



Scaling *Mobile Millennium* with BDAS

Timothy Hunter,
Teodor Moldovan, Matei Zaharia, Samy Merzgui, Justin Ma,
Michael J. Franklin, Pieter Abbeel, Alexandre M. Bayen

UC Berkeley

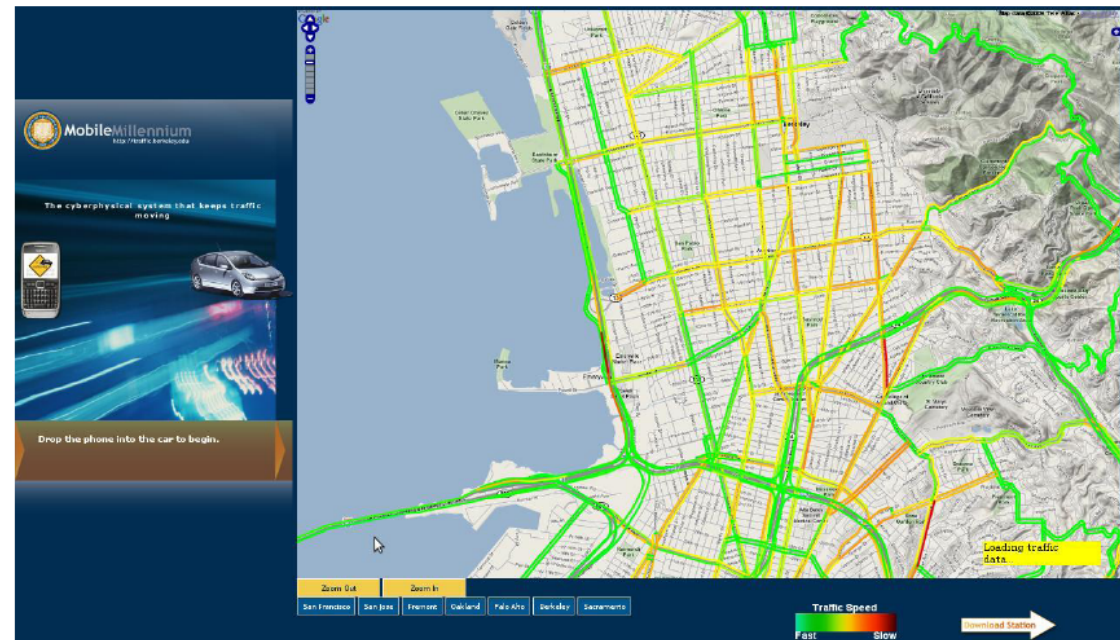


August 22, 2012

- Combining A, M and P in a real application:
 - Complex models (car traffic estimation)
 - Crowd-sourced data (mobile phones)
 - Computations on the cloud
- How we run Spark inside *Mobile Millennium*

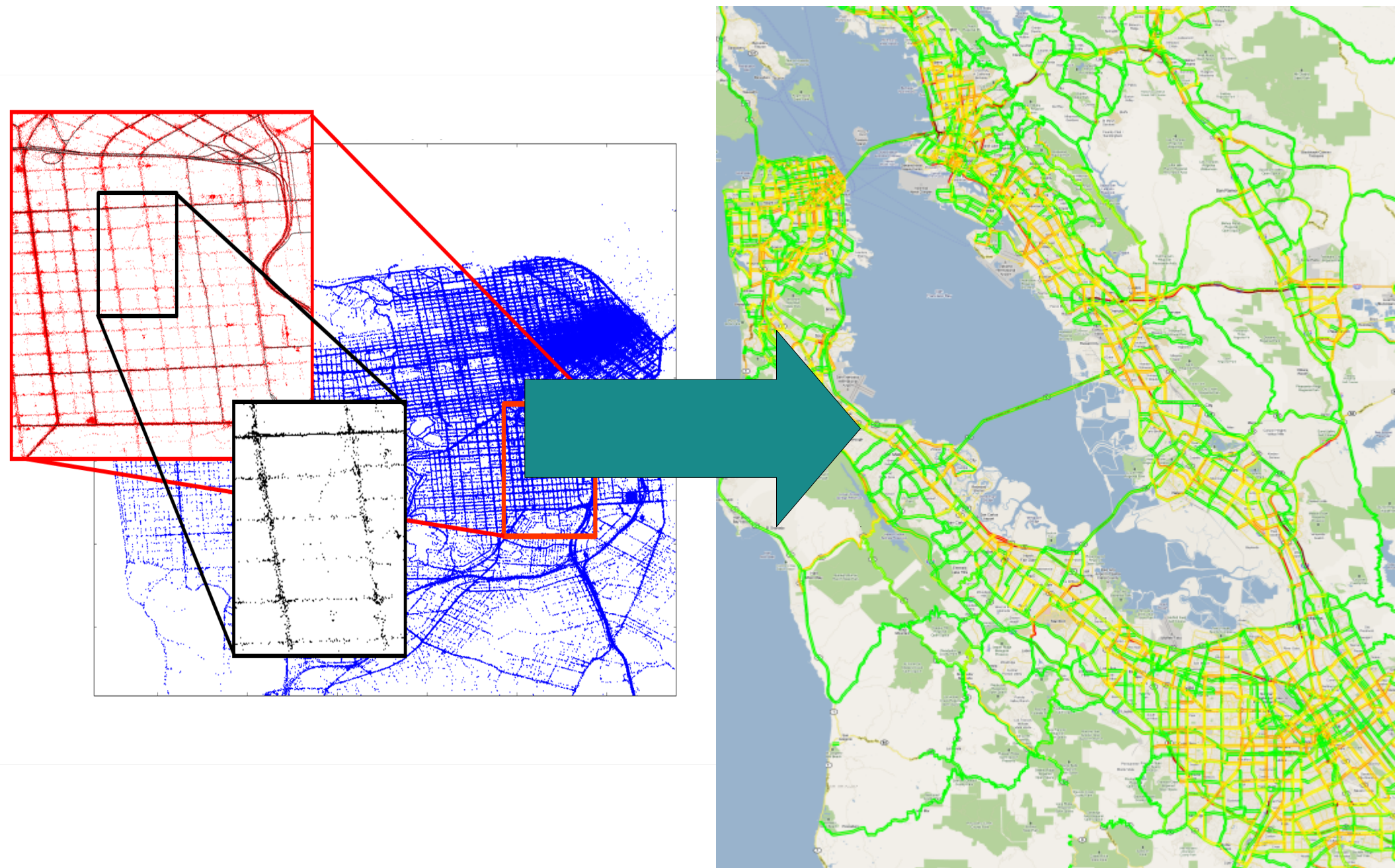
- Why *car* traffic estimation
- Overview of *Mobile Millennium*
- 2 minutes of applied Machine Learning
- Programming with the Spark framework
- Conclusion: the good, the bad, the not so beautiful

- Traffic congestion affects everyone
- Up-to-date estimation is critical
- Complex for urban streets (*arterial roads*)





