traveling to destination $d$, that are on link $\ell$ at time $t$ but are not on it at time $t + 1$ (i.e. these vehicles left $\ell$ during $(t, t + 1]$).

Observe that the above definitions imply that the input at time $t$ will affect the state at time $t + 1$. The output, as we shall see later, is simply that part of the state required to maintain the conservation of flow.

The last prerequisite for formulating a state equation is a functional relationship between input and state. This is where we resort to the impedance function. For all $\ell \in E$ we define

$f_\ell(z)$ = impedance of link $\ell$ when loaded with $z$ vehicles;

      = travel time of a vehicle joining $\ell$ when $z$ vehicles are traveling on this link.

The impedance function is assumed to be an integer-valued staircase function (as sketched in Fig. 1) that is right continuous. Formally the impedance function may be defined as

$$f_\ell(z) = \sum_{k=T_0}^{T_\ell} k\chi_{[a_k, b_k)}(z),$$

where $a_{T_0} = 0$, $b_{T_\ell} = \infty$, $b_{T_\ell - 1} < \infty$ and $b_k = a_{k+1}$ for $T_0 \leq k \leq T_\ell - 1$, and $\chi$ is the usual characteristic function.

We adopt for all $x = (x_1, \ldots, x_n)^T \in \mathbb{R}^n$ the notation $\|x\| = |x_1| + |x_2| + \ldots + |x_n|$.

### 3.2. The link dynamical equation

We define a dynamical equation for the vehicles on link $\ell$ and traveling to various destinations $d$. It is assumed that the initial conditions are stated as the state $x_k^{d\ell}(t_0)$, $1 \leq k \leq T_\ell - 1$ where $t_0$ is the initial time. To begin, let the link have no input (i.e. no new vehicles entering the link). We assume that the exit time of a vehicle from link $\ell$ is fixed at the time of entry by the impedance function as follows:

$$x_k^{d\ell}(t + 1) = x_{k+1}^{d\ell}(t) \qquad 1 \leq k \leq T_\ell - 1$$

and

$$x_{T_\ell}^{d\ell}(t + 1) = 0.$$

The first equation states that the vehicles that were on the link at time $t$ and were supposed to stay until just before $t + k + 1$ (corresponding to $x_{k+1}^{d\ell}(t)$) will still be on the link at time $t + 1$ and will stay on the link until just before $t + 1 + k$ (corresponding to $x_k^{d\ell}(t + 1)$). The second equation expresses the fact that none of the vehicles present on the link at time $t$ can take more than $T_\ell$ time instants to travel the link. This will be generalized in Sections 6.2–6.4 where blocking of vehicles at the end of a link and link outflow functions will be incorporated into the model.

Thus we obtain a linear homogeneous first order difference equation
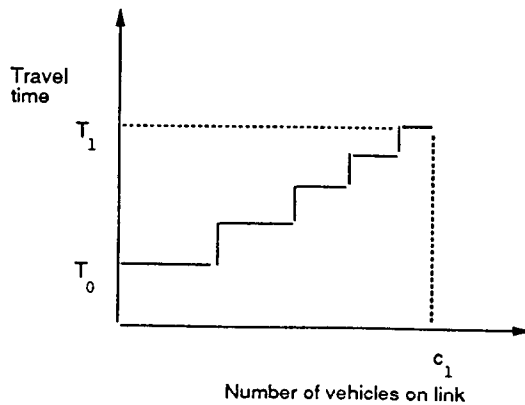
$$x^{d\ell}(t + 1) = A^{d\ell}x^{d\ell}(t),$$



Fig. 1. Impedance function. $c_\ell$ = capacity; $T_\ell$ = Maximum travel time.

where

$$x^{d\ell}(t) = [x_1^{d\ell}(t) \ldots x_{T_\ell}^{d\ell}(t)]^T,$$

and

$$A^{d\ell} = \begin{bmatrix} 0 & 1 & 0 & 0 & \ldots & 0 & 0 \\ 0 & 0 & 1 & 0 & \ldots & 0 & 0 \\ 0 & 0 & 0 & 1 & \ldots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ldots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \ldots & 0 & 1 \\ 0 & 0 & 0 & 0 & \ldots & 0 & 0 \end{bmatrix}_{T_\ell \times T_\ell}$$

Therefore the time evolution of the homogeneous (zero-input) dynamical equation reflects the exit of vehicles from a link. The modeling of vehicle exit as an internal dynamic of the link gives us the following output equation

$$y^{d\ell}(t + 1) = C^{d\ell}x^{d\ell}(t),$$

where

$$C^{d\ell} = [1 \ 0 \ \ldots \ 0]_{1 \times T_\ell}.$$

Now, to incorporate the input $u^{d\ell}(t)$, a simple extension of the state and input is necessary. Let

$$x^\ell(t) = [x^{1\ell}(t)^T \ldots x^{|D|\ell}(t)^T]_{T_\ell |D| \times 1}^T$$

and

$$u^\ell(t) = [u^{1\ell}(t) \ldots u^{|D|\ell}(t)]_{|D| \times 1}^T.$$

Observe that

$$x^\ell(t + 1) = A^\ell x^\ell(t),$$

where

$$A^\ell = \text{diag}\{A^{1\ell} \ldots A^{|D|\ell}\}_{(T_\ell |D|) \times (T_\ell |D|)}.$$

**Remark:** For given matrices $M^j \ j = 1, \ldots, J$, we shall use the notation $\text{diag}\{M^1 \ldots M^J\}$ to represent the block-diagonal matrix

$$\begin{bmatrix} M^1 & 0 & \ldots & 0 \\ 0 & M^2 & & \vdots \\ \vdots & & \ldots & 0 \\ 0 & \ldots & 0 & M^J \end{bmatrix}$$

For each link and destination we define for all $k$ such that $1 \leq k \leq T_\ell$

$$g_k^{d\ell}\left(x^\ell,u^\ell\right) = \begin{cases} 1 & \text{if } f_\ell(\|A^\ell x^\ell\| + \|u^\ell\|) = k \\ 0 & \text{otherwise.} \end{cases}$$

The function $f_\ell$ computes the time spent on link $\ell$ by vehicles joining the link at a given time. The function $g_k^{d\ell}(\cdot)$ is used to schedule events in the future based on the current state. Therefore the state dynamics are deterministically event-driven and the appropriate state variable may be incremented. Accordingly we obtain the first-order difference equation

$$x^{d\ell}(t + 1) = A^{d\ell}x^{d\ell}(t) + g^{d\ell}(x^\ell(t),u^\ell(t)) \cdot u^{d\ell}(t),$$

where

$$g^{d\ell}(x^\ell,u^\ell) = [g_1^{d\ell}(x^\ell,u^\ell) \cdots g_{T_\ell}^{d\ell}(x^\ell,u^\ell)]^T.$$

