

Peak-Period Control of a Freeway System— Some Theoretical Investigations

JOSEPH A. WATTLEWORTH*, Assistant Research Engineer, Texas Transportation Institute, and
DONALD S. BERRY, Chairman, Department of Civil Engineering, Northwestern University

This paper contains a series of theoretical considerations pertaining to the control of a freeway system during peak traffic periods. It attempts to answer such questions as the following: What are the objectives of operation of a highway system? What type of control technique should be used? What type of traffic detection is required? Where is it to be located? How should the entire freeway system be considered and controlled to produce optimal operation?

An arbitrary street and/or freeway system is analyzed to determine the objective function or goal of operation for the system. An input-output analysis is used.

A theory of flow at bottlenecks is developed to explain the reduction of flow rate at some bottlenecks during congestion while the flow rate at other bottlenecks remains at its capacity level during congestion. This is a macroscopic flow model based on basic continuity equations.

Other macroscopic models of traffic behavior at and upstream of a bottleneck are to determine what traffic variables are to be detected to (a) predict congestion and (b) indicate congestion and how far from the critical sections the detection must be made in order to allow control decisions to be made.

Several criteria are established for control techniques and several control techniques are examined in light of these criteria. The possible role of each in a final control system is also discussed.

Finally, a linear programming model of the operation of a freeway system is presented. This can be used as a descriptive model, but with some modifications could probably be used to provide the control actions required for optimal system operation. Interpretation of the dual variables and a sensitivity analysis are included and these provide many valuable insights into the operation of a freeway system. The linear programming approach suggested a method of determining demand at a certain freeway location.

•PEAK-PERIOD congestion is a frequent occurrence on many urban freeways. This is partly due to partial completion of planned freeway systems. It is doubtful, however, if enough freeways can be built (or should be built) to eliminate completely peak-period congestion. Hence, an operational or control means is needed to provide relief from congestion during the peak periods at each stage of freeway construction as well as for the completed system.

*Formerly with Expressway Surveillance Project, Illinois Division of Highways.

Paper sponsored by Committee on Theory of Traffic Flow and presented at the 43rd Annual Meeting.

There is evidence to indicate that congestion on a freeway can decrease the flow rate on the freeway, whereas congestion on an arterial street does not decrease the flow rates on those streets. For this and other reasons, prevention of congestion on freeways seems to be more important than its prevention on surface streets.

Congestion develops at a bottleneck¹ when the demand exceeds the capacity there. The controls then must increase the capacity of some bottlenecks and/or shift the demand, either spatially or temporally. The development of congestion in one part of the system may have quite different effects on various other parts of the system.

1. If the queue which develops at a bottleneck backs upstream past one or more exit ramps, the output rate at these ramps may be decreased.
2. The unusually early development of congestion at one bottleneck may cause the downstream flow rate to decrease earlier than on a "typical" day. This can either (a) decrease the output rate of the system of interest, or (b) delay or eliminate congestion at a downstream bottleneck, thereby increasing the output rate of the system of interest.

The interdependence of critical locations suggests that a systems analysis would be the only adequate analytical approach to the problem of reduction or elimination of congestion on a freeway system.

Because traffic conditions can change rapidly with time, freeway controls should be traffic-adjusted to respond to traffic conditions in a large area rather than only to conditions at individual locations. Thus, a control system is needed rather than a series of independently operated controls. A control system, properly designed and operated, could result in the optimal performance of at least a portion of the entire system, that is to say, a subsystem.

Before the development of a freeway control system, it is necessary to: (a) determine the objectives of the control, (b) understand traffic behavior so that the responses to the control system can be predicted, (c) understand traffic behavior since it affects detector locations, and (d) develop an analytical means of describing and considering a freeway system or subsystem both for predicting the effects of the control and for operating the controls.

SCOPE

This study was not intended to be the complete development of a new and different control system to be applied to streets and freeways. The intention is, rather, to undertake some of the theoretical investigations on which the final development depends. The controls which are investigated have as their purpose the improvement of traffic conditions on urban freeways. Controls to improve operation of arterial streets are not considered, and only special peak-period freeway controls are considered. Primary attention is given to normal operating conditions on the freeways.

Every control scheme has some philosophy behind it, even though it is usually not explicitly stated. Because only controls for the improvement of peak-period freeway traffic flow are considered, the basic philosophy of these controls is that improving operation on the freeways will result in a net improvement in the operation in the entire system of streets and freeways. There can be little doubt of this if the freeway can be operated optimally without diverting a significant portion of the demand to other parts of the system and if the freeway controls do not otherwise affect the operation of the streets. Hence, an emergent philosophy is that as little as possible of the freeway demand will be diverted to the arterial street system. The problem then becomes one of operating the freeway in an optimal manner with the given demand.

¹ In this paper, a location at which the capacity is lower than the possible upstream flow rate is called a bottleneck. Bottlenecks can be placed in two broad categories, permanent or geometric bottlenecks and temporary bottlenecks such as accidents and disabled vehicles. Two types of freeway operation result from these types of bottlenecks. Under "normal operation," traffic behavior is affected primarily by the geometric features of the freeway; under "reduced capacity operation" the capacity of one or more sections of the freeway has been reduced by a temporary bottleneck.

This type of optimal freeway operation does not assure the optimal operation of the entire street and freeway system. It does, however, assure an improvement. Perhaps the diversion of part of the freeway demand to the streets would improve the operation of the total system of streets and freeways. To determine the optimal control for the total system would require a much larger and more difficult study than is undertaken here.

The first part of this study is the development of a criterion function by which the performance of a freeway control system can be evaluated. This is accomplished by means of some theoretical, input-output considerations of an arbitrary highway system.

Also included is a discussion of the macroscopic behavior of traffic at a bottleneck and at points upstream of and downstream of the bottleneck. The applications and implications of these theoretical discussions to peak-period traffic control as well as to the type and location of traffic sensing or detection that would be required for the proper operation of a freeway control system are also discussed.

Finally, a linear programming model of the operation of a one-directional freeway subsystem is developed. This model presents many insights into the operation and control of the system. With some additional refinements it could possibly be used for control purposes to indicate the input rates at the entrance ramps which would yield optimal system operation for the reduced capacity situation.

In summary, this study is a series of theoretical investigations which are meant to provide some of the foundations upon which a traffic control system can be developed.

OBJECTIVES OF FREEWAY CONTROL

The discussions which follow are centered around the behavior of an arbitrary urban highway system. The arbitrary "system of interest" discussed could be the entire street and freeway system, the entire freeway system or any subsystem which does not violate the assumptions which will be presented. The system of interest, which is simply called the system in the following discussions, is visualized as being cut by a cordon line (or cordon lines) and is a closed system. Because it is a closed system, the cordon line(s) defines the locations of all of the system inputs and outputs.

Continuity Characteristics of Traffic Flow

In a closed system, the flow of traffic at all times satisfies the basic continuity equations. In terms of instantaneous rates, the input rate to the system equals the system output rate plus the rate of storage or accumulation within the system. Expressed mathematically, this becomes $i(t) = o(t) + s(t)$.

When used to describe a given time period, the continuity equation states that the number of vehicles entering the system equals the number leaving the system plus the change in the number within the system. For the time period from t_0 to t_1 , $I(t_1) = O(t_1) + \int_{t_0}^{t_1} s(t)dt$. (See footnote 2.)

By similar reasoning, the number of vehicles in the system at time t equals the number in the system at time t_0 plus the difference between the number entering the system (between times t_0 and t) and the number leaving the system (between times t_0 and t). Expressed mathematically, $S(t) = S_0 + I(t) - O(t)$. (See footnote 3.)

System Travel Time

The number of vehicles in the system yields many interesting insights when considered over a period of time. If this number is known for each time, t , either as a graph

$${}^2I(t) = \int_{t_0}^t i(\tau)d\tau = \text{cumulative input from time } t_0.$$

$$O(t) = \int_{t_0}^t o(\tau)d\tau = \text{cumulative output from time } t_0.$$

$${}^3S_0 = S(t_0) = \text{number of vehicles in system at time } t_0.$$

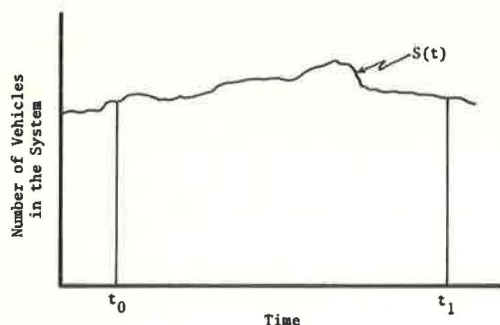


Figure 1. Number of vehicles in the system vs time.

or as a function, $S(t)$, some interesting and important results can be obtained. Figure 1 shows the number of vehicles in a system versus time. The area under this curve between two times, t_0 and t_1 , is the number of vehicle-minutes (if time is in minutes) accumulated in the system during this time. For the time period considered, this is the total "time of occupancy" (1) accruing to all vehicles while they are in the system. If there is no parking in the system, the "time of occupancy" is also the total travel time. The travel time in the system could also have been obtained by integrating $S(t)$ between t_0 and t_1 , i. e., travel time in the system from t_0 to $t_1 = \int_{t_0}^{t_1} S(t)dt$.

From this it can be seen that the total travel time in a given system can theoretically be determined from volume measurements alone. The number of vehicles in the system at the time t_0 (beginning of the period of interest) must be known. From that time until the end of the period an accurate record of all the vehicles entering and leaving the system must be known—according to time increments (such as one minute). For each time period then, the number of vehicles in the system, $S(t)$, is known. This function can be integrated graphically or numerically to yield the total travel time in the system during the period of interest.

The preceding discussions were general and apply to any system or subsystem.