- Input: sampled position of taxicabs
- Observed *every minute*

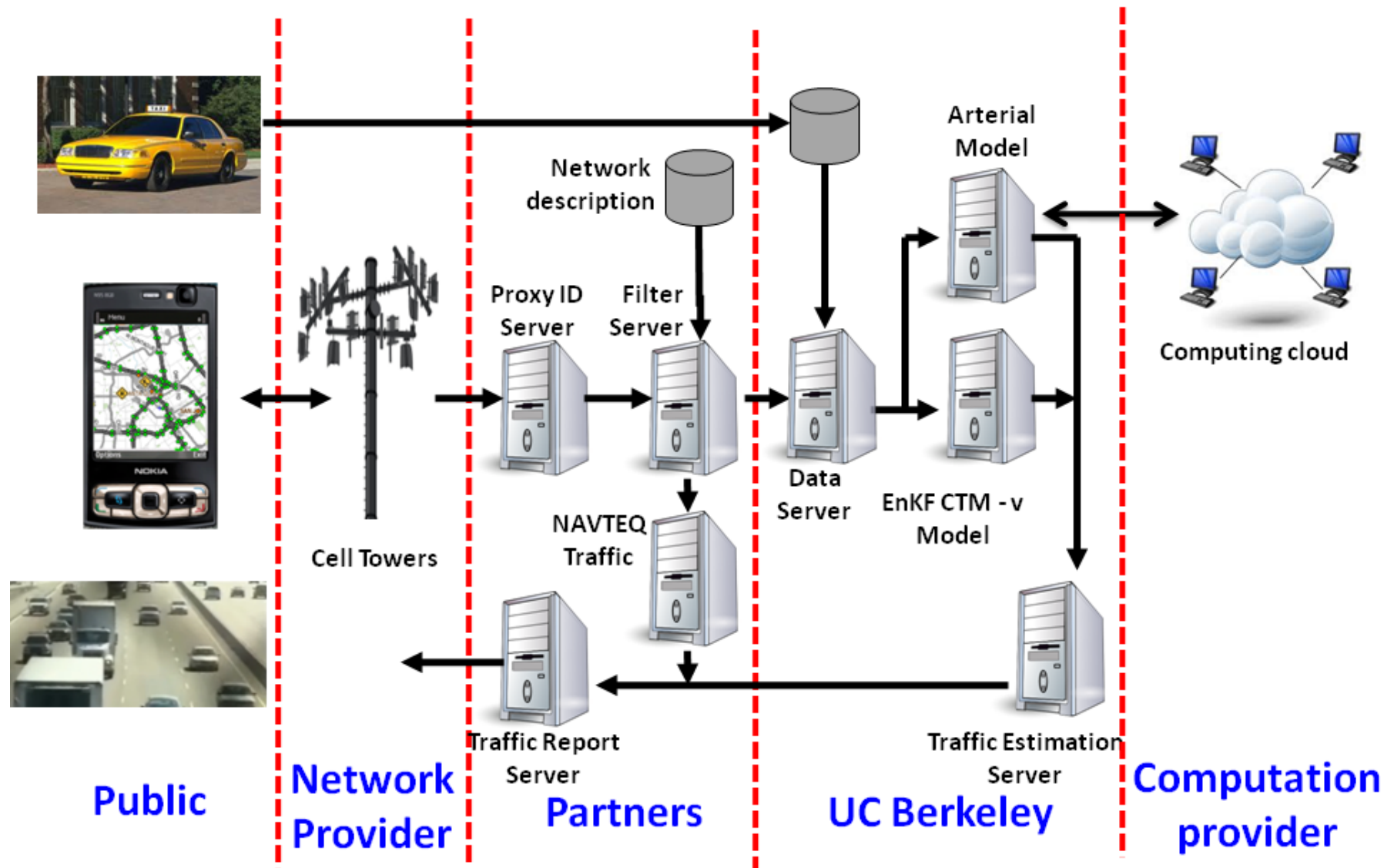




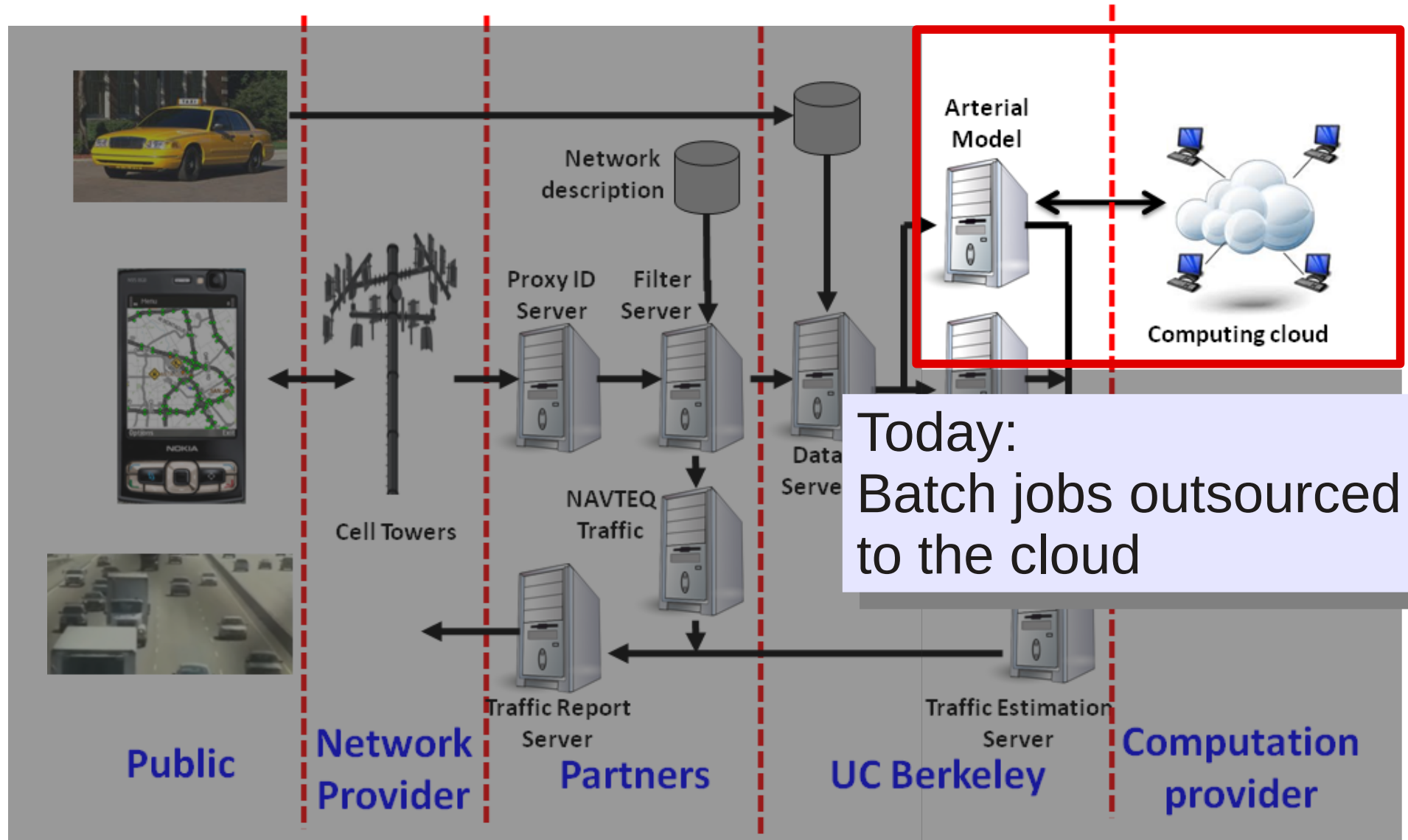
Preprocessing:

- Recovering trajectories from GPS points

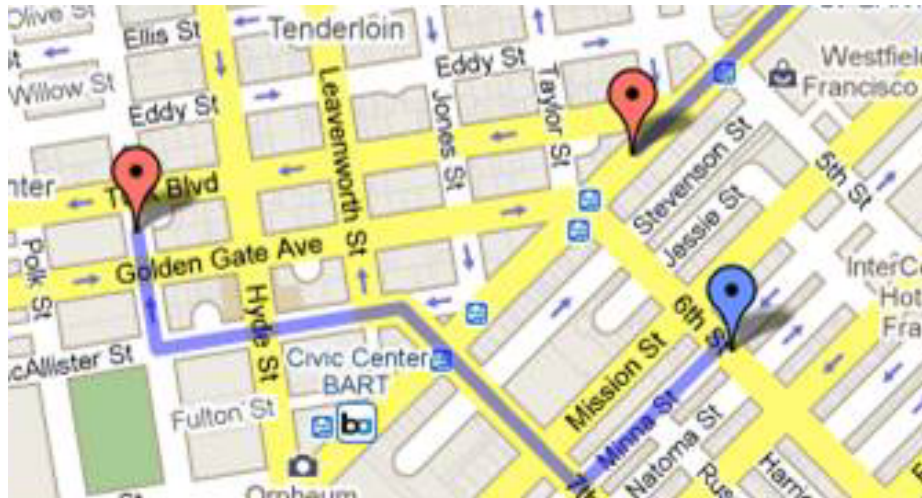
- A cyberphysical system for participatory sensing



