

# **Large-scale dynamic transportation network simulation: a space-time-event parallel computing approach**

Yunchao Qu

Xuesong Zhou\* (Corresponding Author)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Abstract:** This paper describes a computationally efficient parallel-computing framework for mesoscopic transportation simulation on large-scale networks. By introducing an overall data structure for mesoscopic dynamic transportation simulation, we discuss a set of implementation issues for enabling flexible parallel computing on a multi-core shared memory architecture. First, we embed an event-based simulation logic to implement a simplified kinematic wave model and reduce simulation overhead. Second, we present a space-time-event computing framework to decompose simulation steps to reduce communication overhead in parallel execution and an OpenMP-based space-time-processor implementation method that is used to automate task partition tasks. According to the spatial and temporal attributes, various types of simulation events are mapped to independent logical processes that can concurrently execute their procedures while maintaining good load balance. We propose a synchronous space-parallel simulation strategy to dynamically assign the logical processes to different threads. The proposed method is then applied to simulate large-scale, real-world networks to examine the computational efficiency under different numbers of CPU threads. Numerical experiments demonstrate that the implemented parallel computing algorithm can significantly improve the computational efficiency and it can reach up to a speedup of 10 on a workstation with 32 computing threads.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Keywords:** synchronous parallel strategy; mesoscopic transportation simulation; space-time-event network; parallel discrete event simulation

# 1. Introduction

Compared to the sequential computing mode utilized in most existing traffic simulation and planning models, parallel computing not only efficiently utilizes widely available distributed computing powers and communication networks, but also redefines what is tractable for time-critical transportation simulation and management strategy optimization. Emerging multi-core computer processor techniques are offering unprecedented available parallel computing resources, through a wide range of high-performance laptops and desktops currently available in the market. This paper aims to develop a parallel algorithm design for transportation simulation to exploit this paradigm change in computing and to facilitate the most efficient use of emergent parallel hardware.

## 1.1. Literature review

At the core of transportation simulation, traffic flow models are interested in the quantitative relationship between flow, density and speed, and modeling the interactions between different agents. In a transportation network, there may be many routes between each origin and destination and agents choose better routes to reduce travel time. Motivated by network-wide traffic management application needs, such as regional traffic mobility analysis and real-time route guidance, dynamic traffic assignment (DTA) models has been increasingly recognized as an important approach for assessing performance of different traffic system management and information provision strategies. There are macroscopic, mesoscopic or microscopic simulation-based methods for generating time-dependent travel time measures in general traffic simulators and DTA models (Mahmassani et al., 1994; Mahmassani, 2001; Ben Akiva 2002; Peeta and Ziliaskopoulos, 2001; Adler and Blue, 2002; Celikoglu and Dell’Orco, 2007; Chen et al., 2009; Di Gangi and Cantarella, 2016; Dell’Orco et al., 2016). In an effort to reach the right balance between representation detail and computational efficiency, this study focuses on how to implement a mesoscopic-based dynamic network loading model, within a parallel computing framework, on medium and large-scale real-world networks.

One key method to achieve computational efficiency while simulating medium and large-scale networks is parallel computing. A parallel simulation implementation is valuable for researchers to quickly examine the interactions of vehicular flow and analyze the complex traffic phenomena on a larger scale. It also helps to improve the computational efficiency of traffic model validation and calibration, as well as on-line traffic state estimation and prediction, e.g., through a simulation-based optimization framework. In early studies, parallel computing techniques have been applied in several transportation simulation systems (Junchaya and Chang, 1993; Wong, 1997; Ziliaskopoulos et al, 1997). PARAMICS (Cameron and Duncan, 1996) was implemented on a Connection Machine CM-2, and its graph partition algorithm divided the set of links into different sequences of queues and each queue contained a certain

1  
2  
3  
4 number of moving vehicles. Based on a shared memory platform, [Nagel and Schmidt \(2000\)](#) and [Cetin et al. \(2003\)](#) studied the parallelization of microscopic transportation simulation based on TRANSIMS (Transportation Analysis and Simulation System) using another sophisticated graph partition algorithm. [Potuzak \(2012\)](#) reported a distributed microscopic discrete time-stepped simulator DUTS (Distributed Urban Traffic Simulator) and performed road traffic simulation on a cluster of computers with multi-core processors. [Kallioras \(2015\)](#) applied a GPU-based accelerated metaheuristics approach to solve the transit stop inspection and maintenance scheduling problem. In the field of traffic assignment, [Florian and Gendreau \(2001\)](#) offered a good review on parallel computing approaches for performing the shortest path algorithms, e.g., through network decomposition and network replication strategies. [Liu and Meng \(2013\)](#) demonstrated a solid effort for accelerating the Monte Carlo simulation method for solving probit based stochastic user equilibrium problems using a distributed computing system. [Ng and Nguyen \(2015\)](#) proposed a spatial partitioning method to implement parallel computing. [Morosan and Florian \(2015\)](#) applied a shared-memory strategy focused on parallel shortest path computation and reported that the speedup could reach up to 20 when solving the traffic assignment problem. [Auld et al. \(2016\)](#) developed an agent-based modeling software development kit POLARIS that contains a parallel discrete event simulation engine.

31 Other early implementation of parallel transportation simulations are also introduced by [Barceló et al., \(1996\)](#) and [Lee & Chandrasekar \(2002\)](#) and a parallel implementation of AIMSUN reported a speedup of 3.5 on 8 CPUs using multiple threads. As a macro-particle model, a parallel version of DYNEMO has been implemented since 2001 ([Nagel and Rickert, 2001](#)). DYNASMART's research team reported their experiment in implementing functional decomposition ([Mahmassani et al., 1994](#)). DynaMIT introduced a parallelization concept of functional decomposition (i.e., task parallelization) ([Sundaram et al., 2011](#)). Based on GPU techniques, [Zhen et al. \(2011\)](#) recently proposed a parallel computing framework to speed up the traffic simulation and optimize the traffic signal timing.

44 A significant amount of attention has been devoted to advancing parallel implementation for traffic simulation models for specific hardware/software architecture. The major efforts are summarized as follows: (1) in a static fashion, partitioning different geographical areas of the studied region to different CPU cores, and (2) for distributed computing, designing sophisticated message passing and efficient synchronization methods to reduce communication overhead among different computing cores. In our research, from a broader perspective of parallel discrete event simulation, we aim to offer a more feasible task decomposition methodology to synchronize inter-correlated space-time simulation events. This space-time-event oriented approach could take advantages of automated coordination programming interfaces (e.g. through OpenMP) between threads, processors, distributed computers, and Graphical Processing Unit (GPU).

1  
2  
3  
4 An important study by [Nie et al. \(2008\)](#) offered a comprehensive discussion on a unified dynamic  
5 network loading/simulation framework in capturing congestion propagation effects. They also clearly  
6 indicated that their double-buffer-based network loading framework could be used for further parallel  
7 simulation prototype development and system implementation. However, to ensure the actual speedup  
8 under specific parallel computing architecture, in-depth research is still critically needed to examine a  
9 number of important system implementation issues and address how to select an appropriate space-time  
10 resolution and simulation execution sequences for mesoscopic or microscopic simulation.  
11  
12  
13  
14  
15

## 16 **1.2. Space-time-event view for parallel computation**

17 Many further developments ([Fujimoto, 1990, 1993](#); [Ferscha and Tripathi, 1998](#); [Liu, 2009](#); [Fujimoto,](#)  
18 [2015](#)) summarize a number of parallel processing algorithms in terms of time-parallel and space-parallel  
19 categories. In a space-time view presented by [Chandy and Sherman \(1989\)](#), a space-time discrete event  
20 simulation can be divided into regions of arbitrary shape and assigned to separate logical processors (LP)  
21 according to the spatial and temporal decomposable features ([Liu, 2009](#)). In a parallel computing method,  
22 a global simulation task is discretized into a set of communicating logical processes (LP), each LP has its  
23 own memory space and maintains its own simulation clock and event-list, which can be concurrently  
24 executed. One LP is only capable of processing events occurring in its sub-system and communications  
25 between different LPs takes place exclusively by exchanging events.  
26  
27  
28  
29  
30  
31  
32

33 In the general field of parallel discrete event simulation, early research proposed some fundamentally  
34 important synchronization strategies, e.g., the CMB protocol ([Chandy and Misra, 1979](#); [Bryant, 1977](#)) and  
35 the time-warping method ([Jefferson, 1985](#)). In the CMB Chandy–Misra–Bryant (CMB) algorithm, LPs  
36 are assumed to be connected statically via directional links. LPs communicate through timestamped  
37 messages, also called event messages, which are transmitted from one LP to another in a non-decreasing  
38 timestamp order ([Jafer et al., 2013](#)). The CMB mechanism avoids deadlocks by introducing null messages.  
39 A null message essentially indicates the future arrival of the next message.  
40  
41  
42  
43  
44

45 Within this modeling framework, time-parallel simulation methods divide the space-time graph  
46 along the time axis into non-overlapping time intervals and assign them to different processors for parallel  
47 processing, while space-parallel simulation aims to partition the graph into a collection of space-  
48 independent subsystems. The space-parallel methods include two types of LP simulation frameworks:  
49 synchronous vs. asynchronous ([Ferscha and Tripathi, 1998](#)). In synchronous LP simulation, all LPs have  
50 the same global simulation clock and they are executed by a unique time-stepped procedure. In contrast,  
51 the asynchronous technique allows each LP to have its own local virtual time with generally different  
52 clock timestamps at a given point. The asynchronous LP simulation may cause causality errors, as some  
53 events (from the other LPs) could arrive later carrying a timestamp earlier than the target LP's current  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 simulation clock. Accordingly, a number of event-wise synchronization methods, including conservative  
5 protocol and optimistic strategy, are used for efficiently avoiding potential causality problems.  
6

