



Cross-regional cooperative multiple-timescale joint traffic control, signal control, and routing for maximum network throughput

Downloaded from: <https://research.chalmers.se>, 2026-07-01 10:04 UTC

Citation for the original published paper (version of record):

Cui, S., Xue, Y., Gao, K. et al (2026). Cross-regional cooperative multiple-timescale joint traffic control, signal control, and routing for maximum network throughput. *Transportation Research Part B: Methodological*, 211. <http://dx.doi.org/10.1016/j.trb.2026.103519>

N.B. When citing this work, cite the original published paper.



Cross-regional cooperative multiple-timescale joint traffic control, signal control, and routing for maximum network throughput

Shaohua Cui ^{a,b}, Yongjie Xue ^c, Kun Gao ^{a,*}, Bin Yu ^{b,c}, Xiaobo Qu ^d

^a Department of Architecture and Civil Engineering, Chalmers University of Technology, Gothenburg, 41296, Sweden

^b Key Laboratory of Intelligent Transportation Technology and System, Ministry of Education, Beijing, 100191, PR China

^c School of Transportation Science and Engineering, Beihang University, Beijing, 100191, PR China

^d State Key Laboratory of Automotive Safety and Energy, School of Vehicle and Mobility, Tsinghua University, Beijing, 100084, PR China

ARTICLE INFO

Keywords:

Lyapunov optimization
Distributed control
Network stability
Max-pressure control

ABSTRACT

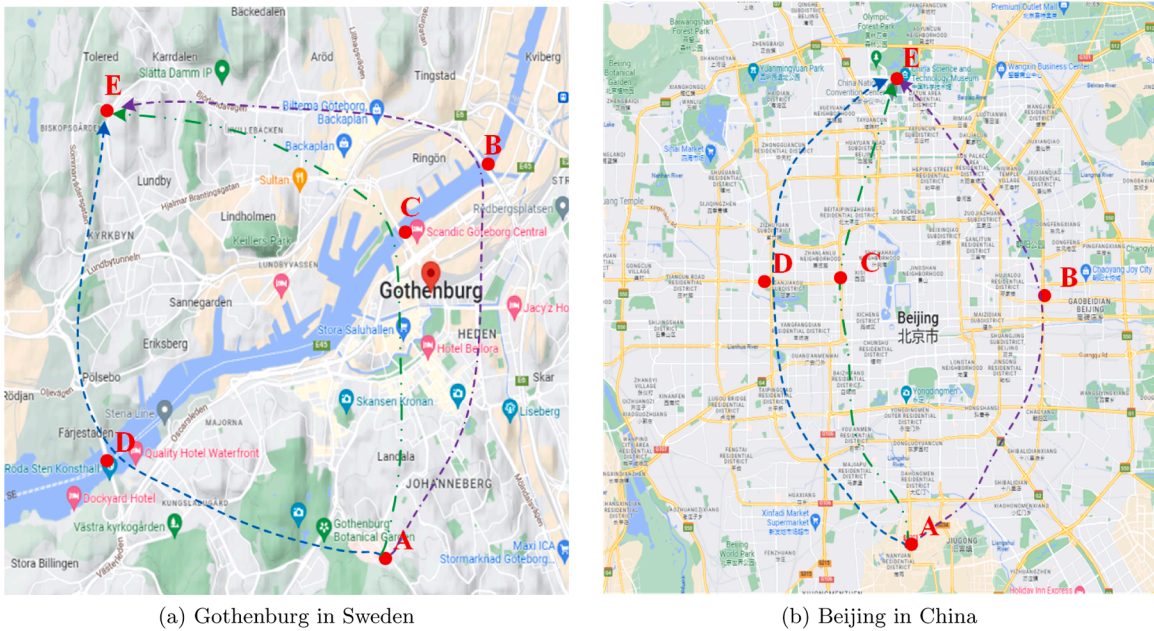
Network-level signal control and routing were independently optimized using centralized algorithms that rely on traffic data from the entire network. These methods cannot guarantee provable network performance under joint control and real-time controller updates. Additionally, they were not integrated with traffic control, which may result in low network performance due to oversaturation. This study proposes single-timescale and dual-timescale joint traffic control, signal control, and routing based on cross-regional cooperation for urban networks of varying scales and communication delays. Urban networks are divided into regions based on road structure and regional functions to narrow down the traffic data required for network-level joint control. Traffic control is implemented collaboratively at entry nodes and inter-regional connection nodes with different timescales to balance the traffic at inter-regional connection nodes via demand allocation and to avoid regional oversaturation by limiting inter-regional traffic transmission rates. Real-time signal control and routing are implemented at intersections to optimize traffic priority and route vehicles efficiently. Both joint control algorithms are fully decomposable and distributed. Using Lyapunov drift functions, this study proves that both the single-timescale and dual-timescale joint control algorithms maximize network throughput, although the latter results in higher average delays. Comparative simulations on a grid-like network and a network with non-uniform link density show that the joint control algorithms improve network throughput by up to 19.3% compared to existing algorithms.

1. Introduction

Traffic congestion has led to significant environmental, economic, and social losses (Palma et al., 2022; Ding et al., 2022; Espadaler-Clapés et al., 2023; Fattah et al., 2022; Huang and Loo, 2022; Han et al., 2021; Marshall and Dumbaugh, 2020). It contributes to 20.3% of all greenhouse gas emissions in Europe (Nugmanova et al., 2019). The economic losses attributed to traffic congestion are estimated to be 0.47% of GDP in Ireland, 2.95% in the Netherlands, and 7% – 9% in Russia (Sakhapov et al., 2016). In Australia's eight capital cities, the total social cost of traffic congestion is estimated at \$16.5 billion (Kuang et al., 2019). Adaptive signal control is an effective way to alleviate traffic congestion because it relies less on historical traffic data and can capture the stochastic behavior of vehicles on links and at intersections (Ampountolas et al., 2020; Barman and Levin, 2023; Gao and Zhang, 2022). Consequently, adaptive

* Corresponding author.

E-mail address: gkun@chalmers.se (K. Gao).



(a) Gothenburg in Sweden

(b) Beijing in China

Fig. 1. The networks consisting of multiple link-dense regions connected by bridges or expressways.

signal control, particularly at the network level, has garnered attention from transportation departments in various countries (Guo and Ban, 2023; Tsitsokas et al., 2023). However, collecting traffic data from numerous urban intersections is already challenging. Utilizing this data for centralized phase optimization in network-level adaptive signal control is difficult, even without integrating real-time routing for tens of thousands of vehicles. Although some map navigation software attempts to integrate the signal phases of all intersections into vehicle routing, they still independently optimize routing for each traveler to minimize travel costs. This makes it challenging to analytically prove network performance under the combined effects of independent vehicle routing and centralized phase optimization.

Alleviating traffic congestion requires collaboration between routing and signal control, cooperation between internal and edge nodes in large-scale networks, effective traffic management, and consideration of diversity in network nodes and links. These issues are particularly severe in rapidly expanding metropolitan areas, where urban networks gradually evolve into multiple dense regions connected through several key bridges (e.g., Gothenburg in Sweden, as shown in Fig. 1(a)) or expressways (e.g., Beijing in China, as shown in Fig. 1(b)). For example,

- if cross-regional traffic is not reasonably distributed across internal bottleneck nodes (e.g., determining whether a vehicle at node A should pass through bridges/expressways B, C, or D to reach destination E in Fig. 1), local control may fail to balance traffic and prevent regional oversaturation;
- if there is no effective transmission rate control on inter-regional connection nodes/links (e.g., bridges/expressways B, C, or D in Fig. 1), local control may fail to improve network performance under regional oversaturation;
- if routing algorithms treat network links and nodes as homogeneous (e.g., dynamic traffic assignment, A^* algorithm, and Dijkstra algorithm), they may not be able to route all vehicles in real time.

Dividing the network into regions of varying scales based on network structure or regional functional differences can effectively reduce computational time and traffic data collection, while better exploiting the heterogeneity of nodes and links.

Hence, this study investigates cross-regional cooperative joint traffic control, signal control, and routing, ensuring provable network performance. Network partitioning can be based on road structure, regional functions, or optimal decomposition algorithms that improve computational efficiency. Traffic control is implemented collaboratively at entry nodes and inter-regional connection nodes. At entry nodes, traffic control allocates cross-regional traffic to inter-regional connection nodes to balance their traffic. At inter-regional connection nodes, it limits the traffic transmission rate between regions to avoid local oversaturation. Signal control and routing are implemented at intersections to optimize signal priority and guide vehicles to their destinations. This study provides three key contributions: 1) It constructs and proves region-based flow conservation constraints, forming a basis for cross-regional cooperation control and the coordination between entry nodes and bottleneck nodes in large-scale networks; 2) It proposes single-timescale and dual-timescale joint traffic control, signal control, and routing algorithms based on cross-regional cooperation for networks of varying scales and communication delays; 3) Using Lyapunov drift functions, it analytically proves that both single-timescale and dual-timescale joint control algorithms maximize network throughput, with the latter resulting in higher average delays.

The remaining sections are organized as follows: [Section 2](#) reviews the relevant literature and summarizes research gaps; [Section 3](#) describes the joint control problem studied in detail; [Section 4](#) proposes the single-timescale and dual-timescale joint control algorithms and provides performance analysis; [Section 5](#) presents comparative numerical simulations for two networks with different scales to verify the effectiveness of the proposed joint control algorithms; and [Section 6](#) concludes the paper and suggests future research directions.

2. Literature review

Adaptive signal control is an effective way to alleviate congestion. Classic optimization-based adaptive signal control algorithms include SCATS (Luk, 1984), RHODES (Mirchandani and Head, 2001), SCOOT (Hunt et al., 1982), TUC (Diakaki et al., 2002), OPAC (Gartner, 1983), and PROLYN (Henry et al., 1983). As the network scale increases, the number of decision variables in optimization models grows exponentially (Chow et al., 2020a,b). Therefore, achieving cooperative control among intersections in a large-scale network using these algorithms becomes challenging. To address this, adaptive signal control with a data-driven solution structure was proposed by (Mo et al., 2022) and Su et al. (2021) to reduce the number of decision variables. However, data-driven algorithms rely heavily on accurate historical and real-time traffic data, which limits scalability (Guo et al., 2019). For detailed discussions on adaptive signal control, readers are referred to Guo et al. (2019) and Papageorgiou et al. (2003). To achieve real-time updates of signal phases and reduce the need for traffic data collection, Varaiya (2013) and Wongpiromsarn et al. (2012) proposed the distributed adaptive max-pressure (MP) signal control, also known as back-pressure signal control. MP signal control relies on the queue lengths on links adjacent to intersections at the current time step for signal phase selection and provably maximizes network throughput (Chen et al., 2020a; Dixit et al., 2020; Gao and Zhang, 2022; Levin et al., 2020, 2019; Rey and Levin, 2019; Tsitsokas et al., 2021, 2022; Wu et al., 2018).

MP signal control requires known mean turn ratios (i.e., routing proportions at intersections). To address this dependency, Gregoire et al. (2014) and Lioris et al. (2016) used the queue lengths on links adjacent to intersections to estimate turn ratios in real time. Le et al. (2017) proposed utility-optimization-based back-pressure traffic control, which directly optimizes the expected turn ratios to stabilize all queues adjacent to intersections and reduce the spatial traffic heterogeneity. Similarly, Gregoire et al. (2016) optimized the turn ratios of controllable vehicles in queues. Unlike optimizing turn ratios, some studies have focused on route optimization for individual vehicles. Liu et al. (2018) and Zaidi et al. (2016) introduced the concept of a virtual network to adaptively route each vehicle. In the virtual network, vehicles are queued separately at intersections based on their destinations. Since virtual queues are used as counters solely for signal optimization and routing, destination-based queues are feasible. Numerical simulations showed that the introduced virtual queues effectively reduce travel delays at high demand but result in high travel delays at low demand. Further simulations indicated that destination-based virtual queues increase travel delays at low demand by causing unnecessary detours rather than lowering average speed. Consequently, destination-based virtual queues do not effectively restrict the feasible path set at low demand.

None of the above extended MP signal control algorithms are integrated with traffic demand control to limit the exogenous demand entering the network and the traffic transmission rate between regions. The macroscopic fundamental diagram (MFD) theory shows that network throughput increases with traffic accumulation up to a critical point, after which additional vehicle accumulation in the network cause a significant decrease in network throughput (Daganzo, 2007). This means that network-level control without integrating traffic control is not sufficient to alleviate congestion and increase network throughput (Moradi et al., 2022). Geroliminis et al. (2013)) and Haddad et al. (2013) proposed perimeter control at the regional level based on MFDs to control the traffic transfer rates between regions. Through perimeter control, they ensured that vehicle accumulation in different regions approaches the critical accumulation level in MFDs to increase network throughput. To improve mobility in networks, route guidance has been integrated with perimeter control to manage outbound traffic from one region to its adjacent regions, with the goal of minimizing total time spent (Guo and Ban, 2020; Hou and Lei, 2022; Ingole et al., 2020; Sirmatel and Geroliminis, 2018), average time spent (Fu et al., 2022; Menelaou et al., 2021), and average delay (Ding et al., 2017), or maximizing trip completion rates (Liu et al., 2018). However, the design of MFDs requires that traffic is spatially homogeneous, which is difficult to guarantee in reality. Moreover, macro-level perimeter control cannot capture the stochastic traffic dynamics at intersections and is not suitable for controlling exogenous demand inputs from entry nodes. Therefore, macro-level perimeter control at inter-regional connection nodes may still cause network oversaturation and local congestion.

To the best of our knowledge, the joint framework of traffic control, signal control, and routing has not been studied in the transportation field. Related studies in wireless communication networks is called joint rate control, scheduling, and routing. This joint control is a decomposition of utility maximization problems subject to rate and scheduling constraints. Eryilmaz and Srikant (2006) and Vo et al. (2011) adopted primal-dual theory, while Aljubayri et al. (2021), Chen et al. (2006), Qu et al. (2015), and Eryilmaz and Srikant (2006) employed dual theory to decompose utility maximization problems into three subproblems: congestion/rate control, scheduling, and routing. This decomposition ensures that the solutions to these optimization problems are completely distributed. The three subproblems are all based on the lengths of destination-based queues at nodes. Scheduling is usually a MP control problem and is similar to MP signal control. Scheduling determines the links that can be activated simultaneously in the wireless network. The routing component then forwards data packets in the queues that are on the activated links and generate the maximum pressure difference at the maximum transmission rate. However, in a traffic network, it is impractical for each intersection to maintain separate queues for all destinations.

According to the above literature, MP signal control is efficient for ensuring provable performance guarantees (e.g., maximizing throughput and utility). However, applying MP signal control to large-scale traffic networks requires addressing the following research

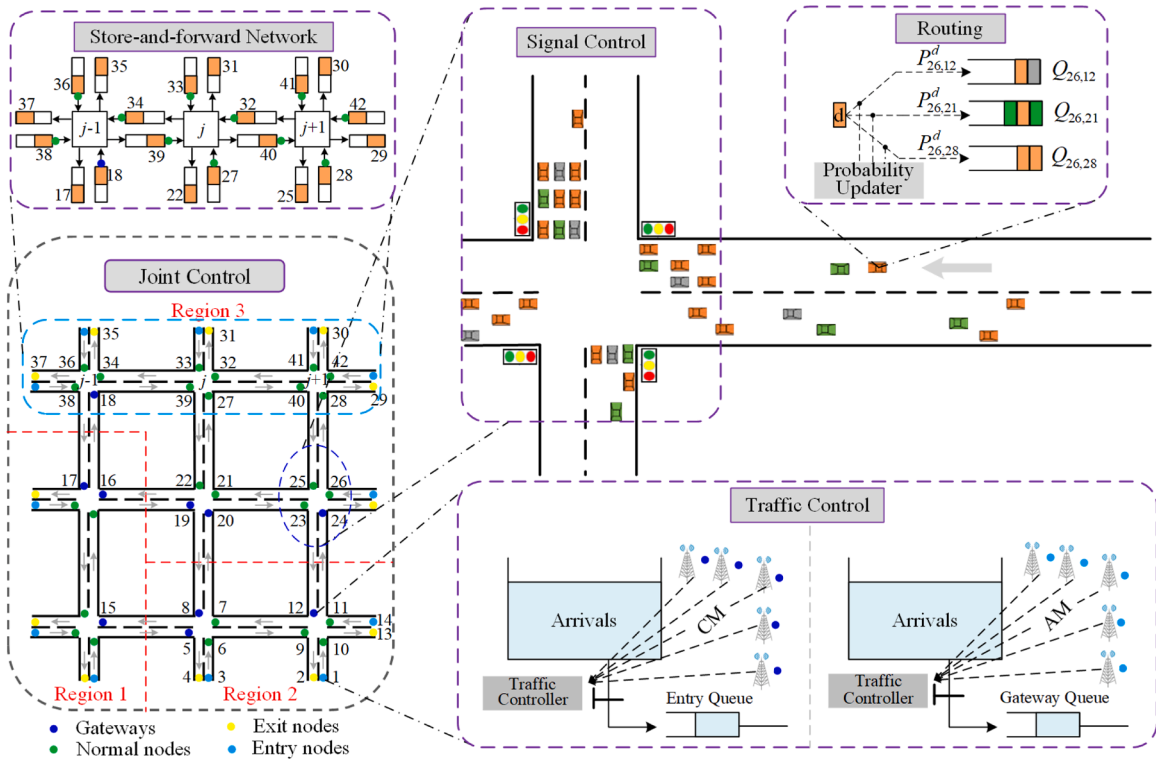


Fig. 2. The cross-regional cooperative joint control consisting of traffic control, signal control, and routing.

gaps: 1) The set of feasible paths should be narrowed down without losing the optimal solution to avoid unnecessary detours for vehicles; 2) In large-scale networks, traffic control coordinated between entry nodes and internal nodes is required to avoid network oversaturation and local congestion; 3) It is necessary to design a micro-level traffic control that can cooperate with signal control and routing without requiring intersections to maintain separate queues for all destinations.

3. Problem description and region-based queueing networks

3.1. Problem description

This study develops the cross-regional cooperative joint control (CRJC) integrating traffic control, signal control, and routing (see Fig. 2) to reduce the required traffic data, ensure real-time controller updates, and achieve provable network performance. According to Cui et al. (2024), Gregoire et al. (2015), and Varaiya (2013), the signalized network is constructed as a store-and-forward queueing network (see the top left of Fig. 2) to capture the stochastic behavior of vehicle arrivals and turn movements at junctions (or intersections). Regional divisions can be based on road structure, regional functions, or optimal network partitioning algorithms that improve computational efficiency. After dividing the network into disjoint regions (see the left side of Fig. 2), CRJC adaptively allocates exogenous demand to inter-regional connection nodes, limits the traffic transmission rate between regions, optimizes signal priority at junctions, and routes vehicles.

Traffic control is performed collaboratively at entry nodes and inter-regional connection nodes. To clearly characterize the role of inter-regional connection nodes in CRJC, these nodes are referred to as gateways in the following text (a concept borrowed from communication networks) to signify that traffic control can be applied here to regulate traffic entering their respective regions. The traffic controller at entry nodes allocates vehicles to the gateway with the lowest congestion level based on the congestion message (CM) sent by all gateways in the region where the vehicle destinations are located. That is, through traffic control at entry nodes, vehicles enter the regions where their destinations are located through the allocated gateways and then proceed to their destinations. Traffic control at entry nodes balances traffic loads at these gateways to improve the utilization of road resources. Traffic control at gateways determines the number of vehicles allowed to enter their regions according to the allocation message (AM) from all entry nodes. Traffic control at gateways prevents local congestion by regulating the inflow of vehicles entering their regions from other regions. Signal control is performed at junctions to optimize signal priority. Routing guides vehicles by determining the probabilities of vehicles choosing the left-turn, straight-through, or right-turn lanes before they reach junctions. The parameters of the three components in the above CRJC are all updated in real time. To distinguish this framework from the multi-timescale joint control algorithm discussed later, we define the above CRJC as single-timescale CRJC (ST-CRJC).

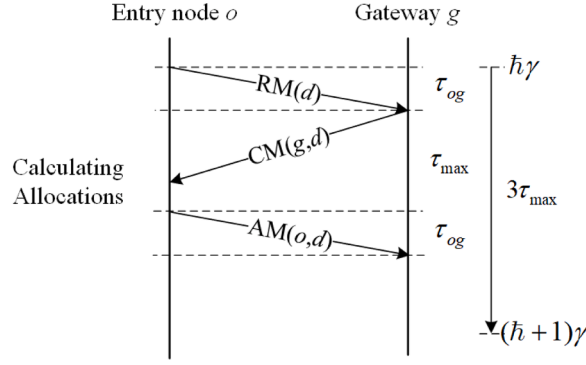


Fig. 3. The three communication delays between entry nodes and gateways.

In a large-scale urban network, there is significant communication delay (e.g., τ_{og}) when transmitting CMs and AMs over long distances between entry nodes and gateways. This delay arises because communication between distant nodes typically uses low-cost standard wireless routers and ordinary 4G modules (Huang et al., 2021). This communication equipment may experience high latency, potentially reaching several seconds, especially during network congestion or under poor signal conditions (Huang et al., 2021). In addition, communication delay also exists when entry nodes send request messages (RMs) or CMs to gateways. These three communication delays are bounded by τ_{max} and illustrated in Fig. 3. As shown in Fig. 3, the maximum update interval of control information for traffic controllers at entry nodes and gateways is less than $3\tau_{max}$. With the evolution of cellular systems from 4G long-term evolution (LTE) to 5G, cellular vehicle-to-everything (C-V2X) is evolving from LTE-V2X to new radio (NR)-V2X (Chen et al., 2020b). Communication delay is reduced from 10 ms in LTE-V2X to 1 ms in NR-V2X (Chen et al., 2020b). Compared with the $3\tau_{max}$ update interval of traffic control, the communication delay associated with transmitting local congestion information between vehicles and junctions for signal control and routing is negligible. The CRJC framework that updates the parameters of traffic control every γ time slots (i.e., $\gamma > 3\tau_{max}$) and updates the parameters of signal control and routing every time slot is called the dual-timescale CRJC (DT-CRJC).

3.2. Assumptions

Based on previous research on MP signal control, such as Cui et al. (2024), Levin et al. (2020), Rey and Levin (2019), and Varaiya (2013), the assumptions adopted in this study are summarized as follows:

- (1) The three components of CRJC algorithms are synchronized independently.
- (2) All vehicles enter the network at entry nodes and leave at exit nodes. Sufficient storage space is available upstream of entry nodes for vehicles that have not yet entered the network to accommodate vehicle queues. Vehicles leave the network immediately upon reaching their destinations.
- (3) All links have separate lanes for straight-through, left-turn, and right-turn movements.
- (4) All flow arrivals are independently and identically distributed (i.i.d.).
- (5) The signalized network operates as a store-and-forward network, ignoring the spatial distribution of queues.
- (6) Communication between vehicles and surrounding signalized intersections is based on LTE-V2X or NR-V2X with a negligible communication delay, while communication between distant nodes uses standard wireless routers or ordinary 4G modules with a non-negligible but bounded communication delay.
- (7) All vehicles comply with the implemented traffic management, without deviation due to personal preferences.