- A cyberphysical system for participatory sensing

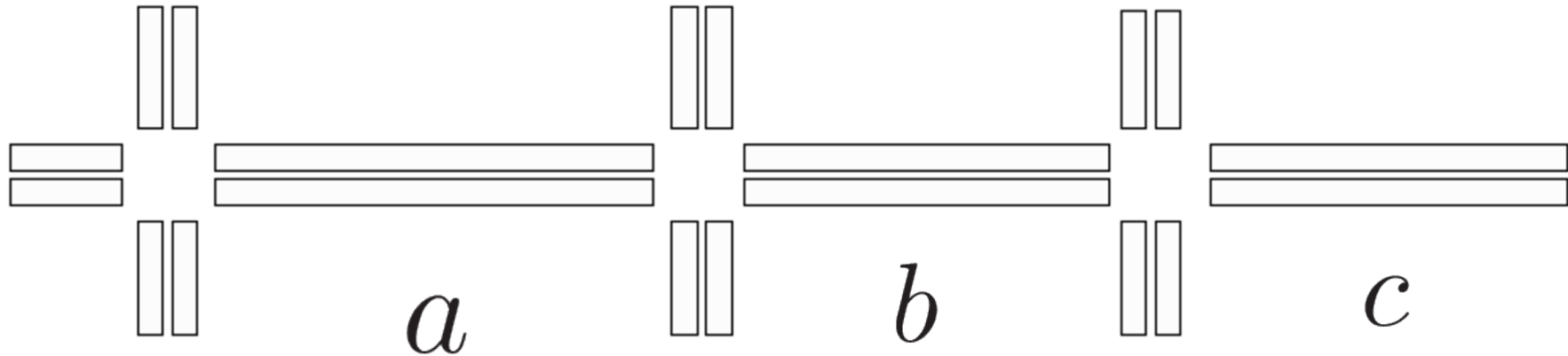


- Input:
 - Pieces of trajectories between GPS points

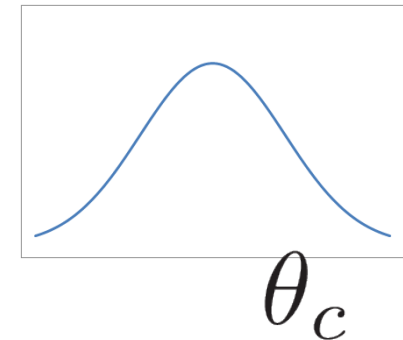
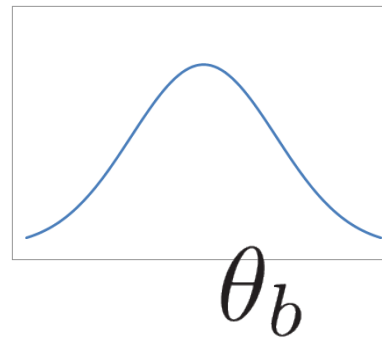
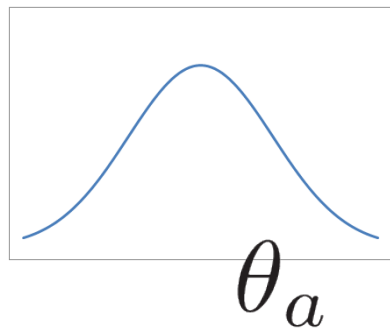


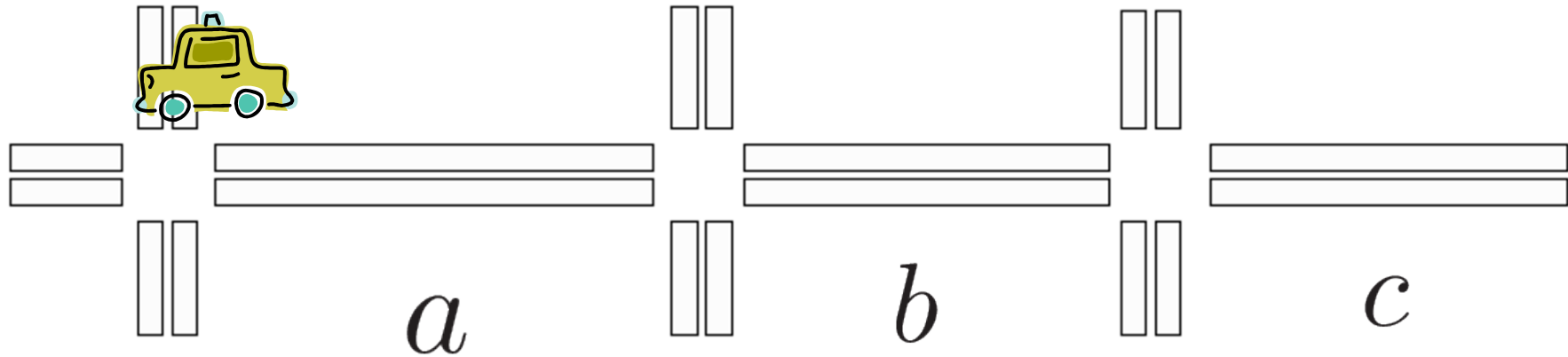
- Output: probability distributions of travel time
 - For each link
 - Parametrized by vector θ (mean and variance of link travel time)

