

# Big Data Analytics

FOSCA GIANNOTTI AND LUCA PAPPALARDO

---

[HTTP://DIDAWIKI.DI.UNIPI.IT/DOKU.PHP/BIGDATAANALYTICS/BDA/](http://didawiki.di.unipi.it/doku.php/bigdataanalytics/bda/)

**DIPARTIMENTO DI INFORMATICA - Università di Pisa**  
**anno accademico 2018/2019**

# Mobility Data Mining

---

CITY DYNAMICS WITH GSM DATA

A solid orange horizontal bar spanning the width of the slide at the bottom.

# Contents

---

- Corporate Users
  - Geomarketing
  - Monitoring Driving-based Segmentation
- Individual Users
  - Self-awareness
  - Proactive Carpooling
- Public Sector
  - Urban Mobility Atlas
  - Borders

---

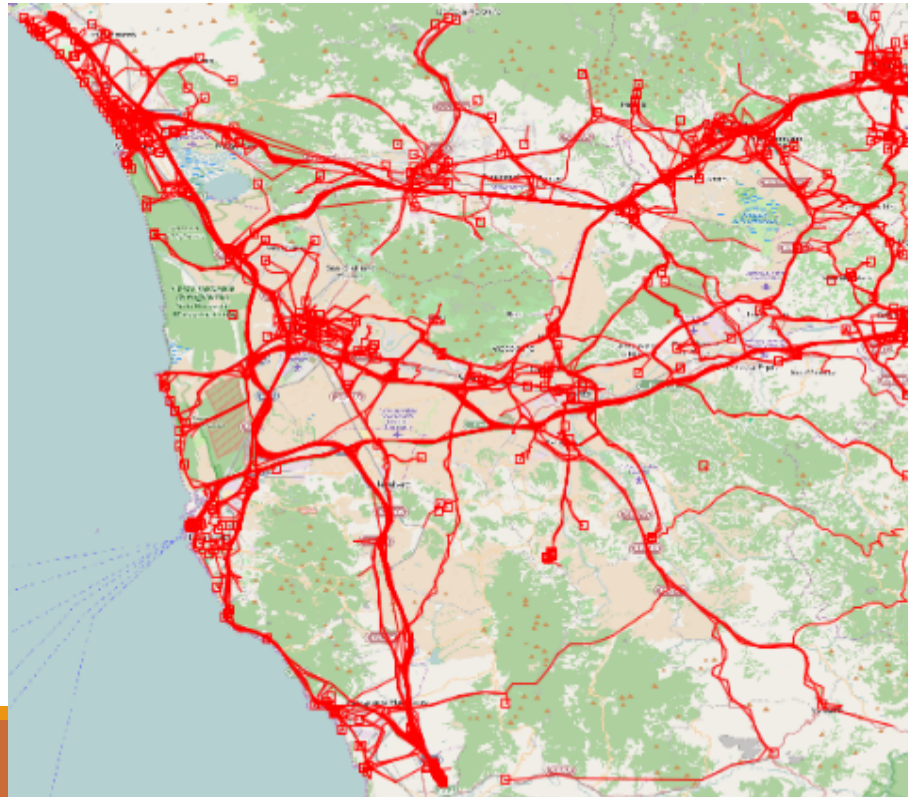
# **Services Towards Corporate Users**

## ***Geomarketing***



# Problem definition

Based on the trajectories of a sample of population, what is the best place to open a new shop / mall ?



# The “best” place

---

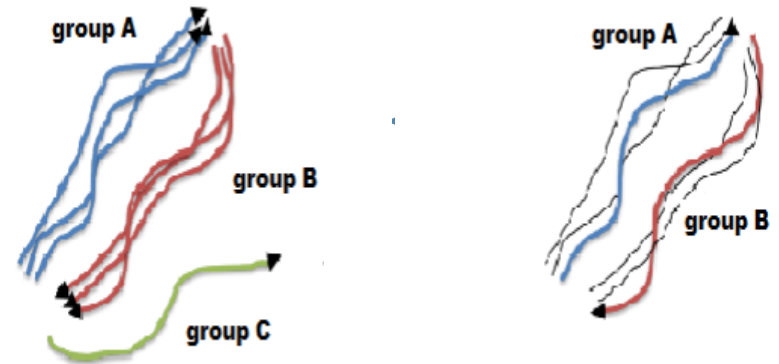
Experts' knowledge: best place to open a mall is where people pass during everyday activities

Area crossed by road segments with a high frequency of systematic travels of people

# Systematic movements

## Step 1: Map-matching

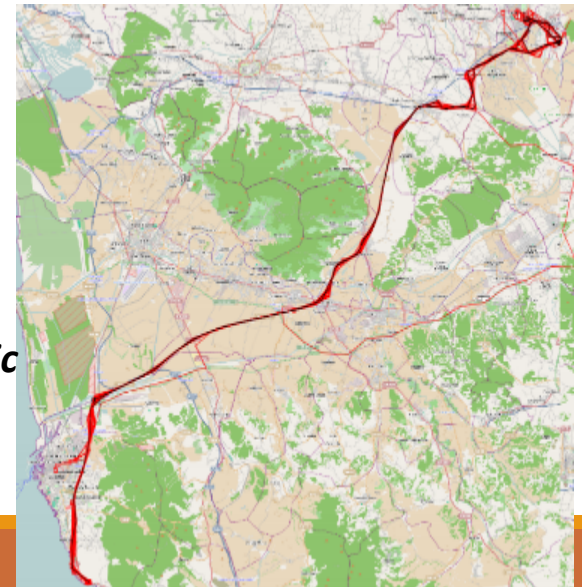
- See users' movements as sequences of road segments.



## Step 2: Mobility profiles

- Select only systematic movements.

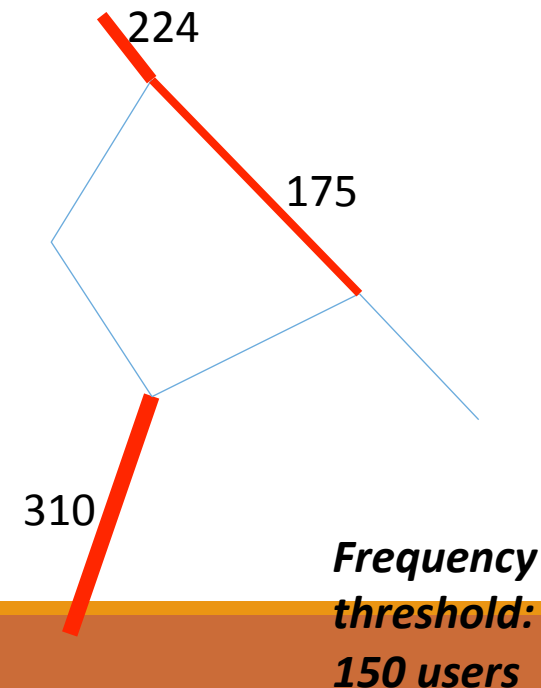
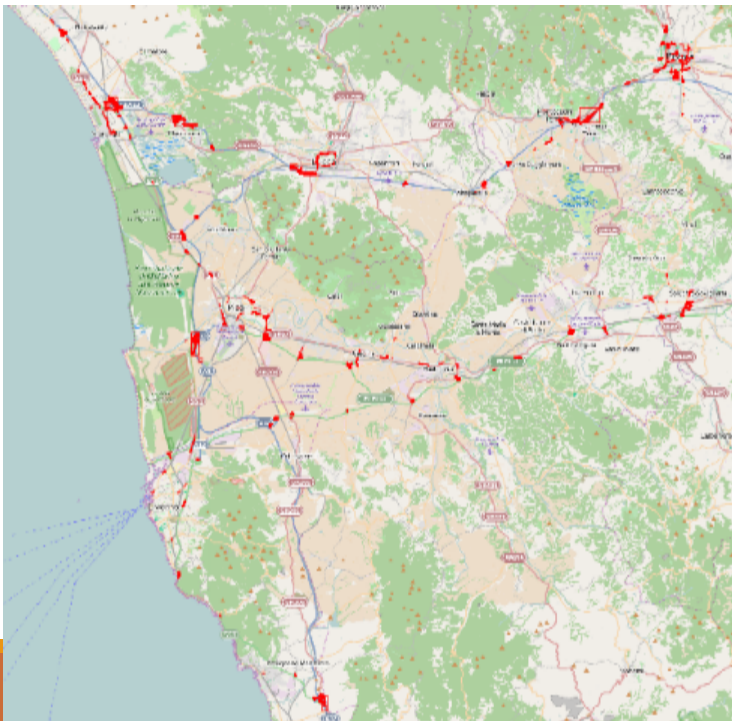
*User's systematic movement:*  
 $L1 \rightarrow L2$



# Frequently visited road segments

---

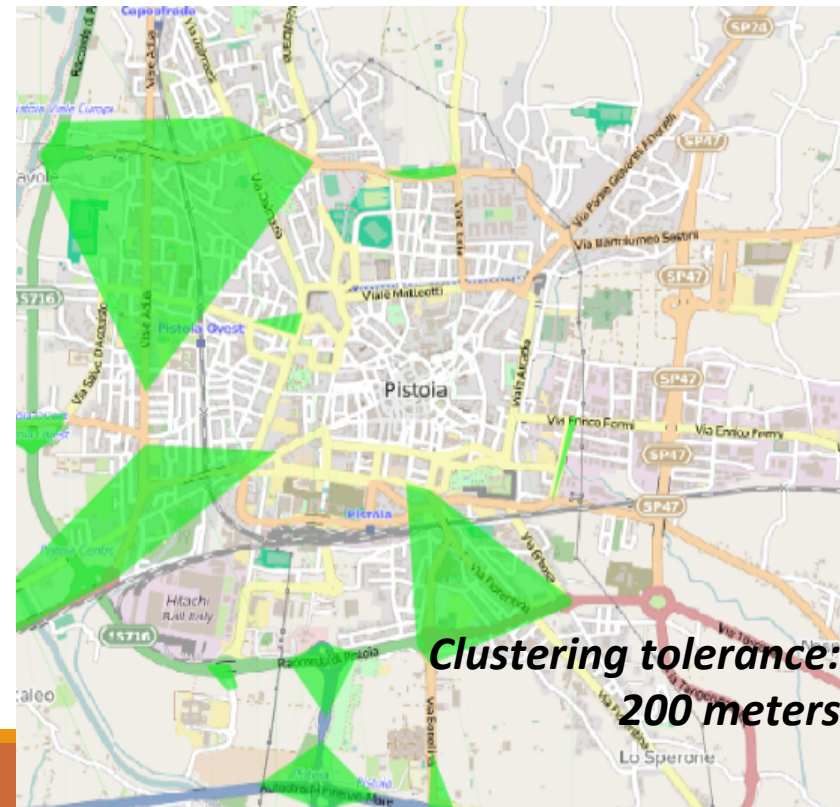
- Aggregate systematic movements by road segments
- Set a threshold to select the frequent ones



# Candidate areas for a mall

Using a spatial clustering we can extract cluster of frequent road segments which are spatially close each other.

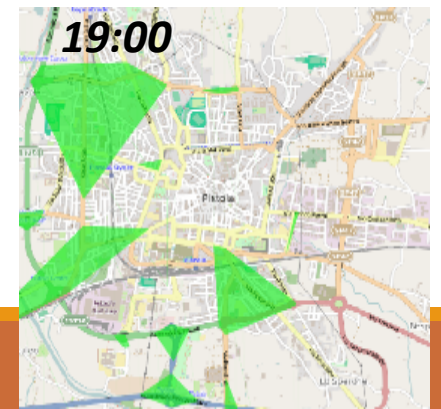
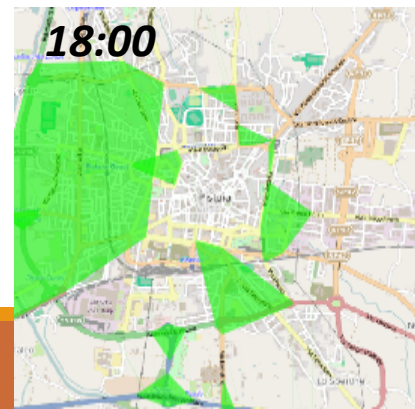
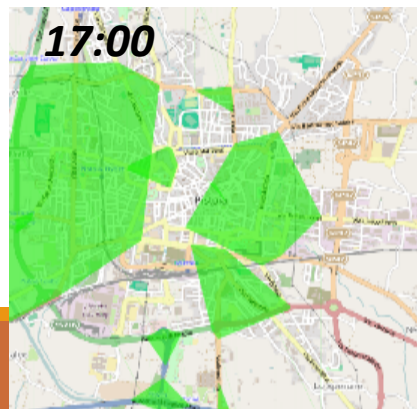
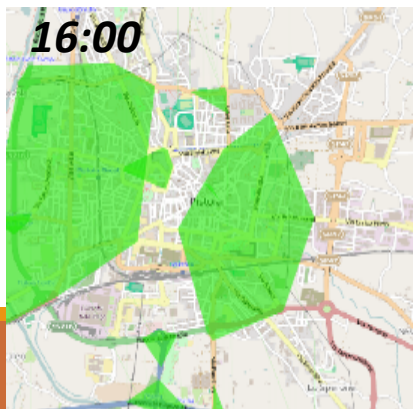
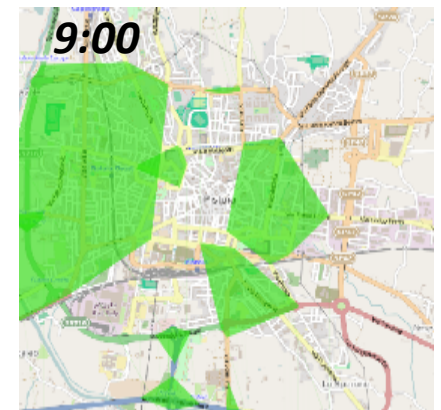
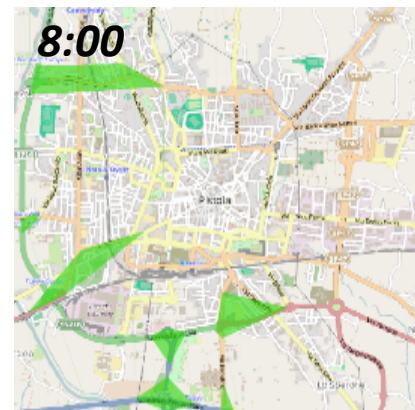
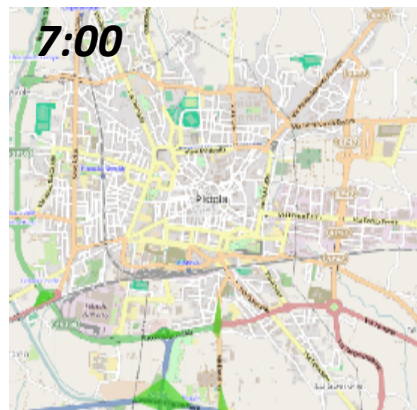
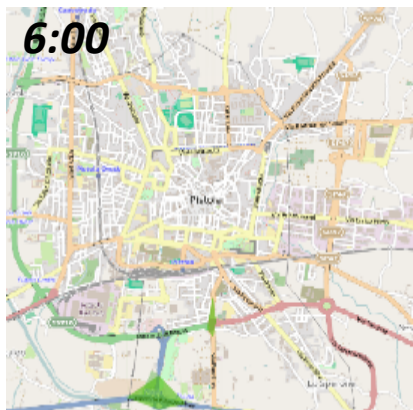
- Distance of 2 segments
  - Compare vertices
- Draw clusters as convex hull





# Temporal evolution

Repeat this process for each hour of the day and analyze how they evolve



---

# **Services Towards Corporate Users**

***Monitoring Driving-based Segmentation***

# Segmentation and monitoring

---

- Mobility application scenario of the LIFT European project



USING LOCAL INFERENCE  
IN MASSIVELY DISTRIBUTED SYSTEMS



- Focused on distributed monitoring technologies



# Scenario context & motivation

- **Customer segmentation:** a marketing strategy that involves dividing a broad target market into subsets of consumers who have common needs

*[http://en.wikipedia.org/wiki/Customer\\_segmentation](http://en.wikipedia.org/wiki/Customer_segmentation)*

- **Needs:** car insurance companies would like to define customer segments that capture different driving profiles
  - Each segment could then be offered suitable contract conditions
- **Opportunities:** the vehicles insured by some companies have on-board GPS devices that can trace their movements
  - They could aggregate such traces into driving habit indicators based on recent history for the driver and transmit them



*dreamstime.com*

# Scenario description

- 
- Driving indicators
    - **Each vehicle** continuously keeps track of recent movements, compute aggregate indicators and sends them to controller
  - Profile extraction
    - **The controller** uses initial indicator values to build clusters of drivers, each corresponding to a “driving profile”
  - Profile monitoring
    - **The controller** continuously checks updates to verify that the driving profiles extracted are still good enough

# Step 1: Features for individual mobility behaviors

---

- Indicators for recent mobility behaviors
- Computed over recent history → sliding window



- Include information derivable from standard GPS devices

# Step 1: Features for individual mobility behaviors

- Which features?

- Superset of those currently used by insurance companies

How fast I drive  
w.r.t. speed limits

Where I drive  
w.r.t. road categories

How dynamic I drive  
w.r.t. acc-/decelerations

Il Quality Level in dettaglio



Livello Prudenza



% Km oltre i limiti di velocità: 5,1%

Il tuo giudizio: \* **Buono**  
Livello Prudenza: 222/450

E' calcolato sulla percentuale di km percorsi nel rispetto dei limiti di velocità, con una tolleranza di 10km/h.

Livello Rischio



Il tuo giudizio: \* **Molto Buono**  
Livello Rischio: 309/450

Misura la percentuale di km percorsi nei diversi tipi di strada durante mattino, pomeriggio/sera e notte. Le combinazioni meno rischiose migliorano il Livello.