3.3. Region-based queueing network topology

The signalized traffic network is modelled as a directed graph $G = (\mathcal{N}, \mathcal{L})$, where \mathcal{N} and \mathcal{L} represent the sets of nodes and directed links, respectively. Nodes $a \in \mathcal{N}$ represent the roads with queued vehicles. Links $(a, b) \in \mathcal{L}$ denote traffic movements from nodes $a \in \mathcal{N}$ to nodes $b \in \mathcal{N}$. Each signalized junction consists of multiple incoming and outgoing roads with interacting traffic movements. The set \mathcal{J} of junctions is introduced to model the interfering traffic movements at junctions. Let \mathcal{A}_j (respectively, \mathcal{B}_j) denote the set of nodes such that there exist roads from the nodes (respectively, junction $j \in \mathcal{J}$) pointing to the junction $j \in \mathcal{J}$ (respectively, the nodes). Fig. 4 shows one standard junction j . For the junction j in Fig. 4, \mathcal{A}_j equals $\{1, 2, 3, 4\}$ and \mathcal{B}_j equals $\{5, 6, 7, 8\}$. We assume that the vehicle queue for each downstream node at a given node is independent. Let Q_{ab} represent the number of vehicles queued at node $a \in \mathcal{N}$ and waiting to join downstream node $b \in \mathcal{N}$. Let C denote a region, with $C(b)$ representing the region containing node b . \mathcal{G}_C and \mathcal{I}_C denote the sets of gateways and interior nodes of region C , respectively. Let \mathcal{H}_b represent the set of all gateways in the network and the interior nodes of the region containing node b , except for node b itself, i.e., $\mathcal{H}_b = \cup_C \mathcal{G}_C \cup \mathcal{I}_{C(b)} \setminus \{b\}$. The region 2 in Fig. 2 can be denoted as $C(1)$, $C(5)$, or $C(12)$. The sets of gateways and interior nodes of region 2 are $\mathcal{G}_{C(1)} = \{5, 8, 12\}$ and $\mathcal{I}_{C(1)} = \{1, 2, 3, 4, 6, 7, 9, 10, 11, 13, 14\}$. Thus, $\mathcal{H}_1 = \{5, 8, 12, 15, 16, 17, 18, 19, 20, 24\} \cup \mathcal{I}_{C(1)} \setminus \{1\}$. Let $[o, d]$ denote a flow with origin o and destination d . Let \mathcal{F} be the set of all flows $[o, d]$. The arrivals $f_o^d(t)$ of flow $[o, d]$ are i.i.d. over time slots with the arrival rate

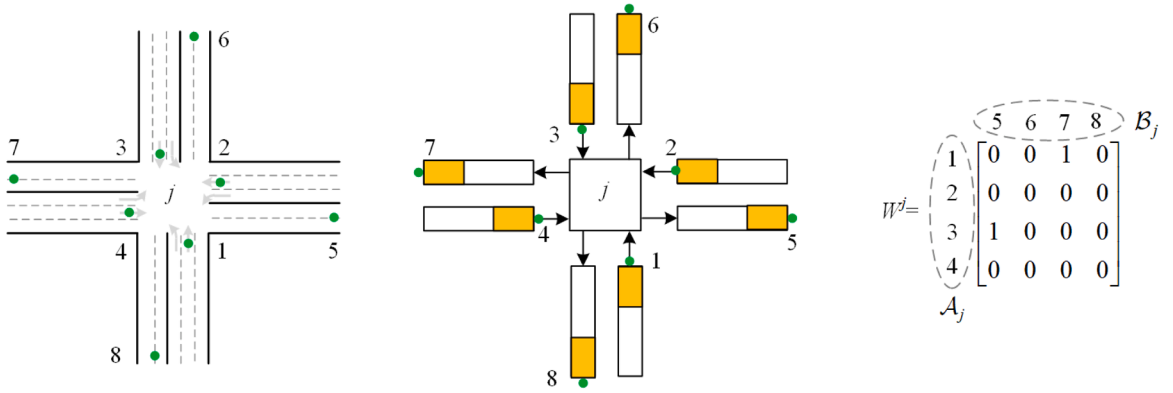


Fig. 4. Signalized traffic network modeling. Intersection layout (left), link-node representation (middle), and the control matrix of the phase with movements (1, 7) and (3, 5).

$E\{f_o^d(t)\} = f_o^d$. Let f denote $\{f_o^d\}$. Let $S(t)$ denote $\{S_{ab}^{[o,d]}(t)\}$, where $S_{ab}^{[o,d]}(t)$ represents the service rate of flow $[o, d]$ over link (a, b) at t .

3.4. Phase and feasible demands

One traffic movement (a, b) at junction $j \in \mathcal{J}$ can be represented as a pair of nodes with $a \in \mathcal{A}_j$ and $b \in \mathcal{B}_j$. Let $\{(a, b) \in \mathcal{A}_j \times \mathcal{B}_j\}$ be the set of traffic movements at junction $j \in \mathcal{J}$. Phase is a term used for the simultaneous actuation of certain subsets of movements at junction $j \in \mathcal{J}$. Each phase is denoted as a control matrix $W^j = \{w_{ab}^j, a \in \mathcal{A}_j, b \in \mathcal{B}_j\}$ with entries $w_{ab}^j = 1$ or 0 depending on whether movement (a, b) is actuated. In Fig. 4, the control matrix of the phase consisting of movements (1, 7) and (3, 5) is shown on the right side. Let \mathcal{W}^j be the set of all such control matrices at junction $j \in \mathcal{J}$. The network control matrix $W = \{w_{ab}\}$, with dimensions $|\mathcal{N}| \times |\mathcal{N}|$, is represented by a collection of all junction control matrices $\{W^j\}$, where $|\mathcal{N}|$ denotes the number of elements in set \mathcal{N} . If the entry w_{ab}^j in the junction control matrix W^j equals 1, then w_{ab} equals 1; otherwise, w_{ab} equals 0. Let \mathcal{W} be the set of all such network control matrices. If movement (a, b) is actuated in time slot t , up to $R_{ab}(t)$ vehicles in queue $Q_{ab}(t)$ can be served at the junction; otherwise, no vehicle in queue $Q_{ab}(t)$ is served. $R_{ab}(t)$ is a stochastic service rate with a mean equal to the saturation/service flow rate r_{ab} and is updated at the beginning of time slot t .

For the infinitely admissible network control sequence $\{W(1), W(2), \dots, W(t), \dots\}$ with $W(t) \in \mathcal{W}$, we define the matrix $\Sigma^W = \{\Sigma_{ab}^W\}$ with entries:

$$\Sigma_{ab}^W = \liminf_T \frac{1}{T} \sum_{t=1}^T w_{ab}(t), \forall a, b \in \mathcal{N} \quad (1)$$

Similarly, for the infinite service rate $S_{ab}^{[o,d]}(t)$ of flow $[o, d]$ over link (a, b) , we define the matrix $S = \{S_{ab}^{[o,d]}\}$ with entries:

$$S_{ab}^{[o,d]} = \liminf_T \frac{1}{T} \sum_{t=1}^T S_{ab}^{[o,d]}(t), \forall (a, b) \in \mathcal{L}, [o, d] \in \mathcal{F} \quad (2)$$

Definition 1. The demand $f = \{f_o^d\}$ is supportable by the admissible network control sequence $\{W(1), W(2), \dots, W(t), \dots\}$ if the following two conditions hold:

$$f_b^d \mathbf{1}_{b=0} + \sum_{a:(a,b) \in \mathcal{L}} S_{ab}^{[o,d]} = \sum_{c:(b,c) \in \mathcal{L}} S_{bc}^{[o,d]}, \forall [o, d] \in \mathcal{F} \quad (3)$$

$$r_{ab} \Sigma_{ab}^W > \sum_{(o,d) \in \mathcal{F}} S_{ab}^{[o,d]}, \forall (a, b) \in \mathcal{L} \quad (4)$$

Eq. (3) denotes the flow-based conservation constraint. Eq. (4) means that the traffic movements $r_{ab} \Sigma_{ab}^W$ served by the control sequence $\{W(1), W(2), \dots, W(t), \dots\}$ per time slot on average are larger than the average demand $\sum_{(o,d) \in \mathcal{F}} S_{ab}^{[o,d]}$ over link (a, b) . By

Definition 1, we know that there exist $\bar{T} < \infty$ and $\chi > 0$ such that:

$$\sum_{t=1}^T r_{ab} w_{ab}(t) - T \sum_{(o,d) \in \mathcal{F}} S_{ab}^{[o,d]} > T\chi, T > \bar{T} \quad (5)$$

Eq. (5) means that the length of queue Q_{ab} at junctions is bounded. The following convex hull of the set \mathcal{W} of network control matrices is defined for the subsequent construction of the set of feasible demands:

$$Co(\mathcal{W}) = \left\{ \sum_{W \in \mathcal{W}} \lambda_W W \mid \lambda_W \geq 0, \sum_{W \in \mathcal{W}} \lambda_W = 1 \right\} \quad (6)$$

The set \mathcal{W} of network control matrices is bounded, so $Co(\mathcal{W})$ is a bounded polytope.

Proposition 1. *The matrix $\Sigma = \{\Sigma_{ab}\}$ belongs to the convex hull $Co(\mathcal{W})$ (i.e., $\Sigma \in Co(\mathcal{W})$) when there is an admissible control sequence $\{W(1), W(2), \dots, W(t), \dots\}$ such that:*

$$\Sigma_{ab} = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T w_{ab}(t), \forall (a, b) \in \mathcal{L} \quad (7)$$

We obtain a matrix $W' = \{w'_{ab}\}$ by setting some entries in $W \in \mathcal{W}$ to zero. W' means that some phases in W are not actuated, so W' still belongs to \mathcal{W} , i.e., $W' \in \mathcal{W}$. If $\Sigma \in Co(\mathcal{W})$ and $0 \leq \Sigma' \leq \Sigma$, where Σ' is obtained by control sequence $\{W'(1), W'(2), \dots, W'(t), \dots\}$ (i.e.,

$$\Sigma'_{ab} = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T w'_{ab}(t)), \text{ then } \Sigma' \in Co(\mathcal{W}).$$

Let D be the set of demand matrices $f = \{f_o^d\}$ such that there exists:

$$\Sigma \in Co(\mathcal{W}) \text{ with } r_{ab} \Sigma_{ab} \geq \sum_{(o,d) \in \mathcal{F}} S_{ab}^{[o,d]}, \forall (a, b) \in \mathcal{L} \quad (8)$$

Let D^0 be the interior of D if the set of demand matrices $f = \{f_o^d\}$ such that there exists:

$$\Sigma \in Co(\mathcal{W}) \text{ with } r_{ab} \Sigma_{ab} > \sum_{(o,d) \in \mathcal{F}} S_{ab}^{[o,d]}, \forall (a, b) \in \mathcal{L} \quad (9)$$

Proposition 2. *The point 0 belongs to D . D is a convex, compact polytope. If $f \in D$ and $0 \leq f' \leq f$, then f' also belongs to D .*

4. Joint control algorithm development and performance analysis

Section 4.1 introduces a virtual network to simplify the analysis of congestion for vehicles arriving at distant junctions. Section 4.2 designs the ST-CRJC algorithm without considering the large communication delay between entry nodes and gateways and conducts a performance analysis. Section 4.3 develops the DT-CRJC algorithm, which integrates the large communication delay between entry nodes and gateways, and conducts a performance analysis.

4.1. Virtual network

Due to limited road capacity in the real network, vehicles at junctions queue based on adjacent downstream junctions. The queue lengths at each junction in the real network can only indicate local congestion of vehicles from the current junction to adjacent downstream junctions. Large-scale collection of queue lengths at junctions for network-level control not only causes high communication overhead and computational complexity but also makes fully distributed control difficult to implement. Similar to Zaidi et al. (2016), a virtual network identical in structure to the real network is constructed to ensure that queue lengths at junctions directly reflect the congestion of vehicles arriving at distant junctions.

They required each junction to maintain separate queues for all destinations. This approach not only wastes computational and storage resources when all nodes in the network serve as destinations, but also causes unnecessary detours of vehicles at low demand, as verified in their numerical simulations. Unlike their approach, we only require each node $b \in \mathcal{N}$ in the virtual network to maintain independent queues for all gateways in the network and the interior nodes of the region containing node b , except for node b itself. More specifically, each node $b \in \mathcal{N}$ maintains independent queues \tilde{Q}_b^i for all nodes i in set H_b .

When a vehicle enters the real network, a corresponding virtual vehicle enters the virtual network with probability 1, and an additional virtual vehicle enters the virtual network with probability $\sigma > 0$. Hence, the arrival rate \tilde{f}_o^d of flow $[o, d]$ in the virtual network is $1 + \sigma$ times the arrival rate f_o^d in the real network. These additional virtual vehicles are introduced to ensure the equivalence of queue stability between the virtual network and the real network. The detailed proof is provided in Theorem 2. Before reaching the regions containing their destinations, all virtual vehicles at junctions join the queues associated with the gateways allocated to them by traffic control at entry nodes. After entering the regions containing their destinations, virtual vehicles at junctions join the queues associated with their destinations. It is worth noting that the vehicles and queues in the virtual network serve only as counters. Unlike the real network, the virtual network permits infinitely long queues.

4.2. ST-CRJC development and performance analysis

The ST-CRJC algorithm includes traffic control, signal control, and routing. The overall flowchart of ST-CRJC is shown in Fig. 5, which illustrates the sequence of control operations and corresponding updates of control variables in the virtual and real networks within a single time slot. ST-CRJC primarily updates the control parameters based on the queue information at nodes in the virtual

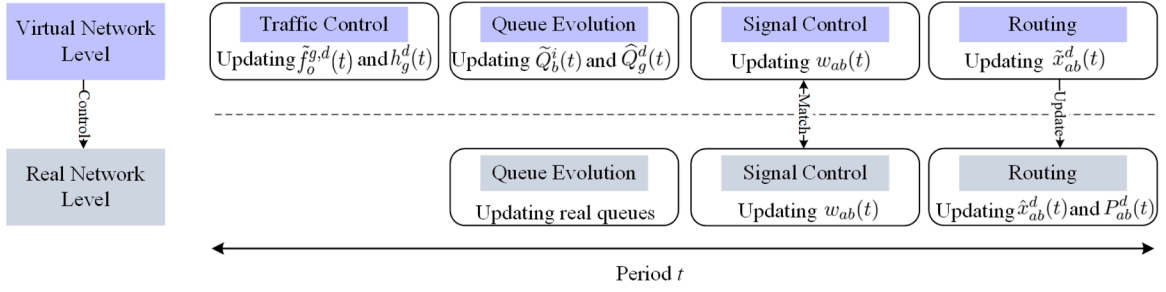
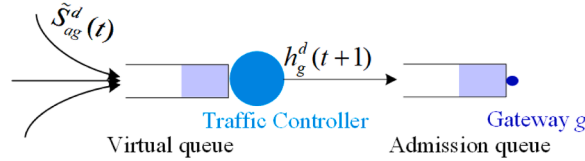


Fig. 5. The process diagram of ST-CRJC.

Fig. 6. Traffic control at gateway g for destination d in the interiors of $C(g)$.

network to direct vehicle movement in the real network. More specifically, the virtual network performs traffic control, signal control, and routing, while the real network follows the virtual network's control by implementing the same signal phases, updating routing probabilities, and updating queues accordingly. In other words, in the real network, vehicles at each junction only enter adjacent downstream junctions after phase actuation and join the left-turn, straight-through, and right-turn queues based on routing probabilities. The three components in ST-CRJC are designed as follows:

Traffic Control: Traffic control is implemented through the cooperation of entry nodes and gateways to balance the traffic at inter-regional connection nodes and avoid region oversaturation:

Traffic Control at Entry Nodes: At the beginning of time slot $(t + 1)$, the traffic controller at entry node o splits the exogenous virtual demand $\tilde{f}_o^d(t + 1)$ into $\{\tilde{f}_o^{g,d}(t + 1)\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$ based on the following:

$$\tilde{f}_o^{g,d}(t + 1) = \begin{cases} \tilde{f}_o^d(t + 1), & g = g^* \\ 0, & g \neq g^* \end{cases} \quad (10)$$

where g^* equals $\arg \min_{g \in \{0\} \cup \mathcal{G}_{C(d)}} \varpi_o^{g,d}(t)$ with $\begin{cases} \varpi_o^{g,d}(t) = \tilde{Q}_o^g(t) + \tilde{Q}_g^d(t), & \text{if } g \in \mathcal{G}_{C(d)} \\ \varpi_o^{0,d}(t) = \tilde{Q}_o^d(t), & \text{if } o \in I_{C(d)}. \\ \varpi_o^{0,d}(t) = \infty, & \text{if } o \notin I_{C(d)} \end{cases}$. Eq. (10) indicates that the traffic controller at entry

nodes o allocates all exogenous demand $\tilde{f}_o^d(t + 1)$ to the node $g \in \{0\} \cup \mathcal{G}_{C(d)}$ with the lowest congestion level to enter the region containing destination d . The case $g = 0$ represents that flow $[o, d]$ can be routed directly to destination d without passing through the gateways of other regions because its destination and origin are in the same region. If origin o is an interior node of the region containing destination d (i.e., $o \in I_{C(d)}$), the congestion $\varpi_o^{0,d}(t)$ of virtual demand $\tilde{f}_o^d(t + 1)$ arriving at destination d equals the queue length $\tilde{Q}_o^d(t)$ for destination d at entry node o ; otherwise, $\varpi_o^{0,d}(t)$ equals infinity. If destination d and node o are in different regions or if destination d is a gateway of the region containing node o , the congestion $\varpi_o^{g,d}(t)$ of virtual demand $\tilde{f}_o^d(t + 1)$ passing through each gateway $g \in \mathcal{G}_{C(d)}$ to destination d equals the queue length $\tilde{Q}_o^g(t)$ at entry node o for gateway $g \in \mathcal{G}_{C(d)}$ plus the queue length $\tilde{Q}_g^d(t)$ at gateway g for destination d . If there is no path from gateway $g \in \cup_c \mathcal{G}_C$ to destination d within its region (i.e., $d \in C(g)$)¹, then $\tilde{Q}_g^d(t)$ is set to infinity. The reason for calculating the congestion of virtual demand $\tilde{f}_o^d(t + 1)$ passing through gateway $g \in \mathcal{G}_{C(d)}$ to destination d based only on the queue lengths $\tilde{Q}_o^g(t)$ and $\tilde{Q}_g^d(t)$ is explained in the following Remark 1.

Traffic Control at Gateways: In the virtual network, each gateway $g \in \cup_c \mathcal{G}_C$ maintains a virtual queue and an admission queue for the destination d in the interiors of region $C(g)$, as shown in Fig. 6. The admission queues in these gateways are used for signal control and routing. At the beginning of time slot $(t + 1)$, the traffic controller allows the following vehicles from the virtual queue to enter

¹ Since we only need to determine network connectivity, we can ignore the travel costs (i.e., weights) between nodes. The depth-first search (DFS) (Yigit et al., 2021) and breadth-first search (BFS) (Jarjan, 1972) algorithms, which have low computational complexity, can be adopted. The complexity of both algorithms is $O(|\mathcal{N}_C| + |\mathcal{L}_C|)$, where \mathcal{N}_C and \mathcal{L}_C represent the sets of nodes and links in region C , respectively, and $|\cdot|$ represents the size of the set. DFS and BFS typically handle large-scale unweighted networks with thousands to hundreds of thousands of nodes within milliseconds to seconds. We only need to check whether there are paths between gateways and the remaining nodes within the same region after regional divisions. DFS and BFS ensure the computational efficiency of this check. Moreover, this connectivity check is performed once after regional divisions, so it does not affect the efficiency of ST-CRJC and DT-CRJC.

the admission queue:

$$h_g^d(t+1) = (1+\epsilon) \sum_{o:|o,d|\in F} \tilde{f}_o^{g,d}(t+1), \forall g \in \cup_C \mathcal{G}_C \text{ and } d \in \mathcal{I}_{C(g)} \quad (11)$$

where ϵ is an arbitrarily small positive number used to adjust the traffic transmission rate between regions. Although there may be some ambiguity, we define the admission queue at gateway g for destination d in the interiors of region $C(g)$ as $\tilde{Q}_g^d(t)$ and the virtual queue as $\hat{Q}_g^d(t)$ for simplicity. The reason why traffic control at gateways can adjust the traffic transmission rate between regions, i.e., how parameter ϵ adjusts the transmission rate between regions, is explained in [Remark 2](#).

Signal Control: At the end of time slot t , the network control matrix $W(t) = \{w_{ab}(t)\}$ is updated for the use in time slot $(t+1)$ by solving the following optimization problem:

$$\max \sum_{(a,b) \in \mathcal{L}} r_{ab} w_{ab}(t) k_{ab}(t) \quad (12)$$

$$\text{s.t. } W(t) = \{w_{ab}(t)\} \in \mathcal{W} \quad (13)$$

where $k_{ab}(t)$ denotes the pressure of traffic movement (a, b) . $k_{ab}(t)$ equals the maximum queue length difference between the upstream and downstream queues and is defined as $k_{ab}(t) = \tilde{Q}_a^{e^*ab}(t) - \tilde{Q}_b^{e^*ab}(t)$, where $e^*_{ab}(t) = \begin{cases} \arg \max_{g \in \cup_C \mathcal{G}_C} (\tilde{Q}_a^g(t) - \tilde{Q}_b^g(t)), & \text{if } b \in \cup_C \mathcal{G}_C \\ \arg \max_{e \in \mathcal{H}_a} (\tilde{Q}_a^e(t) - \tilde{Q}_b^e(t)), & \text{otherwise} \end{cases}$. As seen

from the above optimization problem, signal control decisions at all junctions are independent and do not interfere with each other. Therefore, the signal control matrix $W^j(t) = \{w_{ab}^j(t)\}$ for junction $j \in \mathcal{J}$ can be updated through the following optimization problem:

$$\max \sum_{(a,b) \in \mathcal{L}, a \in A_j} r_{ab} w_{ab}^j(t) k_{ab}(t) \quad (14)$$

$$\text{s.t. } W^j(t) = \{w_{ab}^j(t)\} \in \mathcal{W}^j \quad (15)$$

Solving the above optimization problem is straightforward because both r_{ab} and $k_{ab}(t)$ are known. Additionally, each junction has a limited number of candidate signal control matrices. Therefore, the above optimization problem can be optimally solved by enumerating all candidate signal control matrices. After obtaining the optimal control matrix, the corresponding phases are actuated. Then, the virtual vehicles in queue $e^*_{ab}(t)$ and the real vehicles constrained to the corresponding phases are served. If $e^*_{ab}(t) = b$ and b is a gateway, virtual vehicles in queue $\tilde{Q}_a^b(t)$ at node a are destined for b , or travel to their destinations through gateway b . In the latter case, the served virtual vehicles join the virtual queue $\hat{Q}_b^d(t)$ at gateway b for destination d .

Routing: Let $\tilde{x}_{ab}^d(t)$ denote the number of virtual vehicles traveling from node a to node b with destination d during time slot t by the above signal control algorithm. \bar{x}_{ab}^d is the expected value of $\tilde{x}_{ab}^d(t)$ when the virtual queue process is in a stationary regime. $\hat{x}_{ab}^d(t)$ is the estimated value of $\tilde{x}_{ab}^d(t)$, calculated at the end of time slot t . When the queues in the virtual network are positive recurrent, "reliable" estimators $\hat{x}_{ab}^d(t)$ can be maintained by simple averaging and remain close to the expected value \bar{x}_{ab}^d . According to [Zaidi et al. \(2016\)](#) and [Athanasopoulou et al. \(2013\)](#), the exponential averaging method is used to update the estimates $\hat{x}_{ab}^d(t)$, i.e., $\hat{x}_{ab}^d(t) = (1-\alpha)\hat{x}_{ab}^d(t-1) + \alpha\tilde{x}_{ab}^d(t)$, where α is a smoothing factor. The routing probabilities are updated at the end of time slot t as follows:

$$P_{ab}^d(t) = \frac{\hat{x}_{ab}^d(t)}{\sum_{c:(a,c) \in \mathcal{L}} \hat{x}_{ac}^d(t)} \quad (16)$$

In the real network, a vehicle arriving node a with destination d will be routed to node b with probability $P_{ab}^d(t)$ in time slot $(t+1)$. In other words, a vehicle entering node a destined for node d joins the real queue Q_{ab} with probability $P_{ab}^d(t)$ in time slot $(t+1)$. We use the routing diagram in the upper right corner of [Fig. 2](#) as an example to illustrate this. If $P_{26,12}^d(t) = 0.2$, $P_{26,21}^d(t) = 0.4$, and $P_{26,28}^d(t) = 0.4$, vehicles with destination d join queue $Q_{26,12}$ with probability 0.2, queue $Q_{26,21}$ with probability 0.4, and queue $Q_{26,28}$ with probability 0.4 when approaching node 26 in time slot $(t+1)$.