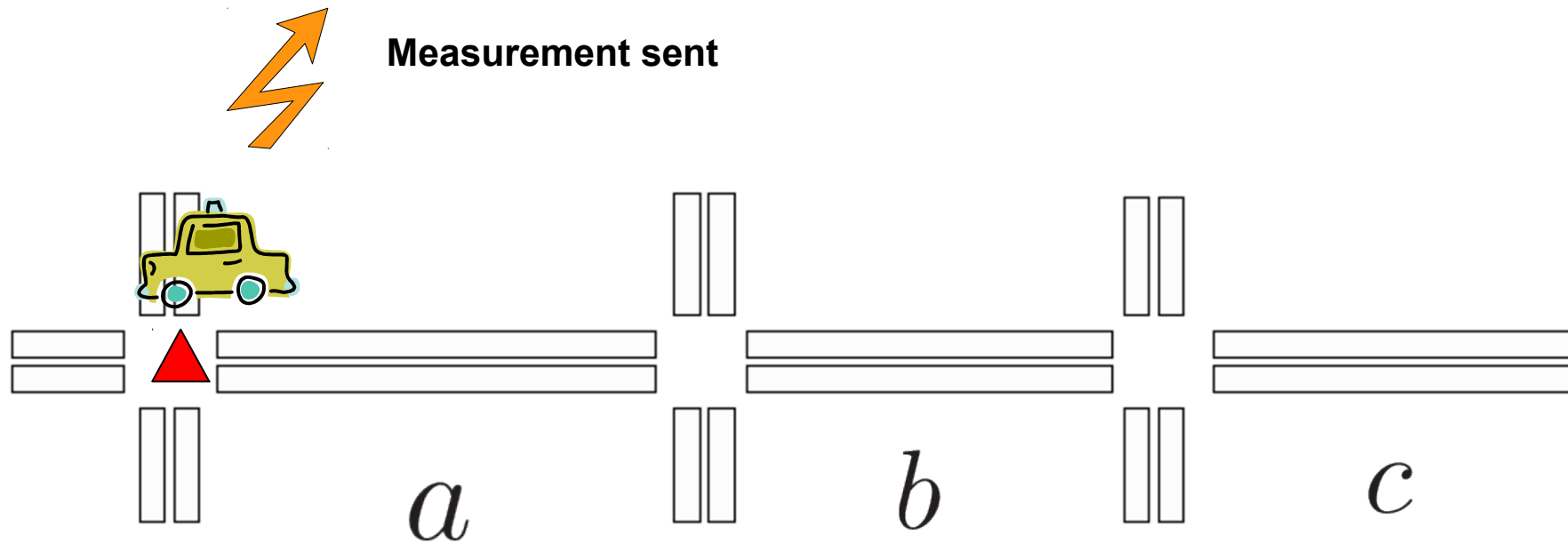
- Example road network

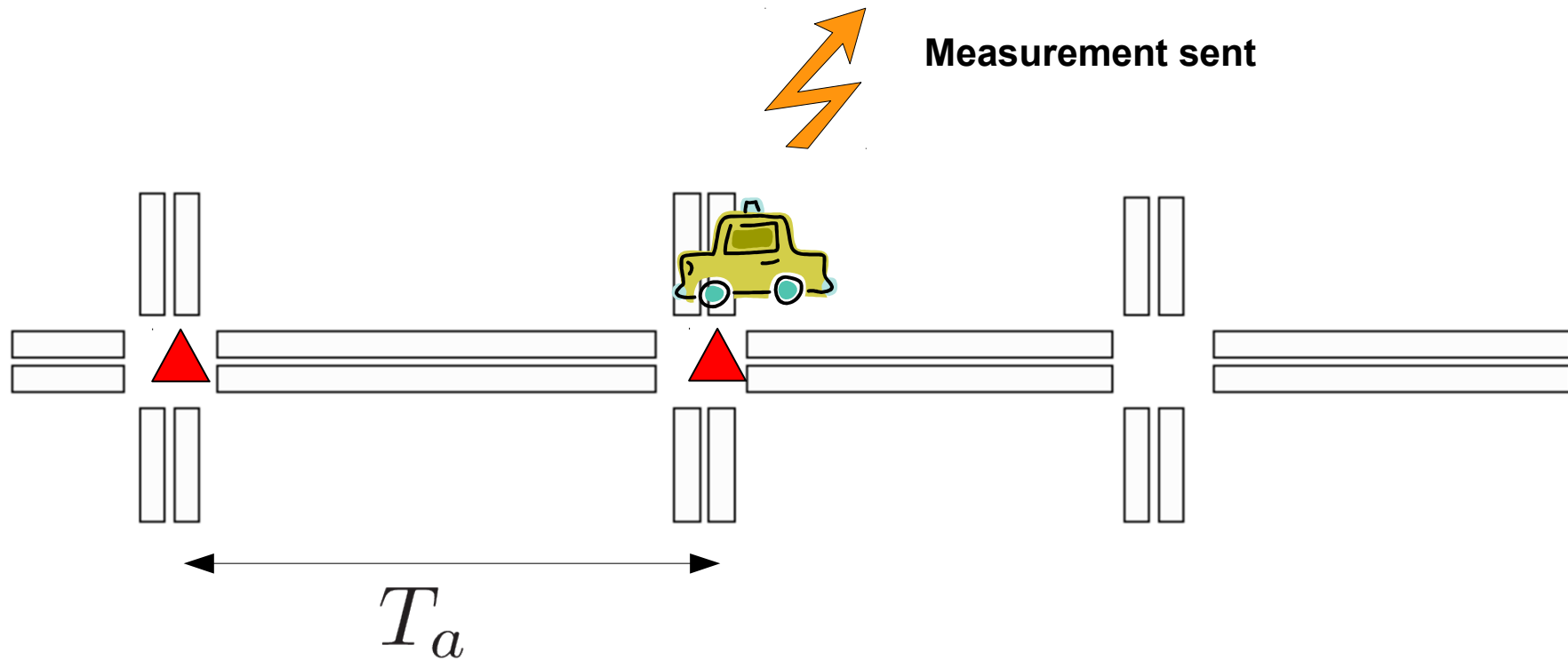


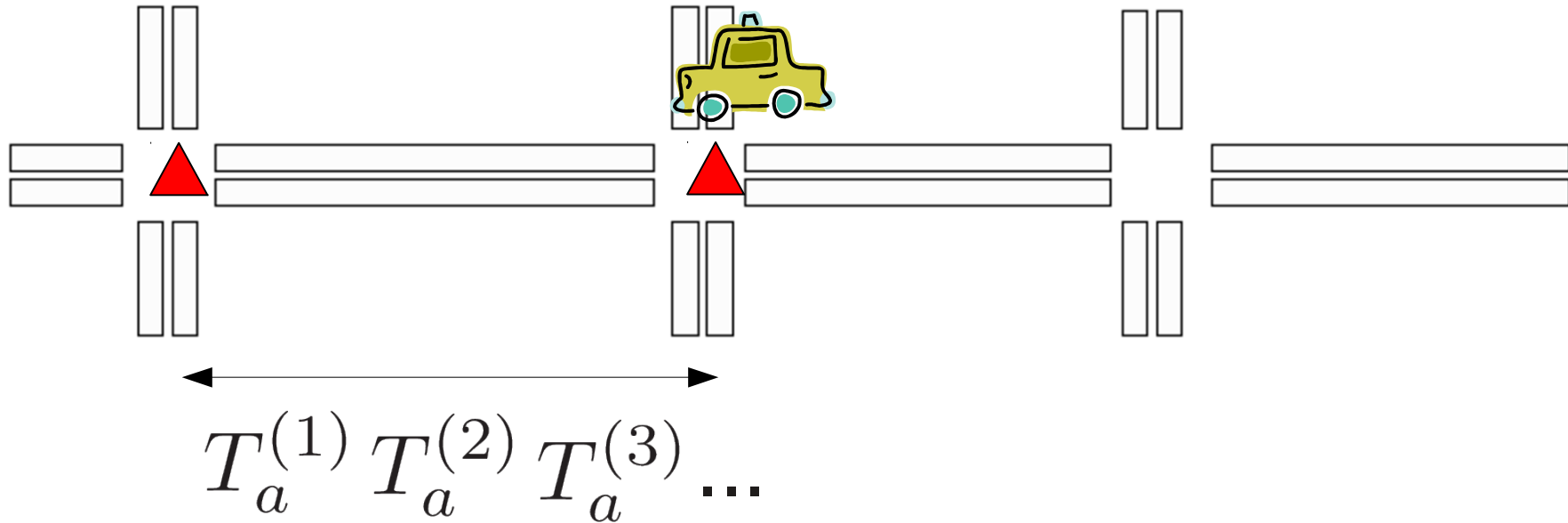
- Associated link travel times:

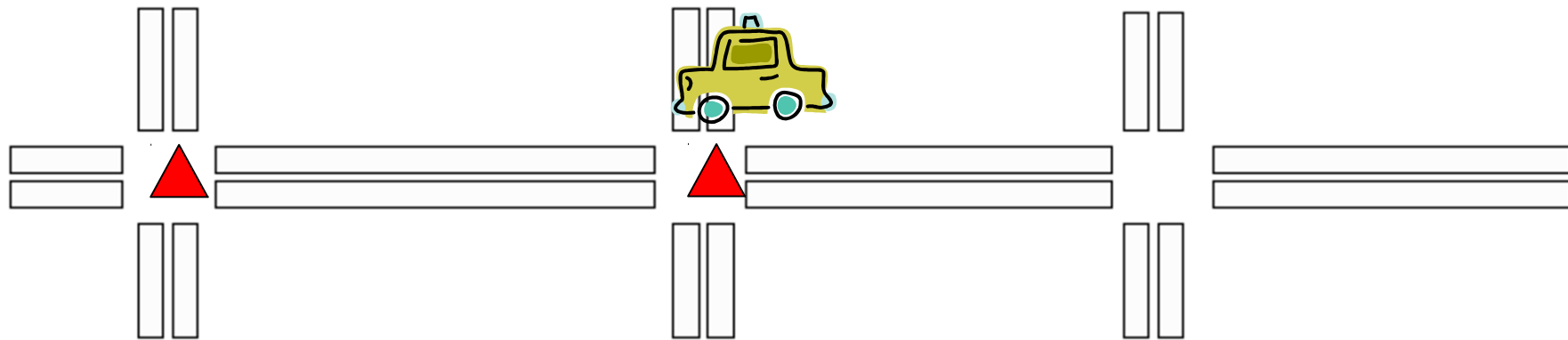




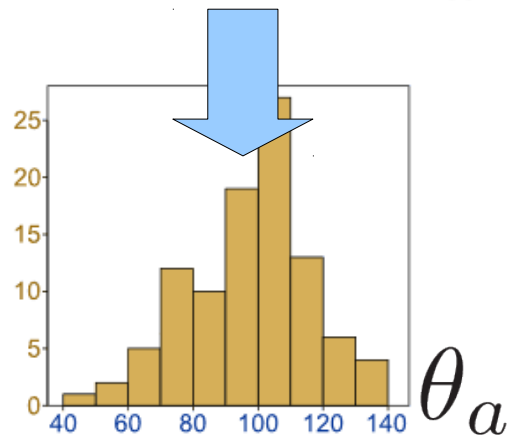


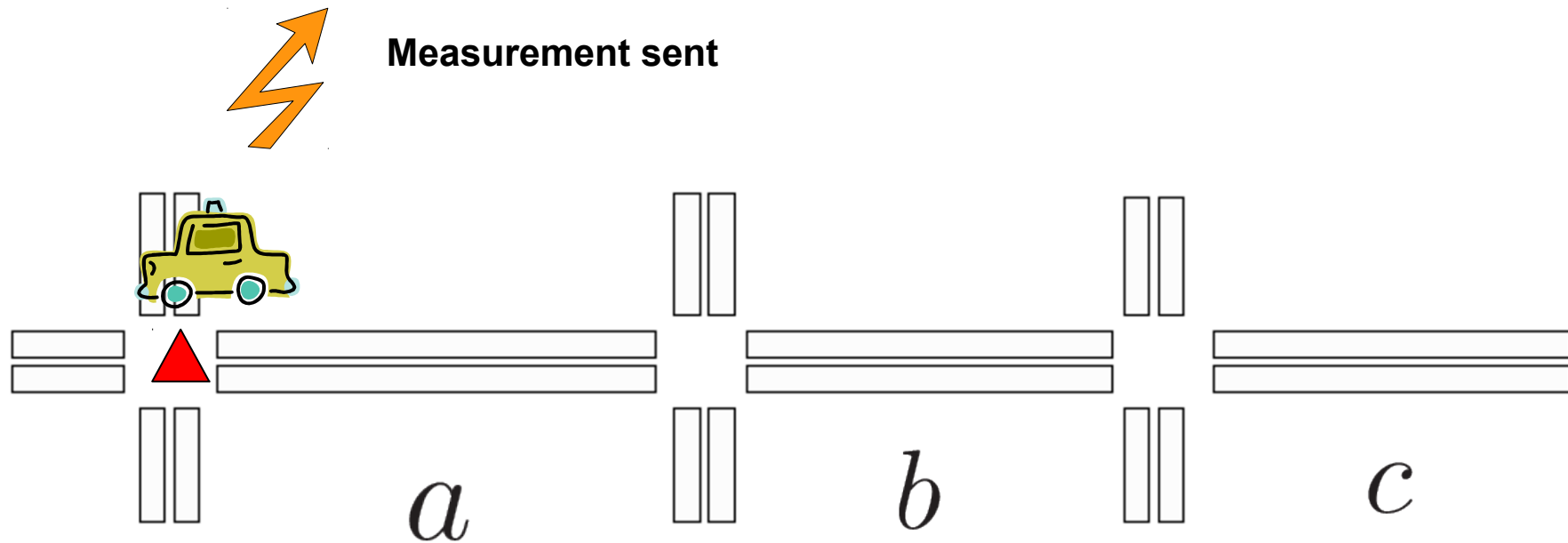






$$T_a^{(1)} T_a^{(2)} T_a^{(3)} \dots$$

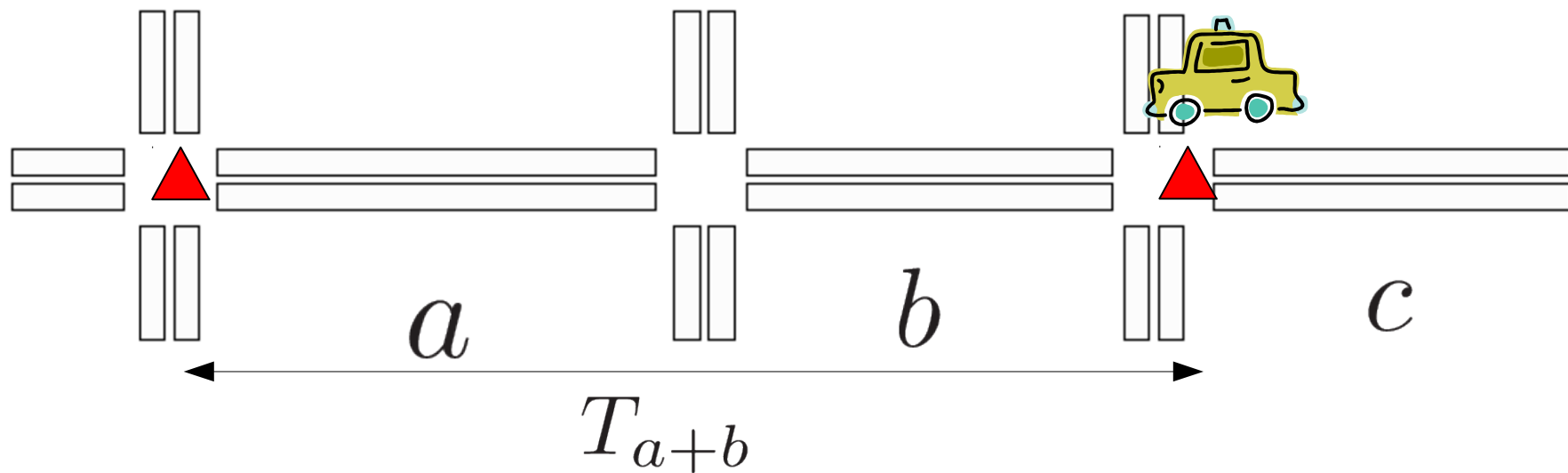




- Long time between observations



Measurement sent

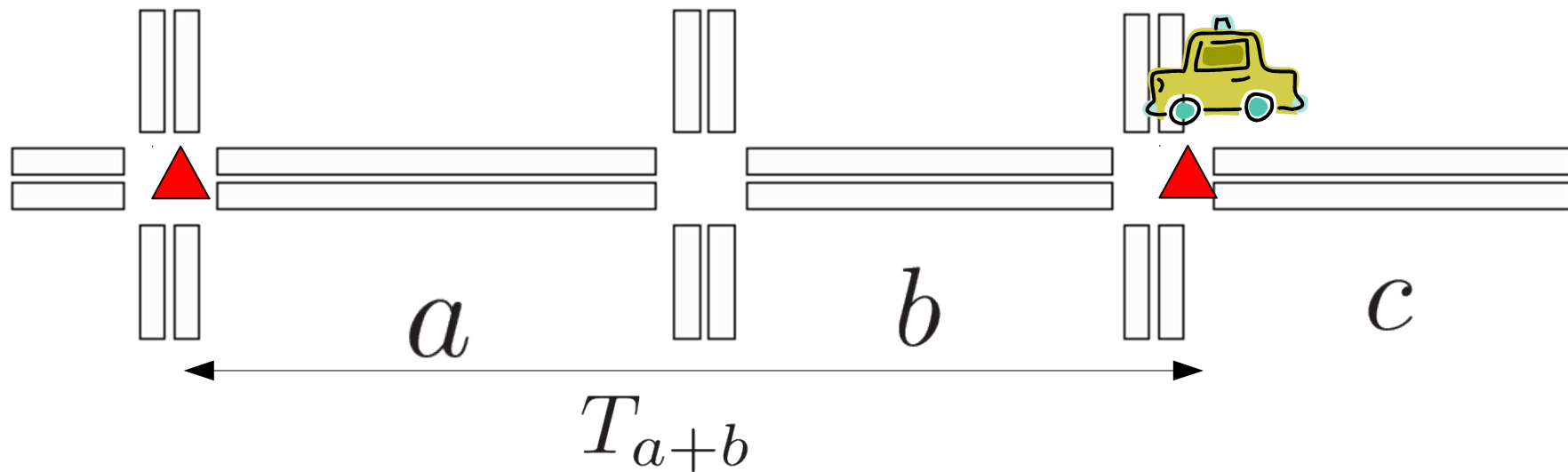


$T_a?$ $T_b?$

- Long time between observations



Measurement sent



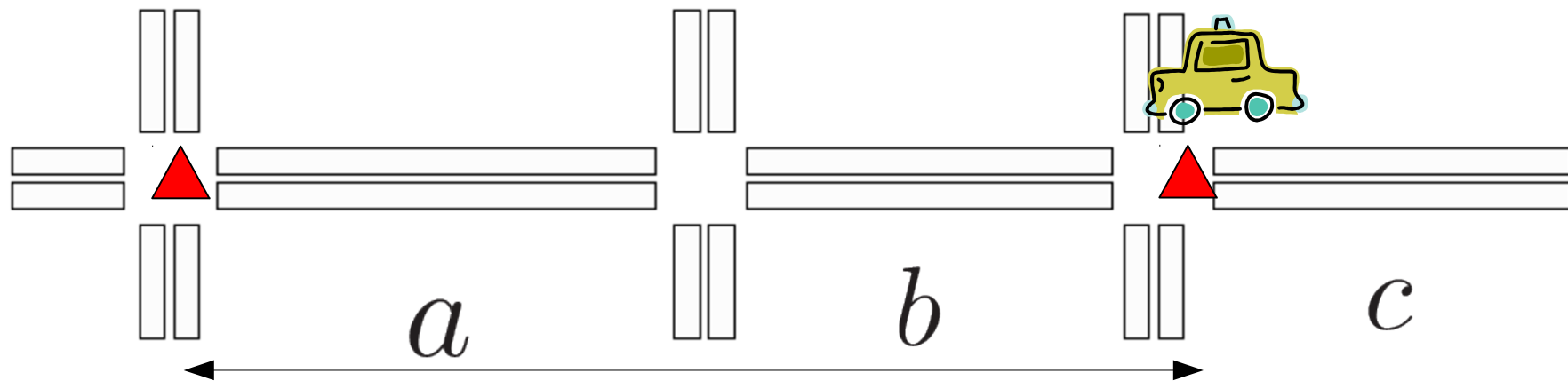
$T_a?$ $T_b?$

- Solution: sample!

- Long time between observations



Measurement sent

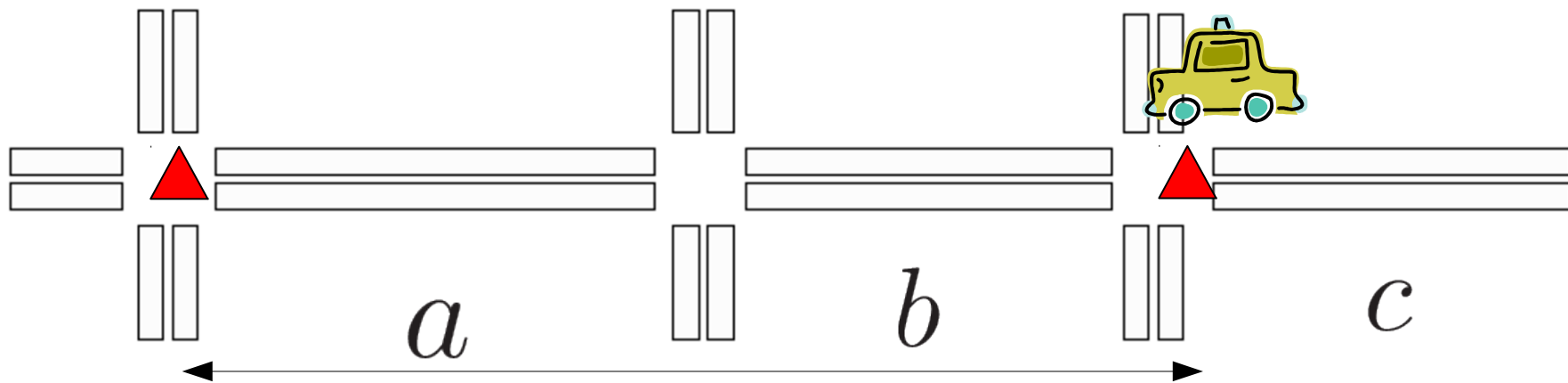


$$t_a^{(1)}, t_a^{(2)}, t_a^{(3)}, \dots \quad T_{a+b} \quad t_b^{(1)}, t_b^{(2)}, t_b^{(3)}, \dots$$