7  
8 For a large-scale transportation simulator, in our view, the synchronous space-parallel LP simulation  
9 approach is a more desirable choice for several reasons. First, a transportation network is spatially  
10 consisting of sets of links and nodes and space-parallel simulation offers a more robust solution to  
11 decompose events. Second, in most mesoscopic transportation simulators, the length of simulation time  
12 intervals should be no shorter than the free-flow travel time of the shortest links in the network (e.g., 6  
13 seconds used in DYNASMART). As opposed to the possibility of extremely short event execution time  
14 intervals such as 0.0001 s between two events in a generic discrete event simulator, the discretized time  
15 interval with reasonably fine resolution (e.g. 6 s) in a mesoscopic simulator enables an efficient use of  
16 barrier synchronization available in parallel processing environment.  
17  
18  
19  
20  
21

### 22 23 **1.3. Approach**

24  
25 A spatial graph partition approach has been implemented in several traffic simulation systems. In  
26 these systems, the events are not completely or logically separated which leads to reduced computational  
27 efficiency. To improve the computational efficiency, asynchronous LP simulation strategies, including a  
28 conservative strategy and optimistic strategy, have been applied in distributed computing traffic  
29 simulation. As an optimistic synchronization protocol, Hunter et al. (2009) proposed an innovative ad hoc  
30 distributed traffic simulation framework with the key elements of space-time memory, state aggregation,  
31 and rollback based synchronization.  
32  
33  
34  
35

36  
37 In this paper, we introduce a space-time-event view to understand a parallel computation mechanism  
38 for large-scale mesoscopic traffic network simulation. Inspired by the classical conservative CMB  
39 protocol (Chandy and Misra, 1979; Bryant, 1977) and a few of early implementations by Mahmassani et  
40 al. (1994) and Nie et al. (2008), we model a direct traffic link as two event buffers. This double-buffer  
41 representation further classifies individual agent's movements as two types of events: arrival events (AE)  
42 and departure events (DE) at entrance buffers and exit buffers. This approach is consistent with the  
43 cumulative flow count-based traffic kinematic wave model proposed by Newell (1993) described in the  
44 Appendix. A unique space-time-event network-based parallel computing method is developed to  
45 schedule the AEs and DEs of all agents at different shared-memory or distributed processors. This study  
46 uses an Open Multi-Processing (OpenMP) Application Programming Interface (API) (Dagum and Enon,  
47 1998) to facilitate our program to distribute computational tasks to different processors in a shared-  
48 memory multiprocessing environment. Without directly dealing with Message Passing Interfaces (MPI),  
49 the OpenMP API also provides a simple and flexible interface for developing parallel computing  
50 programs for Graphical Computing Units (e.g., Nvidia Tesla) developed by Intel Inc.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



The rest of the paper is organized as follows. The next section presents the problem statement, network representation and a space-time-event view to decompose the event-based simulation to independent logical processes. Section 3 discusses the synchronous parallel computing implementation in a processor-space-time scheduling network. Finally, the proposed method is applied to four medium-scale and large-scale real-world networks with an examination of different speedup ratios under different CPU number configurations in Section 4.

## 2. Problem statement and framework of parallel computing

The parallel discrete event simulation system consists of events, actions (including movement actions and waiting actions), buffers, logical processes, and processors. Notations and definitions are listed below.

- $n$  Index of a physical node.
- $l$  Index of a physical link.
- $e$  Event: agent's entering or leaving event with spatial and temporal attributes.
- $b$  Buffer: upstream (entrance buffer) and downstream (exit buffer) segments of a link. A buffer stores a group of events occurring in this buffer.
- $a$  Action: connecting two adjacent agent events. There are two types of actions: movement action and waiting action.
- $m$  Movement action: A movement action connects two events in different buffers and it represents that an agent moves from a buffer to its spatially adjacent buffer. There are two categories of movement actions: link transfer and node transfer.
- $w$  Waiting action: A waiting action connects two events in the same buffer and it represents that an agent is waiting in a buffer for a certain period of time. There are two categories of waiting actions: link waiting and node waiting.
- $LP$  Logical Process: A LP contains events, buffers and actions. It independently executes a series of movements.
- $P$  Processor: the logic circuitry that responds to and processes the basic instructions, e.g., CPU and GPU. One processor may contain one or more LPs. If there are many LPs, the processor could execute the LPs sequentially, e.g., LP1->LP2->LP3.

### 2.1. Problem statement

The proposed parallel computing method aims to perform a dynamic network loading (DNL) process. Given a set of time-dependent OD or path flow on a congested network, the DNL problem determines the time-dependent link/path travel times, and other traffic states such as link flow and density over a fixed time period. Consider a transportation network with a set of nodes  $N = \{n\}$  and a set of links  $L = \{l\}$ , and the simulation time horizon is discretized to  $t = 1, \dots, T$ . In traditional methods, the given data are the demand and supply data, which include (1) the initial time-dependent OD demand matrices between activity locations or traffic analysis zones, (2) the network supply in terms of time-dependent link capacity  $q^{\max}(l, t)$ ,  $k_{jam}$  and lane miles  $\Delta X(l)$  on each link  $l$ , as well as certain capacity distribution rules around intersections and freeway bottlenecks. Given a set of path inflow patterns, the dynamic traffic network simulation problem aims to find the cumulative arrival and departure curves on each link,

1  
2  
3  
4 which also leads to the simulated time-dependent link density and travel time along each used path. The  
5 output data include (1) individual time-space trajectory in the network, and (2) aggregated time-dependent  
6 link/path travel times. These data enable transportation researchers and software developers to expand its  
7 range of capabilities to various traffic management application, e.g., traffic prediction and analysis,  
8 emission estimation, traffic demand calibration.  
9

10  
11 Specifically, the input of our core simulation model is a set of agents with given paths, which can be  
12 regarded as discretized path flows. For one OD pair, we first distribute the OD demand to path flows, then  
13 discretize the path flows to integer values, and finally generate a number of agents with these paths.  
14 Because we have assigned a certain path for each agent, the first-in first-out (FIFO) principle is adopted  
15 to describe the queue behavior and resolve capacity request conflicts at merge nodes. If there are several  
16 vehicles with the same entering time, the road capacity distributes to different paths or movements  
17 according to the path or movement flow proportion.  
18

19 [insert Fig. 1 here]

20  
21 Fig. 1a illustrates a simple freeway corridor, which consists of the link set  $\{a, b, c, d, e\}$  and the node  
22 set  $\{1,2\}$ . To clearly describe the event-based traffic simulation, a link is divided into two parts, namely  
23 entrance buffer (ENB) and exit buffer (EXB), to record agents' arriving and departing events. As shown  
24 in Fig. 1b, ENB (represented by a square) is located at the upstream of a link and EXB (represented by a  
25 triangle) is located at the downstream end. It should be noted that the exit buffer can be further  
26 decomposed to several exit buffers according to different link movements (directed to different  
27 downstream links) (Nie et al., 2008). Without a loss of generality, we only consider the single exit buffer  
28 situation in our paper.  
29

30  
31 In a parallel simulation system, each logical process (LP) possesses its own local simulation clock  
32 and local memory for private data. According to the spatial structure of a transportation network, the  
33 events of all agents can be spatially assigned into node-based LPs and link-based LPs. A node-based  $LP_n$   
34 contains the EXBs of its upstream links and ENBs of its downstream links, while a link-based  $LP_l$   
35 contains the ENB and EXB of the current link. For example, in Fig. 1c, there are two node-based LPs,  
36 that are  $LP_1$  and  $LP_2$ . The node-based  $LP_1$  consists of two exit buffers  $EXB_a$  and  $EXB_b$  and one entrance  
37 buffer  $ENB_c$  for vehicles to be loaded into the physical network. In Fig. 1d, there are five link-based LPs  
38 and each  $LP_l$  contains the  $ENB_l$  and  $EXB_l$  of link  $l$ .  
39

## 40 2.2. Framework of transportation simulation

41  
42 Fig. 2 illustrates the procedures of an overall parallel transportation assignment and simulation system.  
43 This process begins by reading the input supply and demand data from external files, and assigning a  
44 route for each agent. This routing task can be concurrently executed for each zone. After setting  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 parameters, the simulation proceeds in a time-stepped strategy and at each time stamp there are five  
5 elementary sequential tasks. A node-based LP executes the node-based actions to perform the waiting  
6 action in an exit buffer or update spatial attributes of agents, while a link-based LP executes the link-  
7 based actions to perform the waiting action in an entrance buffer or update the timestamp attributes of  
8 agents.

9  
10  
11  
12 *Task 1:* All vehicles with departure time  $t$  are loaded to the entrance buffers of their first links and some  
13 vehicles move forward to the entrance buffer if there is available space.

14  
15  
16 *Task 2:* The inflow capacity of each link can be calculated simultaneously using a traffic flow model, such  
17 as the simplified Kinematic Wave Model presented in the Appendix.

18  
19  
20 *Task 3:* Synchronize capacity distribution at each node (i.e., a bottleneck or queue server). Final outgoing  
21 capacity values for each incoming link are assigned, according to the inflow capacities of a  
22 bottleneck, using node models (e.g., [Daganzo's \(1995b\)](#) priority-based merge model).

23  
24  
25 *Task 4:* Node transfer from exit buffer to entrance buffer of connected links. Given assigned capacity, this  
26 procedure moves vehicles from the exit buffers of inbound links to the entrance buffers of  
27 outbound links or vehicles complete their trips at the destination nodes.

28  
29  
30 *Task 5:* Link transfer from the entrance buffer to exit buffer of the same link. Agents transverse from the  
31 link entrance buffers to link exit buffers by updating the travel times at the corresponding events.  
32 The arrival event at time  $t$  is deleted from the entrance buffer and replaced by a ready-to-depart  
33 event at time  $s$  at the exit buffer of the link. Here, the travel time  $s - t$  depends on the traffic flow  
34 models and Newell's simplified kinematic wave model that uses free-flow-travel-time (FFTT) to  
35 move agents from the cumulative arrival curve  $A(t)$  to the virtual cumulative departure curve  $V(t)$ .  
36 The cumulative departure curve  $D(t)$  is finally updated when the agent moves out of this link with  
37 capacity quotas assigned from Task 3.