EQUIVALENCY OF MINIMIZING TRAVEL TIME AND MAXIMIZING SYSTEM OUTPUT RATE

The purpose of this section is to determine a criterion function or figure of merit to be used to guide in the operation of a freeway control system. Such a criterion function is essential in that it provides both a goal to strive for and a means of evaluating the operation of the freeway system. To optimize system operation, it is necessary to know what is to be optimized. Ideally, such a criterion function should be easily understandable and amenable to continuous measurement in the system.

Traditionally, travel time has been used to evaluate operation in transportation systems, subsystems and individual links. Because most drivers are trying to go between their origins and destinations in the shortest time, there is little doubt that the total street and freeway system travel time is a good figure of merit in examining operation of the entire system. Subject to the constraints of safety, comfort, convenience, etc., travel time in the system would ideally be minimized.

The use of total travel time in the evaluation of the operation of a subsystem or a link may not be quite so meaningful. In these, the total travel time accumulated in a given time period can be reduced by decreasing the per-vehicle travel time or by decreasing the number of vehicles using the facility during the time period. Even the average travel time on a link can be decreased by decreasing the volume on the link. This can be done by diverting some vehicles but, in this case, the problem may merely be shifted to adjacent links or subsystems. Total travel time is meaningful when no attempt is made to alter the input. Thus, the following discussions are limited to any system or subsystem for which the input rate is unaffected by congestion internal to the system. As pointed out previously, these are the only cases in which travel time comparisons are meaningful.

It was shown previously that the total travel time in a system or subsystem is the integral over time of the number of vehicles in the system. That is to say that, for the period from t_0 to t_2 , the total travel time = $\int_{t_0}^{t_2} S(t)dt$. Since $S(t) = S_0 + I(t) - O(t)$,

the total travel time = $\int_{t_0}^{t_2} [S_0 + I(t) - O(t)]dt = S_0(t_2 - t_0) + \int_{t_0}^{t_2} I(t)dt - \int_{t_0}^{t_2} O(t)dt$.

Since S_0 , t_2 and t_0 are constants, the first term of the travel time expression is also a constant. On a given day, the cumulative input to the system, $I(t)$, is a fixed function of time and is assumed unchanged by any controls which are exerted on the freeway system. Thus, the second term in the travel time expression is also a constant. Hence, the total travel time = a constant - $\int_{t_0}^{t_2} O(t)dt$. This is the total travel time accumulated in the system under consideration in the time period from t_0 to t_2 .

One objective of control or operation of a transportation system is to minimize the total travel time accumulated in the system in a certain critical period of time. Due to the previous considerations, for a fixed input function, minimizing the total travel time is equivalent to maximizing $\int_{t_0}^{t_2} O(t)dt$.

Figure 2 shows a hypothetical cumulative system output function plotted between times t_0 and t_2 . If the time period under consideration is defined properly, there will be no appreciable congestion in the system at either time t_0 or t_2 . The total system output in the time period, then, will be equal to the total demand for the period and will be a fixed value as shown in Figure 2. Since $O(t)$ is the cumulative system output, $\int_{t_0}^{t_2} O(t)dt$ is the area under this output curve in the period t_0 to t_2 . The problem then becomes one of maximizing the area under a curve which passes through zero at t_0 and through a fixed point at t_2 .

There are two major constraints which are placed upon this maximization. There is an upper limit on the slope of the curve because the slope of the curve is the output rate of the system. Hence, the maximum possible slope of the curve is the capacity rate of output. The height of the curve or the cumulative system output can never exceed the cumulative system input by more than the number of vehicles originally in the system.

Within these constraints, the area under this curve would be maximized by, starting at time t_0 , maximizing the slope of the curve at each instant of time. If one were considering discrete time intervals, the object would be to increase the height of the curve by as much as possible in each time interval. The traffic interpretation of this is that in most cases the control strategy of maximizing the output rate at each moment of time (or the output in a given time period) is equivalent to minimizing the total travel time in the system for a fixed-system input function.

Several conclusions can be drawn from these considerations. In most cases in which the system input is not affected by controls or other system changes, maximizing the

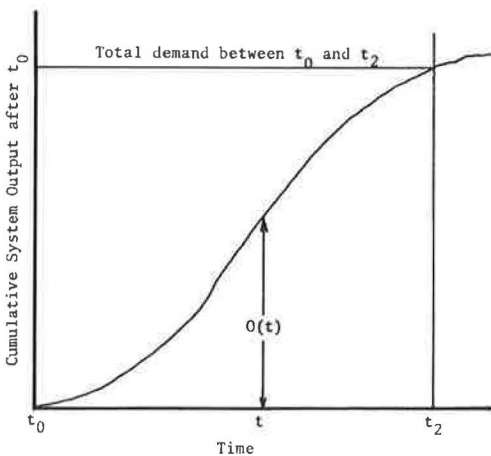


Figure 2. Cumulative system output vs time.

output rate at each moment of time is equivalent to minimizing the travel time in the system. If the controls or other changes cannot increase the output rate at some time during the period of interest, no decrease in travel time can be produced. Also, what happens within the system is of importance to total travel time only as it changes the output characteristics. Hence, in an analysis of this type, such items as speed are of importance (to system travel time) due only to their effect on the output rate of the system.

So far, all discussions have been made for large systems. Most of the same discussions hold for smaller subsystems or individual facilities, but a little more care must be exercised so that some of the assumptions are not violated. In some instances it is possible that an increase in

the output rate of a freeway subsystem (a one-directional length of freeway with some on-and-off ramps) can cause downstream congestion which will restrict the output rate later, resulting in an increase in system travel time. (This is partly a problem of proper system selection.)

It is also possible that the optimal method of operation of the freeway and surrounding arterial streets is to change the freeway input rate by diverting some traffic. This paper concerns itself primarily with determining the optimal operation of a freeway subsystem for fixed inputs. The effect of this operation on the arterial streets is not considered. The philosophy which has been adopted in this paper is to accept the demand for freeway use as it occurs (or at least to store the excess demand for a period of time until it can be allowed on the freeway without causing congestion) and operate the freeway system optimally for that demand. This results in a minimum total travel time for all vehicles which demand use of the freeway in the time period being considered. It is not necessarily the optimal operating procedure when the surface street operation is considered along with the freeway operation. However, optimal operation of the freeways will almost certainly produce a substantially improved operation of the overall system.

EFFECT OF CONGESTION

There is considerable evidence which indicates that the development of congestion at a permanent freeway bottleneck can decrease the flow rate there (2 through 5). Since congestion can decrease the output rate of a freeway system and since one objective of freeway operation is to maximize the system output rate, another goal of freeway operation is to prevent congestion on the freeway. This is consistent with the overall objective of minimizing system travel time.

TYPE OF CONTROL

Criteria for Controls

The controls must be able to prevent or alleviate congestion in freeway system or subsystem without causing inefficiently low flows on the freeway. They must be flexible enough to respond to traffic conditions in the system. They must also be quite positive and firm so that the desired traffic behavior can be obtained and so that a given result can definitely be associated with a given control action. The controls must also be acceptable to the drivers, must not be hazardous and must fall within the limits of economic feasibility.

Entrance Ramp Metering

After considering many types of special peak-period freeway controls, ramp metering seemed to best meet the criteria which were established for such controls. Each of the other types of controls had some distinct advantages in some particular situations, but in the general situation, ramp metering seemed to hold the most promise.

Entrance ramp metering is a system by which the maximum flow rate at each entrance ramp is set based on freeway traffic conditions. This is done by releasing vehicles from the ramp to the freeway at the chosen time headways. In this way it is at least theoretically possible to maintain but not exceed the merging capacity. Ramp metering has been tested by the Congress Expressway Surveillance Project (6).

The possible benefits of metering are as follows:

1. Reduction of freeway congestion,
2. Increase of merging capacities,
3. Making merging maneuver easier for ramp vehicles, and
4. Diversion of some short trips from the freeways due to the time delay caused by the metering.

One potentially serious problem is the storage of queued vehicles on the ramps. This problem must be given a great deal of consideration in the design of a metering system.

TRAFFIC DETECTION IN A FREEWAY CONTROL SYSTEM

Purposes of Detection

Prediction of Congestion.—The first purpose of the traffic detection for a peak-period control system is to predict traffic conditions that will occur at bottleneck locations while there is still time to take corrective action. The detection system must have the capability of predicting the development of congestion at the bottlenecks so that the controls can be applied to prevent this congestion. The lead time of prediction must be sufficient to allow the ramp metering to respond so as to prevent the congestion.

Indication of Congestion.—Even with a peak-period control system in operation congestion will, at times, develop in the freeway system due to accidents or other unusual events. Congestion must be detected so that remedial controls can be initiated at once to minimize the effect of the unusual event. Thus, the second purpose of the detection system is that of providing an indication of congestion.

Variables to Be Detected.—An understanding of the behavior of each of the numerous freeway traffic variables (volume rate, speed, density, lane occupancy, etc.) prior to and during the peak period is important in the design of the detection system. One or more of these variables must be detected in order to accomplish the two objectives: (a) predicting traffic conditions at critical sections, and (b) indicating congestion.

The portion of freeway shown in Figure 3 is used as a framework for discussion, and a descriptive model of the behavior of the traffic variables in this zone is presented. Section B is the bottleneck or critical section and section U is assumed to be one-half mile upstream of section B. The assumed volume-density curves for these two sections are shown in Figure 4.

The behavior of three variables, volume rate (q), speed and density, is examined. These variables were chosen partially because of this interrelationship, $q = V \times k$. Lane occupancy behaves very much like density and, in fact, it could be considered a time-based density. Because of the similarity of behavior of lane occupancy and density, only density is considered.

