

# Building a Large-Scale Microscopic Road Network Traffic Simulator in Apache Spark

Zishan Fu

Jia Yu

Mohamed Sarwat

Arizona State University



# Traffic Congestion



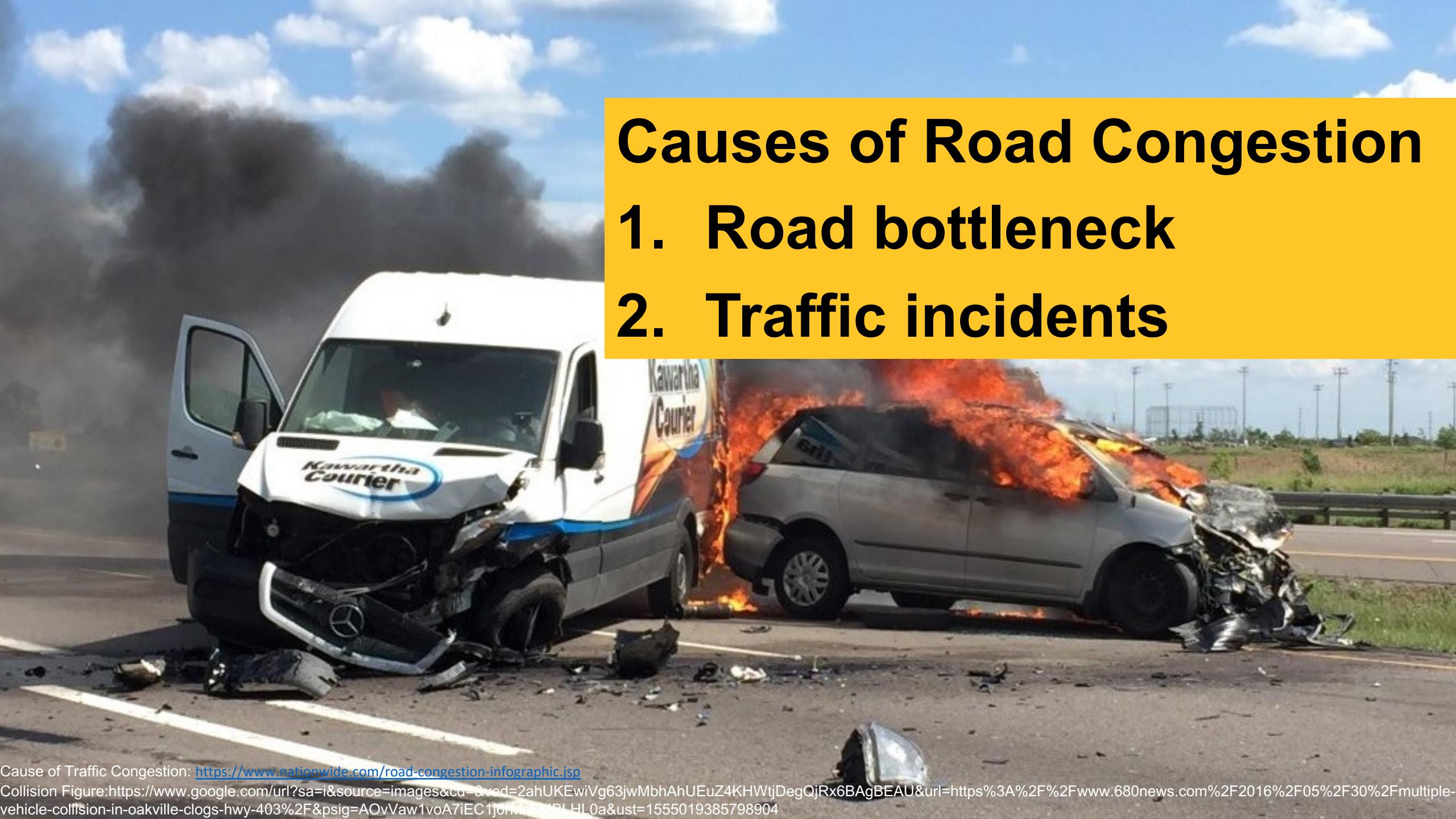
Traffic Figure: <https://www.google.com/search?sa=i&source=images&cd=&ved=2abQKEwjRuLkfwmhbAhWRu54KHaDiCl8QjRx6BAgBEAU&url=https%3A%2F%2Fwww.straitstimes.com%2Fasia%2Fsea-asia%2Fthailand-has-worlds-most-congested-roads-survey&psig=AOvVaw3IU1ztCJ-gcsvD3y5f9rk1&ust=155501924029581>

97

hours were lost in traffic congestion on average

87

billion total cost for traffic in 2018



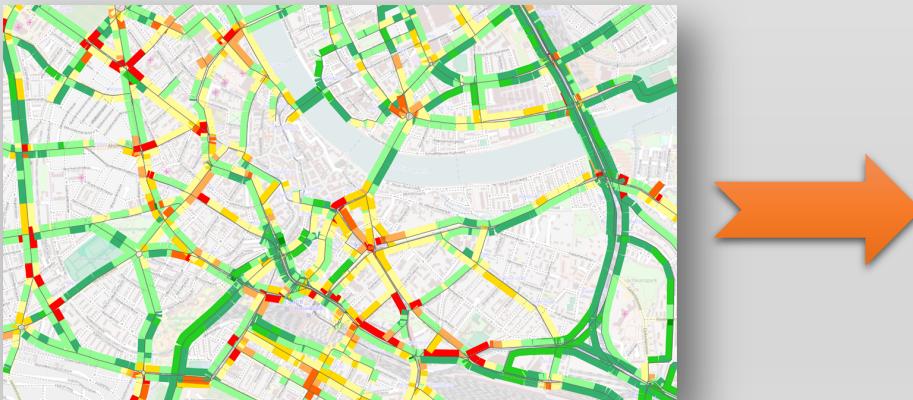
# Causes of Road Congestion

1. Road bottleneck
2. Traffic incidents

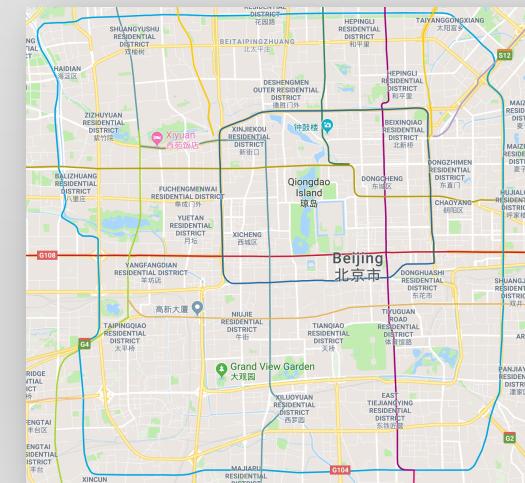
# Traffic data can help

- Urban planning
- Traffic prediction
- Spatial-temporal databases

## Traffic data



## Road network



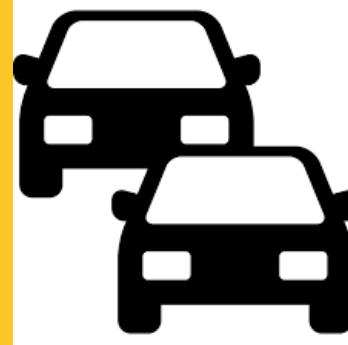
Where is the bottleneck?

What if we add a new road?

What if a road is blocked?

# Collecting real traffic data is NOT EASY

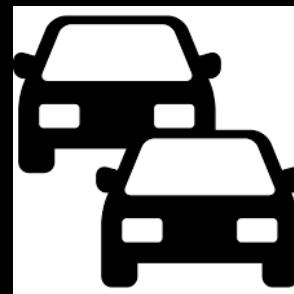
- Problems
  - Small-scale: requires extra monitoring devices / privacy
  - Inaccurate: many trajectory data points deviate from roads
  - Road planning: the road is not built yet
- What we want
  - Large-scale traffic data: various driving behaviors and road situations
  - High-precision: follow on the road network
  - Test planned roads
- Synthetic but simulates realistic driving behaviors and road situations



# GeoSparkSim

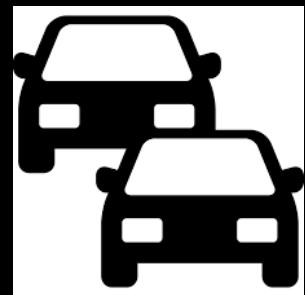
# GeoSparkSim: Large-Scale Microscopic Road Network Traffic Simulator

- Scalable
  - big traffic trajectories
  - big road network
- Microscopic traffic simulation:
  - driving behaviors: speed, lane changing, car following
  - traffic signals
  - traffic rules
- Extensible: customizable driving behaviors
- Traffic visualization



# GeoSparkSim

- Related Works
- Architecture
- Simulation
- Performance
- Demo

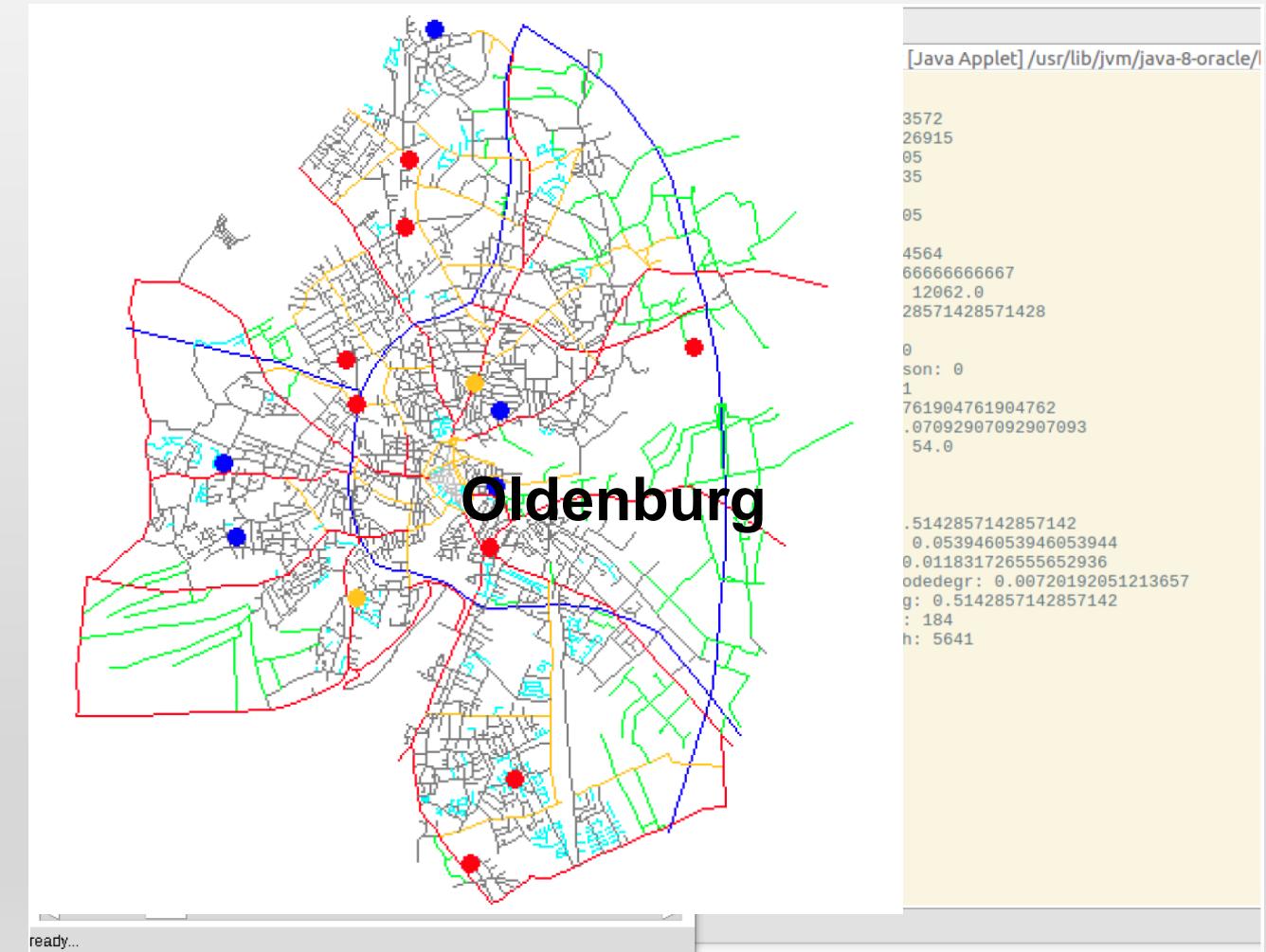


# GeoSparkSim

- **Related Works**
- Architecture
- Simulation
- Performance
- Demo

# Brinkhoff [1] and BerlinMOD [2]

- Macroscopic Traffic Simulator
- Driving Behaviors



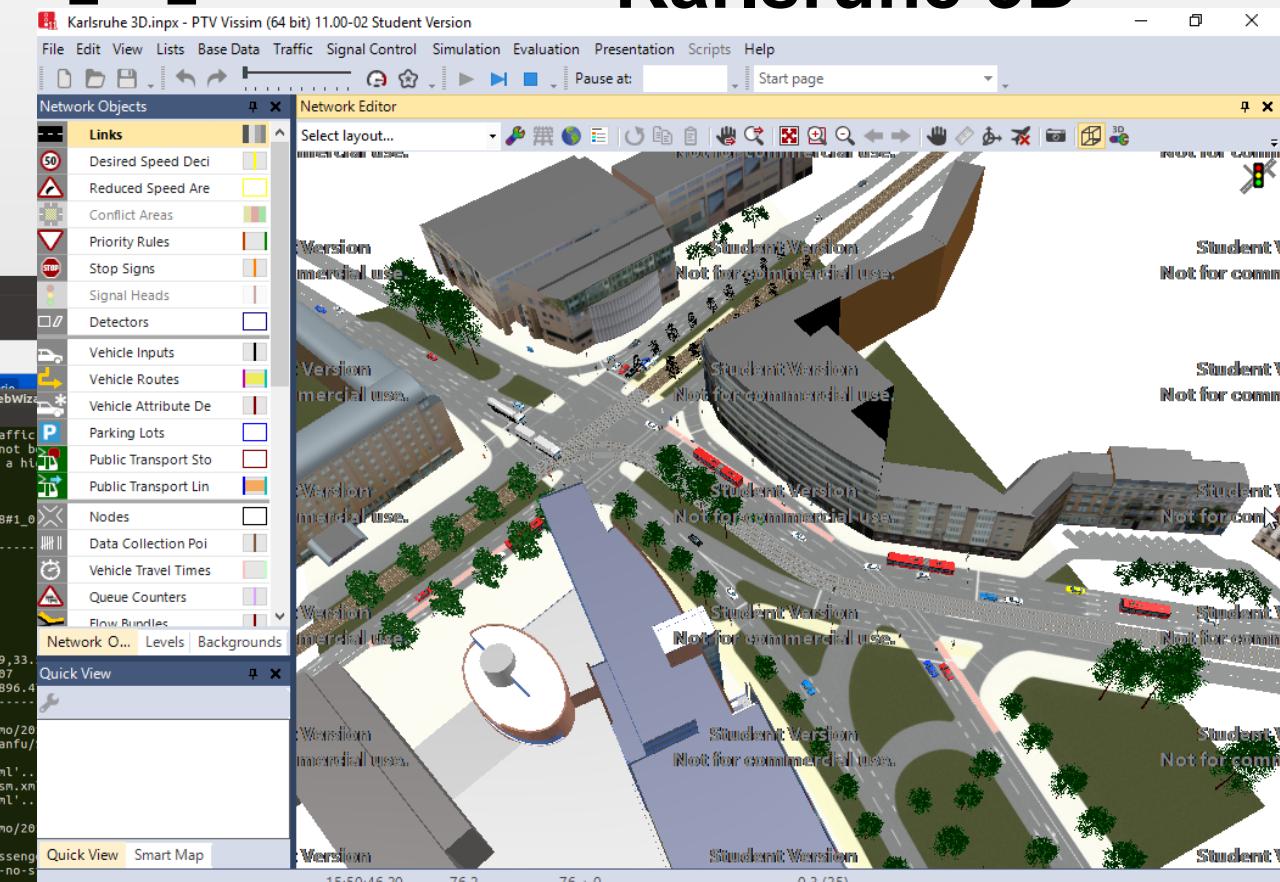
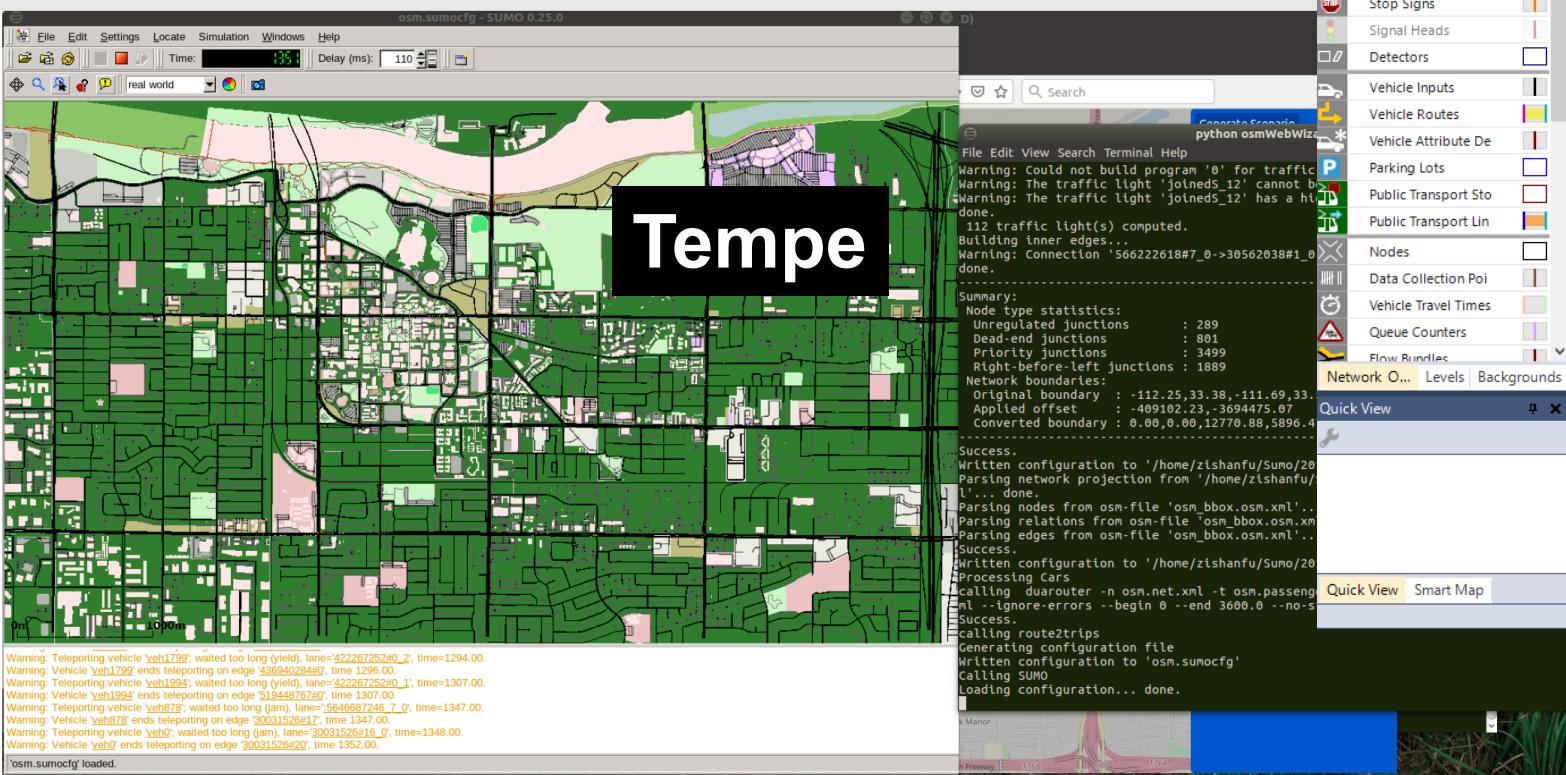
[1] Brinkhoff T. A framework for generating network-based moving objects[J]. Geoinformatica, 2002, 6(2): 153-180.

