# A Dynamic-Zone-Based Coordinated Ramp-Metering Algorithm With Queue Constraints for Minnesota's Freeways

Nikolas Geroliminis, Anupam Srivastava, and Panos Michalopoulos

*Abstract*—Following about 40 years of successful deployment of coordinated traffic-responsive ramp control, a new generation is being developed for Minnesota's freeways based on density measurements, rather than flow rates. This was motivated from recent research indicating that the critical value of density at which capacity is observed is less sensitive and more stable than capacity, thereby allowing the opportunity for more effective control. The main goals of the new approach are to delay the onset of the breakdown and accelerate system recovery when ramp metering is disabled due to the violation of maximum allowable ramp waiting times. This is obtained by a dynamic zone partitioning of the freeway network to identify critical bottleneck locations and coordinated balancing of ramp delays, which aims to avoid mainline breakdown. The effectiveness of the new control strategy is assessed by comparison with the currently deployed version of the stratified zone metering algorithm through microscopic simulation of a real 12-mi 17-ramp freeway section. Simulations show a decrease in delays of mainline and ramp traffic and an improvement of 8% in overall system delays while avoiding maximum ramp delay violations.

*Index Terms*—Coordination, queues, ramp metering, traffic congestion, traffic control.

## I. INTRODUCTION

**O**N-RAMP metering has been widely employed as an effective strategy to reduce congestion and increase freeway operational efficiency. It is one of the most efficient tools to mitigate congestion, other than adding capacity to infrastructures. The fundamental philosophy of ramp metering is that a corridor can maintain its optimal operation by regulating the freeway input to be under its capacity. Over the years, a number of traffic-responsive ramp control strategies have been developed to regulate the entrance ramp demand to freeways. These strategies use basic principles of feedback and feedfor-

N. Geroliminis is with the Urban Transport Systems Laboratory and the School of Architecture, Civil and Engineering, Ecole Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland (e-mail: nikolas.geroliminis@epfl.ch).

A. Srivastava and P. Michalopoulos are with the Department of Civil Engineering, University of Minnesota, Minneapolis, MN 55455-0116 USA (e-mail: anupam83@gmail.com; micha001@umn.edu).

ward controls with minor modifications and can be classified as isolated or coordinated.

In isolated ramp metering, the rate for an on-ramp is determined based on its local traffic conditions. The most prominent example of the isolated metering strategy with multiple successful field applications is ALINEA, which targets a set point for the downstream occupancy or density [1]. A review of other isolated algorithms is presented in Papageorgiou and Kotsialos [2]. The main inefficiencies of isolated ramp algorithms are manifested in three cases: 1) high demand, where inequity issues arise, as travelers at the on-ramps experience significantly higher delays than those in the mainline or adjacent ramps; 2) short on-ramps, where long queues can spill back to the arterial network and create significant delays in adjacent streets; and 3) multiple active bottlenecks.

Coordinated ramp-metering strategies utilize system-wide traffic measurements from an entire region of the network to control all on-ramps within the region. Some operational examples of coordinated algorithms are the Bottleneck [3], Zone [4], Metaline [5], Helper [6], and Swarm [7]. A variety of optimal control theories coupled with dynamic traffic simulation models have been tested in [8]–[13]. The objective in most of these studies is to minimize the total system travel time. An extensive review of heuristic coordinated traffic-responsive ramp-metering algorithms is presented in [14] and [15].

A more advanced concept in freeway control (similar to the coordinated one) employs a hierarchical control structure where overall freeway control is decomposed into several components or layers, such as prediction, optimization, and direct control. The main objective of hierarchical control is to achieve computational feasibility and robustness of the control solutions by determining the systemwide nominal metering rates first and then adjusting them according to real traffic conditions. Previous efforts have employed the hierarchical control concept based on two, three, or four layers. Papageorgiou [16] proposed a three-layer control system consisting of an optimization, a direct control, and an adaptation layer. Four-layer hierarchical control, as proposed by May [17], includes initialization, estimation, optimization, and tactics processors. The Stratified Zone Metering (SZM) algorithm, which has been deployed by Mn/DOT since 2003, is essentially of three-layers, including ramp, zone, and system layer (i.e., broken zone identification layer) design. Total ramp volume is distributed over all metered ramps in the zone in proportion to their demands. SZM targets the flow capacity in the merge area and is based on mainstream and

ramp measurements of flow and occupancy. Further extensions to this algorithm include improvements for the determination of the ramp minimum release rate [18] and the queue size estimation [19].

The aforementioned Minnesota's strategies rely on capacity estimations from flow measurements. However, analysis of real traffic data from freeway merge areas indicates that the recurrent traffic breakdown during peak hours can occur at different flow values, even under the same weather and lighting conditions [20]. Furthermore, the probabilistic characteristics of capacity have been broadly observed in the literature, e.g., [21], due to the breakdown phenomena or variability of weather conditions. On the other hand, earlier research [22] indicated that the critical value of occupancy at which capacity is observed is less sensitive and quite stable. Our own research on the subject in Minnesota's freeways confirms the preceding findings and indicates a capacity drop of as high as 15%, resulting in substantial miscalculation of the optimal metering rates. Similar results for capacity drop have been reported for many locations in the same network [23]. This suggests that a control strategy based on flow thresholds is likely to underload the freeway or lead to traffic congestion. Therefore, it is desirable to develop a new methodology to overcome this uncertainty for making the stratified ramp control strategy more robust and adaptive to real-time traffic conditions.

The strategy presented here resulted from discussions with Mn/DOT, which first identified the need to move to a new generation based on density considerations rather than flows. It was also recognized that the principles of the current control scheme were developed in the late 1960's; albeit the most recent ramp queue improvements were only deployed in 2003, the main line modeling was not changed. Thus, the new generation should have a major impact on both the state of the practice and the state of the art. The main goals of the strategy are to delay the onset of the breakdown and to accelerate system recovery when ramp metering is unable due to the violation of the maximum allowable ramp waiting time. To achieve this, the strategy consists of dynamic zone partitioning of the freeway network (which is used to identify critical congestion locations), followed by coordinated balancing of ramp delays and avoidance of mainline breakdown. In this paper, we present some empirical observations for the current algorithm from the Twin Cities Metropolitan area. This is followed by a description of the new density-based coordinated algorithm and an implementation to a 12-mi corridor at the H-169 freeway through microsimulation. In addition to the preliminary simulation results presented here, a field implementation and testing of the control strategy is currently being planned for refinements and large-scale deployment on the Twin Cities freeway.

## II. EMPIRICAL OBSERVATIONS

The control objectives of the current SZM algorithm are regulating the zone inputs so that the total entering volumes do not exceed the zone capacity and limiting maximum ramp wait times below a predetermined value (4 min) during the control period. The algorithm is built on the basic philosophy of balancing the volumes entering and leaving a section of freeway based on flow conservation. A unique feature of SZM is that zones are grouped into multiple layers from 0.5 to 3 mi in length. A layer is defined as a sequence of successive zones comprising an equal number of stations. As zones overlap, ramps are assigned multiple release rates (one for each layer), and the most restrictive is selected as the dominant one. A more detailed description of the SZM algorithm can be found in [24].

With respect to the current algorithm, by carefully analyzing empirical data of active bottlenecks in the Twin Cities freeways, we noticed that there are many cases where four conditions hold.