Livello Attenzione

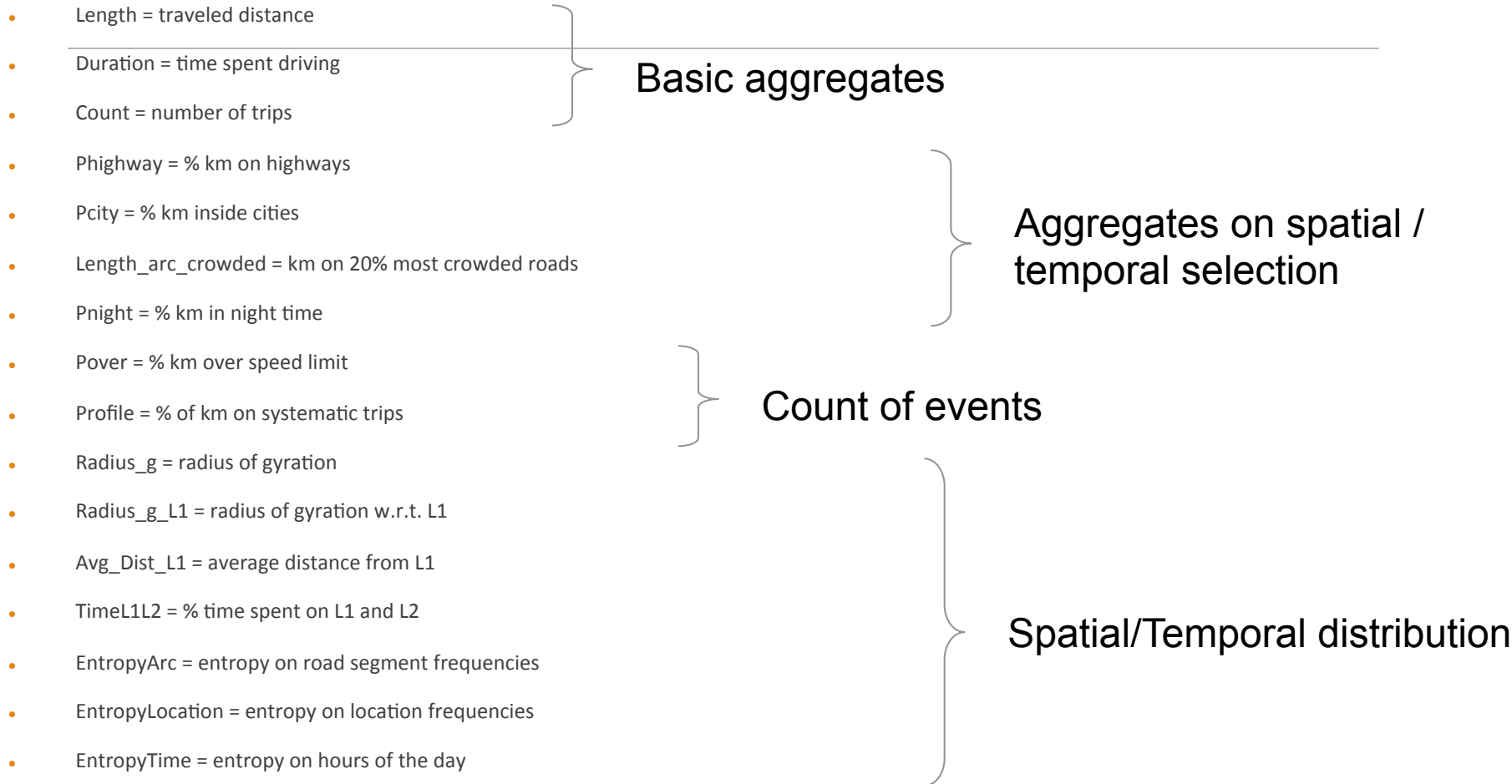


% Km oltre i limiti di velocità: 5,1%

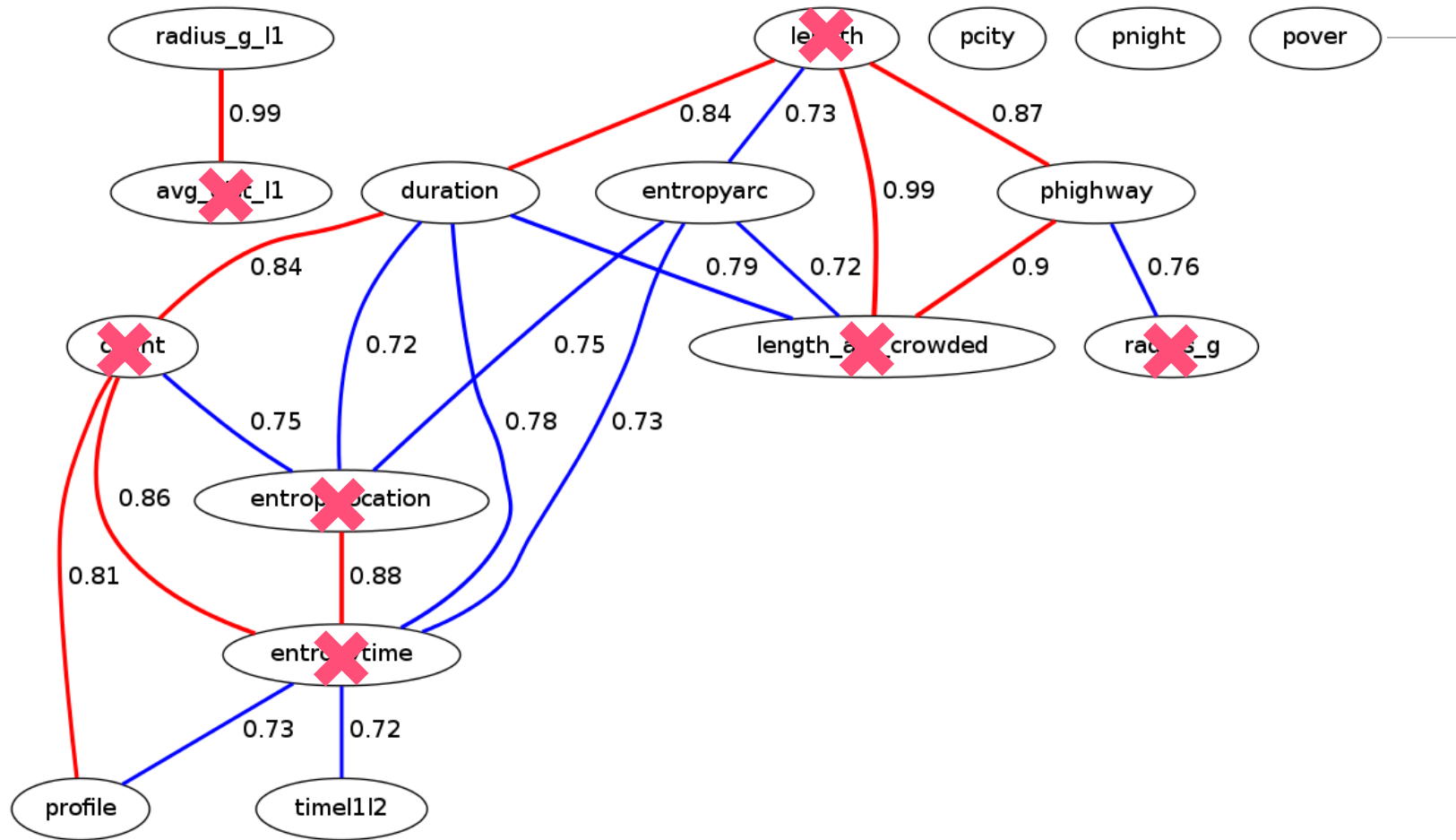
Il tuo giudizio: \* **Buono**  
Livello Attenzione: 49/100

Considera l'intensità delle accelerazioni e decelerazioni durante la guida. Al momento questo livello viene calcolato in proporzione al Livello Prudenza.

# Features over sliding window



# Correlation analysis

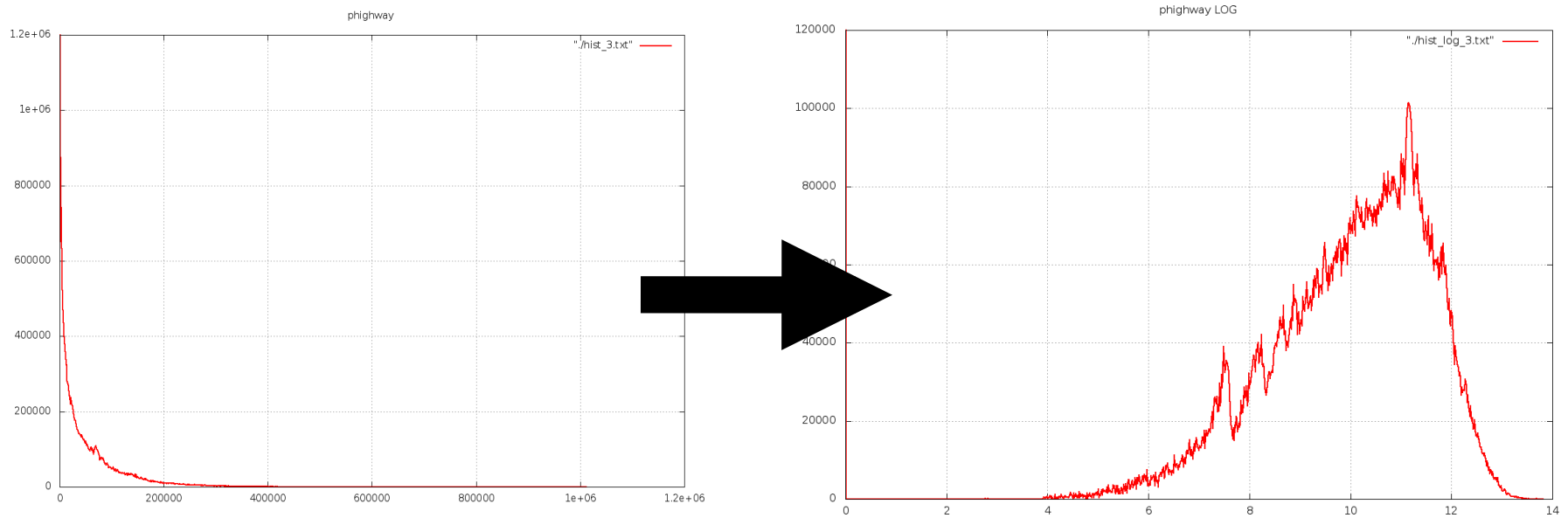


# Features over sliding window

- Length = traveled distance
  - Duration = time spent driving
  - Count = number of trips
  - Phighway = % km on highways
  - Pcity = % km inside cities
  - Length\_arc\_crowded = km on 20% most crowded roads
  - Pnight = % km in night time
  - Pover = % km over speed limit
  - Profile = % of km on systematic trips
  - Radius\_g = radius of gyration
  - Radius\_g\_L1 = radius of gyration w.r.t. L1
  - Avg\_Dist\_L1 = average distance from L1
- Basic aggregates
- Aggregates on spatial / temporal selection
- Count of events
- Spatial/Temporal distribution

# Features normalization

- Log transformation for features with skewed distribution



- Z-score normalization for all features



## (2) Compute driving profiles

- Clustering-based definition

---

  - Profile = representative set of indicators for a large group of drivers with similar behaviors (i.e. similar indicator values)
- Clustering method
  - **K-means** – a partitional, center-based clustering algorithm
  - **Euclidean distance** over driving indicators
  - Refinements: Iterated K-means & select best solution + Noise removal
- Profile = average point of each cluster

# Cluster refinement

---

- Iterated K-means
  - Run clustering multiple times (→ initial random seeding)
  - Select output with best quality
    - Based on clusters compactness (→ SSE – see definition later)
- Noise removal
  - Performed at postprocessing
  - From each cluster, remove points  $p$  such that
$$d(p,c) > 2 \text{ median } \{ d(x,c) \mid x \text{ in cluster} \}$$
where  $c$  is the cluster center
  - Alternative solutions are possible
    - e.g.: density-based noise removal

# Experimental setting

- GSP traces of an insurance

company customers

- 35 days monitoring

- Sample of ~11k vehicles

moving in the area

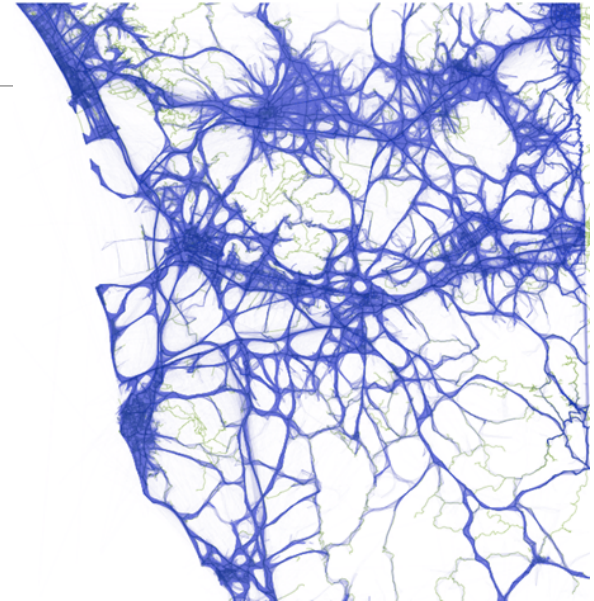
- Short temporal thresholds for

testing purposes

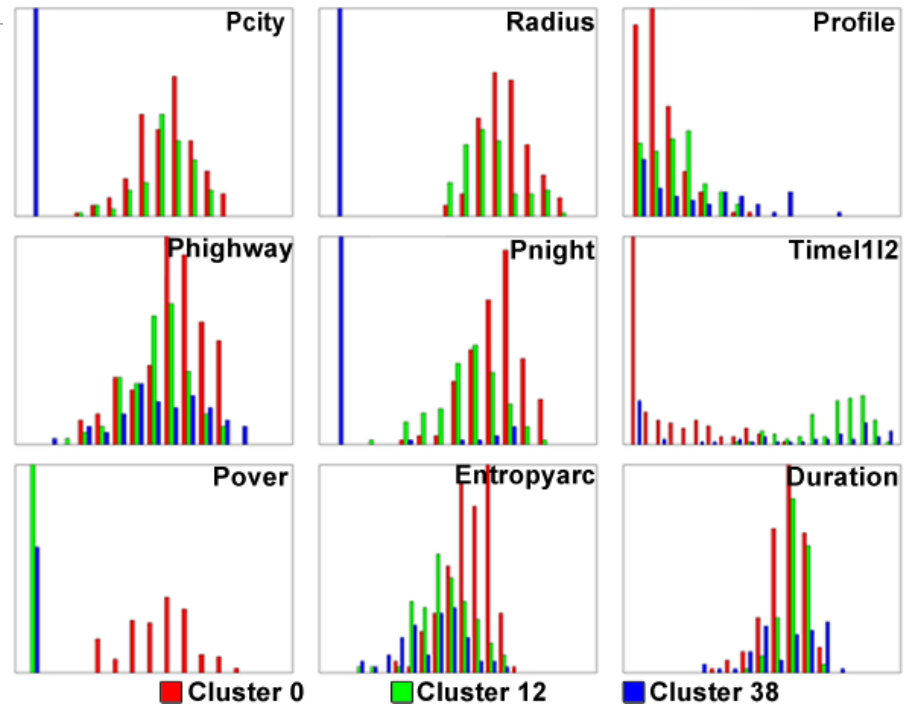
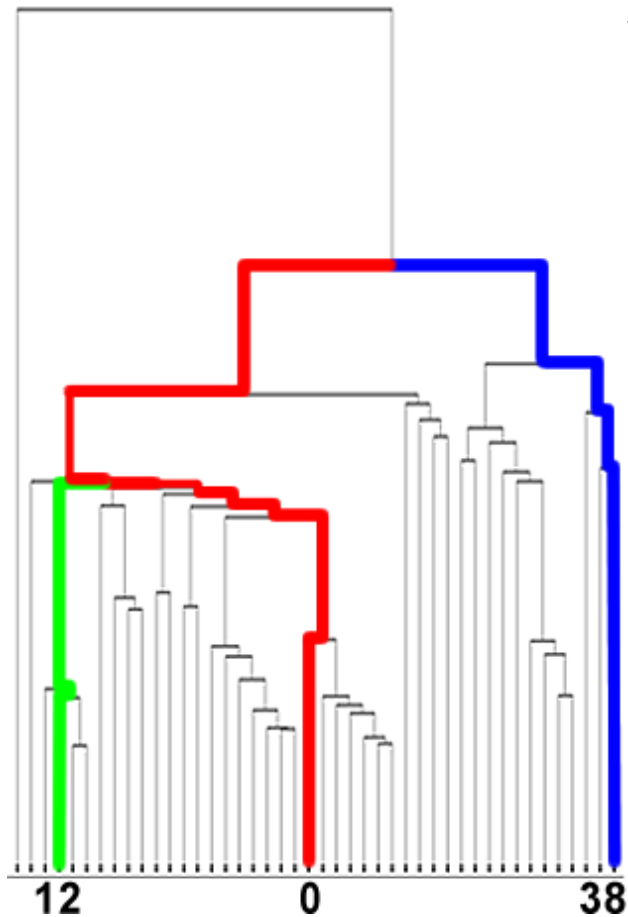
- Compute driving indicators over a sliding window of 3 days

- Update indicators every 15'

- Most likely larger in a real application – parameter tuning to be done with domain experts



# Experiments: clusters inspection



**Explorers**

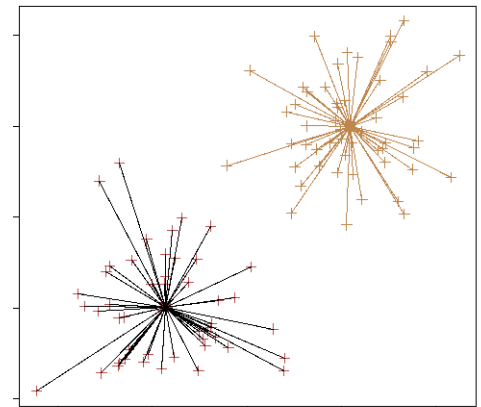
**Long-range commuters**

**Sunday drivers**

# (3) Driving profiles monitoring

---

- Translated to “cluster quality monitoring”
- Quality measure: SSE = Sum of Squared Errors
  - Given a clustering  $C = \{ C_1, \dots, C_k \}$ , and average points  $m_i$  for each cluster  $C_i$



### (3) Driving profiles monitoring

---

DEFINITION 1 (CLUSTER MONITORING PROBLEM).

*Given a clustering  $C = \{C_1, \dots, C_k\}$  having initial SSE equal to  $SSE_0$ , and given a tolerance  $\alpha \in \mathcal{R}^+$ , we require to ensure that at each time instant  $t$  the following holds for the SSE of the (dynamic) dataset  $D_t$ :*

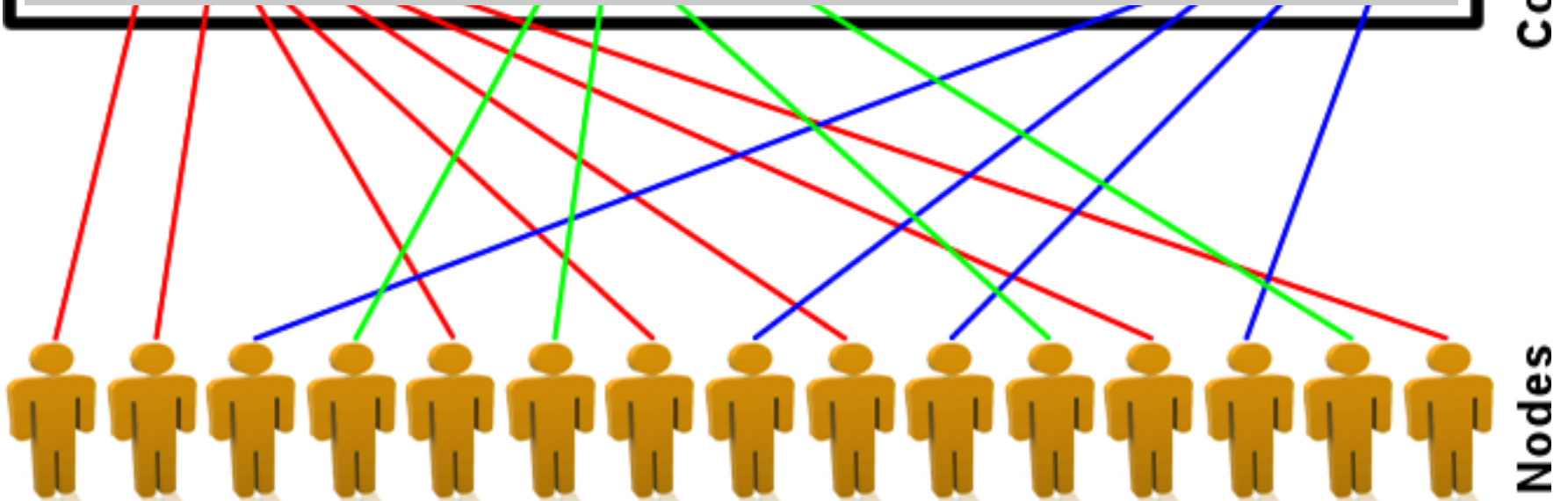
$$SSE_t \leq (1 + \alpha)SSE_0$$

*When that does not happen, a recomputation/update of cluster assignments should be performed.*

# Monitoring process

**Initialization:** compute clusters, cluster centers (used as reference points for Safe Zones) and distribute SSE thresholds to clusters