[2] Düntgen C, Behr T, Güting R H. BerlinMOD: a benchmark for moving object databases[J]. The VLDB Journal—The International Journal on Very Large Data Bases, 2009, 18(6): 1335-1368.

# SUMO[1] and Vissim[2]

Karlsruhe 3D

- Microscopic Traffic Simulator



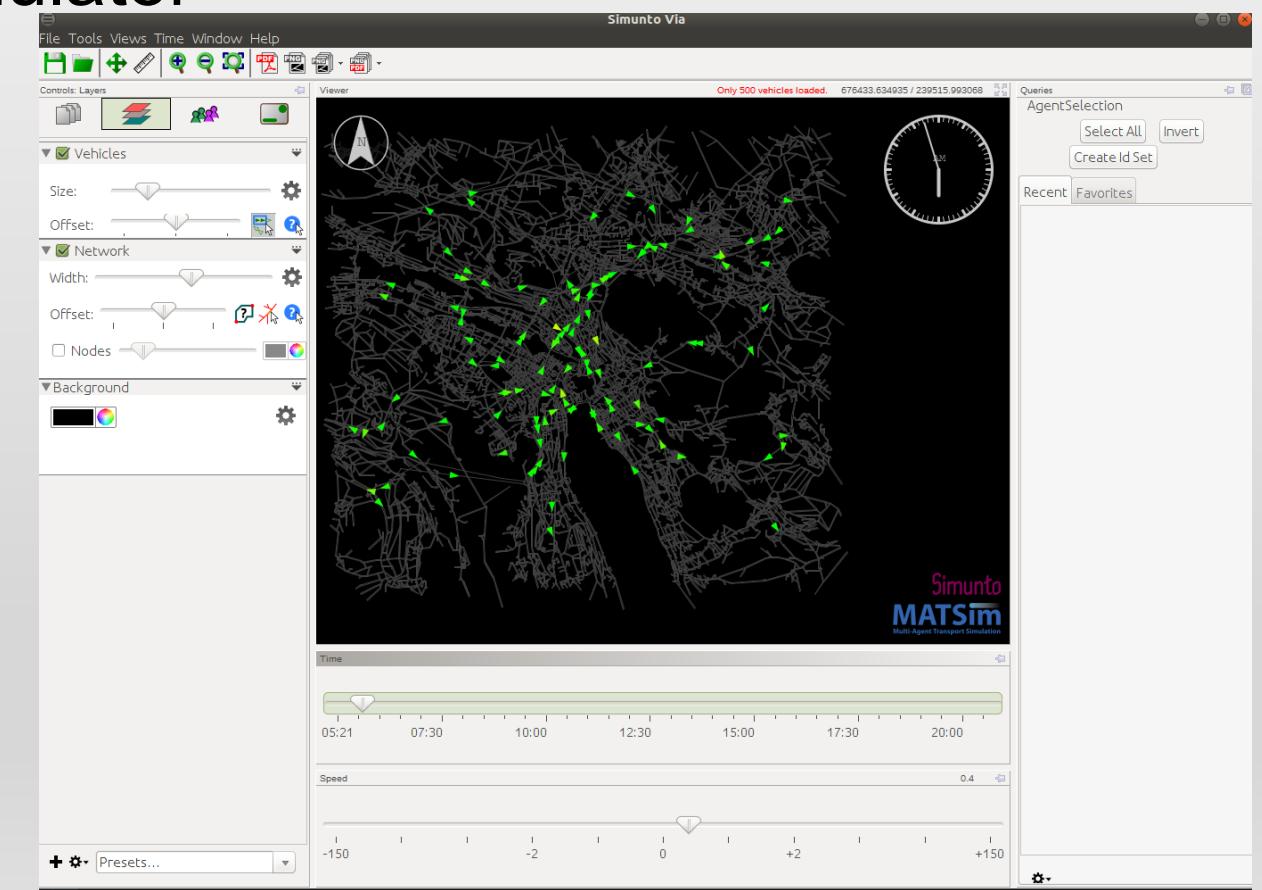
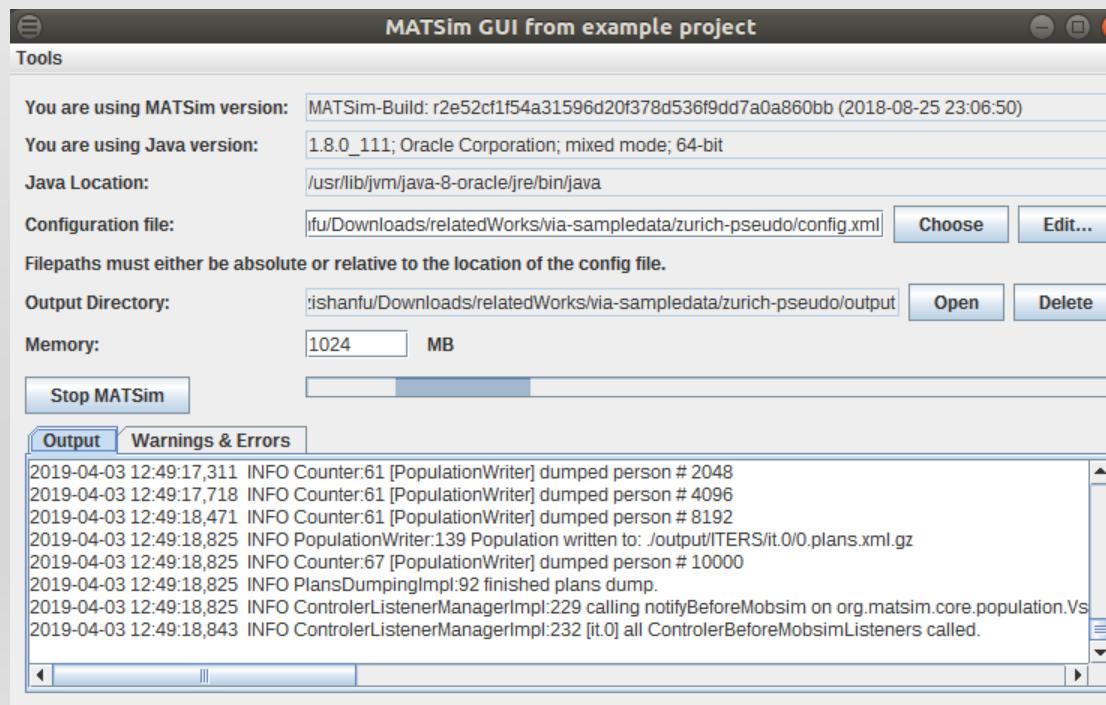
[1] Krajzewicz D, Hertkorn G, Rössel C, et al. SUMO (Simulation of Urban MObility)-an open-source traffic simulation[C]//Proceedings of the 4th middle East Symposium on Simulation and Modelling (MESM2002). 2002: 183-187.

[2] Vissim, <http://vision-traffic.ptvgroup.com/en-us/products/ptv-vissim/>

# MATSim[1]

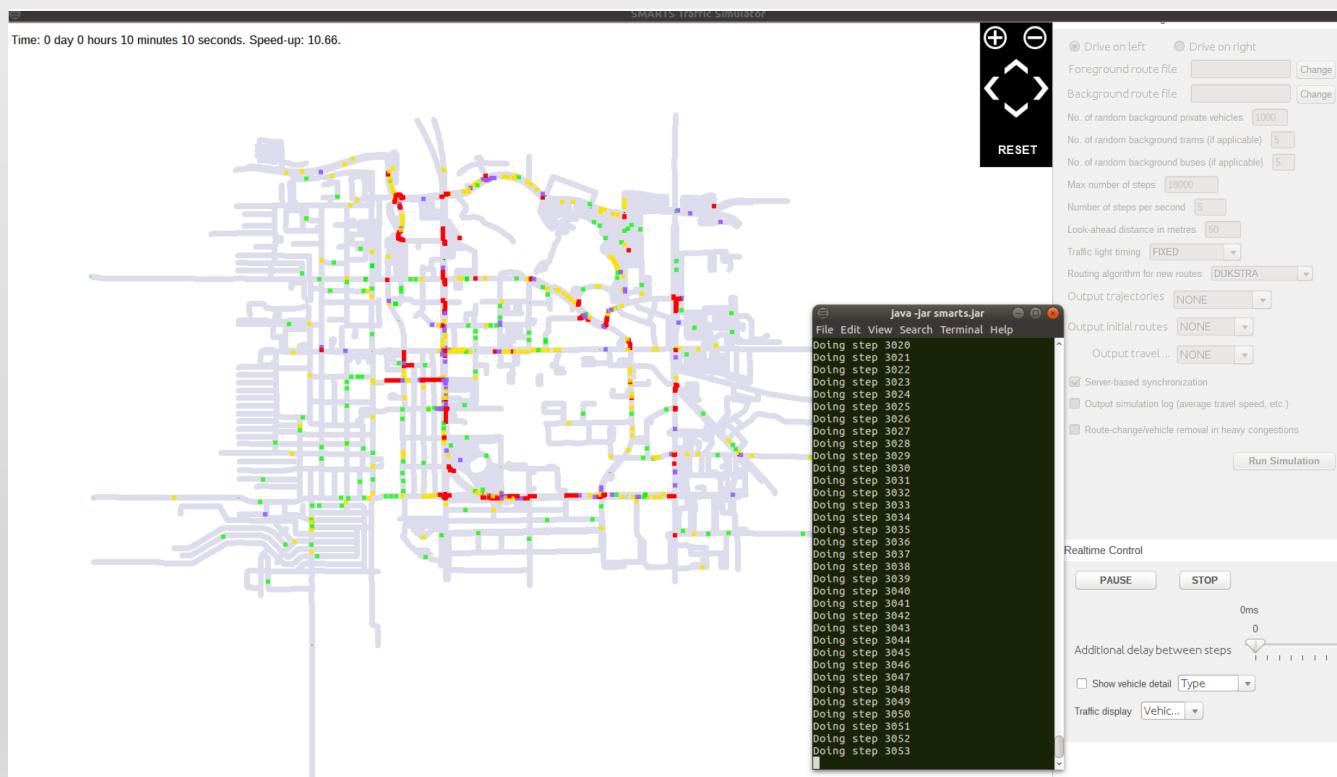
- Parallel traffic microscopic simulator

Zurich



# Transims [1], ParamGrid [2] and SMARTS [3]

- Distributed Microscopic Simulator



[1] Nagel, K. and M. Rickert, "Parallel implementation of the transims micro-simulation", *Parallel Computing* 27, 12, 1611–1639 (2001).

[2] Klefstad, R., Y. Zhang, M. Lai, R. Jayakrishnan and R. Lavanya, "A distributed, scalable, and synchronized framework for large-scale microscopic traffic simulation", in "Intelligent Transportation Systems, 2005. Proceedings. 2005 IEEE", pp. 813–818 (IEEE, 2005).

[3] Ramamohanarao, K., H. Xie, L. Kulik, S. Karunasekera, E. Tanin, R. Zhang and E. B. Khunayn, "Smarts: Scalable microscopic adaptive road traffic simulator", *ACM Transactions on Intelligent Systems and Technology (TIST)* 8, 2, 26 (2017).

- Computing Resource

Number of workers required  Apply  
Number of connected workers

- Map

Show road network graph  
 Show map image in background

Road map areas

Afghanistan  
Albania  
Algeria  
American Samoa  
Andorra  
Angola  
Anguilla  
Antigua and Barbuda  
Argentina  
Armenia

Import roads from OpenStreetMap file

- Miscellaneous Settings

Drive on left  Drive on right  
Foreground route file   
Background route file   
No. of random background private vehicles   
No. of random background trams (if applicable)   
No. of random background buses (if applicable)   
Max number of steps   
Number of steps per second   
Look-ahead distance in metres   
Traffic light timing    
Routing algorithm for new routes    
Output trajectories     
Output initial routes    
Output travel...    
 Server-based synchronization  
 Output simulation log (average travel speed, etc.)  
 Route-change/vehicle removal in heavy congestions

# Comparison

- Brinkhoff, BerlinMOD
- SUMO, Vissim
- MATSim
- Transims, ParamGrid, SMARTS

GeoSparkSim

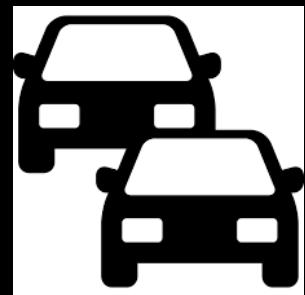
Driving behaviors

Scalable

Workload balancing

Spatial distribution over time

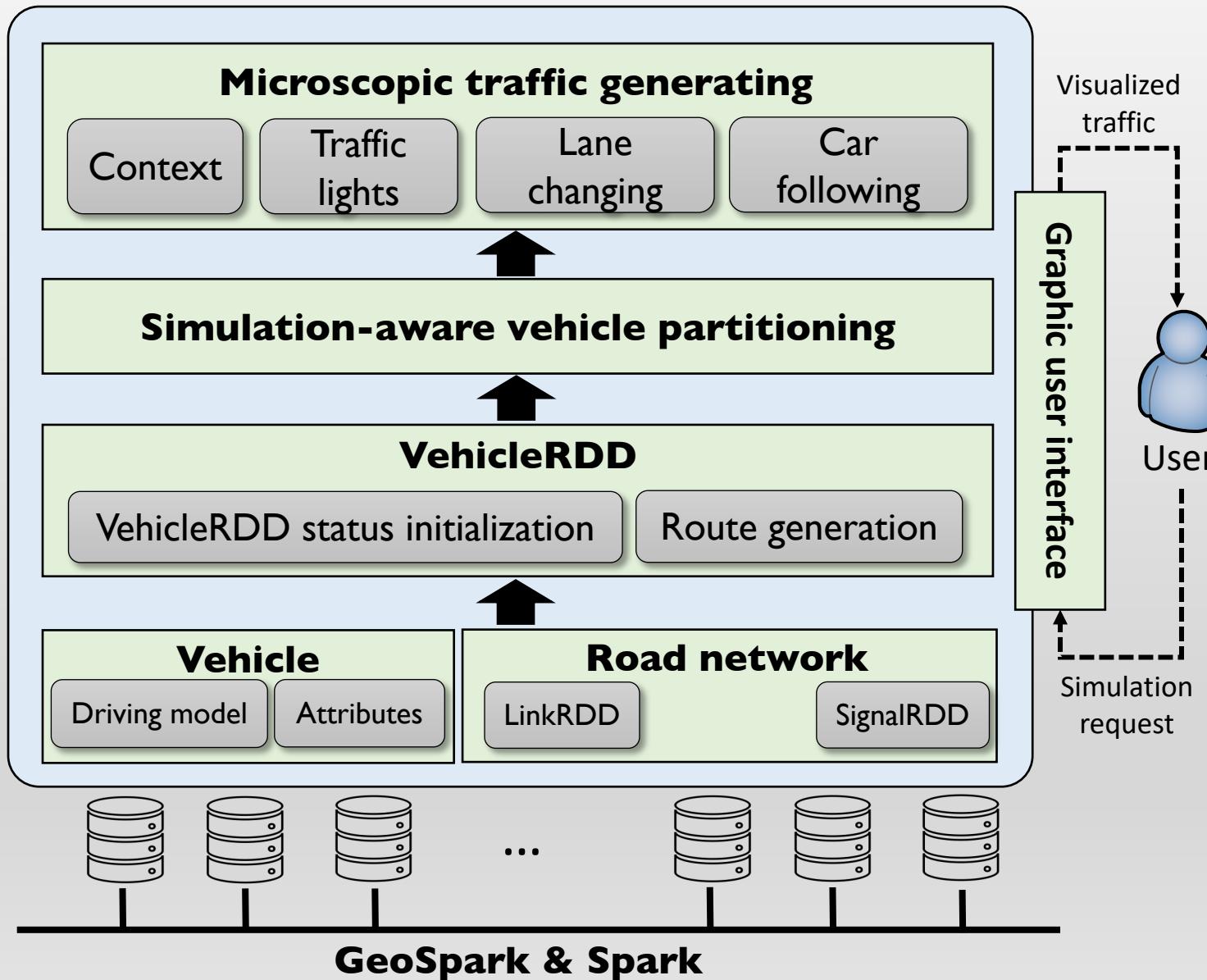
Large-scale data analytics in Apache Spark and GeoSpark