1) Capacity is underutilized because its value before the occurrence of the breakdown is usually higher than that implemented in the current algorithm. In addition, there are cases where the 4-min ramp constraint is miscalculated by the control logic.
2) Once the 4-min ramp constraint is violated, the metering rates significantly fluctuate, e.g., in a 10-min period, ramp rates rise and fall between very low and very high many times. These high fluctuations increase the severity of stop-and-go traffic.
3) The coordination is not efficient with respect to the ramp waiting times. Many ramps close to that with long delays exhibit short delays, even when the demand is high.
4) Once a freeway is congested, the capacity considerably drops, and the system is unable to return to a state of "capacity before the first breakdown" for too long, thereby prolonging the congestion period and increasing system delays.

Fig. 1 summarizes results from real-data analysis using loop detectors for the current SZM algorithm. Fig. 1(a) shows the queue size in an active bottleneck in TH-169 (Plymouth on-ramp location) and the ramp waiting times in the three upstream ramps. Note that, when the active bottleneck of Plymouth experiences very long ramp queues (queue detector occupancy of 100%), the maximum waiting times at the three upstream ramps are much less than 4 min, and coordination is inefficient. In a well-operated coordinated algorithm, upstream ramps should store more vehicles for relieving the load of the critical ramp. Fig. 1(b) plots the ramp-metering rate and the ramp queue detector occupancy in the Plymouth location at 30-s intervals. The graph clearly shows the strong variability of metering rates when ramp constraint is violated. Note that the metering rate jumps from 240 to 1200 vh/h within 1 min to decrease the queue length and then decreases to a very low value. Fig. 1(c) shows an example of capacity drop at Plymouth active bottleneck. Note a 15% range of capacity drop after the breakdown (from 4600 to 4000 vh/hr) for a density of 35 vh/mi/ln.

Based on these observations, we identified six shortcomings of the SZM, which we addressed to improve the overall freeway performance.

1) Capacity is considered constant during all times at all bottlenecks. We have observed a significant capacity drop after the breakdown in many locations (varying from 10% to 20%).
2) The total capacity of an active bottleneck (mainline + on-ramp) depends on the ratio of the two flows. More
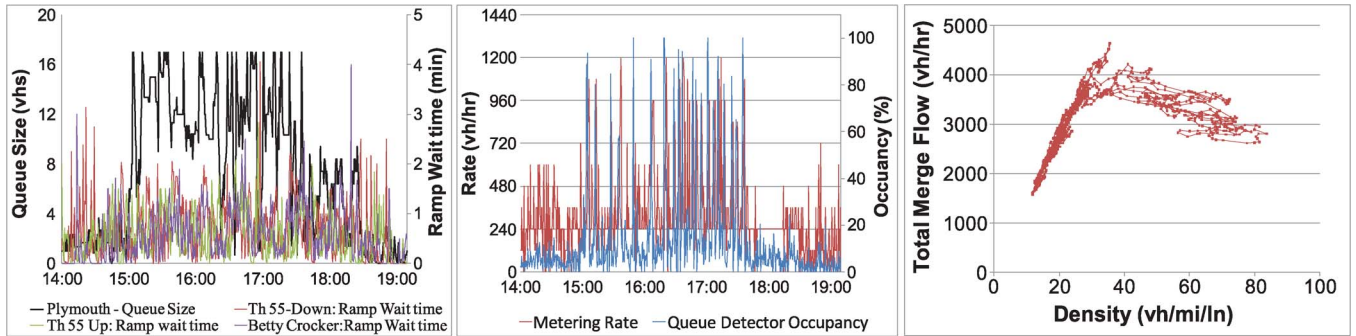
Fig. 1.   Empirical observations for Plymouth bottleneck. (a) Inefficient coordination: queue size at Plymouth (bold) versus ramp wait times at three upstream on-ramps. (b) Strong ramp rate variations: Plymouth's on-ramp metering rate versus occupancy at the ramp queue detector. (c) Capacity drop: The capacity before the breakdown is 15% larger than that after the breakdown.

specifically, when ramp flows are higher, the capacity is smaller ($\sim$5%–10%). This is related to the overreaction of the algorithm when the 4-min violation occurs and the metering rates "jump." The effect of ramp flows has been reported by many researchers, e.g., [20] and [25].

3) While values of occupancy near capacity are quite stable, bottleneck capacity has stochastic variations, and a control strategy based on constant predetermined flow thresholds is likely to underload the freeway or, conversely, lead to traffic congestion.

4) The current algorithm cannot postpone the occurrence of the 4-min violation, which is directly related to the initiation of congestion.

5) The multilayer approach cannot always balance the ramp rates between successive ramps because applied rates are often estimated at different layers.

6) A careful consideration of the "maximum delay at ramps" parameter could lead to better system performance.

The first three issues are directly addressed by switching from a flow- to a density-based algorithm. Point 4 is addressed by better balancing of ramp queues. Dynamic zone partitioning resolves point 5, whereas some sensitivity analysis enhances point 6.

## III. ALGORITHM'S DESCRIPTION

Instead of a layer-based algorithm as in the current SZM algorithm, we proceed with a dynamic-zone-based algorithm. The whole freeway system is divided into zones, where the *length of each zone is dynamic and is estimated in real time*, i.e., it is embedded in the algorithm. As a fully coordinated feedback control algorithm is not a feasible solution (from both a computational and an implementation point of view), the role of zones is to decompose the problem without significant loss of optimality. Within each zone, the metering rates are chosen independently of the conditions in other zones. The zone boundaries and length mainly depend on the following: 1) the level of congestion on the mainline direction; 2) the demand and queue lengths at individual on-ramps; and 3) the location of the active bottleneck(s). An *active* bottleneck arises when vehicles discharge from an upstream queue (to guarantee that the bottleneck served vehicles at a maximum rate) and vehicles are unimpeded by traffic conditions emanating from

further downstream. The location of an active bottleneck is either at freeway merges, where there is an additional inflow, or in areas where the capacity changes (e.g., a change in topology from three to two lanes). The algorithm sequentially identifies the traffic states of freeway sections, partitions the freeway to zones, and determines the metering rates for each section for all zones.

Let us first define a section of a freeway. The section is (in most of the cases) a portion of the freeway between two consecutive mainline detector locations. Each section contains one on-ramp and one off-ramp, with a few exceptions (e.g., freeway-to-freeway connections, two successive ramps, etc.). Each section cannot contain more than one on-ramp. If there are more than one on-ramp between two consecutive mainline detectors, this part of the freeway consists of $x$ sections, where $x$ is the number of on-ramps.

When freeway traffic is free flowing and on-ramp demands are low, there is no need for a coordinated algorithm, and rates are estimated based on local conditions (zone length = *one* section). A coordinated algorithm is used (which results in longer zones) when there is a series of high-demand ramps upstream from an active bottleneck, or there are two closely spaced active bottlenecks and the downstream one is propagating upstream. Thus, we need to provide the algorithm i) with a procedure to identify the active bottlenecks in the freeway, ii) with quantitative criteria to choose the appropriate type (local or coordinated), and iii) with a methodology to estimate boundaries and lengths of zones (starting section and ending section or zone length).

We now provide the definition of a zone. A zone comprises an active bottleneck location (or a location that is a threat for an active bottleneck if demand conditions continue to be the same) at its downstream end and the upstream sections that are influenced by the bottleneck/threat. The associated on-ramp of this downstream bottleneck is the *controlling ramp* for this zone. All rates in the other ramps will be estimated to postpone the occurrence of the breakdown at the controlling ramp.

Each zone should tail off if possible with two to three sections with low on-ramp demands that are historically uncongested. The preceding segmentation may lead to zones with more than one bottleneck location. We will address later the issue of choosing the controlling ramp among different active bottlenecks. We also integrate a constraint for the

maximum length of a zone ($\sim$5 mi) as the controllability of long zones is difficult when conditions are not stationary.