- Long time between observations



Measurement sent



$$T_{a+b}^{(1)}, T_{a+b}^{(2)}, \dots$$

$$t_a^{(1,1)}, t_a^{(1,2)}, t_a^{(1,3)}, \dots$$

$$t_b^{(1,1)}, t_b^{(1,2)}, t_b^{(1,3)}, \dots$$

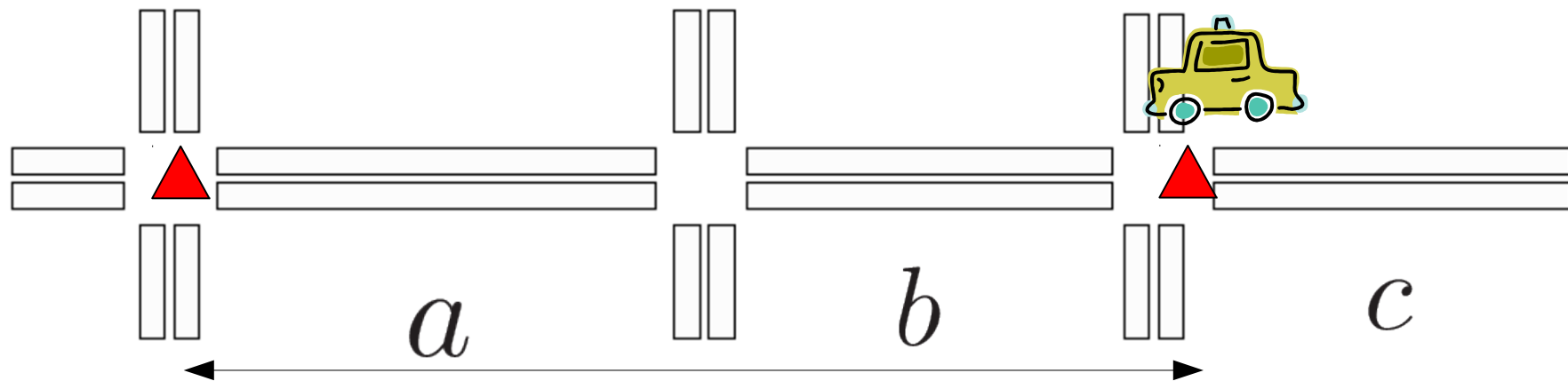
$$t_a^{(2,1)}, t_a^{(2,2)}, t_a^{(2,3)}, \dots$$

$$t_b^{(2,1)}, t_b^{(2,2)}, t_b^{(2,3)}, \dots$$

- Long time between observations



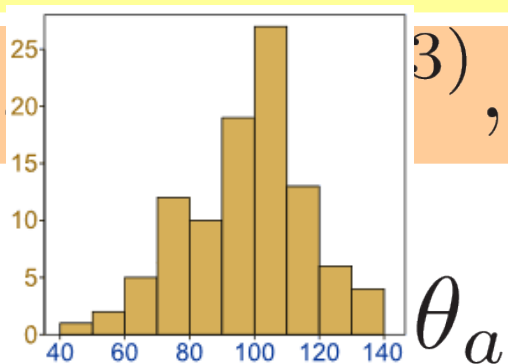
Measurement sent



$$T_{a+b}^{(1)}, T_{a+b}^{(2)}, \dots$$

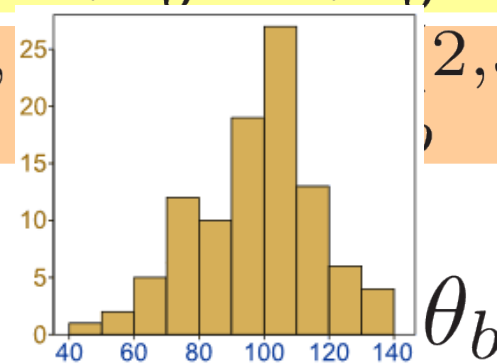
$$t_a^{(1,1)}, t_a^{(1,2)}, t_a^{(1,3)}, \dots$$