Next let

$$G^\ell(x^\ell,u^\ell) = \text{diag}\,\{g^{1\ell}(x^\ell,u^\ell) \cdots g^{|D|\ell}(x^\ell,u^\ell)\},$$

which is of dimension $(T_\ell|D| \times |D|)$. Then it is evident that

$$x^\ell(t + 1) = A^\ell x^\ell(t) + G^\ell(x^\ell(t),u^\ell(t))u^\ell(t).$$

The extension to link output is made as usual. Let

$$y^\ell(t) = [y^{1\ell}(t) \cdots y^{|D|\ell}(t)]_{|D| \times 1}^T$$

and

$$C^\ell = \text{diag}\{C^{1\ell} \cdots C^{|D|\ell}\},$$

which is of dimension $|D| \times (T_\ell|D|)$. Then

$$y^\ell(t) = C^\ell x^\ell(t).$$

### 3.3. The network dynamical equation

The network dynamical equation is very easily written by exploiting the modular structure of the system. We define the following notation:

1. $X(t) = [x^1(t)^T \cdots x^{|E|}(t)^T]^T$.
   Thus $X(t)$ is a $\left(\Sigma_{\ell=1}^{|E|}T_\ell|D|\right) \times 1$ dimensional vector.
2. $U(t) = [u^1(t)^T \cdots u^{|E|}(t)^T]^T$.
   Thus $U(t)$ is a $(|D||E|) \times 1$ dimensional vector.
3. $A = \text{diag}\{A^1 \cdots A^{|E|}\}$.
   Thus A is of dimension $\left(\Sigma_{\ell=1}^{|E|}T_\ell|D|\right) \times \left(\Sigma_{\ell=1}^{|E|}T_\ell|D|\right)$.
4. $G(X,U) = \text{diag}\{G^1(x^1,u^1) \cdots G^{|E|}(x^{|E|},u^{|E|})\}$.
   Thus $G(X,U)$ is of dimension $\left(\Sigma_{\ell=1}^{|E|}T_\ell|D|\right) \times (|D||E|)$.
5. $C = \text{diag}\{C^1 \cdots C^{|E|}\}$.
   Thus C is of dimension $(|E||D|) \times (\Sigma_{\ell=1}^{|E|}T_\ell|D|)$.

The required dynamical equation for the whole network is then

$$X(t + 1) = AX(t) + G(X(t),U(t))U(t).$$

It is assumed that the initial conditions are stated as the quantity $X(t_0) \equiv X_0$ where $t_0$ is the initial time. The homogeneous system remains linear first order and the overall system is also first order. The network output equation is

$$Y(t) = CX(t).$$

### 3.4. Definition of the feasible region for routing decisions

The following are the input and state constraints of the dynamical system:

*1. Flow conservation equations.*

For all nodes $v \in V$ we have the following equations for each destination $d \in D$:

$$d \neq v \Rightarrow \sum_{\ell \in S_v} u^{d\ell}(t) = \sum_{\ell \in P_v} y^{d\ell}(t) + r^{dv}(t), \, t = 0,1, \ldots$$

$$d = v \Rightarrow \forall \ell \in S_v, \, u^{d\ell}(t) = 0, \, t = 0,1, \ldots,$$

where

$S_v$ = set of successor links of node $v$;
$P_v$ = set of predecessor links of node $v$;
$r^{dv}(t)$ = number of new vehicles entering the network at node $v$ in the interval $(t,t + 1]$ and traveling to destination $d \in D$.
$R(t) = [r^{11}(t)r^{21}(t)\cdots r^{|D|1}(t)r^{12}(t)\cdots r^{|D||V|}]^T$ is the $(|V||D| \times 1)$ vector of travel demands for all nodes $v \in V$ and all destinations $d \in D$.

The first equation is the usual nodal flow conservation equation. The inflows and outflows to a node are balanced. The $r^{dv}(t)$ term and the second equation arise because each node is a possible source or sink for traffic. It may be noted that the flow at each node is conserved by destination and not only as total nodal flow.

*2. Headway constraints.*

We define for each $\ell \in E$ the quantity $K_\ell$ as

$$K_\ell = \frac{\Delta t}{t_{\text{headway}}} \times \text{number of lanes on } \ell$$

where $\Delta t$ is the real-time value of the sampling interval and $t_{\text{headway}}$ is the specified minimum separation time between two vehicles. The constraint equation is

$$\sum_{d \in D} x_k^{d\ell}(t) \leq K_\ell, \, 1 \leq k \leq T_\ell, \, t = 0,1, \ldots$$

*3. Capacity constraints.*

For all $\ell \in E$ we require

$$\sum_{d \in D} \|x^{d\ell}(t)\| \leq c_\ell, \, t = 0,1 \ldots$$

where $c_\ell$ is the capacity of link $\ell$.

We remind the reader of our remark made at the beginning of Section 3, about the existence of $T_\ell$. Since the impedance function is bounded on the interval $[0, c_\ell]$ the capacity constraint ensures that the range of the impedance function is confined to the mapping of the interval $[0, c_\ell]$. Thus we may assume that $T_\ell = f_\ell(c_\ell)$.

The three sets of equations in 1–3 together with a given state $X(t)$ and demand $R(t)$ define a feasible region for the input $U(t)$. We refer to this set of admissible controls by the function $\Omega(X(t),R(t))$. Formally,

$$\Omega(X,R) = \left\{ U \in \mathbb{N}^{|D||E|} : \left[ (\forall v \in V)(\forall d \in D)\left( d \neq v \Rightarrow \left[ \sum_{\ell \in S_v} u^{d\ell} = \sum_{\ell \in P_v} C^{d\ell} x^{d\ell} + r^{dv} \right] \right) \right. \right.$$

$$\wedge \left( d = v \Rightarrow \left[ (\forall \ell \in S_v) u^{d\ell} = 0 \right] \right) \right]$$

$$\wedge \left[ (\forall \ell \in E)\left( \left[ 1 \leq k \leq T_\ell \Rightarrow \| M_k(A^\ell x^\ell + G^\ell(x^\ell,u^\ell)u^\ell) \| \leq K_\ell \right] \right. \right.$$

$$\left. \left. \left. \wedge \left[ \| A^\ell x^\ell + G^\ell(x^\ell,u^\ell)u^\ell \| \leq c_\ell \right] \right) \right] \right\}.$$

The matrix $M_k$ extracts all the $x_k^{d\ell}$ components from $x^\ell$ and may be defined as follows. Let

$$m_k = [0\ 0\ \cdots\ 0\ 1\ 0\ \cdots\ 0]_{1 \times T_\ell},$$

where the 1 is in the $k$th position, and

$$M_k = \text{diag}\{m_k, \cdots, m_k\}_{|D| \times |D|T_\ell}.$$

**Remark:** For simplicity of notation, we have omitted possible time dependencies of certain variables in the above presentation. In general, $K_\ell, c_\ell,$ and even the impedance function $f_\ell(\cdot)$ could be time-varying.