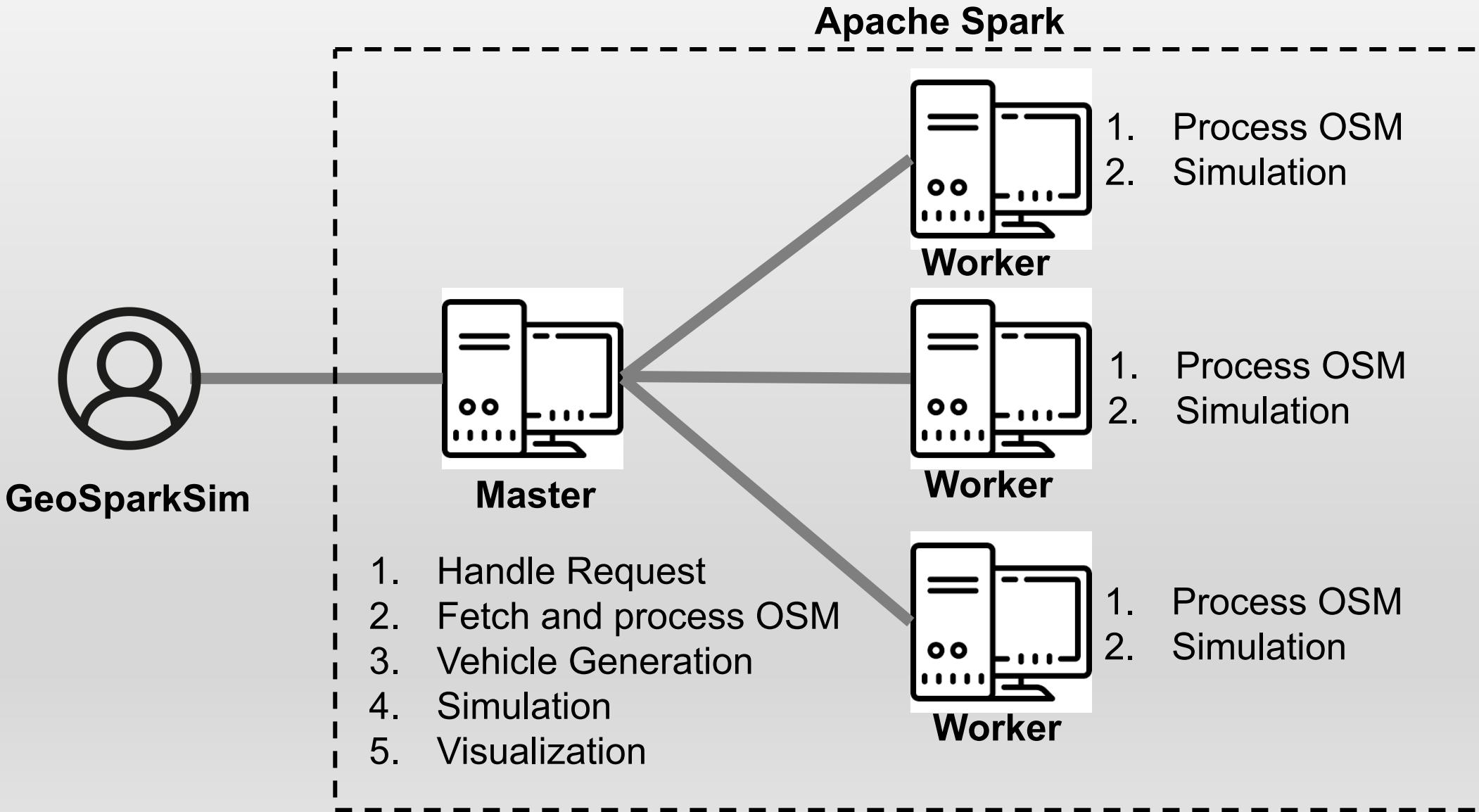
# GeoSparkSim

- Related Works
- **Architecture**
- Simulation
- Performance
- Demo

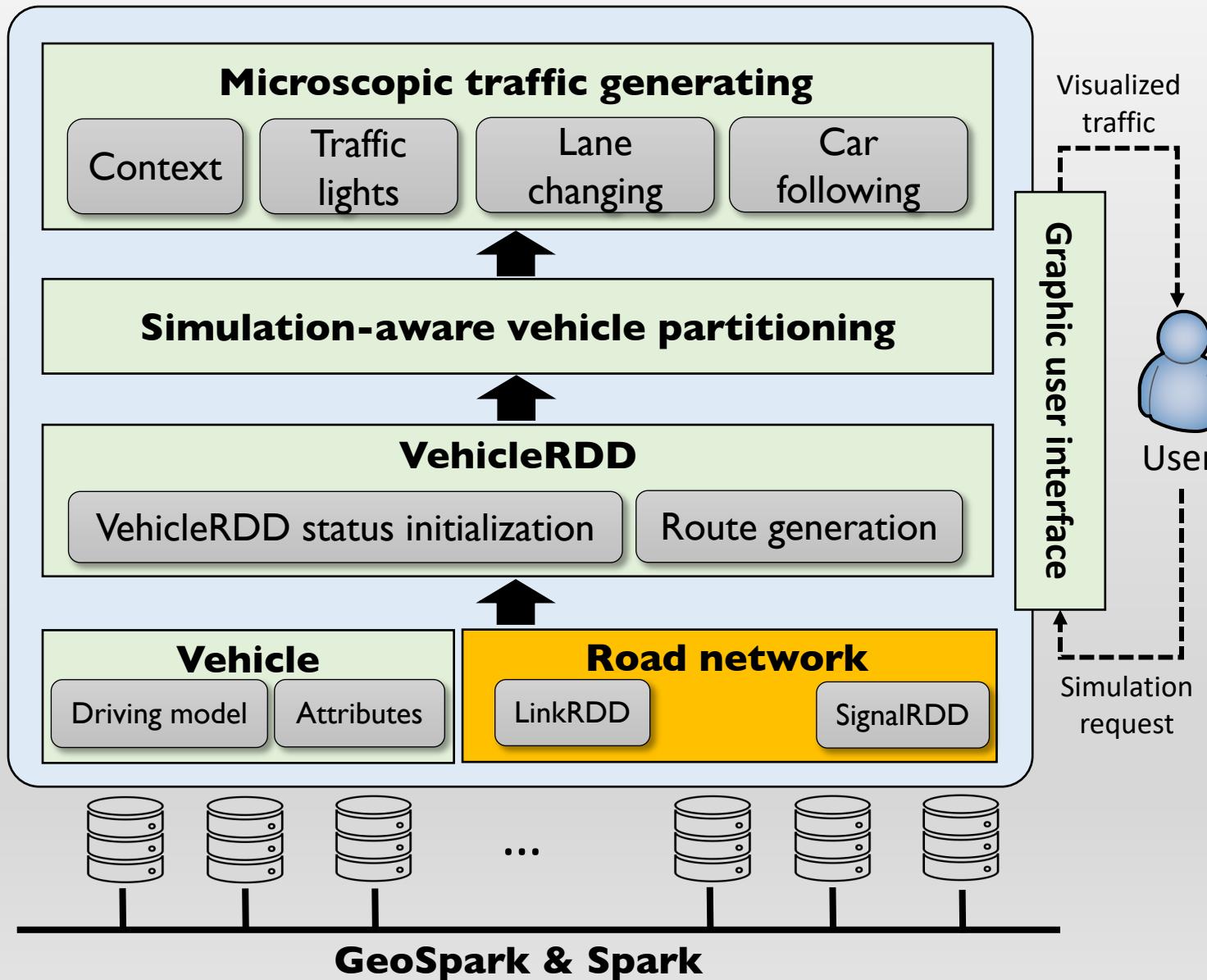
# Architecture



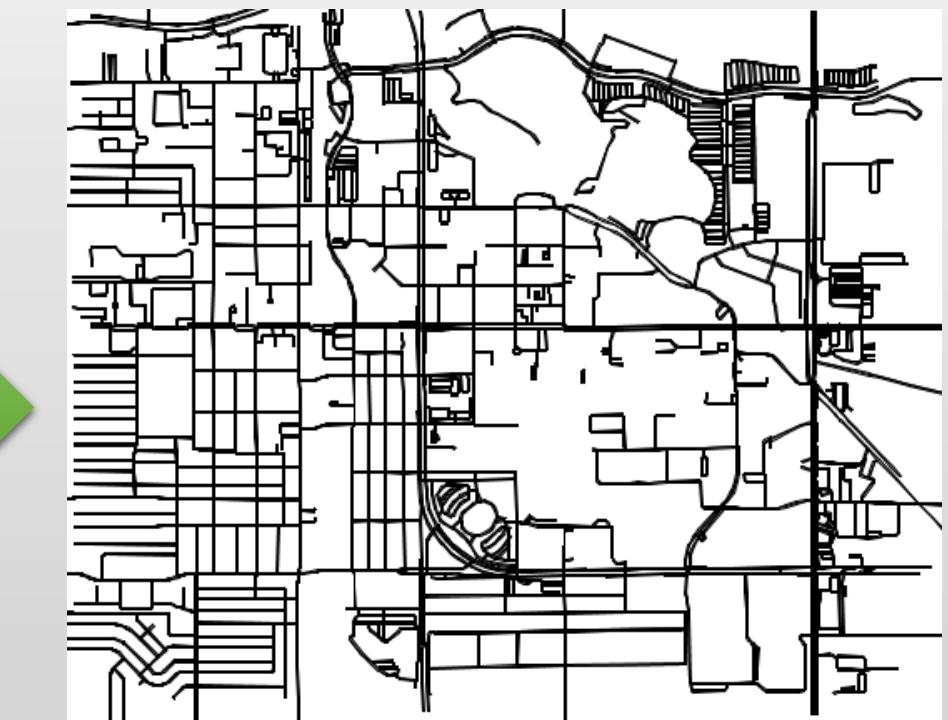
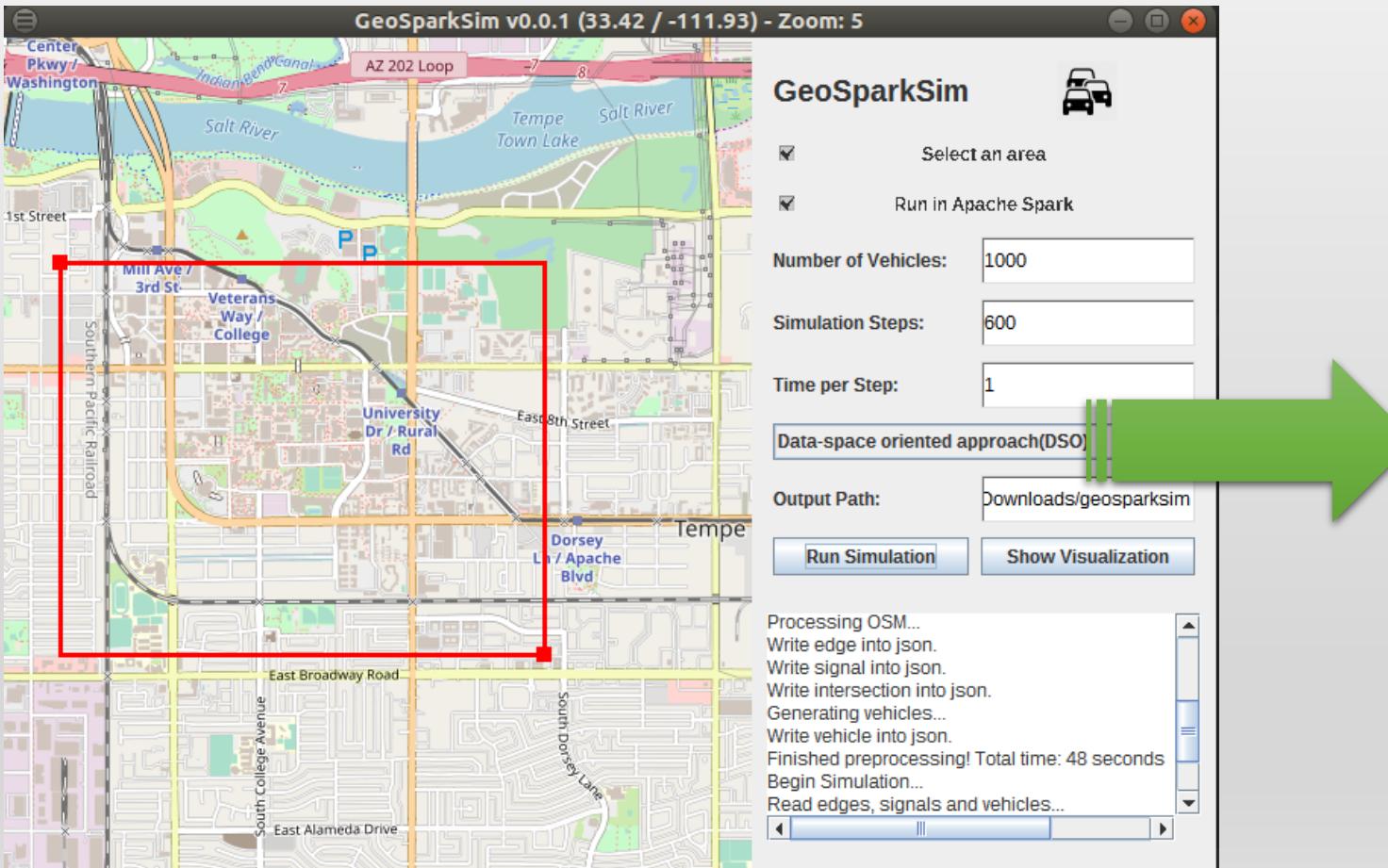
# Distributed Architecture



# Architecture



# Road Network



# Road Network

```

<?xml version="1.0" encoding="UTF-8"?>
<osm version="0.6" generator="Overpass API">

<node id="1" lat="33.4231268" lon="-111.9304264">
    <tag k="highway" v="traffic_signals"/>
</node>
<node id="2" lat="33.4245016" lon="-111.9275146" />

<way id="1">
    <nd ref="1"/>
    <nd ref="2"/>
    <tag k="highway" v="primary"/>
    <tag k="lanes" v="8"/>
    <tag k="oneway" v="no"/>
</way>

</osm>

```

```

Node: {
    ID: 1,
    Latitude: 33.4231268,
    Longitude: -111.9304264,
    Tags: {
        tag1: {
            key: highway,
            value: traffic_signals
        },
        tag2: ...
    }
}

```

Diagram illustrating the structure of a Node:

```

graph LR
    SignalRDD[SignalRDD] --- Node[Node]
    subgraph Node
        direction TB
        ID[ID: long]
        Tags[Tags: array]
        element[element: struct]
        key[key: string]
        value[value: string]
        Latitude[Latitude: double]
        Longitude[Longitude: double]
    end

```

```

Way: {
    ID: 1,
    Nodes: {
        node1: {
            index: 1,
            nodeID: 1
        },
        node2: {
            index: 2,
            nodeID: 2
        }
    },
    Tags: {
        tag1: {
            key: highway,
            value: traffic_signals
        },
        tag2: ...
    }
}