Remark 1. At the beginning of time slot t , traffic control at entry nodes only uses the queue lengths $\tilde{Q}_o^g(t-1)$ and $\tilde{Q}_g^d(t-1)$ to determine the congestion associated with demand $\tilde{f}_o^d(t)$ passing through gateway $g \in \mathcal{G}_{C(d)}$ to reach destination d . This is because, in signal control, the pressure weight associated with each turn movement (a, b) in time slot t depends on the maximum difference $k_{ab}(t-1)$ between upstream and downstream queue lengths. Let us take the simple network shown in [Fig. 7](#) as an example. The network in [Fig. 7](#) consists of two regions, two gateways (i.e., 3 and 4), and two flows (i.e., [1, 5] and [1, 6]). The queues at each node are marked above it. At the beginning of time slot t , traffic control at entry node 1 allocates demand $\tilde{f}_1^4(t)$ into $\tilde{f}_1^{3,5}(t)$ and $\tilde{f}_1^{4,5}(t)$ based on the minimum of $\tilde{Q}_1^3(t-1) + \tilde{Q}_3^5(t-1)$ and $\tilde{Q}_1^4(t-1) + \tilde{Q}_4^5(t-1)$, and splits demand $\tilde{f}_1^6(t)$ into $\tilde{f}_1^{3,6}(t)$ and $\tilde{f}_1^{4,6}(t)$ based on the minimum of $\tilde{Q}_1^3(t-1) + \tilde{Q}_3^6(t-1)$ and $\tilde{Q}_1^4(t-1) + \tilde{Q}_4^6(t-1)$. At the end of time slot $t-1$, signal control at node 1 calculates the pressure $k_{12}(t-1)$ of traffic movement $(1, 2)$ by calculating $e_{12}^*(t-1) = \arg \max_{g \in \{3,4\}} \{\tilde{Q}_1^g(t-1) - \tilde{Q}_2^g(t-1)\}$, and allows virtual vehicles in queue $\tilde{Q}_1^{e_{12}^*(t-1)}(t)$ to join queue $\tilde{Q}_2^{e_{12}^*(t-1)}(t)$ after phase actuation. Usually, when the pressure $k_{12}(t-1)$ is positive (i.e., $\tilde{Q}_1^{e_{12}^*(t-1)}(t-1)$ is greater than $\tilde{Q}_2^{e_{12}^*(t-1)}(t-1)$), the corresponding phase can be actuated. Therefore, the congestion of virtual vehicles passing through node 2 to reach gateway $e_{12}^*(t)$ is effectively shown in the queue length $\tilde{Q}_1^{e_{12}^*(t)}(t)$ at node 1. Each gateway also has a traffic controller.

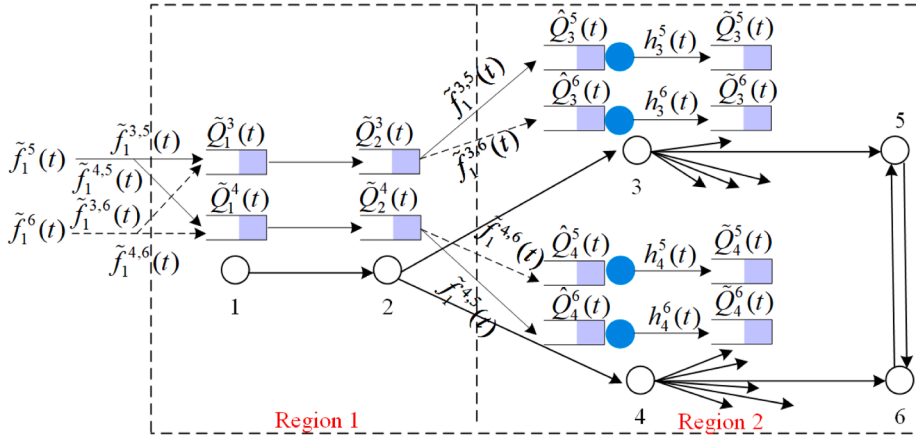


Fig. 7. Splitting flows [1, 5] and [1, 6] in time slot t .

Similarly, the congestion of vehicles at gateway $e_{12}^*(t)$ traveling toward destination 5 or 6 is effectively reflected in the queue length $\tilde{Q}_{e_{12}^*(t)}^5(t)$ or $\tilde{Q}_{e_{12}^*(t)}^6(t)$ at the gateway, respectively. Therefore, the queue lengths $\tilde{Q}_o^g(t-1)$ and $\tilde{Q}_g^d(t-1)$ are sufficient for traffic control at entry nodes to estimate the congestion associated with demand $\tilde{f}_o^d(t)$ passing through gateway $g \in \mathcal{G}_{C(d)}$ toward destination d .

Remark 2. In the virtual network, signal control is based on the MP algorithm, which typically selects signal phases that maximize the queue differential between upstream and downstream queues. Gateways connect their regions to other regions, which implies that all traffic destined for nodes within their regions must pass through them. Typically, there is heavy upstream traffic at gateways, which often causes signal control to activate phases that serve the upstream movements of gateways, thereby releasing a large volume of traffic from neighboring regions into the corresponding region. The traffic controller in Fig. 6 uses the parameter ϵ to reconstruct an admission queue. This shorter admission queue is then used for signal phase optimization at gateways. This reconstruction controls the activation of phases that limit traffic from other regions, thereby regulating the rate of traffic transmission between regions.

As shown in the above ST-CRJC algorithm, in the real network, vehicles at each junction only enter adjacent downstream junctions after phase actuation and join the left-turn, straight-through, and right-turn queues based on routing probabilities. Therefore, the queue updates at junctions in the real network are omitted. The queues at each node in the virtual network are updated as follows:

$$\begin{aligned} \tilde{Q}_b^i(t+1) = & \tilde{Q}_b^i(t) + \left[(1+\epsilon) \sum_{o: [o,d] \in \mathcal{F}} \tilde{f}_o^{b,i}(t+1) \wedge \tilde{Q}_b^i(t) \right] \\ & - \left[\sum_{c: (b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(i)}} R_{bc}(t+1) \mathbf{1}_{\{e_{bc}^*(t)=i\}} \mathbf{w}_{bc}(t) \wedge \tilde{Q}_b^i(t) \right], \forall b \in \cup_c \mathcal{G}_C \text{ and } i \in \mathcal{I}_{C(b)} \end{aligned} \quad (17)$$

$$\begin{aligned} \hat{Q}_b^i(t+1) = & \hat{Q}_b^i(t) + \sum_{d \in \mathcal{C}(b)} \tilde{f}_b^{i,d}(t+1) + \left[\sum_{a: (a,b) \in \mathcal{L}} R_{ab}(t+1) \mathbf{1}_{\{e_{ab}^*(t)=i\}} \mathbf{w}_{ab}(t) \wedge \tilde{Q}_a^i(t) \right] \\ & - \left[(1+\epsilon) \sum_{o: [o,d] \in \mathcal{F}} \tilde{f}_o^{b,i}(t+1) \wedge \hat{Q}_b^i(t) \right], \forall b \in \cup_c \mathcal{G}_C \text{ and } i \in \mathcal{I}_{C(b)} \end{aligned} \quad (18)$$

$$\begin{aligned} \tilde{Q}_b^i(t+1) = & \tilde{Q}_b^i(t) + \tilde{f}_b^{0,i}(t+1) + \left[\sum_{a: (a,b) \in \mathcal{L}} R_{ab}(t+1) \mathbf{1}_{\{e_{ab}^*(t)=i\}} \mathbf{w}_{ab}(t) \wedge \tilde{Q}_a^i(t) \right] \\ & - \left[\sum_{c: (b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(i)}} R_{bc}(t+1) \mathbf{1}_{\{e_{bc}^*(t)=i\}} \mathbf{w}_{bc}(t) \wedge \tilde{Q}_b^i(t) \right], \forall b \in \cup_c \mathcal{I}_C \text{ and } i \in \mathcal{I}_{C(b)} \end{aligned} \quad (19)$$

$$\begin{aligned} \tilde{Q}_b^i(t+1) = & \tilde{Q}_b^i(t) + \sum_{d \in \mathcal{C}(b)} \tilde{f}_b^{i,d}(t+1) + \left[\sum_{a: (a,b) \in \mathcal{L}} R_{ab}(t+1) \mathbf{1}_{\{e_{ab}^*(t)=i\}} \mathbf{w}_{ab}(t) \wedge \tilde{Q}_a^i(t) \right] \\ & - \left[\sum_{c: (b,c) \in \mathcal{L}} R_{bc}(t+1) \mathbf{1}_{\{e_{bc}^*(t)=i\}} \mathbf{w}_{bc}(t) \wedge \tilde{Q}_b^i(t) \right], \forall b \in \mathcal{N} \text{ and } i \in \cup_c \mathcal{G}_C \end{aligned} \quad (20)$$

where $z \wedge y$ is a function defined as $\min[z, y]$. Eq. (17) indicates that the admission queue $\tilde{Q}_b^i(t)$ at gateway $b \in \cup_c \mathcal{G}_C$ for its interior node $i \in \mathcal{I}_{C(b)}$ increases by up to $(1+\epsilon) \sum_{o: [o,d] \in \mathcal{F}} \tilde{f}_o^{b,i}(t+1)$ vehicles from the output of the traffic controller at gateway $b \in \cup_c \mathcal{G}_C$ in

time slot $(t + 1)$; it decreases by up to $R_{bc}(t + 1)$ vehicles in time slot $(t + 1)$ if traffic movement (b, c) is actuated and the queue for i generates the maximum queue length difference (i.e., $e_{bc}^*(t) = i$). Eq. (18) indicates that the virtual queue $\tilde{Q}_b^i(t)$ at gateway $b \in \cup_C \mathcal{G}_C$ for its interior node $i \in \mathcal{I}_{C(b)}$ increases by up to $\sum_{d \in \mathcal{C}(b)} \tilde{f}_b^{i,d}(t + 1)$ vehicles from exogenous demand and $R_{ab}(t + 1)$ vehicles from upstream node a in time slot $(t + 1)$; it decreases by up to $(1 + \epsilon) \sum_{o: [o,d] \in \mathcal{F}} \tilde{f}_o^{b,i}(t + 1)$ vehicles after passing through the traffic controller at gateway $b \in \cup_C \mathcal{G}_C$ to join the virtual queue $\tilde{Q}_b^i(t + 1)$ in time slot $(t + 1)$. Eq. (19) indicates that the virtual queue $\tilde{Q}_b^i(t)$ at interior node $b \in \cup_C \mathcal{I}_C$ for another interior node $i \in \mathcal{I}_{C(b)}$ increases by up to $\tilde{f}_b^{0,i}(t + 1)$ vehicles from exogenous demand and $R_{ab}(t + 1)$ vehicles from upstream node a in time slot $(t + 1)$; it decreases by up to $R_{bc}(t + 1)$ vehicles in time slot $(t + 1)$ as they join the downstream queue $\tilde{Q}_c^i(t)$. Eq. (20) indicates that the virtual queue $\tilde{Q}_b^i(t)$ at any node $b \in \mathcal{N}$ for gateway $i \in \cup_C \mathcal{G}_C$ decreases as described in Eq. (19) in time slot $(t + 1)$, but increases by $\sum_{d \in \mathcal{C}(i)} \tilde{f}_b^{i,d}(t + 1)$ from exogenous demand in time slot $(t + 1)$.

Before analyzing the performance of ST-CRJC, we introduce Lemma 1 to reformulate flow conservation constraints and feasible demand in the virtual network. The flow conservation constraints and feasible demand in Definition 1 are based on the origin o and destination d of $\mathbf{f} = \{f_o^d\}$. According to the queue updates in Eqs. (17)–(20), the flow conservation constraints and feasible demand in the virtual network should be based on the split demand $\left\{ \tilde{f}_o^{g,d}(t) \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$. Let $\tilde{\mathbf{f}}_s$ be $\left\{ \tilde{f}_o^{g,d} \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$ with entries

$$\tilde{f}_o^{g,d} = \liminf \frac{1}{T} \sum_{t=1}^T \tilde{f}_o^{g,d}(t). \text{ Let } \tilde{\mathcal{S}}(t) \text{ be } \left\{ \tilde{\mathcal{S}}_{ab}^g(t) \right\}, \text{ where } \tilde{\mathcal{S}}_{ab}^g(t) \text{ represents the service rate of the virtual queue for } g \text{ over link } (a, b) \text{ in } t.$$

$$\text{Let } \tilde{\mathcal{S}} \text{ be } \left\{ \tilde{\mathcal{S}}_{ab}^g \right\} \text{ with entries } \tilde{\mathcal{S}}_{ab}^g = \liminf \frac{1}{T} \sum_{t=1}^T \mathcal{S}_{ab}^g(t).$$

Lemma 1. *The virtual demand $\tilde{\mathbf{f}} = \{\tilde{f}_o^d\}$ is supportable by joint control such that the infinitely admissible signal control sequence $\{W(1), W(2), \dots, W(t), \dots\}$ accommodates $\tilde{\mathbf{f}}_s = \left\{ \tilde{f}_o^{g,d} \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$:*

$$\tilde{f}_o^d = \tilde{f}_o^{0,d} \mathbf{1}_{o \in \mathcal{I}_{C(d)}} + \sum_{g \in \mathcal{G}_{C(d)}} \tilde{\mathbf{f}}_o^{g,d}, \forall [o, d] \in \mathcal{F} \quad (21)$$

$$\sum_{o: [o,d] \in \mathcal{F}} \tilde{f}_o^{b,d} = \sum_{c: (b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(d)}} \tilde{\mathcal{S}}_{bc}^d, \forall b \in \cup_C \mathcal{G}_C \text{ and } d \in \mathcal{I}_{C(b)} \quad (22)$$

$$\sum_{d \in \mathcal{C}(g)} \tilde{f}_b^{g,d} + \sum_{a: (a,b) \in \mathcal{L}} \tilde{\mathcal{S}}_{ab}^g = \sum_{c: (b,c) \in \mathcal{L}} \tilde{\mathcal{S}}_{bc}^g, \forall b \in \mathcal{N} \text{ and } g \in \cup_C \mathcal{G}_C \quad (23)$$

$$\tilde{f}_b^{0,d} + \sum_{a: (a,b) \in \mathcal{L}} \tilde{\mathcal{S}}_{ab}^d = \sum_{c: (b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(d)}} \tilde{\mathcal{S}}_{bc}^d, \forall b \in \cup_C \mathcal{I}_C \text{ and } d \in \mathcal{I}_{C(b)} \quad (24)$$

$$r_{ab} \Sigma_{ab} > \sum_{g \in \mathcal{H}_a} \tilde{\mathcal{S}}_{ab}^g, \forall (a, b) \in \mathcal{L} \quad (25)$$

where Σ_{ab} equals $\liminf \frac{1}{T} \sum_{t=1}^T w_{ab}(t)$.

Proof. The proof is provided in Appendix A. \square

An intuitive interpretation of the flow conservation constraints in Eqs. (21)–(24) is shown in Fig. 8. Eq. (21) indicates that the exogenous demand \tilde{f}_o^d is divided into $\left\{ \tilde{f}_o^{g,d} \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$ by the traffic controller at entry node o , depending on whether destination d and origin o are in the same region. Eq. (22) implies that the amount of traffic routed to gateway b for destination d equals the amount of traffic departing from gateway b for d . Eq. (23) indicates that the amount of traffic leaving the queue for gateway g at any node b equals the incoming traffic from adjacent upstream nodes plus the traffic generated by node b itself to be routed to gateway g . Eq. (24) implies that if destination d and node b are in the same region, the traffic leaving the queue equals the traffic coming from upstream nodes plus the traffic generated by node b itself and routed directly to destination d without passing through the gateways of other regions. Similar to Eqs. (4) and (25) means that the traffic $r_{ab} \Sigma_{ab}$ served by signal control per time slot on average is greater than the average traffic movement $\sum_{g \in \mathcal{H}_a} \tilde{\mathcal{S}}_{ab}^g$ generated by all demands per time slot.

According to Levin et al. (2020), network throughput equals the sum of the input rates of all entry nodes when the queueing network is stable (i.e., all queue lengths have finite means). Intuitively, all vehicles entering the network can be served if the queueing network is stable. This means that the control algorithm achieving network stability is equivalent to maximizing network throughput. Following Neely (2010), network stability is defined as follows:

Definition 2. The queueing network $Q(t) = \{Q_{ab}(t), (a, b) \in \mathcal{L}\}$ is strongly stable in the mean if the sum of the time-averaged expected queue lengths is bounded, i.e.,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{(a,b) \in \mathcal{L}} E\{Q_{ab}(t)\} < \infty \quad (26)$$

where $E\{\cdot\}$ is the expectation function.

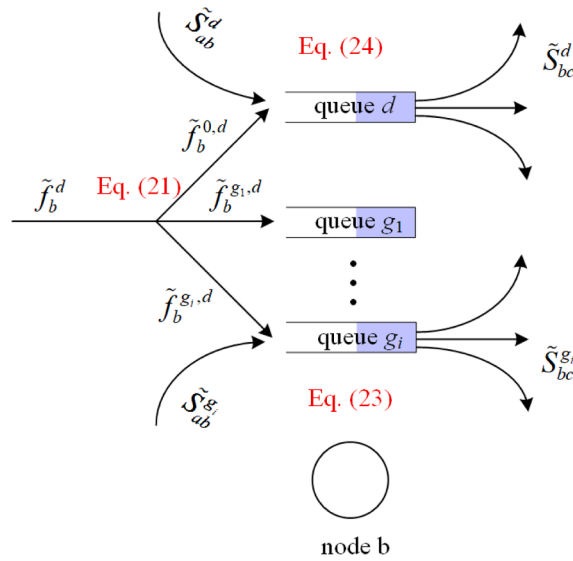


Fig. 8. Flow conservation constraints at node b .

There are virtual queues and admission queues in the virtual network. To simplify subsequent proof of network stability, we define all queues in the virtual network into matrix $\tilde{Q}(t) = \{\tilde{Q}_b^i(t), b \in \mathcal{N}, i \in H_b\}$ and matrix $\hat{Q}(t) = \{\hat{Q}_g^d(t), g \in \cup_C \mathcal{G}_C, d \in I_{C(g)}\}$. According to Definition 2, we conclude that the virtual network is strongly stable in the mean if the following two equations hold:

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{b \in \mathcal{N}} \sum_{i \in H_b} E\{\tilde{Q}_b^i(t)\} < \infty \tag{27}$$

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{g \in \cup_C \mathcal{G}_C} \sum_{d \in I_{C(g)}} E\{\hat{Q}_g^d(t)\} < \infty \tag{28}$$

Theorem 1. If the exogenous virtual demand \tilde{f} is feasible (i.e., $\tilde{f} = \{\tilde{f}_o^d\} \in D^0$), the virtual network under ST-CRJC is strongly stable in the mean.

Proof. Let $|\tilde{Q}(t)|^2$ be the sum of squares of all the queue lengths in matrix $\tilde{Q}(t)$, i.e., $|\tilde{Q}(t)|^2 = \sum_{b \in \mathcal{N}} \sum_{i \in H_b} [\tilde{Q}_b^i(t)]^2$. Under ST-CRJC, there exist $M_1 < \infty$ and $\vartheta > 0$ such that the one-step conditional Lyapunov drift function satisfies the following equation:

$$E\left\{|\tilde{Q}(t+1)|^2 - |\tilde{Q}(t)|^2 \mid \tilde{Q}(t)\right\} \leq M_1 - \vartheta |\tilde{Q}(t)| \tag{29}$$

The proof of Eq. (29) is shown in Appendix B. Taking an expectation with $\tilde{Q}(t)$ and summing over $t = 1, \dots, T$ yield:

$$E|\tilde{Q}(T+1)|^2 - E|\tilde{Q}(1)|^2 \leq M_1 T - \vartheta \sum_{t=1}^T E|\tilde{Q}(t)| \tag{30}$$

Since $E|\tilde{Q}(T+1)|^2$ is larger than 0, we have:

$$\vartheta \frac{1}{T} \sum_{t=1}^T E|\tilde{Q}(t)| \leq M_1 + \frac{1}{T} E|\tilde{Q}(1)|^2 - \frac{1}{T} E|\tilde{Q}(T+1)|^2 \leq M_1 + \frac{1}{T} E|\tilde{Q}(1)|^2 \tag{31}$$

Taking a limit as $T \rightarrow \infty$ yields:

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T E|\tilde{Q}(t)| \leq \frac{M_1}{\vartheta} \tag{32}$$

The proof of Eq. (27) ends here. Let $\mathcal{K}_g^d(t)$ be the total number of vehicles destined for d that pass through gateway g and have not yet reached the admission queue for node d at gateway g . We assume that there exists $t \geq t_0$ such that:

$$\mathcal{K}_g^d(t) > |\tilde{Q}(t)| + M \tag{33}$$

where M is the maximum change in any queue length in one time slot and is defined in Appendix B. Eq. (33) implies that:

$$\hat{Q}_g^d(t) > M \tag{34}$$

Through the traffic control at gateways in Eq. (11), we have:

$$\mathcal{K}_g^d(t+1) = \sum_{o:[o,d] \in \mathcal{F}} \tilde{f}_o^{g,d}(t+1) + \mathcal{K}_g^d(t) - h_g^d(t+1) \leq \mathcal{K}_g^d(t) \quad (35)$$

This implies that for any $t > t_0$:

$$\mathcal{K}_g^d(t) \leq \max \left\{ \mathcal{K}_g^d(t_0), \left| \tilde{Q}(t) \right| + 2M \right\} \quad (36)$$

Hence, all queue lengths in the matrix $\hat{Q}(t) = \left\{ \hat{Q}_g^d(t), g \in \cup_c \mathcal{G}_C, d \in \mathcal{I}_{C(b)} \right\}$ are bounded. This directly implies that Eq. (28) holds. Therefore, the sum of the time-averaged expected lengths of all queues in the virtual network is bounded. Based on Definition 2, we know that the virtual network is strongly stable in the mean. This completes the proof of Theorem 1. \square

Theorem 1 implies that the virtual network is stable in the mean. According to Varaiya (2013), network stability in the mean implies that all queues in the network are positive recurrent. Therefore, reliable estimates $\hat{x}_{ab}^d(t)$ can be maintained by simple averaging (e.g., exponential averaging). The routing probabilities in Eq. (16) remain close to their ideal values:

$$\bar{P}_{ab}^d = \frac{\bar{x}_{ab}^d}{\sum_{c:(a,c) \in \mathcal{L}} \bar{x}_{ac}^d} \quad (37)$$

Theorem 2. *If the exogenous virtual demand \tilde{f} is feasible (i.e., $\tilde{f} = \{\tilde{f}_o^d\} \in D^0$), ST-CRJC ensures that the real network is strongly stable in the mean.*

Proof. The detailed proof is provided in Appendix C. In Appendix C, we show that the increase rate of each real queue is strictly smaller than the average service rate allocated to it under the ideal probabilities in Eq. (37). This implies that the lengths of all real queues have bounded means (see Eq. (5) for a similar analysis). Based on Definition 2, we conclude that the real network is strongly stable in the mean. \square

Theorem 3. *For a network with random demand arrival rates, full connectivity, and a finite state space, ST-CRJC ensures that all vehicles eventually reach their destinations and do not cycle infinitely.*

Proof. Consider the movement of vehicles in the network as a finite state space Markov chain. Each state represents vehicles being at a certain node in the network at a specific time. Vehicles with destinations d transitioning from node a to node b with probability $P_{a,b}^d(t)$ indicate that this process is stochastic. According to Theorem 2, the queuing process in the real network is stable in the mean. Stability in the mean implies that the Markov chain describing the real queuing process is positive recurrent and has a unique steady-state probability distribution (Varaiya, 2013). In a finite state space Markov chain, if every state is positive recurrent, the mean return time to any state is finite. Positive recurrence ensures that no states in the network are visited with zero probability in the long run. This means that vehicles do not become trapped in any part of the network indefinitely. The full connectivity of the network and random routing choices at each node together ensure the irreducibility of the Markov chain, meaning all states are reachable from each other. Random routing breaks any fixed periodicity, thus fulfilling aperiodicity. Positive recurrence, irreducibility, and aperiodicity together ensure that the Markov chain is ergodic (Gallager, 1996; Ghahramani, 2005). Ergodicity implies that, in the long run, the system visits all states with the probability of stationary distribution (Meshalkin, 1958; Sandrić, 2014), ensuring that vehicles can eventually reach any node, including their final destinations. Therefore, vehicles do not cycle infinitely in the network. The behavior of the vehicles covers all possible states without being confined to any subset or falling into infinite cycling paths. \square

4.3. DT-CRJC development and performance analysis

We recall that τ_{ab} is the information transmission delay from node $a \in \mathcal{N}$ to node $b \in \mathcal{N}$, with a bound τ_{\max} , i.e., $\tau_{\max} = \max_{a,b \in \mathcal{N}} \tau_{ab}$. According to the problem description of DT-CRJC in Section 3.1, traffic control updates control parameters every γ time slots (i.e., $\gamma > 3\tau_{\max}$), while signal control and routing update control parameters every time slot. Therefore, for DT-CRJC, we keep the signal control and routing from ST-CRJC and only develop the low-timescale traffic control.