The algorithm's goal is to keep the car density levels at all ramps below the congestion thresholds and not to allow low speeds to occur in the mainline by constraining the ramp delays. The ramp rates become stricter when the mainline density is close to the congestion threshold, and the ramp rates increase when the ramp waiting times are close to the ramp delay threshold. When it is not possible to keep both uncongested because of high on-ramp and mainline demands, the algorithm seeks to delay the violation of both thresholds as much as possible. The metering rates are decided based on the following: 1) how close the mainline density and/or the ramp delay is to the congestion thresholds and 2) how fast they are approaching these thresholds (rate of change). For example, if two ramps have the same ramp delays but these delays increase to different rates, the applied ramp flows should be different as well. To account for both conditions, two key variables are introduced for each merge section. These are time durations $T_k$ for the mainline and $T_w$ for the ramp, until the mainline and the ramp reach the congestion thresholds, respectively.[1]

By identifying the controlling bottleneck in a freeway section, which is also the downstream end of one zone, all the upstream ramp rates are estimated in such a way that *the ramp violation will simultaneously occur in all ramps*. As the zones and controlling bottlenecks are identified in real time, the algorithm is able to capture changes in demand in different locations of the network.

## A. Associated Variables and Parameters

Let us first define the required variables for the development of the algorithm. Their values are based on historical or real-time data from loop detectors. For each ramp $i$, we estimate the following variables, which are distinguished in historical and real-time data: The first set of variables should be estimated before the implementation of the algorithm and is not time dependent, whereas the second set of variables is dynamically estimated every 30 s. Note that all variables are estimated based on 30-s flow and density measures at a moving window of *5-min moving averages every 30 s*. This procedure smoothens random or unpredictable (difficult to estimate) phenomena such as short fluctuations in demand, percentage of trucks in the sample of vehicles, detector errors, behavioral characteristics of drivers, etc. Some of these variables have almost identical values among all ramps, and in this case, global values are assigned. We also define some parameters that are used to identify different traffic states and decisions.

1) Critical density for the mainline upstream from ramp $i$, i.e., $k_{\text{crit}}(i)$. This variable is used to identify the threshold between congested and uncongested conditions in the mainline. It is estimated from historic data as the average density among 2% of the flow-density pairs observed (over multiple days) with the highest flow for the closest upstream detector to a merge. In most current ramp-

metering algorithms based on density (e.g., ALINEA [1]), a slightly smaller value than the critical density is applied, e.g., $0.95 \times k_{\text{crit}}(i)$. We also follow a similar approach.[2]

2) Capacity flows before and after the occurrence of congested conditions $c^h(i)$ and $c^l(i)$. $c^h(i)$ is estimated as the average flow from all the flow-density pairs with values of density between $0.95 \times k_{\text{crit}}(i)$ and $1.00 \times k_{\text{crit}}(i)$, whereas $c^l(i)$ is that for values of density between $1.00 \times k_{\text{crit}}(i)$ and $1.05 \times k_{\text{crit}}(i)$.[3]

3) Average density between the merge locations of the acceleration lane associated with the ramp and the upstream mainline detector or 300-ft upstream, whichever is shorter at time interval $t$. This estimation extends the kinematic wave theory modeling with finite differences of [26] by using internal boundaries for the mainline detector measures and fundamental diagrams with capacity drops estimated from empirical data. More details about the density estimation can be found elsewhere [27]. The length of 300 ft can also be overwritten for special cases. The rate of change of the mainline density for an interval with duration $T$ is then calculated as

$$\Delta k_t = [k_t - k_{t-1}(i)]/T. \tag{1}$$

4) Maximum waiting time for vehicles on the ramp at interval $t$ (based on a combination of arrival and departure patterns and on occupancy of queue detectors) $w_t(i)$. The rate of change of the delay for the ramps is then calculated as

$$\Delta w_t(i) = [w_t(i) - w_{t-1}(i)]/T. \tag{2}$$

5) The time duration until the occurrence of congested conditions on the mainline, i.e., the time to congestion on the mainline, is estimated as

$$T_t^k(i) = (k_{\text{crit}}(i) - k_t(i))/\Delta k_t(i). \tag{3}$$

6) The time duration until the violation of ramp wait time constraint $T_{\text{crit}}$, i.e., the time to ramp congestion, is estimated as

$$T_t^w(i) = [T_{\text{crit}} - w_t(i)]/\Delta w_t(i). \tag{4}$$

We now assign state conditions for each merge location. 0-state is considered safest, 1-state is near congestion, and 2-state is in the completely congested regime. These states will contribute in the estimation of zone lengths and the type of ramp metering employed (local or coordinated). State 0 occurs when the mainline density is small, and the time to congestion is large for both mainline and ramp. State 1 occurs when the mainline density is medium and the time to congestion is positive for both mainline and ramp, whereas state 2 occurs

---

[1]This does not mean that the algorithm uses a predictive approach, but these key variables are introduced to facilitate the control decisions, based on the level of congestion and the rate of its change.

[2]By analyzing historical data from many different locations, we have observed that this value is close to 40 veh/mi for most locations. Special attention should be given to secondary or inactive bottlenecks, which cannot reach values close to their capacity because of demand restrictions. For both reasons, we propose a constant value of critical density $k_{\text{crit}}(i)$, with a few exemptions if needed.

[3]As detectors are placed slightly upstream of the merge and not at the active bottleneck locations, capacity flow is estimated as the sum of mainline and ramp detector flows. In addition, for locations where flow cannot reach the true value of capacity, common values are chosen $c^h$ and $c^l$.

when congestion is observed in the mainline or the ramp, i.e., $T_t^k(i)$ or $T_t^w(i)$ is negative. Different states $S_t(i)$ at each merge location are identified as follows (where $\delta$, $\tau_k$, and $\tau_w$ are thresholds to quantitatively distinguish between states; $\delta$ is a low density margin for the mainline to separate state 0 from 1, e.g., 70%–80% of the critical density; and $\tau_k$ and $\tau_w$ are the mainline and ramp-area safe time intervals before the mainline (or the ramp) is in risk of congestion (or overflow), e.g., 10 min):

$$S_t(i) = \begin{cases} 0, & \text{if } [k_t(i) < \delta \times k_{\mathrm{crit}}(i)] \; AND \\ & [T_t^w(i) > \tau_w] \; AND \; [T_t^k(i) > \tau_k] \\ 2, & \text{if } [T_t^k(i) < 0] \; OR \; [T_t^w(i) < 0] \\ 1, & \text{otherwise.} \end{cases} \quad (5)$$

When there exist two problematic locations close to congestion thresholds that are neighboring (two to three sections apart), the algorithm has to choose which one will be the controlling section. There are two options.

1) If the upstream location experiences more severe demand conditions, then the freeway is separated into two zones, with one ending at the upstream location and another one between the two problematic locations.

2) If the downstream location experiences more severe demand conditions, then there is only one zone for this freeway, and the controlling ramp is the downstream one. All ramp rates are estimated with respect to the controlling location; the goal is that all ramps simultaneously reach the ramp waiting constraint.