```

Diagram illustrating the structure of a Way:

```

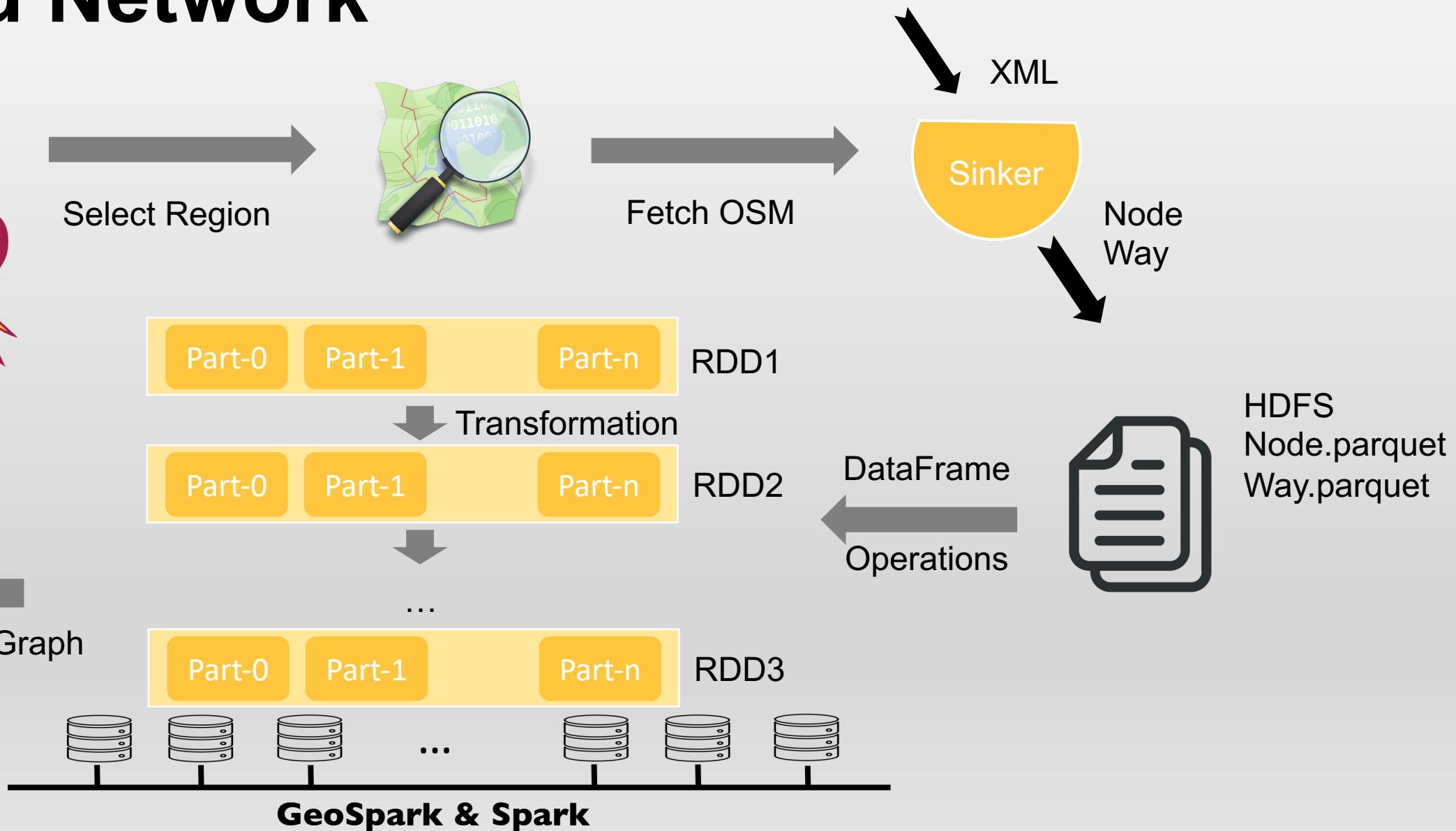
graph LR
    LinkRDD[LinkRDD] --- Way[Way]
    subgraph Way
        direction TB
        ID[ID: long]
        Tags[Tags: array]
        element[element: struct]
        key[key: string]
        value[value: string]
        Nodes[Nodes: array]
        elementN[element: struct]
        indexN[index: integer]
        nodeIDN[nodeID: long]
    end

```

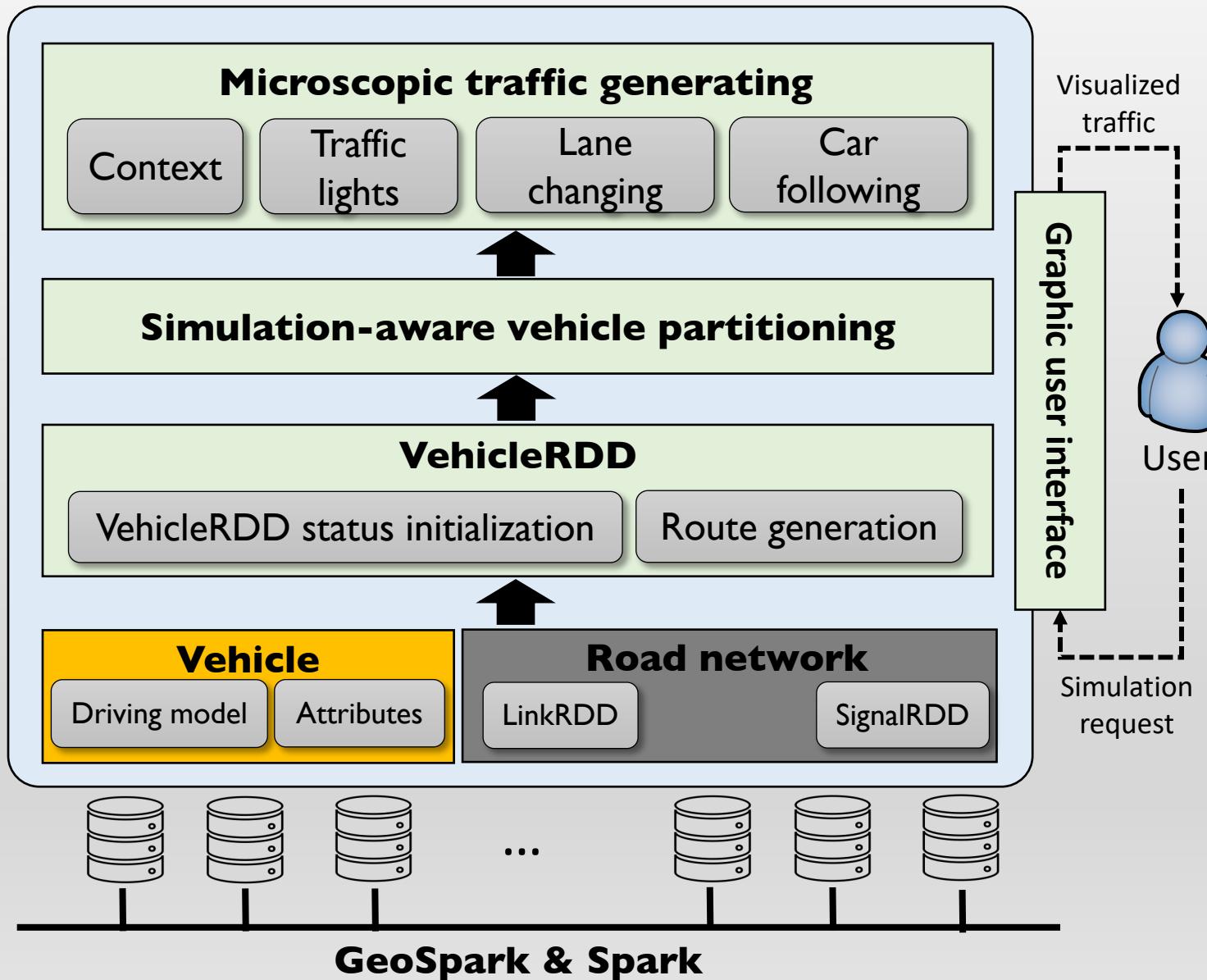
# Road Network



Road Network Graph  
linkRDD  
signalRDD



# Architecture



# Vehicle

**Black box: field**

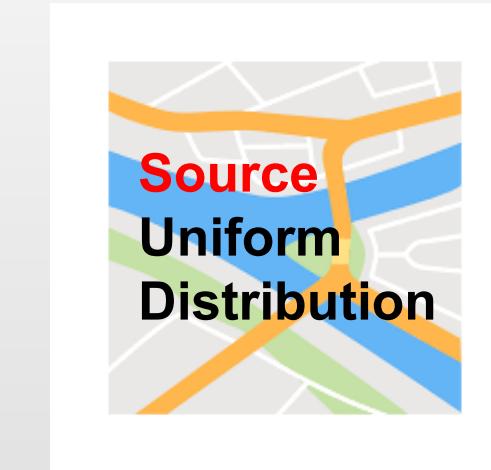
**Red box: method**

## Vehicle

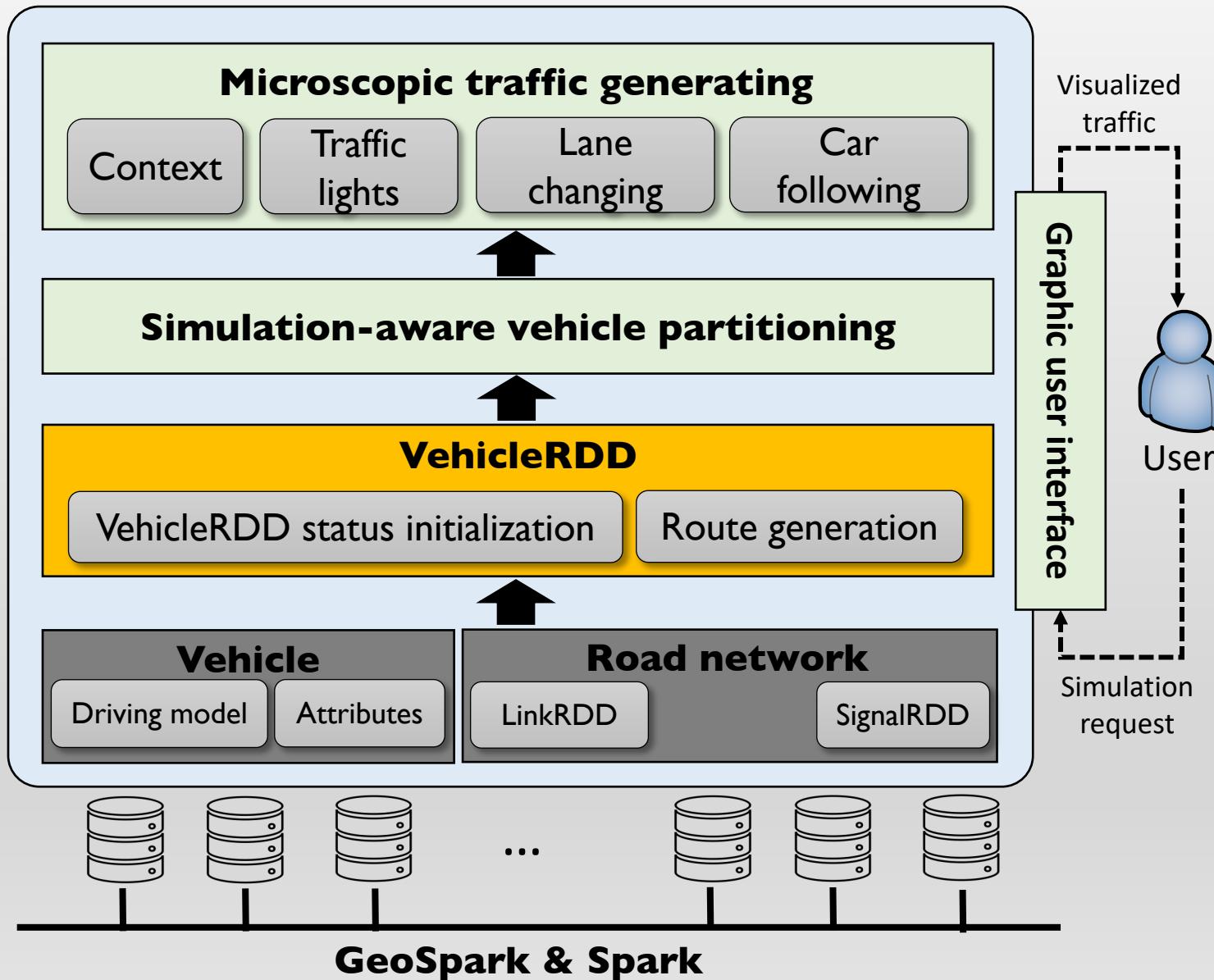
ID: String  
Source: Coordinate  
Target: Coordinate  
EdgePath: Long Array  
Costs: Double Array  
FullPath: Coordinate Array  
Front: Coordinate  
Rear: Coordinate  
edgeIndex: int  
carLength: 3  
currentLane: int  
isArrive: boolean  
currentLink: Link

getLinkByHead  
headWayCheck  
basicMovement  
born

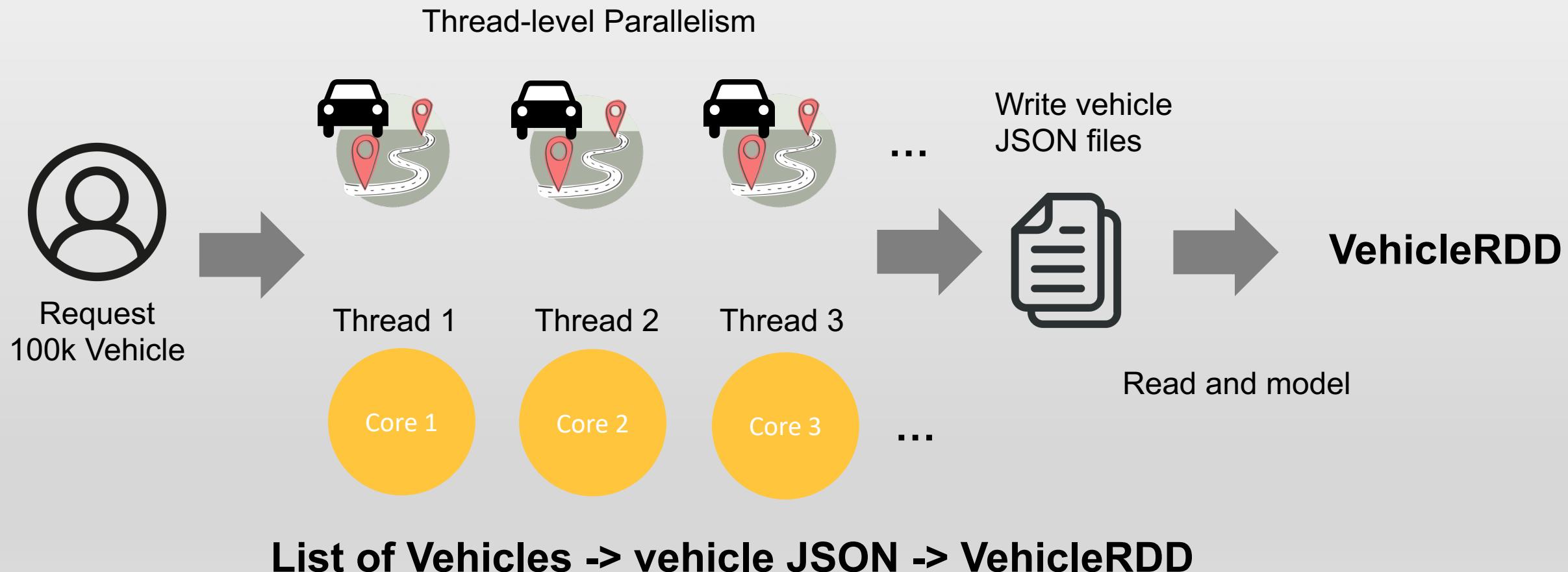
- Vehicle ID
- Source coordinate
- Destination coordinate
- Compute path trajectory
- Array of link ID
- Array of link distance cost
- Other attributes (car length, etc)