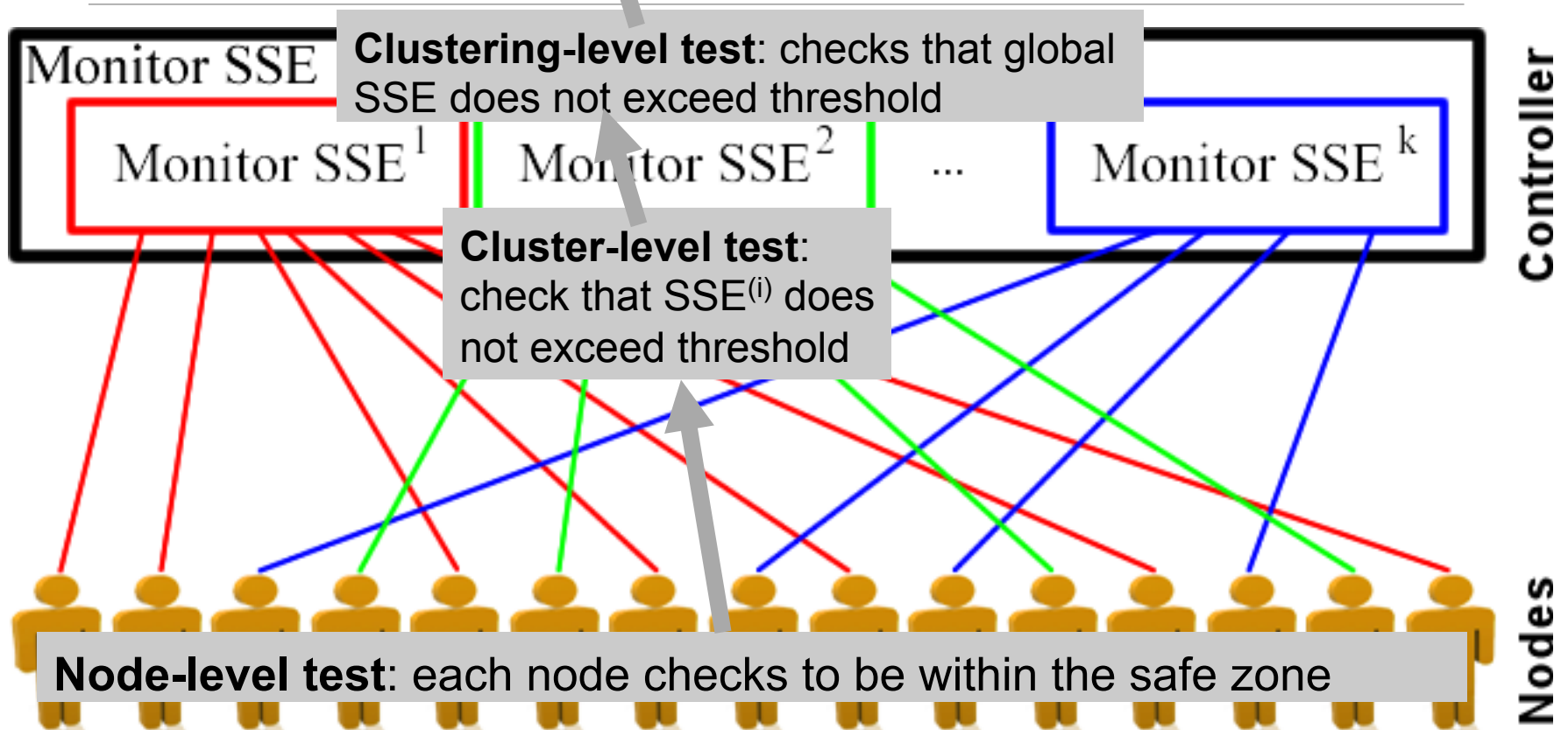
Controller



Nodes

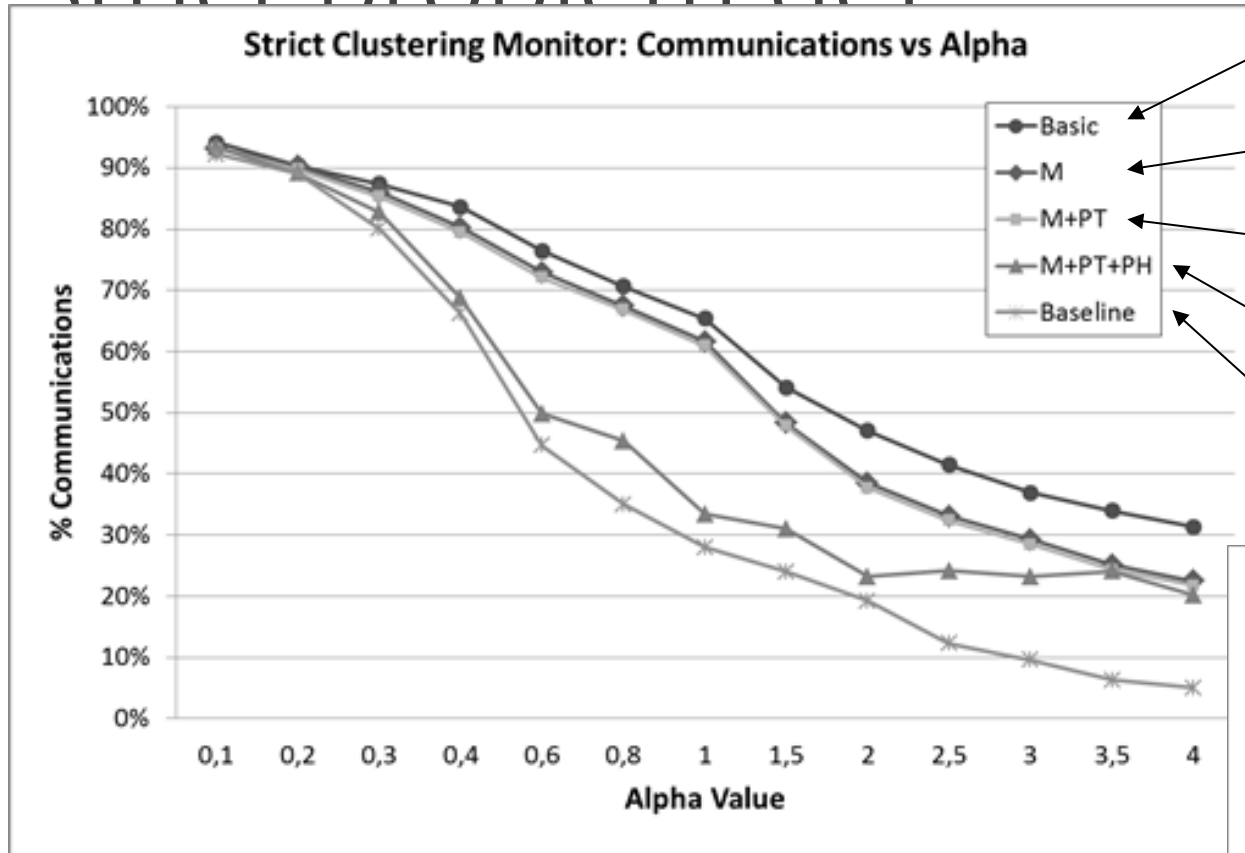
# Monitoring process

Re-clustering





# Experiments: communications / strict problem def



Balancing/memoryless

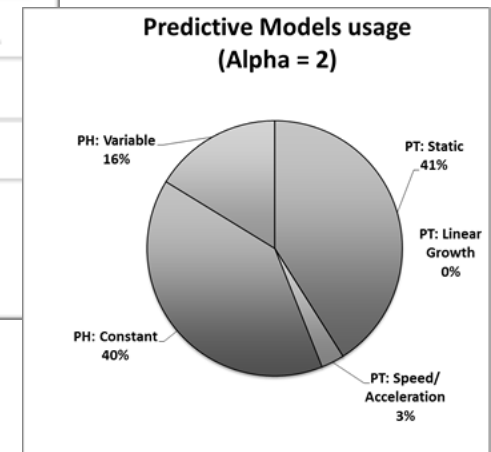
Balancing/memory

Trend Predictive Ms

History Predictive Ms

Oracle (no false alarms)

Communications from controller w/ broadcasting:  
between 1.23% and 2.34%, dominated by balancing



---

# **Services Towards Individual Users**

***Self-awareness***

# Self-awareness services

- Mobility-based specialization of self-awareness services for generic users
  - Provide summary of activity of the user
  - Provide comparison against collectivity

# Self-awareness services

- Summaries based on

---

  - Temporal statistics
  - Spatial statistics / distributions
  - Movement aggregates

# User's activity summaries

- An example within Generali



## Chilometraggio mensile

## Il Quality Level in dettaglio

### Livello Prudenza

### Livello Rischio

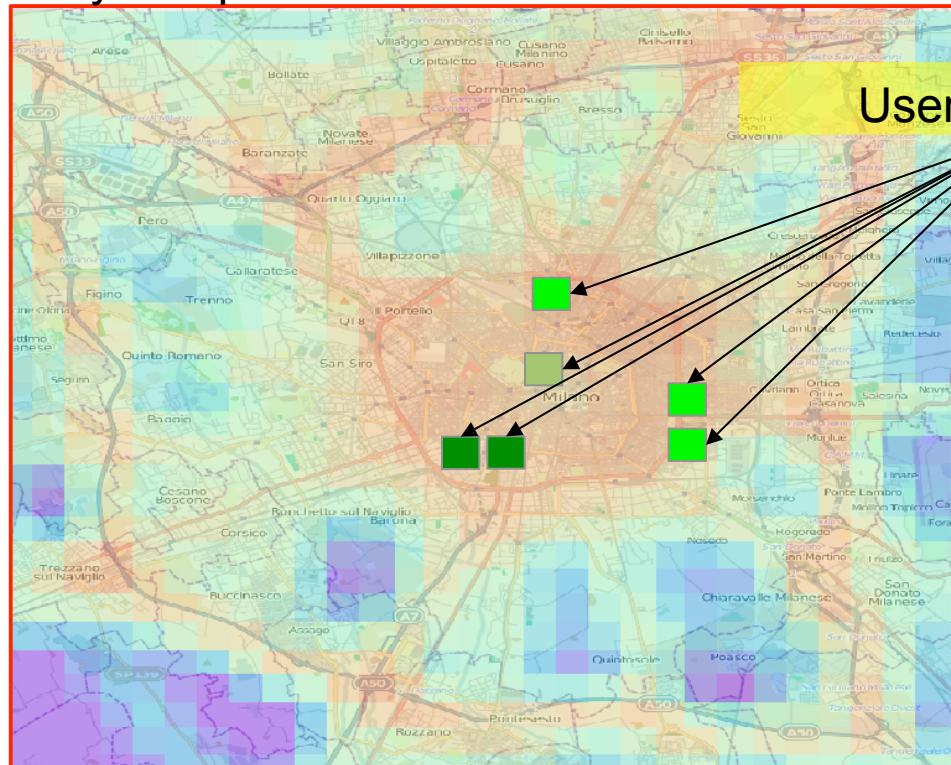
### Livello Attenzione

## Il tuo chilometraggio per Marzo 2013

# Comparison against collectivity

- In space

City hotspots

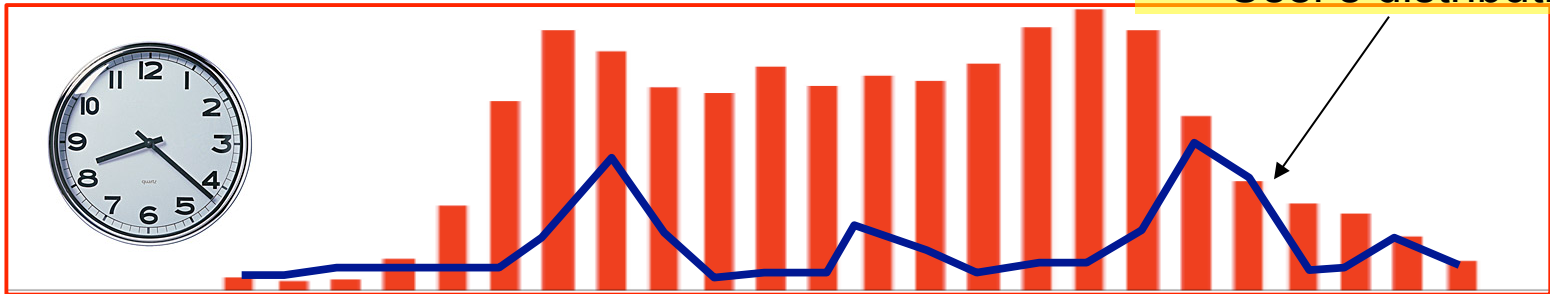


User's hotspots

# Comparison against collectivity

- In time

City time distribution

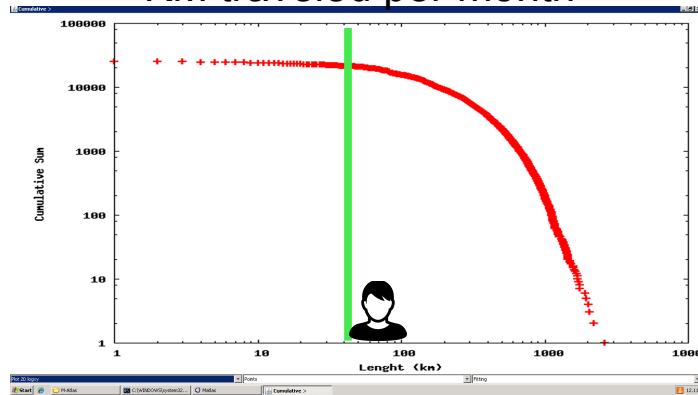


User's distribution

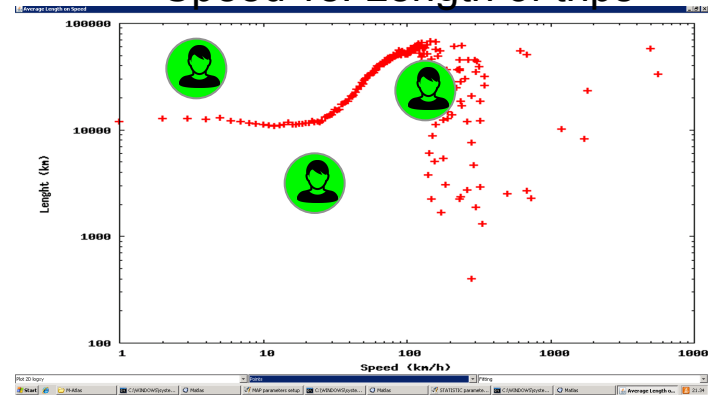
# Comparison against collectivity

- On general statistics

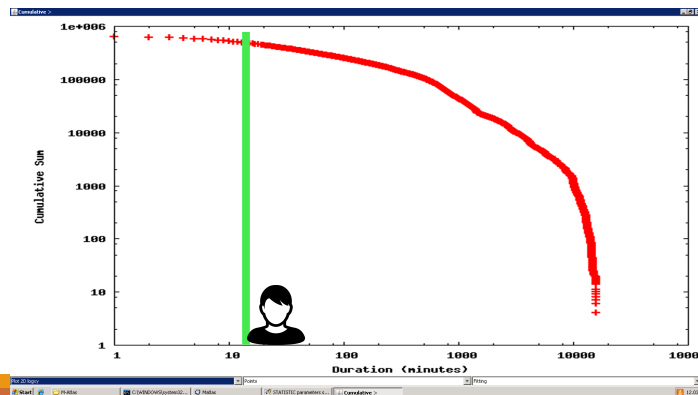
KM traveled per month



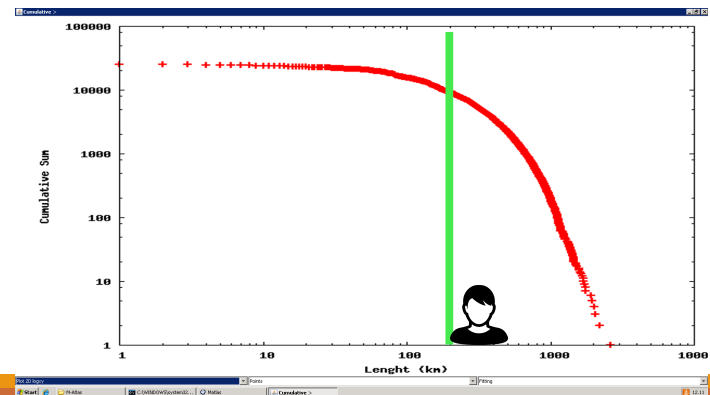
Speed vs. Length of trips



Total duration of travels



Radius of gyration





---

# Services Towards Individual Users

## *Proactive Carpooling*



# Proactive car pooling

# Carpooling cycle

## Context

- Several initiatives, especially on the web



Raggiungi  
**Fieracavalli**  
in carpooling!



# Carpooling cycle

## Distinctive features

---

### Traditional approach vs. Data-driven cycle

- |  |   |   |
|--|---|---|
| • Users manually insert and update their rides | → | • System autonomously detect systematic trips |
| • Users search and contact candidate pals      | → | • System automatically suggest pairings       |
| • Users make individual, “local” choice        | → |   |

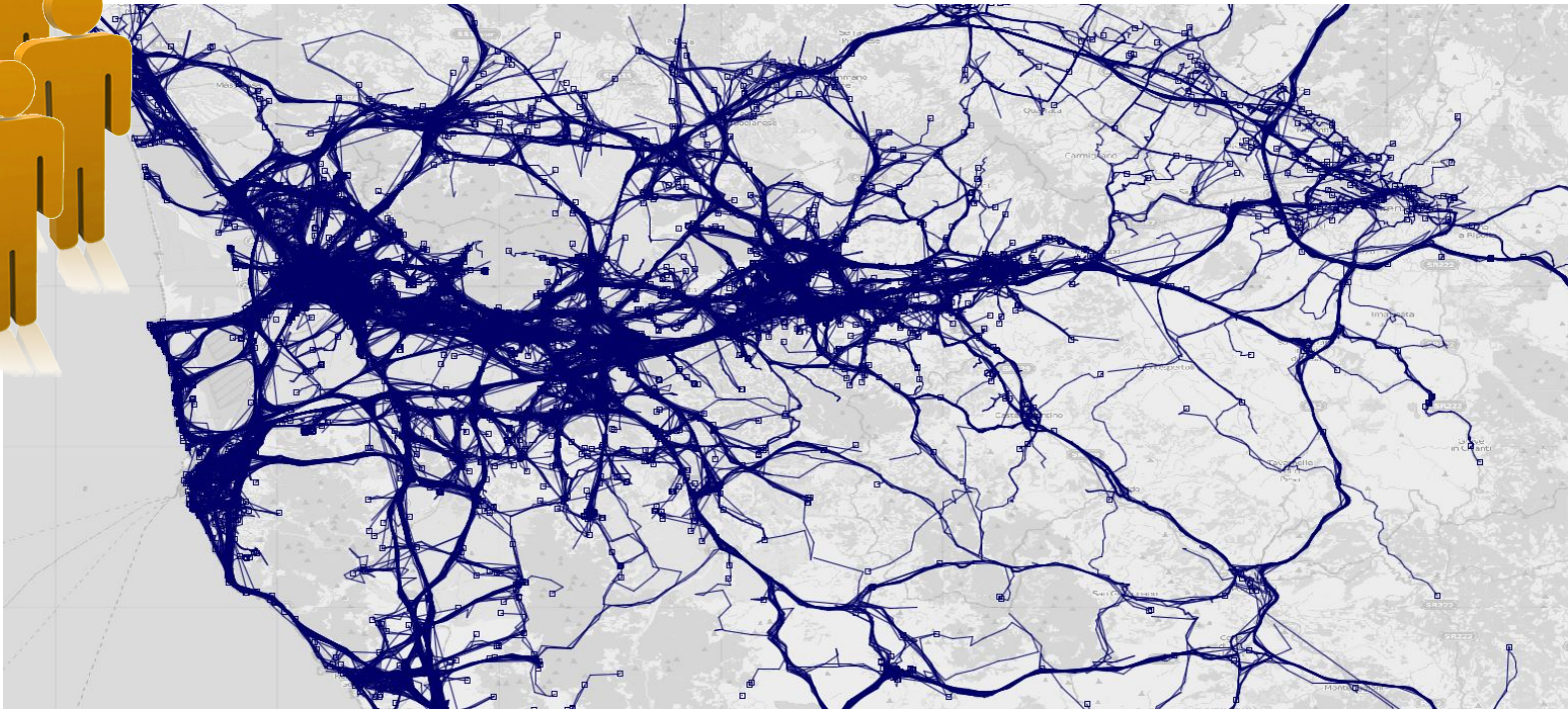
System cooks

# Carpooling cycle

## Assumptions

---

- Users provide access to their mobility traces

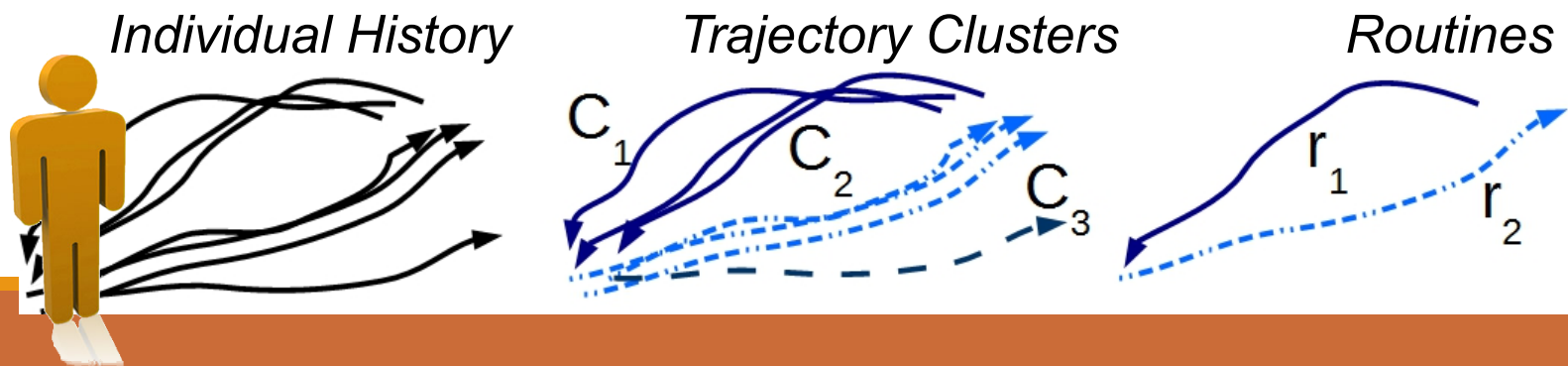




# Carpooling cycle

## Step 1: Inferring Individual Systematic Mobility

- Extraction of Mobility Profiles
  - Describes an abstraction in space and time of the systematic movements of a user.
  - Exceptional movements are completely ignored.
  - Based on trajectory clustering with noise removal