Low-Timescale Traffic Control: Similar to the real-time traffic control described in Section 4.2, low-timescale traffic control is also implemented under the cooperation of entry nodes and gateways, as shown below:

Low-Timescale Traffic Control at Entry Nodes: During time slots $\{(h-1)\gamma, \dots, h\gamma-1\}$, the origin o of flow $[o,d]$ receives $\left\{ \tilde{Q}_g^d((h-1)\gamma) \right\}_{g \in \mathcal{G}_{C(d)}}$ from the gateways of the region containing destination d and calculates $\left\{ \hat{f}_o^{g,d}((h-1)\gamma) \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$ according to the following:

$$\hat{f}_o^{g,d}((h-1)\gamma) = \begin{cases} \tilde{f}_o^d((h-1)\gamma), & g = g^* \\ 0, & g \neq g^* \end{cases} \quad (38)$$

where $g^* = \arg \min_{g \in \{0\} \cup \mathcal{G}_{C(d)}} \tilde{\omega}_o^{g,d}((h-1)\gamma)$. $\tilde{\omega}_o^{g,d}(t)$ is defined in Eq. (10). During time slots $\{h\gamma, \dots, (h+1)\gamma-1\}$, origin o splits the virtual demand $\tilde{f}_o^d(t)$ based on $\left\{ \hat{f}_o^{g,d}((h-1)\gamma) \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$.

Low-Timescale Traffic Control at Gateways: During time slots $\{(h-1)\gamma, \dots, h\gamma-1\}$, gateways receive the information $\left\{ \hat{f}_o^{g,d}((h-1)\gamma) \right\}_{g \in (0) \cup \mathcal{G}_{C(d)}}$ from all origins and calculate $\left\{ \hat{h}_g^d((h-1)\gamma) \right\}_{g \in \cup \mathcal{G}_{C(d)}}$ according to the following:

$$\hat{h}_g^d((h-1)\gamma) = (1 + \epsilon) \sum_{o: |o,d| \in \mathcal{F}} \hat{f}_o^{g,d}((h-1)\gamma), \text{ if } g \in \cup \mathcal{G}_C \text{ and } d \in \mathcal{I}_{C(g)} \quad (39)$$

During time slots $\{h\gamma, \dots, (h+1)\gamma-1\}$, the number of vehicles entering the admission queue at gateway g for destination d from the virtual queue at gateway g for destination d is given by $h_g^d(t) = \hat{h}_g^d((h-1)\gamma)$.

Theorem 4. *If the exogenous virtual demand $\tilde{\mathbf{f}}$ is feasible (i.e., $\tilde{\mathbf{f}} = \{\tilde{f}_o^d\} \in \mathcal{D}^0$), the virtual network is strongly stable in the mean under DT-CRJC.*

Proof. Under DT-CRJC, there exist $M_2 < \infty$ and $\varphi > 0$ such that the one-step conditional Lyapunov drift function satisfies the following equation:

$$E \left\{ \left| \tilde{\mathcal{Q}}((h+1)\gamma) \right|^2 - \left| \tilde{\mathcal{Q}}(h\gamma) \right|^2 \mid \tilde{\mathcal{Q}}(h\gamma) \right\} \leq M_2 - \varphi \left| \tilde{\mathcal{Q}}(h\gamma) \right| \quad (40)$$

The proof of Eq. (40) is shown in Appendix D. Similar to the process in Eqs. (30)–(32), we have:

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{h=1}^T E \left| \tilde{\mathcal{Q}}(h\gamma) \right| \leq \frac{M_2}{\varphi} \quad (41)$$

Hence, all queue lengths in the matrix $\tilde{\mathcal{Q}}(h\gamma)$ have bounded means. Similarly, using the process described in Eqs. (33)–(35), we have for any $h > h_0$:

$$\mathcal{K}_g^d(h\gamma) \leq \max \left\{ \mathcal{K}_g^d(h_0\gamma), \left| \tilde{\mathcal{Q}}(h\gamma) \right| + 2M \right\} \quad (42)$$

This means that all queue lengths in the matrix $\left\{ \hat{\mathcal{Q}}_g^d(h\gamma), g \in \cup \mathcal{G}_C, d \in \mathcal{I}_{C(b)} \right\}$ are bounded. Based on Definition 2, the virtual network is strongly stable in the mean. This completes the proof of Theorem 4. \square

Theorem 5. *If the exogenous virtual demand $\tilde{\mathbf{f}}$ is feasible (i.e., $\tilde{\mathbf{f}} = \{\tilde{f}_o^d\} \in \mathcal{D}^0$), DT-CRJC ensures that the real network is strongly stable in the mean.*

The proof of Theorem 5 is similar to that of Theorem 2 and is therefore omitted in this study. According to Levin et al. (2020), both ST-CRJC and DT-CRJC maximize network throughput.

Theorem 6. *For a network with random demand arrival rates, full connectivity, and a finite state space, DT-CRJC ensures that all vehicles do not cycle infinitely and eventually reach their destinations.*

The proof of Theorem 6 is similar to the proof of Theorem 3 and therefore is ignored.

Theorem 7. *DT-CRJC with low-timescale traffic control results in a higher average delay compared to ST-CRJC.*

Proof. Based on the definitions of M_1 in Eq. (B-32) and M_2 in Eq. (D-1), we know $M_1 < M_2$. Both ϑ and φ are proportional to the distance of the exogenous virtual demand $\tilde{\mathbf{f}}$ from the boundary of the feasible demand set \mathcal{D} , so ϑ and φ are equal for the same virtual demand $\tilde{\mathbf{f}}$. From Eqs. (32) and (41), we can see that the average queue length under DT-CRJC is greater than that under ST-CRJC. The Little's theorem shows that the average queue length is proportional to the average delay. Therefore, based on Little's theorem, the average delay under DT-CRJC is higher than that under ST-CRJC. This concludes the proof of Theorem 7. \square

Remark 3. The formation of gateways stems from regional divisions. However, this study does not impose any specific restrictions on regional divisions. Regardless of the regional divisions, ST-CRJC and DT-CRJC are both proven to maximize network throughput (see Theorems 2 and 5), with the latter resulting in a higher average delay (see Theorem 7). The most intuitive method for regional divisions is based on network structural or functional differences. This approach can adjust inflow rates according to the characteristics or requirements of local regions and efficiently guide flow from low-congestion nodes into local regions to utilize network capacity. Therefore, the average number of gateways generally correlates positively with the number of network structural or functional deviations. The upper bound on the number of gateways is the number of nodes in the network (if region division lines cut across all links), while the lower bound is zero (if the entire network is treated as a single region).

Remark 4. In terms of time and space, ST-CRJC and DT-CRJC do not increase complexity compared to previous algorithms that require maintaining queues for each destination. Instead, they reduce the storage space for queue data and improve computational efficiency by reducing the number of virtual queues maintained at each node. This is because both the joint control algorithms based on different timescales and previous algorithms use these virtual queues to determine the optimal signal phase and routing in each time slot. In signal phase optimization, ST-CRJC and DT-CRJC require maintaining fewer queues than previous algorithms, which reduces the computation time for weights of each candidate phase, thus improving computational efficiency. Routing only requires counting changes in the number of virtual vehicles in each queue in each time slot. Due to the smaller number of queues maintained at each node, ST-CRJC and DT-CRJC allow faster counting of these changes compared to previous algorithms. The introduction of traffic control at entry nodes and gateways does not add computational complexity. Traffic control at entry nodes uses queue lengths at entry nodes and gateways to perform all-or-nothing allocation. Gateways use traffic allocation results from entry nodes to determine

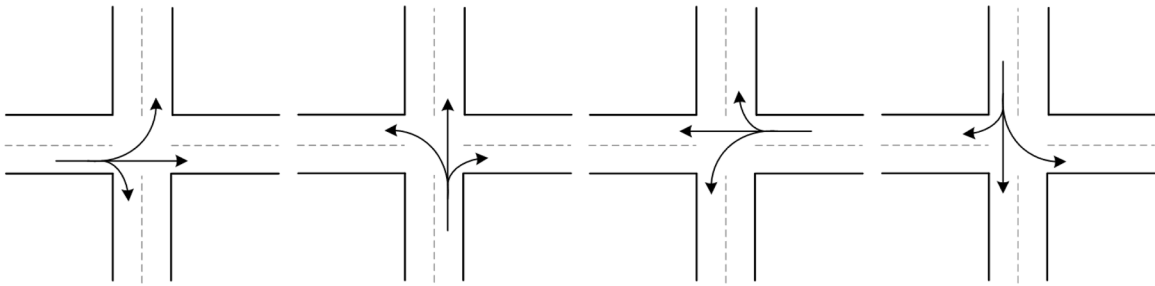


Fig. 9. The four-phase signal scheme.

the amount of traffic allowed into their regions. These are distributed computations and do not introduce additional computational complexity. Traffic control at entry nodes balances the traffic of inter-regional connection nodes, and at gateways, it prevents local oversaturation, which previous algorithms do not address. In terms of network space, ST-CRJC and DT-CRJC do not add complexity as they do not change the network topology and only deploy traffic control at inter-regional connection nodes.

5. Numerical simulations

To demonstrate the effectiveness of ST-CRJC and DT-CRJC, Sections 5.1 and 5.2 present comparative numerical simulations using a small-scale grid-like network and a large-scale network with non-uniform link density, respectively. Both networks are based on real-world data from Beijing, China. Each road in the networks contains three separate lanes, with each lane designed for either going straight, turning left, or turning right. Each junction employs the four-phase scheme from Le et al. (2017), as shown in Fig. 9. According to Le et al. (2017), the duration of each time slot is set to 25 s, and there is no phase switching time. The algorithms used for comparative simulations include ST-CRJC, DT-CRJC, fixed-time signals combined with shortest path control (FTSP), and adaptive routing back-pressure control (AR-BP) from Zaidi et al. (2016). In FTSP, the signal phase is fixed at all junctions, and all vehicles are routed using the shortest path algorithm. In AR-BP, both signal phase updates and routing are adaptively optimized through the destination-based queues in the virtual network. The smoothing factor α is set to 0.25.

5.1. Numerical simulations on the grid-like network

The small-scale grid-like network tested in this section is shown in Fig. 10. The network consists of 18 entry nodes (labeled in blue) and 18 exit nodes (labeled in yellow). An origin-destination (O-D) pair exists between each entry node and each exit node, resulting in 324 (i.e., 18×18) O-D pairs in the network. The exogenous demand from each entry node is evenly distributed to all exit nodes. Before analyzing the effect of low-timescale congestion information on network performance, the low-timescale traffic control in DT-CRJC updates control parameters every 50 s. Before analyzing the effect of region size on network performance, the network is divided into four regions using orange, yellow, blue, and green blocks. The gateways under this region division are labeled purple in Fig. 10.

5.1.1. Network capacity region

Following Li and Jabari (2019), we set up uniform demand for all entry nodes to measure the capacity region of the grid-like network in Fig. 10. Using uniform demand allows us to measure the capacity region of the network with a single parameter (i.e., demand). We perform comparative simulations with demand varying from 50 veh/h to 1300 veh/h in 50 intervals. Each simulation runs for three hours. The simulation results for average trip delay and network throughput are shown in Fig. 11. As shown in Fig. 11(a), the average trip delay under FTSP is consistently the worst. When demand is less than or equal to 700 veh/h, the average trip delay performance under AR-BP is better than that under ST-CRJC and DT-CRJC. However, when demand exceeds 700 veh/h, the average trip delay performance is reversed, with ST-CRJC and DT-CRJC outperforming AR-BP. The average trip delay performance under ST-CRJC and DT-CRJC is similar, though ST-CRJC performs slightly better. This means that the low-timescale network congestion information in traffic control causes an increase in trip delay, verifying Theorem 7. Further comparing the average trip delay results under ST-CRJC and DT-CRJC, we find that the low-timescale network congestion information has a greater impact on the average trip delay at medium demand (i.e., 800 veh/h - 1050 veh/h). This is a reasonable result. The network is in free flow at low demand, so there are almost no long queues at each entry node and gateway. At high demand, the network is congested or even blocked, leading to relatively stable long queues at any gateway and entry node. In contrast, the queue lengths at each gateway and entry node are unstable and vary at medium demand. Therefore, low-timescale traffic control has only a weak impact on delay performance at low demand and high demand, but has a significant impact on delay performance at medium demand.

Fig. 11(b) shows that the maximum network throughput under FTSP is 43,875 veh/h at a demand of 850 veh/h, while the maximum network throughput under AR-BP is 57,319 veh/h at a demand of 1100 veh/h. The maximum network throughput under ST-CRJC and DT-CRJC is similar, with values of 62,777 veh/h and 62,903 veh/h, respectively. These maximum values are achieved at a demand of 1200 veh/h. Compared to AR-BP, ST-CRJC and DT-CRJC increase the maximum network throughput by about 8.7%. At

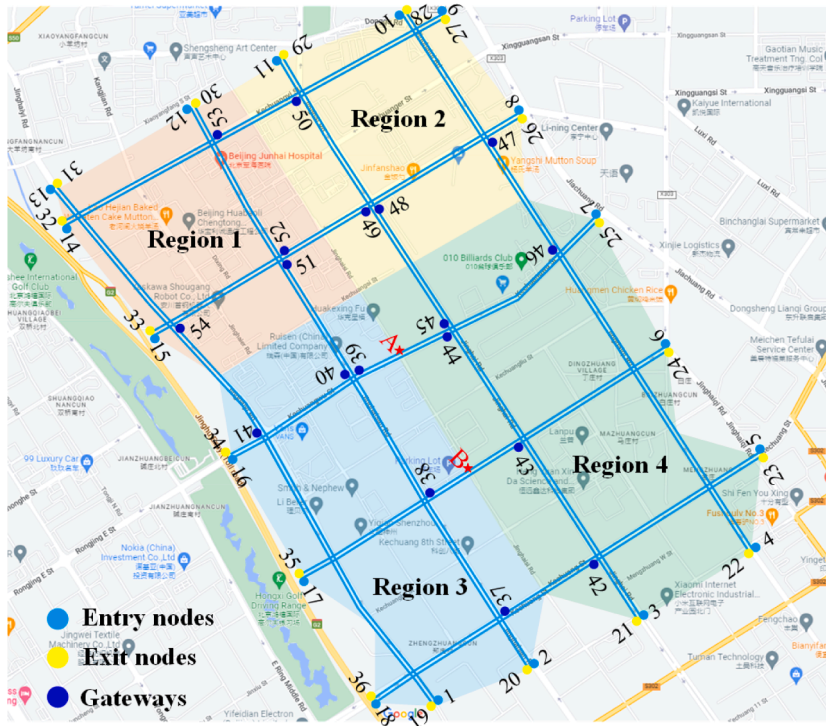


Fig. 10. The grid-like network in Beijing.

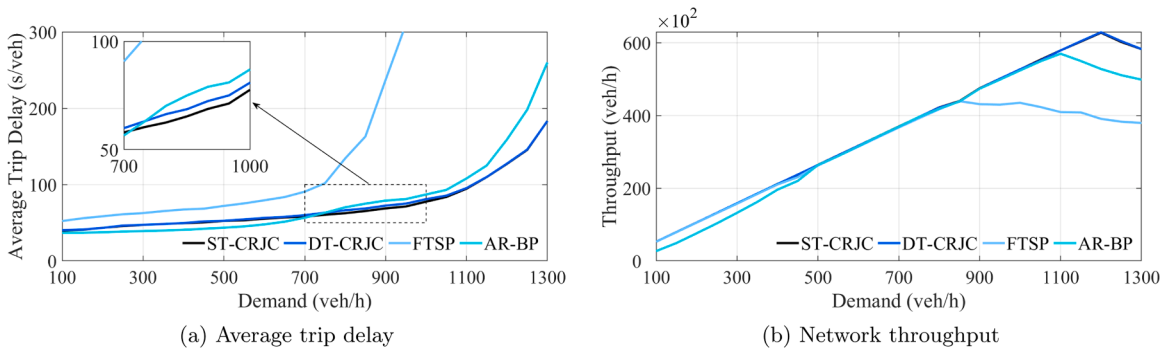


Fig. 11. The average trip delay and network throughput under different control algorithms.

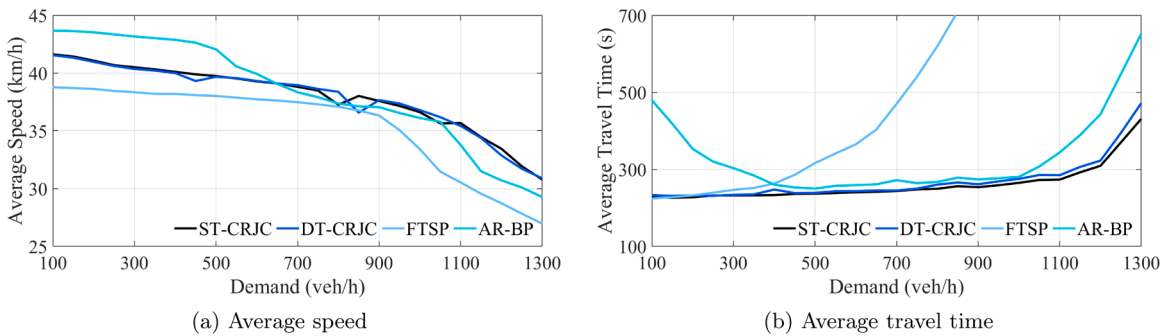


Fig. 12. The average speed and average travel time under different control algorithms.

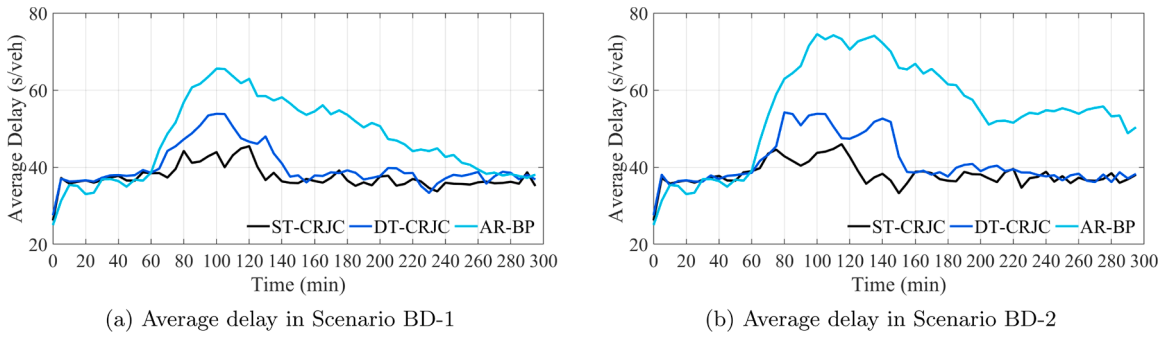


Fig. 13. Average delay under bursty demand.

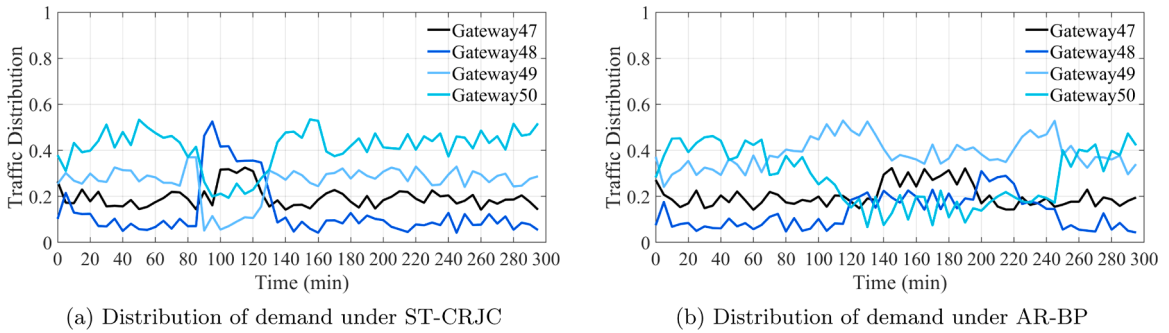


Fig. 14. Distribution of demands from region 3 to region 2 among the gateways of region 2.

any demand level, the network throughput under ST-CRJC and DT-CRJC is similar, with deviations ranging between -0.8% and $+0.3\%$. Therefore, DT-CRJC with low-timescale traffic control can still maximize network throughput, verifying [Theorem 5](#). We observe that when network throughput reaches its maximum value, network throughput decreases and average trip delay increases dramatically as demand continues to increase. According to [Li and Jabari \(2019\)](#), the demand at which the network achieves its maximum throughput can be defined as the capacity region bound. The capacity region bounds under FTSP, AR-BP, and ST-CRJC (DT-CRJC) are 850 veh/h, 1100 veh/h, and 1200 veh/h, respectively. Therefore, the proposed ST-CRJC and DT-CRJC widen the capacity region bound.

From [Fig. 11](#), we further observe that when demand is less than 500 veh/h, AR-BP, which causes the minimum average trip delay, results in the lowest network throughput. In fact, the network throughput under AR-BP is even lower than that under FTSP. To analyze the reason for this phenomenon, we examine the average speed and average travel time shown in [Fig. 12](#). In [Fig. 12](#), when demand is less than 500 veh/h, AR-BP leads to the highest average speed but also the longest average travel time. This indicates that at low demand, AR-BP reduces average trip delay but causes unnecessary detours. Hence, at low demand, using only the congestion information of the current junction for long-distance routing leads to evenly distributed vehicles, resulting in unnecessarily long paths. The proposed ST-CRJC and DT-CRJC allocate gateways to vehicles, allowing them to use short-distance gateways as intermediate destinations before reaching their final long-distance destinations. This allocation of gateways narrows the set of feasible paths, thus avoiding unnecessary detours at low demand.

5.1.2. Response to bursty demand

We set bursty demand at entry nodes 10 and 11 in region 2 to analyze the ability of ST-CRJC and DT-CRJC to respond to bursty demand. The total simulation time is five hours. The demand at entry nodes 10 and 11 from the 60th min to the 120th min is set to 1450 veh/h in Scenario BD-1 and 1550 veh/h in Scenario BD-2. The demand at entry nodes 10 and 11 during the remaining simulation time, as well as the demand at other entry nodes during the entire simulation, are set to 900 veh/h. We count the average delay every 300 s. When the average delay approaches 300 s, the network is fully blocked. The average delay results are shown in [Fig. 13](#). As can be seen from [Fig. 13](#), ST-CRJC and DT-CRJC perform well in both network recovery time and average delay. In Scenario BD-1, ST-CRJC and DT-CRJC recover the network within 10 min and 25 min, respectively, while AR-BP takes 140 min to recover the network. In Scenario BD-2, ST-CRJC and DT-CRJC take 15 min and 35 min to recover the network, respectively, while AR-BP fails to recover the network even by the 180th min after the end of bursty demand.