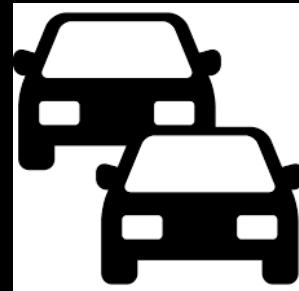
# Architecture



# Vehicle to VehicleRDD

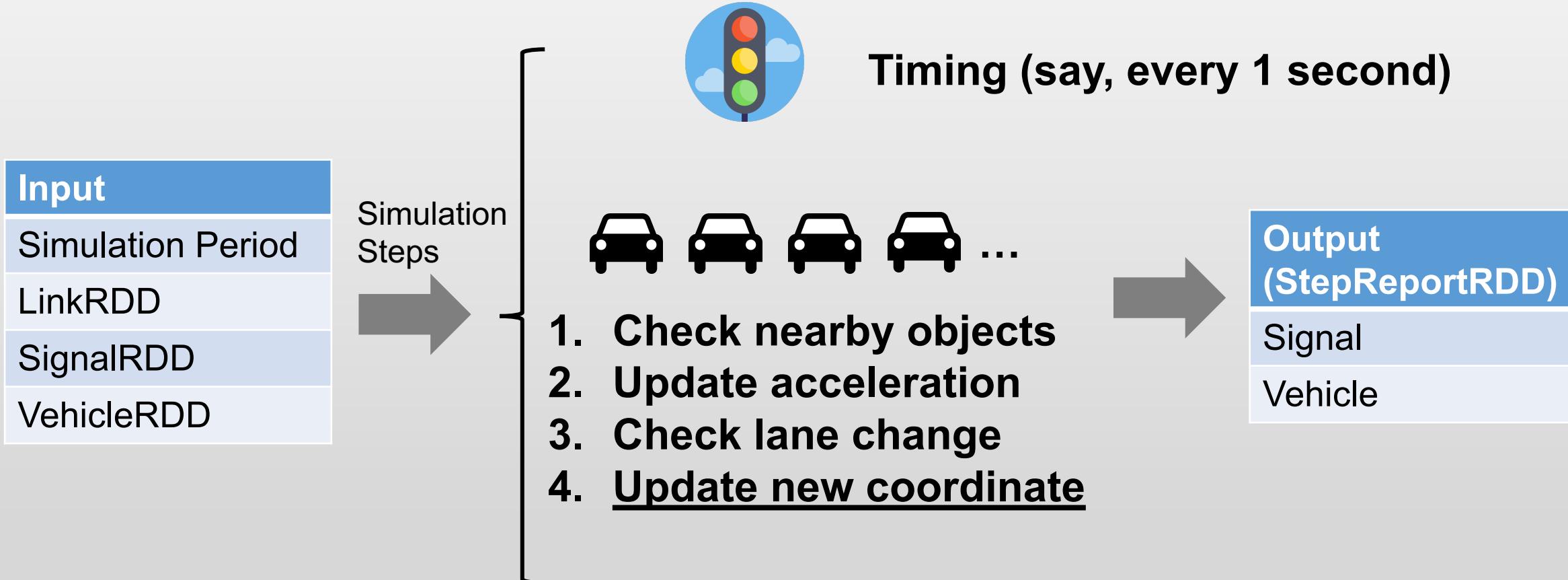


# GeoSparkSim

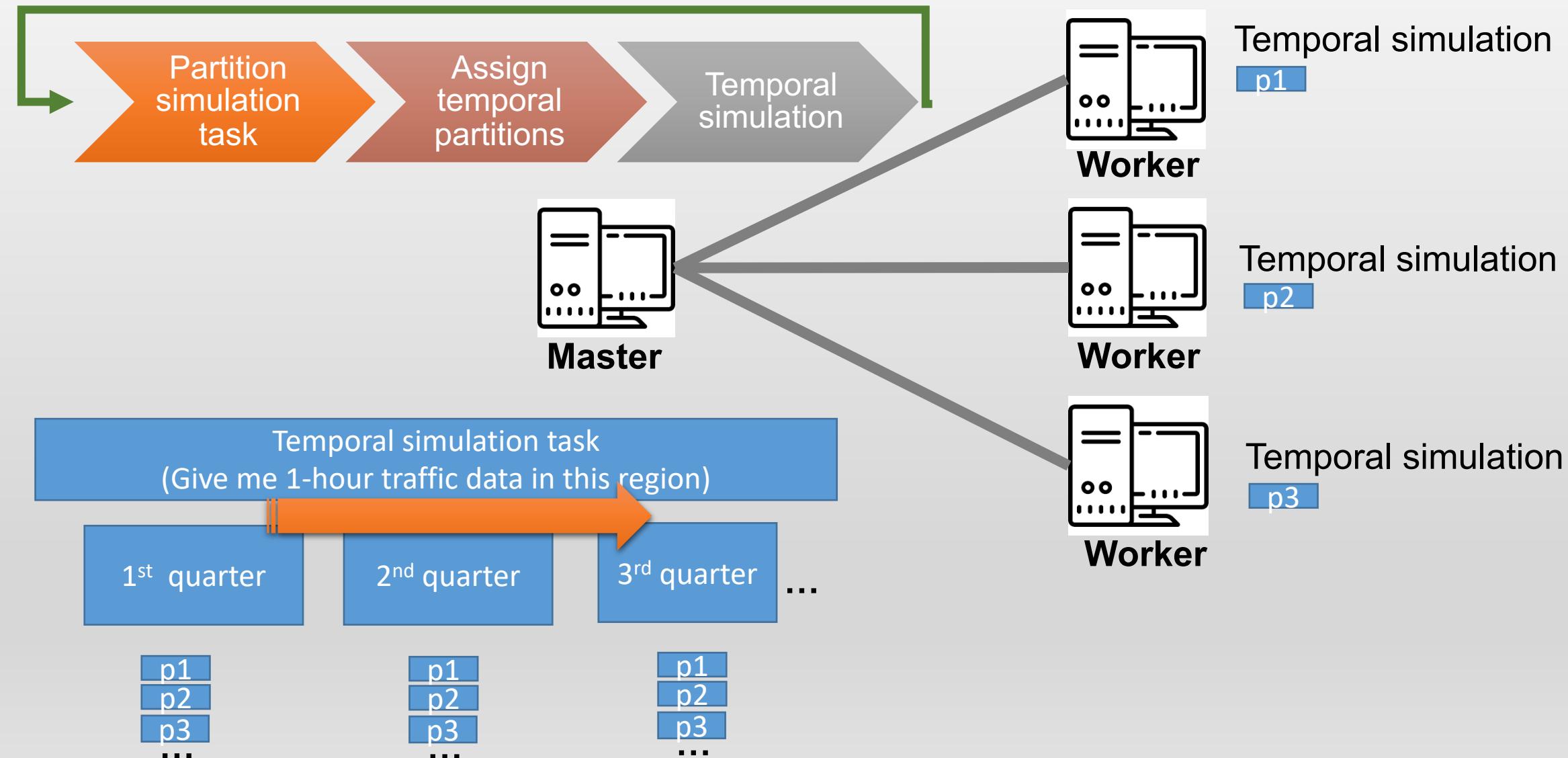


- Related Works
- Architecture
- **Simulation**
- Performance
- Demo

# Simulation

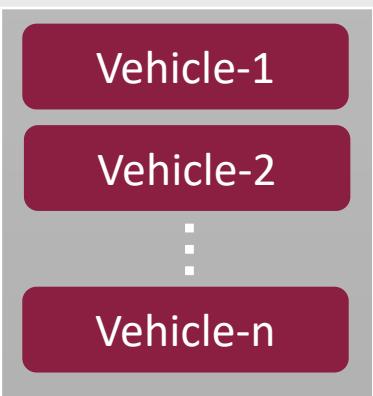


# Workload Distribution



# Spatial Resilient Distributed Dataset

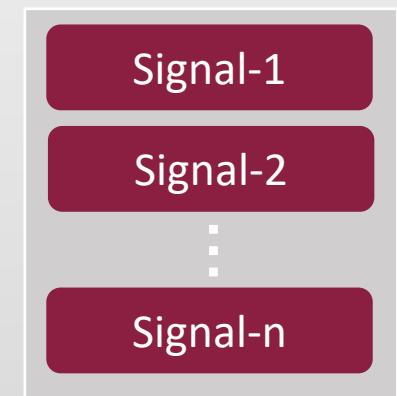
**VehicleRDD**



**LinkRDD**



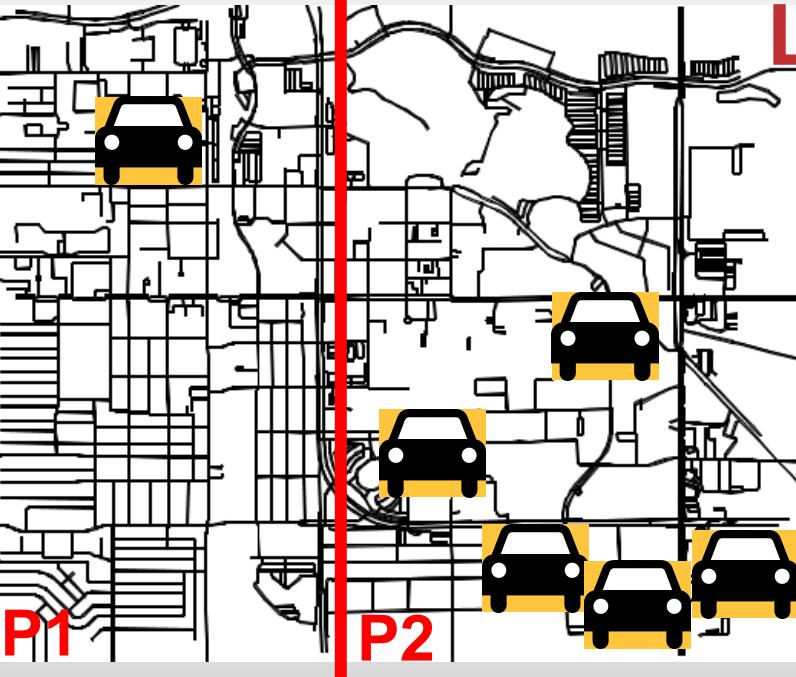
**SignalRDD**



Data Space -> Space-partitioning Tree

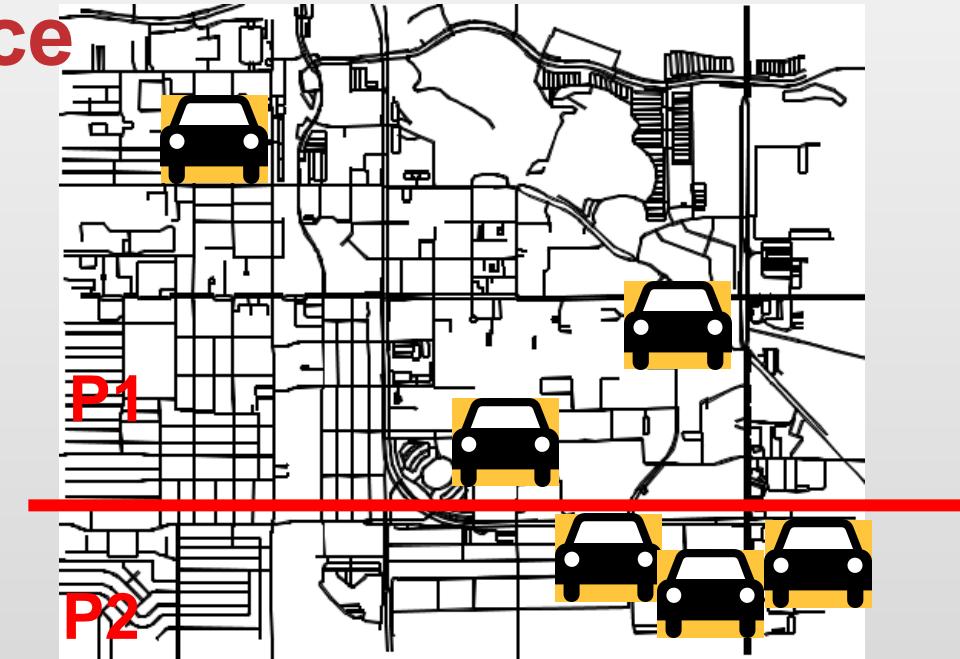
**Apply same partition to linkRDD and signalRDD**

# Spatial Workload Distribution



Load balance

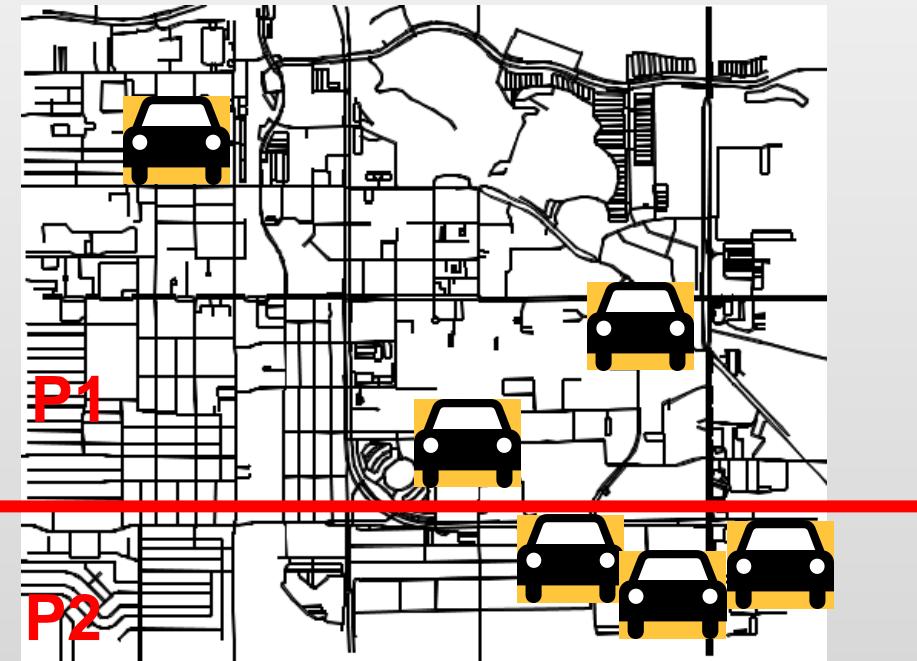
Partition by road network



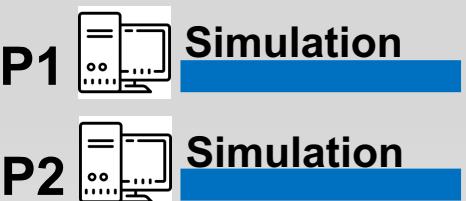
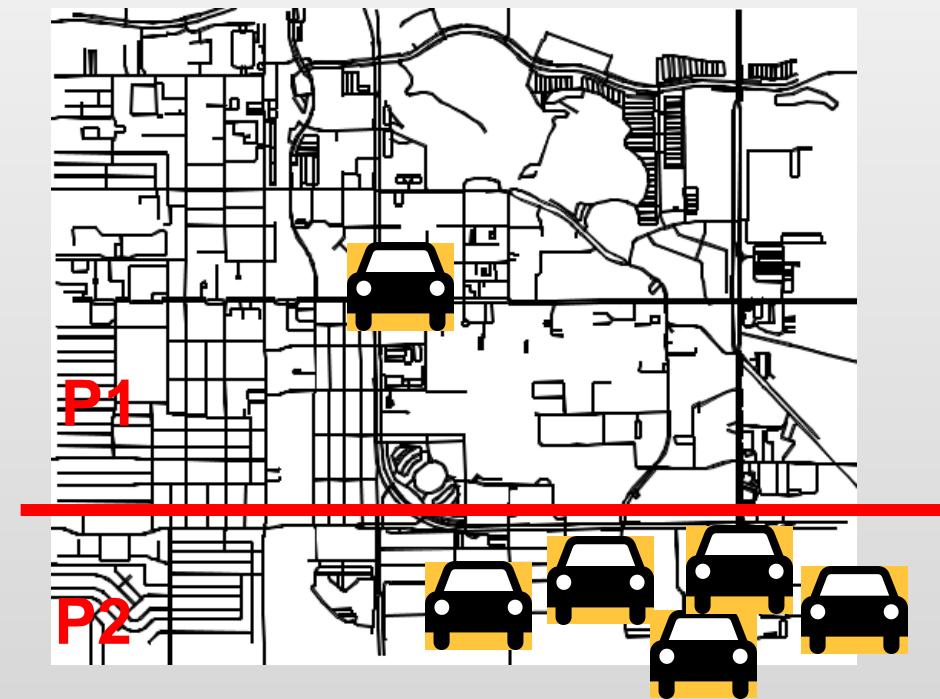
Partition by vehicle density



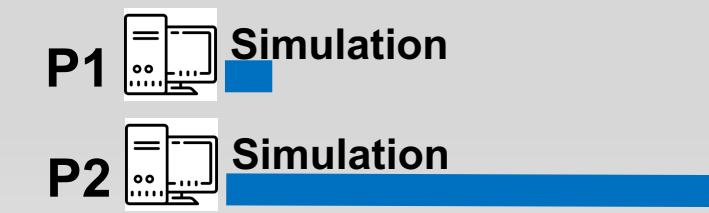
# Spatial Workload Balancing



2 min



Uneven!