38  
39  
40  
41  
42 [insert Fig. 2 here]

43  
44  
45  
46 The core of the traffic flow model used in this study is a discretized (vehicular) kinetic wave (KW)  
47 model. In the standard macroscopic KW traffic flow model and the node capacity distribution model, the  
48 values (i.e., traffic volume) are non-integers. In our simulation, an improved long-period pseudo-random  
49 number generator (Panneton et al., 2006) is used to generate uniform random numbers and then round  
50 floating-point capacity values to the nearest integers in terms of the number of vehicles.

51  
52  
53  
54 After the simulation process, path-specific travel times are needed for the traffic assignment module.  
55 A standard way for approximating the link and path travel time is through tracking the link cumulative  
56 flow count curves, but it may involve interpolation and back-tracking and constitute a significant portion  
57 of the computational overhead. Snelder (2009) has proposed a method to reduce the error caused by the  
58  
59  
60  
61

1  
2  
3  
4 discretization rounding off errors. Different from the above flow-based method, we use vehicle  
5 trajectories from all agents to compute the time-dependent (experienced) path travel time and waiting  
6 time, based on the time stamps along individual paths. Define  $t_p^a(\tau)$  as the travel time of agent  $a$  travels  
7 along path  $p$  with departure time  $\tau$ ,  $t(a, l_p^{last}, EXB)$  as the time interval of agent  $a$  arrives the exit buffer  
8 (EXB) of the last link of path  $l_p^{last}$ . During the simulation, the individual entering time and leaving time  
9 of each buffer of each link are dynamically recorded, so the individual time-dependent path travel time  
10  $t_p^a(\tau)$  and the aggregated average path travel time  $\bar{t}_p(\tau)$  can be easily calculated by Eq. (1).

$$11 \quad t_p^a(\tau) = t(a, l_p^{last}, EXB) - \tau, \quad \bar{t}_p(\tau) = \frac{\sum_a t_p^a(\tau)}{N} \quad (1)$$

12  
13 By following the time discretization approach presented by Lu et al. (2009), we generate time-dependent  
14 aggregated link travel cost using the shortest path algorithm and dynamic equilibrium gaps.

### 15 2.3. Space-time-event view for parallel computing

16  
17 In this paper, we use a four-indexed notation  $(a, l, t, B)$  to record the space-time event (or state) of  
18 each agent, where  $a$  represents the index of an agent,  $l$  represents the index of a link,  $t$  represents the  
19 simulation time interval, and  $B$  represents the buffer type of the current link  $l$ . The space-time trajectory  
20 of one agent is recorded as an event list. For example, in Fig. 3, the space-time trajectory of one agent  $a$   
21 in the corridor  $l_1 \rightarrow l_2 \rightarrow l_3$  can be represented by the event list  $(a, l_1, t_1, ENB) \rightarrow (a, l_1, t_3, EXB) \rightarrow$   
22  $\dots \rightarrow (a, l_3, t_9, EXB)$ . Accordingly, the input buffers store the events  $(a, l, t, ENB)$  that represent agents  
23 entering the upstream node of link  $l$  while the output buffers store the events  $(a, l, t, EXB)$  that represent  
24 agents entering and leaving the downstream node of link  $l$ . For a single agent/vehicle, the execution  
25 sequence of its events is dynamically driven by three categories of actions

- 26 1) Link transfer:  $(a, l, t, ENB) \rightarrow (a, l, s, EXB)$ , with a travel time of  $s - t$  ( $s > t$ ).
- 27 2) Node transfer:  $(a, l, t, EXB) \rightarrow (a, m, t, ENB)$ . The downstream link  $m$  is predefined ahead of the  
28 agent's trip or determined by real-time routing behavior. In our study, we assume that the  
29 operation time of node transfer movement is zero to establish a clear activity-on-the-link network  
30 structure.
- 31 3) Node waiting:  $(a, l, t, EXB) \rightarrow (a, l, s, EXB)$ . Agent  $a$  is waiting at the exit buffer  $EXB$  of link  $l$   
32 from time interval  $t$  and is ready to be considered for a move to the next link at time interval  $s$ .  
33 The waiting time  $s - t$  depends on the inflow capacities of the adjacent downstream links and  
34 outflow capacity of the current link  $l$ .

35 [insert Fig. 3 here]

36  
37 In a multi-agent simulation, tasks 2 and 3 proposed in Section 2.2 calculate capacity requests and  
38 allocation prior to moving agents, and during this process the agents' spatial and temporal attributes are

1  
2  
3  
4 not updated. After these two tasks, agents with permission perform node-transfer movements and link-  
5 transfer movements while other agents execute node-waiting actions/commands.  
6  
7  
8  
9

### 10 **3. Scheduling logical processors for parallel computing implementation**

11 Given a space-time-event representation presented above for an overall transportation simulation, we  
12 then need to execute all LPs and events synchronously or asynchronously in parallel. An important local  
13 causality constraint (Fujimoto, 1990) principle requires that every LP should process the events in a non-  
14 decreasing timestamp order. This section systematically describes an overall control/coordination  
15 mechanism to schedule the parallel processors to attain maximum speed-up.  
16  
17  
18  
19  
20

#### 21 **3.1. Synchronous LP simulation**

22 In most mesoscopic transportation simulators, the length of simulation time intervals are not shorter  
23 than the free-flow travel time of the shortest links in the network, e.g., 6 seconds used in DYNASMART.  
24 This discretized time horizon is very useful for conducting synchronous space-parallel LP simulation.  
25 Problems of deadlock and memory request are avoided with the help of the barrier synchronization  
26 mechanisms available in most every parallel processing environments (Ferscha and Tripathi, 1998). On  
27 the other hand, LPs with shorter execution times should wait for all other LPs to complete their  
28 executions and thus, this synchronous strategy could lead to idle time. If events are spatially and  
29 temporally distributed to all LPs in an even manner, the CPU time differences could be reduced, leading  
30 to less system-wide idle time.  
31  
32  
33  
34  
35  
36

37 While there are many advantages of adopting the synchronous LP strategy, we also need to  
38 recognize the following limitations. First, the time discretization might result in temporal precision errors  
39 as the occurrence time of every event should be integral multiples of the simulation time interval. Higher  
40 time precision requires finer simulation time resolution. Second, events are spatially decomposed into  
41 different LPs and each logical process can simultaneously manage non-overlapping sets of buffers. This is  
42 shown in Fig. 1-(c) for a node-based partition and Fig. 1-(d) for a link-based partition. It should be noted  
43 that under the space-parallel strategy, each LP can only communicate with their neighbors, which requires  
44 the distance  $d$  traveled by an agent in one time step  $\Delta t$  not to exceed  $\Delta l$  (condition (2)).  
45  
46  
47  
48  
49  
50

$$51 \quad d \leq \Delta l = v_{\max} \cdot \Delta t \quad (2)$$

52 The LP parallel computing strategy can be further considered to a possible application in the  
53 microscopic traffic or pedestrian flow models that satisfy condition (1). Examples along line could  
54 include the floor field cellular automata model (e.g., Kirchner and Schadschneider, 2002) and force-based  
55 model (e.g., Qu et al., 2014; Qu et al., 2015) for pedestrian simulation, as well as the cellular  
56 transmission model for traffic simulation (Daganzo, 1994). In these models, an individual can only move  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 to one of its neighbor cells/links in a single time stamp. However, in some other microscopic models,  
5 such as the cellular automata model for traffic simulation by Nagel and Schreckenberg (1992), a vehicle  
6 can move more than one cell/link at a time interval. In such a situation, condition (1) is not satisfied and  
7 the parallel strategy could be invalidated because it cannot exactly predict the potential conflict area and  
8 accordingly resolve the conflicts. To deal with such situations, some special microscopic models, such as  
9 simple linear car-following model (CF (L)), two cellular automata models (CA (L, M)), proposed by  
10 Daganzo (2006), can be converted to the kinematic wave model (KW) with a triangular fundamental  
11 diagram in order to enable the use of our proposed parallel computing strategy with simplified traffic flow  
12 models. For more complicated models considering lane changing and overtaking behaviors, it is still a  
13 very challenging task to design an effective parallel computing method.  
14  
15  
16  
17  
18  
19  
20  
21

22 The pseudo-code of the parallel simulation process is described as follows:

23 For time  $t = 0$  to  $T$  (e.g., 24 hours)

24 Step 1: Load newly generated agents to the loading buffer of their first links. If there is sufficient space  
25 available for each agent entering the entrance buffer, move it from the entrance buffer to the exit  
26 buffer.  
27

28 Step 2: Capacity request estimation for each link

29 #pragma omp parallel for

30 For link  $l = 1$  to  $L$

31 Based on the external user input and traffic flow model, calculate inflow capacity for link  $l$

32 Use a link model and outflow capacities, adjust inflow capacity for link  $l$

33 End for // link

34 Step 3: Synchronize capacity distribution at each bottleneck (merge and diverge)

35 #pragma omp parallel for

36 For each node  $i = 1$  to  $N$

37 Total Incoming Demand = 0;

38 For each link  $l$  in incoming link set  $IL(i)$

39 Total Incoming Demand += outflow demand from link  $l$

40 End for each incoming link

41 If (Total Incoming Demand > Inflow Link Capacity  $cap^{out}(l')$ ) //bottleneck

42 For each incoming link  $IL(i)$

43 Distribute inflow capacity to outflow capacity for each incoming link at bottleneck

44 End for each incoming link

45 EndIf

46 End for each node

47 Step 4: Node transfer from exit buffer to entrance buffer of connected links

48 #pragma omp parallel for

49 For each node  $i = 1$  to  $N$

50 For each link  $l$  in the incoming link set  $IL(i)$

51 While (fetch the front vehicle from the exit buffer)

52 If current simulation time  $t$  equals to or is greater than  $t_v$ , and both inflow and  
53 outflow capacity are available (for kinematic wave based queuing models) //

54 Check exit buffer

55 Move the vehicle from the exit buffer of link  $l$  to the entrance buffer of link  $l'$

56 Update cumulative arrival agent list on outgoing link  $l'$   
57  
58  
59  
60  
61  
62  
63  
64  
65

```

1
2
3
4         Reduce the inflow capacity on  $l'$  at time  $t$  by 1
5         Reduce the outflow capacity on  $l$  at time  $t$  by 1
6     Else
7         Break
8     End while
9 End for each incoming link
10 End for each node
11 Step 5: Link traversal from the entrance buffer to exit buffer of the same link
12 #pragma omp parallel for
13 For each link  $l=1$  to  $L$ 
14     Arrival events at time  $t_a$ , departure events at time  $t_v = t_a + FFTT(l)$ , create the event into exit
15     buffer,
16     Update link cumulative arrival curve  $A(l, t)$ 
17 End for
18 End for //time
19
20
21

```

Fig. 4 briefly explains the fork-joint parallelism mechanism in OpenMP. A parallel region is a block of one or more statements with one point of entry at the top (*fork point*) and one point of exit at the bottom (*joint point*). Beginning from the code structure with hashtag “#pragma omp parallel for”, a number of threads are created at the fork point to concurrently execute the decomposed tasks and then threads wait here for all threads to end at the joint point. Typically, a processor can only work on one thread per core. A CPU can have multiple cores and the hyper threading technology can allow one CPU to create and work on up to two threads per core.