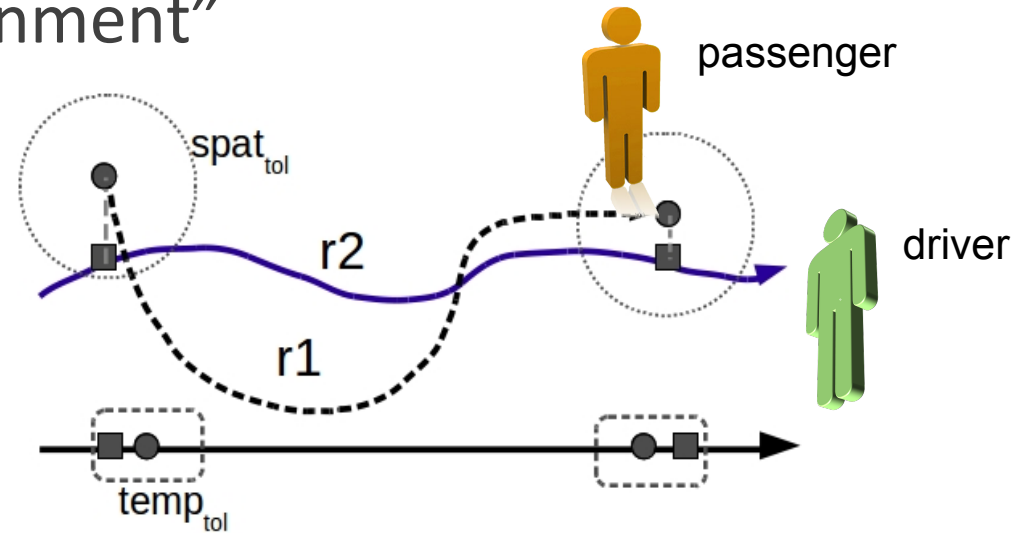


# Carpooling cycle

Step 2: Build Network of possible carpool matches

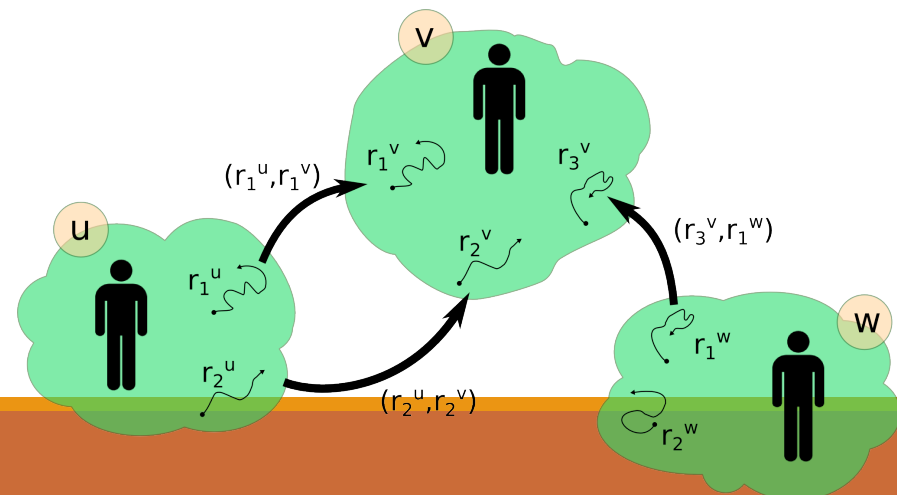
- Based on “routine containment”

- One user can pick up the other along his trip



- Carpooling network

- Nodes = users
- Edges = pairs of users with matching routines



# Carpooling cycle

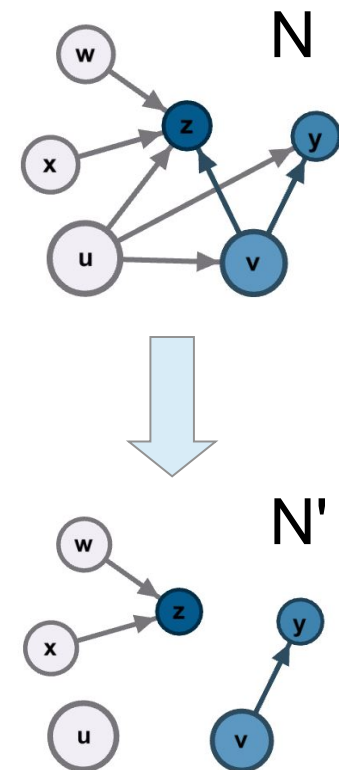
## Step 3: Optimal allocation of drivers-passengers

- Given a Carpooling Network  $N$ , select a subset of edges that minimizes  $|S|$

- $S$  = set of circulating vehicles

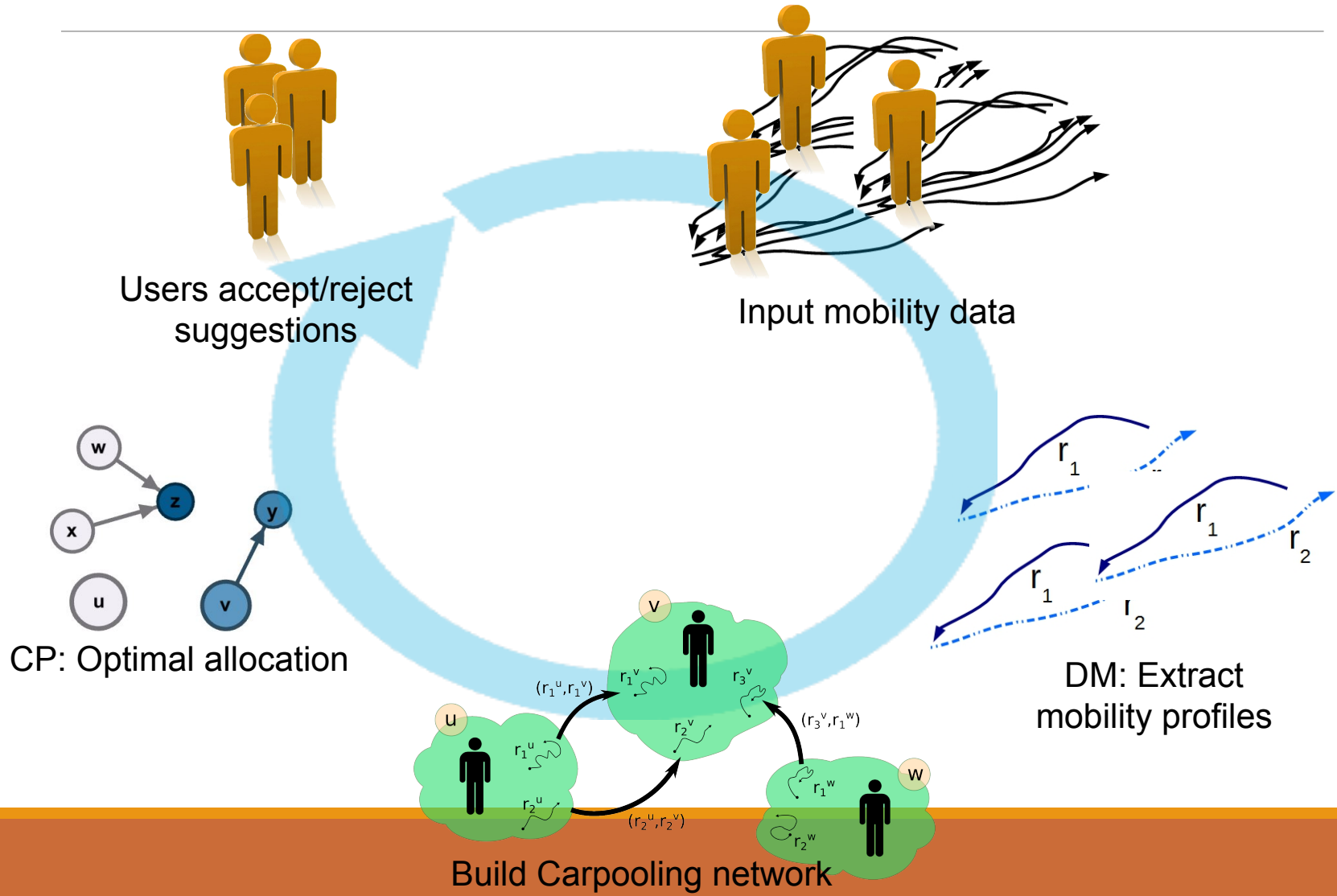
provided that the edges are coherent, i.e.:

- $\text{indegree}(n)=0$  OR  $\text{outdegree}(n)=0$  (a driver cannot be a passenger)
  - $\text{indegree}(n) \leq \text{capacity}(n)$





# Carpooling cycle



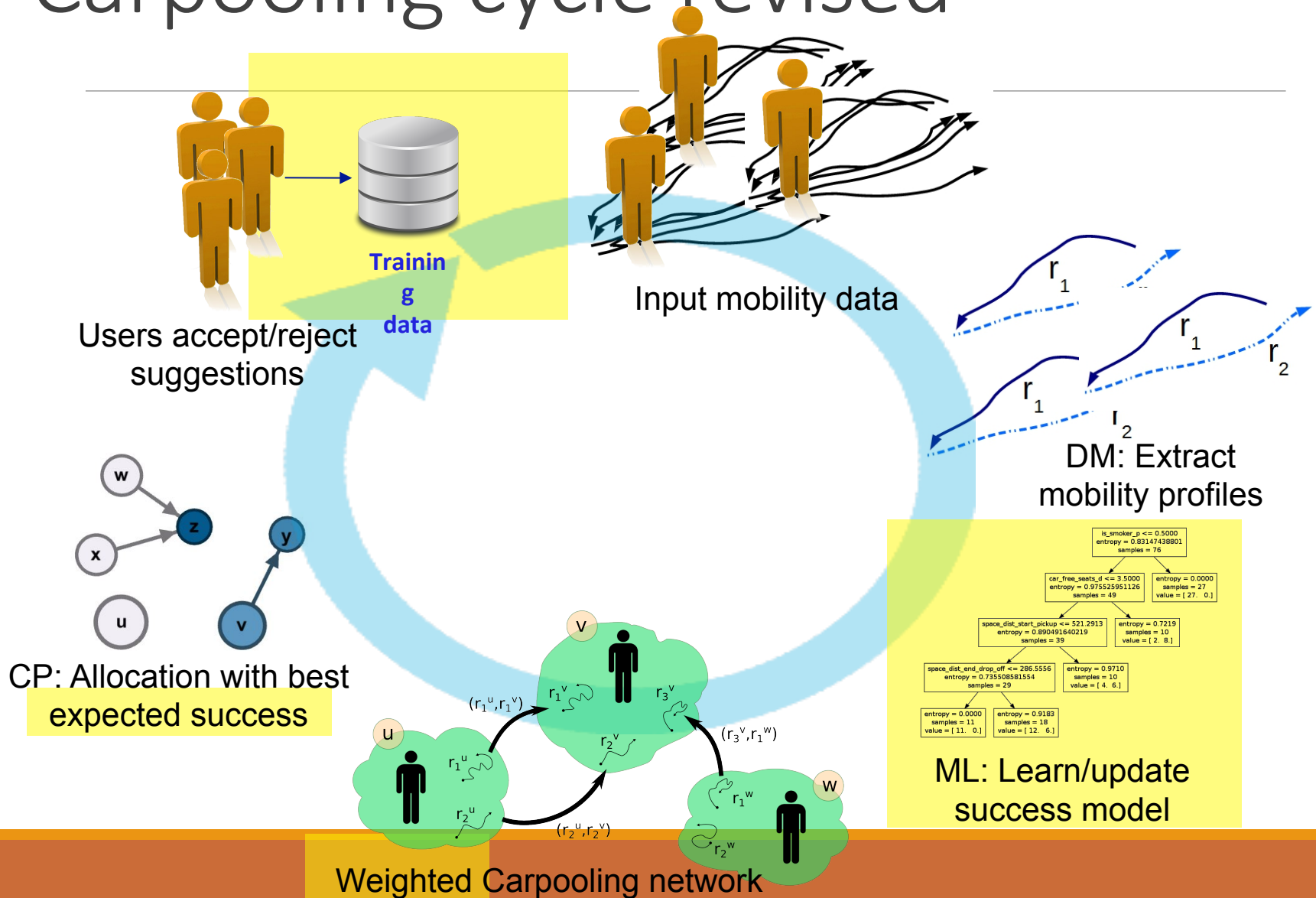
# Carpooling cycle

## Improvement

---

- In carpooling (especially if proactive) users might not like the suggested matches
  - Impossible to know who will accept a given match
  - Modeling acceptance might improve results
- Two new components
  - **Learning** mechanism to guess success probability of a carpooling match
  - **Optimization** task exploits it to offer solution with best expected overall success

# Carpooling cycle revised



# Carpooling cycle

## Learning a success model

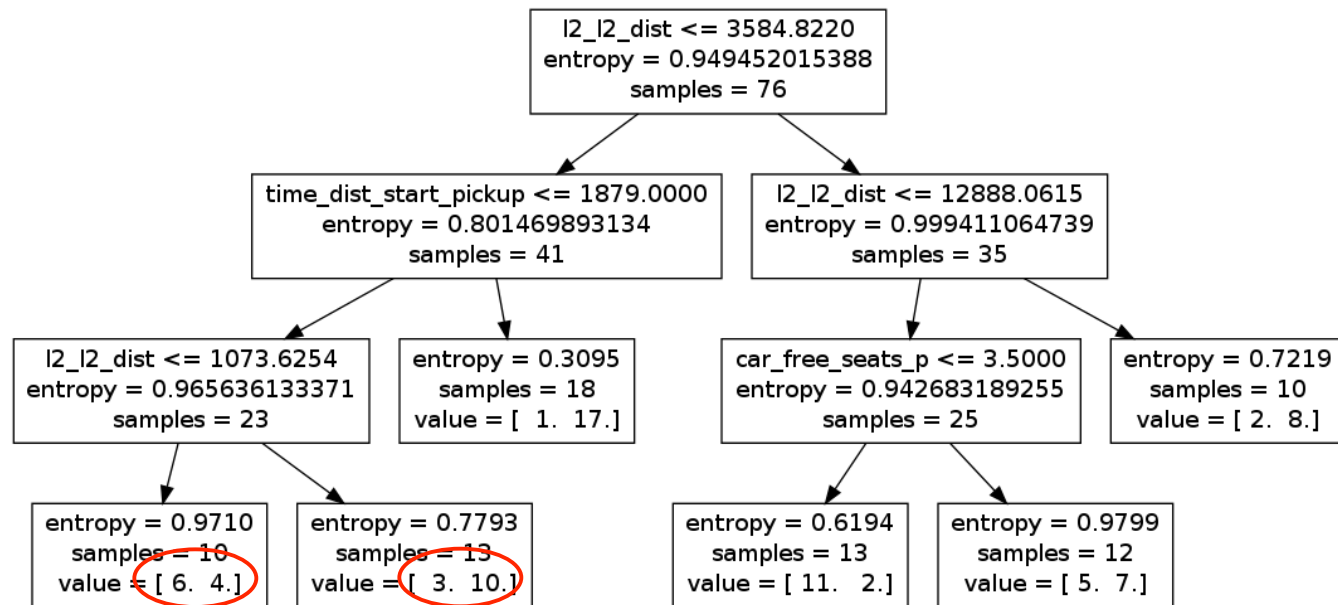
---

- **Input:** set of features describing a single carpooling pair
- **Output:** success probability  $p$  in  $[0,1]$
- 36 Features adopted
  - **Ease of carpooling:** space\_dist\_start\_pickup, space\_dist\_end\_drop\_off, time\_dist\_start\_pickup, time\_dist\_end\_drop\_off, time\_pick\_up\_get\_off, start\_together, end\_together, distance\_between\_homes, dist\_between\_works
  - **Personal features** (of both driver and passenger): age, gender, marital\_status, occupation, is\_smoker, has\_children, has\_animals, car\_free\_seats →  
Cannot be inferred, need external data
  - **Past personal history in the service** (of both driver and passenger): last\_driver\_accepted, last\_passenger\_accepted, %\_acceptance\_driver, %\_acceptance\_passenger
  - **History of the two users together** (if any): last\_accepted\_pair, last\_rejected\_pair, %\_accepted\_pair

# Carpooling cycle

## Learning a success model

- Model selected: “probability estimation tree”
  - simple decision tree with assigned probabilities of prediction in the leaves



$P(\text{Yes}) = 6/10 = 60\%$        $P(\text{Yes}) = 3/13 = 23\%$

# Carpooling cycle

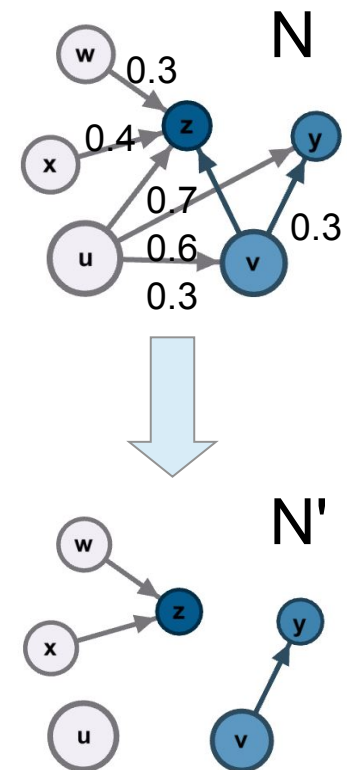
## Revised optimization model

- Given a Carpooling Network  $N$ , select a subset  $W$  of edges that maximize

- $\sum p(w) \mid w \text{ in } W$

provided that the edges are coherent, i.e.:

- $\text{indegree}(n)=0$  OR  $\text{outdegree}(n)=0$  (a driver cannot be a passenger)
  - $\text{indegree}(n) \leq \text{capacity}(n)$



# Carpooling cycle

## Two usage scenarios

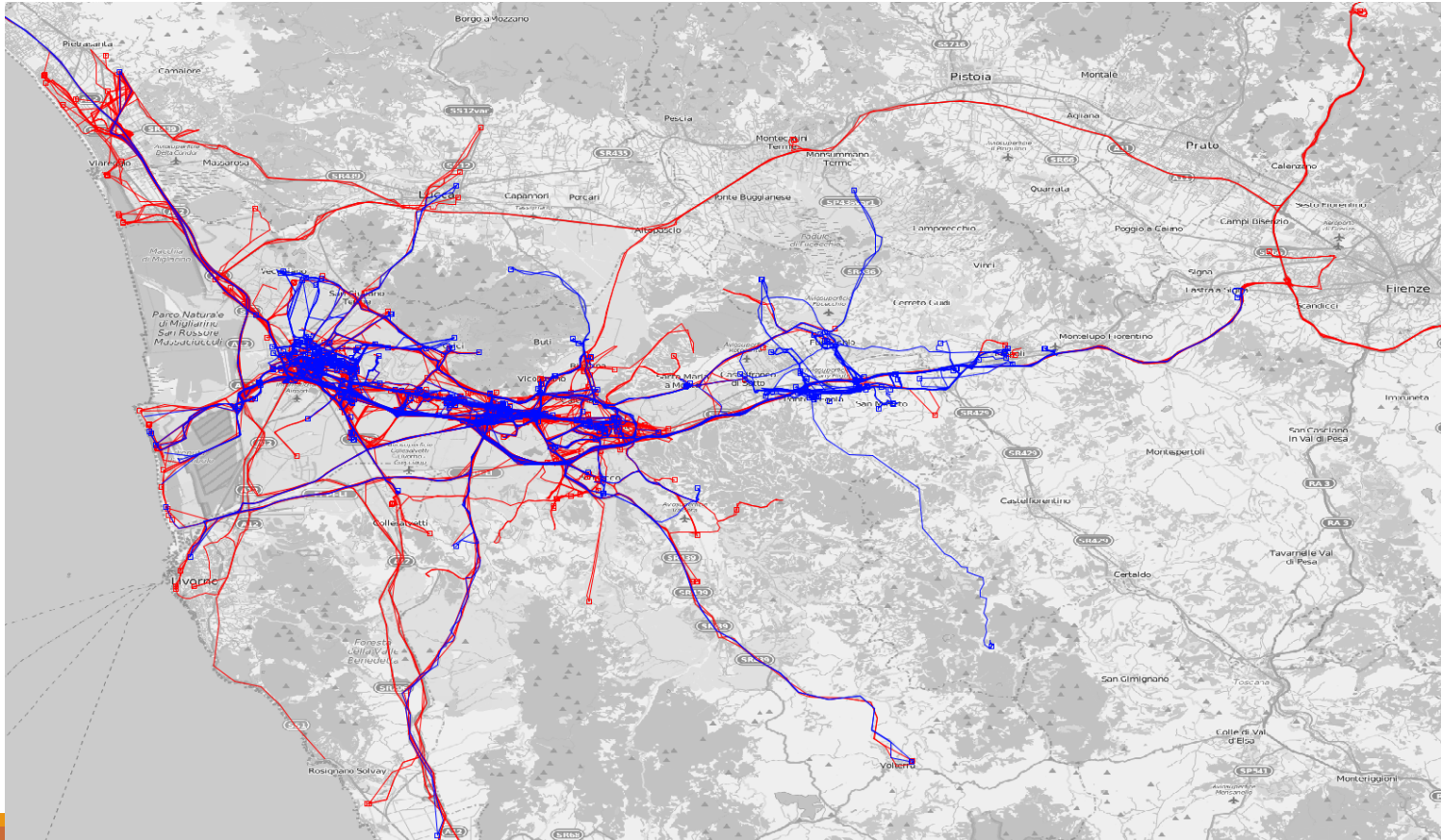
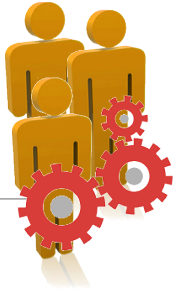
- Scenario 1:
  - Real service is implemented, with real users interacting (accept/reject suggestions)
- Scenario 2:
  - Simulation environment where the users' behaviour is simulated through a model
  - Mobility data is taken from historical traces
  - Useful to perform what-if analyses on
    - (i – social) effects of different users' behaviours
    - (ii – performances) effects of different learning strategies