Since ST-CRJC and DT-CRJC are similar, we take ST-CRJC as an example to analyze why the proposed joint control enables the network to recover quickly from bursty demand. For ST-CRJC and AR-BP, [Fig. 14](#) shows the proportion of demands originating in region 3 and destined for region 2, passing through gateways 47, 48, 49, and 50 in region 2 to reach their destinations in Scenario BD-1. [Fig. 14\(a\)](#) shows that under ST-CRJC, the proportion of demands passing through gateways 50 and 49 decreases in the 5th and 15th min of bursty demand, respectively. The bursty demand from entry nodes 10 and 11 causes congestion at gateway 50. ST-CRJC

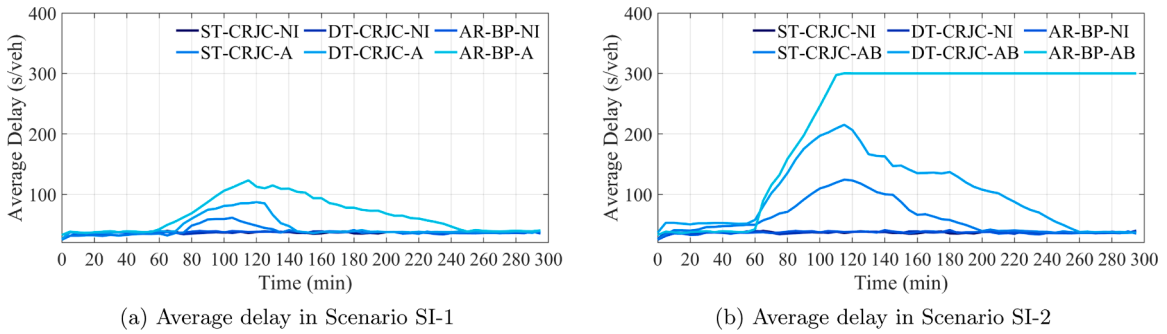


Fig. 15. Average delay under sudden incidents.

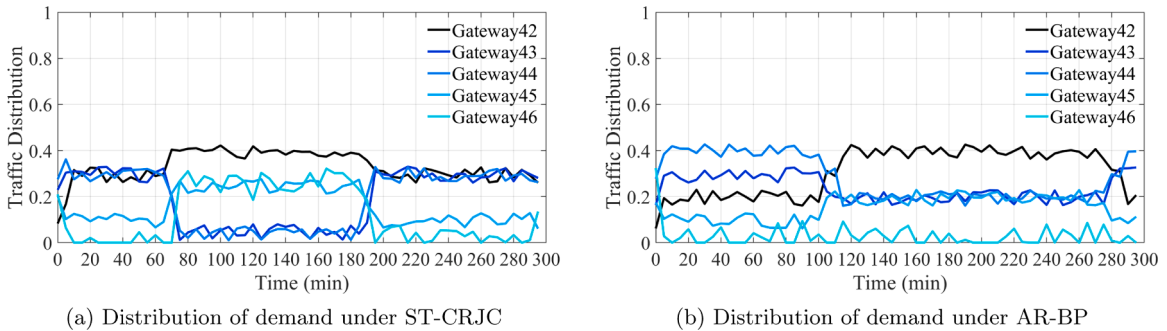


Fig. 16. Distribution of demands from region 1 to region 4 among the gateways in region 4.

reallocates the demands originating in region 3 and destined for region 2 to gateways 48 and 47, which are farther from entry nodes 10 and 11, to alleviate congestion at gateway 50. In contrast, AR-BP, which relies solely on local congestion information for routing, transfers more demands from gateway 50 to gateway 49 at the 35th min of bursty demand (see Fig. 14(b)). Therefore, ST-CRJC and DT-CRJC demonstrate strong robustness in coping with bursty demand.

5.1.3. Response to sudden incidents

We set two scenarios to test the ability of ST-CRJC and DT-CRJC to respond to sudden incidents. In Scenario SI-1, an incident occurs at point A (see Fig. 10) from the 60th min to the 120th min. In Scenario SI-2, two incidents occur simultaneously at points A and B (see Fig. 10) from the 60th min to the 120th min. The numerical simulations in both scenarios are performed for five hours, with the demand at all entry nodes set to 1100 veh/h. We measure the average delay every 300 s. The average delay results are shown in Fig. 15. “NI” denotes the simulation results with no incidents. In Scenario SI-1, ST-CRJC, DT-CRJC, and AR-BP recover the network at the 15th, 25th, and 130th min, respectively, after the incident at point A ends. In Scenario SI-2, ST-CRJC and DT-CRJC recover the network at the 80th and 140th min, respectively, after the incidents at points A and B end, while AR-BP fails to recover the network even by the 180th min after the incidents at points A and B.

Similar to Section 5.1.2, we take ST-CRJC as an example to analyze why the proposed joint control allows the network to recover quickly from sudden incidents. For ST-CRJC and AR-BP, Fig. 16 shows the proportion of demands originating in region 1 and destined for region 4, passing through gateways 42, 43, 44, 45, and 46 in region 4 to reach their destinations in Scenario SI-2. As seen in Fig. 16, when sudden incidents occur simultaneously at points A and B, ST-CRJC reallocates more demands from gateways 43 and 44 to gateways 45 and 46 to quickly respond to congestion at points A and B. AR-BP does not allocate gateways to demands in advance. Consequently, under AR-BP, demands cannot effectively adjust their routes after reaching gateways 43 and 44, resulting in significant, unrecoverable congestion.

5.1.4. Response to low-timescale congestion information exchange

To analyze the effect of traffic control with low-timescale congestion information on trip delay and network throughput, we further test the congestion information exchange delays of 30 s, 70 s, and 90 s. The simulation results for average trip delay and network throughput are shown in Fig. 17. As seen in Fig. 17(a), compared to ST-CRJC, DT-CRJC-30, DT-CRJC-50, DT-CRJC-70, and DT-CRJC-90 cause the average trip delay to increase by more than 2% in the demand intervals of 650 veh/h - 800 veh/h, 550 veh/h - 1000 veh/h, 300 veh/h - 1100 veh/h, and 300 veh/h - 1150 veh/h, respectively, and lead to an increase of 2.4%, 5.2%, 10.9%, and 11.2% in the maximum average trip delay, respectively. Therefore, as the congestion information exchange delay increases, the average trip delay increases. The growth margin of average trip delay decreases after the congestion information exchange delay exceeds 70 s. Fig. 17(b) shows that at the capacity region bound of 1200 veh/h tested in Section 5.1.1, the maximum network throughput

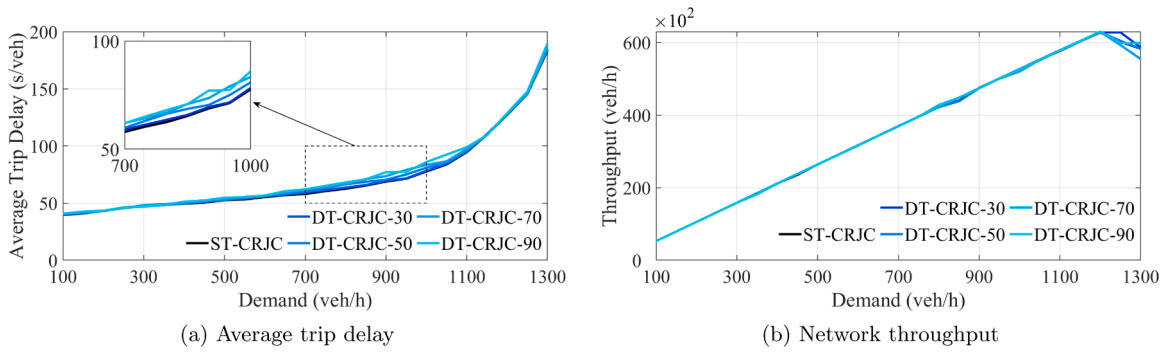


Fig. 17. Average trip delay and network throughput under different information exchange delays.

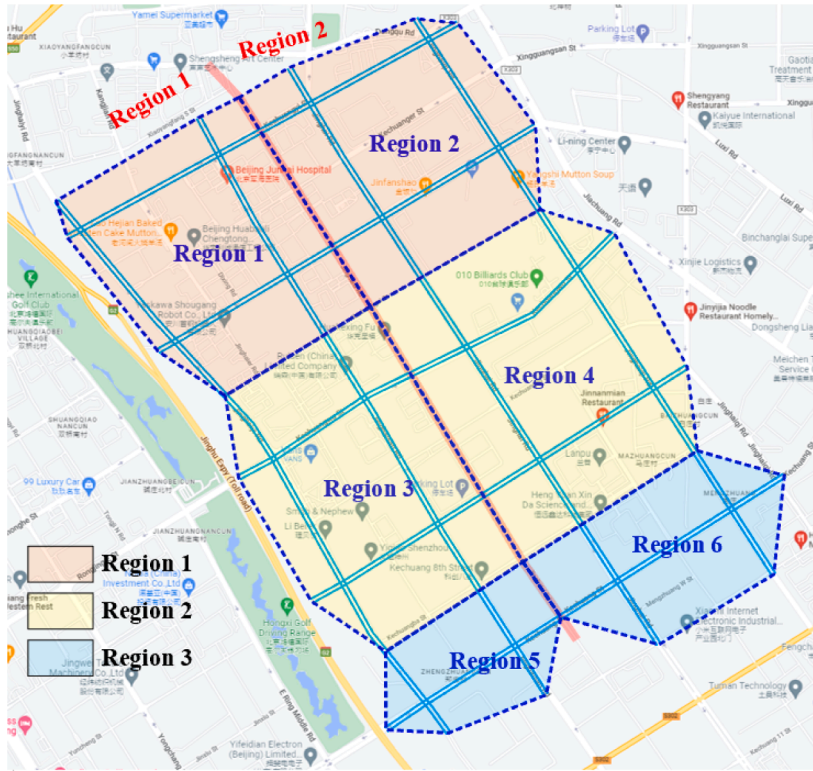


Fig. 18. Three new region divisions for the grid-like network.

under ST-CRJC, DT-CRJC-30, DT-CRJC-50, DT-CRJC-70, and DT-CRJC-90 is 62,777 veh/h, 62,793 veh/h, 62,903 veh/h, 62,945 veh/h, and 62,849 veh/h, respectively. For demand within the capacity region (i.e., demand less than or equal to 1200 veh/h), the deviation of network throughput under different congestion information exchange delays is between -1.1% and $+0.52\%$. Therefore, while congestion information exchange delay leads to an increase in trip delay, it does not reduce network throughput within the capacity region, further verifying [Theorems 2, 5, and 7](#).

5.1.5. Response to cluster size

In [Fig. 18](#), we redivide the grid-like network from [Fig. 10](#) into two regions (using the bold red line), three regions (using the orange, yellow, and blue blocks), and six regions (using the blue dotted lines) to analyze the robustness of the proposed joint control algorithms. The simulation results for network throughput and average trip delay are shown in [Fig. 19](#). The average trip delay results in [Fig. 19\(a\)](#) are similar to those in [Fig. 11\(a\)](#). Specifically, the average trip delay under ST-CRJC is higher than that under AR-BP when demand is low (i.e., less than or equal to 750 veh/h), while the opposite is true when demand is high. As seen in [Fig. 19\(b\)](#), if demand is within the capacity region (i.e., less than 1200 veh/h), ST-CRJC keeps the deviation of maximum network throughput under different region divisions to less than 2%. Therefore, the proposed joint control algorithms demonstrate strong robustness. [Fig. 19\(a\)](#) further shows that the results for average trip delay under ST-CRJC are similar when dividing the network into three

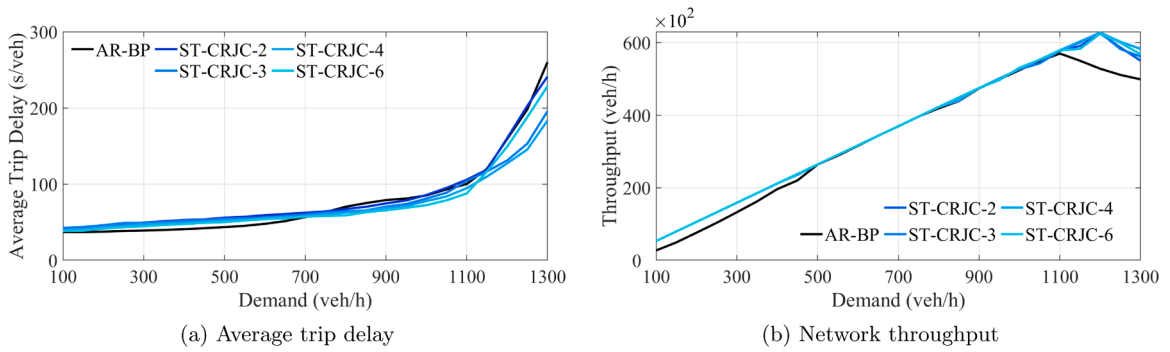


Fig. 19. Average trip delay and network throughput under different region divisions.

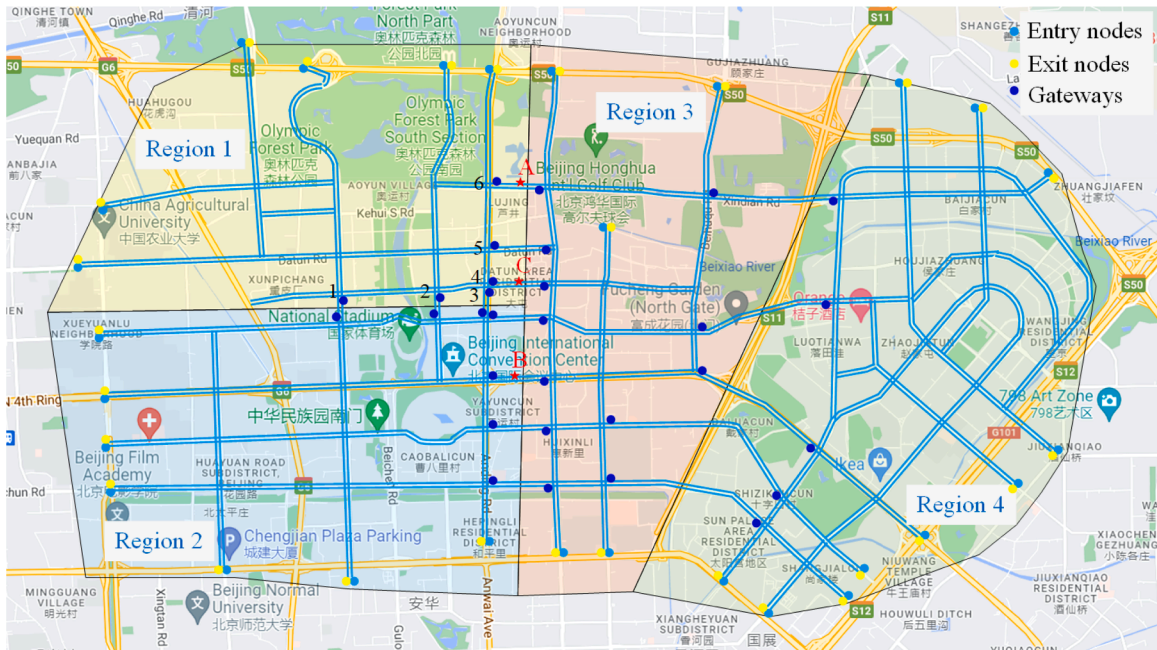


Fig. 20. The large-scale network with non-uniform link density in Beijing.

regions and four regions. The network divided into six regions under ST-CRJC performs optimally in terms of average delay when demand is less than 1100 veh/h. However, the average trip delay increases sharply when demand exceeds 1100 veh/h, similar to the results under AR-BP. When the network is divided into six regions, gateways are close to destinations (i.e., there is only one junction between gateways and destinations), causing the set of feasible paths for gateways to almost overlap with the set of feasible paths for destinations. When the network is divided into two regions, the average delay performance under ST-CRJC is consistently the worst. This is because the lengths of virtual queues for destinations are short at gateways far from the destinations. Consequently, more demands must pass through gateways far from their destinations, leading to unnecessary detours and delays.

5.2. Numerical simulations on the network with non-uniform link density

This section conducts numerical simulations on a large-scale network with non-uniform link density (see Fig. 20) to verify the robustness of the proposed joint control algorithms in response to irregular networks. The network contains 28 entry nodes (labeled in blue) and 28 exit nodes (labeled in yellow). There is an O-D pair between every entry and exit node, resulting in 784 O-D pairs in the network. The network is divided into four regions by three black lines, marked with yellow, orange, blue, and green blocks, respectively. Under this division, gateways are labeled purple in Fig. 20. The congestion information exchange delay for DT-CRJC is set to 120 s. We still set uniform demand for all entry nodes to measure the capacity region bound of the large-scale network. We test average trip delay and network throughput with demand varying from 100 veh/h to 1400 veh/h in intervals of 100. The comparative simulation results are shown in Fig. 21. Fig. 21(a) shows that the average trip delay under FTSP is consistently the worst. When demand is low, the average trip delay under AR-BP is lower than that under ST-CRJC and DT-CRJC, while the opposite

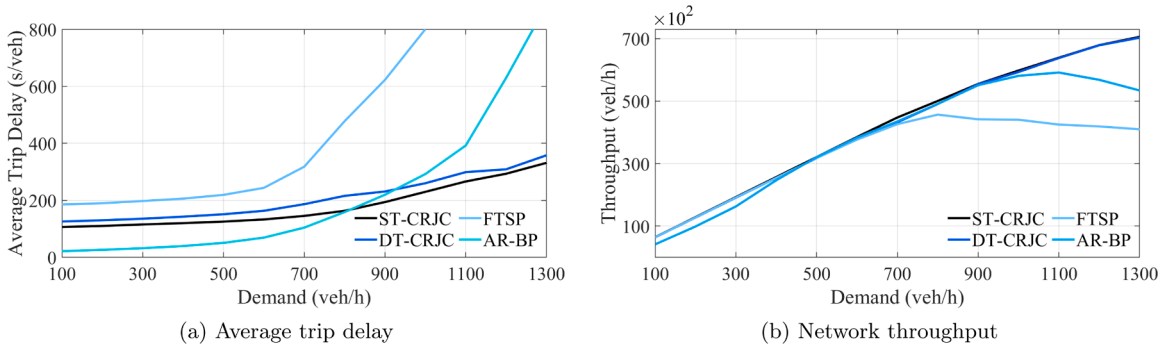


Fig. 21. The average trip delay and network throughput for the large-scale network.

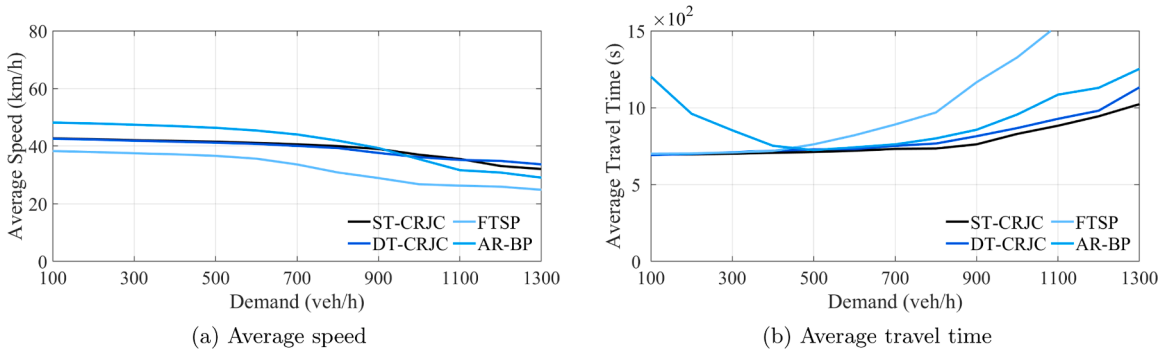


Fig. 22. The average speed and average travel time for the large-scale network.

is true when demand is high. These trip delay results are consistent with the simulation results of the grid-like network in Fig. 10. Fig. 21(b) shows that under FTSP, ST-CRJC, DT-CRJC, and AR-BP, the capacity region bounds are 800 veh/h, 1300 veh/h, 1300 veh/h, and 1100 veh/h, respectively, and the maximum network throughputs are 45,681 veh/h, 70,572 veh/h, 70,302 veh/h, and 59,159 veh/h, respectively. Compared to AR-BP, ST-CRJC and DT-CRJC widen the capacity region bound and increase the maximum network throughput by about 19.3%. Fig. 22 shows average speed and average travel time. When demand is low, AR-BP ensures that vehicles run at the highest average speed but causes the longest travel time due to unnecessary detours. The proposed ST-CRJC and DT-CRJC narrow the set of feasible paths through gateway allocation to avoid long-distance detours at low demand. Therefore, the proposed joint control algorithms are still effective for large-scale networks with non-uniform link density.

We set up two bursty demand scenarios to further test the high capacity and throughput under the joint control algorithms proposed in this study. In Scenarios BD-1 and BD-2, the demand at all entry nodes is set to 1200 veh/h and 1400 veh/h, respectively, from the 60th min to the 180th min. The total simulation runtime is five hours. During the remaining time in both scenarios, the demand at all entry nodes is set to 900 veh/h. We record the average delay every 600 s. When the average delay approaches 600 s, it indicates that the network is completely congested. The test results are shown in Fig. 23. In Scenario BD-1, ST-CRJC and DT-CRJC recover the network 10 min and 30 min, respectively, after the bursty demand ends. In Scenario BD-2, the two algorithms recover the network 30 min and 60 min, respectively, after the bursty demand ends. Conversely, in Scenario BD-1, AR-BP causes the network to deadlock during the bursty demand and only manages to reduce the average delay to 380 s/veh within 120 min after the bursty demand ends. In Scenario BD-2, AR-BP causes the network to completely collapse, with no recovery even 120 min after the bursty demand ends. This is because, under bursty demand, the network controlled by ST-CRJC and DT-CRJC does not reach saturation. Therefore, although the average delay increases under high bursty demand, the network throughput also increases, allowing the network to recover quickly. However, AR-BP causes the network to become oversaturated during bursty demand, resulting in reduced throughput and significantly increased delay. This leads to widespread network paralysis. Therefore, the proposed joint control algorithms in this study improve network throughput.

We set up two incident scenarios to further verify the robustness of our proposed algorithms through effective traffic allocation. In Scenario SI-1, incidents occur simultaneously at nodes A and B, while in Scenario SI-2, incidents occur simultaneously at nodes A, B, and C (see Fig. 20). All incidents occur from the 60th min to the 180th min, occupying the left-turn and straight-through lanes. The demand at all entry nodes is set to 900 veh/h. The total simulation runtime is five hours. We record the average delay every 600 s. The test results of the average delay are shown in Fig. 24. In Scenario SI-1, ST-CRJC, DT-CRJC, and AR-BP restore the network 10 min, 30 min, and 50 min, respectively, after the incidents end. Although DT-CRJC causes higher average delays than AR-BP before the incidents, it results in lower average delays during the incidents. In Scenario SI-2, ST-CRJC and DT-CRJC restore the network 30 min and 50 min, respectively, after the incidents end, while AR-BP still fails to restore the network 120 min after the incidents end.

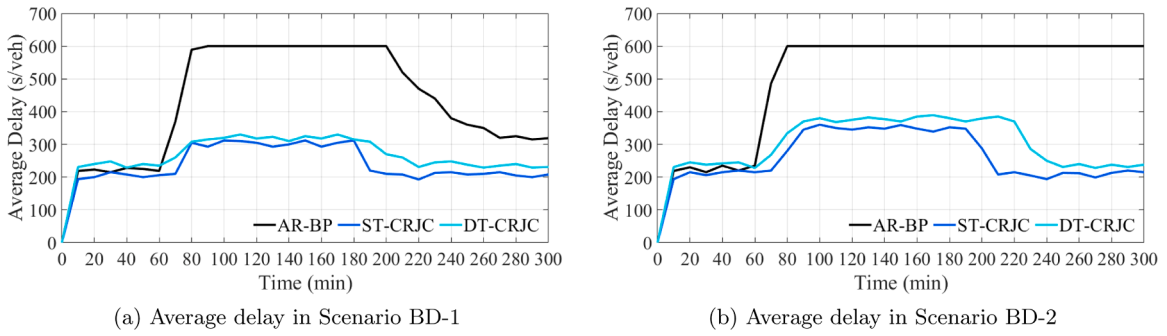


Fig. 23. Average delay under bursty demand.