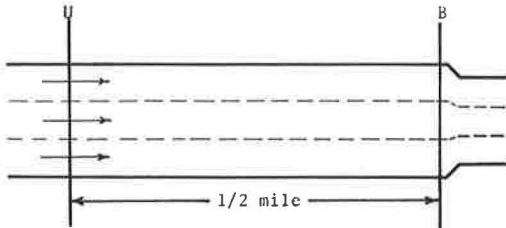


Figure 3. Freeway with bottleneck section.

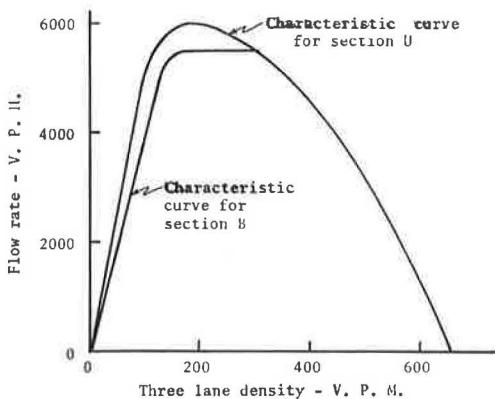
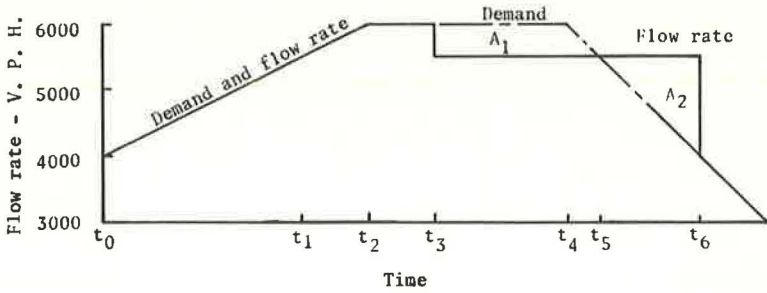


Figure 4. Flow rate-density curves for sections B and U.

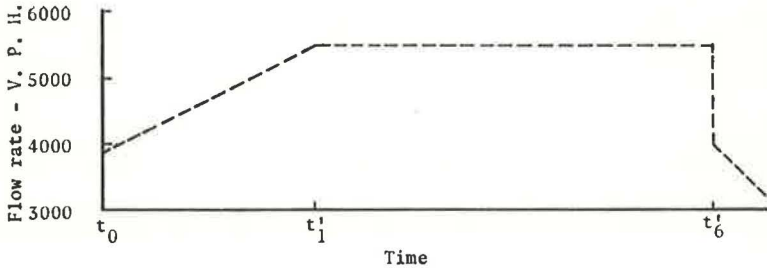
Figure 5a shows the assumed demand rate at section U. The following discussions do not depend on the exact shape of this curve so the somewhat linearized demand curve does not alter the results.

Because of the volume-density curves which were chosen (Fig. 4), the bottleneck capacity is assumed to be 5,500 veh/hr, while the capacity at section U is 6,000 veh/hr. It is also assumed that the density during congestion equals 300 veh/mi so the bottleneck flow rate does not in this case (by assumption) decrease due to congestion. This assumption was made to simplify the example. A similar, somewhat more difficult, analysis could be made in which the density during congestion could exceed 300 veh/mi thereby decreasing the flow rate at the bottleneck. For the purposes for which this analysis is used, however, the assumption that the density during congestion equals 300 veh/mi does not alter the results, mainly because the period of primary interest is that prior to congestion.

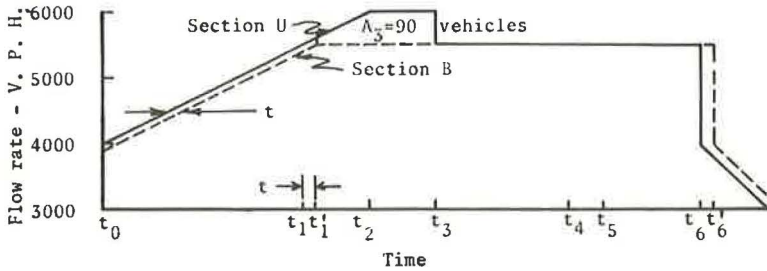
At time t_0 the volume rate at section U equals the demand rate of 4,000 veh/hr while the speed is 50 mph and the density is 80 veh/mi. The volume rate is also



a. Flow rate at section U versus time.



b. Flow rate at section B versus time.



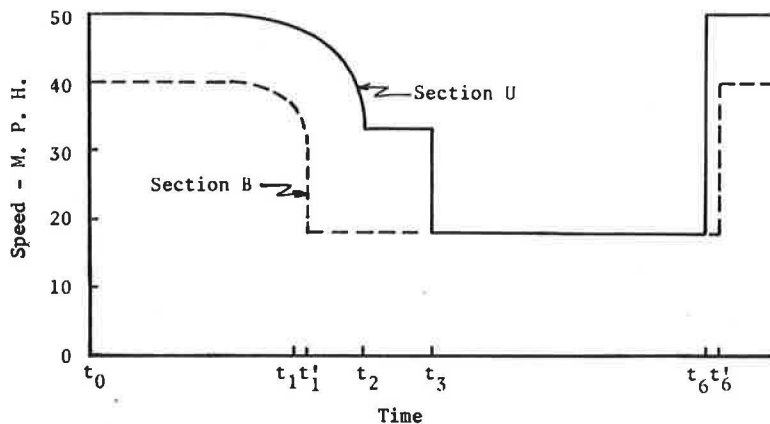
c. Flow rates at sections U and B versus time.

Figure 5. Flow rates at two freeway sections vs time.

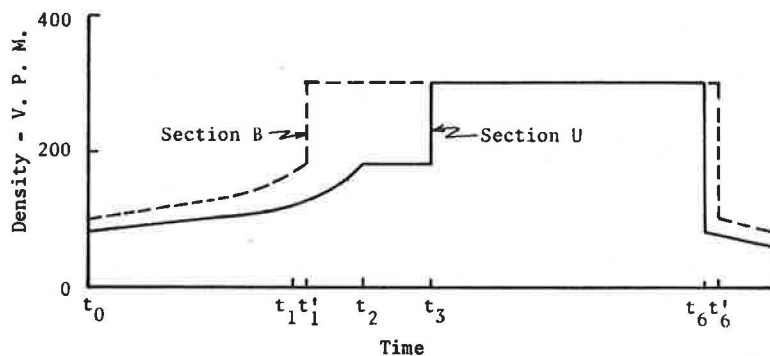
increasing because the peak period (it is assumed) is approaching. (Short-time variations in flow rate are not considered, only the overall volume rate trend which is increasing.) Due to the distance ($\frac{1}{2}$ mi) between sections U and B, there is a time lag between an increase (or decrease) at B. During the volume buildup prior to congestion, the volume rate curve at B will parallel the corresponding curve for section U but will lag by a time t (Fig. 5c). At time $t_0 + t$ the volume rate at B equals 4,000 veh/hr, the speed equals 40 mph and the density equals 100 veh/mi. Figure 6a and b shows the speeds and density plots at the two locations.

At time t_1 the volume rate at U reaches 5,500 veh/hr, the bottleneck capacity. The corresponding flow rate at section B is slightly lower than this (Figure 5a and c). The speeds at U and B are, respectively, 48 and 37 mph (Figure 6a) while the densities are 120 and 170 veh/mi, respectively, at U and B.

A short time after t_1 , the flow at U exceeds 5,500 veh/hr. It was seen previously that the input rate minus the output rate equals the storage rate. In this case if the input exceeds the output capacity there is certain to be a positive storage rate.



a. Speeds at sections B and U versus time.



b. Density at sections B and U versus time.

Figure 6. Speed and density at sections B and U vs time.

At time t'_1 ($t_1 + t$) the volume rate at B reaches the capacity flow rate there. At this time storage or queuing occurs upstream of B and the density at B climbs to the (assumed) congested density of 300 veh/mi. The corresponding (Fig. 6a) speed at B decreases to 18 mph. The behavior of the speed and density at time t'_1 at the bottleneck can be seen in Figure 6a and b (dashed lines). At this time, which is the beginning of congestion at B, the speeds drop sharply and the density rises sharply. The volume rate, however, remains constant after time t'_1 . Thus, a measurement of volume alone would be a poor indicator of congestion, whereas speed and density are quite sensitive to congestion. Since lane occupancy is similar to density, it too, would be very sensitive to congestion. Hence, speed, density or lane occupancy could be used as indicators of congestion.

Due to the time lag between vehicles passing an upstream section and the same vehicles passing a bottleneck, the traffic behavior at the upstream section should provide predictive information on what is going to occur at the bottleneck. (To keep this description relatively simple, a constant time lag, t , is used; in reality, however, the time lag would vary, depending on the speeds in zone UB.)

It has been shown that when the volume rate at U exceeded the capacity flow rate at B, storage had to take place. Also, the excess of flow rate at U over the capacity flow rate at B led to the formation of the queue at B which caused the density to increase sharply and the speed to decrease sharply there. Thus, this congestion was caused by

the flow rate at U exceeding the capacity flow rate at B. The volume rate at U exceeded 5,500 veh/hr for a time period t prior to the development of congestion at B. Thus the upstream flow rate, when related to the bottleneck capacity, can be used to predict the development of congestion at the bottleneck. High volume rate can be considered the cause of congestion, whereas high density and low speed can be considered the effects of congestion.