# Carpooling cycle

## Scenario 2 – sample results

- Profiles involved in carpooling network



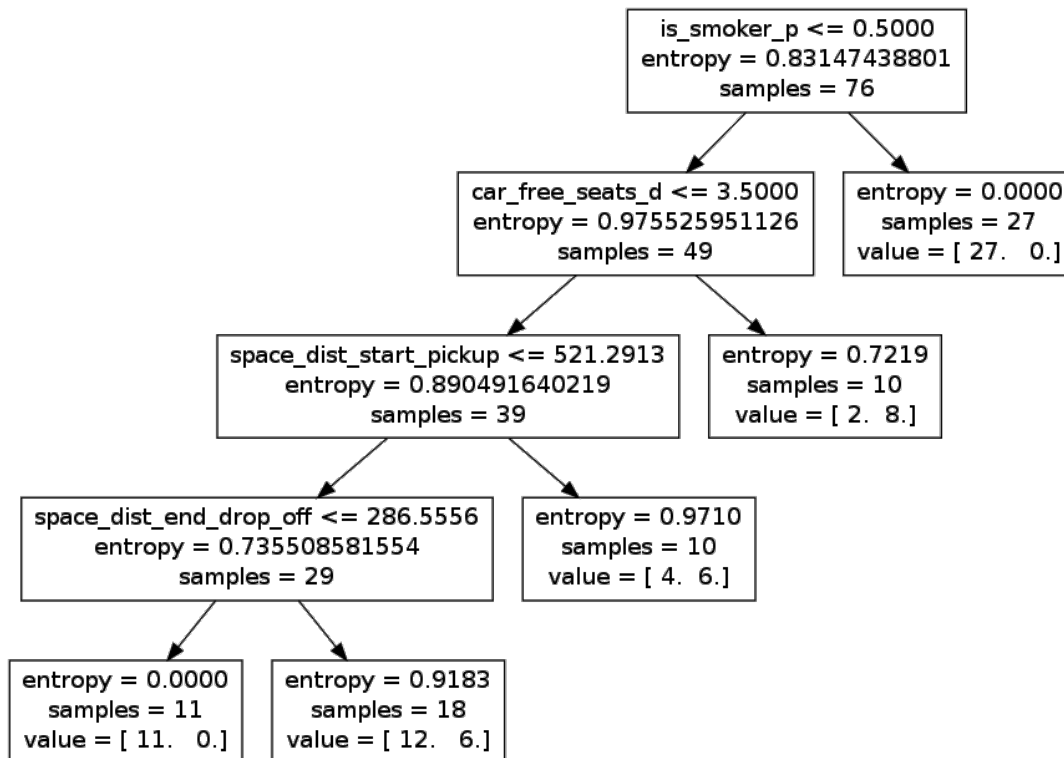


# Carpooling cycle

## Scenario 2 – sample results



- Prediction models



### Iteration 0

is\_smoker\_p : 0.51763342041  
car\_free\_seats\_d : 0.196822768067  
space\_dist\_end\_drop\_off : 0.161445930025  
space\_dist\_start\_pickup : 0.124097881498  
time\_dist\_start\_pickup : 0.0  
last\_accepted\_pair : 0.0  
l1\_l1\_dist : 0.0  
age\_d : 0.0  
gender\_p : 0.0  
has\_children\_p : 0.0

### Iteration 4

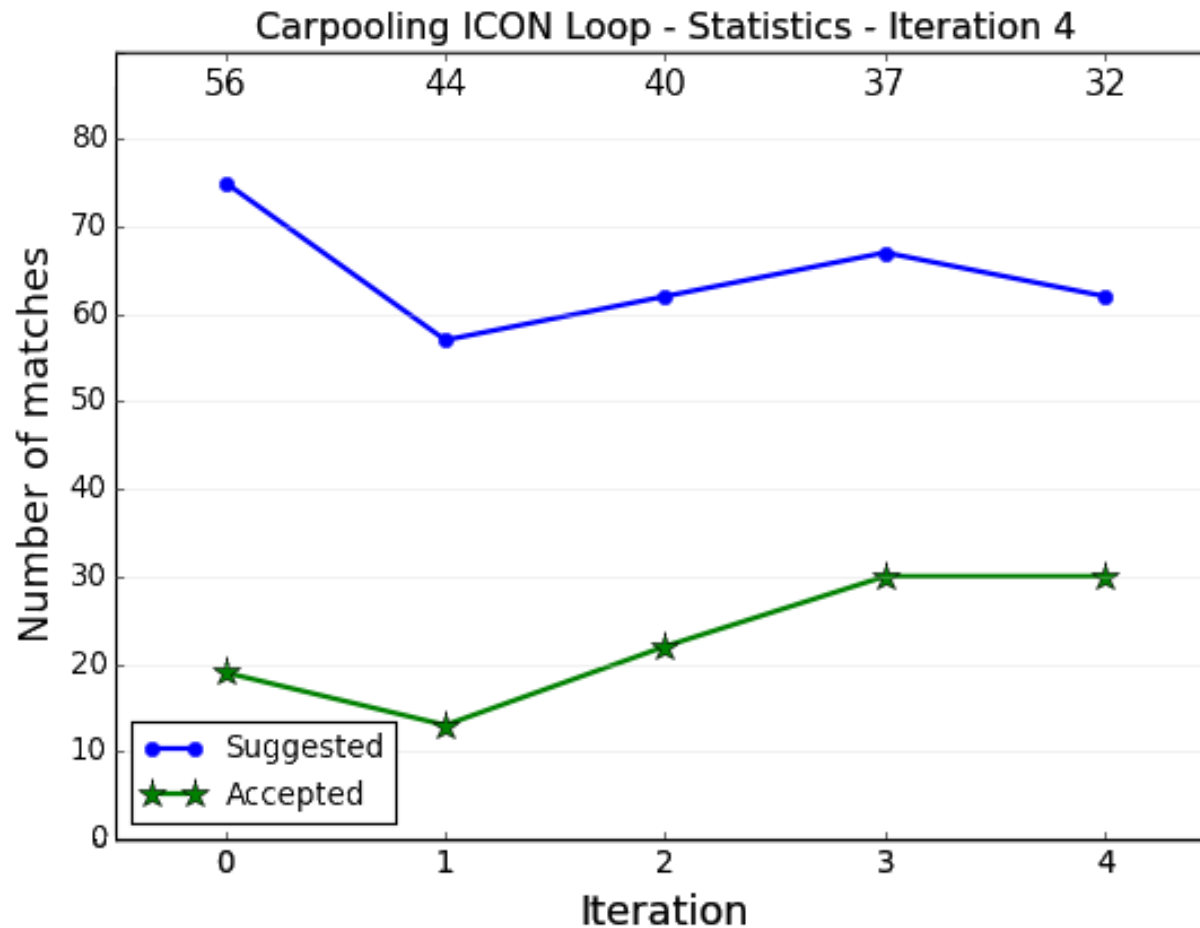
last\_accepted\_pair : 0.300609683595  
%\_accepted\_pair : 0.18422352604  
gender\_d : 0.121782490916  
is\_smoker\_d : 0.096830535215  
l1\_l1\_dist : 0.0947711528021  
is\_smoker\_p : 0.0921934235296  
age\_p : 0.0549409842076  
gender\_p : 0.0396236591312  
time\_dist\_start\_pickup : 0.00874162379163  
car\_free\_seats\_d : 0.00628292077177

# Carpooling cycle

## Scenario 2 – sample results



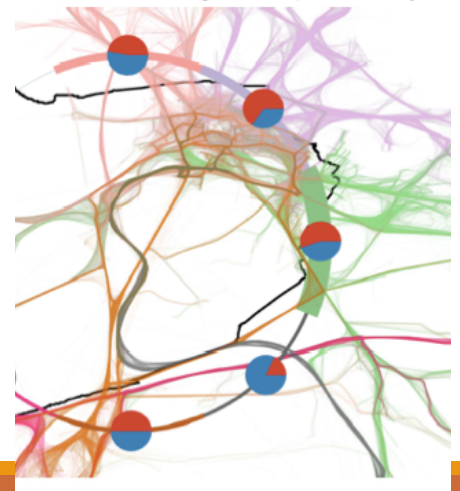
- Performances



---

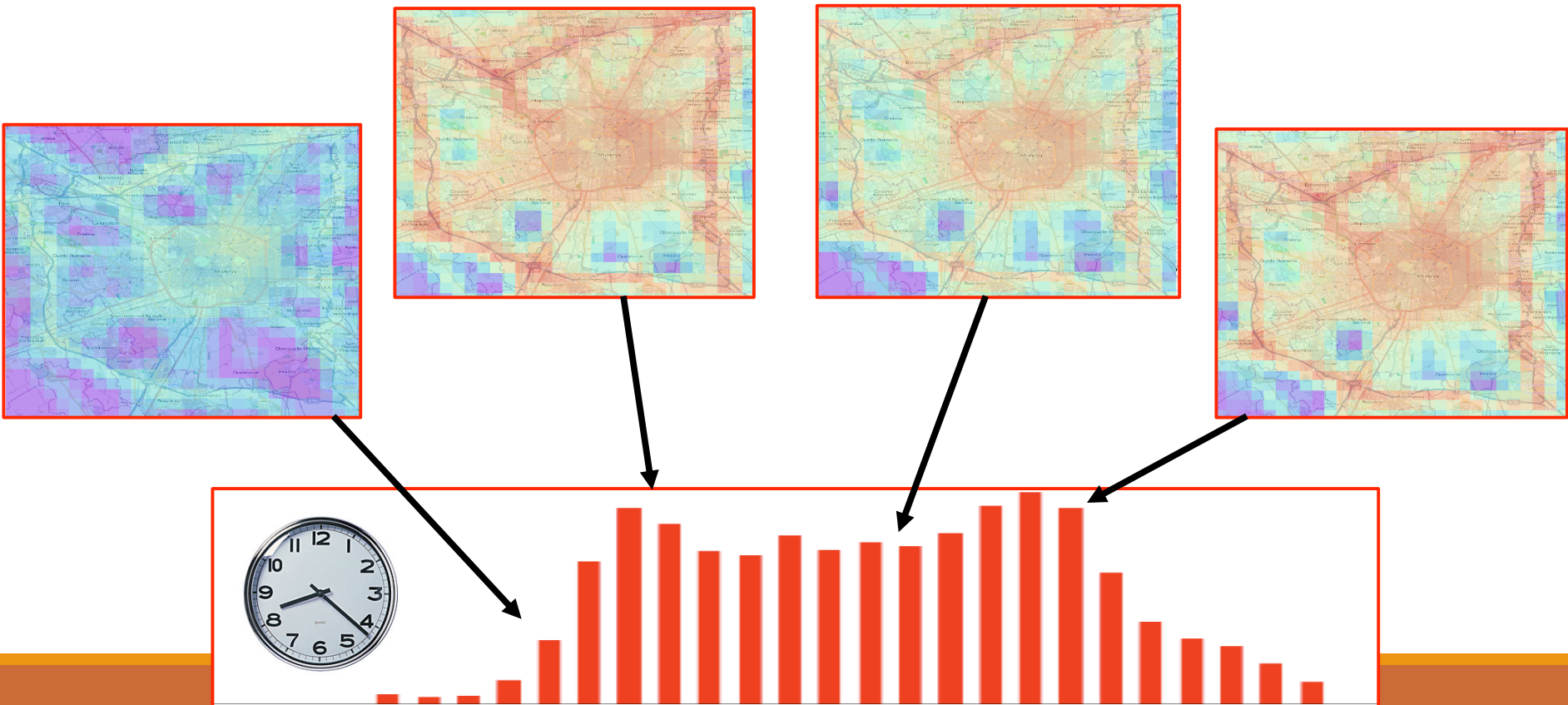
# Services Towards Public Sector

*Urban Mobility Atlas*

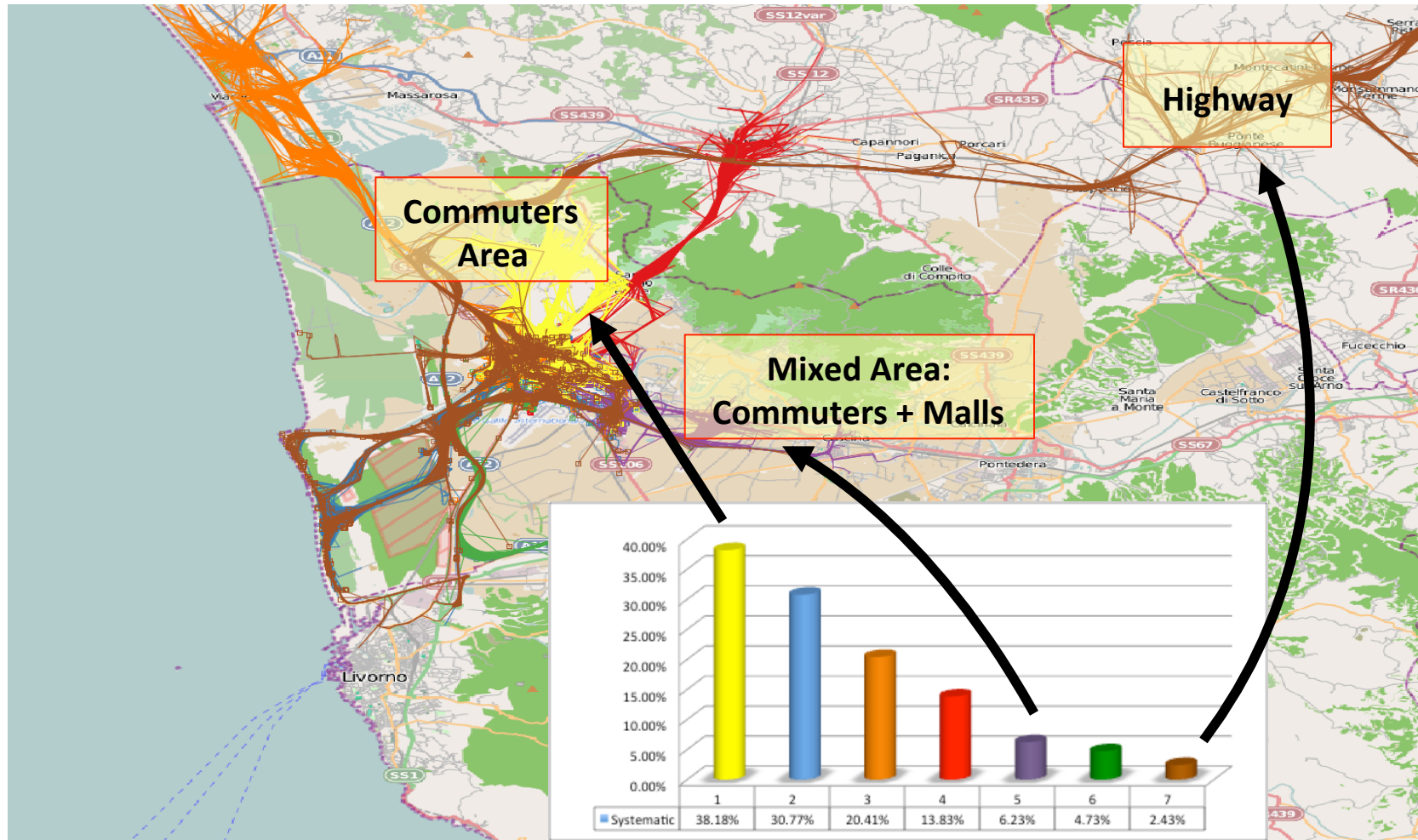


# Dynamics of urban mobility

---



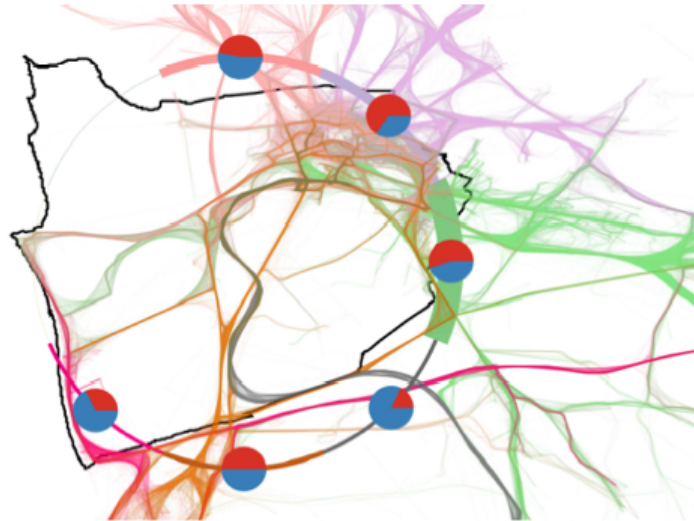
# Impact of Systematic Mobility



Access Routes  
Systematic Mobility (%)

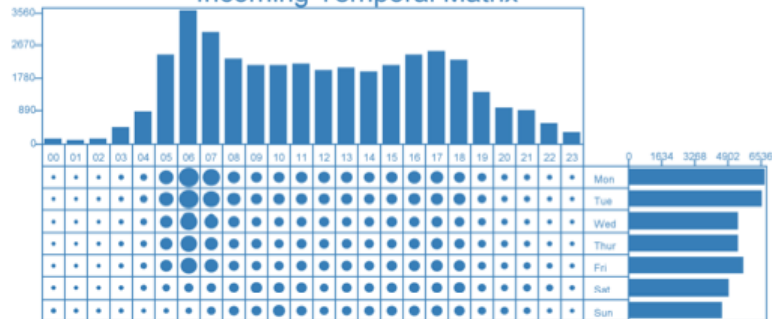
# Pisa – Incoming traffic

Incoming Traffic (38.464 Trajectories)



	City	Traj	Perc
NORD 32%	San Giuliano T...	4.816	62%
	Vecchiano	1.425	94%
	Viareggio	1.142	99%
	Lucca	862	67%
OVEST 0%	Camaione	358	94%
SUD 12%	Livorno	2.843	92%
	Collesalveti	565	50%
	Rosignano Mar...	140	41%
	Fauglia	137	19%
	Cecina	124	45%
EST 54%	Cascina	7.078	97%
	San Giuliano T...	2.881	37%
	Pontedera	1.350	95%
	Calci	795	79%
	Calcinaia	693	92%