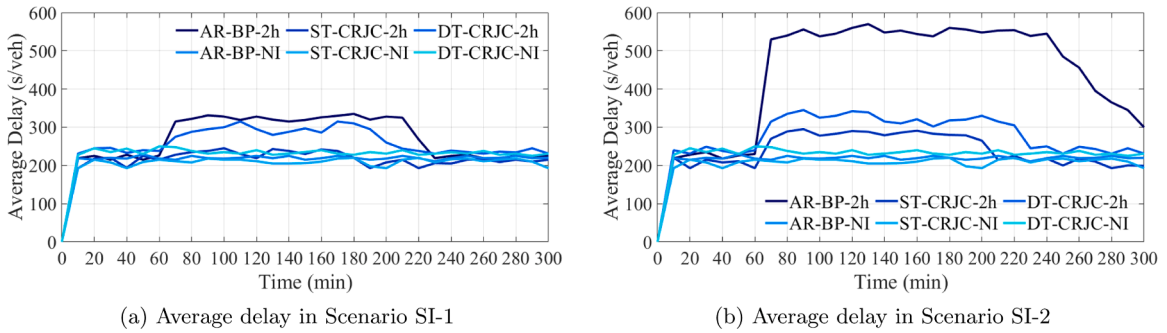


Fig. 24. Average delay under sudden incidents.

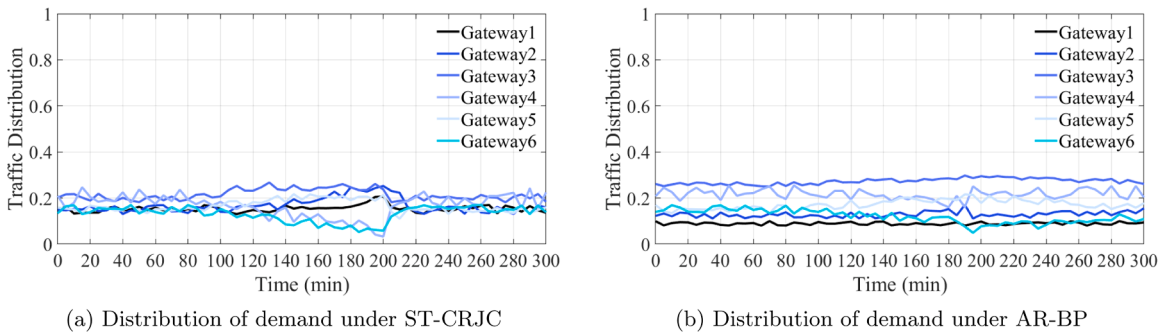


Fig. 25. Distribution of demands from regions 3 and 4 to region 1 among the gateways of region 1.

In Scenario SI-1, after incidents occur at nodes A and B, there are short-distance candidate gateways (i.e., 3, 4, and 5) for the demand from regions 3 and 4 to regions 1 and 2. Therefore, ST-CRJC, DT-CRJC, and AR-BP can all adjust the traffic in a relatively short time to respond to the incidents.

In Fig. 25, we take Scenario SI-2 as an example to show the traffic distribution of demands from regions 3 and 4 to region 1 through gateways 1 to 6 under ST-CRJC and AR-BP algorithms. The figure illustrates that when incidents occur simultaneously at nodes A, B, and C, most of the demand from regions 3 and 4 to region 1 first enters region 2 and then bypasses through gateways 1, 2, and 3 to reach region 1. The AR-BP algorithm fails to use local information to make quick judgments, causing a significant amount of traffic to wait on the right-turn lanes of gateways 4 and 6, where the incidents occur, to enter region 1. This leads to a large number of vehicles waiting at the nodes around the incident area, inducing large-scale local congestion. As a result, the network deadlocks and fails to recover in a short time.

6. Conclusion and future work

To achieve real-time control of large-scale networks and ensure provable network performance, this study divides the network into regions of different scales based on road structure and regional functions, and proposes the cross-regional cooperative joint traffic control, signal control, and routing. Traffic control implemented at entry nodes and inter-regional connection nodes balances the traffic at these connection nodes and prevents region oversaturation, respectively. Signal control and routing are applied at each

junction to optimize traffic priority and route vehicles efficiently. All three control mechanisms rely on the lengths of virtual queues at junctions and can be fully decomposed and distributed. Due to information transmission delays, traffic control using queue length information from distant nodes updates control parameters at a lower timescale than signal control and routing. Using Lyapunov drift functions, this study proves that both single-timescale and dual-timescale joint control algorithms maximize network throughput, with the latter causing higher travel delays. Comparative simulations on a small-scale grid-like network and a large-scale network with non-uniform link density from Beijing show that both joint control algorithms maximize network throughput and widen the capacity region. Additionally, both joint control algorithms demonstrate high robustness against bursty demand, sudden incidents, varying regional division scales, and special road structures.

The single-timescale and dual-timescale joint control algorithms presented in this study aim to maximize network throughput while ensuring bounded average travel delays. However, they do not account for travelers' anticipatory behavior or the fairness of delays across different flows. When travelers attempt to minimize their travel costs (i.e., dynamic user equilibrium) or demand fair/equal travel delays, equilibrium constraints are introduced into the stochastic optimization problem of maximizing network throughput, including traffic control and signal control. These non-convex and non-smooth equilibrium constraints make solving the stochastic optimization problem challenging, especially for large-scale networks. For future work, it is interesting that the space-phase-time hypernetwork from Li et al. (2015) and Li and Zhou (2017) may be introduced to consider travelers' anticipatory behavior and delay fairness. Additionally, the signal control in our joint control algorithms is based on the basic MP algorithm. However, the basic MP algorithm can result in signal phases forming unstable and unpredictable sequences, which can frustrate drivers and lead to chaotic and potentially dangerous behaviors. It is interesting future work to incorporate cyclic phases from Levin et al. (2020) and Le et al. (2015) to ensure that phases follow a fixed, orderly sequence and guarantee non-zero service times for each phase. Furthermore, it is worth to investigate asynchronous signal control and routing to form a green wave, further improving local traffic efficiency and reducing energy consumption. Lastly but not at least, it is promising to consider heterogeneous vehicles with different priorities to increase the traffic priority of buses and rescue vehicles, thereby improving social welfare.

CRedit authorship contribution statement

Shaohua Cui: Writing – review & editing, Writing – original draft, Visualization, Validation, Resources, Methodology, Investigation, Formal analysis, Data curation, Conceptualization; **Yongjie Xue:** Writing – review & editing, Writing – original draft, Visualization, Investigation, Formal analysis, Data curation, Conceptualization; **Kun Gao:** Writing – review & editing, Validation, Resources, Project administration, Investigation, Funding acquisition, Formal analysis, Data curation, Conceptualization; **Bin Yu:** Writing – review & editing, Supervision, Resources, Investigation, Data curation, Conceptualization; **Xiaobo Qu:** Writing – review & editing, Validation, Investigation, Conceptualization.

Data availability

No data was used for the research described in the article.

Acknowledgement

The research is supported by Area of Advance Transport at Chalmers University of Technology and JPI Urban Europe and Energimyndigheten (e-MATS, P2023-00029). Any opinions, findings, and conclusions or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of the sponsors.

Appendix A. Proof of Lemma 1

For any feasible virtual demand $\tilde{f} = \{\tilde{f}_o^d\}$ (i.e., $\tilde{f} \in D^0$), we can construct a service rate matrix $\{\tilde{S}_{ab}^{[o,d]}\}$ that satisfies the basic flow-based conservation constraints in Definition 1:

$$\tilde{f}_b^d \mathbf{1}_{b=0} + \sum_{a:(a,b) \in \mathcal{L}} \tilde{S}_{ab}^{[o,d]} = \sum_{c:(b,c) \in \mathcal{L}} \tilde{S}_{bc}^{[o,d]}, \forall [o, d] \in \mathcal{F} \quad (\text{A-1})$$

$$r_{ab} \sum_{ab} > \sum_{(o,d) \in \mathcal{F}} \tilde{S}_{ab}^{[o,d]}, \forall (a, b) \in \mathcal{L} \quad (\text{A-2})$$

Next, we use Eqs. (A-1) and (A-2) to prove that Eqs. (21)–(25) hold for any feasible virtual demand $\tilde{f} = \{\tilde{f}_o^d\} \in D^0$. We construct a network $G^{[o,d]}$ for each flow $[o, d]$. Let $\tilde{S}_{ab}^{[o,d]}$ be the capacity of link (a, b) . Let $L = (l_1, \dots, l_k)$ be a loop in network $G^{[o,d]}$ if $l_1 = l_k$ and $\tilde{S}_{ab}^{[o,d]} > 0$. Let $\mathcal{L}\mathcal{O}$ be the set of all loops in network $G^{[o,d]}$ and \tilde{S}_{\min}^L be $\min_{i=1, \dots, k-1} \tilde{S}_{l_i l_{i+1}}^{[o,d]}$ for $L \in \mathcal{L}\mathcal{O}$. We then construct a new network $\tilde{G}^{[o,d]}$ and define the capacity of each link as follows:

$$P_{l_i l_{i+1}}^{[o,d]} = \tilde{S}_{l_i l_{i+1}}^{[o,d]} - \sum_{L:(l_i, l_{i+1}) \in L, L \in \mathcal{L}\mathcal{O}} \tilde{S}_{\min}^L \quad (\text{A-3})$$

Hence, the new network $\tilde{G}^{[o,d]}$ is loop-free. Let $\tilde{\mathcal{L}}^{[o,d]}$ be the set of links in the new network $\tilde{G}^{[o,d]}$. The input and output rates of each node b in network $\tilde{G}^{[o,d]}$ are defined as follows:

$$p_{\text{in}(b)}^{[o,d]} = \begin{cases} \sum_{a:(a,b) \in \tilde{\mathcal{L}}^{[o,d]}} p_{ab}^{[o,d]}, & \text{if } b \neq o \\ \tilde{f}_b^d, & \text{if } b = o \end{cases} \quad (\text{A-4})$$

$$p_{\text{out}(b)}^{[o,d]} = \sum_{c:(b,c) \in \tilde{\mathcal{L}}^{[o,d]}} p_{bc}^{[o,d]} \quad (\text{A-5})$$

From Eqs. (A-1) and (A-3), we have:

$$p_{\text{out}(b)}^{[o,d]} = p_{\text{in}(b)}^{[o,d]} \quad (\text{A-6})$$

Let $R = (l_1, \dots, l_k)$ be one origin-destination route in network $\tilde{G}^{[o,d]}$ if $l_1 = o$, $l_k = d$, and $p_{l_i l_{i+1}}^{[o,d]} > 0$ for all $i = 1, \dots, k-1$. Let $\mathcal{R}^{[o,d]}$ be the set of all such origin-destination routes in network $\tilde{G}^{[o,d]}$. For the loop-free network $\tilde{G}^{[o,d]}$, the set $\mathcal{R}^{[o,d]}$ is finite. Let m_R be the last gateway on route R . m_R can divide route R into two parts $R_1 = (o, \dots, m_R)$ and $R_2 = (m_R, \dots, d)$. For each route $R = (l_1, \dots, l_k)$, we set:

$$\tilde{f}_{(l_1, \dots, l_k)}^{[o,d]} = \tilde{f}_o^d \prod_{i=1, \dots, k-1} \frac{p_{l_i l_{i+1}}^{[o,d]}}{p_{\text{out}(l_i)}^{[o,d]}} \quad (\text{A-7})$$

It is easy to verify:

$$\sum_{R \in \mathcal{R}^{[o,d]}} \tilde{f}_R^{[o,d]} = \tilde{f}_o^d \quad (\text{A-8})$$

$$\sum_{R \in \mathcal{R}^{[o,d]}, (a,b) \in R} \tilde{f}_R^{[o,d]} = p_{ab}^{[o,d]}, \forall (a,b) \in \tilde{\mathcal{L}}^{[o,d]} \quad (\text{A-9})$$

Next, $\tilde{f}_s = \left\{ \tilde{f}_o^{g,d} \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$ and $\tilde{\mathcal{S}} = \left\{ \tilde{\mathcal{S}}_{ab}^g \right\}$ are constructed based on $\left\{ \tilde{f}_R^{[o,d]} \right\}_{R \in \mathcal{R}^{[o,d]}}$. For each flow $[o,d] \in \mathcal{F}$, we have:

$$\tilde{f}_o^{g,d} = \sum_{R \in \mathcal{R}^{[o,d]}, m_R = g} \tilde{f}_R^{[o,d]}, \forall g \in \mathcal{G}_{C(d)} \quad (\text{A-10})$$

$$\tilde{f}_o^{0,d} = \tilde{f}_o^d - \sum_{g \in \mathcal{G}_{C(d)}} \tilde{f}_o^{g,d} \quad (\text{A-11})$$

For each link $(a,b) \in \mathcal{L}$, we have:

$$\tilde{\mathcal{S}}_{ab}^g = \sum_{[o,d] \in \mathcal{F}} \left(\sum_{R:(a,b) \in R_1, m_R = g} \tilde{f}_R^{[o,d]} \right), \forall g \in \mathcal{G}_{C(d)} \quad (\text{A-12})$$

$$\tilde{\mathcal{S}}_{ab}^d = \sum_{[o,d] \in \mathcal{F}} \left(\sum_{R:(a,b) \in R_2} \tilde{f}_R^{[o,d]} \right), \forall b \in \mathcal{I}_{C(d)} \quad (\text{A-13})$$

From Eqs. (10)-(A-13), we have:

$$\tilde{\mathcal{S}}_{\text{out}(b)}^g = \sum_{[o,d] \in \mathcal{F}} \left(\sum_{R:b \in R_1, m_R = g} \tilde{f}_R^{[o,d]} \right) = \tilde{\mathcal{S}}_{\text{in}(b)}^g + \sum_d \tilde{f}_b^{g,d}, \forall b \neq g \quad (\text{A-14})$$

$$\tilde{\mathcal{S}}_{\text{out}(b)}^d = \sum_{[o,d] \in \mathcal{F}} \left(\sum_{R:b \in R_2} \tilde{f}_R^{[o,d]} \right) = \tilde{\mathcal{S}}_{\text{in}(b)}^d + \mathbf{1}_{b \in \cup_c \mathcal{G}_C} \sum_o \tilde{f}_o^{b,d}, \forall b \neq d \quad (\text{A-15})$$

Hence, $\tilde{f}_s = \left\{ \tilde{f}_o^{g,d} \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$ and $\tilde{\mathcal{S}} = \left\{ \tilde{\mathcal{S}}_{ab}^g \right\}$ satisfy the flow conservation constraints in Eqs. (21)–(24). For each link $(a,b) \in \mathcal{L}$, we have:

$$\sum_{i \in \cup_c \mathcal{G}_C \cup \mathcal{I}_{C(b)}} \tilde{\mathcal{S}}_{ab}^i = \sum_{[o,d] \in \mathcal{F}} \left(\sum_{R:(a,b) \in R} \tilde{f}_R^{[o,d]} \right) = \sum_{[o,d] \in \mathcal{F}} p_{ab}^{[o,d]} \leq \sum_{[o,d] \in \mathcal{F}} \tilde{\mathcal{S}}_{ab}^{[o,d]} \quad (\text{A-16})$$

Based on Eq. (A-2), we know that $\tilde{f}_s = \left\{ \tilde{f}_o^{g,d} \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$ is feasible if $\tilde{f} = \left\{ \tilde{f}_o^d \right\}$ is supportable. The proof of Lemma 1 ends. \square

Appendix B. Proof of Eq. (29)

We convert the flow conservation constraints in Lemma 1 into the following notations to simplify proof:

$$\tilde{f}_{\text{in}(b)}^i = \begin{cases} \sum_{o:[o,d] \in \mathcal{F}} \tilde{f}_o^{b,i}, & \text{if } b \in \cup_c \mathcal{G}_C \text{ and } i \in \mathcal{I}_{C(b)} \\ \sum_{d \in \mathcal{C}(i)} \tilde{f}_b^{i,d}, & \text{if } i \in \cup_c \mathcal{G}_C \\ \tilde{f}_b^{0,i}, & \text{if } b \in \cup_c \mathcal{I}_C \text{ and } i \in \mathcal{I}_{C(b)} \end{cases} \quad (\text{B-1})$$

$$\tilde{S}_{in(b)}^i = \begin{cases} 0, & \text{if } b \in \cup_C \mathcal{G}_C \text{ and } i \in \mathcal{I}_{C(b)} \\ \sum_{a:(a,b) \in \mathcal{L}} \tilde{S}_{ab}^i, & \text{otherwise} \end{cases} \tag{B-2}$$

$$\tilde{S}_{out(b)}^i = \begin{cases} \sum_{c:(b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(i)}} \tilde{S}_{bc}^i, & \text{if } i \notin \cup_C \mathcal{G}_C \\ \sum_{c:(b,c) \in \mathcal{L}} \tilde{S}_{bc}^i, & \text{otherwise} \end{cases} \tag{B-3}$$

Based on Lemma 1, we have:

$$\tilde{f}_{in(b)}^i + \tilde{S}_{in(b)}^i = \tilde{S}_{out(b)}^i, \forall b \in \mathcal{N}, i \in \mathcal{H}_b \tag{B-4}$$

Let the matrix $\kappa(t) = \{\kappa_b^i(t)\}$ be defined as $\kappa(t) = \tilde{\mathcal{Q}}(t+1) - \tilde{\mathcal{Q}}(t)$. We have:

$$\left| \tilde{\mathcal{Q}}(t+1) \right|^2 - \left| \tilde{\mathcal{Q}}(t) \right|^2 = \left| \tilde{\mathcal{Q}}(t) + \kappa(t) \right|^2 - \left| \tilde{\mathcal{Q}}(t) \right|^2 = 2\tilde{\mathcal{Q}}^T(t)\kappa(t) + \kappa^T(t)\kappa(t) \tag{B-5}$$

Next, we determine the bounds of $\tilde{\mathcal{Q}}^T(t)\kappa(t)$ and $\kappa^T(t)\kappa(t)$ separately.

B.1. Bound on $\tilde{\mathcal{Q}}^T(t)\kappa(t)$

Based on Eqs. (17), (19), and (20), we have:

$$\tilde{\mathcal{Q}}^T(t)\kappa(t) \leq (1 + \epsilon)\lambda_1(t) + \lambda_2(t) \tag{B-6}$$

where

$$\begin{aligned} \lambda_1(t) = & \sum_{b \in \cup_C \mathcal{G}_C} \sum_{i \in \mathcal{I}_{C(b)}} \tilde{\mathcal{Q}}_b^i(t) \left[\sum_{o: [o,d] \in \mathcal{F}} \tilde{f}_o^{b,i}(t+1) \right] + \sum_{b \in \cup_C \mathcal{I}_C} \sum_{i \in \mathcal{I}_{C(b)}} \tilde{\mathcal{Q}}_b^i(t) \tilde{f}_b^{0,i}(t+1) \\ & + \sum_{b \in \mathcal{N}} \sum_{i \in \cup_C \mathcal{G}_C} \tilde{\mathcal{Q}}_b^i(t) \left[\sum_{d \in \mathcal{C}(i)} \tilde{f}_b^{i,d}(t+1) \right] \end{aligned} \tag{B-7}$$

and

$$\begin{aligned} \lambda_2(t) = & \sum_{b \in \cup_C \mathcal{I}_C} \sum_{i \in \mathcal{I}_{C(b)}} \tilde{\mathcal{Q}}_b^i(t) \left[\sum_{a:(a,b) \in \mathcal{L}} R_{ab}(t+1) \mathbf{1}_{\{e_{ab}^*(t)=i\}} \mathbf{W}_{ab}(t) \wedge \tilde{\mathcal{Q}}_a^i(t) \right] \\ & + \sum_{b \in \mathcal{N}} \sum_{i \in \cup_C \mathcal{G}_C} \tilde{\mathcal{Q}}_b^i(t) \left[\sum_{a:(a,b) \in \mathcal{L}} R_{ab}(t+1) \mathbf{1}_{\{e_{ab}^*(t)=i\}} \mathbf{W}_{ab}(t) \wedge \tilde{\mathcal{Q}}_a^i(t) \right] \\ & - \sum_{b \in \mathcal{N}} \sum_{i \in \mathcal{I}_{C(b)}} \tilde{\mathcal{Q}}_b^i(t) \left[\sum_{c:(b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(i)}} R_{bc}(t+1) \mathbf{1}_{\{e_{bc}^*(t)=i\}} \mathbf{W}_{bc}(t) \wedge \tilde{\mathcal{Q}}_b^i(t) \right] \\ & - \sum_{b \in \mathcal{N}} \sum_{i \in \cup_C \mathcal{G}_C} \tilde{\mathcal{Q}}_b^i(t) \left[\sum_{c:(b,c) \in \mathcal{L}} R_{bc}(t+1) \mathbf{1}_{\{e_{bc}^*(t)=i\}} \mathbf{W}_{bc}(t) \wedge \tilde{\mathcal{Q}}_b^i(t) \right] \end{aligned} \tag{B-8}$$

Taking the expectation of Eq. (B-7) with respect to $\tilde{\mathcal{Q}}(t)$ yields:

$$\begin{aligned} & E \left\{ \lambda_1(t) \tilde{\mathcal{Q}}(t) \right\} \\ = & \sum_b \sum_d E \left\{ \sum_{i \in \mathcal{G}_{C(b)}} \tilde{f}_b^{i,d}(t+1) \left(\tilde{\mathcal{Q}}_b^i(t) + \tilde{\mathcal{Q}}_i^d(t) \right) + \mathbf{1}_{b \in \mathcal{I}_{C(b)}} \tilde{f}_b^{0,d}(t+1) \tilde{\mathcal{Q}}_b^d(t) \mid \tilde{\mathcal{Q}}(t) \right\} \\ = & \sum_b \sum_d E \left\{ \sum_{i \in \mathcal{G}_{C(b)}} \tilde{f}_b^{i,d^*}(t+1) \min_{i \in \{0\} \cup \mathcal{G}_{C(d)}} \varpi_b^{i,d}(t) + \mathbf{1}_{b \in \mathcal{I}_{C(b)}} \tilde{f}_b^{0,d^*}(t+1) \min_{i \in \{0\} \cup \mathcal{G}_{C(d)}} \varpi_b^{i,d}(t) \mid \tilde{\mathcal{Q}}(t) \right\} \\ = & \sum_b \sum_d E \left\{ \tilde{f}_b^d(t+1) \min_{i \in \{0\} \cup \mathcal{G}_{C(d)}} \varpi_b^{i,d}(t) \mid \tilde{\mathcal{Q}}(t) \right\} \\ = & \sum_b \sum_d \tilde{f}_b^d \min_{i \in \{0\} \cup \mathcal{G}_{C(d)}} \varpi_b^{i,d}(t) \\ \leq & \sum_b \sum_d \left(\sum_{i \in \mathcal{G}_{C(b)}} \tilde{f}_b^{i,d} \left(\tilde{\mathcal{Q}}_b^i(t) + \tilde{\mathcal{Q}}_i^d(t) \right) + \mathbf{1}_{b \in \mathcal{I}_{C(b)}} \tilde{f}_b^{0,d} \tilde{\mathcal{Q}}_b^d(t) \right) \end{aligned}$$

$$= \sum_{b \in \mathcal{N}'} \sum_{i \in \mathcal{H}_b} \tilde{Q}_b^i(t) \tilde{f}_{\text{in}(b)}^i \quad (\text{B-9})$$

where $\left\{ \tilde{f}_o^{g,d*}(t+1) \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$ in the second equation is obtained by the traffic controller at entry nodes in Eq. (10) and satisfies $\tilde{f}_o^d(t+1) = \tilde{f}_o^{0,d*}(t+1) \mathbf{1}_{o \in \mathcal{I}_{C(d)}} + \sum_{g \in \mathcal{G}_{C(d)}} \tilde{f}_o^{g,d*}(t+1) \cdot \left\{ \tilde{f}_o^{g,d} \right\}_{g \in \{0\} \cup \mathcal{G}_{C(d)}}$ in the fifth equation is obtained through any stationary traffic controller at entry nodes and satisfies $\tilde{f}_o^d = \tilde{f}_o^{0,d} \mathbf{1}_{o \in \mathcal{I}_{C(d)}} + \sum_{g \in \mathcal{G}_{C(d)}} \tilde{f}_o^{g,d}$. Eq. (B-1) is applied in the last equation. Taking the expectation of Eq. (B-8) with respect to $\tilde{Q}(t)$ yields:

$$\begin{aligned} E \left\{ \lambda_2(t) | \tilde{Q}(t) \right\} &= - \sum_{(a,b) \in \mathcal{L}} E \left\{ \sum_{g \in \mathcal{U}_C \mathcal{G}_C} \left(R_{ab}(t+1) \mathbf{1}_{\{e_{ab}^*(t)=g\}} \mathbf{w}_{ab}(t) \wedge \tilde{Q}_a^g(t) \right) \left(\tilde{Q}_a^g(t) - \tilde{Q}_b^g(t) \right) \right. \\ &\quad \left. + \mathbf{1}_{b \in \mathcal{U}_C \mathcal{I}_C} \sum_{d \in \mathcal{I}_{C(b)}} \left(\mathbf{R}_{ab}(t+1) \mathbf{1}_{\{e_{ab}^*(t)=d\}} \mathbf{w}_{ab}(t) \wedge \tilde{Q}_a^d(t) \right) \left(\tilde{Q}_a^d(t) - \tilde{Q}_b^d(t) \right) \middle| \tilde{Q}(t) \right\} \\ &= \lambda_3(t) + \lambda_4(t) \end{aligned} \quad (\text{B-10})$$

where

$$\begin{aligned} \lambda_3(t) &= - \sum_{(a,b) \in \mathcal{L}} \left[\sum_{g \in \mathcal{U}_C \mathcal{G}_C} \left(r_{ab} \mathbf{1}_{\{e_{ab}^*(t)=g\}} \mathbf{w}_{ab}(t) \right) \left(\tilde{Q}_a^g(t) - \tilde{Q}_b^g(t) \right) \right. \\ &\quad \left. + \mathbf{1}_{b \in \mathcal{U}_C \mathcal{I}_C} \sum_{d \in \mathcal{I}_{C(b)}} \left(\mathbf{r}_{ab} \mathbf{1}_{\{e_{ab}^*(t)=d\}} \mathbf{w}_{ab}(t) \right) \left(\tilde{Q}_a^d(t) - \tilde{Q}_b^d(t) \right) \right] \end{aligned} \quad (\text{B-11})$$

and

$$\begin{aligned} \lambda_4(t) &= \sum_{(a,b) \in \mathcal{L}} \left[\sum_{g \in \mathcal{U}_C \mathcal{G}_C} \left(r_{ab} - E \left\{ R_{ab}(t+1) \wedge \tilde{Q}_a^g(t) | \tilde{Q}(t) \right\} \right) \mathbf{1}_{\{e_{ab}^*(t)=g\}} \mathbf{w}_{ab}(t) \left(\tilde{Q}_a^g(t) - \tilde{Q}_b^g(t) \right) \right. \\ &\quad \left. + \mathbf{1}_{b \in \mathcal{U}_C \mathcal{I}_C} \sum_{d \in \mathcal{I}_{C(b)}} \left(\mathbf{r}_{ab} - E \left\{ \mathbf{R}_{ab}(t+1) \wedge \tilde{Q}_a^d(t) | \tilde{Q}(t) \right\} \right) \mathbf{1}_{\{e_{ab}^*(t)=d\}} \mathbf{w}_{ab}(t) \left(\tilde{Q}_a^d(t) - \tilde{Q}_b^d(t) \right) \right] \end{aligned} \quad (\text{B-12})$$

Lemma B1. For all $(a, b) \in \mathcal{L}$ and t ,

$$\lambda_4(t) \leq \sum_{(a,b) \in \mathcal{L}} r_{ab} \bar{R}_{ab} \quad (\text{B-13})$$

where \bar{R}_{ab} represents the maximum of $R_{ab}(t)$ for all t .