To identify the most problematic location, we define the net inflow between the two locations $i$ and $j$, i.e., $M_t(i, j)$, which is given as the sum of on-ramp volumes minus the sum of off-ramp volumes for all ramps between the two locations. This justifies whether the demand at the upstream problematic location is higher or lower than the demand at the downstream location, thus representing the gain of vehicles within the given section for the specified interval. This variable will help in choosing the critical bottleneck, which will define the metering rates in the coordinated strategy. If the capacities of the two sections are different, we account for this by adding the loss in capacity from the net input obtained from the ramp flows. Thus, the net inflow ($I_t^k$ and $0_t^k$ are the on- and off-ramp flows associated with section $k$ at time $t$) is

$$M_t(i, j) = \sum_i^j \left( I_t^k - O_t^k \right) + c^h(i) - c^h(j). \quad (6)$$

### B. Bottleneck and Zone Identification

The state conditions identify the problematic locations that can be expected to become the active bottleneck(s). The zone identifies the associated sequence of sections that get affected by a given active bottleneck (or are expected to get affected by a potential active bottleneck before it is actually formed). Thus, the definition of zones also relies on the identification of the active bottleneck locations/threats. Usually, a zone would have congested/near-congested location(s) at its downstream end and two to three uncongested sections at the upstream end.
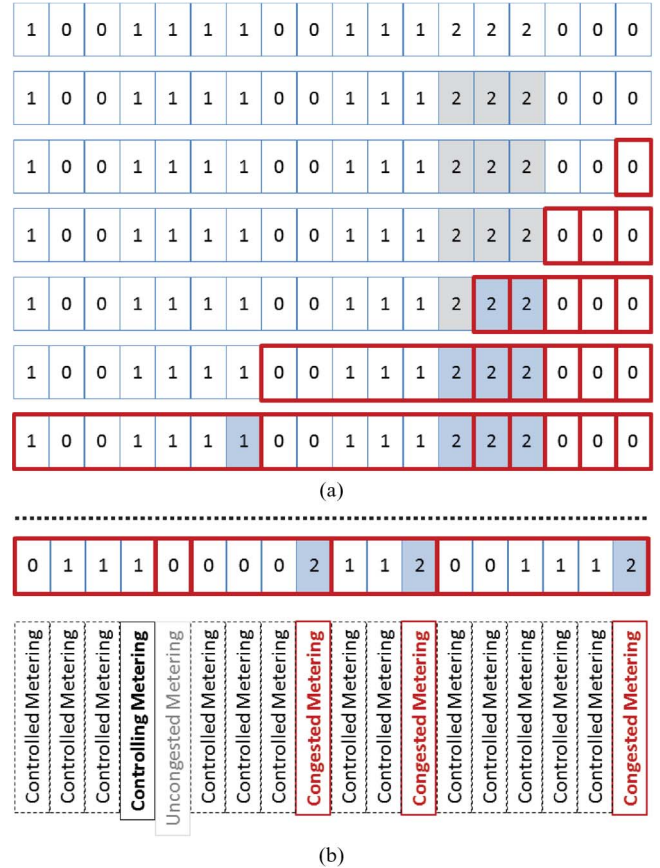


Fig. 2. (a) Example of zone identification (traffic moves from left to right). Blue cells show the controlling ramps. (b) Assignment of control strategies following zone identification.

To identify the controlling ramps and the zone length when conditions are near or at congestion, we start from the most downstream location of the freeway and move upstream by analyzing each location individually at each step. We then use the following rules to mark a merge as the controlling location. Each zone has only one controlling section/ramp at its downstream end, and all the upstream ramp rates are calculated as a function of the controlling one.

For clarity, we adopt a naming convention for the sections such that the sections are numbered from 1 to $N$ ($N$ is the total number of freeway sections) starting at the downstream end of the stretch. (Section $j$ is located upstream from $j-1$.) We first estimate the state of all sections, $S_t(j)$ for all $j = 1, 2, \ldots, N$. Then, zone identification is given as follows:

If section $j-1$ is not controlling, and $S_t(j) = 1$ or 2, then mark $j$ as controlling.

If section $(j-k)$ is controlling and $S_t(j) = 1$ for $k < 3$,
    then, if $M_t(j, j-k) > 0$, mark $j$ as controlling.
    else, mark as noncontrolling.[4]

If current state = 2, then mark $j$ as controlling.

If none of above, then mark $j$ as noncontrolling.

---

[4]The role of $k$: If there are a few on-ramps with low demand at state 0 between two on-ramps at state 1 with high demand, the upstream ramp may be controlled by the downstream one, i.e., they belong in the same zone. The criterion is the net inflow between the ramps.

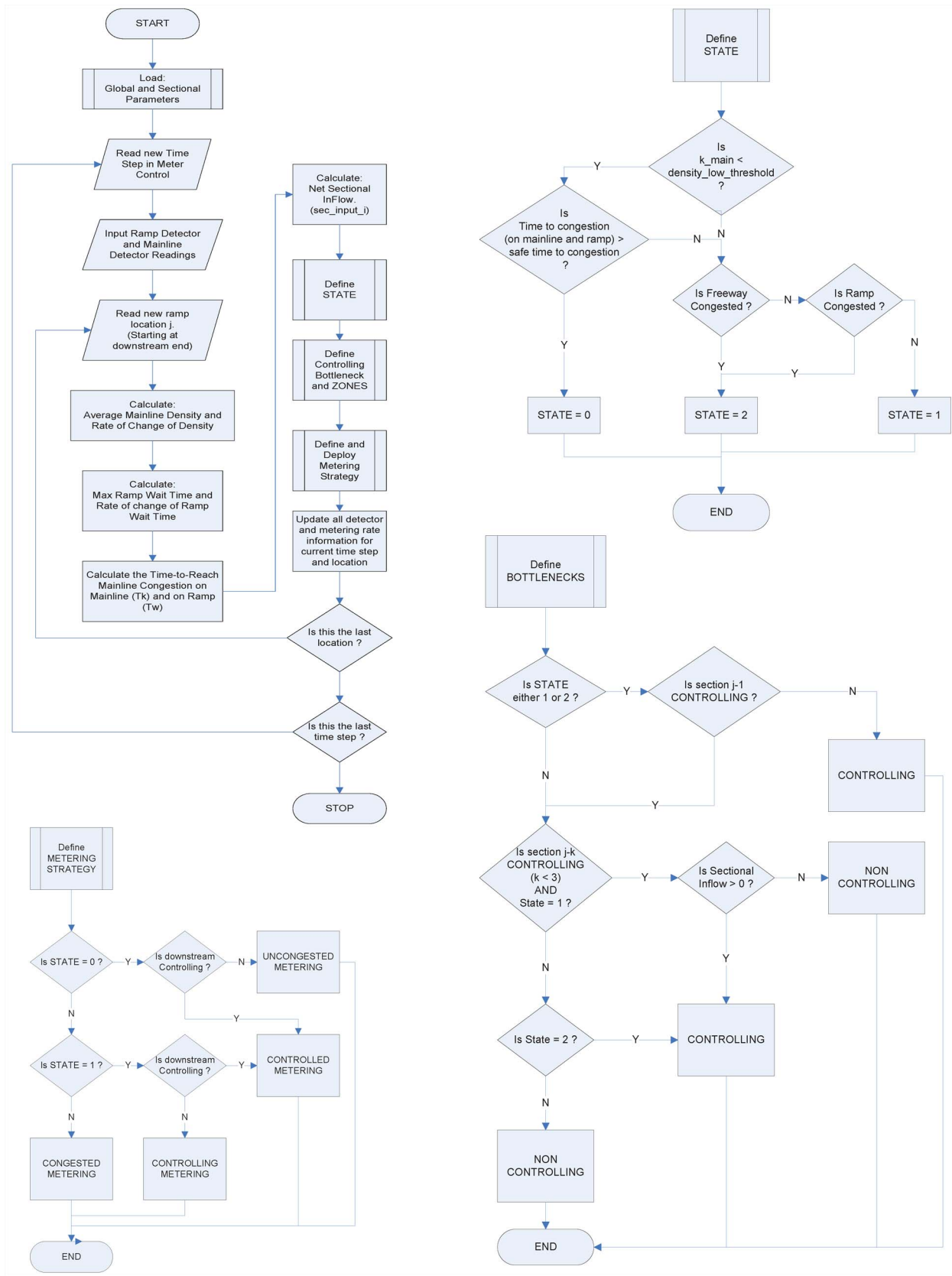Fig. 3. Decision charts for the ramp-metering algorithm.