[insert Fig. 4 here]

All CPU threads share the memory during the execution procedure in the structured parallel block, and OpenMP is able to dynamically create a private memory block for each task in a thread. Define  $t_f$  as the time stamp at the fork point,  $t_j$  as the time stamp at the joint point,  $t_{i,s}$  as the time stamp of thread  $i$  being created, and  $t_{i,c}$  as the time stamp of thread  $i$  completing its executing. Here,  $i$  is the total number of parallel executed threads. It should be noted that the time stamps are real computation time stamps, which are not simulation time stamps. The relationships between the time stamps can be expressed as Eq. (3).

$$t_f = \min\{t_{1,s}, \dots, t_{i,s}\}, t_j = \max\{t_{1,c}, \dots, t_{i,c}\} \quad (3)$$

It should be noted that OpenMP automates the distribution of the LPs (e.g. for different nodes, links or zones) to different processors during the simulation, and thus, the LP distributions may not be the same at different computing instances. Consider Fig. 5 as an example where there are three LPs (LP<sub>1</sub>, LP<sub>2</sub>, LP<sub>3</sub>) and four LPs (LP<sub>1</sub>, LP<sub>4</sub>, LP<sub>6</sub>, LP<sub>7</sub>) allocated in CPU thread 1 at two time instances, respectively.

1  
2  
3  
4 [insert Fig. 5 here]

5  
6 An important property for parallel computing is that different numbers of CPU threads and different  
7 LP distributions should not influence the simulation results. Thanks to the flexible coding structure in  
8 OpenMP, our program attaches its own random number seed for each node, link, or zone shown in Fig. 2.  
9 Thus, the computing process and results for each object (to be parallelized) remain the same, regardless of  
10 which CPU thread it is assigned to or how many CPU threads created by the system. We have proved this  
11 concept according to our simulation results.  
12  
13  
14  
15

### 16 **3.2. Understanding the execution of simulation activities at multiple processors**

17  
18 The atomic traffic flow simulation logic at each time stamp consists of five sequential tasks, each of  
19 which is executed by a link-based LP or a node-based LP, as shown in Fig. 2. Focusing on the link-  
20 transfer and node-transfer tasks, Fig. 6 illustrates the space-time view of the dynamic LP decomposition  
21 procedure for a simple road segment example, which is a traffic corridor with two merge and diverge  
22 ramps. To simply describe the execution procedure, we only focus on nodes 1 and 2 with a little  
23 complicated topological structure. The simulation procedure has five link-based LPs and two node-based  
24 LPs, e.g., node-based LP<sub>1</sub> (node 1) contains the exit buffer of links *a* and *b* and the entrance buffer of link  
25 *c*.  
26  
27  
28  
29  
30  
31

32 [insert Fig. 6 here]

33  
34 Constructing a space-time-processor graph (Chandy and Sherman, 1989) is useful for us to  
35 understand how independent link-based and node-based LPs are assigned to different processors. This  
36 graph consists of a set of activity and barrier vertexes and a set of schedule arcs. Specifically, an activity  
37 vertex  $v(i, t, p)$  represents a LP being executed at space *i* (a link *l* or a node *n*) at time *t* on processor *p*.  
38 A barrier vertex  $v^*(i^*, t, p^*)$  represents a forking or joint activity at a virtual space *i*\* at time *t* on  
39 processor *p*\*, and it represents the intersection of results before re-distribution back to the CPUs. A  
40 schedule arc  $(v, v^*)$  connects an activity vertex and a barrier vertex. The schedule network is sequentially  
41 executed along the time dimension (in terms of simulation steps) and parallel computing activities are  
42 concurrently executed in the space dimension inside each link-based LP or node-based LP stage.  
43  
44  
45  
46  
47  
48

49 [insert Fig. 7 here]

50  
51 As illustrated in Fig. 7, we consider an example in a corridor (1->2->3) including 3 nodes (node #1,  
52 2, 3) and 2 links (link #a, b), on a system with 2 CPUs. During the link-transfer step, link-based LP<sub>a</sub> and  
53 LP<sub>b</sub> are concurrently executed on CPU#1 and CPU#2, respectively; while inside the node-transfer step,  
54 the node-based LP<sub>1</sub> and LP<sub>2</sub> are performed on CPU#1, while LP<sub>3</sub> is performed on CPU#2. Between every  
55 two steps, there is a barrier vertex connecting the activity vertexes to synchronize all LPs. Fig. 8  
56 illustrates the execution time flow chart of an instance of Fig. 7. Because fewer LPs are assigned to CPU2  
57  
58  
59  
60  
61  
62  
63  
64  
65



1  
2  
3  
4 compared to CPU1, the execution time of CPU2 is shorter than CPU1 and the overall speedup depends on  
5 the percentage of the idle time during each full execution step.  
6

7 [insert Fig. 8 here]  
8  
9

## 10 11 12 **4. Numerical experiments**

13 The proposed parallel computing method is implemented on a workstation with an Intel Xeon E5-  
14 2680 2.80GHz multi-core CPU and 192 GB memory. This workstation has a total of 40 threads. To  
15 ensure the normal running of the workstation, we only use up to 32 threads for the simulation program in  
16 the experiments and reserve the remaining 8 threads for critical operating systems and other system-level  
17 threads. The source codes were implemented using Visual Studio C++ 2013 and its built-in OpenMP API  
18 library. Based on the proposed synchronous LP spatial-parallel computing method, our research group  
19 developed an open-source traffic simulation software [DTALite \(Zhou and Taylor, 2014; Zhou et al., 2015;](#)  
20 [Ruan et al., 2016\)](#) that was used to measure the total execution time and speedup under different numbers  
21 of threads.  
22  
23  
24  
25  
26  
27  
28

### 29 **4.1. Network and experiment design**

30 A set of experiments were performed on four different scales of networks. These include West  
31 Jordan network, Clackamas, Salt Lake City, and Maryland statewide network, with increasing scale  
32 respectively. [Table 1](#) lists the supply and demand summaries of all the four networks. In scenario (a), the  
33 West Jordan area extracted from Salt Lake City (UDOT Report No. UT-15.09 by Zhou et al., 2015) is a  
34 small-scale network with 0.1K nodes and 0.4K links and the traffic demand loading period is 2 hours  
35 (15:30-17:30). In scenario (b), the Clackamas subarea network extracted from the Portland metropolitan  
36 network ([Nevers et al., 2013](#)) is a medium-scale network with 1.5K nodes and 4.1K links and the traffic  
37 demand loading period is 5 hours (14:00-19:00). In scenario (c), the Salt Lake City network is a large-  
38 scale regional network with 13.9K nodes and 26.8K links and the demand loading period is 5 hours. In  
39 scenario (d), the Maryland statewide network ([Erdoğan et al., 2015](#)) is a large-scale metropolitan/mega  
40 region network, which includes all of Maryland, Delaware and the District of Columbia, along with  
41 adjacent portions of Virginia, Pennsylvania and West Virginia. This network has a total of 26.6 M  
42 travelers during a 24-hour simulation horizon. [Fig. 9](#) plots the corresponding networks and approximate  
43 peak hour demand values.  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

54 [insert Table 1 here]  
55  
56

57 [insert Fig. 6 here]  
58  
59  
60  
61  
62  
63  
64  
65

## 4.2. Simulation results

The dynamic traffic assignment program includes both an assignment/routing stage and simulation stage. The dynamic network loading stage is able to produce time dependent link volume, speed, density, queue length, and has the ability to track trajectories of individual vehicles on each link. Focusing on the execution time of the simulation stage, each scenario is executed 10 iterations to reduce sampling errors for getting a tally of average execution time of dynamic network loading at each iteration. Taking scenario b) for example, the cumulative execution time of each iteration is listed in Table 2. The execution time of single-thread computing is 291.825 s. As the number of threads increases, the CPU usage is improved with shorter execution time, for example, with a speedup of about 9.30 in a 32-thread computing environment.

[insert Table 2 here]

Recall that, a speedup is measured by comparing the execution time of a simulation using multiple threads and one obtained using one thread. The speedups of the aforementioned four test cases are plotted in Fig. 10, which indicate a possible nonlinear speedup as a function of the number of threads used in computing. Overall, both the scale of network (in terms of number of links and nodes) and travel demand (in terms of number of agents) could affect the overall parallel computing efficiency.

[insert Fig. 10 here]

The parallel computing speedup values vary from 4.1 to 10.7 in simulating the four real-world networks. The speedup is higher for larger road traffic networks, because a large number of links and nodes allows a more balanced computing resource allocation to each processor, or more precisely, link-based or node-based LPs. In contrast, for smaller networks, the events have to be distributed to only a few of the processors, and therefore, the idle time of unused processors could be significant. The higher demand or a large number network could introduce a more balanced computing task assignment across nodes and links (at different CPU cores), which leads to a potentially higher speed up rate. The average execution time of scheduling events could dramatically increase with more frequent interactions and data exchanges between LPs. As a result, the overall communication overhead consumes more computation resources, leading to a slower increasing trend of the speedup curve. Overall, we have demonstrated that the speed-up in the parallel computing environment is significant regardless of the number of CPUs available in our example.