$$t_a^{(2,1)}, t_a^{(2,2)}, t_a^{(2,3)}, \dots$$

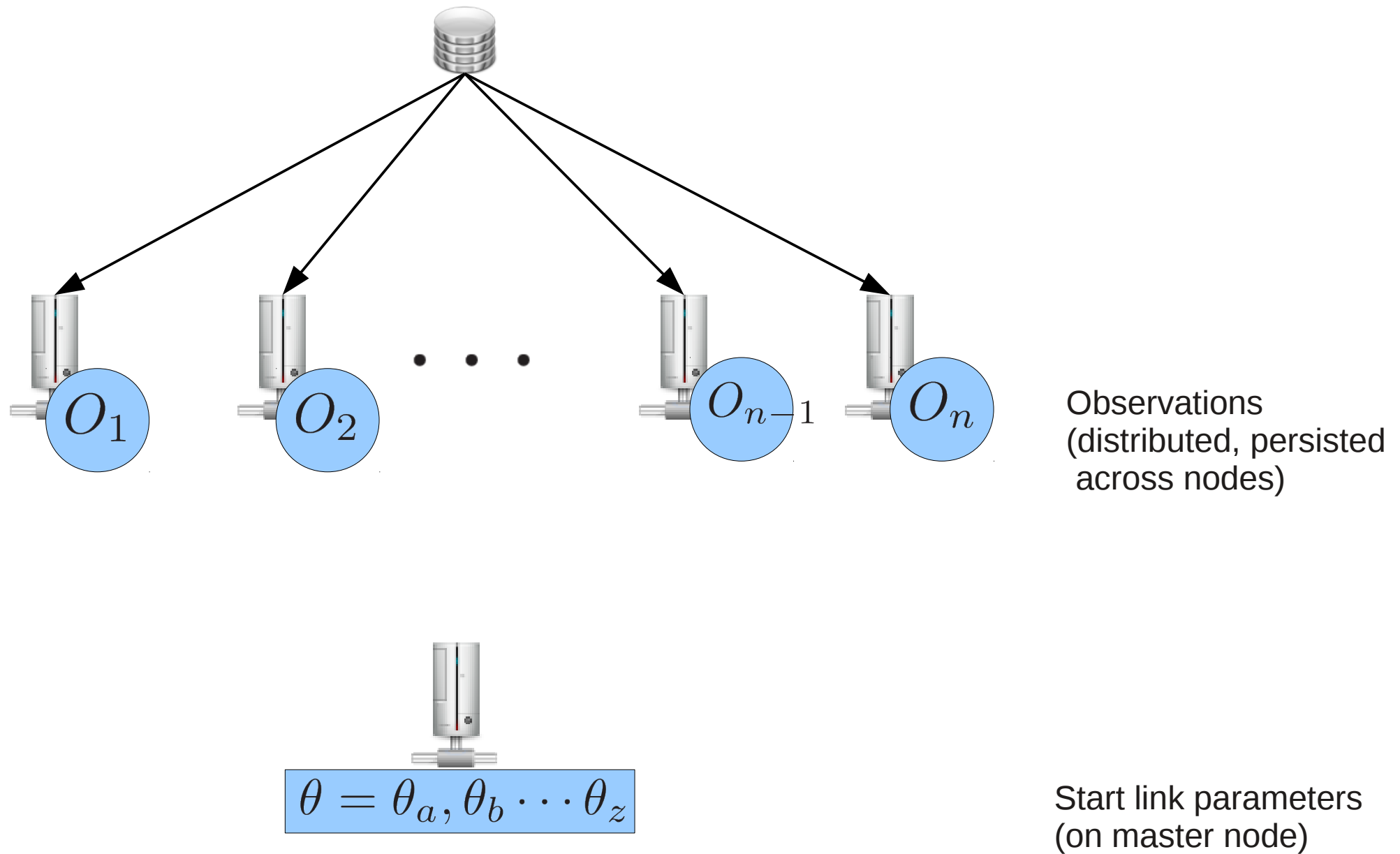


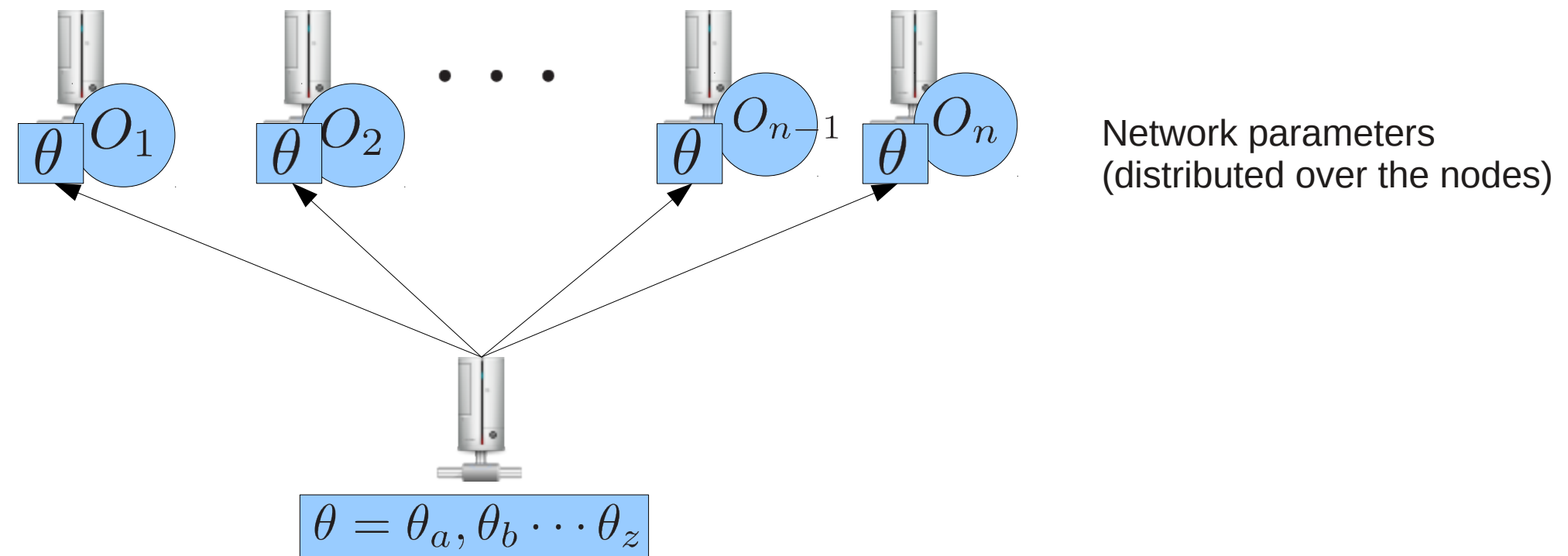
$$t_b^{(1,1)}, t_b^{(1,2)}, t_b^{(1,3)}, \dots$$

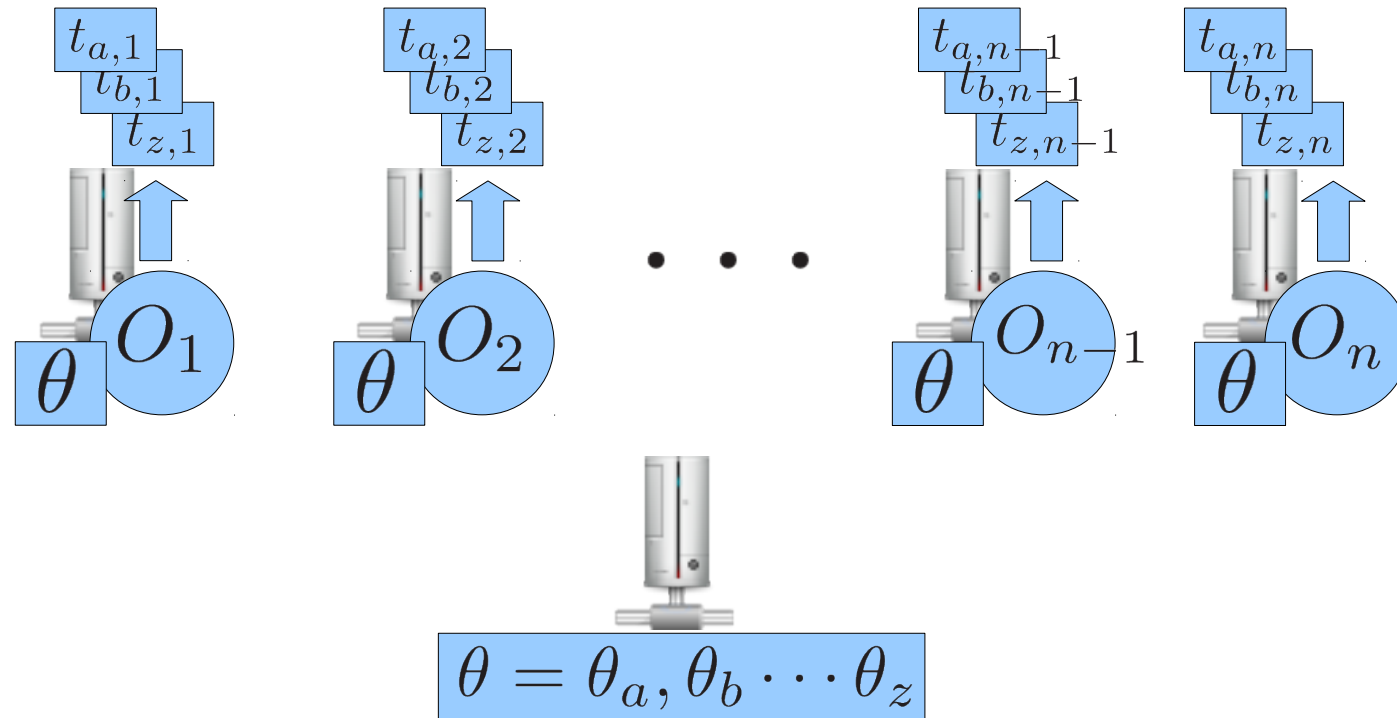
$$t_b^{(2,1)}, t_b^{(2,2)}, t_b^{(2,3)}, \dots$$



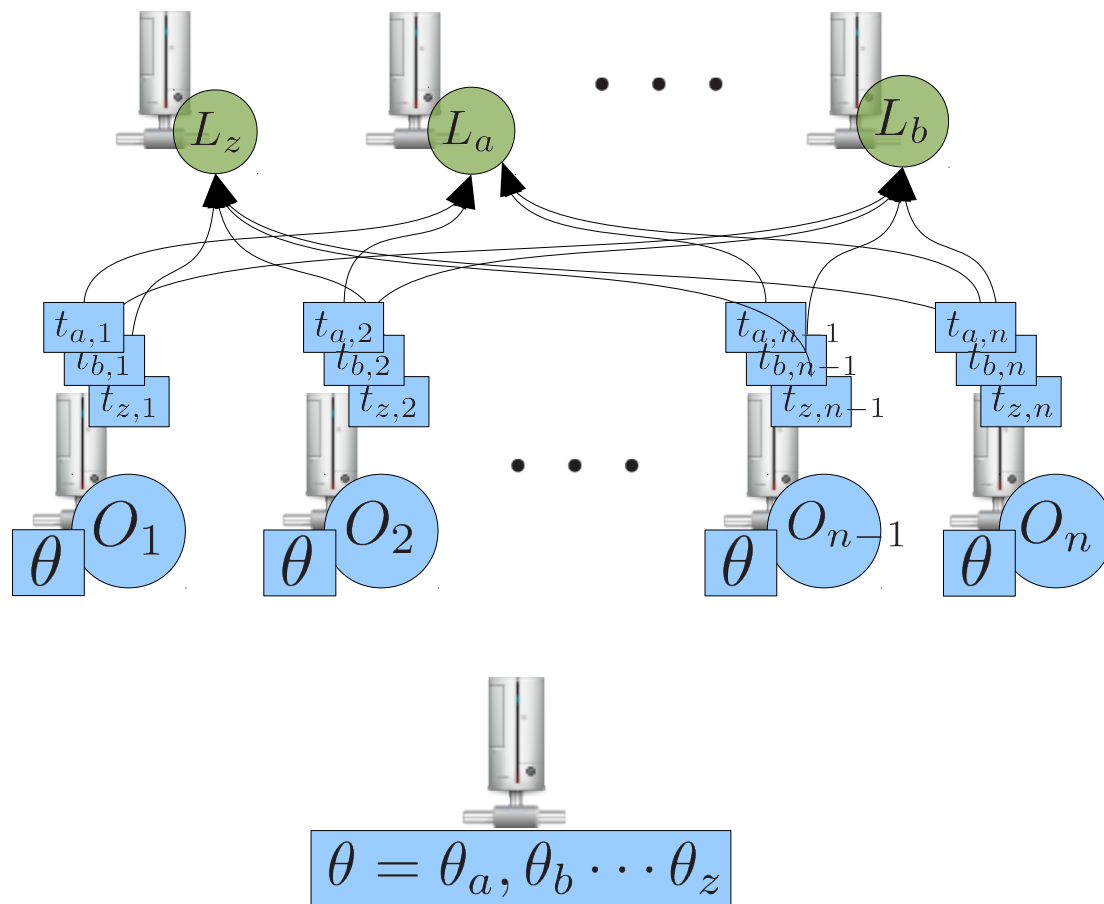
- Procedure called *Expectation Maximization*
- Iterative in nature:
 - Alternates between sampling (E step) and learning (M step)
- Some figures:
 - 50k road links (parameters)
 - 50M observations (15GB, avg. 4 links / observation)
 - 200M partial travel times
 - x1000 samples per partial travel times