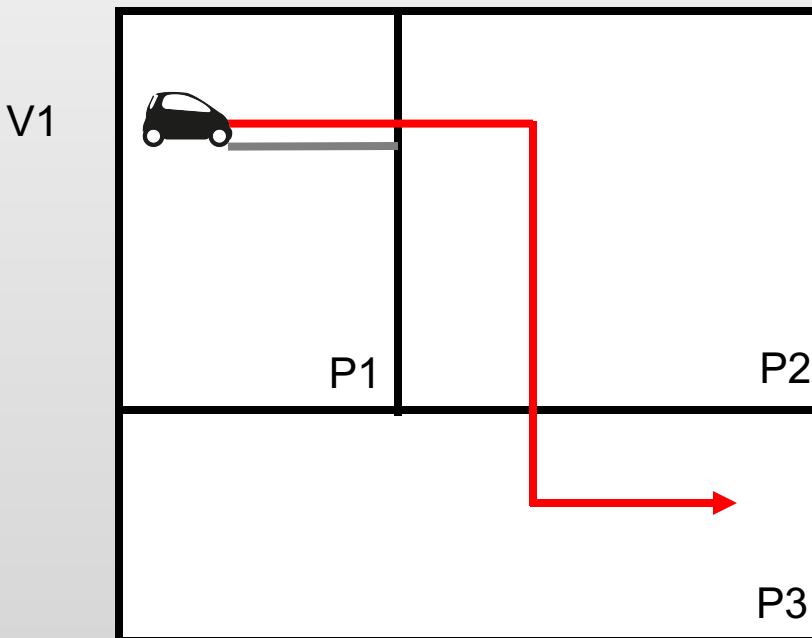
**How to do spatial partition in  
the simulation?**

**Vehicles are moving!**

**Source coordinate?**

**Planned route?**

# Partition by Source coordinate

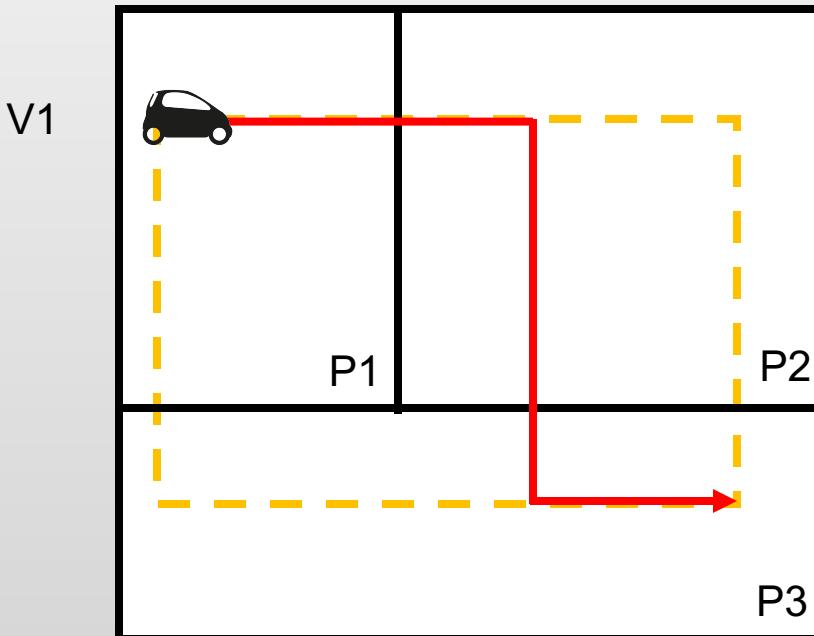


**Update**

P1: V1 —————→ P2: V1

**Message exchange, inefficient!**

# Partition by planned routes

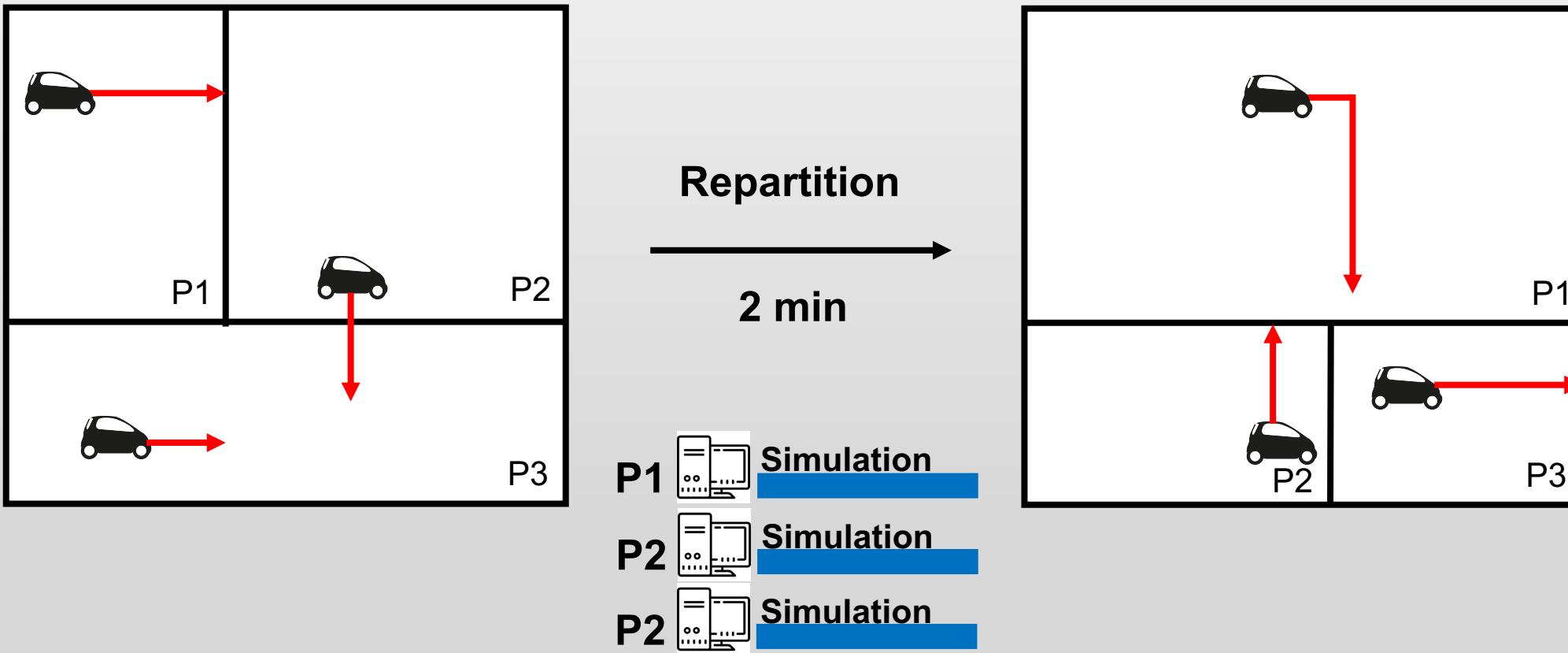


Minimum Bounding Rectangle  
(MBR)

P1: V1  
P2: V1  
P3: V1

Routes are too long  
Overlaps, intersections!

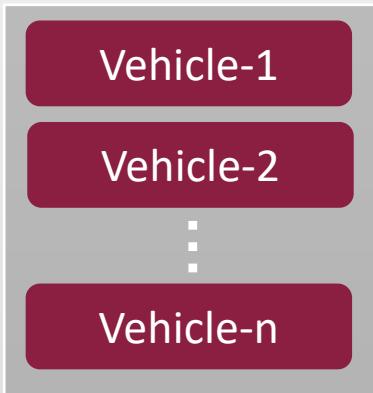
# Spatial Workload Balancing



**Update partition by spatial distribution over time!**

# Spatial Resilient Distributed Dataset

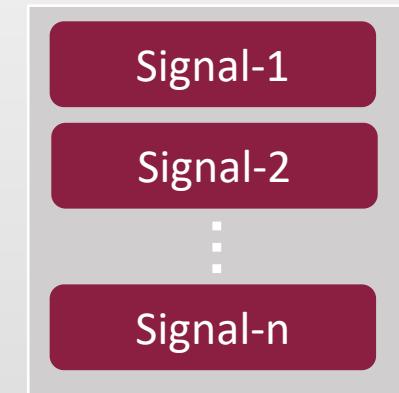
## VehicleRDD



## LinkRDD



## SignalRDD



Data Space -> Space-partitioning Tree

Part-0   Part-1   ...   Part-n

Spatial Proximity

Part-0   Part-1   ...   Part-n

Zip Partitions

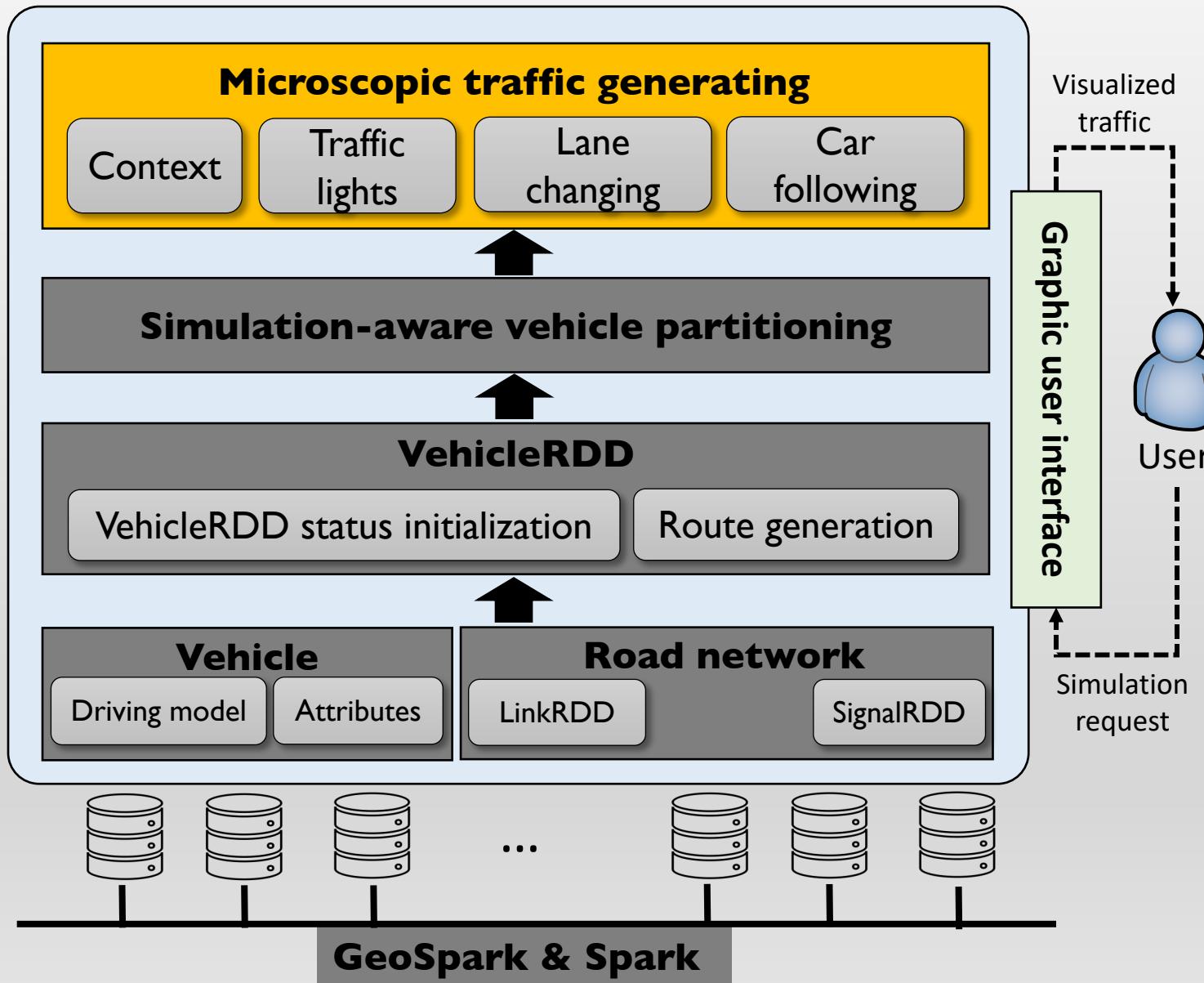
Part-0   Part-2   Part-3   Part-4   ...   Part-p



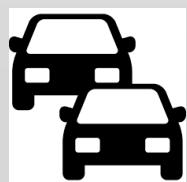
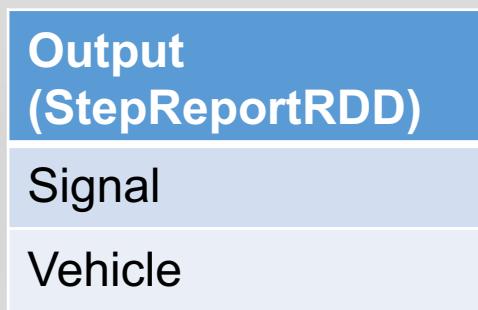
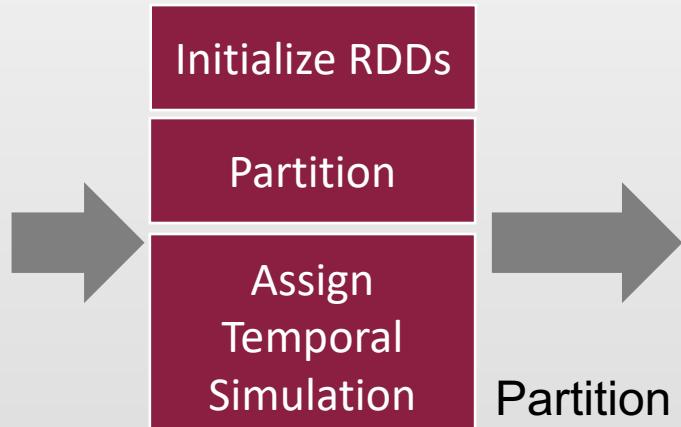
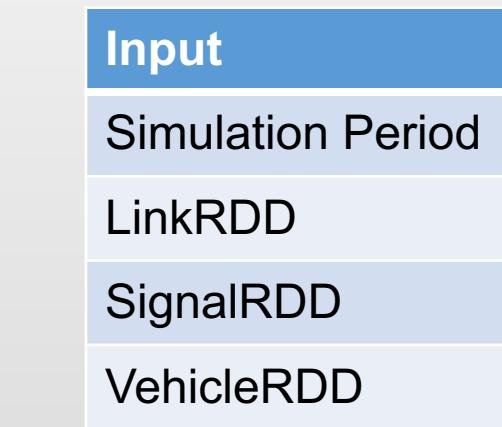
Cluster