### 4. MATHEMATICAL STATEMENT OF THE OPTIMAL CONTROL PROBLEM

From Sections 3.3 and 3.4 we obtain the equation of motion of the system and the control constraint set function, i.e. the set of feasible inflows. These are

$$X(t + 1) = AX(t) + G(X(t),U(t))U(t)$$

$$U(t) \in \Omega(X(t),R(t)).$$

(For simplicity of notation, we will abbreviate $G(X(t),U(t))U(t)$ as $GU(t)$ in this section.)

The problems are modeled as decision networks in which the decision nodes are pairs $(t, X(t))$ such that $X(t)$ is a system state reachable at time $t$. The decision arcs represent the control variables (link inputs) $U(t)$, with arc $U(t)$ present and connecting $(t, X(t))$ with $(t + 1, X(t + 1))$ if and only if $U(t) \in \Omega(X(t), R(t))$ and $X(t + 1) = AX(t) + GU(t)$. By imposing a cost structure and boundary conditions on the decision network we formulate two problems, one fixed-endtime free-endpoint problem and the other free-endtime fixed-endpoint.

#### 4.1. Fixed-endtime free-endpoint problem

Let the study horizon be a fixed integer terminal time denoted by $T_f$. The boundary conditions are then

$$X(t_0) = X_0$$

$$X(T_f) \in \mathbb{N}^{\Sigma_{\ell=1}^{|E|} T_\ell |D|}.$$

Having obtained our two point boundary value problem we proceed to define the cost. For each state trajectory in the interval $[t_0, T_f]$, denoted $X(t; t_0, T_f)$ the associated cost $J$ is defined as

$$J(X(t; t_0, T_f)) = \sum_{j=t_0}^{T_f} \|X(j)\| + F(X(T_f)).$$

The rationale for our definition is as follows. The norm of the state represents the number of vehicles on the network during one sampling interval. Thus, if the sampling time is one hour, then the first term of the cost function gives us the total number of vehicle-hours incurred within the time horizon. The second term represents a penalty for not clearing the network within the allotted time $T_f$, since all nonzero terms in $X(T_f)$ represent vehicles that have not yet reached their destinations.

The optimal cost as a function of the initial state may be defined as

$$J^*(X_0) = \min \Big\{ J(X(t; t_0, T_f)) : X(t_0) = X_0, X(j + 1) = AX(j) + GU(j),$$

$$\text{and } U(j) \in \Omega(X(j), R(j)), t_0 \leq j < T_f \Big\}.$$

Pick points $t_1, \ldots, t_p$ such that $t_0 < t_1 < \cdots < t_p < T_f$. From the definition of $J$ it is evident that

$$J(X(t; t_0, T_f)) = J(X(t; t_0, t_1)) + \cdots + J(X(t; t_{p-1}, t_p)) + J(X(t; t_p, T_f)),$$

where if $t_k < T_f$, then $J(X(t; t_0, t_k)) = \sum_{j=t_0}^{t_k} \|X(j)\|$. Thus by the above additivity $J$ obeys the principle of optimality. Furthermore the costs associated with trajectories from a state and those to the state are independent of each other. Consequently we develop recursive dynamic programming equations to solve the optimization problem.

### 4.2. Dynamic programming recursion equation for the fixed-endtime free-endpoint problem

*4.2.1. Forward recursion equation.* For a given initial state $X_0$, we define the cost $J_F$ to reach state $X$ at time $t$ along trajectory $X(j; t_0, t)$ to be

$$J_F(X, X(j; t_0, t), t) = \begin{cases} \sum_{j=t_0}^{t} \|X(j)\| & \text{if } t < T_f \\ \sum_{j=t_0}^{t} \|X(j)\| + F(X) & \text{if } t = T_f \end{cases},$$

where $X(j; t_0, t)$ must satisfy $X(t_0) = X_0$ and $X(t) = X$. Then

$$J(X(t; t_0, T_f)) = J_F(X(T_f), X(j; t_0, T_f), T_f).$$

For each feasible $X_t$ at a time instant $t$ we define the following forward value function $V_F$:

$$V_F(X_t, t) = \min \Big\{ J_F(X_t, X(j; t_0, t), t) : X(t_0) = X_0, X(t) = X_t, X(j + 1)$$

$$= AX(j) + GU(j), \text{ and } U(j) \in \Omega(X(j), R(j)), t_0 \leq j < t \Big\}.$$

Thus $V_F(X_t, t)$ represents the lowest cost to reach state $X_t$ at time $t$ among all admissible state trajectories between $X_0$ and $X_t$. It is evident that

$$J^*(X_0) = \min\{V_F(X_{T_f}, T_f): X_{T_f} \in X_{T_f}\},$$

where

$$X_{T_f} = \Big\{X : X(t_0) = X_0, X(T_f) = X, X(j + 1) = AX(j) + GU(j),$$

$$\text{and } U(j) \in \Omega(X(j), R(j)), t_0 \leq j < T_f\Big\}$$

is the set of all feasible terminal states for this fixed endtime problem.

The initial condition on $V_F(X, t)$ is

$$V_F(X_0, t_0) = \|X_0\|.$$

By the principle of optimality the recursion equation is as follows.

**Case** $t_0 < t < T_f$:

$$V_F(X_t, t) = \|X_t\| + \min_{X_{t-1} \in X_{t-1}} V_F(X_{t-1}, t - 1),$$

where

$$X_{t-1} = \{X : X_t = AX + GU, U \in \Omega(X, R(t - 1))\}.$$

**Case** $t = T_f$:

$$V_F(X_{T_f}, T_f) = F(X_{T_f}) + \|X_{T_f}\| + \min_{X_{T_f-1} \in X_{T_f-1}} V_F(X_{T_f-1}, T_f - 1),$$

where $X_{T_f-1}$ is defined as above.

*4.2.2. Backward recursion equation.* Analogous to the prior case we define the cost to complete from a given state $X$ at a given time $t$ along a trajectory $X(j; t, T_f)$ as

$$J_B(X, X(j; t, T_f), t) = \sum_{j=t}^{T_f} \|X(j)\| + F(X(T_f)),$$

where $X(j; t, T_f)$ must satisfy $X(t) = X$. Then

$$J(X(t; t_0, T_f)) = J_B(X(t_0), X(j; t_0, T_f), t_0).$$

For each feasible $X_t$ at time instant $t$ we define

$$V_B(X_t, t) = \min \Big\{J_B(X_t, X(j; t, T_f), t) : X(t) = X_t, X(j + 1) = AX(j) + GU(j)$$

$$\text{and } U(j) \in \Omega(X(j), R(j)), t \leq j < T_f\Big\}$$

as the backward value function at state $X_t$ at time $t$. Accordingly,

$$J^*(X_0) = V_B(X_0, X(j; t_0, T_f), t_0).$$

For each $X \in X_{T_f}$ (as defined in Section 4.2.1) the boundary condition is

$$V_B(X, T_f) = \|X\| + F(X).$$

By the principle of optimality, for all $t$ such that $t_0 \le t < T_f$, the backward recursion equation is

$$V_B(X,t) = \|X\| + \min_{U \in \Omega(X,R(t))} V_B(AX + GU, t + 1).$$

While the backward recursion equation appears more elegant in its formulation, we consider the forward recursion equation to be more useful. In the absence of a definite endpoint the set of feasible terminal states is too large to allow computation. The well defined initial condition, on the other hand, allows forward chaining through the solution space in a well defined recursive manner. An example illustrating the pruning of forward search using Dijkstra's implementation of dynamic programming is presented in Section 5.