## 5. Conclusions

This paper has presented a parallel computing implementation approach for vehicular traffic simulation in real-world networks. According to the properties of individual events, a space-time-event

1  
2  
3  
4 view is proposed to allocate agent’s entering and leaving events to different link buffers. Based on the  
5 double-buffer representation, we decompose the general mesoscopic transportation simulation into five  
6 sequential tasks that can be concurrently executed on link-based and node-based logical processes. The  
7 logical processes are then dynamically distributed to different processors by the proposed synchronous  
8 space-parallel LP strategy. As a unique contribution, this paper establishes a space-time-processor  
9 network to illustrate the spatial and temporal LP scheduling and synchronization mechanism for  
10 microscopic transportation system simulation. The proposed parallel computing method is then  
11 implemented and applied to simulate four real-world transportation networks. The simulation results  
12 show significant speedup using our method and that that speedup magnitude is dependent on the scale of  
13 the network and the travel demand. In the examples provided, our method attained a speedup of up to  
14 10.7 on a workstation with 32 computing CPU threads.  
15  
16  
17  
18  
19  
20  
21  
22

23 In our future work, we will perform further testing of the proposed parallel computing algorithm  
24 using a more efficient communication protocol for inter-process communication. We will further evaluate  
25 a wide range of other parallelization techniques, such as distributed shared memory-based methods and  
26 GPU computing technology, to further improve the computational efficiency of large-scale mesoscopic  
27 traffic simulations. The proposed space-time parallel framework is computationally efficient in traffic  
28 simulation, and we will also plan to extend this parallel computing method to other large and complex  
29 multi-agent simulation systems, e.g., pedestrian simulation (Qu et al., 2014; Qu et al, 2015) and urban rail  
30 transit networks, by considering boundedly rational route choice behavior (Liu and Zhou, 2016).  
31  
32  
33  
34  
35  
36  
37  
38

### 39 **Acknowledgement**

40 This paper is mainly supported by National Science Foundation – United States under Grant No.  
41 CMMI 1538105 “Collaborative Research: Improving Spatial Observability of Dynamic Traffic Systems  
42 through Active Mobile Sensor Networks and Crowdsourced Data”, and U.S. Department of Energy’s  
43 (DOE) Advanced Research Projects Agency – Energy (ARPA-E), “Traveler Response Architecture using  
44 Advanced Novel Signaling for Network Efficiency in Transportation (TRANSNET)”. The second and  
45 corresponding author especially thanks Mr. Brian J Gardner at Federal Highway Administration’s Office  
46 of Planning, for his encouragement in our development process of the open-source traffic simulation  
47 package DTALite. We also thank Jiangtao Liu and Dr. Taylor Li at Arizona State University, and Dr.  
48 Jinjin Tang at Beijing Jiaotong University, Jeff Taylor at University of Utah and Dr. Jeffrey Stempihar at  
49 Arizona State University for their useful comments and suggestions. The work presented in this paper  
50 remains the sole responsibility of the authors.  
51  
52  
53  
54