Incoming Temporal Matrix



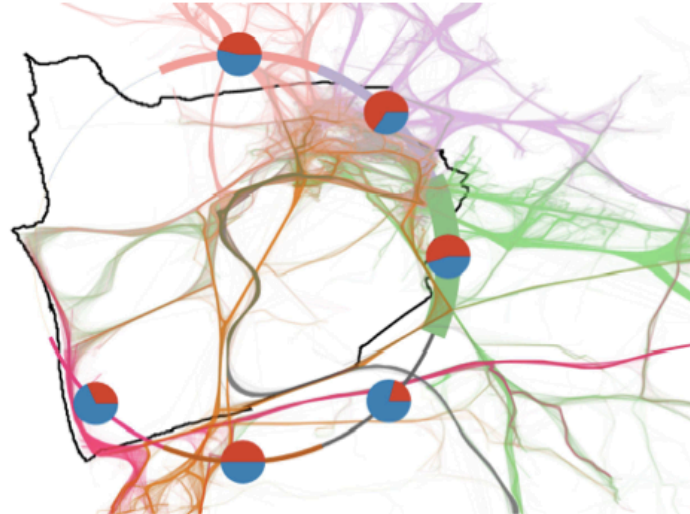
Regular VS Occasional





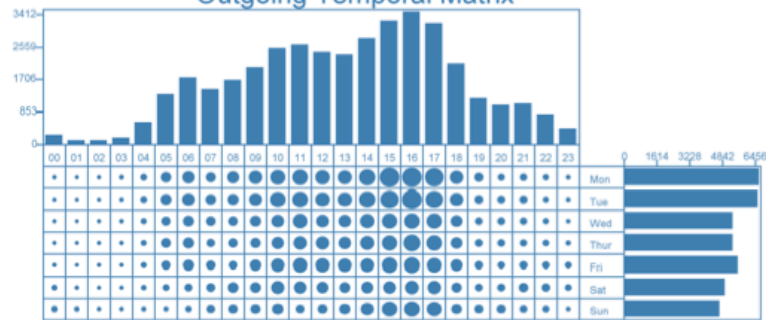
# Pisa – Outgoing Traffic

Outgoing Traffic (38.271 Trajectories)



	City	Traj	Perc
NORD 32%	San Giuliano T.	4.842	62%
	Vecchiano	1.418	93%
	Viareggio	1.117	99%
	Lucca	886	67%
OVEST 0%	Camaione	329	96%
SUD 13%	Livorno	2.812	92%
	Collesalveti	565	51%
	Rosignano Mar.	143	44%
	Fauglia	130	19%
	Cecina	123	45%
EST 54%	Cascina	7.253	97%
	San Giuliano T.	2.860	37%
	Pontedera	1.326	95%
	Calci	798	82%
	Calcinaia	704	93%