Once the controlling ramps are identified, the zones are implicitly obtained as the complete lengths between any two controlling ramps. Thus, we ensure that the downstream end is a controlling bottleneck (either a single congested/near-congestion section or a series of sections in the congested/near-congestion state) and that it has the highest state value among

all sections within the zone. (We cannot have a noncontrolling section in state 2 if the controlling section of this specific zone is at state 1.) We also ensure that the upstream end is uncongested since it is just downstream of another bottleneck that defines the boundary for the upstream zone. A zone may sometimes have free-flow sections at its downstream boundary, if zone length = 1 section. This scenario occurs when the section is considered independent of any congestion due to being far from any downstream congestion threats. Fig. 2 shows an example of zone partitioning. (Blue cells are the controlling ramps, and red squares are the zones.)

### C. Action Matrix

Given the zone identification, the metering actions are then governed for each zone by the state of the various ramp locations and controlling bottlenecks within the zone, as defined by the following rules:

*1) Merge in State 0 With no Downstream Controlling Ramp* (*zone length, $z = 1$*) $=>$ **Uncongested_Metering**()*:* Isolated ramp metering is applied with metering rate $r_t(i)$ equal to before-the-breakdown capacity minus mainline demand $q_t(i)$. Alternatively, any simple isolated strategy (e.g., ALINEA) could be applied, i.e.,

$$r_t(t) = c^h(i) - q_t(i). \tag{7}$$

*2) Merge in State 0 or 1 With Downstream Controlling Bottleneck* (*$z > 1$*) $=>$ **Controlled_Ramp_Metering**()*:* The goal of the coordinated strategy is to balance all ramp delays to reach the ramp wait-time constraint simultaneously (balance the delay to ramp wait constraint with that for controlling bottlenecks). The following expression guarantees that all ramps within a zone will simultaneously reach the ramp wait constraint with the controlling downstream ramp:

$$r_t(i) = d_i - \frac{\Delta t_i * d_i \left[d_j - r_t(j)\right]}{\Delta t_j d_j} \tag{8}$$

where $d_i$ is the demand at ramp $i$, $j$ is the bottleneck or the controlling ramp, and $\Delta t_i = T_{\text{crit}} - w_t(i)$. Another approach for balancing queues is reported in [28].

    *Proof:* Consider at time $t$ a ramp $i$ of a specific zone **Z**, with demand $d_i$, ramp metering rate $r_t(i)$, and the ramp waiting time of the vehicle just entering the freeway $w_t(i)$ If $d_i > r_t(i)$, the time to ramp congestion $T_t^w(i)$ should be equal to $T_t^w(j) \forall i \in \mathbf{Z}$, where $j$ is the controlling ramp of zone **Z**. Time $T_t^w(j)$ can be estimated $\forall i \in \mathbf{Z}$ from the following:

$$T_{cr} + \frac{T_t^w(j)r_t(i)}{d_j} = w_t(i) + T_t^w(j). \tag{9}$$

After some manipulations, (9) leads to

$$T_t^w(j) = \frac{(T_{cr} - w_t(i))\, dj}{d_j - r_t(i)} = \frac{\Delta t_i d_j}{d_j - r_t(i)} \forall i \in \mathbf{Z}. \tag{10}$$

By substituting $i = j$ in (10) and considering that $T_t^w(i) = T_t^w(j)$, we get (8). ∎

Note that, from (8), $r_t(i)$ is positive if and only if $T_t^w(j) > \Delta t_i$. This means that the coordinated algorithm should be activated some time before $T_t^w(j)$ reaches $T_{cr}$. This fact explains
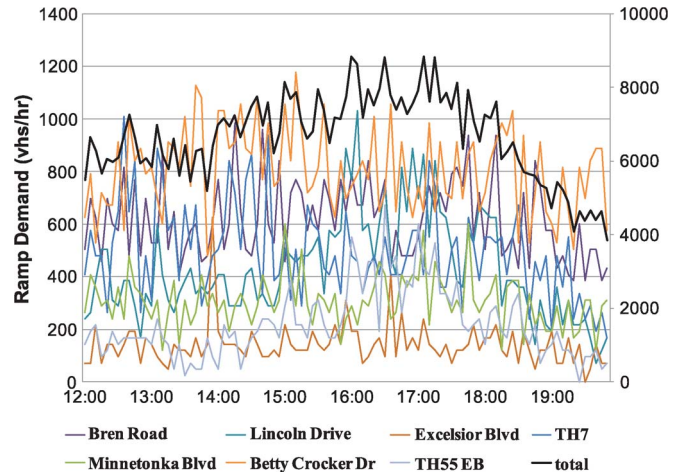


Fig. 4.    Selected test site (TH-169 NB).



Fig. 5.    On-ramp demands for the study site.

why there is a need for an intermediate state $S = 1$ in the developed philosophy.

*3) Merge in State 1 With no Downstream Controlling ramp* $=>$ **Controlling_Ramp_Metering**()*:* The goal of the strategy for the controlling ramp is to postpone the congestion occurrence as much as possible (override ramp delay threshold or mainline critical density). This is achieved by balancing the mainline congestion with delay to ramp wait time violation on the specific location

$$r_t(i) = r_{t-i}(i) - K_1 \left(T_t^w(i) - \tau_w\right) + K_2 T_t^k(i) \tag{11}$$

where $K_1$ and $K_2$ are contribution parameters for the importance of breakdown at ramp and mainline, respectively in units of vh/min$^2$.

*4) Merge in state 2*      $=>$ **Congested_Meteringl**()*:* Priority is given to the ramp violation. If there is a ramp violation, the goal is to balance metering to "just"-avoid ramp wait time violation. Otherwise, the metering rate decreases to improve conditions in the mainline. That is

$$if\ [T_t^w(i) < 0]\, then\{r_t(i) = r_{t-1}(i) + K_1 T_t^w(i)$$
$$else,\ r_t(i) = r_{t-1}(i) + K_2 T_t^k(i). \tag{12}$$

To avoid rapid increases in the metering rate (one of the main drawbacks of the current algorithm), we define a maximum change of metering $a$ when $T_t^w(i) < 0$ so that $r_t(i) - r_{t-1}(i) < a$. Note that, since a good estimator of $T_t^w(i)$ is not easily obtainable when the queue detector is occupied, $K_1 T_w$ can be replaced by $a$. Fig. 3 summarizes the control logic.