### 4.3. Free-endtime fixed-endpoint problem

Here the control objective is to clear the network in as short a time as possible. Thus the target state is the zero vector. The dynamical equation and control constraint set $\Omega(X, R(t))$ are as usual.

We define the cost associated with a state trajectory to be

$$J(X(t;t_0,T)) = \sum_{j=t_0}^{T} \|X(j)\| \quad \text{where} \quad T \in \mathbb{N}.$$

The optimal cost $J^*$ as a function of the initial state $X_0$ at time $t_0$ is

$$J^*(X_0,t_0) = \min \left\{ J(X(t;t_0,T)) : X(t_0) = X_0, X(T) = 0, X(j + 1) \right.$$

$$\left. = AX(j) + GU(j), \text{ and } U(j) \in \Omega(X(j),R(j)), t_0 \le j < T, T \in \mathbb{N} \right\}.$$

However, if this problem is viewed as an infinite horizon one, i.e. all state trajectories are assumed defined on $[t_0, \infty)$, then the problem may be viewed as free-endpoint free-endtime. We define

$$J(X(t;t_0)) = \sum_{j=t_0}^{\infty} \|X(j)\|$$

and

$$J^*(X_0,t_0) = \min \left\{ J(X(t;t_0)) : X(t_0) = X_0, X(j + 1) = AX(j) + GU(j), \right.$$

$$\left. \text{and } U(j) \in \Omega(X(j),R(j)), t_0 \le j \right\}.$$

Thus if $J^*(X_0,t_0)$ exists then $\lim_{j \to \infty} \|X(j)\| = 0$, which implies that the target is implicitly achieved and the exogenous demand $R$ terminates in finite time. Moreover, the requirements of integrality force $X(j) = 0$ in finite time.

### 5. AN EXAMPLE ILLUSTRATING THE FORWARD DP RECURSION

The triangle network of Fig. 2 is used. $A$ and $B$ are valid origin nodes and $B$ and $C$ are valid destination nodes. Accordingly,

$$V = \{A,B,C\}$$

Fig. 2. Triangle network.

$$E = \{1 = (A,B), 2 = (B,C), 3 = (A,C)\}$$

$$D = \{B,C\}.$$

The impedance functions are presented in tabular form in Table 1.

We wish to find a system-optimal routing for the following simple demand pattern:

$$r^{BA}(0) = 1$$

$$r^{CA}(0) = 1$$

$$r^{CA}(10) = 3.$$

Our objective is to clear the network, as in Section 4.3. To better illustrate the working of the DP algorithm, an equivalent but different formulation of the cost function is used. We define:

$$J(X(t;0,T)) = \sum_{t=0}^{T-1} \sum_{d \in D} \sum_{\ell \in E} \left[ \sum_{k=1}^{T_\ell} k g_k^{d\ell}(x^\ell(t), u^\ell(t)) \right] u^{d\ell}(t).$$

Thus the total cost of traveling a link is incurred as soon as the vehicle is routed to a link.

The flow of the algorithm is presented as a tree in Fig. 3 and the important computations are shown in Table 2. (By "Terminal" we mean a node that will reach the terminal state without any further inputs. The "not expanded" nodes are ones that are not expanded further since they cannot yield optimal solutions.)

Let

$$n(t) = (n_{AB}(t), n_{BC}(t), n_{AC}(t)).$$

Then

Table 1. Impedance functions

| Number of Vehicles | Travel Time-$AC$ | Travel Time-$AB$ and $BC$ |
|---|---|---|
| 1 | 20 | 15 |
| 2 | 20 | 17 |
| 3 | 30 | 22 |
| 4 | 40 | 33 |
| 5 | 60 | 45 |
| 6 | 100 | 70 |

Capacity = 6; Maximum travel time = 100 ($AC$); 70 ($AB$, $BC$).

| | | 011: J=169; (5,0,0) | Not expanded |
|---|---|---|---|
| | 01; J=34; (2,0,0) | 012: J=124; (2,0,3) | 0121: J=139; Terminal |
| | | 013: J=120; (4,0,1) | 0131: J=135; Not Expanded |
| | | 014: J=116; (3,0,2) | 0141: J=131; Not Expanded |
| 0: J=0; (0,0,0) | | 021: J=134; (4,0,1) | Not expanded |
| | 02; J=35; (1,0,1) | 022: J=170; (1,0,4) | Terminal |
| | | 023: J=104; (3,0,2) | 0231: J=138; Terminal |
| | | 024: J=112; (2,0,3) | 0241: J=127; Terminal and OPTIMAL |

Fig. 3. DP tree $(n_{AB}, n_{BC}, n_{AC})$ = (number of vehicles on $(A,B)$, number of veh. on $(B,C)$, number of veh. on $(A,C)$).