At t_1 there are also certain values of speed (48 mph) and density (120 veh/mi) at section U which correspond to the 5,500-veh/hr volume rate. It could be argued that an average speed of less than 48 mph or a density greater than 120 veh/mi could also be used to predict congestion at B. However, in this case speed and density are used to predict volume and would have no significance in themselves. A speed or density at U cannot be specifically related to a present or future speed or density anywhere else. The continuity equations hold true only for volumes. Indeed, Barker (7) studied the propagation of discontinuities of volume, speed and density and found that under free-flow conditions discontinuities in all three variables were propagated downstream. These discontinuities were followed through a series of detector stations about 400 ft apart. For a given variable, the downstream propagation of "waves" or discontinuities is a necessary but not a sufficient condition for use as a predictor of congestion at a downstream location. It is also necessary that the variable satisfy the continuity equations. These equations can be written for volume (or volume rate) for a closed system, namely the input equals the output plus storage. For a "straight pipe" length of free-way (such as in Fig. 3) this means the volume past the upstream section equals the volume past the downstream section plus the additional number of vehicles between the sections. Similar equations cannot be written for speed, density or lane occupancy.

Since congestion develops at a bottleneck when the upstream flow rate exceeds the bottleneck capacity, it seems even more logical to use the upstream flow rate to predict congestion. Capacity values of speed and density of lane occupancy have not yet been developed.

For the preceding discussions, one-to-one transformations between speed or density and volume rate were assumed—meaning that a given speed or density has one and only one corresponding volume rate. In reality, however, a range of volume rates would be associated with a fixed average speed or density and this would make prediction of volume rate from other variables more difficult. The relationships between volume rate and speed or density could also change from one upstream location to another, further contributing to the problem of predicting volume rate from the other variables.

It is also interesting that noncongested flow prevailed at the upstream location U after the bottleneck (section B) flow became congested. The time during which this holds true equals the time for the rear of the queue to reach the upstream location. The speed of the rear of the queue is a function of (a) the change in density from noncongested to congested operation, and (b) the storage rate. The storage rate equals the flow rate at U minus the flow rate at B. Hence, the speed of the rear of the queue equals $(q_U - q_B)/(k_B - k_U)$. On the volume rate-density curve, this would be the slope of a vector drawn between the operating points for sections U and B, as shown by Lighthill and Whitham (8).

If the density in zone UB were 120 veh/mi at time t'_1 , a total storage of 90 vehicles would be required to place the rear of the queue at U. This is because it was already assumed that the steady-state congested density is 300 veh/mi and zone UB has $\frac{1}{2}$ -mi length. Thus, the time T between the development of congestion at sections B and U is such that $\int_{t'_1}^{t'_1 + T} (q_U - 5,500)dt = 90$. This is the area A_3 in Figure 5c. The queue reaches U at time t_3 ($t_3 = t'_1 + T$).

At time t_3 zone UB is in a steady-state condition. That is, the density in zone UB is a constant 300 veh/mi. This requires that $q_U = q_B = 5,500$ veh/hr. Thus at t_3 the volume rate at U drops sharply from 6,000 to 5,500 veh/hr, the speeds drop sharply from 33 to 18 mph and the density increases from 180 to 300 veh/mi. At time t_6 all the congestion clears at U and at time t'_6 it clears at B.

This analysis assumed that the flow rate at the bottleneck did not decrease due to congestion. If this assumption were not made it would mean that the steady-state con-

gested operating conditions could be found in Figure 4 at the point corresponding to the congested density. Speeds and volume rates during congestion would be lower since the density is now greater than 300 veh/mi. Since the rear of the queue would travel upstream faster, the time required to travel the $\frac{1}{2}$ -mi distance would be less than T . All of these differences occur after congestion has set in and have no effect on the conclusions regarding the prediction of congestion.

This descriptive model was used to explain the behavior of three traffic variables—volume, speed and density—prior to and during congestion. Since it is a descriptive model, several assumptions can be made to facilitate the discussions. Probably the most important of these is the assumption of a known, fixed bottleneck capacity. Many things can change the capacity of a bottleneck or cause a bottleneck where none previously existed. The effect on the capacity of a freeway roadway of a disabled vehicle, an accident, adverse weather or other factor is an area needing a great deal of research. Perhaps, there is no one fixed capacity at a given location but rather there may be a probability of congestion developing due to a given flow rate. This probability may also change with different drivers, weather, etc.

The detection system which has been described is, perhaps, somewhat idealized and no consideration has been given to the economy of such a system. Such considerations might require that some compromises be made. For example, it may be necessary to sample a variable, such as volume, in one lane to estimate the variable across all lanes. It may also be necessary to estimate volume from speed, density or lane occupancy. In this way a one-variable system could possibly be developed. The accuracy of this sampling and estimating one variable from another should be carefully examined to determine if such procedures can be used successfully.

LOCATIONS OF DETECTORS FOR PREDICTION OF TRAFFIC CONDITIONS AT A SECTION

The traffic detection system must be capable of determining the proper volume rates to be allowed to enter at each entrance ramp in order to keep the freeway system operating in the best possible manner under existing conditions—either normal or reduced capacity operation. The volume rate on each ramp can be determined by the capacity of the merging section or by a downstream bottleneck. In either case, the detection of the upstream freeway volume can be used to determine the maximum allowable ramp volume. The metering rates at the ramps can then be set to allow no more than the predetermined rates of flow on the ramps.

The preceding section discussed a time lag between a certain volume passing an upstream location and its passing the bottleneck location. The purpose of this section is to examine this lag time with respect to other critical times of detection and control. The situation considered is that of a metered entrance ramp. Vehicles are detected

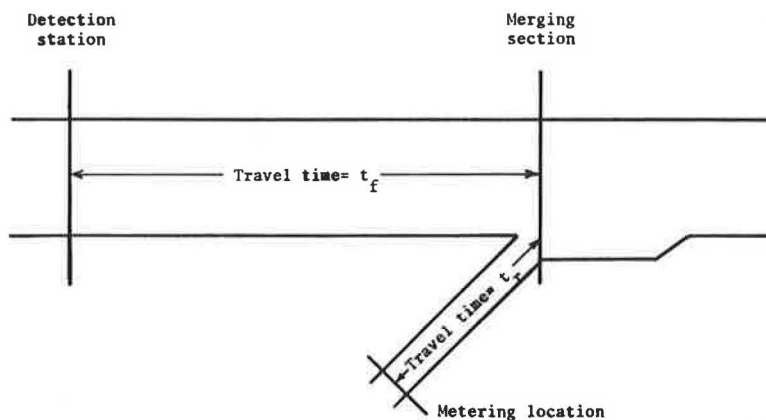


Figure 7. Schematic of metering and detection locations.

In all probability control decisions would not be made on instantaneous happenings at any detector location. It is more likely that the results of detection in a length of time would be used as a basis for control decisions. Short observation times would probably produce quite erratic results (for example, for a 1-sec time period the detection of no vehicles corresponds to a 0 veh/hr flow rate while the detection of one vehicle corresponds to 3,600 veh/hr). Lengthening the observation time or averaging over a longer time would damp out much of the random fluctuation.

In the third case considered, the effect on the required freeway travel time of the time period, t_d , is used for detection or averaging of detector data. Figure 8c shows the time relationships in this case. A time t_d is required from the beginning to the end of the detection period. When the detection period is over a time t_c is needed to assimilate the data, make a control decision, and initiate the control. After a time t_r the first ramp vehicles released under the new control reach the merging section. They should be merging with the first vehicles to be detected on the freeway at a time t_f earlier. Hence, in this case $t_f = t_d + t_c + t_r$. For example, if the detection and averaging time is 45 sec, computation time is 15 sec and ramp travel time is 10 sec, t_f has to be greater than or equal to 70 sec.

It was seen that extra time required between the start of detection and the adjustment of the controls produces an increase in the required freeway travel time between the detection station and the merging section. Other time requirements would similarly increase t_f .

So far the considerations of various time requirements have been used to determine t_f . A fixed t_f or an upper limit on t_f can also be used to establish limits on other times. For example, if $t_f = 45$ sec and if $t_r = 10$ sec, $t_d + t_c$ must be less than or equal to 35 seconds.

Because congestion develops at a bottleneck and is propagated upstream, the ideal location for prompt detection of congestion would be at or slightly upstream of a bottleneck. As the distance upstream of the bottleneck increases the time lag between the development and detection of congestion also increases. Again this location immediately upstream of the bottleneck is the ideal location of the detectors. For reasons of economy it may be necessary to use the same detectors for prediction of traffic behavior and for detection of congestion. In some cases, one detector station can be located at one bottleneck to detect congestion and could also be used to obtain volume data for predicting traffic behavior at a downstream location.

In summary, it appears that short-period volumes upstream of a bottleneck provide the best prediction of impending congestion and are probably the best variable to measure for control purposes. Volume alone cannot be used to differentiate between congested and noncongested operation. Hence, speed, density or lane occupancy must also be measured at or upstream of the bottleneck in order to provide this information.

It is fortunate that volume is one variable to be measured for control purposes. It was previously indicated that maximizing the system output volume rate will lead to optimal operation of the system. Volume is also the only variable for which continuity equations can be written. Volume measurements are susceptible to point measurements that are much easier to accomplish than measurements over a length of roadway. The volume rate is also the most easily controllable variable and volume capacities can be established. Individual time headways at the input sources can be controlled and there is a one-to-one relationship between time headways and volume rate. Other advantages will appear later.