TABLE I
MEASUREMENTS OF EFFECTIVENESS FOR SEPT. 25, 2008, 2:00–8:00 P.M.

| PERFORMANCE MEASURES | | | SZM | New Metering | Perc Incr over SZM |
|---|---|---|---|---|---|
| Total Travel Time | veh-hrs | Mainline | 5692 | 5605 | -1.53% |
| | | Ramp | 430 | 341 | -20.63% |
| | | Entire Site | 6164 | 5980 | -2.99% |
| Total Travel | veh-mi | Entire Site | 290187 | 288864 | NA |
| Total Delay | veh-hrs | Mainline | 1427 | 1359 | -4.77% |
| | | Ramp | 276 | 196 | -29.01% |
| | | Entire Site | 1714 | 1566 | -8.64% |
| Average Delay | min/veh | Mainline | 1.27 | 1.26 | -0.79% |
| | | Ramp | 0.39 | 0.28 | -28.21% |
| | | Entire Site | 1.66 | 1.54 | -7.23% |
| Total Stop Time | veh-hrs | Mainline | *181* | *159* | *-12.15%* |
| Total No of Stops | number | Mainline | 114015 | 97770 | -14.25% |
| No of Stops per Veh | #/veh | Mainline | 1.75 | 1.56 | -10.86% |
| Speed | mph | Mainline | 54.69 | 54.79 | 0.18% |

TABLE II
RAMP MEASUREMENTS OF EFFECTIVENESS FOR THE PEAK OF SEPT. 25, 2008, 4:00–6:30 P.M.

| 4:00 - 6:30 | Avg Ramp Delay (minutes) | | Max Ramp Wait (minutes) | | Total Ramp Delay (vehicle-hours) | | Average Queue Size (vehicles) | | Maximum Queue (vehicles) | |
|---|---|---|---|---|---|---|---|---|---|---|
| Ramp | SZM | New | SZM | New | SZM | New | SZM | New | SZM | New |
| Valley View | 0.29 | 0.43 | 1.65 | 1.43 | 8.78 | 13.10 | 2 | 3 | 12 | 17 |
| TH 62 | 0.50 | 0.68 | 2.12 | 5.03 | 29.23 | 39.59 | 4 | 5 | 25 | 27 |
| Bren Road | 0.86 | 1.06 | 4.51 | 3.88 | 24.42 | 30.11 | 5 | 6 | 24 | 30 |
| Lincoln Drive | 0.43 | 0.07 | 1.10 | 0.52 | 2.77 | 0.45 | 3 | 1 | 7 | 1 |
| Excelsior Blvd | 0.73 | 0.22 | 3.47 | 0.82 | 16.89 | 4.99 | 4 | 1 | 25 | 8 |
| TH7 | 0.14 | 0.17 | 1.56 | 0.63 | 3.48 | 4.35 | 3 | 1 | 15 | 7 |
| 36th Street | 0.29 | 0.07 | 2.00 | 0.10 | 3.84 | 0.91 | 2 | 1 | 11 | 1 |
| Minnetonka Blvd | 0.27 | 0.07 | 2.33 | 0.23 | 2.61 | 0.70 | 2 | 1 | 9 | 1 |
| I-394 | 3.66 | 1.82 | 3.75 | 2.55 | 19.04 | 9.49 | 3 | 1 | 13 | 5 |
| Betty Crocker Dr | 0.56 | 0.52 | 4.00 | 4.10 | 9.02 | 8.38 | 5 | 2 | 18 | 21 |
| TH55 | 0.51 | 0.12 | 4.06 | 0.70 | 12.32 | 2.91 | 3 | 1 | 13 | 5 |
| Plymouth Ave | 1.29 | 0.07 | 4.10 | 0.37 | 20.83 | 1.12 | 5 | 1 | 17 | 6 |
| Medicine Lake Rd | 1.11 | 0.12 | 4.29 | 0.50 | 12.17 | 1.30 | 3 | 1 | 20 | 5 |
| TOTAL | | | | | 165.39 | 117.41 | | | | |

It contains decision charts to define the state, the metering strategy, and the steps of the algorithm.

## IV. APPLICATION AND RESULTS

The developed strategy is tested through an AIMSUN microsimulator on a 12-mi segment of Trunk Highway 169 northbound (TH-169 NB), starting from the I-494 interchange and ending at the 63rd Avenue North (see Fig. 2). This site is a circumferential freeway traversing the Twin Cities west metropolitan region. It includes ten weaving sections, four high-occupancy-vehicle bypass ramps, 24 entrance ramps (17 metered), and 25 exit ramps. Among the metered ramps, 15 local access ramps and two freeway-to-freeway ramps connect TH-62 and I-394, respectively (see Fig. 4). The upstream and downstream boundaries are uncongested. The 1-min-interval traffic demand data used in the simulation were extracted from the Mn/DOT loop detector database. For the purpose of this study, only one peak period was tested, specifically the afternoon peak on September 25, 2008, between 2 P.M.

and 8 P.M.. The parameters of the simulator for this network have been calibrated in previous studies [29]. Fig. 5 summarizes the demand for a subset of the on-ramps of the study site. The measures of effectiveness (MOEs) selected for the preliminary assessment are given as follows:

1) delay (in vehicle hours) on the mainline, ramps, and entire system (freeway and ramps);
2) total travel (in vehicle miles) on the mainline, ramps, and entire system (freeway and ramps);
3) average speed on the mainline;
4) total number of stops on the mainline;
5) average number of stops on the mainline;
6) mean queue size on ramps;
7) maximum queue size on ramps;
8) mean ramp wait times;
9) maximum ramp wait times.

The simulation results are summarized in Tables I and II, which present the aforementioned MOEs. In both tables, the base case for the comparison is the current SZM strategy. From Table I, it can be seen that, under the new strategy, both

TABLE III
MOEs OBSERVED FOR VARIOUS VALUES OF $T$ FOR ALL RAMPS AND FOR TH62 RAMP

| Ramp wait constraint | | TH62 | 5 mins | 5 mins | 5 mins | 5 mins | 3 mins | 4 mins | 5 mins | 6 mins |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Other Ramps | 3 mins | 4 mins | 5 mins | 6 mins | 4 mins | 4 mins | 4 mins | 4 mins |
| Total Travel Time | veh-hrs | Mainline | 5726 | 5605 | 5675 | 5750 | 5464 | 5598 | 5605 | 5669 |
| | | Ramp | 305 | 341 | 367 | 502 | 705 | 367 | 341 | 314 |
| | | Entire Site | 6064 | 5980 | 6075 | 6286 | 6202 | 5998 | 5980 | 6017 |
| Total Travel | veh-mi | Mainline | 280962 | 281113 | 280782 | 281250 | 280474 | 280614 | 281113 | 280846 |
| | | Ramp | 6539 | 6539 | 6539 | 6539 | 6539 | 6539 | 6539 | 6539 |
| | | Entire Site | 288713 | 288864 | 288533 | 289002 | 288225 | 288366 | 288864 | 288597 |
| Total Delay | veh-hrs | Mainline | 1484 | 1359 | 1433 | 1500 | 1228 | 1359 | 1359 | 1430 |
| | | Ramp | 159 | 196 | 221 | 356 | 559 | 221 | 196 | 168 |
| | | Entire Site | 1654 | 1566 | 1665 | 1867 | 1798 | 1590 | 1566 | 1609 |
| Average Delay | min/veh | Mainline | 1.37 | 1.26 | 1.32 | 1.38 | 1.13 | 1.25 | 1.26 | 1.32 |
| | | Ramp | 0.23 | 0.28 | 0.32 | 0.51 | 0.8 | 0.32 | 0.28 | 0.24 |
| | | Entire Site | 1.6 | 1.54 | 1.64 | 1.89 | 1.93 | 1.57 | 1.54 | 1.56 |
| Total Stop Time | veh-hrs | Mainline | 182 | 159 | 169 | 180 | 120 | 153 | 159 | 167 |
| | | Ramp | 108 | 146 | 172 | 312 | 523 | 172 | 146 | 118 |
| | | Entire Site | 294 | 309 | 344 | 496 | 647 | 329 | 309 | 288 |
| Total Stops | # | Mainline | 114598 | 97770 | 108502 | 114628 | 78511 | 97002 | 97770 | 104908 |
| | | Ramp | 21024 | 21941 | 21926 | 22717 | 23318 | 21925 | 21941 | 21478 |
| | | Entire Site | 136957 | 121024 | 131725 | 138665 | 103129 | 120235 | 121024 | 127695 |
| Stops per Veh | #/veh | Mainline | 1.82 | 1.56 | 1.72 | 1.8 | 1.25 | 1.54 | 1.56 | 1.67 |
| | | Ramp | 0.502 | 0.524 | 0.523 | 0.542 | 0.557 | 0.523 | 0.524 | 0.513 |
| | | Entire Site | 2.322 | 2.084 | 2.243 | 2.342 | 1.807 | 2.063 | 2.084 | 2.183 |
| Speed | mph | Mainline | 54.32 | 54.79 | 54.49 | 54.29 | 55.46 | 54.78 | 54.79 | 54.49 |
| | | Ramp | 28.37 | 27.65 | 27.44 | 26.7 | 26.13 | 27.44 | 27.65 | 28.11 |
| | | Entire Site | 52.17 | 52.57 | 52.31 | 52.11 | 53.11 | 52.55 | 52.57 | 52.32 |