Table 2. Computational data for nodes of DP tree

| Node # | Time | State | Demand | Inputs | $g$ | Cost |
|---|---|---|---|---|---|---|
| 0 | 0 | Start State | | | | |
| 01 | 0, 1 | $x_{17}^{B1}=1, x_{17}^{C1}=1$ | $r^{BA}=1, r^{CA}=1$ | $u^{B1}=1, u^{C1}=1$ | $g_{17}^{B1}=1, g_{17}^{C1}=1$ | 34 |
| 02 | 0, 1 | $x_{15}^{B1}=1, x_{20}^{C3}=1$ | $r^{BA}=1, r^{CA}=1$ | $u^{B1}=1, u^{C3}=1$ | $g_{15}^{B1}=1, g_{20}^{C3}=1$ | 35 |
| 011 | 10, 11 | $x_8^{B1}=1, x_8^{C1}=1$; $x_7^{B1}=1, x_7^{C1}=1$; $x_{45}^{C1}=3$ | $r^{CA}=3$ | $u^{C1}=3$ | $g_{45}^{C1}=1$ | 169 |
| 012 | 10, 11 | $x_8^{B1}=1, x_8^{C1}=1$; $x_7^{B1}=1, x_7^{C1}=1$; $x_{30}^{C3}=3$ | $r^{CA}=3$ | $u^{C3}=3$ | $g_{30}^{C3}=1$ | 124 |
| 013 | 10, 11 | $x_8^{B1}=1, x_8^{C1}=1$; $x_7^{B1}=1, x_7^{C1}=1$; $x_{33}^{C1}=2, x_{20}^{C3}=1$ | $r^{CA}=3$ | $u^{C1}=2, u^{C3}=1$ | $g_{33}^{C1}=1, g_{20}^{C3}=1$ | 120 |
| 014 | 10, 11 | $x_8^{B1}=1, x_8^{C1}=1$; $x_7^{B1}=1, x_7^{C1}=1$; $x_{22}^{C1}=1, x_{25}^{C3}=2$ | $r^{CA}=3$ | $u^{C1}=1, u^{C3}=2$ | $g_{22}^{C1}=1, g_{25}^{C3}=1$ | 116 |
| 021 | 10, 11 | $x_6^{B1}=1, x_{11}^{C3}=1$; $x_5^{B1}=1, x_{10}^{C3}=1$; $x_{33}^{C1}=3$ | $r^{CA}=3$ | $u^{C1}=3$ | $g_{33}^{C1}=1$ | 134 |
| 022 | 10, 11 | $x_6^{B1}=1, x_{11}^{C1}=1$; $x_5^{B1}=1, x_{10}^{C3}=1$; $x_{45}^{C3}=3$ | $r^{CA}=3$ | $u^{C3}=3$ | $g_{45}^{C3}=1$ | 170 |
| 023 | 10, 11 | $x_6^{B1}=1, x_{11}^{C3}=1$; $x_5^{B1}=1, x_{10}^{C3}=1$; $x_{22}^{C1}=2, x_{25}^{C3}=1$ | $r^{CA}=3$ | $u^{C1}=2, u^{C3}=1$ | $g_{22}^{C1}=1, g_{25}^{C3}=1$ | 104 |
| 024 | 10, 11 | $x_6^{B1}=1, x_{11}^{C3}=1$; $x_5^{B1}=1, x_{10}^{C3}=1$; $x_{17}^{C1}=1, x_{30}^{C3}=2$ | $r^{CA}=3$ | $u^{C1}=1, u^{C3}=2$ | $g_{17}^{C1}=1, g_{30}^{C3}=1$ | 112 |
| 0111 | Not expanded by algorithm | | | | | |
| 0121 | 17 | $x_1^{B1}=1, x_1^{C1}=1$; $x_{39}^{C1}=3$ | | $u^{C2}=1$ | $g_{15}^{C2}=1$ | 139 |
| 0131 | 17 | $x_1^{B1}=1, x_1^{C1}=1$; $x_{27}^{C1}=2, x_{14}^{C3}=1$ | | $u^{C2}=1$ | $g_{15}^{C2}=1$ | 135 |
| 0141 | 17 | $x_1^{B1}=1, x_1^{C1}=1$; $x_{16}^{C1}=2, x_{18}^{C3}=1$ | | $u^{C2}=1$ | $g_{15}^{C2}=1$ | 131 |
| 0211 | Not expanded by algorithm | | | | | |
| 0221 | Terminal state | | | | | |
| 0231 | 32 | $x_1^{C1}=2, x_4^{C3}=1$ | | $u^{C2}=2$ | $g_{17}^{C2}=1$ | 138 |
| 0241 | 27 | $x_1^{C1}=1, x_{14}^{C3}=2$ | | $u^{C2}=1$ | $g_{15}^{C2}=1$ | 127 |

All data not stated in the table is zero

$$n(0) = (0,0,0)$$
$$n(1) = (1,0,1)$$
$$n(11) = (2,0,3)$$
$$n(41) = (0,0,0).$$

We use $n(t)$ as a simplified representation of the state trajectory of this example. The detailed description of the state trajectory using the notation developed in prior sections may be found in Table 2.

## 6. IMPORTANT EXTENSIONS OF THE BASIC MODEL

### 6.1. Background traffic

The model of Section 3 assumes that all vehicles are to be routed during their travel in the network. We now discuss how to include background traffic into the model. By background traffic we mean vehicles whose routing is not part of the optimization but rather is a known deterministic function of time. This function could for instance be based on historical data, on shortest path calculations, or on some traffic equilibrium solution.

We need to distinguish between *guided* and *background* vehicles. Let $X_g(\cdot)$ denote the state vector (as defined in Section 3.2) for guided vehicles and $X_b(\cdot)$ that for background vehicles. Similarly for $U_g(\cdot)$ & $U_b(\cdot)$, $Y_g(\cdot)$ & $Y_b(\cdot)$ and $R_g(\cdot)$ & $R_b(\cdot)$. Finally, define

$$X_{tt}(t) = X_g(t) + X_b(t)$$
$$U_{tt}(t) = U_g(t) + U_b(t)$$
$$R_{tt}(t) = R_g(t) + R_b(t)$$

for all $t$, where $tt$ stands for "total."

For the background traffic, by assumption,

$$u_b^{dt}(t) = f_b(t, Y_b(t), r_b^{dv}(t)),$$

where $f_b$ is a known function and where $v$ is the origin node of $\ell$. In order to obtain a meaningful problem formulation, we assume that

$$U_b(t) \in \Omega(X_b(t), R_b(t)),$$

i.e. the routing of the background traffic should satisfy the flow conservation equations and the headway and capacity constraints in the absence of guided vehicles ($X_g(\cdot) \equiv 0$ and $R_g(\cdot) \equiv 0$).

The network dynamical equation now consists of two parts:

$$X_g(t + 1) = AX_g(t) + G(X_{tt}(t), U_{tt}(t))U_g(t)$$
$$X_b(t + 1) = AX_b(t) + G(X_{tt}(t), U_{tt}(t))U_b(t)$$
$$Y_g(t) = CX_g(t)$$
$$Y_b(t) = CX_b(t),$$

with appropriate initial conditions $X_b(t_0)$ and $X_g(t_0)$. The feasible region for routing decisions $U_g(\cdot)$ is defined by the flow conservation equations

$$\sum_{\ell \in S_v} u_g^{dt}(t) = \sum_{\ell \in P_v} y_g^{dt}(t) + r_g^{dv}(t) \qquad [d \neq v]$$

$$\forall \ell \in S_\nu, \ u_g^{d\ell}(t) = 0 \qquad [d = \nu],$$

and by the headway and capacity constraints ($M_k$ below is as defined in Section 3.4)

$$\sum_{d \in D} M_k x_{\ell\ell}^{d\ell}(t) \le K_t, \ 1 \le k \le T_t, \ t = 0, 1, \ldots$$

$$\sum_{d \in D} \|x_{\ell\ell}^{d\ell}(t)\| \le c_t, \ t = 0, 1, \ldots$$

In other words, the headway and capacity constraints are in terms of the total traffic, while the flow conservation equations are in terms of the guided traffic, since by assumption the routing of the background traffic satisfies its own set of flow conservation equations.