### 55 **References**

56  
57 Adler, J. L., & Blue, V. J. (2002). A cooperative multi-agent transportation management and route  
58 guidance system. *Transportation Research Part C: Emerging Technologies*, 10(5), 433-454.  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4 Auld, J., Hope, M., Ley, H., Sokolov, V., Xu, B., & Zhang, K. (2016). POLARIS: Agent-based modeling  
5 framework development and implementation for integrated travel demand and network and  
6 operations simulations. *Transportation Research Part C: Emerging Technologies*, 64, 101-116.  
7 Barceló, J, Ferrer, J.L, García D., Florian, M. and E. Le Saux, The Distributedization of AIMSUN2  
8 Microscopic Simulator for ITS Applications, Proc. 3rd. World Congress on Intelligent Transport  
9 Systems, Orlando, 1996.  
10 Ben-Akiva, M. E., Bierlaire, M., Burton, D., Koutsopoulos, H. N., & Mishalani, R. (2002). Network state  
11 estimation and prediction for real-time transportation management applications. Paper presented  
12 at the Transportation Research Board 81st Annual Meeting.  
13 Bliemer, M.C.J., H.H. Versteegt and R.J. Castenmiller (2004) INDY: A New Analytical Multiclass  
14 Dynamic Traffic Assignment Model, Proceedings of the TRISTAN V conference, Guadeloupe.  
15 Bryant, R. E.. Simulation of packet communication architecture computer systems. Technical Report  
16 MIT-LCS-TR-188, MIT, 1977.  
17 Cameron, G. D., & Duncan, G. I. (1996). PARAMICS—Parallel microscopic simulation of road traffic.  
18 *The Journal of Supercomputing*, 10(1), 25-53. K. Nagel and M. Schmidt. Parallel DYNEMO:  
19 Mesoscopic traffic flow simulation on large networks. preprint, 2000.  
20 Celikoglu, H. B., & Dell’Orco, M. (2007). Mesoscopic simulation of a dynamic link loading process.  
21 *Transportation Research Part C: Emerging Technologies*, 15(5), 329-344.  
22 Cetin N., Burri A., and Nagel K. 2003. A large-scale agent-based traffic microsimulation based on queue  
23 model. Swiss Transport Research Conference, Monte Verita, CH.  
24 Chandy, K. M., & Misra, J. (1979). Distributed simulation: A case study in design and verification of  
25 distributed programs. *IEEE Transactions on software engineering*, (5), 440-452.  
26 Chandy, K. M., Sherman, R. Space-time and simulation. University of Southern California, Information  
27 Sciences Institute, 1989.  
28 Chen, B., Cheng, H. H., & Palen, J. (2009). Integrating mobile agent technology with multi-agent systems  
29 for distributed traffic detection and management systems. *Transportation Research Part C:*  
30 *Emerging Technologies*, 17(1), 1-10.  
31 Comfort, J. C. The simulation of a master-slave event set processor. *Simulation*, 1984, 42(3): 117-124.  
32 Daganzo, C. F. The cell transmission model: A dynamic representation of highway traffic consistent with  
33 the hydrodynamic theory. *Transportation Research Part B: Methodological*, 1994, 28(4): 269-287.  
34 Daganzo, C. F. (1995a). A finite difference approximation of the kinematic wave model of traffic flow.  
35 *Transportation Research Part B: Methodological*, 29(4): 261-276.  
36 Daganzo, C. F. (1995b). The cell transmission model, part II: network traffic. *Transportation Research*  
37 *Part B: Methodological*, 29(2), 79-93.  
38 Daganzo, C. F. In traffic flow, cellular automata= kinematic waves. *Transportation Research Part B:*  
39 *Methodological*, 2006, 40(5): 396-403.  
40 Dagum, L., Enon, R. (1998). OpenMP: an industry standard API for shared-memory programming.  
41 *Computational Science & Engineering, IEEE*, 5(1), 46-55.  
42 Dell’Orco, M., Marinelli, M., & Silgu, M. A. (2016). Bee Colony Optimization for innovative travel time  
43 estimation, based on a mesoscopic traffic assignment model. *Transportation Research Part C:*  
44 *Emerging Technologies*, 66, 48-60.  
45 Di Gangi, M., Cantarella, G. E., Di Pace, R., & Memoli, S. (2016). Network traffic control based on a  
46 mesoscopic dynamic flow model. *Transportation Research Part C: Emerging Technologies*, 66, 3-  
47 26.  
48 Erdoğan, S., Zhou, X., & Liu, J. (2015). A Simplified Dynamic Traffic Assignment Framework for  
49 Statewide Traffic Modeling. In *Transportation Research Board 94th Annual Meeting (No. 15-6090)*.  
50 Ferscha, A., Tripathi, S. K. Parallel and distributed simulation of discrete event systems. 1998.  
51 Florian, M. and Gendreau, M., (2001). Applications of parallel computing in transportation. *Parallel*  
52 *Computing*, 27(12), 1521-1522.  
53 Fujimoto, R. M. Parallel discrete event simulation. *Communications of the ACM*, 1990, 33(10): 30-53.  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4 Fujimoto, R. M. Parallel discrete event simulation: Will the field survive?. *ORSA Journal on Computing*,  
5 1993, 5(3): 213-230.
- 6 Fujimoto, R. (2015). Parallel and distributed simulation. In *Proceedings of the 2015 Winter Simulation*  
7 *Conference* (pp. 45-59). IEEE Press.
- 8 Hunter, M., Kim, H. K., Suh, W., Fujimoto, R., Sirichoke, J., & Palekar, M. (2009). Ad hoc distributed  
9 dynamic data-driven simulations of surface transportation systems. *Simulation*, 85(4), 243-255.
- 10 Jafer, S., Liu, Q., & Wainer, G. (2013). Synchronization methods in parallel and distributed discrete-event  
11 simulation. *Simulation Modelling Practice and Theory*, 30, 54-73.
- 12 Jefferson, D. R.. Virtual time. *ACM Transactions on Programming Languages and Systems*, 7(3):404–  
13 245, July 1985.
- 14 Junchaya, T., & Chang, G. L. (1993). Exploring real-time traffic simulation with massively parallel  
15 computing architecture. *Transportation Research Part C: Emerging Technologies*, 1(1), 57-76.
- 16 Kallioras, N. A., Kepaptsoglou, K., & Lagaros, N. D. (2015). Transit stop inspection and maintenance  
17 scheduling: A GPU accelerated metaheuristics approach. *Transportation Research Part C:*  
18 *Emerging Technologies*, 55, 246-260.
- 19 Kirchner, A., & Schadschneider, A. (2002). Simulation of evacuation processes using a bionics-inspired  
20 cellular automaton model for pedestrian dynamics. *Physica A: Statistical Mechanics and its*  
21 *Applications*, 312(1), 260-276.
- 22 Lee, D. H., & Chandrasekar, P. (2002). A framework for parallel traffic simulation using multiple  
23 instancing of a simulation program. *ITS Journal*, 7(3-4), 279-294.
- 24 Liu, J. *Parallel Discrete - Event Simulation*[M]. John Wiley & Sons, Inc., 2009.
- 25 Liu, J. & Zhou, Z. (2016) Capacitated transit service network design with boundedly rational agents.  
26 *Transportation Research Part B: Methodological*, 93, 225-250
- 27 Liu, Z., Meng, Q. (2013). Distributed computing approaches for large - scale probit - based stochastic  
28 user equilibrium problems. *Journal of Advanced Transportation*, 47(6), 553-571.
- 29 Lu, C.C., Mahmassani, H.S., Zhou, X., (2009). Equivalent gap function-based reformulation and solution  
30 algorithm for the dynamic user equilibrium problem. *Transportation Research Part B:*  
31 *Methodological*, 43(3), 345-364.
- 32 Lu, C. C., Zhou, X., & Zhang, K. (2013). Dynamic origin–destination demand flow estimation under  
33 congested traffic conditions. *Transportation Research Part C: Emerging Technologies*, 34, 16-37.
- 34 Mahmassani, H. S., Hu, T-Y, Peeta, S., and Ziliaskopoulos, A. Development and testing of dynamic  
35 traffic assignment and simulation procedures for ATIS/ATMS applications. Report DTFH61-90-  
36 R-00074-FG, U.S. DOT, Federal Highway Administration, McLean, Virginia, 1994.
- 37 Mahmassani, H. S. Dynamic network traffic assignment and simulation methodology for advanced  
38 system management application. *Networks and Spatial Economics*, 2001, Vol. 1, 267-292.
- 39 Morosan, C. D., Florian, M. (2015). The benefits of parallel computing for large-scale network  
40 equilibrium models. Downloaded from [https://www.inrosoftware.com/assets/pres-](https://www.inrosoftware.com/assets/pres-pap/TRB2015/P15-6591.pdf)  
41 [pap/TRB2015/P15-6591.pdf](https://www.inrosoftware.com/assets/pres-pap/TRB2015/P15-6591.pdf).
- 42 Nagel, K., & Schreckenberg, M. (1992). A cellular automaton model for freeway traffic. *Journal de*  
43 *physique I*, 2(12), 2221-2229.
- 44 Nagel, K., & Rickert, M. (2001). Parallel implementation of the TRANSIMS micro-simulation. *Parallel*  
45 *Computing*, 27(12), 1611-1639.
- 46 Nevers, Brandon L., et al. *The Effective Integration of Analysis, Modeling, and Simulation Tools*. No.  
47 FHWA-HRT-13-036. 2013.
- 48 Newell, G. F. (1993). A simplified theory of kinematic waves in highway traffic, part I: general theory.  
49 *Transportation Research Part B: Methodological*, 27(4), 281-287
- 50 Ng, M. W., Duc. T. Nguyen. (2015). Domain Decomposition, Parallel Computing and Traffic  
51 Assignment. Downloaded from [http://amonline.trb.org/trb57535-2015-1.1793793/t006-](http://amonline.trb.org/trb57535-2015-1.1793793/t006-1.1818822/116-1.1809928/p15-6590-1.1820620/p15-6590-1.1955721?qr=1)  
52 [1.1818822/116-1.1809928/p15-6590-1.1820620/p15-6590-1.1955721?qr=1](http://amonline.trb.org/trb57535-2015-1.1793793/t006-1.1818822/116-1.1809928/p15-6590-1.1820620/p15-6590-1.1955721?qr=1)
- 53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4 Nie, Y., Ma, J., & Zhang, H. M. (2008). A polymorphic dynamic network loading model. *Computer -*  
5 *Aided Civil and Infrastructure Engineering*, 23(2), 86-103.
- 6 Panneton, F., L'ecuyer, P., & Matsumoto, M. (2006). Improved long-period generators based on linear  
7 recurrences modulo 2. *ACM Transactions on Mathematical Software (TOMS)*, 32(1), 1-16.
- 8 Peacock, J. K., J. W. Wong, and E. G. Manning. *Distributed Simulation using a Network of Processors.*  
9 *Computer Networks*, Vol. 3, No. 1, pp. 44, 1979.
- 10 Peeta, S., & Ziliaskopoulos, A. K. (2001). Foundations of dynamic traffic assignment: The past, the  
11 present and the future. *Networks and Spatial Economics*, 1(3-4), 233-265.
- 12 Potuzak, T. *Distributed-Parallel Road Traffic Simulator for Clusters of Multi-core Computers.* in  
13 *Distributed Simulation and Real Time Applications (DS-RT)*, 2012 IEEE/ACM 16th  
14 International Symposium on. 2012.
- 15 Qu, Y., Gao, Z., Xiao, Y., Li, X., (2014). Modeling the pedestrian's movement and simulating evacuation  
16 dynamics on stairs. *Safety science*, 70, 189-201.
- 17 Qu, Y., Gao, Z., Orenstein, P., Long, J. and Li, X., (2015). An effective algorithm to simulate pedestrian  
18 flow using the heuristic force-based model. *Transportmetrica B: transport dynamics*, 3(1), 1-26.
- 19 Ruan, J. M., Liu, B., Wei, H., Qu, Y., Zhu, N., & Zhou, X. (2016). How Many and Where to Locate  
20 Parking Lots? A Space-time Accessibility-Maximization Modeling Framework for Special Event  
21 Traffic Management. *Urban Rail Transit*, 1-12.
- 22 Snelder, M. (2009). A comparison between dynameq and indy. CIRRELT.
- 23 Sundaram, S., Koutsopoulos, H.N., Ben-Akiva, M., Antoniou, C., Balakrishna, R., 2011. Simulation-  
24 based dynamic traffic assignment for short-term planning applications. *Simulation Modelling*  
25 *Practice and Theory* 19(1), 450-462.
- 26 Wong, S. C. (1997). Group-based optimisation of signal timings using parallel computing. *Transportation*  
27 *Research Part C: Emerging Technologies*,5(2), 123-139.
- 28 Yperman, I. *The link transmission model for dynamic network loading.* 2007.
- 29 Ziliaskopoulos, A., Kotzinos, D., & Mahmassani, H. S. (1997). Design and implementation of parallel  
30 time-dependent least time path algorithms for intelligent transportation systems applications.  
31 *Transportation Research Part C: Emerging Technologies*, 5(2), 95-107.
- 32 Zhen, S., W. Kai, and Z. Fenghua. Agent-based traffic simulation and traffic signal timing optimization  
33 with GPU. in *Intelligent Transportation Systems (ITSC)*, 2011 14th International IEEE  
34 Conference on. 2011.
- 35 Zhou, X., Taylor, J. (2014). DTALite: A queue-based mesoscopic traffic simulator for fast model  
36 evaluation and calibration. *Cogent Engineering*, 1(1): 961345.
- 37 Zhou, X., Zlatkovic, M., Farhan, M. (2015). Simplified web-based decision support method for traffic  
38 management and work zone analysis. Utah Department of Transportation, Report No. UT-15.09.
- 39 Zhou, X., Tanvir, S., Lei, H., Taylor, J., Liu, B., Roupail, N. M., & Frey, H. C. (2015). Integrating a  
40 simplified emission estimation model and mesoscopic dynamic traffic simulator to efficiently  
41 evaluate emission impacts of traffic management strategies. *Transportation Research Part D:*  
42 *Transport and Environment*, 37, 123-136.
- 43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6 **List of Figure**  
7  
8  
9

10  
11 **Fig. 1 Network representation and logical process in parallel discrete event simulation**

12 **Fig. 2 Overall framework of mesoscopic transportation simulation**

13  
14 **Fig. 3 Space-time trajectory and vehicular event list**

15  
16 **Fig. 4 Parallelism mechanism in OpenMP with multiple CPU threads**

17  
18 **Fig. 5 Different LP distributions at different parallel computing instances, a CPU label means a**  
19 **CPU thread**  
20

21 **Fig. 6 Space-time-processor scheduling graph**

22 **Fig. 7 Processor-space-time network representation**

23  
24 **Fig. 8 CPU usage time series/Gantt chart for time-dependent discretized event simulation using**  
25 **OpenMP**  
26

27 **Fig. 9 Overall maps of four real-world transportation networks**

28  
29 **Fig. 10 Speedup of four networks**  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

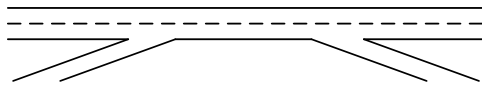
**List of Table**

**Tab. 1 Details of four real networks used for simulations**

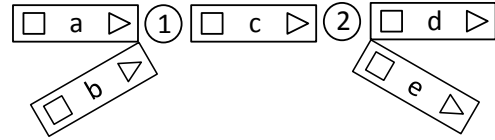
**Tab. 2 Cumulative execution time of each iteration (Scenario b)**



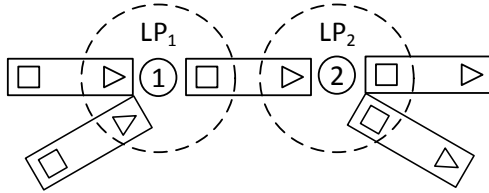
FIGURES



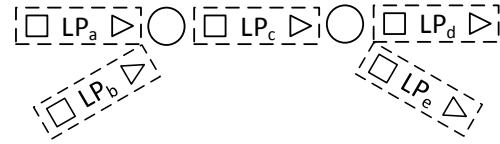
a) freeway corridor



b) double-buffer representation



c) node-based LP



d) link-based LP

Fig. 1 Network representation and logical process in parallel discrete event simulation