Outgoing Temporal Matrix



Regular VS Occasional



# ... and Comparison

---

**Florence**

**Montepulciano**

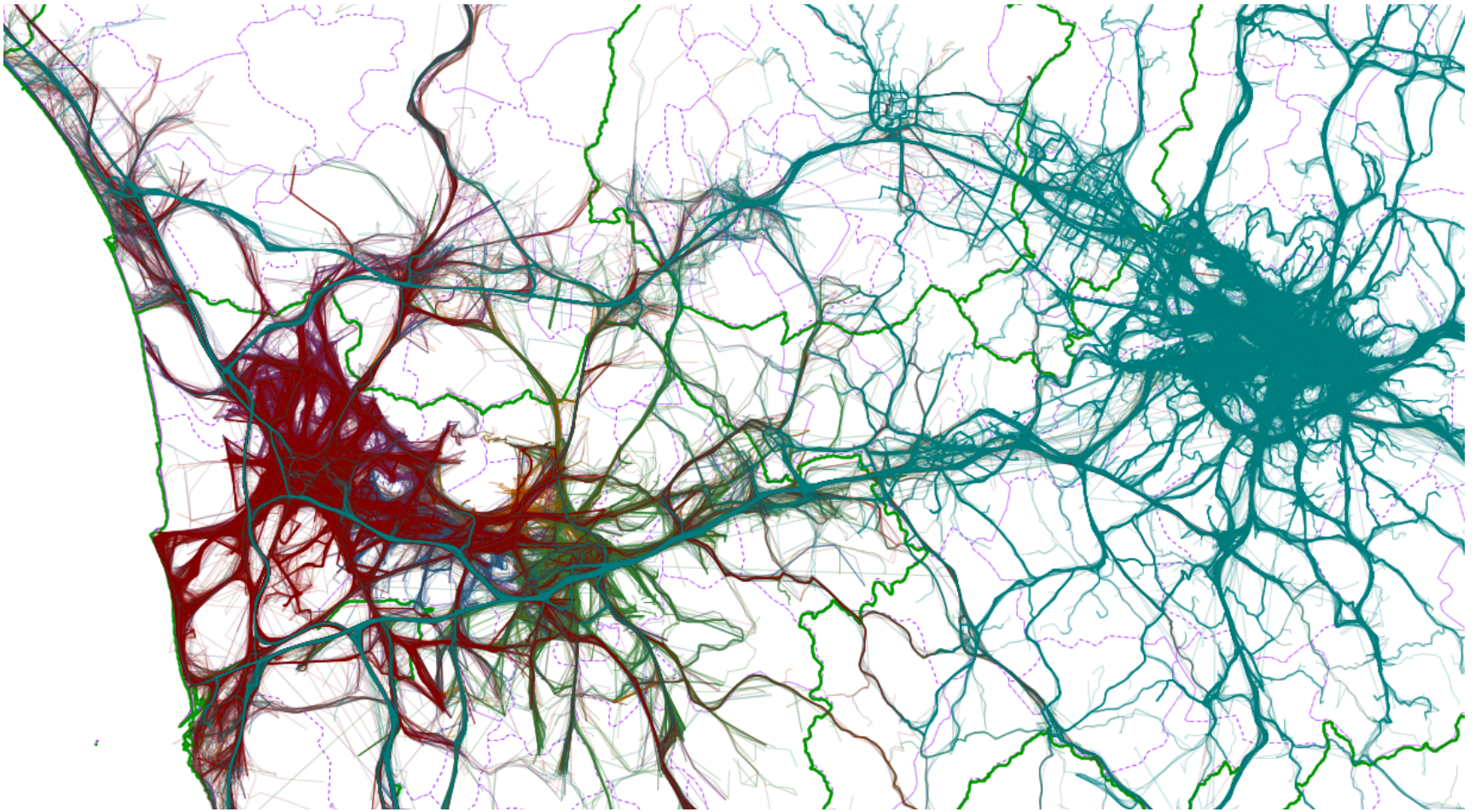


---

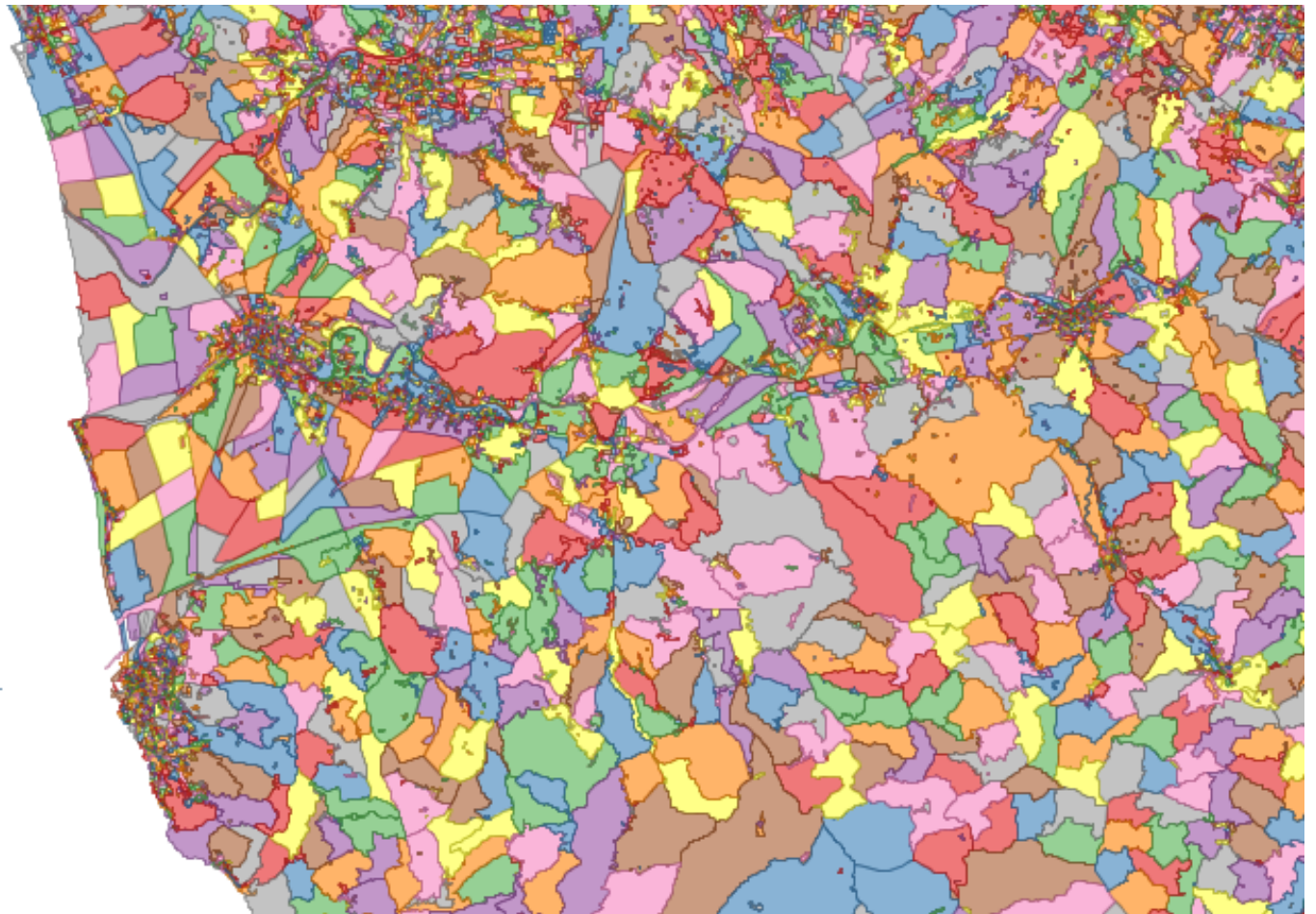
# **Services Towards Public Sector**

***Mobility-based Redefinition of Borders***

# Mobility coverages

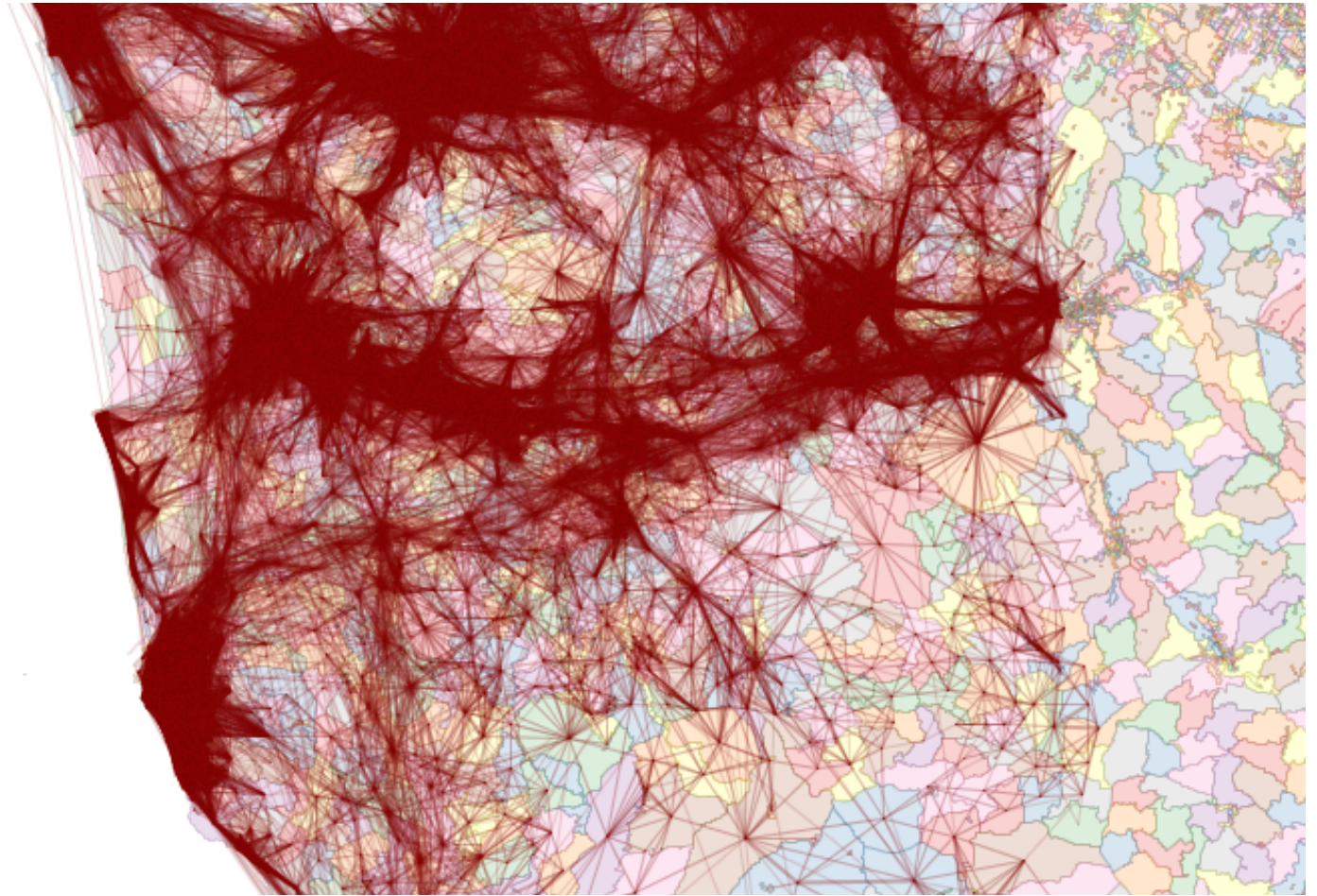


# Step 1: spatial regions

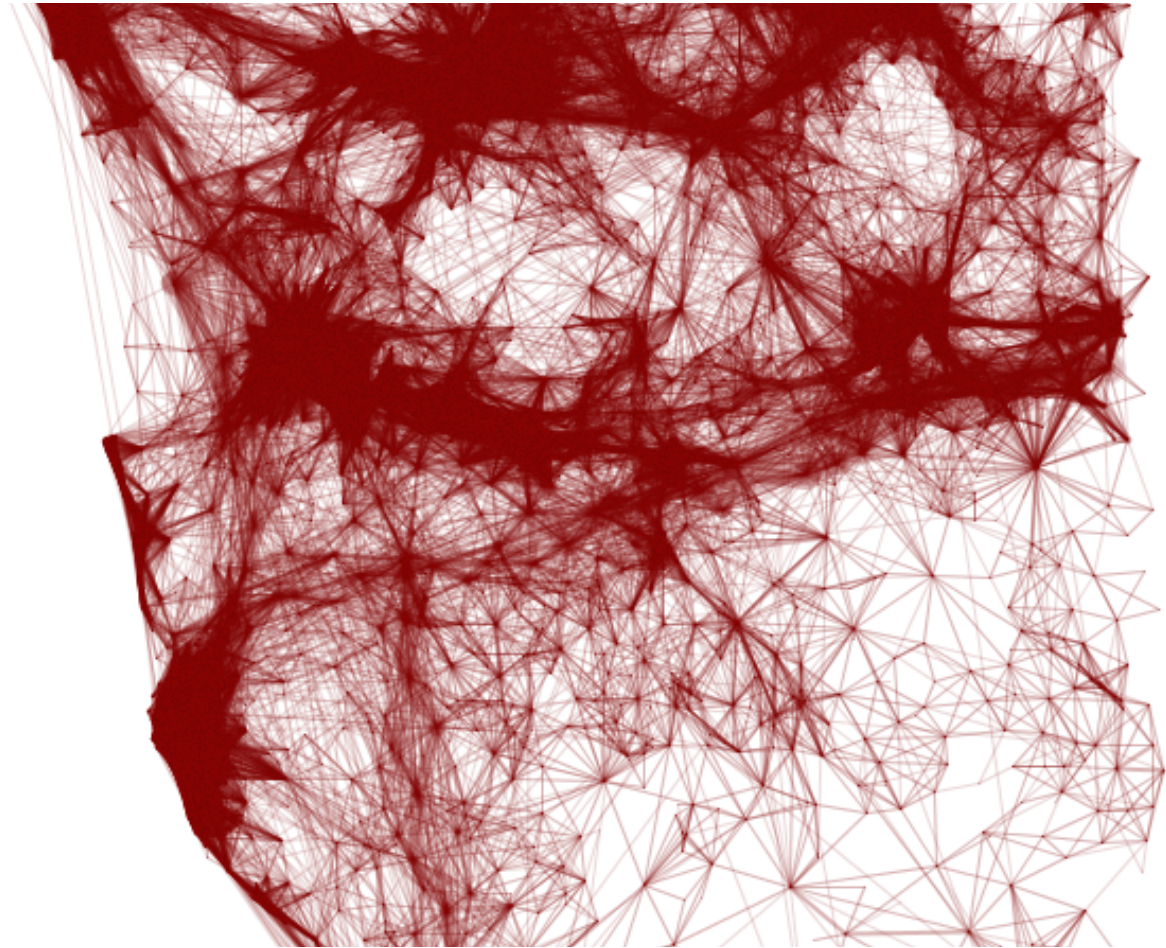




## Step 2: evaluate flows among regions

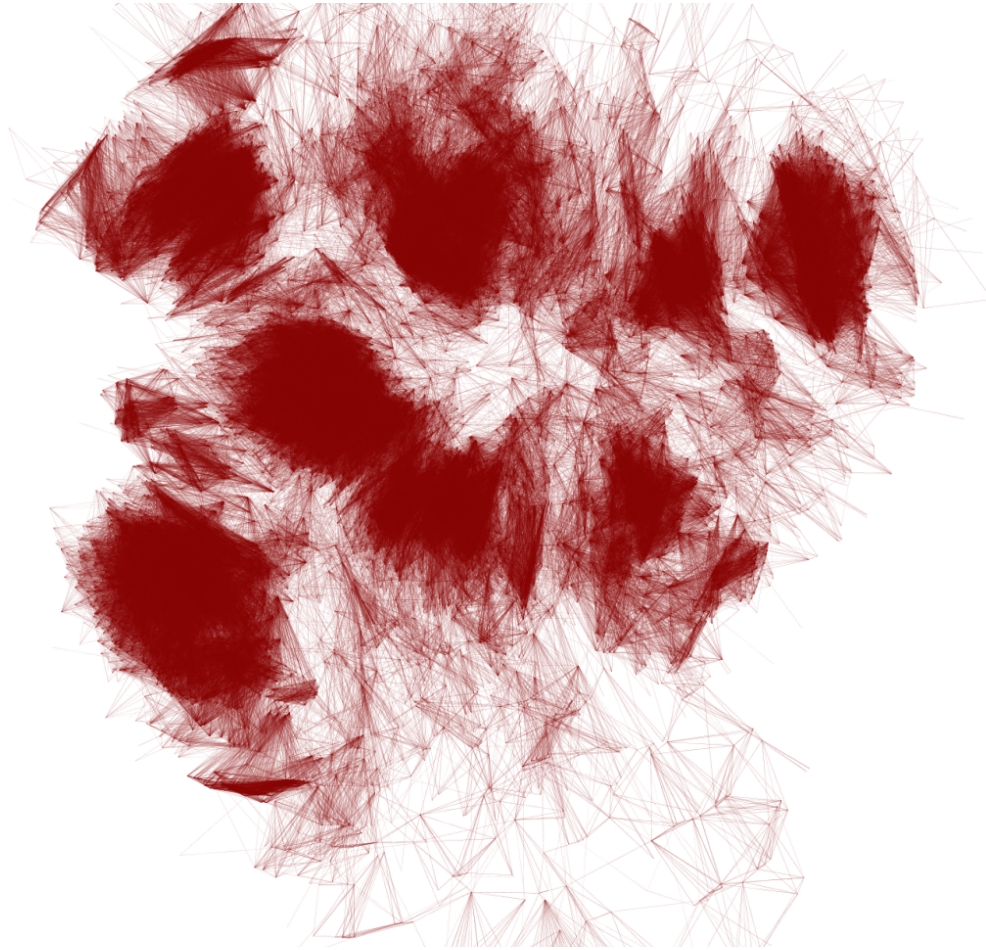


# Step 3: forget geography

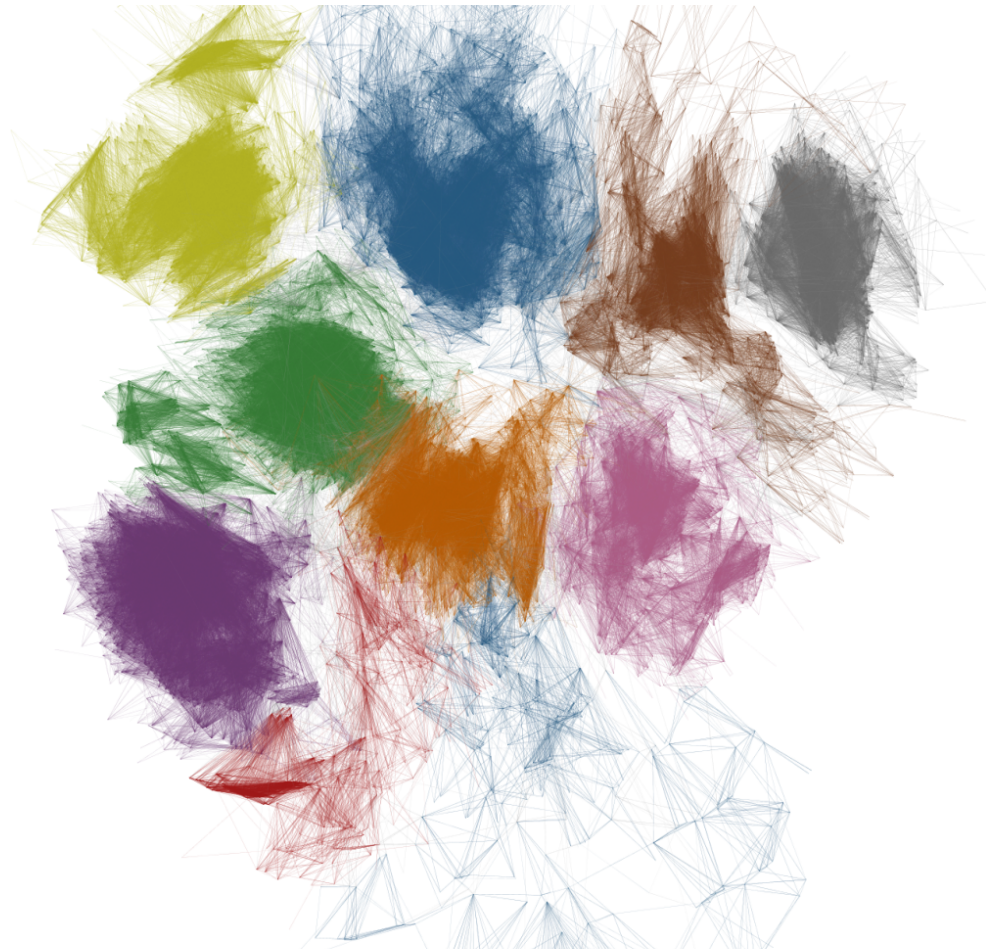




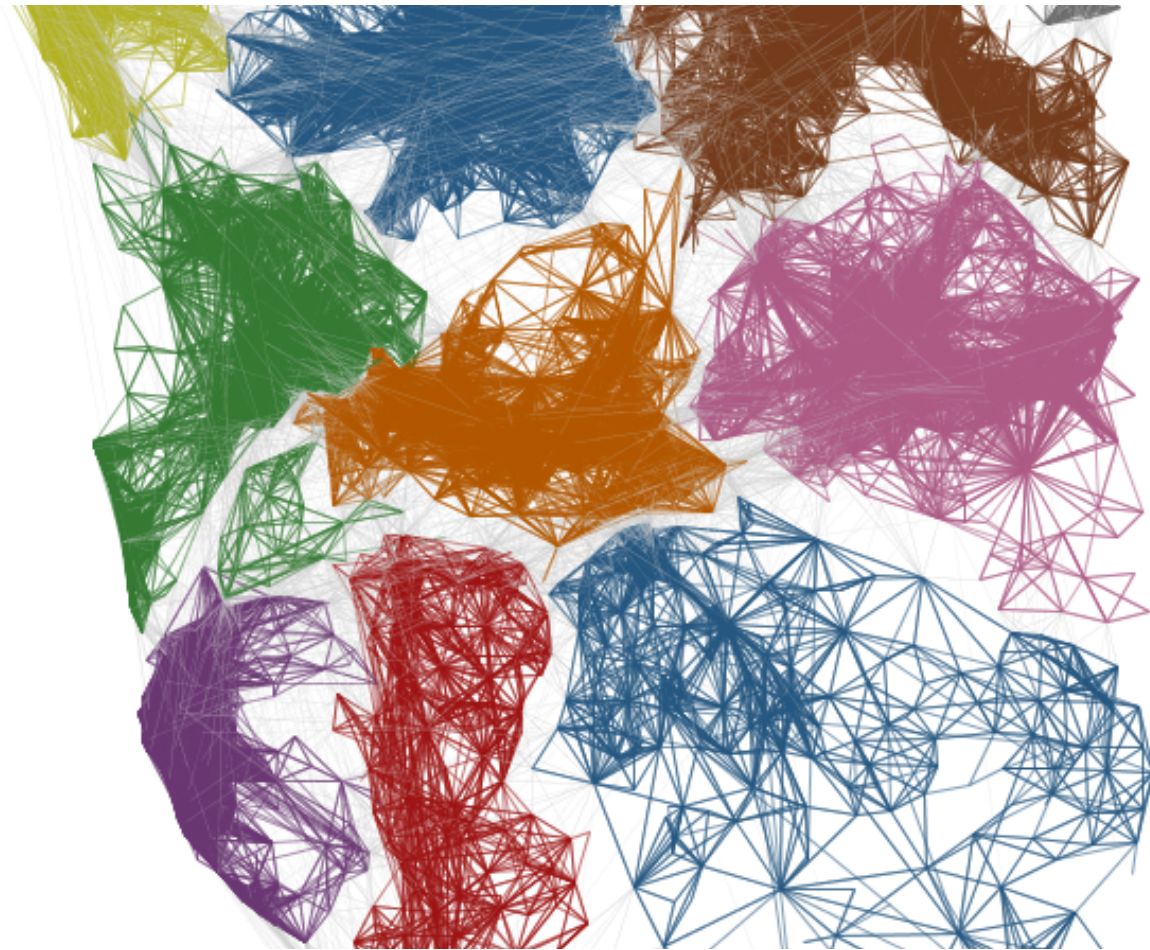
# Step 4: perform community detection



# Step 4: perform community detection

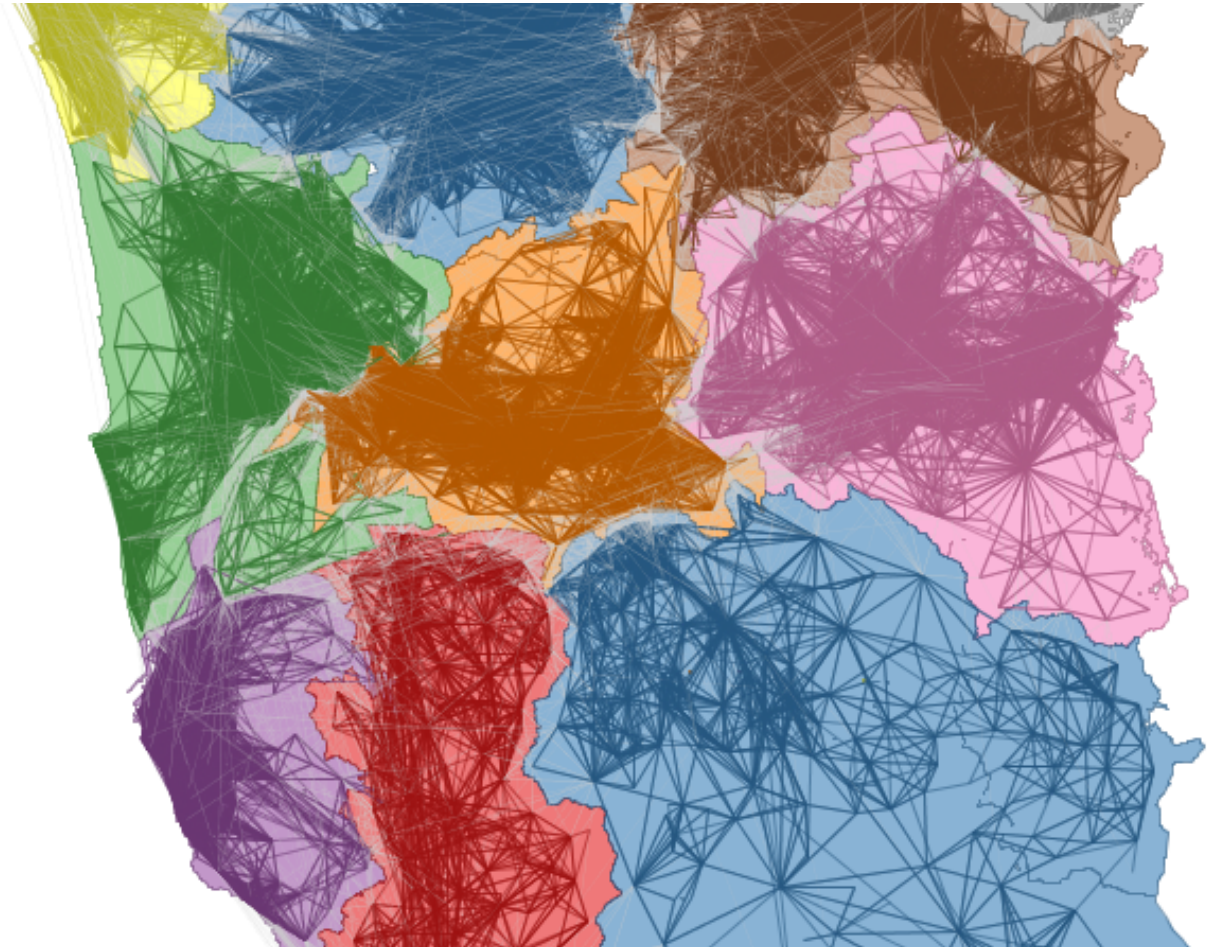


# Step 5: map back to geography

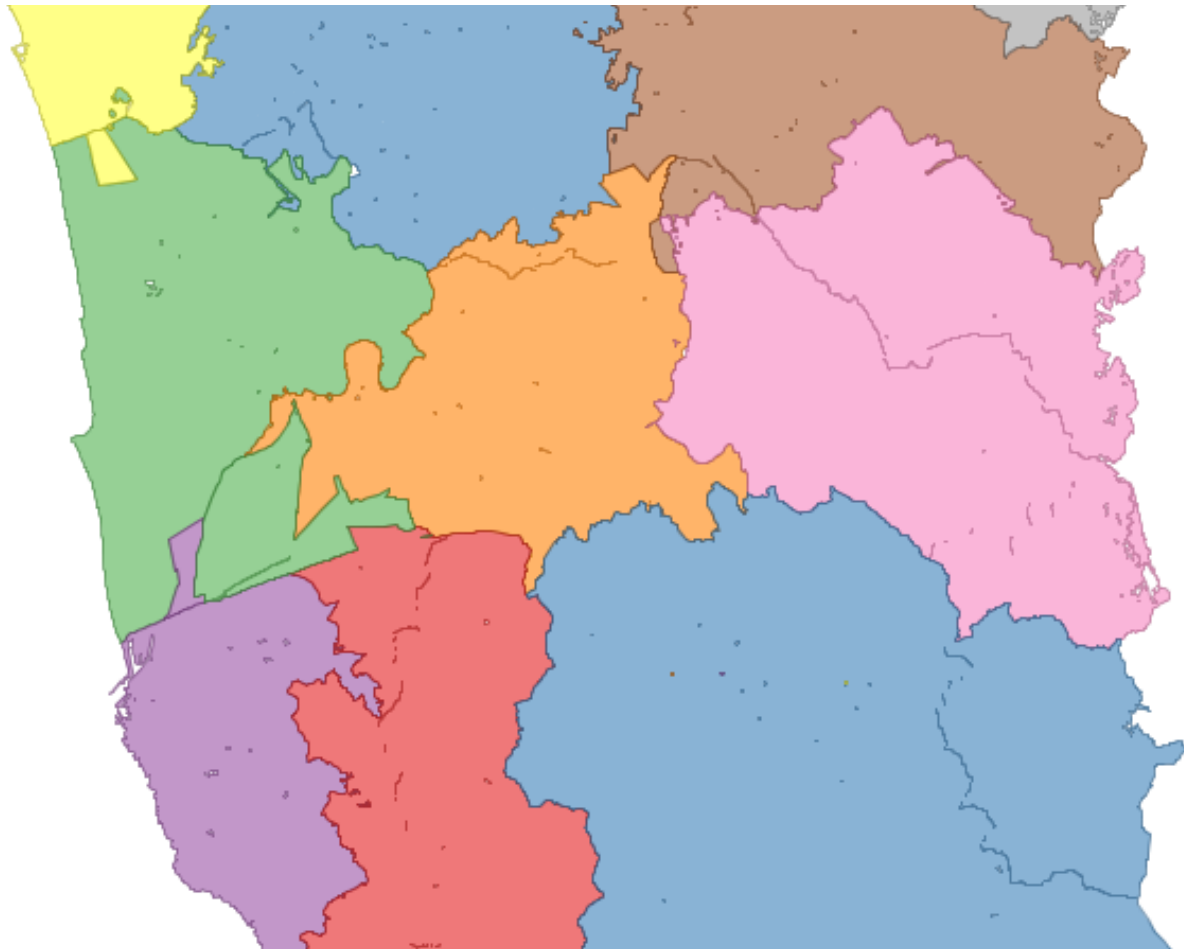




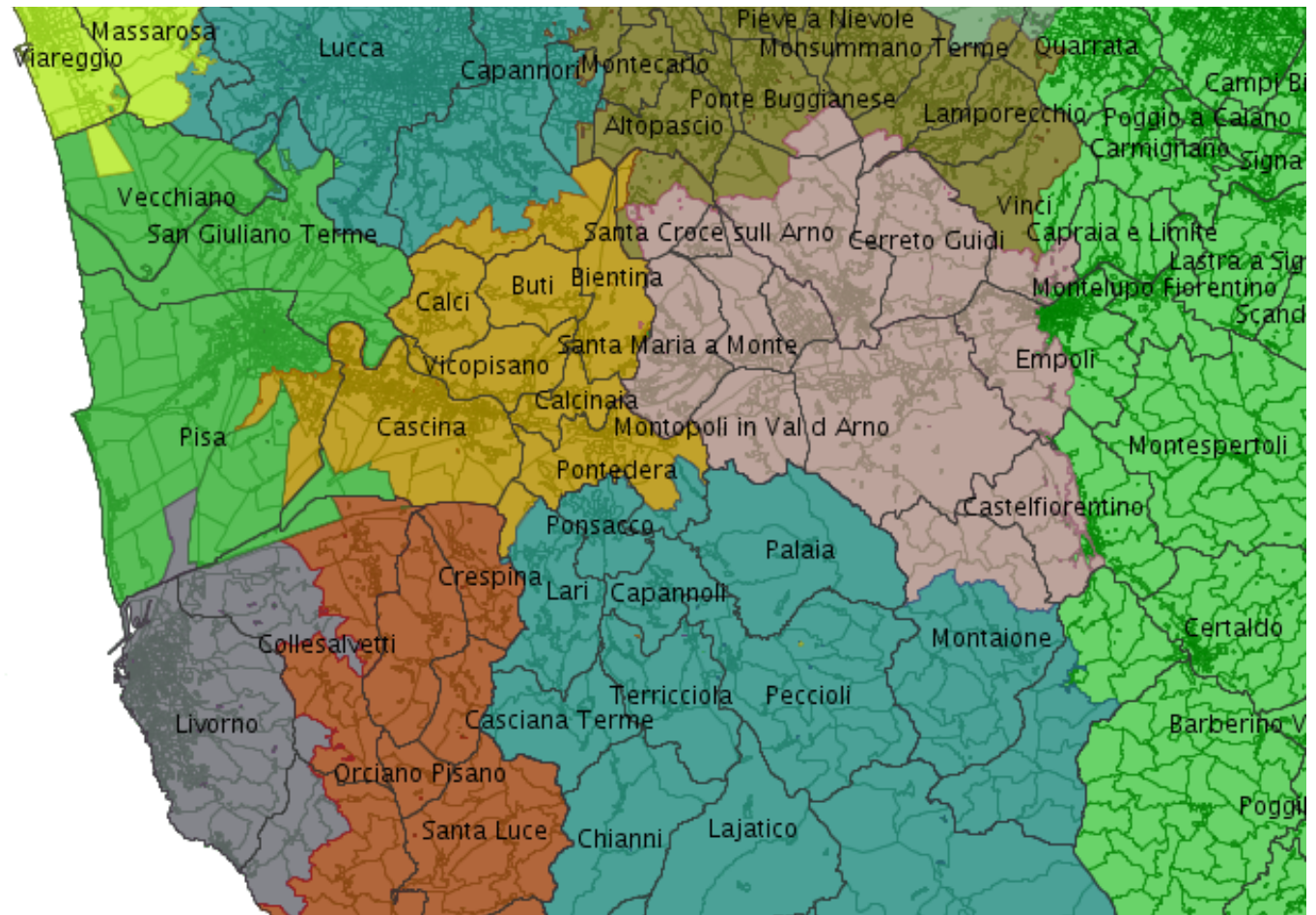
# Step 6: draw borders



# Final result

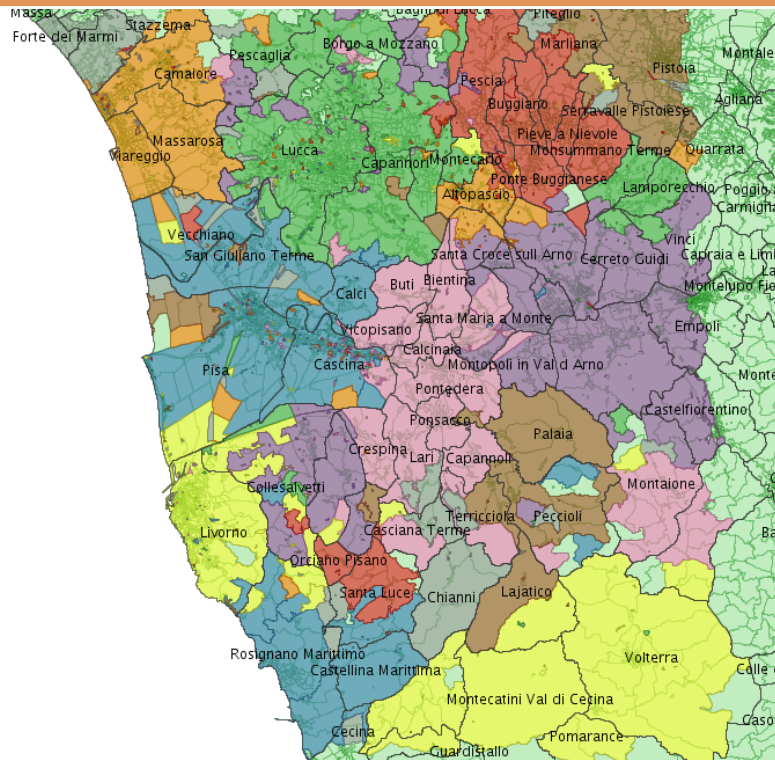


# Final result: compare with municipality borders



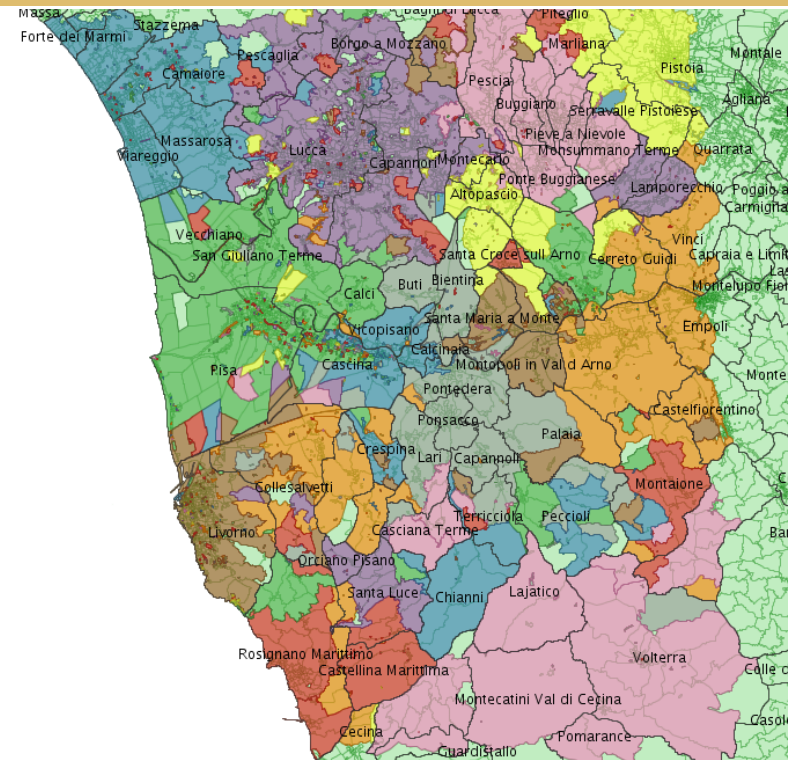
# Borders in different time periods

Only weekdays movements



Similar to global clustering: strong influence of systematic movements