Finally, the cost of a state trajectory

$$X(t;t_0 T_f) = \begin{bmatrix} X_g(t;t_0,T_f) \\ X_b(t;t_0,T_f) \end{bmatrix}$$

could in general depend on both types of traffic:

$$J(X(t;t_0,T_f)) = \sum_{j=t_0}^{T_f} \alpha_1 \|X_g(j)\| + \beta_1 \|X_b(j)\| + \alpha_2 F(X_g(T_f)) + \beta_2 F(X_b(T_f)),$$

where $\alpha_1$, $\alpha_2$, $\beta_1$, and $\beta_2$ are weighting factors.

### 6.2. Blocking controls

We have assumed that link travel time is strictly determined by an impedance function of the number of vehicles on the link. In reality, the travel time on one link may be affected by loadings on nearby links. For instance, suppose that a vehicle assigned to depart a link at time $(t + 1)^-$ cannot enter any of the links immediately downstream due to capacity constraints. With the current model as formulated in Section 3.4, this would imply infeasibility of the state $X(t)$ under consideration. However, an actual vehicle blocked in this way would simply remain on its current link until downstream capacity becomes available, affecting the vehicles on its link and in some cases vehicles on upstream links. Link interactions may be an essential part of congestion modeling, particularly with respect to real-time traffic management in the presence of unforeseen incidents.

In order to address the above limitation, we enlarge the set of admissible controls by (a) interpreting the impedance function $f_t$ not as an actual but as a minimum link travel time; and (b) including into the model routing decisions that actually "block" or "stall" a vehicle at the end of a link instead of immediately routing this vehicle on a successor link (if feasible). We will refer to these decisions as *blocking controls* and use the notation

$$U_{\text{block}}(t) = \left( u_{\text{block}}^{d\ell}(t) \right)_{|D||E| \times 1}$$

to differentiate them from the normal routing decisions $U(t)$. More precisely,

$u_{\text{block}}^{d\ell}(t) =$ number of vehicles traveling to destination $d$ that are on link $\ell$ at time $t$ and are due to exit that link at time $(t + 1)^-$, but will be rescheduled to exit $\ell$ at time $(t + 2)^-$ instead (i.e. they are blocked at the end of $\ell$ for one time period).

The new dynamical equation for this "two-input" system is of the form

$$X(t + 1) = AX(t) + G(X(t),U(t))U(t) + G_{block}U_{block}(t).$$

The new term represents the effect of the blocking controls. When $U(t) \equiv 0$, we have

$$x_1^{d\ell}(t + 1) = x_2^{d\ell}(t) + u_{block}^{d\ell}(t),$$

which can be written more compactly as

$$x^{d\ell}(t + 1) = A^{d\ell}x^{d\ell}(t) + G_{block}^{\ell}u_{block}^{d\ell}(t),$$

where $G_{block}^{\ell} = [1\ 0\ \cdots\ 0]_{T_\ell \times 1}^T$.

Similar to the structure of $G(X,U)$, $G_{block}$ is of dimension $\left(\Sigma_{\ell=1}^{|E|}T_\ell|D|\right) \times (|D||E|)$ and is of the form

$$G_{block} = \begin{bmatrix} G_{block}^1 & & & & & \\ & \ddots & & & & \\ & & G_{block}^1 & & & \\ & & & G_{block}^2 & & \\ & & & & \ddots & \\ & & & & & G_{block}^{|E|} \end{bmatrix},$$

where the block-diagonal form contains $|D|$ copies of each $G_{block}^i$ vector.

Capacity and headway constraints are unchanged. However, the flow conservation equations become

$$\sum_{\ell \in P_v} u_{block}^{d\ell}(t) + \sum_{\ell \in S_v} u^{d\ell}(t) = \sum_{\ell \in P_v} y^{d\ell}(t) + r^{dv}(t) \qquad [d \neq v]$$

$$\forall \ell \in S_v, \, u^{d\ell}(t) = 0 \qquad [d = v]$$

$$\forall \ell \in P_v, \, u_{block}^{d\ell}(t) \leq y^{d\ell}(t) \qquad [d \neq v]$$

$$\forall \ell \in P_v, \, u_{block}^{d\ell}(t) = 0 \qquad [d = v].$$

Observe that no modification of the cost $J$ is necessary since the extra travel time incurred by the blocked vehicles is automatically summed up in $\|X(t)\|$.

Without any further constraints, the increase in the dimensionality of the decision space due to the consideration of all feasible blocking controls is likely to render the problem intractable even for small networks. A reasonable heuristic would be to allow blocking controls *only* when the feasible region $\Omega(X(t))$ (as defined in Section 3) is empty, i.e. when there is no routing decision $U(t)$ that satisfies (a) the flow conservation equations, (b) the headway constraints, and (c) the capacity constraints.

### 6.3. Link outflow functions

So far we have modeled link travel times as being fixed by an impedance function applied at the time of link entry, possibly subject to a link exit delay through the use of blocking controls. In this section we demonstrate that the model is easily altered to incorporate another widely used form of vehicle dynamics, and in the next section we generalize the model to include both dynamical forms.

The models of Merchant and Nemhauser (1978) and Friesz *et al.* (1989) use link outflow functions $h$, which give link output as a function of total link volume, without reference to impedance functions. In the single destination case without blocking, we can change our model to be similarly free of impedance functions by replacing $y_\ell(t) = C^\ell x^\ell(t)$ (cf. Section 3.2) by

$$y^\ell(t) = h_\ell(\|x^\ell(t)\|) \, \ell \in E.$$

The link outflow $y^\ell(t)$ depends solely on the total link volume at time $t$, hence, in contrast to what was done in Section 3.2, we need not divide the link volume into classes $x_k^\ell(t)$ by time $k$ until link departure (the superscript $d$ is omitted since there is a single destination). To complete the alternate model, we rewrite the dynamical equation of Section 3.3 and the constraints of Section 3.4 with this reduction in state space.

Carey (1987) considers the outflow function as an upper bound on outflow rather than actual outflow, writing $y_\ell(t) = C^\ell x^\ell(t)$ in order to obtain a convex mathematical program; this is equivalent to the outflow function version of our model with the addition of blocking controls. The remark concluding the previous section applies here as well.

The multidestination case presents the issue of how to balance link outflows across destination classes. In this case we constrain total outflows over all destination classes by

$$\sum_{d \in D} y^{d\ell}(t) = h_\ell(\|x^\ell(t)\|) \quad \ell \in E.$$

In some cases it may be desirable to prevent the system optimality criterion from delaying one destination class in favor of another. Carey suggests balancing outflows through constraints of the form

$$\frac{y^{1\ell}(t)}{y^{1\ell}(t)} = \frac{x^{d\ell}(t)}{x^{1\ell}(t)} \quad d \in D,$$

but the integral solutions we explore in the dynamic programming solution procedure would in general fail to satisfy this constraint. We can attempt to satisfy this criterion by placing a penalty on violation of the balancing constraints into the objective function. Another alternative for balancing would be to define two-argument outflow functions $h'$ that take as inputs both the total traffic volume on the link and volume associated with the particular destination. Then

$$y^{d\ell}(t) = h_\ell'(\|x^\ell(t)\|, x^{d\ell}(t)) \quad \ell \in E.$$

We remark that link outflow functions seem to model nonrecurring congestion, i.e. capacity reductions due to traffic incidents, more closely than impedance functions or changes in link capacity constraints. It is not the physical capacity of the link to hold vehicles that is reduced, but the maximum rate at which the end of the link can release vehicles.