Proof. Since the function $c \rightarrow c \wedge x$ is concave in c , we have:

$$E \left\{ R_{ab}(t+1) \wedge \tilde{Q}_a^g(t) | \tilde{Q}(t) \right\} \leq E \left\{ R_{ab}(t+1) | \tilde{Q}(t) \right\} \wedge \tilde{Q}_a^g(t) = r_{ab} \wedge \tilde{Q}_a^g(t) \leq r_{ab} \quad (\text{B-14})$$

Similarly, we have:

$$E \left\{ R_{ab}(t+1) \wedge \tilde{Q}_a^d(t) | \tilde{Q}(t) \right\} \leq E \left\{ R_{ab}(t+1) | \tilde{Q}(t) \right\} \wedge \tilde{Q}_a^d(t) = r_{ab} \wedge \tilde{Q}_a^d(t) \leq r_{ab} \quad (\text{B-15})$$

Based on Eqs. (B-14) and (B-15) and the definition of $k_{ab}(t)$ in Eq. (12), we have:

$$\begin{aligned} \lambda_4(t) &\leq \sum_{(a,b) \in \mathcal{L}} \left[\sum_{g \in \mathcal{U}_C \mathcal{G}_C} \left(r_{ab} - E \left\{ R_{ab}(t+1) \wedge \tilde{Q}_a^g(t) | \tilde{Q}(t) \right\} \right) \mathbf{1}_{\{e_{ab}^*(t)=g\}} \mathbf{w}_{ab}(t) \mathbf{k}_{ab}(t) \right. \\ &\quad \left. + \mathbf{1}_{b \in \mathcal{U}_C \mathcal{I}_C} \sum_{d \in \mathcal{I}_{C(b)}} \left(\mathbf{r}_{ab} - E \left\{ \mathbf{R}_{ab}(t+1) \wedge \tilde{Q}_a^d(t) | \tilde{Q}(t) \right\} \right) \mathbf{1}_{\{e_{ab}^*(t)=d\}} \mathbf{w}_{ab}(t) \mathbf{k}_{ab}(t) \right] \end{aligned} \quad (\text{B-16})$$

It is easy to verify for any $(a, b) \in \mathcal{L}$:

$$0 \leq r_{ab} - E \left\{ R_{ab}(t+1) \wedge \tilde{Q}_a^g(t) | \tilde{Q}(t) \right\} = \begin{cases} r_{ab} - E \left\{ R_{ab}(t+1) \wedge \tilde{Q}_a^g(t) | \tilde{Q}(t) \right\}, & \text{if } \tilde{Q}_a^g(t) \leq \bar{R}_{ab} \\ 0, & \text{if } \tilde{Q}_a^g(t) > \bar{R}_{ab} \end{cases} \quad (\text{B-17})$$

where \bar{R}_{ab} represents the maximum of $R_{ab}(t)$ for all t . This yields:

$$\left(r_{ab} - E \left\{ R_{ab}(t+1) \wedge \tilde{Q}_a^g(t) | \tilde{Q}(t) \right\} \right) \mathbf{1}_{\{e_{ab}^*(t)=g\}} \mathbf{w}_{ab}(t) \mathbf{k}_{ab}(t) \leq \mathbf{r}_{ab} \bar{R}_{ab}, \forall (a, b) \in \mathcal{L} \quad (\text{B-18})$$

Similarly, we have:

$$\left(r_{ab} - E \left\{ R_{ab}(t+1) \wedge \tilde{Q}_a^d(t) \middle| \tilde{Q}(t) \right\} \right) \mathbf{1}_{\{e_{ab}^*(t)=d\}} \mathbf{W}_{ab}(t) \mathbf{k}_{ab}(t) \leq r_{ab} \bar{R}_{ab}, \forall (a, b) \in \mathcal{L} \quad (\text{B-19})$$

Substituting Eqs. (B-18) and (B-19) into Eq. (B-16) yields:

$$\lambda_4(t) \leq \sum_{(a,b) \in \mathcal{L}} r_{ab} \bar{R}_{ab} \quad (\text{B-20})$$

The proof of Lemma B1 ends. \square

Lemma B2. Under the signal control in Eqs. (12) and (13), there exists $\vartheta > 0$ such that:

$$\lambda_3(t) + (1 + \epsilon) E \left\{ \lambda_1(t) \middle| \tilde{Q}(t) \right\} \leq -\vartheta |\tilde{Q}(t)| \quad (\text{B-21})$$

Proof. Let $W^*(t) = \{w_{ab}^*(t)\}$ be the network control matrix solved by the optimization problem in Eqs. (12) and (13). We have:

$$\sum_{(a,b) \in \mathcal{L}} w_{ab}^*(t) r_{ab} k_{ab}(t) = \max_{W \in \mathcal{W}} \sum_{(a,b) \in \mathcal{L}} w_{ab}(t) r_{ab} k_{ab}(t) = \max_{\Sigma \in Co(\mathcal{W})} \sum_{(a,b) \in \mathcal{L}} \Sigma_{ab} r_{ab} k_{ab}(t) \quad (\text{B-22})$$

Based on Lemma 1, we know that if $\tilde{f} = \{\tilde{f}_a^d\}$ is feasible, there exist $\varsigma > 0$ and $\Sigma' = \{\Sigma'_{ab}\} \in Co(\mathcal{W})$ such that $r_{ab} \Sigma'_{ab} \geq \sum_{g \in \mathcal{H}_a} \tilde{S}_{ab}^g + \varsigma$ for all $(a, b) \in \mathcal{L}$. By Proposition 1, we know that if $0 \leq \Sigma \leq \Sigma'$, then $\Sigma \in Co(\mathcal{W})$. Therefore, there exists $\Sigma \in Co(\mathcal{W})$ such that:

$$\Sigma_{ab} r_{ab} = \begin{cases} \sum_{g \in \mathcal{H}_a} \tilde{S}_{ab}^g + \varsigma, & \text{if } k_{ab} > 0 \\ 0, & \text{if } k_{ab} \leq 0 \end{cases} \quad (\text{B-23})$$

Based on the network control matrix $W^*(t) = \{w_{ab}^*(t)\}$ solved by the optimization problem in Eqs. (12) and (13), we have $\lambda_3(t) = -\sum_{(a,b) \in \mathcal{L}} w_{ab}^*(t) r_{ab} k_{ab}(t)$. By Eq. (B-9), we have:

$$\begin{aligned} \lambda_3(t) + (1 + \epsilon) E \left\{ \lambda_1(t) \middle| \tilde{Q}(t) \right\} &\leq - \sum_{(a,b) \in \mathcal{L}} w_{ab}^*(t) r_{ab} k_{ab}(t) + (1 + \epsilon) \sum_{b \in \mathcal{N}} \sum_{i \in \mathcal{H}_b} \tilde{Q}_b^i(t) \tilde{f}_{in(b)}^i \\ &\leq - \sum_{(a,b) \in \mathcal{L}} \Sigma_{ab} r_{ab} k_{ab}(t) + \sum_{b \in \mathcal{N}} \sum_{i \in \mathcal{H}_b} \tilde{Q}_b^i(t) \left(\tilde{S}_{out(b)}^i - \tilde{S}_{in(b)}^i \right) \\ &= - \sum_{(a,b) \in \mathcal{L}} \Sigma_{ab} r_{ab} k_{ab}(t) + \sum_{(a,b) \in \mathcal{L}} \left[\sum_{g \in \mathcal{U}_C \mathcal{G}_C} \tilde{S}_{ab}^g \left(\tilde{Q}_a^g(t) - \tilde{Q}_b^g(t) \right) \right. \\ &\quad \left. + \mathbf{1}_{b \in \mathcal{U}_C \mathcal{I}_C} \tilde{S}_{ab}^d \left(\tilde{Q}_a^d(t) - \tilde{Q}_b^d(t) \right) \right] \\ &\leq \sum_{(a,b) \in \mathcal{L}} \left(\sum_{g \in \mathcal{H}_a} \tilde{S}_{ab}^g - \Sigma_{ab} r_{ab} \right) k_{ab}(t) \\ &= \sum_{(a,b) \in \mathcal{L}} k_{ab}^-(t) S_{ab} - \sum_{(a,b) \in \mathcal{L}} \varsigma k_{ab}^+(t) \\ &\leq -\vartheta \sum_{(a,b) \in \mathcal{L}} |k_{ab}(t)| \end{aligned} \quad (\text{B-24})$$

where $k_{ab}^-(t) = \begin{cases} k_{ab}(t), & \text{if } k_{ab}(t) < 0 \\ 0, & \text{if } k_{ab}(t) \geq 0 \end{cases}$ and $k_{ab}^+(t) = \begin{cases} k_{ab}(t), & \text{if } k_{ab}(t) > 0 \\ 0, & \text{if } k_{ab}(t) \leq 0 \end{cases}$. Eq. (B-4) is applied in the second term of the second equation.

In the fourth equation, the fact that $\sum_{g \in \mathcal{H}_a} \tilde{S}_{ab}^g = \sum_{g \in \mathcal{U}_C \mathcal{G}_C} \tilde{S}_{ab}^g + \mathbf{1}_{b \in \mathcal{U}_C \mathcal{I}_C} \tilde{S}_{ab}^d$ and the definition of $k_{ab}(t)$ in Eq. (12) are applied. The fifth equation is obtained by Eq. (B-23). The matrix $\{k_{ab}(t)\}$ is a linear function of $\tilde{Q}(t) = \left\{ \tilde{Q}_b^i(t), b \in \mathcal{N}, i \in \mathcal{H}_b \right\}$, so there exists β such that:

$$\sum_{(a,b) \in \mathcal{L}} k_{ab}(t) \geq \beta |\tilde{Q}(t)| \quad (\text{B-25})$$

By Eqs. (B-24) and (B-25), we have:

$$\lambda_3(t) + (1 + \epsilon) E \left\{ \lambda_1(t) \middle| \tilde{Q}(t) \right\} \leq -\vartheta |\tilde{Q}(t)| \quad (\text{B-26})$$

The proof of Lemma B2 ends. \square

By Lemmas B1 and B2, we have:

$$E \left\{ \tilde{Q}^T(t) \kappa(t) \middle| \tilde{Q}(t) \right\} \leq \sum_{(a,b) \in \mathcal{L}} r_{ab} \bar{R}_{ab} - \vartheta |\tilde{Q}(t)| \quad (\text{B-27})$$

B.2. Bound on $\tilde{\mathbf{Q}}^T(t)\kappa(t)$

Based on Eq. (17), we obtain:

$$\begin{aligned} \kappa_b^i(t) &= \left[(1+\epsilon) \sum_{o:[o,d] \in \mathcal{F}} \tilde{f}_o^{b,i}(t+1) \wedge \tilde{Q}_b^i(t) \right] - \left[\sum_{c:(b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(i)}} R_{bc}(t+1) \mathbf{1}_{\{e_{bc}^*(t)=i\}} \mathbf{w}_{bc}(t) \wedge \tilde{Q}_b^i(t) \right] \\ &\leq \max_{b \in \cup_C \mathcal{G}_C, i \in \mathcal{I}_{C(b)}} \left\{ (1+\epsilon) \sum_{o:[o,d] \in \mathcal{F}} \tilde{f}_o^{b,i}, \sum_{c:(b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(i)}} \bar{R}_{bc} \right\}, \forall b \in \cup_C \mathcal{G}_C \text{ and } i \in \mathcal{I}_{C(b)} \end{aligned} \quad (\text{B-28})$$

where $\tilde{f}_o^{b,i}$ is the maximum of $\tilde{f}_o^{b,i}(t)$ for all t . Based on Eq. (19), we have:

$$\begin{aligned} \kappa_b^i(t) &= \tilde{f}_b^{0,i}(t+1) + \left[\sum_{a:(a,b) \in \mathcal{L}} R_{ab}(t+1) \mathbf{1}_{\{e_{ab}^*(t)=i\}} \mathbf{w}_{ab}(t) \wedge \tilde{Q}_a^i(t) \right] \\ &\quad - \left[\sum_{c:(b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(i)}} R_{bc}(t+1) \mathbf{1}_{\{e_{bc}^*(t)=i\}} \mathbf{w}_{bc}(t) \wedge \tilde{Q}_b^i(t) \right] \\ &\leq \max_{b \in \cup_C \mathcal{I}_C, i \in \mathcal{I}_{C(b)}} \left\{ \tilde{f}_b^{0,i} + \sum_{a:(a,b) \in \mathcal{L}} \bar{R}_{ab}, \sum_{c:(b,c) \in \mathcal{L}, c \in \mathcal{I}_{C(i)}} \bar{R}_{bc} \right\}, \forall b \in \cup_C \mathcal{I}_C \text{ and } i \in \mathcal{I}_{C(b)} \end{aligned} \quad (\text{B-29})$$

where $\tilde{f}_b^{0,i}$ is the maximum of $\tilde{f}_b^{0,i}(t)$ for all t . By Eq. (20), we have:

$$\begin{aligned} \kappa_b^i(t) &= \sum_{d \in \mathcal{C}(i)} \tilde{f}_b^{i,d}(t+1) + \left[\sum_{a:(a,b) \in \mathcal{L}} R_{ab}(t+1) \mathbf{1}_{\{e_{ab}^*(t)=i\}} \mathbf{w}_{ab}(t) \wedge \tilde{Q}_a^i(t) \right] \\ &\quad - \left[\sum_{c:(b,c) \in \mathcal{L}} R_{bc}(t+1) \mathbf{1}_{\{e_{bc}^*(t)=i\}} \mathbf{w}_{bc}(t) \wedge \tilde{Q}_b^i(t) \right] \\ &\leq \max_{b \in \mathcal{N}, i \in \cup_C \mathcal{G}_C} \left\{ \sum_{d \in \mathcal{C}(i)} \tilde{f}_b^{i,d} + \sum_{a:(a,b) \in \mathcal{L}} \bar{R}_{ab}, \sum_{c:(b,c) \in \mathcal{L}} \bar{R}_{bc} \right\}, \forall b \in \mathcal{N} \text{ and } i \in \cup_C \mathcal{G}_C \end{aligned} \quad (\text{B-30})$$

Let $M = \max_{b \in \mathcal{N}} \left\{ (1+\epsilon) \sum_{o:[o,d] \in \mathcal{F}} \tilde{f}_o^d + \sum_{a:(a,b) \in \mathcal{L}} \bar{R}_{ab}, \sum_{c:(b,c) \in \mathcal{L}} \bar{R}_{bc} \right\}$ be the maximum change in a queue length that can occur in one time slot. By Eqs. (B-28)–(B-30), we have:

$$|\kappa(t)|^2 = \sum_{b \in \mathcal{N}} \sum_{i \in \mathcal{H}_b} (\kappa_b^i(t))^2 \leq N^2 M^2 \quad (\text{B-31})$$

Substituting Eqs. (B-27) and (B-31) into Eq. (B-5) yields:

$$E \left\{ \left| \tilde{\mathbf{Q}}(t+1) \right|^2 - \left| \tilde{\mathbf{Q}}(t) \right|^2 \middle| \tilde{\mathbf{Q}}(t) \right\} = E \left\{ 2\tilde{\mathbf{Q}}^T(t)\kappa(t) + \kappa^T(t)\kappa(t) \middle| \tilde{\mathbf{Q}}(t) \right\} \leq M_1 - \vartheta \left| \tilde{\mathbf{Q}}(t) \right| \quad (\text{B-32})$$

where M_1 equals $2 \sum_{(a,b) \in \mathcal{L}} r_{ab} \bar{R}_{ab} + N^2 M^2$. The proof of Eq. (29) ends. \square

Appendix C. Proof of Theorem 2

Through Appendix A, we know that for any feasible virtual demand $\tilde{\mathbf{f}} = \{\tilde{f}_o^d\} \in D^0$, there exists a service rate matrix $\{\tilde{S}_{ab}^{[o,d]}\}$ such that:

$$\tilde{f}_b^d \mathbf{1}_{b=0} + \sum_{a:(a,b) \in \mathcal{L}} \tilde{S}_{ab}^{[o,d]} = \sum_{c:(b,c) \in \mathcal{L}} \tilde{S}_{bc}^{[o,d]}, \forall [o, d] \in \mathcal{F} \quad (\text{C-1})$$

$$r_{ab} \Sigma_{ab} > \sum_{(o,d) \in \mathcal{F}} \tilde{S}_{ab}^{[o,d]}, \forall (a, b) \in \mathcal{L} \quad (\text{C-2})$$

Since the virtual network is strongly stable, all virtual queues are positive recurrent. Hence, reliable estimates $\bar{x}_{ab}^d(t)$ can be maintained by simple averaging (e.g., exponential averaging). The routing probabilities in Eq. (16) remain near their ideal values:

$$\bar{P}_{ab}^d = \frac{\bar{x}_{ab}^d}{\sum_{c:(a,c) \in \mathcal{L}} \bar{x}_{ac}^d} \quad (\text{C-3})$$

Let S_{bc}^d be the mean arrival rate of real traffic at link (b, c) destined for d . Through the flow conservation constraint in Eq. (3), we obtain:

$$S_{bc}^d = f_b^d \mathbf{1}_{b=0} \bar{\mathbf{P}}_{bc}^d + \sum_{a:(a,b) \in \mathcal{L}} S_{ab}^d \bar{\mathbf{P}}_{bc}^d, \forall (b, c) \in \mathcal{L}, \mathbf{d} \in \mathcal{N} \quad (\text{C-4})$$

If $S_{bc}^d = \frac{\bar{x}_{bc}^d}{1+\sigma}$ for all $(b, c) \in \mathcal{L}$, Eq. (C-4) is the same as Eq. (C-1), where the fact that $\tilde{f}_o^d = (1+\sigma)f_o^d$ is applied. Then, by Eq. (C-2), we have:

$$\frac{1}{r_{bc} \sum_{bc} S_{bc}^d} = \frac{1}{(1+\sigma)r_{bc} \sum_{bc} \bar{x}_{bc}^d} < 1, \forall (a, b) \in \mathcal{L} \quad (\text{C-5})$$

Through Definition 1, we know that the demand $f = \{f_o^d\}$ in the real network is feasible. The proof of Theorem 2 ends. \square

Appendix D. Proof of Eq. (40)

Let the matrix $\hat{\kappa}(\hbar\gamma) = \{\hat{\kappa}_b^i(\hbar\gamma)\}$ be defined as $\tilde{\mathbf{Q}}((\hbar+1)\gamma) - \tilde{\mathbf{Q}}(\hbar\gamma)$. We have:

$$\left| \tilde{\mathbf{Q}}((\hbar+1)\gamma) \right|^2 - \left| \tilde{\mathbf{Q}}(\hbar\gamma) \right|^2 = \left| \tilde{\mathbf{Q}}(\hbar\gamma) + \hat{\kappa}(\hbar\gamma) \right|^2 - \left| \tilde{\mathbf{Q}}(\hbar\gamma) \right|^2 = 2\tilde{\mathbf{Q}}^T(\hbar\gamma)\hat{\kappa}(\hbar\gamma) + \hat{\kappa}^T(\hbar\gamma)\hat{\kappa}(\hbar\gamma) \quad (\text{D-1})$$

Similar to Appendix B, we determine the bounds of $\tilde{\mathbf{Q}}^T(\hbar\gamma)\hat{\kappa}(\hbar\gamma)$ and $\hat{\kappa}^T(\hbar\gamma)\hat{\kappa}(\hbar\gamma)$ separately.

D.1. Bound on $\tilde{\mathbf{Q}}^T(\hbar\gamma)\hat{\kappa}(\hbar\gamma)$

Based on the definitions of $\lambda_1(t)$ and $\lambda_2(t)$ in Eqs. (B-7) and (B-8), we have:

$$\tilde{\mathbf{Q}}^T(\hbar\gamma)\hat{\kappa}(\hbar\gamma) \leq (1+\epsilon)\eta_1 + \eta_2 \quad (\text{D-2})$$

where η_1 equals $\sum_{t=\hbar\gamma}^{(\hbar+1)\gamma-1} \lambda_1(t)$ and η_2 equals $\sum_{t=\hbar\gamma}^{(\hbar+1)\gamma-1} \lambda_2(t)$.