OPTIMAL OPERATION OF A FREEWAY SUBSYSTEM

The subsystem considered here includes storage areas for vehicles waiting to enter the freeway. Several entrance and exit ramps increase and decrease the volumes along the freeway. Figure 9 is a schematic of the freeway system (a portion of the westbound Congress Street Expressway in the Chicago area) which is used for the development of a prototype model. Four lanes of traffic enter at the Cicero Ave. end of the system and three lanes exit downstream of the Des Plaines Ave. entrance ramp. The demand at the (Cicero Ave.) freeway input source is not controlled, and thus, is one further re-

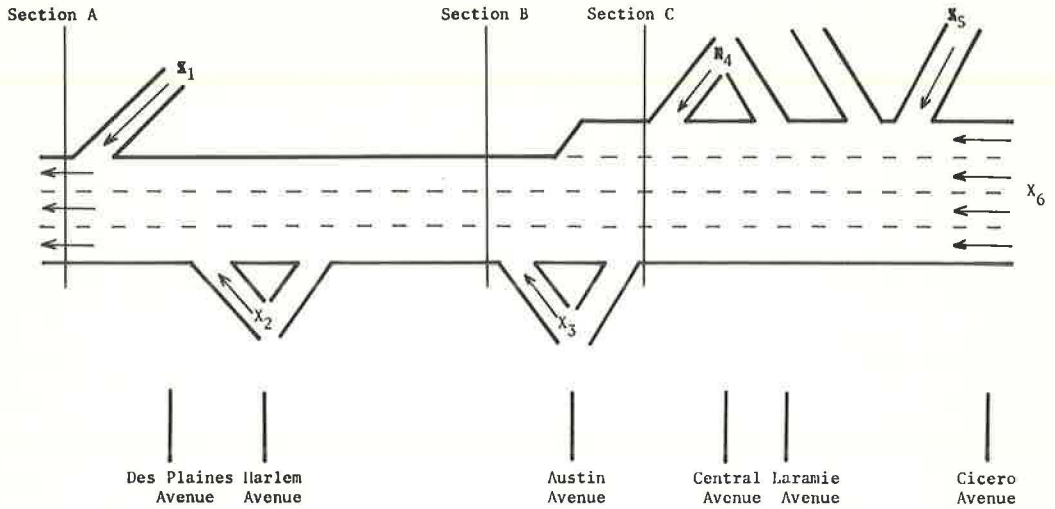


Figure 9. Schematic of freeway subsystem used in the development of the prototype linear programming model.

striction to the problem. Generally such a restriction will not affect the optimality of the operation, however. As used in this section, input refers to flows onto the freeway roadway and is a different use of the term than was made previously.

The controls which are considered are entrance ramp metering controls which can limit the inputs to the freeway from the various entrance ramps. The mathematical model yields the volume which is required or permitted on each of these ramps during the analysis period in order to obtain optimum performance of the portion of the freeway system under consideration. Metering devices are assumed to be in operation on the ramps to limit the ramp flow rate to the required level. Besides limiting the ramp flow rate, the metering system also serves to damp out large variations in demand on the entrance ramps. Hopefully, this will lead to smoother merging operations in the vicinity of the ramp and perhaps to higher merging capacity rates.

The objective function which was selected for optimization is the output of the system in a given time period. This is, of course, to be maximized. Since it is assumed that congestion can decrease the flow rate at a bottleneck, the controls must be operated so as to prevent the development of congestion at all bottleneck locations in the system in order to keep the output rate at its maximum level. Hence, critical points or potential bottlenecks in the section must be identified so that the demand on these sections can be kept below capacity levels.

The flow rate upstream of each bottleneck section must be so controlled as to prevent the development of congestion. The model which is developed tells how much flow from each source can be accommodated under these conditions during the analysis period. (The actual ramp flow rates would not be constant during the period.) The remainder is the amount which must be diverted or stored until the time period is over. The entire demand from the freeway input source (in this case at the Cicero Ave. four-lane section) is accepted into the system.

Origin-Destination Information

Since the volume at each critical or bottleneck section is composed of vehicles from several origins, the effect of altering the flow at one or more of the origins must be known. Therefore, for each input the percentage of its inflow vehicles crossing each critical section must be known.

Existing O-D data (9) for each entrance ramp provided part of this information. For the peak period at each entrance ramp the percent of entering vehicles destined for each of the exit ramps was available for the system of interest. These data are shown in lines 1 to 5 of Table 1.

TABLE 1
ORIGIN-DESTINATION DATA FOR THE FREEWAY SUBSYSTEM
(FIG. 9)

Input Source	% Destined for:				
	Laramie	Central	Austin	Harlem	Through
1. Des Plaines	-	-	-	-	100.0
2. Harlem	-	-	-	-	100.0
3. Austin	-	-	-	5.1	94.9
4. Central	-	-	-	6.7	93.3
5. Cicero	0.9	2.2	4.7	9.8	82.4
6. Cicero	13.7	8.6	15.8	10.0	51.9

When combined with volume counts of vehicles entering and leaving the system, both on the mainline (the freeway roadway) and on ramps, these data were used to determine the destinations of the vehicles entering at Cicero Ave. on the freeway. At each output location, the total output volume during the period was known. From the input volume data for each location and the knowledge of the percent of this volume destined for each output location, it was possible to estimate, for each output, the volume coming from all the input sources except the mainline input. When this subtracted from the total volume at a particular output, the number of vehicles at the output coming from the mainline input is determined. For example, assume that the Laramie Ave. exit ramp volume (for the period considered) was 700 vehicles and the Cicero Ave. entrance ramp volume was 1,000 vehicles. It is known that 0.9 percent of the vehicles which enter the freeway from the Cicero on-ramp exit at Laramie Avenue. This means that the expected number of vehicles entering at the Cicero ramp and exiting at Laramie during the period is 9. Thus the other 691 exiting vehicles must have come from the freeway input. Since the total freeway input is known, the percent of this volume leaving at Laramie can be computed. The calculated percents of vehicles entering on the freeway at Cicero and destined for each of the outputs are shown in line 6 of Table 1.

Deterministic Linear Programming Model

The deterministic linear programming model discussed here optimizes the operation of a freeway subsystem or system subject to several constraints. The period considered in the optimization is one hour and it is selected for several reasons:

1. The ends of the section are temporally separated by approximately 6 to 8 minutes (the travel time from one end to the other); hence, extremely short time periods are not satisfactory.
2. It is a convenient time period for many data measurements.
3. It is a period that is shorter than the period of congestion in the specific instance under consideration.

Even at locations at which the congestion lasts for a shorter time, demand can be conveniently expressed as hourly rates.

All volumes, capacities, and demands used in the model are for a 1-hr period. This assumption can be, as it turns out, quite revealing. If the capacity of a critical location that now regularly experiences congestion is determined and the corresponding capacity restraint is not exceeded in the linear programming model, it simply means that metering the flows would have prevented congestion at this location and that the entire ramp demand would be satisfied within the hour. In other words, if the hourly capacity constraint is not exceeded, the congestion is caused by short-time surges of traffic or a downstream bottleneck.

Objective Function.—The model has as its objective the maximization of the output of the system of interest in the time period considered. The output is considered to be the volume leaving the system via the freeway mainline output and all of the exit ramps. The number of vehicles entering a closed freeway system during some time period equals the number leaving plus the number stored in the system during the same time. Thus, the input to the system equals the output plus storage. The storage rate can be positive or negative depending on whether vehicles are entering or leaving storage. As congestion develops and as congestion diminishes the absolute value of the storage rate increases. During steady-state conditions, either in free flow or congestion, the change in storage approaches zero, i.e., the number of vehicles in the system over the time period remains relatively constant. Since the prevention of congestion is incorporated into the linear programming model, the change in storage is considered to be zero. Hence, the input of the system is equal to its output, so input can be substituted for output in the criterion function. The objective of the model, then, is to maximize the input to the system. This can be interpreted on an intuitive basis; the model is maximizing the number of vehicles which can enter the system without encountering congestion before leaving. The variables of the model are the volumes at the input sources.

[It will be noticed that two different systems have been discussed so far. The first is the system in which the travel time is being minimized. This is the system consisting of the freeway and the areas where vehicles wait to be allowed to enter the freeway. The second system, which is considered in the linear programming model, consists of only the freeway and the entrance ramps up to the metering devices; thus, it does not include the waiting areas. If we call the first system the larger system and the second system the smaller system, the two are related as follows. (See Fig. 10.) We have seen that minimizing the travel time in the larger system is equivalent to maximizing the output (rate) in the larger system. Since the outputs of the larger and smaller systems are identical, this is equivalent to maximizing the output of the smaller system. If no congestion develops in the smaller system the input equals (very nearly) its output. Hence, maximizing the input (rate) to the smaller system without causing congestion is equivalent to minimizing the total travel time in the larger system.]

N - Metering Device

W - Waiting Areas - in Large System Only

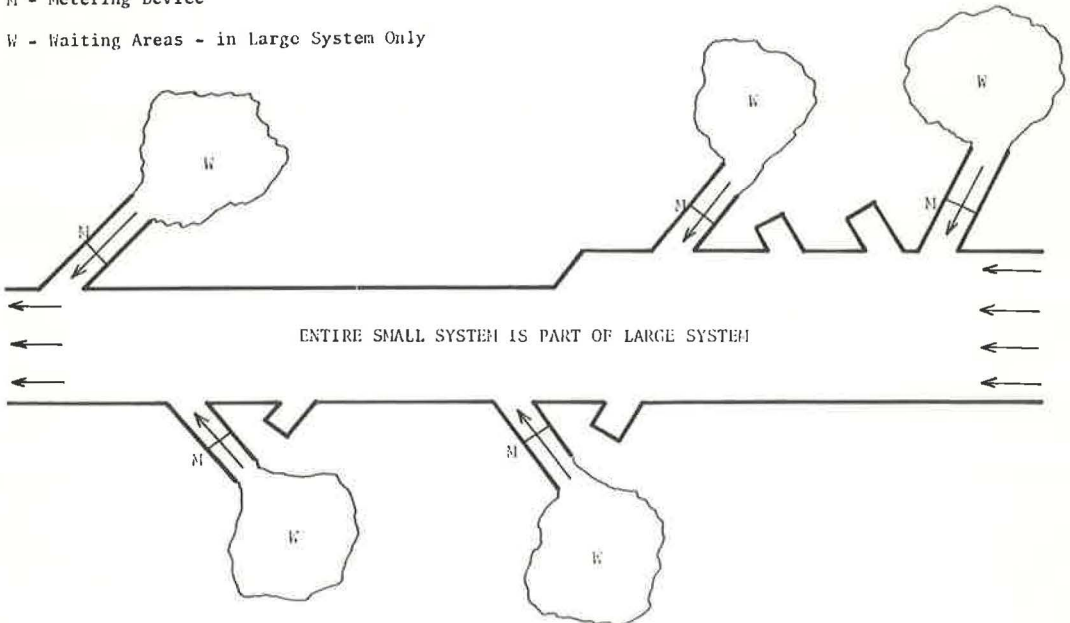


Figure 10. Two systems considered in the analyses.