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

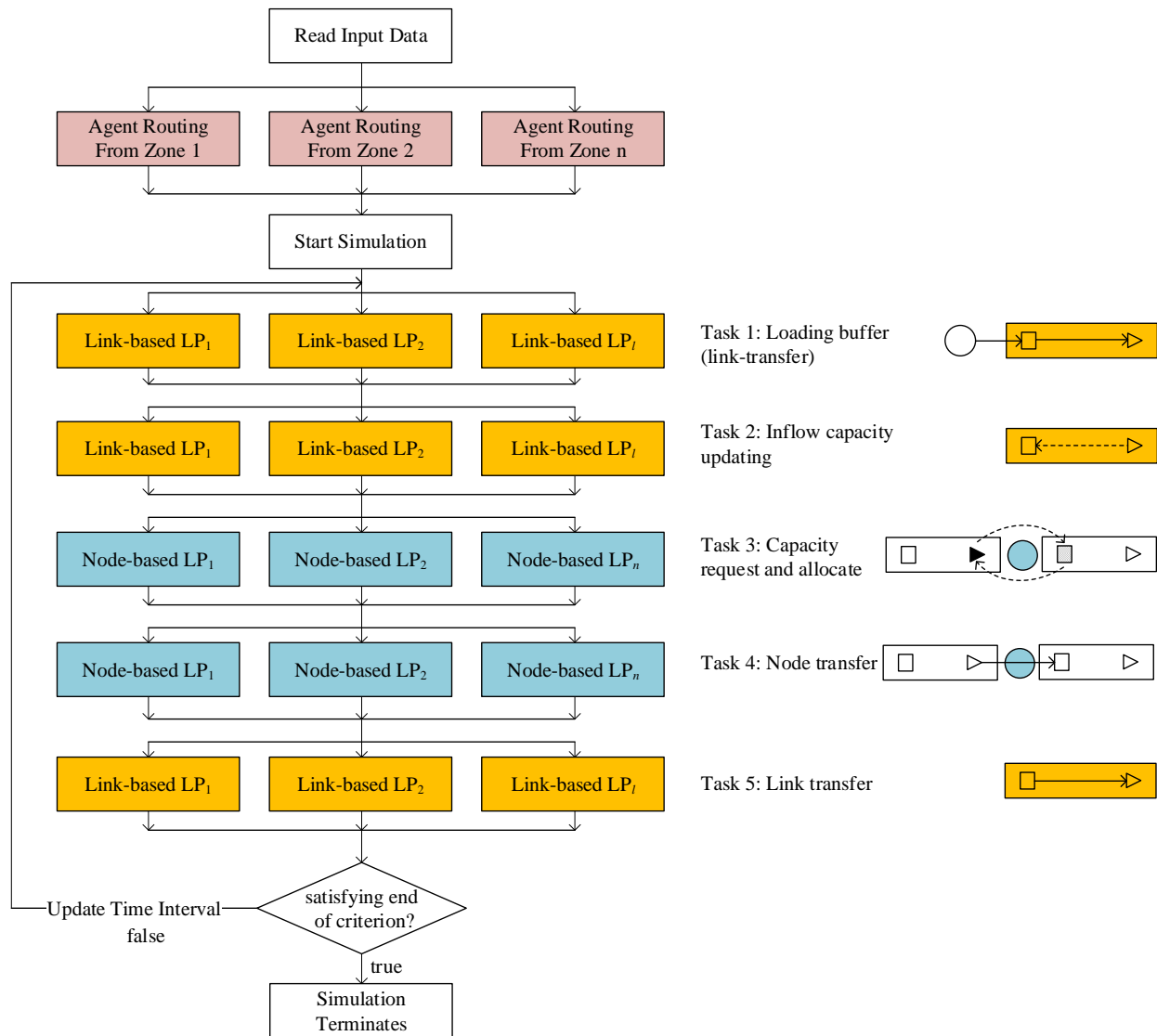


Fig. 2 Overall framework of mesoscopic transportation simulation

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

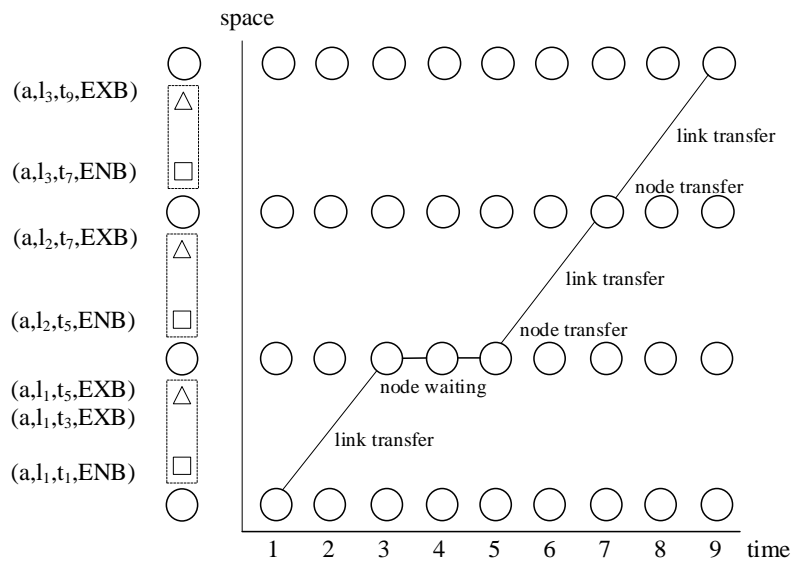


Fig. 3 Space-time trajectory and vehicular event list

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

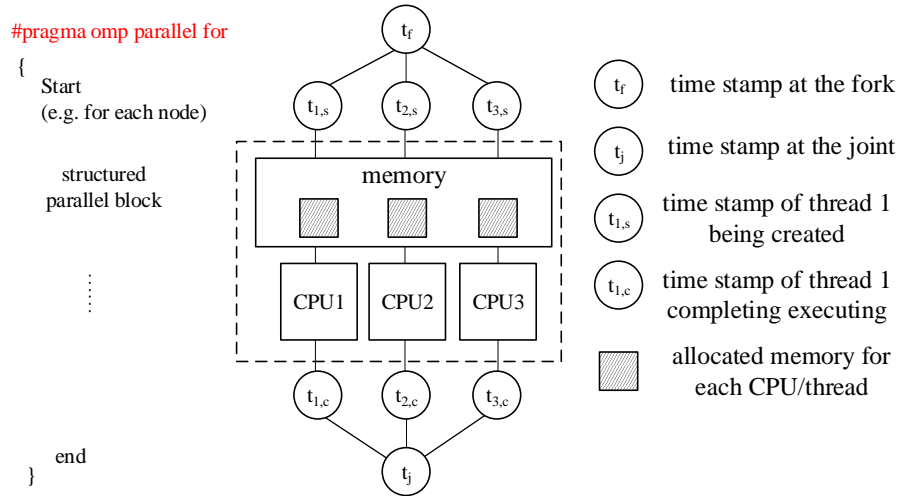


Fig. 4 Parallelism mechanism in OpenMP with multiple CPU threads

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

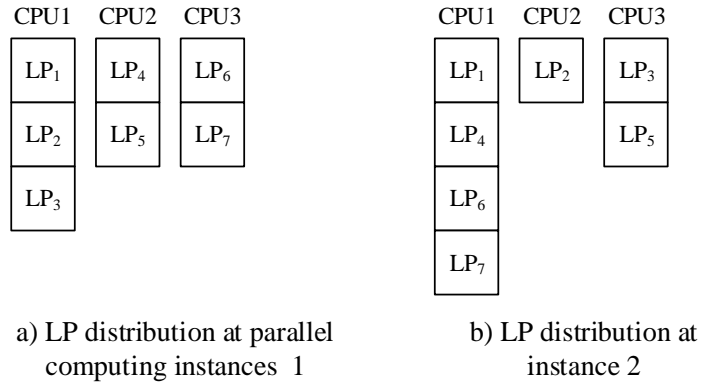


Fig. 5 Different LP distributions at different parallel computing instances, a CPU label means a CPU thread

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

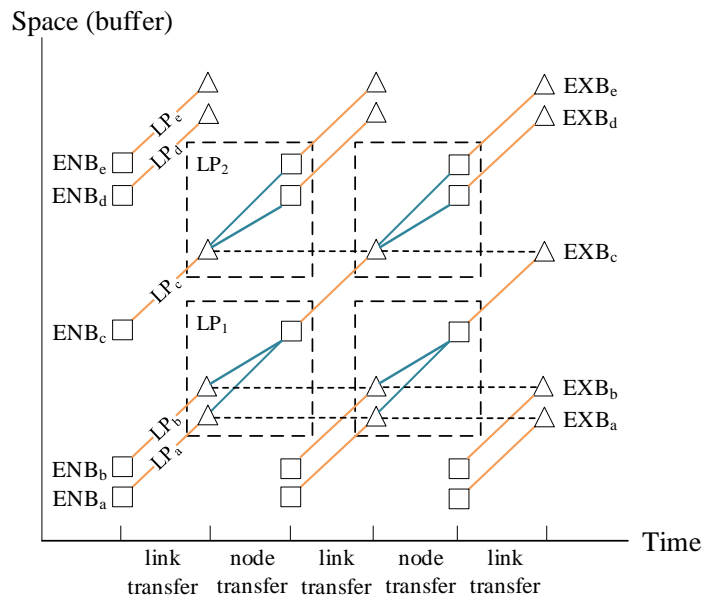
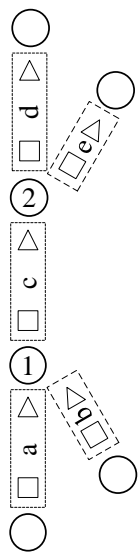