Run local Simulation  
in parallel

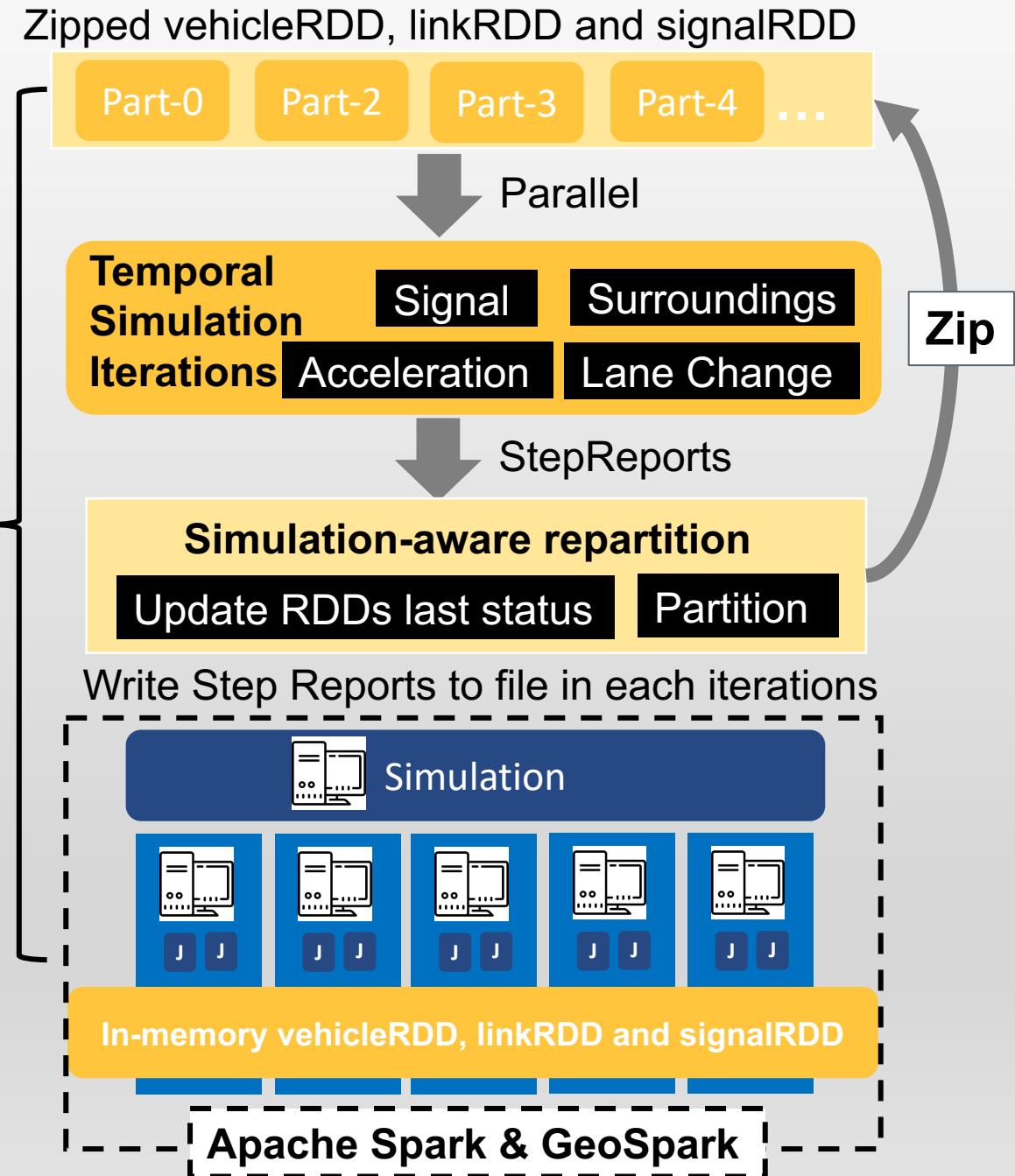
# Architecture



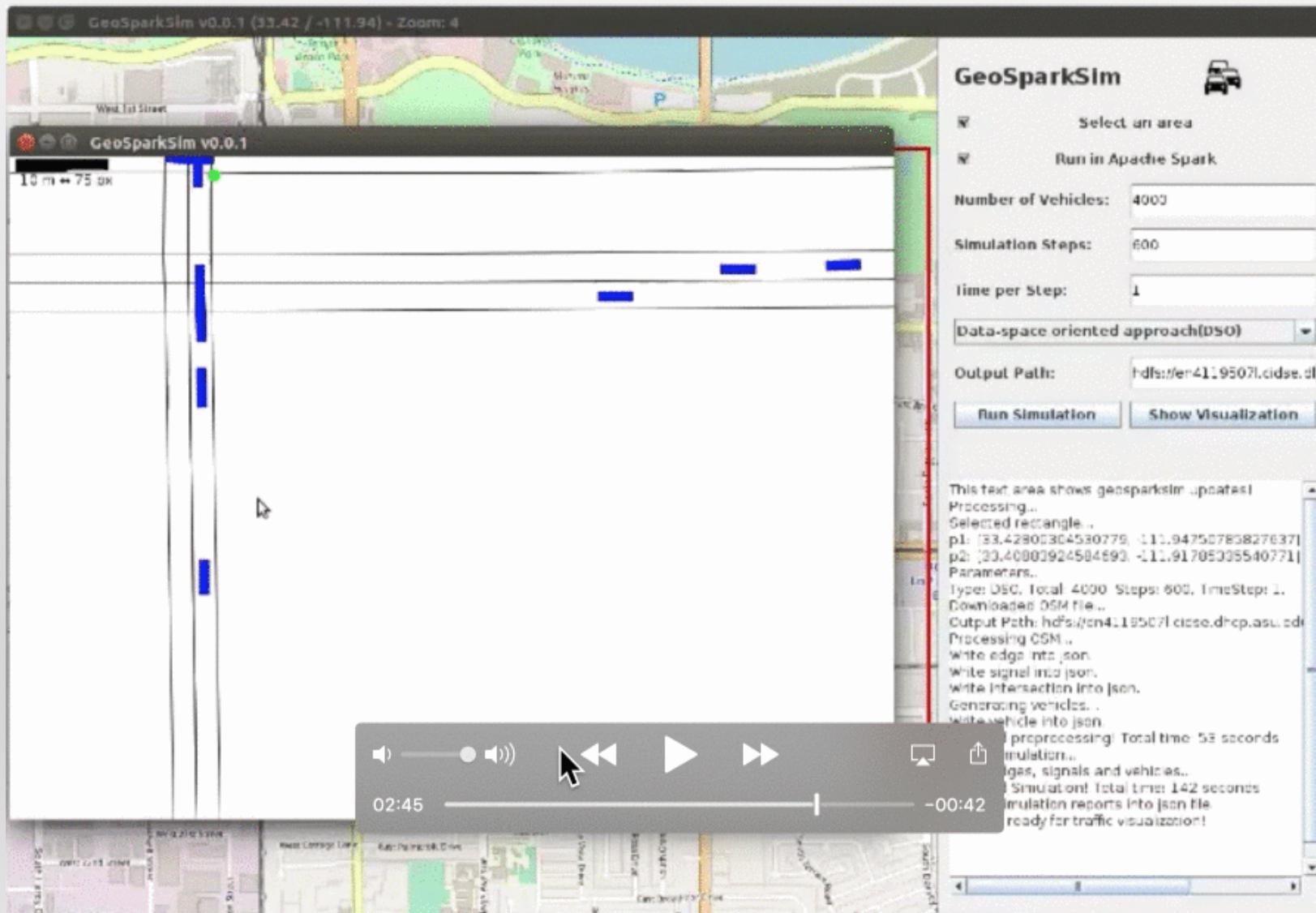
# Simulation

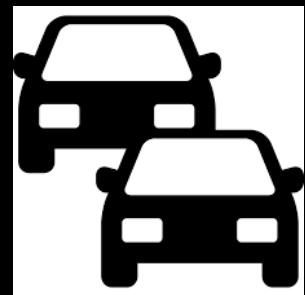


**GeoSparkSim**



# Visualized Traffic Data

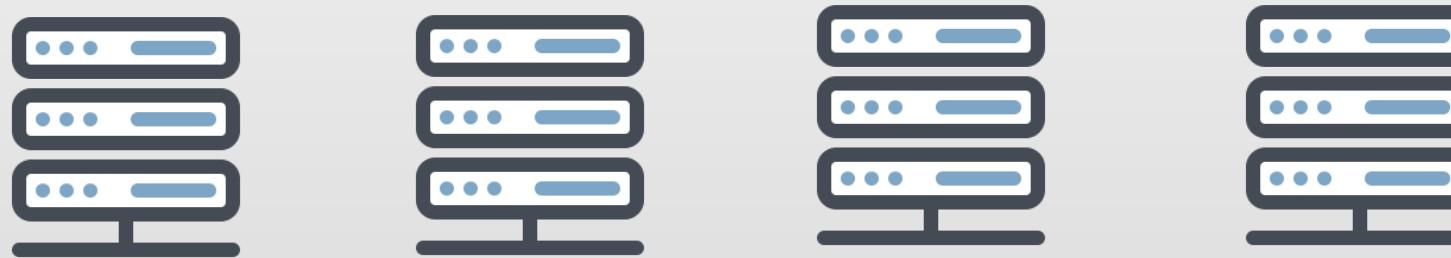




# GeoSparkSim

- Related Works
- Architecture
- Simulation
- **Performance**
- Use

## Cluster: One master and four worker



Machine: 48 cores, 100 GB memory, 4TB HDD

# Parameters

Parameters	Range
Number of vehicles (thousand)	<u>100</u> , 200, 300
Time step (second)	1, 0.8, 0.6, 0.4, 0.2
Simulation period (minute)	<u>10</u> , 30, 60, 120
Number of Partitions	1000, <u>1500</u> , 3000
Repartition period (minute)	1, <u>2</u> , 4, 8

# The Number of Vehicles

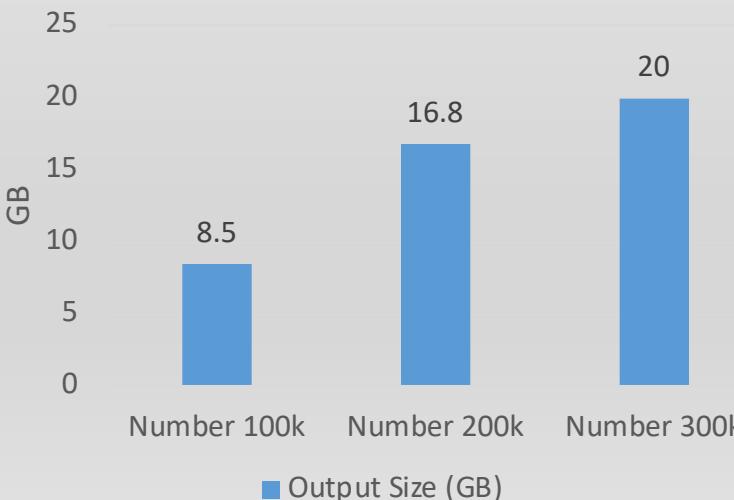
Number of Vehicles: **100k, 200k, 300k**

Time Step: **1 sec**

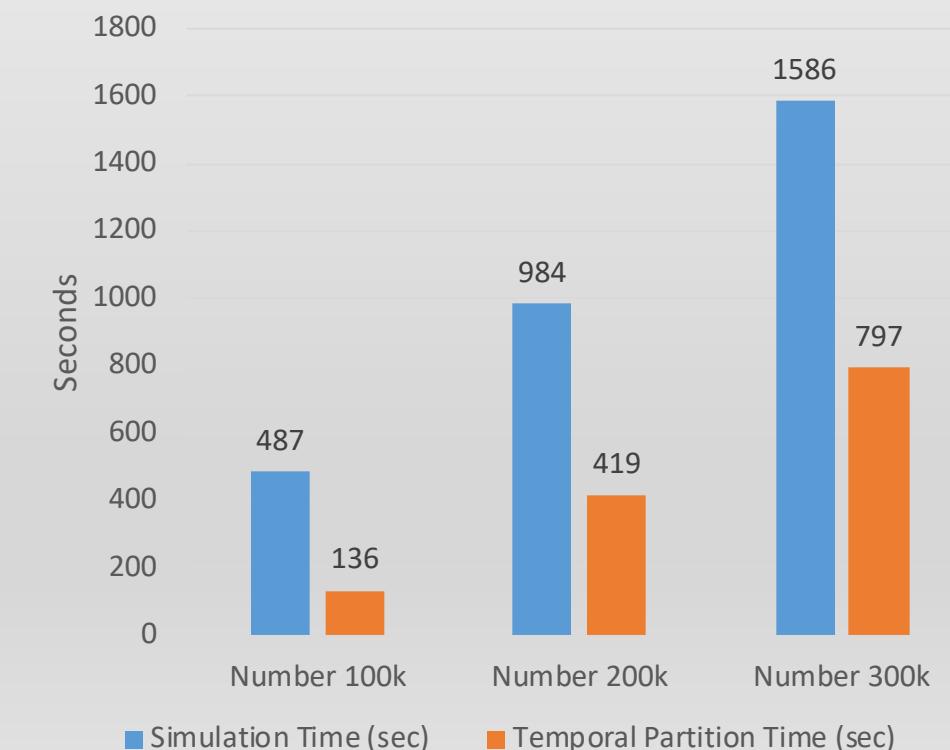
Simulation Period: **10 min**

Repartition Period: **2 min**

Number of Partition: **1500**



**Simulation, temporal partition, output size linear increase**



# Simulation Period

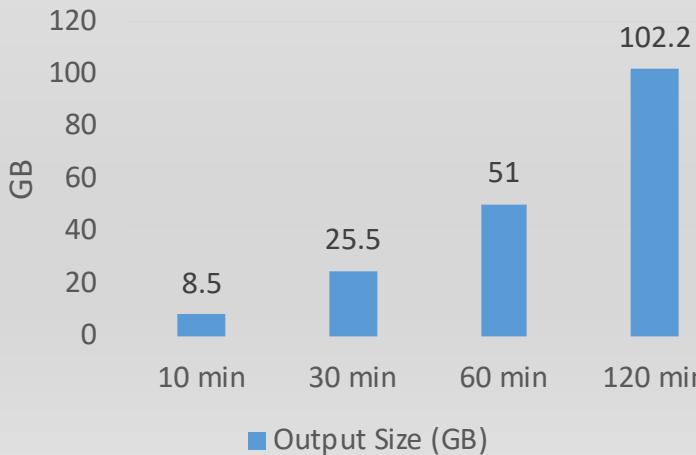
Number of Vehicles: **100k**

Time Step: **1 sec**

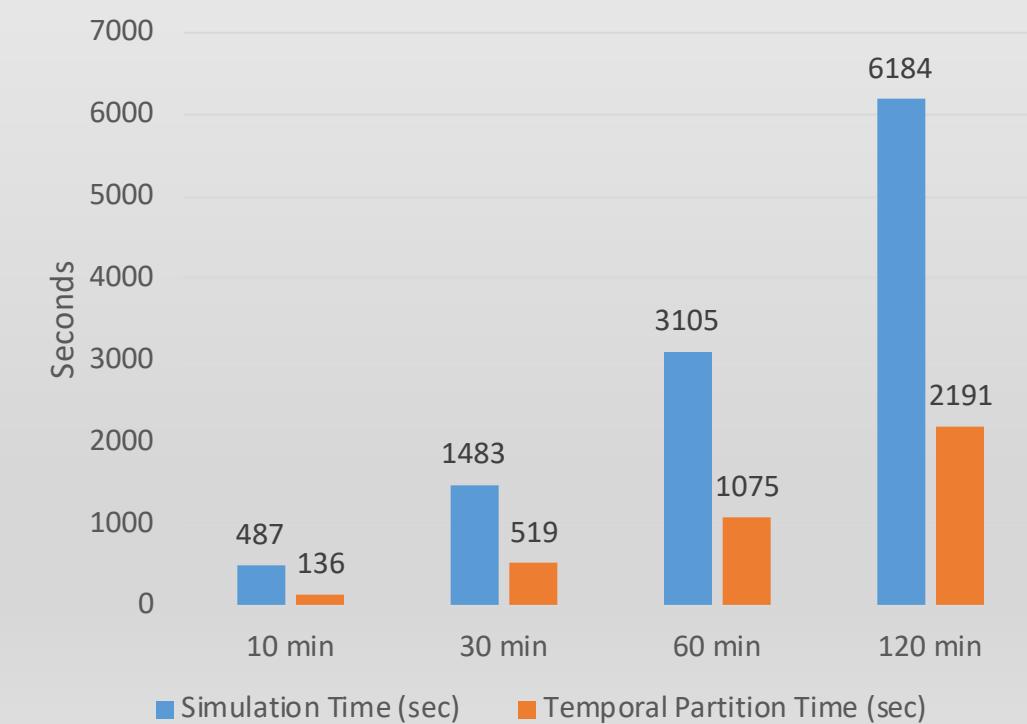
Simulation Period: **10, 30, 60, 120 min**