Constraints.—The model has two types of restraints. Since the development of congestion at many freeway sections reduces the flow rate at these sections, one of the restraints is that congestion will not be allowed to develop at any location on the freeway. Alternately stated, the flow upstream of all critical sections on the freeway will be constrained so that the capacity flow rates at these critical locations will not be exceeded. In applying the model for normal operating conditions, constraints have to be established for only those locations which are, a priori, known to be likely sources of congestion.

The second restraint is that there are certain upper limits placed on the input volumes, because there is a limited demand or number of people desiring to use each ramp during any time period. This demand could, however, increase or decrease when a metering system is put into operation. If travel times on the freeway were lowered, the freeway would become more attractive to some motorists. This generated traffic must enter the freeway somewhere, so some ramp volumes could increase. However, some of the increased traffic on a given ramp might have formerly used another ramp in the system. It is also possible that many vehicles will be diverted from the freeway altogether, because of the delays at the entrance ramps which are caused by the metering operation. Thus an increase or a decrease in a ramp's volume is possible. In this model, it is assumed that the maximum demand is known for all ramp inputs as well as for the freeway input.

Statement of Model.—The deterministic model, then, yields the volume at each input source which maximizes the total input to the system subject to two types of constraints. First, a set of constraint equations is required to assure that congestion will not develop at any location. A second set of constraint equations restricts the inputs from each source so as not to exceed the demand at the source.

Development of Prototype of Deterministic Model

Objective Function.—The variables of the linear programming problem are the input volumes from each of the input sources. The variable corresponding to each input is shown in Table 2. The objective is to maximize $X_1 + X_2 + X_3 + X_4 + X_5 + X_6$.

Constraints.—The first set of constraint inequalities, those which require that demand not exceed the free-flow (or possible) capacity at each critical location on the freeway, utilizes the data from Table 1. It is necessary to know the percent of each input volume that will appear at each bottleneck location on the freeway. For example, the capacity at section A is 5,900 veh/hr, so the total demand at this section must be kept lower than 5,900 vehicles for the hour. The percent of vehicles crossing this section from each input is listed in the "through" column (Table 1). This volume is the freeway mainline output. One hundred percent of the Des Plaines and Harlem ramp traffic crosses this section, while 94.9 percent of the Austin ramp traffic, 93.3 percent of the Central traffic, etc., pass through this bottleneck. Expressed decimally, the first constraint is $1.00 X_1 + 1.00 X_2 + 0.949 X_3 + 0.933 X_4 + 0.824 X_5 + 0.519 X_6 \leq 5,900$. This assures that congestion will not develop at the Des Plaines Avenue entrance ramp merging section.

Another potential bottleneck location is the merge of the Austin Ave. entrance ramp (section B in Fig. 9); the capacity here is 6,000 veh/hr. All of the Austin and Central entrance ramp traffic crosses this section, while only the portion of the Cicero ramp and freeway input traffic which does not exit at Laramie, Central or Austin would pass section B. For the freeway input traffic, 38.1 percent leaves the freeway without reaching section B (13.7 percent at Laramie, 8.6 percent at Central, and 15.8 percent at Austin). The remaining 61.9 percent of this traffic

TABLE 2
VARIABLES AND CORRESPONDING
INPUT VOLUMES

Variable	Represents Input Volume at
X_1	Des Plaines entrance ramp
X_2	Harlem entrance ramp
X_3	Austin entrance ramp
X_4	Central entrance ramp
X_5	Cicero entrance ramp
X_6	Cicero mainline

TABLE 3
MAXIMUM HOURLY DEMAND
AT EACH INPUT

Input	Variable	Max. Hourly Demand
Des Plaines entrance ramp	X ₁	600
Harlem entrance ramp	X ₂	475
Austin entrance ramp	X ₃	450
Central entrance ramp	X ₄	500
Cicero entrance ramp	X ₅	825
Cicero mainline	X ₆	6,800

(0.619 X₆) crosses section B. Similarly, 92.2 percent of the Cicero ramp traffic (0.922 X₅) passes this bottleneck location. The constraint for the Austin Avenue merge is 1.00 X₃ + 1.00 X₄ + 0.922 X₅ + 0.619 X₆ ≤ 6,000.

The third and final bottleneck in the freeway system is the approach to the Austin Ave. exit ramp. This location is critical because of the transition from four lanes to three lanes. The capacity of this location (section C) is 6,450 veh/hr, so the constraint for section C is 1.00 X₄ + 0.969 X₅ + 0.777 X₆ ≤ 6,450.

These three constraints, when met, assure that congestion will not develop at any of the three bottleneck locations.

Another assumption was implicitly made in the formulation of these equations: there is a one-to-one tradeoff between ramp vehicles and freeway vehicles; that is, a ramp vehicle "uses only as much capacity" as a vehicle already on the freeway. If it

is determined that this is not true (i.e., that the reduction of ramp traffic by one vehicle would allow more than one vehicle increase on the freeway) this can easily be put into the model provided the tradeoff is a constant. In the three constraints discussed so far, unity has been the coefficient of the variable corresponding to the traffic on the merging ramp. This would have to be changed to the correct value (the number of additional freeway vehicles which can be passed due to a one-vehicle decrease in ramp traffic) in the constraints. The capacity would also have to be increased to compensate for this change since the ramp volume is weighted more heavily. If 1.5 is found to be the correct coefficient, the second restraint (Austin entrance ramp merge) would be re-written 1.5 X₃ + 1.00 X₄ + 0.922 X₅ + 0.619 X₆ = C, where C is the modified capacity.

Table 3 gives the values of the maximum hourly demands at each input location. The constraint inequalities which prevent an input from exceeding the demand are

$$\begin{aligned}
 X_1 &\leq 600 \\
 X_2 &\leq 475 \\
 X_3 &\leq 450 \\
 X_4 &\leq 500 \\
 X_5 &\leq 825 \\
 X_6 &\leq 6,800
 \end{aligned}$$

Statement of the Prototype Model.—The prototype model, which maximizes the input to the system subject to constraints which (a) assure that congestion will not develop, and (b) assure that ramp volumes do not exceed the demand, is stated as follows:

$$\begin{aligned}
 &\text{Maximize } X_1 + X_2 + X_3 + X_4 + X_5 + X_6 \\
 &\text{Subject to } X_1 + X_2 + 0.949 X_3 + 0.933 X_4 + 0.824 X_5 + 0.519 X_6 \leq 5,900 \\
 &\qquad\qquad\qquad X_3 + X_4 + 0.922 X_5 + 0.619 X_6 \leq 6,000 \\
 &\qquad\qquad\qquad X_4 + 0.969 X_5 + 0.777 X_6 \leq 6,450 \\
 &\qquad\qquad\qquad X_1 \qquad\qquad\qquad \leq 600 \\
 &\qquad\qquad\qquad X_2 \qquad\qquad\qquad \leq 475 \\
 &\qquad\qquad\qquad X_3 \qquad\qquad\qquad \leq 450 \\
 &\qquad\qquad\qquad X_4 \qquad\qquad\qquad \leq 500 \\
 &\qquad\qquad\qquad X_5 \qquad\qquad\qquad \leq 825 \\
 &\qquad\qquad\qquad X_6 \leq 6,800
 \end{aligned}$$

This assumes a one-to-one tradeoff between ramp and freeway vehicle in the merge constraints.

TABLE 4
ORIGINAL SIMPLEX TABLEAU

C_j		1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	b_i
C_i	X_j	X_1	X_2	X_3	X_4	X_5	X_6	S_1	S_2	S_3	S_4	S_5	S_6	S_7	S_8	S_9	
0	S_1	1	1	.949	.933	.824	.519	1									5900
0	S_2			1	1	.922	.619		1								6000
0	S_3				1	.969	.777			1							6450
0	S_4	1									1						600
0	S_5		1									1					475
0	S_6			1									1				450
0	S_7				1									1			500
0	S_8					1									1		825
0	S_9						1									1	6800
$Z_j - C_j$		-1	-1	-1	-1	-1	-1	0	0	0	0	0	0	0	0	0	$Z_0 = 0$

TABLE 5
OPTIMAL SIMPLEX TABLEAU

C_j		1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	
C_i	X_j	X_1	X_2	X_3	X_4	X_5	X_6	S_1	S_2	S_3	S_4	S_5	S_6	S_7	S_8	S_9	
1	X_1	1						1		-.933		-1	-.949		.080	.206	447
0	S_2								1	-1			-1		.047	.158	213
1	X_4				1					1					-.969	-.777	367
0	S_4							-1		.933	1	1	.949		-.080	-.206	153
1	X_2		1									1					475
1	X_3			1									1				450
0	S_7									-1				1	.969	.777	133
1	X_5					1										1	825
1	X_6						1										6800
$Z_j - C_j$		0	0	0	0	0	0	1	0	.067	0	0	.051	0	.111	.429	$Z_0 = 9364$

Computational Solution of Prototype Problem. — This problem can be solved quite easily by the well-known simplex computational technique (10). In order to do this, the inequalities must be converted to equalities. This is done by adding a slack variable to each constraint. The set of these slack variables are S_1, \dots, S_9 . Table 4 shows the original simplex tableau.