During time slots $\{\hbar\gamma, \dots, (\hbar+1)\gamma-1\}$, it is easy to verify that:

$$\left| \tilde{\mathcal{Q}}_b^i(t) - \tilde{\mathcal{Q}}_b^i((\hbar-1)\gamma) \right| \leq 2\gamma M \quad (\text{D-3})$$

where M is defined in Eq. (B-31). By Eqs. (D-3) and (B-7), we have:

$$\begin{aligned} \eta_1 &\leq \sum_{t=\hbar\gamma}^{(\hbar+1)\gamma-1} \sum_b \sum_d \left[\sum_{i \in \mathcal{G}_{C(b)}} \tilde{f}_b^{i,d}(t) \left(4\gamma M + \tilde{\mathcal{Q}}_b^i((\hbar-1)\gamma) + \tilde{\mathcal{Q}}_i^d((\hbar-1)\gamma) \right) \right. \\ &\quad \left. + \mathbf{1}_{b \in \mathcal{I}_{C(b)}} \tilde{\mathbf{f}}_b^{0,d}(t) \left(2\gamma M + \tilde{\mathcal{Q}}_b^d((\hbar-1)\gamma) \right) \right] \\ &\leq 4\gamma M \sum_{t=\hbar\gamma}^{(\hbar+1)\gamma-1} \sum_b \sum_d \tilde{f}_b^d(t) + \sum_{t=\hbar\gamma}^{(\hbar+1)\gamma-1} \sum_b \sum_d \left[\sum_{i \in \mathcal{G}_{C(b)}} \tilde{f}_b^{i,d}(t) \left(\tilde{\mathcal{Q}}_b^i((\hbar-1)\gamma) + \tilde{\mathcal{Q}}_i^d((\hbar-1)\gamma) \right) \right. \\ &\quad \left. + \mathbf{1}_{b \in \mathcal{I}_{C(b)}} \tilde{\mathbf{f}}_b^{0,d}(t) \left(\tilde{\mathcal{Q}}_b^d((\hbar-1)\gamma) \right) \right] \\ &\leq 4\gamma M \sum_{t=\hbar\gamma}^{(\hbar+1)\gamma-1} \sum_b \sum_d \tilde{f}_b^d(t) + \sum_{t=\hbar\gamma}^{(\hbar+1)\gamma-1} \sum_b \sum_d \left[\sum_{i \in \mathcal{G}_{C(b)}} \tilde{f}_b^{i,d^{**}}((\hbar-1)\gamma) \min_{i \in \{0\} \cup \mathcal{G}_{C(d)}} \varpi_b^{i,d}((\hbar-1)\gamma) \right. \\ &\quad \left. + \mathbf{1}_{b \in \mathcal{I}_{C(b)}} \tilde{\mathbf{f}}_b^{0,d^{**}}((\hbar-1)\gamma) \min_{i \in \{0\} \cup \mathcal{G}_{C(d)}} \varpi_b^{i,d}((\hbar-1)\gamma) \right] \\ &= 4\gamma M \sum_{t=\hbar\gamma}^{(\hbar+1)\gamma-1} \sum_b \sum_d \tilde{f}_b^d(t) + \gamma \sum_b \sum_d \min_{i \in \{0\} \cup \mathcal{G}_{C(d)}} \varpi_b^{i,d}((\hbar-1)\gamma) \tilde{f}_b^d((\hbar-1)\gamma) \end{aligned} \quad (\text{D-4})$$

where $\left\{ \tilde{f}_o^{i,d^{**}}((\hbar-1)\gamma) \right\}_{i \in \{0\} \cup \mathcal{G}_{C(d)}}$ in the third equation is obtained by the low-timescale traffic controller in Eq. (38) and satisfies

$\tilde{f}_o^d((\hbar-1)\gamma) = \tilde{f}_o^{0,d^{**}}((\hbar-1)\gamma) \mathbf{1}_{o \in \mathcal{I}_{C(d)}} + \sum_{i \in \mathcal{G}_{C(d)}} \tilde{\mathbf{f}}_o^{i,d^{**}}((\hbar-1)\gamma)$. Taking an expectation of Eq. (D-4) with respect to $\tilde{\mathbf{Q}}(\hbar\gamma)$ yields:

$$\begin{aligned} &E \left\{ \eta_1 \left| \tilde{\mathbf{Q}}(\hbar\gamma) \right. \right\} \\ &\leq 4\gamma^2 M \sum_b \sum_d \tilde{f}_b^d + \gamma \sum_b \sum_d \min_{i \in \{0\} \cup \mathcal{G}_{C(d)}} \varpi_b^{i,d}((\hbar-1)\gamma) \tilde{f}_b^d \end{aligned}$$

$$\begin{aligned}
&\leq 4\gamma^2 M \sum_b \sum_d \bar{f}_b^d + \gamma \sum_b \sum_d \left(\sum_{i \in \mathcal{C}(b)} \bar{f}_b^{i,d} \left(\tilde{Q}_b^i((h-1)\gamma) + \tilde{Q}_i^d((h-1)\gamma) \right) + \mathbf{1}_{b \in \mathcal{I}_{C(b)}} \bar{f}_b^{0,d} \tilde{Q}_b^d((h-1)\gamma) \right) \\
&= 4\gamma^2 M \sum_b \sum_d \bar{f}_b^d + \gamma \sum_{b \in \mathcal{N}} \sum_{i \in \mathcal{H}_b} \tilde{Q}_b^i((h-1)\gamma) \bar{f}_{in(b)}^i
\end{aligned} \tag{D-5}$$

where $\left\{ \bar{f}_o^{i,d} \right\}_{i \in \{0\} \cup \mathcal{C}(d)}$ in the second equation is obtained through any stationary traffic controller at entry nodes and satisfies $\bar{f}_o^d = \bar{f}_o^{0,d} \mathbf{1}_{o \in \mathcal{I}_{C(d)}} + \sum_{i \in \mathcal{C}(d)} \bar{f}_o^{i,d}$. Taking an expectation of η_2 with respect to $\tilde{Q}(\hat{h}\gamma)$ yields:

$$E\left\{ \eta_2 \middle| \tilde{Q}(\hat{h}\gamma) \right\} = \sum_{t=\hat{h}\gamma}^{(h+1)\gamma-1} \lambda_3(t) + \sum_{t=\hat{h}\gamma}^{(h+1)\gamma-1} \lambda_4(t) \leq \sum_{t=\hat{h}\gamma}^{(h+1)\gamma-1} \lambda_3(t) + \gamma \sum_{(a,b) \in \mathcal{L}} r_{ab} \bar{R}_{ab} \tag{D-6}$$

where the right equation is obtained based on Lemma B1. Based on Eqs. (D-3) and (B-11), we have during time slots $\{h\gamma, \dots, (h+1)\gamma-1\}$:

$$\lambda_3(t) = - \sum_{(a,b) \in \mathcal{L}} \Sigma_{ab} r_{ab} k_{ab}(t) \leq 4\gamma M \sum_{(a,b) \in \mathcal{L}} \Sigma_{ab} r_{ab} - \sum_{(a,b) \in \mathcal{L}} \Sigma_{ab} r_{ab} k_{ab}((h-1)\gamma) \tag{D-7}$$

Based on Lemma B2 and Eq. (D-7), we have:

$$\sum_{t=h\gamma}^{(h+1)\gamma-1} \lambda_3(t) + (1+\epsilon) E\left\{ \eta_1(t) \middle| \tilde{Q}(t) \right\} \leq -\varphi \left| \tilde{Q}(h\gamma) \right| + 4\gamma^2 M \left(\sum_b \sum_d \bar{f}_b^d + \sum_{(a,b) \in \mathcal{L}} r_{ab} \right) \tag{D-8}$$

By Eqs. (D-6) and (D-8), we have:

$$E\left\{ \tilde{Q}^T(h\gamma) \hat{\kappa}(h\gamma) \middle| \tilde{Q}(h\gamma) \right\} \leq -\varphi \left| \tilde{Q}(h\gamma) \right| + 4\gamma^2 M \left(\sum_b \sum_d \bar{f}_b^d + \sum_{(a,b) \in \mathcal{L}} r_{ab} \right) + \gamma \sum_{(a,b) \in \mathcal{L}} r_{ab} \bar{R}_{ab} \tag{D-9}$$

D.2. Bound on $\hat{\kappa}^T(h\gamma) \hat{\kappa}(h\gamma)$

As $\hat{\kappa}(h\gamma)$ equals $\tilde{Q}((h+1)\gamma) - \tilde{Q}(h\gamma)$, we have:

$$\hat{\kappa}_b^i(h\gamma) = \tilde{Q}_b^i((h+1)\gamma) - \tilde{Q}_b^i(h\gamma) \leq 4\gamma M + \sum_{t=h\gamma}^{(h+1)\gamma-1} \kappa_b^i(t) \leq 5\gamma M \tag{D-10}$$

where Eq. (B-30) is applied. Hence, the bound of $\hat{\kappa}^T(h\gamma) \hat{\kappa}(h\gamma)$ satisfies:

$$\hat{\kappa}^T(h\gamma) \hat{\kappa}(h\gamma) \leq 25\gamma^2 M^2 N^2 \tag{D-11}$$

Substituting Eqs. (D-9) and (D-11) into Eq. (D-1) yields:

$$E\left\{ \left| \tilde{Q}((h+1)\gamma) \right|^2 - \left| \tilde{Q}(h\gamma) \right|^2 \middle| \tilde{Q}(h\gamma) \right\} \leq M_2 - \varphi \left| \tilde{Q}(h\gamma) \right| \tag{D-12}$$

where M_2 equals $8\gamma^2 M \left(\sum_b \sum_d \bar{f}_b^d + \sum_{(a,b) \in \mathcal{L}} r_{ab} \right) + 2\gamma \sum_{(a,b) \in \mathcal{L}} r_{ab} \bar{R}_{ab} + 25\gamma^2 M^2 N^2$.

The proof of Eq. (40) ends. ■

References

- Aljubayri, M., Yang, Z., Shikh-Bahaei, M., 2021. Cross-layer multipath congestion control, routing and scheduling design in ad hoc wireless networks. *IET Commun.* 15, 1096–1108.
- Amputoulas, K., Santos, J. A.D., Carlson, R.C., 2020. Motorway tidal flow lane control. *IEEE Trans. Intell. Transp. Syst.* 21 (4), 1687–1696.
- Barman, S., Levin, M.W., 2023. Throughput properties and optimal locations for limited deployment of max-pressure controls. *Transp. Res. Part C* 150, 104105.
- Chen, L., Low, S.H., Chiang, M., Doyle, J.C., 2006. Cross-layer congestion control, routing and scheduling design in ad hoc wireless networks. In: 25th IEEE Int. Conf. Comput. Commun. Proc, pp. 676–688.
- Chen, R.S., Hu, J., Levin, M.W., Rey, D., 2020a. Stability-based analysis of autonomous intersection management with pedestrians. *Transp. Res. Part C* 114, 463–483.
- Chen, S., Hu, J., Shi, Y., Zhao, L., Li, W., 2020b. A vision of c-v2x: technologies, field testing, and challenges with chinese development. *IEEE Internet Things J* 7 (5), 3872–3881.
- Chow, A. H.F., Sha, R., Li, S., 2020a. Centralised and decentralised signal timing optimisation approaches for network traffic control. *Transp. Res. Part C* 113, 108–123.
- Chow, A. H.F., Sha, R., Li, Y., 2020b. Adaptive control strategies for urban network traffic via a decentralised approach with user-optimal routing. *IEEE Trans. Intell. Transp. Syst.* 21 (4), 1697–1704.
- Cui, S.H., Xue, Y.J., Gao, K., Wang, K., Yu, B., Qu, X.B., 2024. Delay-throughput tradeoffs for signalized networks with finite queue capacity. *Transp. Res. Part B* 180, 102876.
- Daganzo, C.F., 2007. Urban gridlock: macroscopic modeling and mitigation approaches. *Transp. Res. Part B* 41, 49–62.
- Diakaki, C., Papageorgiou, M., Aboudolas, K., 2002. A multivariable regulator approach to traffic-responsive network-wide signal control. *Control Eng. Pract.* 10 (2), 183–195.
- Ding, H., Guo, F., Zheng, X.Y., Zhang, W.H., 2017. Traffic guidance-perimeter control coupled method for the congestion in a macro network. *Transp. Res. Part C* 81, 300–316.

- Ding, S., Zhang, M., Xing, Y.Y., Lu, J., 2022. Revealing urban community structures by fusing multisource transportation data. *J. Transp. Eng. Part A* 10 (2), 183–195.
- Dixit, V., Nair, D.J., Chand, S., Levin, M.W., 2020. A simple crowdsourced delay-based traffic signal control. *Plos One* 15 (4), e0230598.
- Eryilmaz, A., Srikant, R., 2006. Joint congestion control, routing, and mac for stability and fairness in wireless networks. *IEEE J. Sel. Areas Commun.* 24 (8), 1514–1524.
- Espadaler-Clapés, J., Bampounakis, E., Geroliminis, N., 2023. Traffic congestion and noise emissions with detailed vehicle trajectories from uavs transp. Res. Part D 121, 103822.
- Fattah, M.A., Morshed, S.R., Kafy, A.A., 2022. Insights into the socio-economic impacts of traffic congestion in the port and industrial areas of chittagong city. *Bangladesh Transp. Eng.* 9, 100122.
- Fu, H., Chen, S., Chen, K., Kouvelas, A., Geroliminis, N., 2022. Perimeter control and route guidance of multi-region mfd systems with boundary queues using colored petri nets. *IEEE Trans. Intell. Transp. Syst.* 23 (8), 12977–12999.
- Gallager, R.G., 1996. *Finite State Markov Chains Discrete Stochastic Process*. Springer, Boston, MA, USA, pp. 103–147.
- Gao, H., Zhang, M., 2022. Arrival-based backpressure traffic signal control. *Transp. Res. Rec. J. Transp. Res. Board* 2676 (9), 172–186.
- Gartner, N., 1983. A demand-responsive strategy for traffic signal control. *Transp. Res. Rec.* 906, 75–81.
- Geroliminis, N., Haddad, J., Ramezani, M., 2013. Optimal perimeter control for two urban regions with macroscopic fundamental diagrams: a model predictive approach. *IEEE Trans. Intell. Transp. Syst.* 14 (1), 348–359.
- Ghahramani, S., 2005. *Fundamentals of Probability with Stochastic Processes*. Prentice Hall, 3rd ed. Upper Saddle River, NJ, USA. Rd ed.
- Gregoire, J., Frazzoli, E., Fortelle, A. D.L., Wongpiromsarn, T., 2014. Back-pressure traffic signal control with unknown routing rates. *IFAC-Papersonline* 47 (3), 11332–11337.
- Gregoire, J., Qian, X.J., Frazzoli, E., Fortelle, A., Wongpiromsarn, T., 2015. Capacity-aware backpressure traffic signal control. *IEEE Trans. Control Netw. Syst.* 2 (2), 164–173.
- Gregoire, J., Samaranyake, S., Frazzoli, E., 2016. Back-pressure traffic signal control with partial routing control. In: *IEEE 55th Conf. Decis. Control*, pp. 6753–6758.
- Guo, Q.Q., Ban, X.G., 2020. Macroscopic fundamental diagram based perimeter control considering dynamic user equilibrium. *Transp. Res. Part B* 136, 87–109.
- Guo, Q.Q., Ban, X.G., 2023. A multi-scale control framework for urban traffic control with connected and automated vehicles. *Transp. Res. Part B* 175, 102787.
- Guo, Q.Q., Li, L., Ban, X.G., 2019. Urban traffic signal control with connected and automated vehicles: a survey. *Transp. Res. Part C* 101, 313–334.
- Haddad, J., Ramezani, M., Geroliminis, N., 2013. Cooperative traffic control of a mixed network with two urban regions and a freeway. *Transp. Res. Part B* 54, 17–36.
- Han, X., Yu, Y., Gao, Z.Y., Zhang, M., 2021. The value of pre-trip information on departure time and route choice in the morning commute under stochastic traffic conditions. *Transp. Res. Part B* 152, 205–226.
- Henry, J., Farges, J., Tuffal, J., 1983. The prodyn real time traffic algorithm. *Proc. IFAC Control Transp. Syst.* 16 (4), 305–310.
- Hou, Z., Lei, T., 2022. Constrained model free adaptive predictive perimeter control and route guidance for multi-region urban traffic systems. *IEEE Trans. Intell. Transp. Syst.* 23 (2), 912–924.
- Huang, L., Zhao, X.Y., Chen, W., Poor, H.V., 2021. Low-latency short-packet transmission over a large spatial scale. *Entropy* 23 (7), 916.
- Huang, Z.R., Loo, B. P.Y., 2022. Urban traffic congestion in twelve large metropolitan cities: a thematic analysis of local news contents 2009–2018. *Int. J. Sustain. Transp.* 17 (6), 592–614.
- Hunt, P.B., Robertson, D.I., Bretherton, R.D., Royle, M.C., 1982. The scoot on-line traffic signal optimisation technique. *Traffic Eng. Control* 23 (4), 190–192.
- Ingole, D., Mariotte, G., Leclercq, L., 2020. Perimeter gating control and citywide dynamic user equilibrium: a macroscopic modeling framework. *Transp. Res. Part C* 111, 22–49.
- Jarjan, R., 1972. Depth-first search and linear graph algorithms. *SIAM J. Comput.* 1 (2), 146–160.
- Kuang, Y., Yen, B. T.H., Suprun, E., Sahin, O.Z., 2019. A soft traffic management approach for achieving environmentally sustainable and economically viable outcomes: an australian case study. *J. Environ. Manage.* 237, 379–386.
- Le, T., Kovács, P., Walton, N., Vu, H.L., Andrew, L. L.H., Hoogendoorn, S. S.P., 2015. Decentralized signal control for urban road networks. *Transp. Res. Part C* 58, 431–450.
- Le, T., Vu, H.L., Hoogendoorn, S.P., Kovács, P., Queija, R.N., 2017. Utility optimization framework for a distributed traffic control of urban road networks. *Transp. Res. Part B* 105, 539–558.
- Levin, M.W., Hu, J., Odell, M., 2020. Max-pressure signal control with cyclical phase structure. *Transp. Res. Part C* 120, 102828.
- Levin, M.W., Rey, D., Schwartz, A., 2019. Max-pressure control of dynamic lane reversal and autonomous intersection management. *Transp. Res. Part B* 118, 1693–1718.
- Li, L., Jabari, S.E., 2019. Position weighted backpressure intersection control for urban networks. *Transp. Res. Part B* 128, 435–461.
- Li, P.F., Mirchandani, P., Zhou, X.S., 2015. Solving simultaneous route guidance and traffic signal optimization problem using space-phase-time hypernetwork. *Transp. Res. Part B* 81, 103–130.
- Li, P.F., Zhou, X.S., 2017. Recasting and optimizing intersection automation as a connected-and-automated-vehicle (cav) scheduling problem: a sequential branch-and-bound search approach in phase-time-traffic hypernetwork. *Transp. Res. Part B* 105, 479–506.
- Lioris, J., Kurzhanskiy, A., Varaiya, P., 2016. Adaptive max pressure control of network of signalized intersections. *IFAC-Papersonline* 49 (22), 19–24.
- Liu, M.F., Han, D.C., Li, D.M., Wang, M., 2018. Route guidance during evacuations integrated with perimeter control strategy in large-scale mixed traffic flow networks. *Int. J. Modern Phys. C* 29 (11), 1850112.
- Luk, J., 1984. Two traffic-responsive area traffic control methods: scat and scoot. *Traffic Eng. Control* 25 (1), 14–22.
- Marshall, W.E., Dumbaugh, E., 2020. Revisiting the relationship between traffic congestion and the economy: a longitudinal examination of U.S. metropolitan areas. *Transportation* 47, 275–314.
- Menelaou, C., Timotheou, S., Kolios, P., Panayiotou, C.G., 2021. A convex reformulation solution approach for the joint perimeter control and route guidance problem. *IEEE Int. Intell. Transp. Syst. Conf.* pp. 2541–2546.
- Meshalkin, L.D., 1958. Limit theorems for Markov chains with a finite number of states theory probab. *Theory Probab. Appl.* 3, 335–357.
- Mirchandani, P., Head, L., 2001. A real-time traffic signal control system: architecture, algorithms, and analysis. *Transp. Res. Part C* 9 (6), 415–432.
- Mo, Z.B., Li, W.Z., Fu, Y.J., Ruan, K.R., Di, X., 2022. CVLight: decentralized learning for adaptive traffic signal control with connected vehicles. *Transp. Res. Part C* 141, 103728.
- Moradi, H., Sasaninejad, S., Wittevrongel, S., Walraevens, J., 2022. The contribution of connected vehicles to network traffic control: a hierarchical approach. *Transp. Res. Part C* 139, 103644.
- Neely, M.J., 2010. Stochastic network optimization with application to communication and queueing systems. *Synth. Lect. Commun. Netw.* 3 (1), 1–211.
- Nugmanova, A., Arndt, W.H., Hossain, M.A., Kim, J.R., 2019. Effectiveness of ring roads in reducing traffic congestion in cities for long run: big almaty ring road case study. *Sustainability* 11 (18), 4973.
- Palma, A.D., Stokkink, P., Geroliminis, N., 2022. Influence of dynamic congestion with scheduling preferences on carpooling matching with heterogeneous users. *Transp. Res. Part B* 155, 479–498.
- Papageorgiou, M., Diakaki, C., Dinopoulou, V., Kotsialos, A., Wang, Y.B., 2003. Review of road traffic control strategies. *Proc. IEEE* 91 (12), 2043–2067.
- Qu, L., He, J., Assi, C., 2015. Congestion control, routing, and scheduling in wireless networks with interference cancelation capabilities. *IEEE Trans. Veh. Technol.* 64 (7), 3108–3119.
- Rey, D., Levin, M.W., 2019. Blue phase: optimal network traffic control for legacy and autonomous vehicles. *Transp. Res. Part B* 130, 105–129.
- Sakhapov, R.L., Nikolaeva, R.V., Gatiyatullin, M.H., Makhmutov, M.M., 2016. Risk management model in road transport systems. *J. Phys. Conf. Ser.* 738, 012008.
- Sandrić, N., 2014. Recurrence and transience criteria for two cases of stable-like Markov chains. *J. Theor. Probab.* 27, 754–788.
- Sirmatel, I.I., Geroliminis, N., 2018. Economic model predictive control of large-scale urban road networks via perimeter control and regional route guidance. *IEEE Trans. Intell. Transp. Syst.* 19 (4), 1112–1121.
- Su, Z.C., Chow, A. H.F., Zhong, R.X., 2021. Adaptive network traffic control with an integrated model-based and data-driven approach and a decentralised solution method. *Transp. Res. Part C* 128, 103154.

- Tsitsokas, D., Kouvelas, A., Geroliminis, N., 2021. Efficient max-pressure traffic signal control for large-scale congested urban networks. 22nd Swiss Transp. Res. Conf. <https://doi.org/10.3929/ethz-b-000504166>
- Tsitsokas, D., Kouvelas, A., Geroliminis, N., 2022. Critical node selection method for efficient max-pressure traffic signal control in large-scale congested networks. 10th symp. Eur. Assoc. Res. Transp. <https://doi.org/10.3929/ethz-b-000550765>
- Tsitsokas, D., Kouvelas, A., Geroliminis, N., 2023. Two-layer adaptive signal control framework for large-scale dynamically-congested networks: combining efficient max pressure with perimeter control. *Transp. Res. Part C* 152. 104128
- Varaiya, P., 2013. Max pressure control of a network of signalized intersections. *Transp. Res. Part C* 36, 177–195.
- Vo, P.L., Tran, N.H., Hong, C.S., ChaeK., 2011. A joint congestion control, routing, and scheduling algorithm in multihop wireless networks with heterogeneous flows. *Int. Conf. Inf. Netw.* pp. 347–351.
- Wongpiromsarn, T., Uthacharoenpong, T., Wang, Y., Frazzoli, E., Wang, D.W., 2012. Distributed traffic signal control for maximum network throughput. 15th IEEE conf. Intell. Transp. Syst. pp. 588–595.
- Wu, J., Ghosal, D., Zhang, M., Chuah, C.N., 2018. Delay-based traffic signal control for throughput optimality and fairness at an isolated intersection. *IEEE Trans. Veh. Technol.* 67 (2), 896–909.
- Yigit, Y., Akram, V.K., Dagdeviren, O., 2021. Breadth-first search tree integrated vertex cover algorithms for link monitoring and routing in wireless sensor networks. *Comput. Netw.* 194, 108144.
- Zaidi, A.A., Kulcsár, B., Wymeersch, H., 2016. Back-pressure traffic signal control with fixed and adaptive routing for urban vehicular networks. *IEEE Trans. Intell. Transp. Syst.* 17 (8), 2134–2143.