### 6.4. Modeling impedance functions and link outflow functions together

We will now show that by combining impedance functions and link outflow functions into one model, we can overcome weaknesses associated with each idea taken alone. The combined model is easily interpreted as a generalization both of the link outflow function model of the previous section and of the impedance function model with blocking controls.

When traffic dynamics are modeled solely by outflow function constraints in the inequality sense, vehicles are allowed to exit lightly loaded links only one period after entering, regardless of link length. This would tend to imply vehicle speeds in excess of the maximum freeflow speed. This defect can be remedied by applying an impedance function at the time of link entry to compute a link travel time $s^\ell(t) = f_\ell(\|x^\ell(t)\|)$ for the vehicles entering link $\ell$ at time $t$, and to enforce this as the *minimum* link travel time by setting

$$D^{d\ell}(t) = \sum_{p=t}^{t+s^\ell(t)-2} y^{d\ell}(p) = \text{number of vehicles departing link } \ell \text{ while vehicles reaching the link at time period } t \text{ must remain}$$

(the −2 term in the above upper bound is due to the fact that, by definition, the output at time $t$ is reflected in the state at time $t + 1$), and constraining the link outflows by

$$D^{dt}(t) \le x^{dt}(t) - u^{dt}(t - 1) \; \ell \in E, d \in D.$$

Then only the vehicles that had entered the link before time $t$ can exit within the minimum impedance time $s^{\ell}(t)$ and vehicles that enter at time $t$ occur on a different link no earlier than $t + s^{\ell}(t)$. Note that this constraint also has the beneficial effect of contributing to balance of outflows across destinations by preventing very short link travel times for destination classes that might otherwise be favored by the system-optimality criterion. Furthermore, it enforces a first-in, first-out condition under which no vehicle can decrease its total travel time by delaying its departure at an intermediate node. Kaufman and Smith (1990) argue that this condition tends to be satisfied by traffic flows.

To see how this generalizes not only the outflow function model but the impedance function model with blocking as well, note that we can track satisfaction of the new constraints by assigning each group of entering vehicles to a data storage area according to its earliest possible link exit time. We originally tracked vehicles by variables $x_k^{dt}(t)$ according to *actual* link exit time $t + k$, and the addition of blocking controls transforms this to *earliest* link exit time, subject to feasibility considerations. Then the headway constraints (Section 3.4) can be interpreted as a special case of the generalized model, with a constant outflow function.

Modeling link travel times as being given with equality by impedance functions ignores the interaction between links that is a necessary component of congestion modeling. The situation was improved by the addition of blocking controls, allowing vehicles to be delayed by downstream congestion. The addition of link outflow functions suggests a mechanism for modeling more complicated link interactions. For example, consider several links with the same exit node at which a traffic signal is located. The outflow function with the addition of a time index can model a preprogrammed variable signal timing plan, and if the link outflow functions on these links are given the flows on each link into the node as arguments, we can model a signal that responds in real time to traffic detectors.

## 7. COMPUTATIONAL ISSUES

We present herein a discussion of the computational complexity of solving the optimal control problems posed in Section 4. As was shown in that section, we are essentially dealing with shortest-path problems in the decision network, whose nodes are pairs, $(t, X(t))$ (recall that the decision network is not the same as the geographical road network). It is well known that an algorithm such as Dijkstra and based upon the dynamic programming equation developed in Section 4.2 will be $O(n^2)$ in the worst case, where $n$ is the number of nodes in the decision network. (The subsequent analysis is all worst-case.)

Let there be $N$ nodes in the geographical roadway network, i.e. $|V| = N$. Thus all roadways terminate or originate within this set of $N$ nodes. Furthermore, all traffic on links flowing into a node may be routed to any of the links flowing out of the node. Accordingly, the following worst-case upper bounds are obtained: number of links = $N^2$, number of destinations = $N$.

We assume next that the maximum number of time periods required to travel a link is denoted by $T_{max}$, i.e. $T_{max} = \max_{\ell \in E} T_\ell$. It is immediate that the quantity $T_{max} N^3$ is an upper bound on the dimension of the state. We next estimate the size of our decision network. Each node in this decision network is a state $X(t)$ for some time instant $t$. Each $X(t)$ is a parent of many $X(t + 1)$. Every $X(t + 1)$ is derived from the parent $X(t)$ by some input $U(t)$, which may be viewed as a label for the link connecting parent and child. The input $U(t)$ gives the increase in cost from parent to child. This cost is the travel time experienced by the vehicles routed in $U(t)$ (cf. example in Section 5).

No exact way of estimating the fanout of a state $X(t)$ is as yet known to us. It is

obviously dependent on the initial conditions and the demand pattern. The clarification of this relationship is the core of the optimization problem and our understanding of the same is embodied in the input space represented by $\Omega(X(t), R(t))$. Thus let us denote by $w^{dv}(t)$ the total number of vehicles (both inflow and external) arriving at a particular node $v$ of the geographical network at time $t$ and traveling to destination $d$. It can be shown that the maximum possible number of integral solutions to the flow conservation equations is given by

$$\binom{w^{dv}(t) + N - 1}{w^{dv}(t)}.$$

Thus the total number of possible inputs for the entire state $X(t)$ is bounded above by

$$\prod_{v=1}^{N} \prod_{d=1}^{N} \binom{w^{dv}(t) + N - 1}{w^{dv}(t)}, \tag{1}$$

since routing decisions must be made for each destination at each node. Consequently the maximum fanout for any state $X(t)$ of the decision network is given by eqn (1). It may be noted that the depth of the tree is the time horizon of optimization.

We suspect that the actual computation in most networks would fall far short of the aforementioned bounds. In lightly loaded networks, $w^{dv}(t)$ is small. In the case that $w^{dv} = 1$, we observe that eqn (1) evaluates to $N^3$. If the network is heavily loaded, it will be congested and much of the search space will be pruned by active capacity and outflow rate constraints.

Further avenues of improvement lie in the use of heuristics to yield near-optimal solutions. It is easy to see that for all vehicles constituting a particular $w^{dv}(t)$, their travel time from $v$ to $d$ must exceed the minimum freeflow travel time from $v$ to $d$. This freeflow time is a conservative estimate of the cost to complete from a given state, and thus it could be used to develop an $A^*$ optimization algorithm (see, e.g. Pearl, 1984) to find a shortest path in the decision network.