In this tableau, the second row (X_j row) contains the designation of the variables corresponding to the particular column. The rate of change of the criterion function for each variable is located in the top row (C_j row) above the variable. The X_i column contains the basis variables and the rate of change of the criterion function for each of

these is located adjacent to them in the C_j column. Hence the subscript designation j refers to any variable, whereas the subscript i refers only to the basis variables.

The slack variables, S_1, \dots, S_9 , constitute the original basis since they form an identity matrix. This can be seen in the X_j column since S_1, \dots, S_9 appear in this column. The last column (b_i column) contains the current values of each basis variable. Originally, for example, $S_1 = 5,900$.

The $Z_j - C_j$ row contains the evaluators which are used to determine whether the introduction of a particular variable into the basis will produce an increase in the value of the criterion function. If a particular $Z_j - C_j$ is negative the introduction of the corresponding variable into the basis will produce this increase. As seen in Table 4, any of the six variables ($X_1, X_2, X_3, X_4, X_5, X_6$) will improve the solution when introduced into the basis.

The value of Z_0 is the value of the criterion function which is produced by the particular set of variables in the basis. Since, in the original tableau, all basis variables are slack variables and do not contribute anything to the criterion function, $Z_0 = 0$.

Six routine simplex iterations (one for each variable introduced into the basis) were required to reach the optimal solution for this model (Table 5). The values of the basis variables which yield the optimal solution appear in the column on the right. All of the non-basis variables equal zero, of course.

Interpretation of Solution.—In the formulation of the linear programming model it is possible to put in constraints which turn out to be redundant. For such equations the slack variables have positive values in the optimal solution. The interpretation given to the value of a slack variable in the optimal solution is the amount that a particular constraint would have to be reduced before it would cease to be redundant (10). In other words it is the excess "capacity" contained in the restrictions.

There are three slack variables in the optimal basis. They are S_2, S_4 and S_7 . The slack variable S_2 is for the second equation, which is the capacity restraint at the Austin on-ramp merging section. Since $S_2 = 213$, 213 more vehicles could be passed through section B without developing congestion. The fourth equation states that the volume on the Des Plaines entrance ramp cannot exceed the demand of 600. However, the capacity of section A provides a greater restriction on the Des Plaines volume than this. This merging capacity restricts X_1 (the Des Plaines ramp volume) to 447 vehicles. This is 153 less than the demand on this ramp, so $S_4 = 153$ and there would be an unsatisfied demand of this amount at the Des Plaines ramp. Similarly, since $S_7 = 133$ there would be an unsatisfied demand of 133 vehicles in the hour period at the Central Ave. on-ramp. The capacity at section C provides a greater restriction to the variable X_4 than does the limit which is placed on the ramp demand. The total number of vehicles which must be prevented from entering the freeway during the hour is $S_4 + S_7 = 286$.

The effect on the value of the criterion function of unit changes in constraints is also interesting (Table 5). These are contained in the $Z_j - C_j$ row. The $Z_j - C_j$ at this stage are the optimal values of the dual variables (10). The dual variables are interpreted as the rate of change of the criterion function for a unit change in the corresponding constraint. For example, the dual variable for the first constraint is contained in the $Z_j - C_j$ row of the slack variable for the first equation (S_1). Its value is 1. This means that for a unit increase in capacity at section A, a unit increase in the system output (or input) is realized. The dual variable for the second constraint is zero. Since section B is not now (optimally) operating at capacity, increasing the capacity at this location would merely add to the overcapacity and would not increase the output of the system.

The dual variable for Eq. 3 must be interpreted as an expected value since its value is 0.067. If the capacity at section C were increased by one, an additional vehicle could be allowed to enter the freeway from the Central Ave. ramp without creating congestion at C. However, the vehicle could only be allowed on the freeway if it were going to exit before reaching section A or congestion would develop at that point. Since 6.7 percent of the vehicles entering the freeway at the Central Ave. ramp leave at Harlem Ave., the expected increase in input is only 0.067 vehicle—the value of the dual variable. This analysis indicates that remedial action in this subsystem should begin at the Des Plaines Ave. bottleneck since the value of its dual variable is higher than that of either of the other two bottlenecks.

For both Eqs. 4 and 5 the value of the dual variable is zero. These equations refer to constraints on demand at the Des Plaines and Harlem on-ramps. Since section A is operating at capacity, any increase in demand at either ramp could not be matched by an increase in the output (or input) volume.

Equation 6 placed an upper limit on the volume input at the Austin entrance ramp and the value of the dual variable for these equations is 0.051. This again is the expected value of the increase in system input with a unit increase in demand at this ramp. Of the vehicles which enter the freeway at Austin, 5.1 percent exit at Harlem. Hence, the probability that the additional vehicle would exit at Harlem is 0.051. Since it would be allowed to enter the freeway only if it were not going to pass through section A, the expected value of the increased input is 0.051.

The demand constraint at Central Ave. is redundant so its dual variable is zero. The Cicero ramp and Cicero mainline have dual variables equal to 0.111 and 0.429, respectively. Again these are interpreted as the expected values of the increases in system input for a unit increase in demand at these points. A closer look at the value of 0.111 at the Cicero ramp location is of some interest. The additional vehicle could only be allowed to enter the freeway if it were not going through section A. The probability of being allowed to enter the freeway is 0.176. However, if the vehicle crosses section C it will decrease by one the allowable number of vehicles from the Central ramp. Each vehicle entering at Central Ave. has only a 0.067 probability of being allowed to enter and the probability of a vehicle entering at Cicero Ave. crossing section C is 0.969. Hence, the expected increase in system input with a unit increase in demand at the Cicero ramp is $0.176 - (0.969)(0.067) = 0.111$ vehicle.

The optimal tableau shows a value of 9,364 for Z_0 . This means that, in the hour considered, 9,364 vehicles could enter the system of freeway without encountering congestion and that this is the maximum number that can do so.

Extension of the Deterministic Model.—One of the potential drawbacks of a metering system is the buildup of queues of vehicles on the metered entrance ramps. One restriction that could be placed in the model is an upper limit on the number of vehicles which are not allowed to enter the freeway in the period considered. This can be alternately viewed as establishing the minimum number of vehicles using a given ramp in the time period or a lower limit on the metering rate. This is not necessarily the maximum queue length but could perhaps be related to this quantity without a great deal of difficulty. In any case it might be meaningful to place a limit on the number of vehicles in the hour period which are not allowed to enter the freeway from any input.

This is simply an upper limit on one or more of the slack variables and could be accomplished by adding inequalities of the type $S_j \leq Q_j$, where S_j is the j th slack variable and Q_j is the maximum number of vehicles on the j th ramp which are not allowed to enter the freeway in the time period considered.

Accidents, disabled vehicles, adverse weather, etc., frequently cause reduced capacity operation at one or more sections on the freeway. Hence, it would be desirable to somehow incorporate the effects of these events into the model so that optimal system operation (under the reduced capacity conditions) can be obtained. The discussions will be concerned with a capacity reduction at one section but this can readily be extended to cover the adverse weather situation.

In order to include the reduced capacity situation a capacity constraint will have to be placed on a section between each successive pair of ramps or, alternately viewed, on a section downstream of each location at which the volume can change. These would be similar to the constraints placed on sections A, B, and C in Figure 9. During normal conditions the normal capacity at each section would be used and many of the constraints would be redundant. However, in case an accident reduced the capacity at a given section, the reduced capacity would be used in the constraint for this section. The solution of this problem would yield the optimal inputs at each ramp under the conditions.

When thinking of using the linear programming model for control, one might wonder how the capacity at an accident location could be determined since an accident can have a wide range of effects—from virtually no effect to the closing of all of the freeway lanes (in one direction). If a detection station is located downstream of the accident it is

possible to measure the capacity flow rate directly when congestion develops at the accident. The flow rate at the downstream detection station would be the capacity flow rate past the accident. This capacity could then be used in the linear programming model if it were being used for control purposes.

Limitations.—The use of the deterministic linear programming model assumes accurate knowledge of the O-D characteristics of each volume input source in the system. These data might change significantly with time on a given day or from one day to another. The implementing of a metering system would almost certainly change these characteristics so it would be necessary to obtain new O-D data for the system.

This model can be used only for time periods which are long compared to the travel time through the system. For this reason, it might be necessary to consider a dynamic model.

Obtaining O-D Data and Estimating Demand at a Section.—In view of the sensitivity of the model to changes in certain O-D data, it is quite important to have an accurate knowledge of these data. Since it might vary by time of day the data should be collected according to short time periods (such as 15 minutes). It is necessary to determine for each input source for each time period the percent of vehicles which exit at each output. This could be done in any one of several ways. The method discussed here consists of a 100 percent sample of vehicles entering each ramp on each of several days. It is quite a laborious method but it provides a great deal of information that more conventional O-D techniques could not provide. Sampling could be confined to time periods of interest.

The method was actually used by Brenner, et al. (11), but for different purposes. It consists of recording the time of arrival and the license number of each vehicle entering the freeway at each entrance ramp in the system, the time of departure and license number of each vehicle leaving the freeway at each exit ramp and counts of all vehicles entering and leaving the system via the freeway input and output. Matching the license numbers and times would yield the O-D data by time of day. As was done by Brenner, et al., these data could be used to determine travel times as well.