Repartition Period: **2 min**

Number of Partition: **1500**



**Linear increase**



# Time per Steps

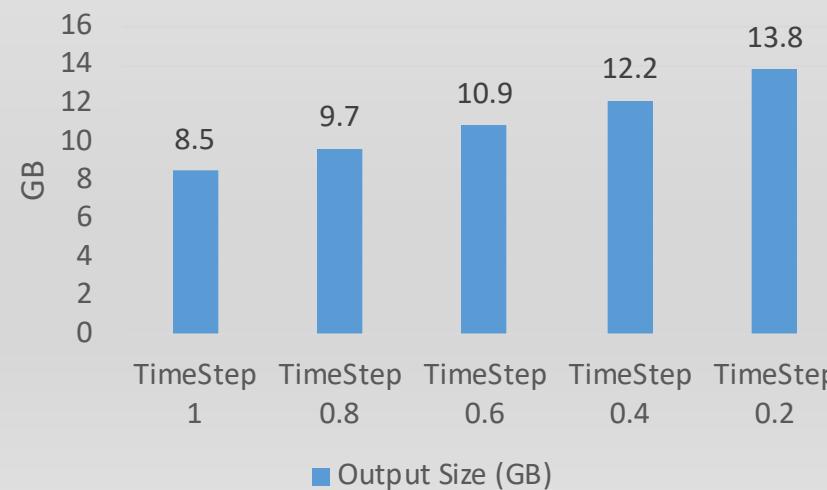
Number of Vehicles: **100k**

Time Step: **1, 0.8, 0.6, 0.4, 0.2 sec**

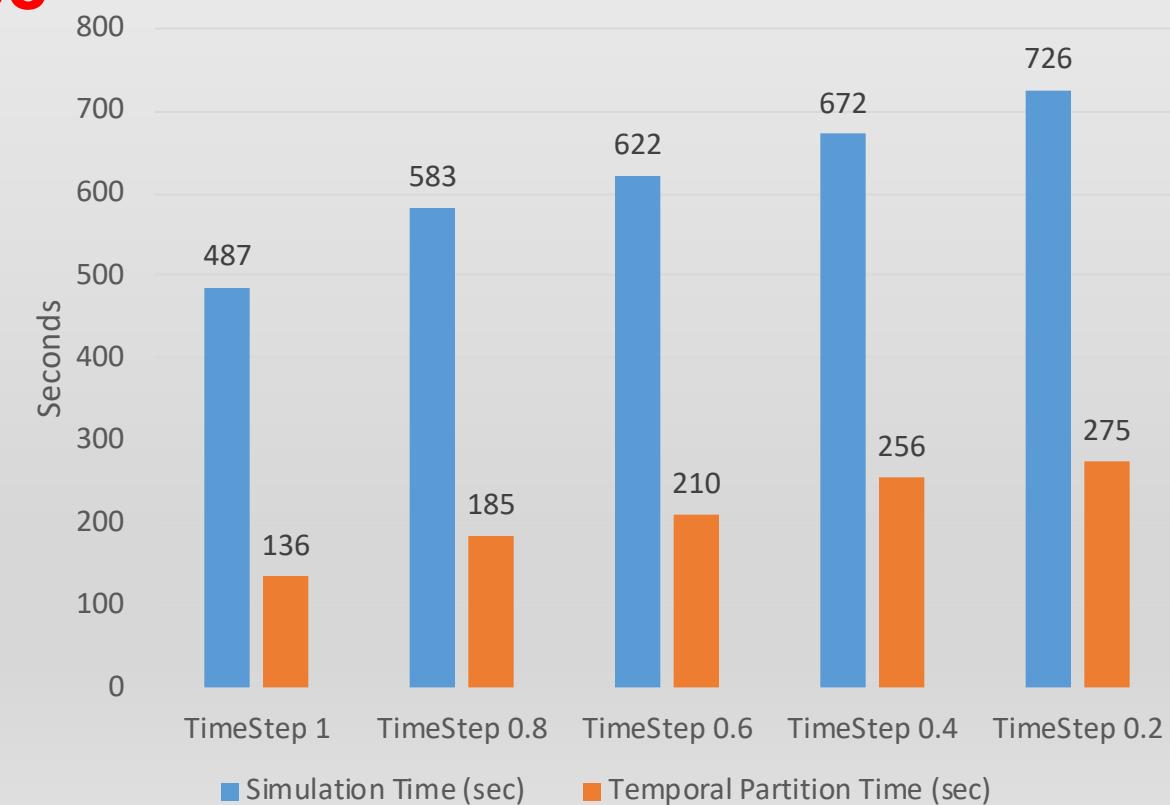
Simulation Period: **10 min**

Repartition Period: **2 min**

Number of Partition: **1500**



**Linear increase**



# The Number of Partitions

Number of Vehicles: **100k**

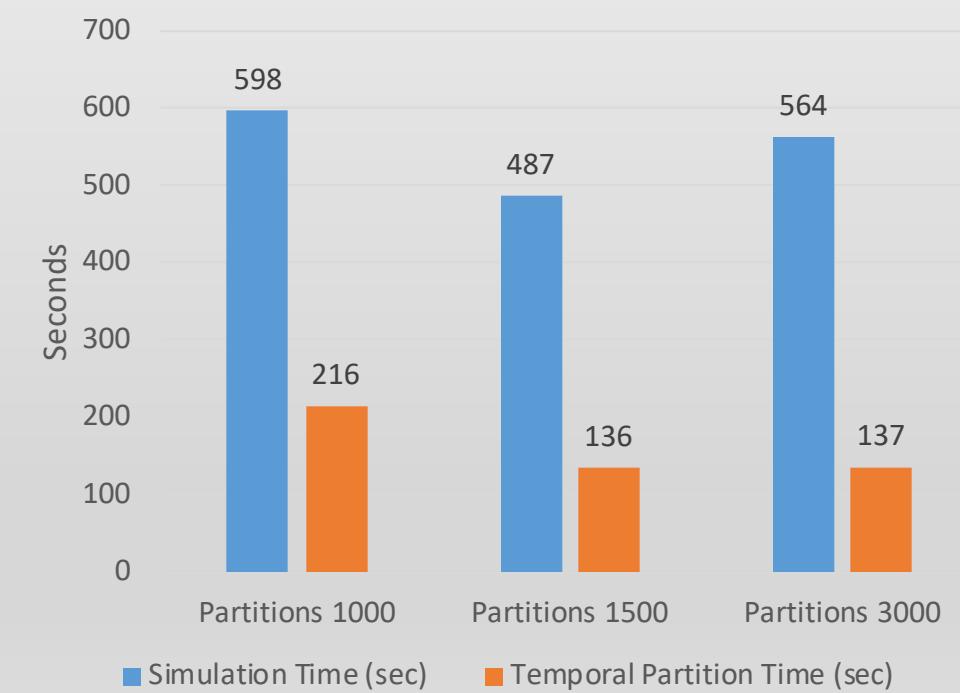
Time Step: **1 sec**

Simulation Period: **10 min**

Repartition Period: **2 min**

Number of Partition: **1000, 1500, 3000**

**More parallelism and  
more data shuffle costs**



# Repartition Period – Temporal

Number of Vehicles: **100k**

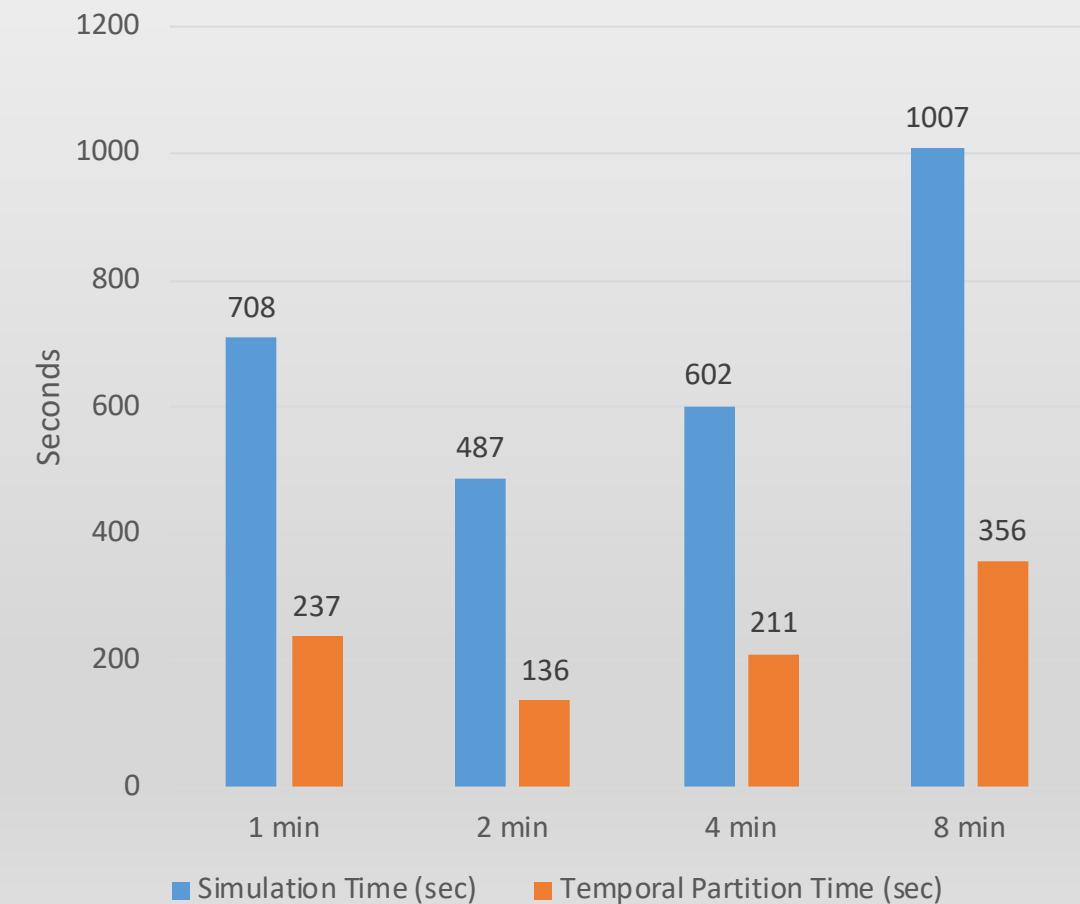
Time Step: **1 sec**

Simulation Period: **10 min**

Repartition Period: **1, 2, 4, 8**

Number of Partition : **1500**

**Increases repartition times  
leads to longer partition time.  
Decreases repartition times  
causes unbalanced workload**



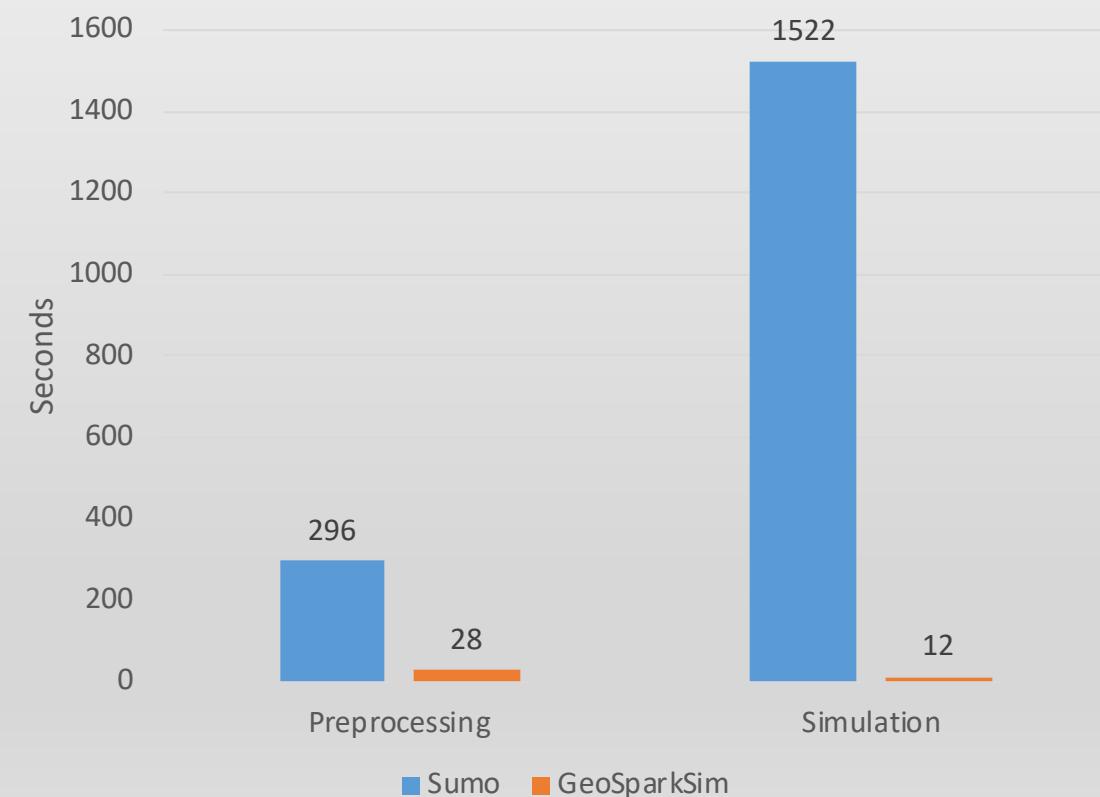
# Sumo and GeoSparkSim

Number of Vehicles: **1000**

Time Step: **1 sec**

Simulation Period: **1 min**

**Scalability**



# Smarts and GeoSparkSim

**100k vehicles, 1 sec time step,**

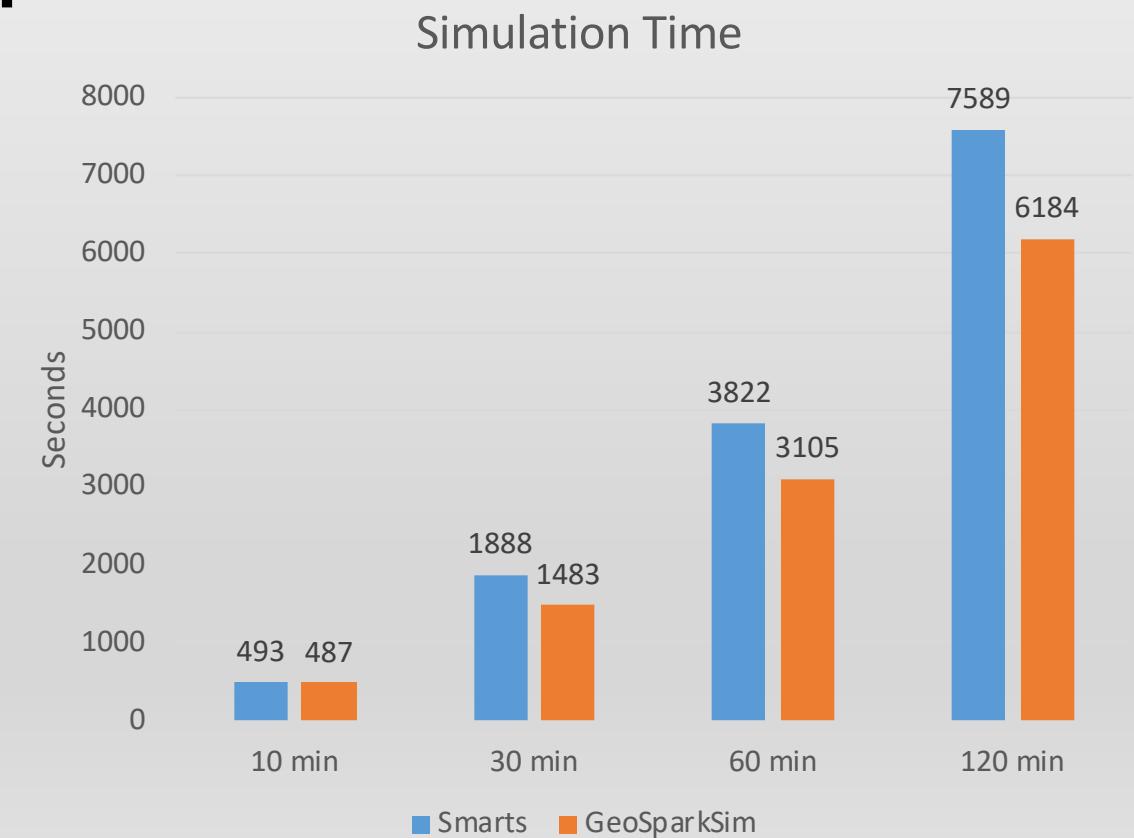
Simulation Period: **10, 30, 60, 120 min**

Setting: **4 machines**

$$\text{Speedup} = \frac{\text{Simulation period}}{\text{Simulation time}}$$



**Workload balancing over time**

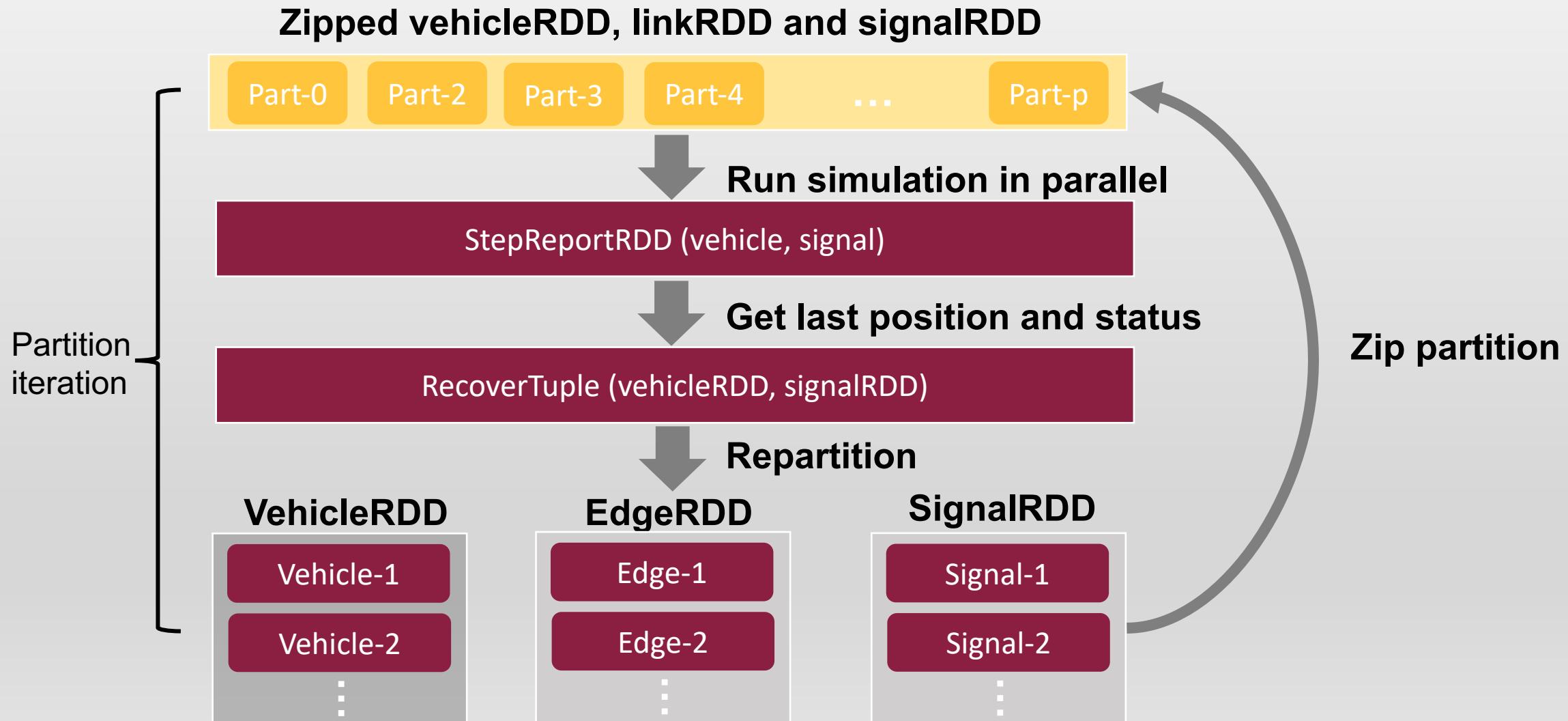


# Conclusion

- Scalable – **100X** faster > single-machine solution, **20%** faster > multiple-machine solution
- Comprehensive – **Easy** to use
- Extensible – Driving model **abstraction**
- Large-scale data analysis – with other GeoSpark components

**DEMO**

# Simulation-aware Partition



# Parallel Microscopic Simulation

**Input:** TrafficLights, Link, Vehicle

1. Initialize light and vehicle
2. Run simulation steps

- Loop simulation steps
  - Update traffic lights
  - Loop vehicles
    - If vehicle not arrive
      - Vehicle head way check
      - Move
    - If arrive, reborn vehicle
- End, collect step reports



Step: N	Time Complexity	Spatial Partition: P
Link: E		
Vehicle: V		
Signal: S	$N(S+CV)$	$\frac{N(S+VC)}{P}$
Driving: C		

**Parallel computing part i from 0 → p**

**N(S+VC)**

- Get closest vehicle and traffic light
- Update velocity (IDM)
- Remove vehicle from lane
- Compute position by velocity and update vehicle info
- Born vehicle
  - Check lane change (MOBIL)
  - Compute coordinate by position and lane
  - Update vehicle to map