Fig. 6 Space-time-processor scheduling graph

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

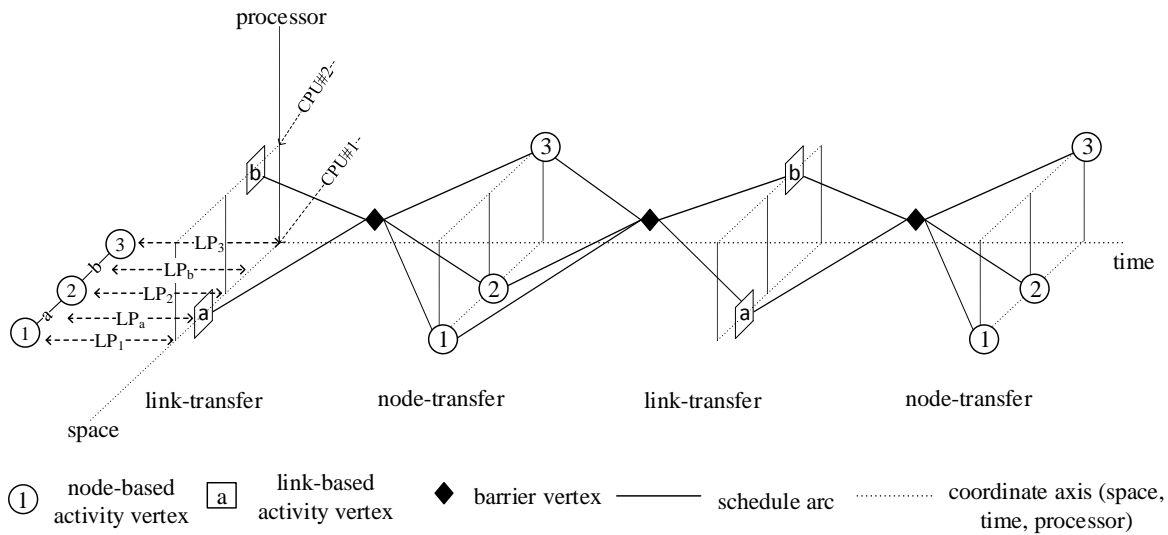


Fig. 7 Processor-space-time network representation

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

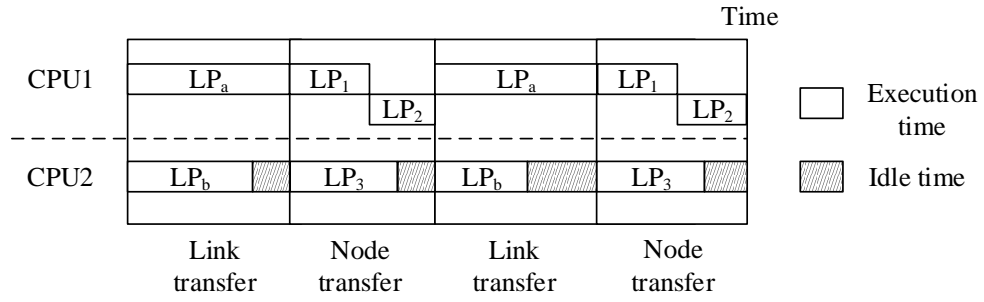


Fig. 8 CPU usage time series/Gantt chart for time-dependent discretized event simulation using OpenMP



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

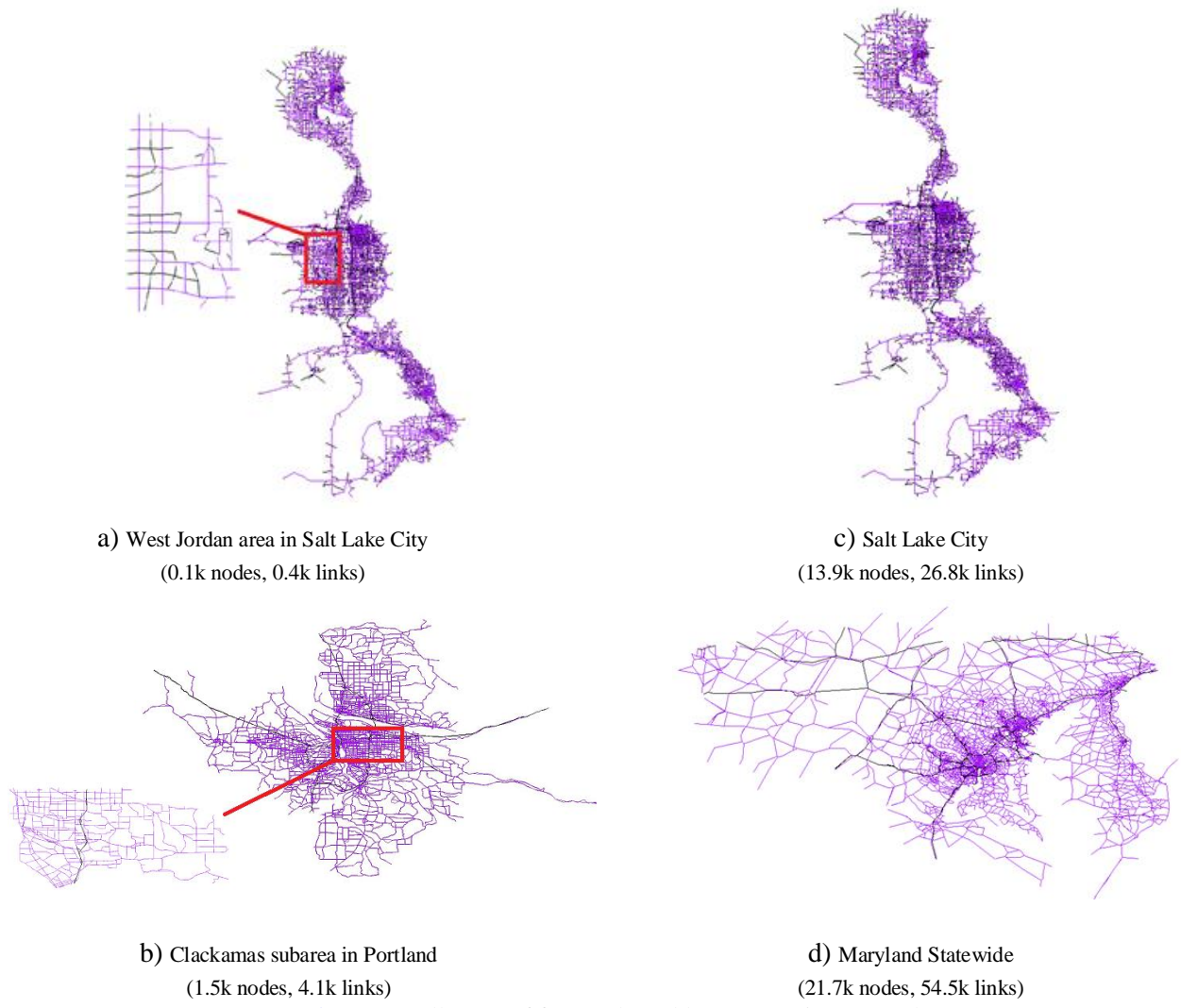


Fig. 9 Overall maps of four real-world transportation networks

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

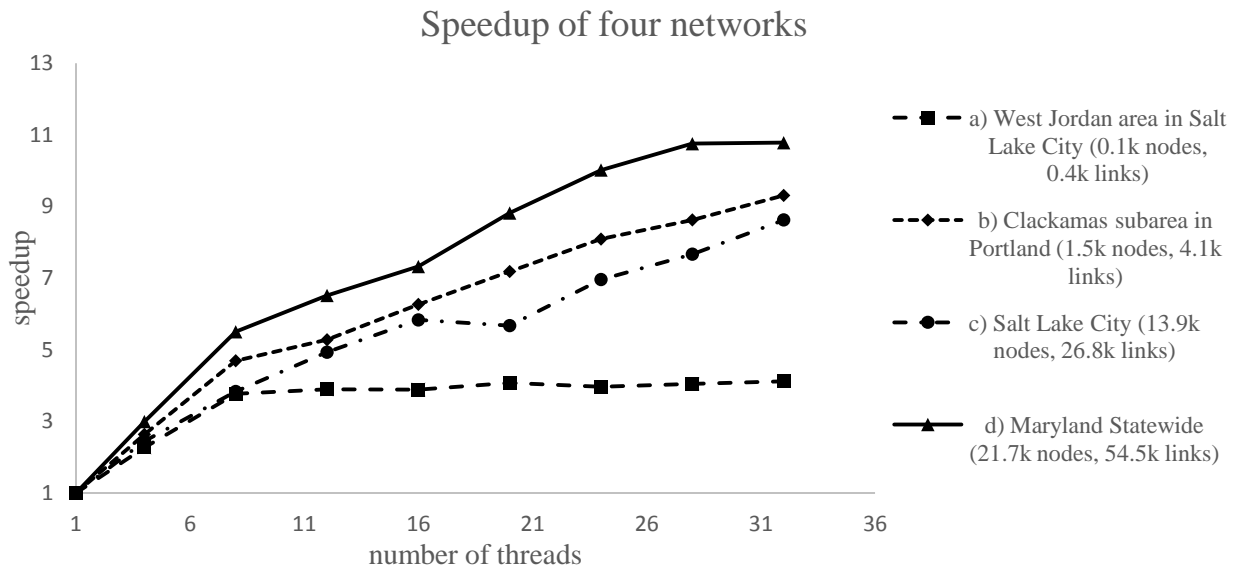


Fig. 10 Speedup of four networks

1  
2  
3  
4 **TABLES**  
5  
6

7 **Tab. 1** Details of four real networks used for simulations

8  
9

Scenario	Network	Number of nodes	Number of links	Number of agents	Demand loading period	Scale
a	West Jordan subarea	0.1K	0.4K	25.3K	2h	Small
b	Clackamas subarea	1.5K	4.1K	491.5K	5h	Medium
c	Salt Lake City region	13.9K	26.8K	1.3M	5h	Large
d	Maryland statewide	21.7K	54.5K	26.6M	24h	Large

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

Tab. 2 Cumulative execution time of each iteration (Scenario b)

Iteration	Number of CPU threads								
	1	4	8	12	16	20	24	28	32
1	28.02	10.529	5.668	4.965	4.668	3.999	3.622	3.217	3.259
2	57.449	21.672	12.319	10.764	9.366	7.996	7.561	6.555	6.205
3	86.273	32.8	18.382	16.571	13.962	12.159	11.06	10.221	9.267
4	115.623	44.054	24.613	22.49	18.698	15.818	14.263	13.534	12.31
5	144.742	55.25	31.438	27.702	23.17	19.811	17.608	16.673	15.499
6	173.806	65.882	37.792	33.519	27.621	24.197	21.058	20.026	18.526
7	203.044	77.119	43.533	38.985	31.93	28.137	24.527	23.69	21.677
8	232.695	87.769	49.712	44.317	36.345	32.282	28.026	27.173	24.646
9	262.601	98.692	56.005	49.489	40.827	36.571	31.577	30.583	27.68
10	291.825	111.61	62.187	55.204	45.432	40.854	35.267	33.559	30.75