For heavily congested networks it is also possible to restrict the resolution to obtain near-optimal solutions. The estimate in eqn (1) is based on the premise that any one vehicle in the group $w^{dv}(t)$ can be routed to any link flowing out of the node $v$. We could decrease the resolution by assuming that vehicles may only be routed in packets of size $p$. Then our fanout estimate would come down to,

$$\prod_{v=1}^{N} \prod_{d=1}^{N} \binom{w^{dv}(t)/p\,(mod\,p) + 1 + N - 1}{w^{dv}(t)/p\,(mod\,p) + 1}$$

$$= \prod_{v=1}^{N} \prod_{d=1}^{N} \binom{w^{dv}(t)/p\,(mod\,p) + N}{w^{dv}(t)/p\,(mod\,p) + 1}.$$

Finally, observe that the necessary computations will be greatly reduced if the optimization is performed over a small time horizon. The optimization problem for the complete time horizon could then be tackled by employing a *rolling-horizon* strategy; the "optimal" control at a given node of the decision network would be determined on the basis of a forward search over a reduced horizon (or limited window into the future). This limited horizon could be as small as one step ahead. In this case the objective of optimization would be to minimize the travel time within the reduced horizon rather than to clear the network of all traffic. The cost function would be as stated at the end of Section 6.1 with $\alpha_2$, $\beta_2$ being chosen to prevent trivial solutions from being accepted as optimal. The control action so determined would then be used and the procedure repeated at the corresponding successor node in the decision network, until construction of a complete

path in the decision network. Such rolling-horizon strategies have been frequently employed in a variety of dynamic optimization problems; see, e.g. Gartner (1982) for their use in dynamic traffic signal control, a problem related to ours.

## 8. CONCLUSION

We have approached the problem of dynamic traffic assignment in networks from the viewpoint of dynamical systems and have proposed a new model for this problem. This model is more detailed than typical macroscopic models, yet it avoids complete microscopic detail by grouping vehicles into platoons irrespective of origin node and time of entry to the network. It has been observed in the literature that traffic models based on impedance functions alone or on link outflow functions alone suffer from certain deficiencies. By using impedance functions to first determine a minimum travel time for a vehicle on a link and then by using blocking controls (Section 6.2) and/or outflow functions (Sections 6.3–6.4) to limit the outflow from a link, our model effectively combines both of these approaches. We believe that this feature of the model is interesting because it closely resembles what happens in a microscopic simulation of traffic (e.g. the INTEGRATION traffic simulator of Van Aerde and Yagar, 1988).

In the context of this dynamical system model, two versions of the problem of optimal traffic assignment were formulated and studied. These problems can be solved by using dynamic programming over the state space. Due to the large size of the state space, the computational complexity is high even for very simple networks. However, the approach employed permits easy introduction of search heuristics, even as the $A^*$ algorithm or rolling-horizon strategies.

The work that we have presented opens several avenues for future investigations. Among these, we wish to mention the following.

1. Refinement of the impedance function. The link impedance function has been assumed to be a function of $z$, which represents the number of vehicles on the link. But this number is a highly aggregated version of the detailed information contained in the state $X(t)$. In fact, the state $X(t)$ is capable of supporting far richer arguments than $z$ for the impedance function. Thus the experienced traffic designer may study the behavior of the model using more complex and generalized impedance functions (see, e.g. Branston, 1976, for a review of impedance functions).
2. Determination of structural properties of optimal routing policies. Such properties could then be used to accelerate the forward search.
3. Investigation of various methodologies, either exact or heuristic, for the reduction of computation. Two such approaches (that in fact could be combined) are the $A^*$ algorithm and rolling-horizon strategies.
4. Coordination with real-time traffic control. Real-time traffic control (signalization, ramp metering, etc.) and anticipatory route guidance are coupled issues. On the one hand, routing decisions should account for real-time traffic control. On the other hand, real-time traffic control should adapt to routing decisions, e.g. when an incident triggers a rapid surge of vehicles off the corridors and onto the surface street network. Ways of modeling these interactions have to be developed.

## REFERENCES

Branston D. (1976) Link capacity functions: A review. *Transpn. Res.*, **10B**, 223–236.
Carey M. (1987) Optimal time-varying flows on congested networks. *Operations Research*, **35**, 58–69.
Friesz T. L., Lugue J., Tobin R. L. and Wie B. (1989) Dynamic network traffic assignment considered as a continuous-time optimal control problem. *Operations Research*, **37**, 893–901.
Gartner N. H. (1982) Development and testing of a demand-responsive strategy for traffic signal control. *Proc. 1982 American Control Conf.*, June, pp. 578–583, Arlington, VA.

Gartner N. H. and Improta G. Guest Editors (1990) Special issue on urban traffic networks: Dynamic control and flow equilibrium. *Transpn. Res.,* **24B**, 407–495.

Janson B. N. (1991a) A convergent algorithm for dynamic traffic assignment. In *Proc. of the 70th Transportation Research Board Meeting,* January, Washington, DC.

Janson B. N. (1991b) Dynamic traffic assignment for urban road networks. *Transpn. Res.,* **25B**, 143–161.

Kaufman D. E., Lee J. and Smith R. L. (1990) Anticipatory traffic modeling and route guidance in intelligent vehicle–highway systems. Technical Report 90-2, IVHS Program, University of Michigan, February.

Kaufman, D. E., and Smith, R. L. (1993) Fastest paths in time-dependent networks for intelligent vehicle-highway systems application. *IVHS Journal,* **1**, 1–11.

Lafortune S., Sengupta R., Kaufman D. and Smith R. L. (1991) A dynamical system model for traffic assignment in networks. In *Proc. 1991 Conf. on Vehicle Navigation & Information Systems,* Dearborn, MI, October, pp. 701–708.

Merchant K. D. and Nemhauser G. L. (1978) A model and an algorithm for the dynamic traffic assignment problem. *Transportation Science,* **12**, 183–199.

Papageorgiou M., Banos J. C. M. and Messmer A. (1990) Optimal control of multidestination traffic networks. In *Proc. 29th IEEE Conf. on Decision and Control,* Honolulu, HI, December, pp. 1355–1361.

Pearl J. (1984) *Heuristics: Intelligent Search Strategies for Computer Problem Solving.* Addison-Wesley, Reading, MA.

Ran, B., Boyce, D. E. and LeBlanc, L. J. (1993) A new class of instantaneous dynamic user-optimal traffic assignment models. *Operations Research,* **37**, 192–202.

Saxton L. Guest Editor (1991) Special issue on Intelligent-Vehicle Highway Systems. *IEEE Trans. Vehicular Technology,* **40**, 1–158.

Van Aerde M. and Yagar S. (1988) Dynamic integrated freeway/traffic signals networks: A routing-based modelling approach. *Transpn. Res.,* **22A**, 445–453.

Wie B., Friesz T. L. and Tobin R. L. (1990) Dynamic user optimal traffic assignment on congested multidestination networks. *Transpn. Res.,* **24B**, 431–442.

Wunderlich K. (1990) Time-variant travel cost calculation under anticipatory routing. Technical Report 90-5, IVHS Program, University of Michigan, August.