mainline and ramp delay during the study period are greatly reduced. Specifically, the total ramp delay decreased by about 29% with the new ramp control strategy, and the mainline delay decreased by about 5%; the total delay for the entire system dropped by almost 9%. This finding suggests that, in this case, the new strategy is effective since it reduces not only ramp delay but total system delay as well. Similarly, under the new control strategy, the total travel time on the mainline decreased by 1.5%, the ramp total travel time dropped by nearly 20%, and the total system (freeway and ramp) travel time decreased by about 3%. Decreased mainline congestion is also evidenced by the decrease in the total number of mainline stops per vehicle and the total stop time, as can be seen from Table I.

Table II presents the effects of the stratified ramp control strategy on all TH169NB metered entrance ramps for the peak period 4:00–6:30 P.M. The results clearly indicate that the new control strategy is very effective in keeping ramp wait times below the maximum allowed and in reducing ramp delay time. Another interesting observation made by analyzing the simulation results is that the new strategy substantially reduces ramp queues, whereas the overall ramp delay for the peak period was reduced by nearly 30%.

A simple sensitivity analysis was carried out on the parameter that indicates the ramp constraint wait times for all the ramps. For the original reported results (see Tables I and II), the ramp constraint wait time for all ramps along the network was set to 4 min (the recommended value used by Mn/DOT). The constraint was relaxed at the TH62 ramp, which reported the highest on-ramp delays in simulation, to 5 min. For the sensitivity analysis, simulation runs were made first for varying allowable ramp delay constraints on all ramps while keeping

the delay constraint constant for the TH62 ramp. Table III lists the system MOEs obtained for this set of test cases. Similarly, a set of simulation runs was also made for varying values of delay constraint on the TH62 ramp while keeping all other ramps constant at 4 min, with Table III listing the resulting MOEs for the cases. From both the test sets, we see that the ramp travel times and wait times increase with an increase in allowable ramp wait times, which follows naturally, while mainline and system-wide performance parameters such as total delay, total stop time, and speed first improve and then worsen with increasing delay constraints. The optimal value for overall system MOEs is $T = 4$ min. This value ensures peak system performance with smooth mainline flow of traffic and not high enough to overly burden the ramps themselves, causing high delays.

The new strategy successfully reaches the objectives of decreasing the length of congested period in the mainline. Fig. 6 shows a comparison of the SZM versus the new strategy through contour plots of mainline density. It is evident that, in many locations, the onset of the breakdown is postponed, and the system recovered faster by 5–10 min.

## V. CONCLUSION

A new-generation coordinated traffic-responsive ramp metering algorithm has been designed for Minnesota's freeways based on density measurements, rather than flow rates. This was motivated in view of empirical results from observations, indicating that, when recurrent traffic breakdowns occur, capacity significantly varies, even under the same operating conditions and that the critical density value for this stochastic
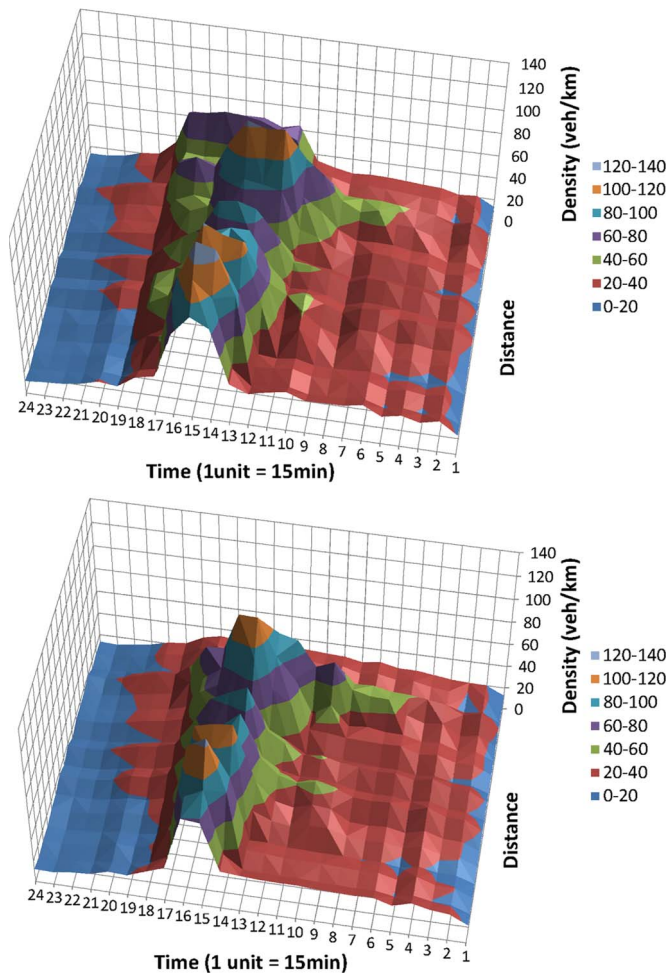
Fig. 6. Density contour plots for (top) SZM and (bottom) new metering strategy.

capacity is more stable. Indeed, our test results have indicated that the density-based feedback coordinated algorithm succeeds in delaying the onset of breakdown, accelerating system recovery when ramp metering is unable to maintain uncongested conditions in the mainline due to the violation of maximum allowable ramp waiting time, and improving freeway and ramp performance, compared with the current deployed algorithm (SZM). To study the effects of varying traffic demand patterns, occurrence of incidents, detector malfunctions, and other relevant issues, additional sites are being selected for testing. This additional testing will allow more in-depth evaluation, parameter optimization, and potential improvements for the new strategy. A field test implementation of the resulting control strategy is currently being planned for refinements and large-scale deployment in the Twin Cities freeways, as simulations are still not accurately replicate real-world conditions.