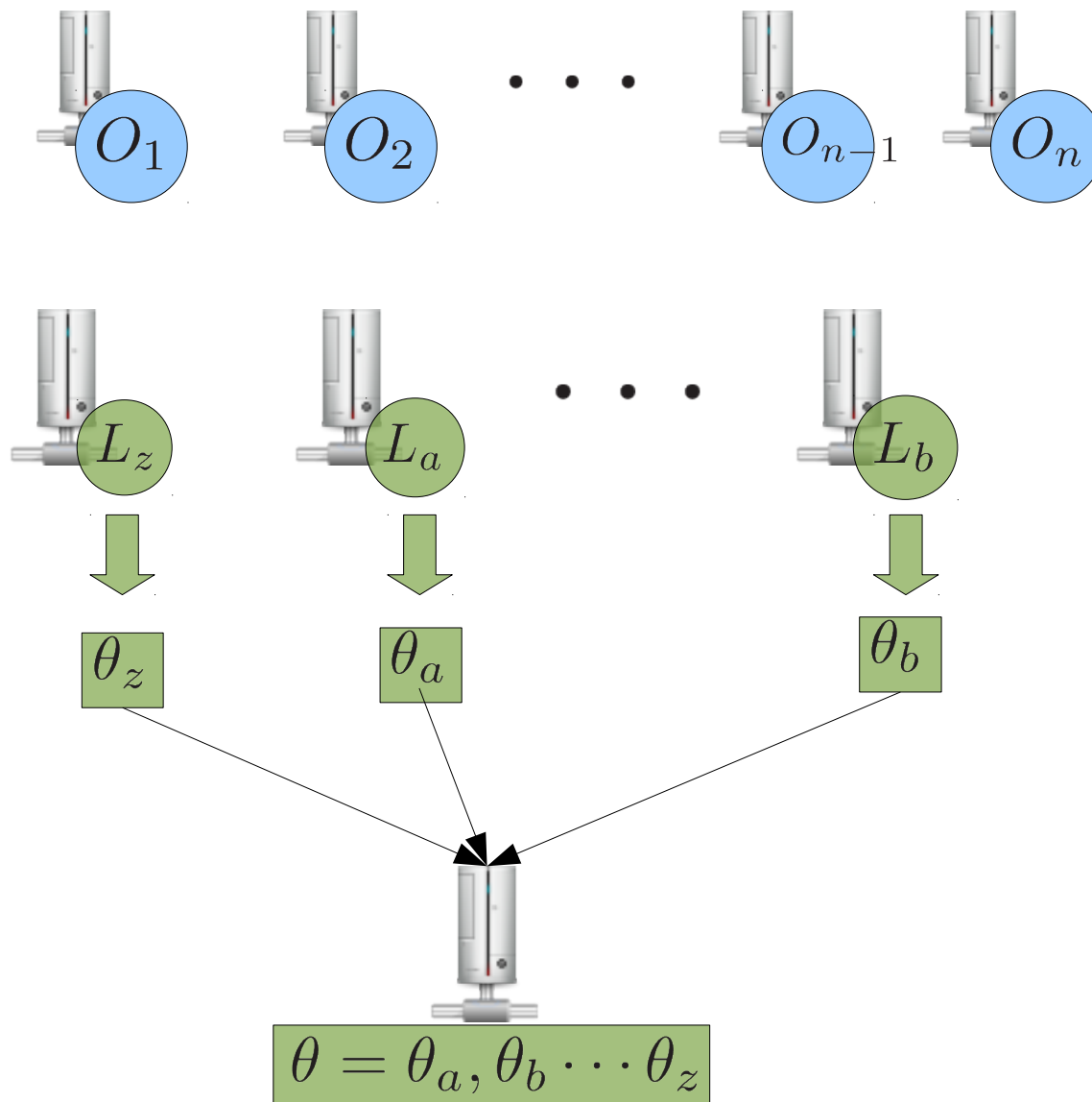




Travel time samples
For each observation link



Travel time samples aggregated
on a link basis



New parameters are generated
The maximize sampled travel times
for each link.

The master collects the vector of
new parameters.

Main loop of the program

```
val observations = spark.textFile("hdfs:...")  
  .map(parseObservation _)  
  .cache()  
var params = // Initialize models parameters  
while (!converged) {
```

Step 1 (E step)

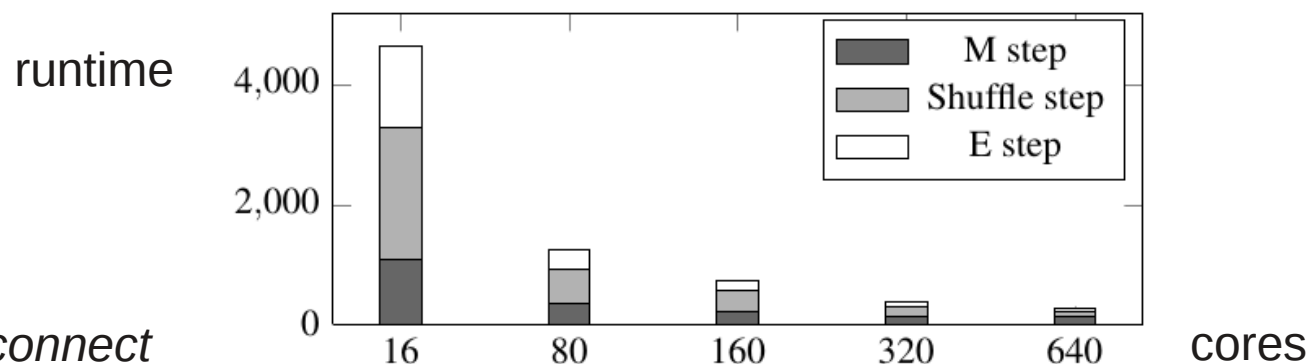
```
val samples = observations.flatMap( obs =>  
  generateSamples(obs, params))
```

Step 2 (M step)

```
params = samples.groupByKey(false).map(  
  case (linkId, vals) =>  
    mostLikelyParam(linkId, vals)  
) .collect()  
}
```

- Before using Spark:
 - 3.5x *slower* than real-time
 - Could not even handle all the data
- With Spark:
 - Similar programming interface (methods on scala collections)
 - Very good scalability (near linear)
 - Each iteration 3x faster than reloading from disk

NERSC cluster:
quad-core Xeon
4X QDR InfiniBand interconnect

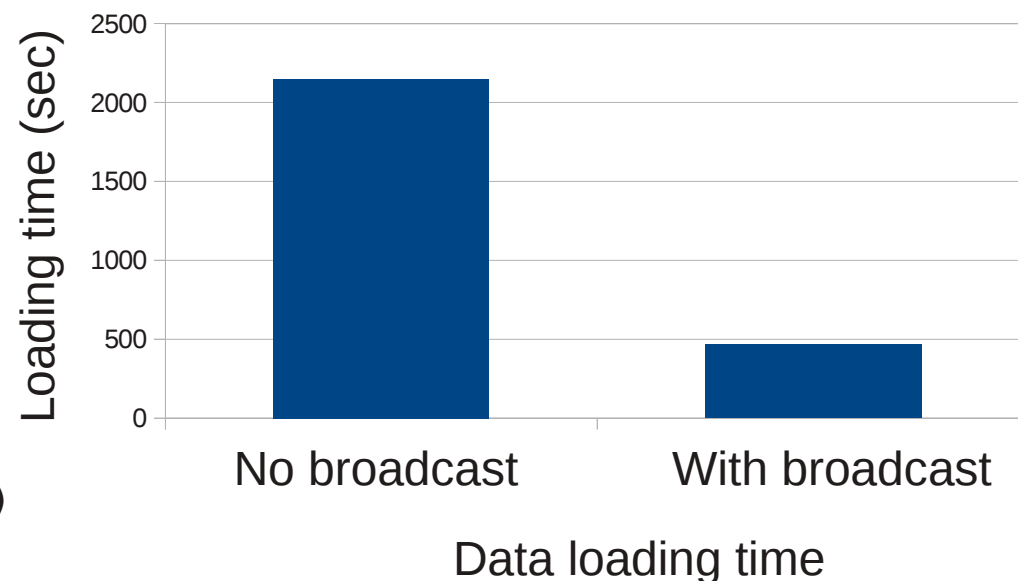


- The observation data is stored in memory:
 - Be careful with the memory footprint
 - Look at logs to monitor GC status
- We cache pointer-based structures
 - Significant overhead in the JVM
- Workaround: use compact collection structures (arrays) and make liberal use of `.toArray()`
- Workaround: RDDs of *serialized data*

- Need to share data between all workers:
 - At the start of the job (network description, > 40MB)
 - Between iterations (updated parameters θ)
- Using Spark's broadcast
- Data loading time reduced by 79%

```
val network = // load network
val observations = spark.textFile("...")
    .map(parseObservation(_, network))

val network = // load network
val bc_net = spark.broadcast(network)
val observations = spark.textFile("...")
    .map(parseObservation(_, bc_net.get()))
```



- An application of Spark:
 - Real-world ML problem
 - Crowd-sourced data
- Implementation now (much) faster than real time
- Not limited by computations:
 - We can use more complex ML tools than before

Thank you