Another valuable by-product of these data would be the ability to estimate the demand on any section in the system. If the free-flow travel time between an entrance ramp and a bottleneck section is t , a vehicle entering the freeway at this entrance ramp at time t_0 represents one unit of demand at the bottleneck at time $t_0 + t$ (providing it does not exit before reaching the bottleneck). This is independent of the effects of intermediate bottlenecks. The sum of these demands over all inputs would yield the actual demand at a given bottleneck.

CONCLUSIONS AND RECOMMENDATIONS

This paper is primarily theoretical and presents many hypotheses which need to be tested. Subject to the validation studies suggested in the section on recommendations which follows, this theoretical study offers the following conclusions applicable to a somewhat idealized urban freeway system of the type which was analyzed.

Conclusions

1. Congestion develops at a bottleneck location when the upstream flow exceeds the bottleneck capacity for a sufficiently long period of time.
2. The development of congestion at a bottleneck causes high-density, low-speed operation upstream of the bottleneck.
3. There is evidence to indicate that the flow rates at many freeway bottlenecks are lower when there is congestion upstream than during some periods of free flow. The reduction in the flow rate at a bottleneck under normal operating conditions may be due primarily to the inability of the congested upstream freeway to supply vehicles to the bottleneck at its capacity flow rate. Under these conditions the start-and-stop flow upstream of the bottleneck may be the factor limiting the flow rate to the bottleneck.
4. Since congestion upstream of a freeway bottleneck can cause the bottleneck flow rate to decrease, the output of the freeway system can be increased (or maintained at its maximum level) by the prevention of congestion at all locations in the system. One goal of a control system, then, is to prevent the development of congestion everywhere in the freeway system.

5. In most cases, for a given demand on a street and/or freeway system, the peak-period objective of maximizing the output of the system is equivalent to minimizing the total travel time in the system.

6. Traffic control on freeway systems holds promise for reducing freeway congestion and reducing travel time in the total street and freeway system. Such controls include both (a) controls on the freeway and (b) control of the inputs to the freeway.

7. Control of certain inputs to the freeway system seems to be the most effective method of preventing congestion on the freeway during normal operating conditions and minimizing the effects of reduced capacity operation. Of the various input controls, ramp metering appears to hold the most promise. By allowing ramp vehicles to enter the freeway at the maximum rate that will not cause congestion, it should be possible to obtain the best use of the freeway system.

8. At most metered entrance ramps, the vehicular storage capacity probably is insufficient to store the maximum queue which develops at the ramps. The storage of queued vehicles is one major problem of metering.

9. Volume measurements should be very useful as predictors of developing congestion, as long as detection takes place upstream of all bottlenecks and entrance ramps. However, measurements of lane occupancy, speed or density are needed in addition as indicators of congestion, such measurements preferably to be made at bottleneck locations.

10. The use of volume measurements in a freeway control system also has other advantages: (a) such measurements provide a check to determine whether the output volume rate of the system is being maximized; and (b) the continuity characteristics of volume make it the only variable which is well suited to theoretical system analyses.

11. Application of a control at one location affects traffic operations at many other locations. The entire system should be studied, not the isolated locations. For this reason, a systems analysis is perhaps the most adequate analytical technique for predicting the effect of a control or control system on the system under consideration.

12. Linear programming provides a valuable tool for describing the operation of a freeway system or subsystem. It is possible that it could be used in reduced capacity situations to determine the proper controlled ramp inputs to provide optimal peak-period operation of a freeway system.

13. Demand on a freeway section can be estimated by sampling the O-D characteristics and volume-time characteristics of free-flowing system inputs upstream of the section.

Recommendations

Probably the most important empirical study that should be undertaken is the study of the behavior of the macroscopic traffic variables at and upstream of various types of bottlenecks in order to determine the effect of control on traffic operation at these locations. The questions of whether or not congestion normally decreases the flow rate, in what situations the volume rate decrease takes place, and how much the volume decreases due to congestion must be answered. In situations in which congestion does not decrease the flow rate, freeway travel time under normal operations can be significantly decreased only by decreasing the inputs. In this case, the output rate of the freeway is little affected by freeway storage, so the prevention of freeway congestion is not necessarily the peak-period objective (although it probably would still be desirable).

The flow rate away from a queue should also be investigated because it will furnish some information on a steady-state congested flow rate. Perhaps there is no single value, but if there is, our knowledge of it will contribute greatly to the evaluation of the possible effects of a control system.

Shock wave development and propagation should also be studied. Queuing forms and dissipates at bottlenecks prior to congestion but finally flow breaks down when excessive queuing takes place. The causes of queuing and the behavior of the shock waves at a bottleneck should be studied, because such study will furnish information on the causes of congestion and for what time period and by how much upstream flow rate can be allowed to exceed the bottleneck capacity. In other words, it will help to determine the probability of congestion which is associated with a particular set of upstream flow conditions.

Since many freeway bottlenecks are at entrance ramp locations, the merging maneuver could be studied along with the shock waves and queues at bottlenecks. The capacity of the merging areas under different conditions must be obtained in order to establish the upper limit of the ramp input rate for a given upstream flow rate. It is frequently assumed that there is a fixed merging capacity and that the proportions of vehicles merging from the ramp and freeway do not affect the capacity value. This assumption of a one-to-one tradeoff between ramp and freeway vehicles should be investigated.

The effect of entrance ramp metering on the merging capacity is also needed. If metered arrivals from the ramp allow a larger upstream freeway flow rate, an additional benefit of metering will have been realized.

A method for the estimation of demand must be evaluated. The method of recording license numbers and arrival times of vehicles at the upstream entrance ramps as well as a free-flowing freeway section should be tested. The duration as well as the severity of the control at a given location depends on the demand function at this location. If it is not known, the selection of the proper control may be difficult. The changes in demand caused by upstream controls must also be determined.

While the license numbers and arrival times are being recorded on the ramps, another important study should be conducted. The linear programming model assumed that the freeway O-D pattern remained constant during the peak period. If license numbers of exiting vehicles are recorded at each exit ramp, the changes in the O-D patterns with time can be obtained.

The control system proposed here would work best under "normal" traffic conditions (i.e., no accidents, disabled vehicles or other "unusual" events), since the full capacity of the freeway could be used. However, the "unusual" situation can also be taken care of since ramp closure is possible as part of a flexible metering system. In this case some vehicles would be diverted around the capacity reductions on the freeway and onto those surface streets which have remaining capacity. The frequency of these events under various volume rates and congestion conditions should be examined. Even more basic and important, the effects of a traffic accident, tire changer, disabled vehicle and other "unusual" events on traffic behavior and especially on the flow rates should be studied. The effect of adverse weather on the capacity flow rate of bottleneck sections also warrants intensive investigation. The complexity of the final control system may depend on the outcome of these studies.

The cost of the final detection system could be substantially reduced if it is possible to use measurements of speed, lane occupancy or density to estimate volume or if it is possible to estimate the traffic variables for all lanes by sampling detection in one lane. These possibilities should be thoroughly investigated to determine the sacrifice of accuracy that accompanies an economic savings.

The philosophy of this paper is to accept the total demand on the freeway and to operate the freeway system in an optimal manner. No vehicles (or at least as few as possible) would be diverted from the freeway. They could be delayed from entering by means of the metering system. This set of conditions permits only the optimal operation of the freeway system but does not assure optimal operation of the total system which includes the streets as well as the freeways. The next logical step is the development of a model which would yield the optimal operation of the entire system. Perhaps the Charnes-Cooper multi-copy model (10) with capacitated entrance ramp links (to produce travel time increases with volume increases) would fill this need.

ACKNOWLEDGMENTS

This paper is based on part of a doctoral dissertation at Northwestern University. Wattleworth was the author and Berry was chairman of the faculty committee in charge of evaluating the dissertation. Other members of this committee, whose contributions certainly deserve mention, are Professors Abraham Charnes, Adolf D. May and Loring G. Mitten.

The work on this material was also part of Wattleworth's duties at the Expressway Surveillance Project of the Illinois Division of Highways. Discussions with the staff of this project, including Project Director Adolf D. May, also proved valuable in the preparation of this material.

REFERENCES

1. Rothrock, C. A., and Keefer, L. E. Measurement of Urban Traffic Congestion. Highway Research Board Bull. 156, pp. 1-13, 1957.
2. Edie, L. C., and Foote, R. S. Traffic Flow in Tunnels. Highway Research Board Proc., Vol. 37, pp. 334-344, 1958.
3. Edie, L. C., and Foote, R. S. Effect of Shock Waves on Tunnel Traffic Flow. Highway Research Board Proc., Vol. 39, pp. 492-505, 1960.
4. Wattleworth, J. A. Bottleneck Throughput Study. Unpublished report presented to the Freeway Operations Comm. of the Highway Research Board, 1963.
5. Highway Capacity Manual. U. S. Dept. of Commerce, Bureau of Public Roads. 1950.
6. May, A. D. Experimentation with Manual and Automatic Ramp Control. Highway Research Record 59, pp. 9-38, 1964.
7. Barker, J. L. Determination of Discontinuities in Traffic Flow as a Factor in Freeway Operational Control. Inst. of Traffic Engineers Proc., pp. 25-40, 1961.
8. Lighthill, M. J., and Whitham, G. B. On Kinematic Waves II. A Theory of Traffic Flow on Long Crowded Roads. Royal Soc. of London Proc., Series A, Vol. 229, pp. 317-345, 1955.
9. Characteristics of Traffic Entering the Westbound Congress Expressway Within the Pilot Detection System Study Section. Unpublished project report of the Expressway Surveillance Project, Chicago, 1962.
10. Charnes, A., and Cooper, W. W. Management Models and Industrial Applications of Linear Programming. 2 Vols., Wiley, 1961.
11. Brenner, R., Telford, E., and Frischer, D. A Quantitative Evaluation of Traffic in a Complex Freeway Network. Highway Research Board Bull. 291, pp. 163-206, 1961.