## ACKNOWLEDGMENT

## REFERENCES

[1] M. Papageorgiou, H. Haj-Salem, and J. Blosseville, "ALINEA: A local feedback control law for on-ramp metering," *Transp. Res. Rec.*, no. 1320, pp. 58–64, 1991.
[2] M. Papageorgiou and A. Kotsialos, "Freeway ramp metering: An overview," *IEEE Trans. Intell. Transp. Syst.*, vol. 3, no. 4, pp. 271–281, Dec. 2002.
[3] L. Jacobson, K. Henry, and O. Mehyar, "Real time metering algorithm for centralized control," *Transp. Res. Rec.*, no. 1232, pp. 17–26, 1989.
[4] Y. Stephanedes, "Implementation of on-line zone control strategies for optimal ramp metering in the minneapolis ring road," in *Proc. 7th Int. Conf. Road Traffic Monitoring Control*, 1994, pp. 181–184.
[5] M. Papageorgiou, H. Haj-Salem, and J. Bloseville, "Modeling and real time control of traffic flow on the southern part of the Boulevard Peripherique in Paris—Part II: Coordinated on-ramp metering," *Transp. Res. Part A.*, vol. 24, no. 5, p. 361, 1990.
[6] L. E. Lipp, L. J. Corcoran, and G. A. Hickman, "Benefits of central computer control for Denver ramp metering system," *Transp. Res. Rec.*, no. 1320, pp. 3–6, 1991.
[7] G. Paesani, J. Kerr, P. Perovich, and F. E. Khosravi, "System wide adaptive ramp metering (SWARM)," in *Proc. Amer. 7th Annu. Meeting ITS*, Washington, DC, 1997 [CDROM].
[8] Y. J. Stephanedes and K. K. Chang, "Optimal control of freeway corridors," *J. Transp. Eng.*, vol. 119, no. 4, pp. 504–514, Jul./Aug. 1993.
[9] H. M. Zhang and W. W. Recker, "On optimal freeway ramp control policies for congested traffic corridors," *Transp. Res. Part B, Method*, vol. 33, no. 6, pp. 417–436, 1999.
[10] A. Kotsialos and M. Papageorgiou, "Nonlinear optimal control applied to coordinated ramp metering," *IEEE Trans. Control Syst. Technol.*, vol. 12, no. 6, pp. 920–933, Nov. 2004.
[11] G. Gomes and R. Horowitz, "Optimal freeway ramp metering using the asymmetric cell transmission model," *Transp. Res. Part C*, vol. 14, pp. 244–262, 2006.
[12] I. Papamichail and M. Papageorgiou, "Traffic-responsive linked ramp-metering control," *IEEE Trans. Intell. Transp. Syst.*, vol. 9, no. 1, pp. 111–121, Mar. 2008.
[13] D. Zhao, X. Bai, F.-Y. Wang, J. Xu, and W. Yu, "DHP method for ramp metering of freeway traffic," *IEEE Trans. Intell. Transp. Syst.*, vol. 12, no. 4, pp. 990–999, Dec. 2011.
[14] K. Bogenberger and A. D. May, "Advanced coordinated traffic responsive ramp metering strategies," California PATH Working Paper, Berkeley, CA, 1999, UCB-ITS-PWP-99-19.
[15] M. Zhang, T. Kim, X. Nie, W. Jin, L. Chu, and W. Recker, "Evaluation of on-ramp control, algorithms," California PATH Working Paper, Berkeley, CA, 2001, UCB-ITS- PRR-2001-36.
[16] M. Papageorgiou, "A hierarchical control system for freeway traffic," *Transp. Res. B*, vol. 17, no. 3, pp. 251–261, Jun. 1983.
[17] A. May, "A proposed dynamic freeway control system hierarchy," in *Proc. 3rd Int. Symp. Control Transp. Syst., IFAC/IFIP/IFORS*, Columbus, OH, 1976, pp. 1–12.
[18] B. Feng, J. Hourdos, and P. Michalopoulos, "Improving minnesota's stratified ramp control strategy," *Transp. Res. Rec.*, vol. 1959, no. 1, pp. 77–83, 2006.
[19] H. Liu, X. Wu, and P. Michalopoulos, "Improving Queue size estimation for Minnesota's stratified zone metering strategy," *Transp. Res. Rec.*, vol. 2012, pp. 38–46, 2007.
[20] L. Eleftheriadou, R. Roess, and W. Shane, "Probabilistic nature of breakdown at freeway merge junctions," *Transp. Res. Rec.*, vol. 1484, pp. 80–89, 1995.
[21] R. Kuhne and R. Mahnke, "Controlling traffic breakdowns," in *Proc. 16th ISTTT*, 2005, pp. 229–244.
[22] M. Cassidy and J. Rudjanakanoknad, "Increasing the capacity of an isolated merge my metering its on-ramp," *Transp. Res. B, Methodol.*, vol. 39, no. 10, pp. 896–913, Dec. 2005.
[23] L. Zhang and D. Levinson, "Some properties of flows at freeway bottlenecks," *Transp. Res. Rec.*, vol. 1883, no. 1, pp. 122–131, 2004.
[24] W. Xin, P. Michalopoulos, J. Hourdakis, and D. Lau, "Minnesota's new ramp control strategy: Design overview and preliminary assessment," *Transp. Res. Rec.*, vol. 1867, pp. 69–79, 2004.
[25] L. Eleftheriadou, W. Brilon, F. Hall, B. Persaud, S. Washburn, and A. Kondyli, "Proactive ramp metering based on breakdown probabilities," in *Proc. 6th Int. Symp. Highway Capacity*, Stockholm, Sweden, 2011.
[26] P. G. Michalopoulos, D. E. Beskos, and J. K. Lin, "Analysis of interrupted traffic flow by finite difference methods," *Transp. Res. B*, vol. 18, pp. 409–421, 1984.

[27] N. Geroliminis, A. Srivastava, and P. G. Michalopoulos, "Experimental observations of capacity drop phenomena in freeway merges with ramp metering control and integration in a first-order model," in *Proc. 90th Trans. Res. Board Conf.*, Washington, DC, 2011.

[28] I. Papamichail and M. Papageorgiou, "Balancing of Queues or waiting times on metered dual-branch on-ramps," *IEEE Trans. Intell. Transp. Syst.*, vol. 12, no. 2, pp. 438–452, Jun. 2011.

[29] B. Feng, J. Hourdos, and P. Michalopoulos, "Improving Minnesota's stratified ramp control strategy," *Transp. Res. Rec.*, vol. 1959, pp. 77–83, 2006.

**Anupam Srivastava** received the Bachelor's degree in civil engineering from Indian Institute of Technology, Kharagpur, India, in 2005. He is currently working toward the M.Sc. degree with the Department of Civil Engineering, University of Minnesota, Minneapolis, where he works under the supervision of Prof. N. Geroliminis.

His research interests include traffic flow modeling, bottleneck analysis, freeway ramp metering, and traffic flow simulation.

**Nikolas Geroliminis** received the Diploma degree in civil engineering from National Technical University of Athens, Athens, Greece, and the M.Sc. and Ph.D. degrees in civil engineering from the University of California, Berkeley.

Before joining the Ecole Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland, he was an Assistant Professor with the faculty of the Department of Civil Engineering, University of Minnesota, Minneapolis. He is a member of the Transportation Research Board's Traffic Flow Theory Committee. In summer 2009, he was an Academic Visitor with the Chair of Sociology, in particular of Modeling and Simulation with the Swiss Federal Institute of Technology Zurich, Switzerland. He is currently the Director of the Urban Transport Systems Laboratory and an Assistant Professor with the School of Architecture, Civil and Engineering, EPFL. His research interests include urban transportation systems, traffic flow theory and control, public transportation, and logistics.

Dr. Geroliminis was the recipient of the University of California Transportation Student of the Year Award in 2007 and the Outstanding Graduate Student Instructor Award in 2006.

**Panos Michalopoulos** received the B.S. degree in civil engineering from the College of Engineering, Athens, Greece, and the M.S. and Ph.D. degrees in civil engineering with specializations in transportation from the University of Florida, Gainesville.

He is currently with the Department of Civil Engineering, University of Minnesota, Minneapolis. He has more than 30 years of experience in traffic engineering, management, operations, and control. His research led to the development of AUTOSCOPE, which is the first and most widely used machine-vision-based vehicle detection and surveillance system, with more than 20 000 installations worldwide since 1993. His recent work has included continuum modeling of traffic flow dynamics in complex large-scale networks; he is a leader in the development of a wide-area vehicle detection system based on machine vision for testing and validating traffic flow models, adaptive (real-time) control, and automatic incident detection. He has presented or published more than 150 papers in proceedings and peer-reviewed journals, in addition to numerous research reports. He is the holder of U.S. and E.U. patents for machine-vision-based vehicle detection systems. His research interests include the modeling of traffic flow dynamics, traffic flow theory, traffic simulation and control, and the development of advanced traffic management systems.