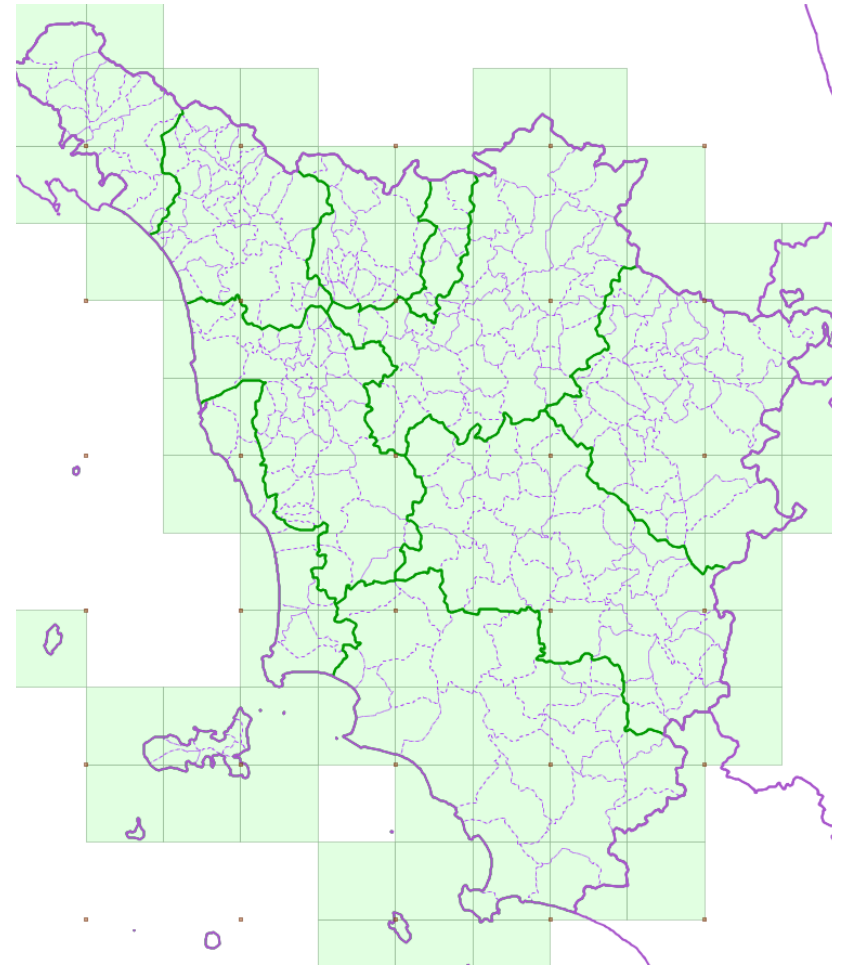
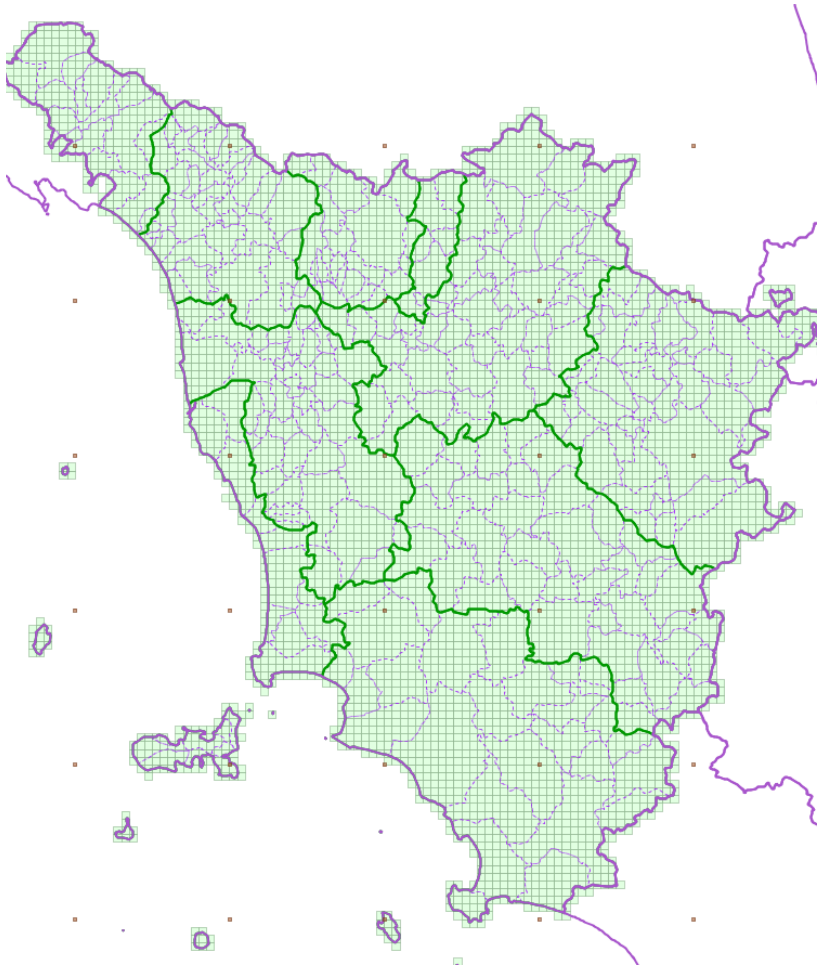
Only weekend movements



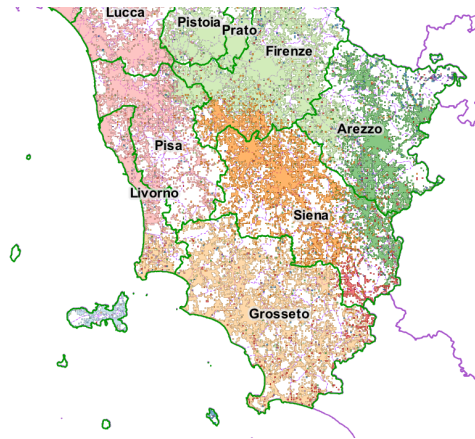
Strong fragmentation: the influence of systematic movements (home-work) is missing



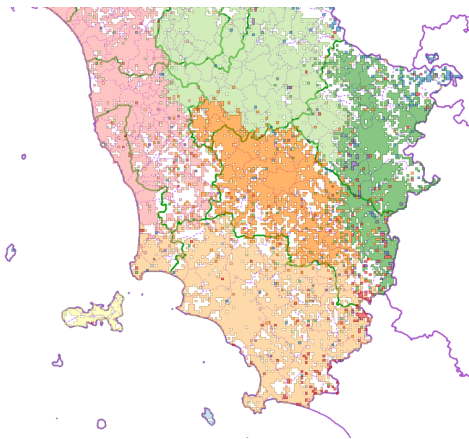
# Borders at regional scale



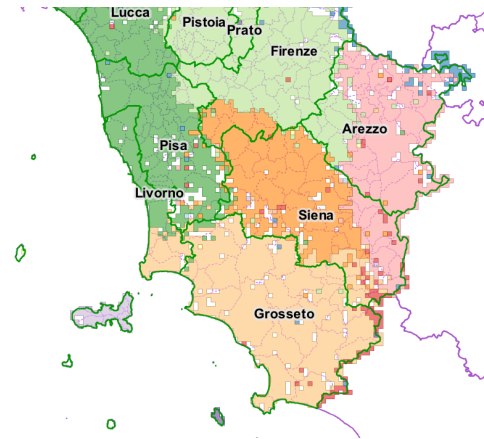
# Final results



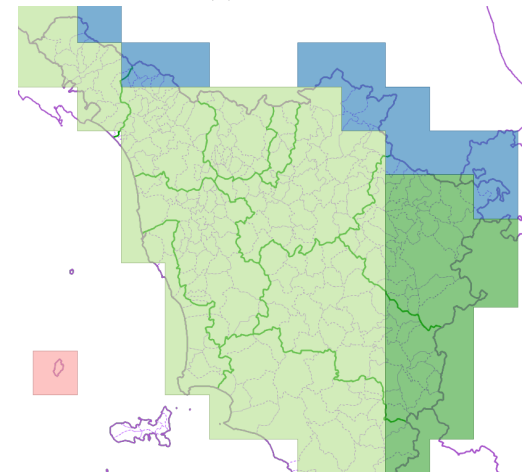
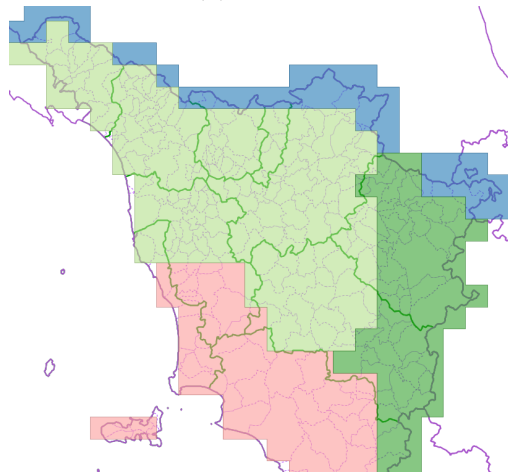
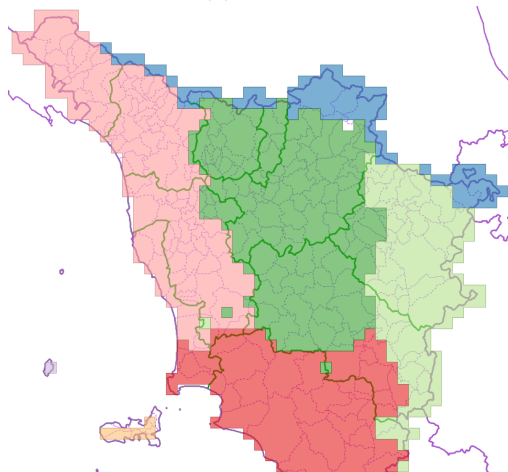
(a) 500m



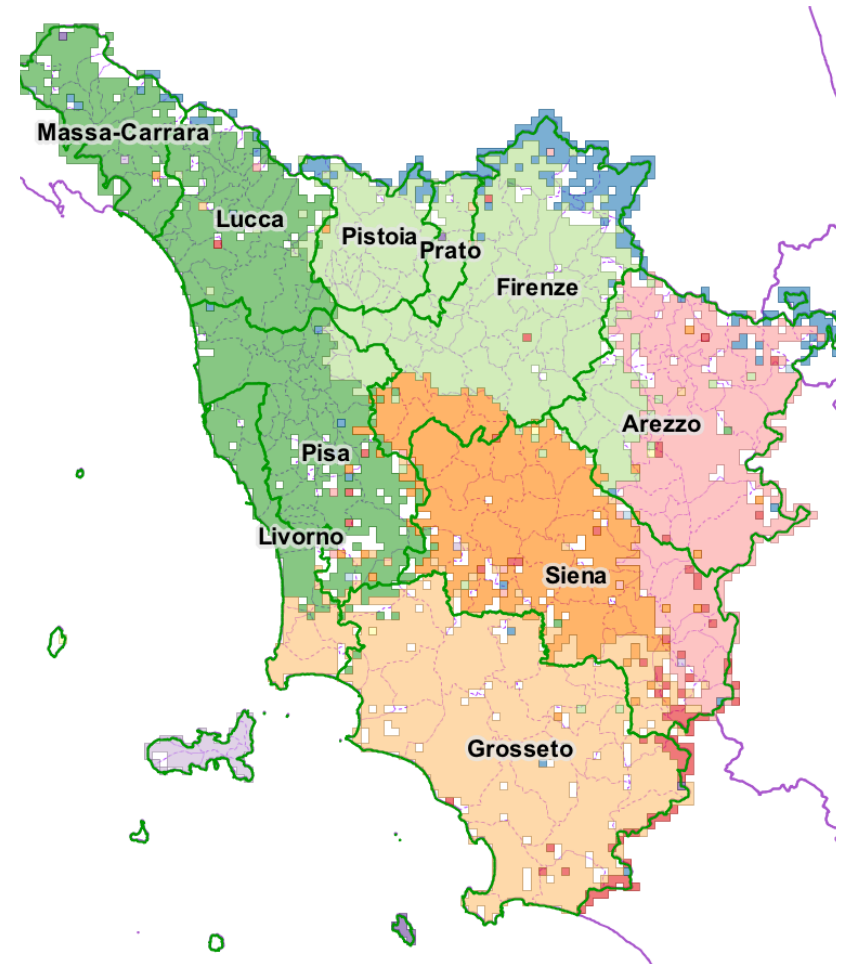
(b) 1000m

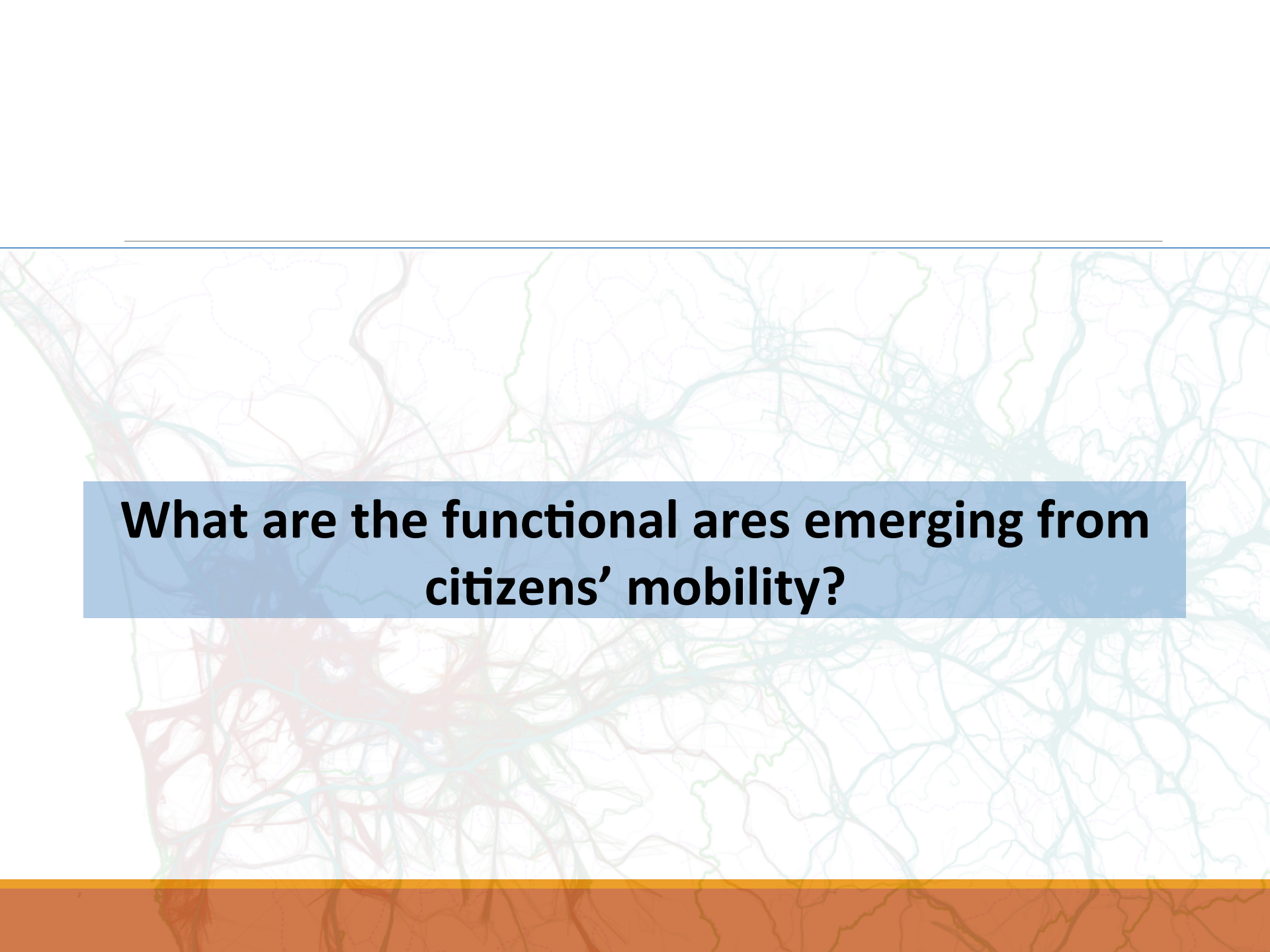


(c) 2000m



# Comparison with “new provinces”

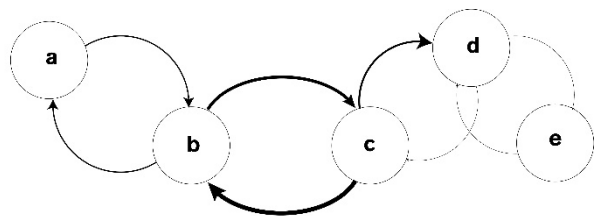


The background of the slide is a map with a complex network of colored lines in shades of red, green, and blue. These lines represent movement patterns or mobility data, with some areas showing higher density. A solid orange horizontal bar is located at the bottom of the slide.

---

**What are the functional areas emerging from citizens' mobility?**

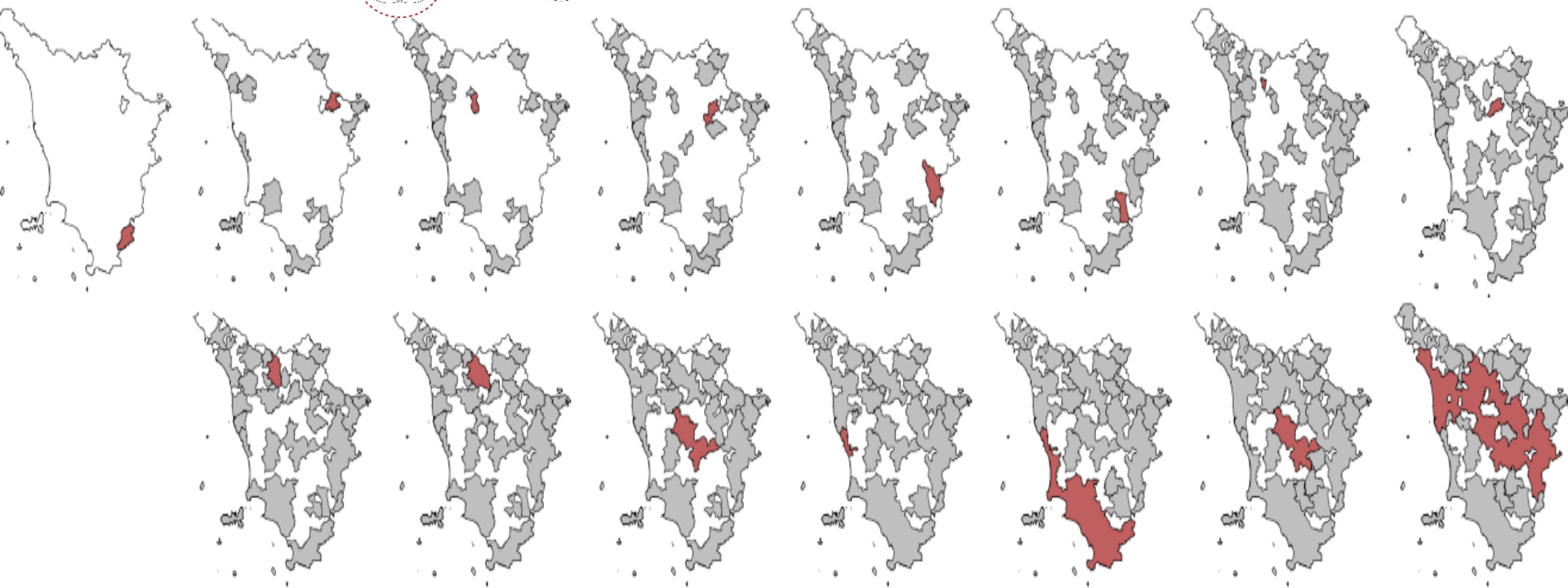
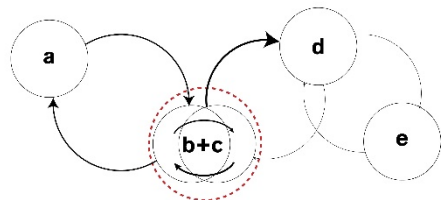




STEP 0

FIND COMMUNITIES THAT MAXIMIZE  
SELFCONTAINMENT OF MOBILITY FLUXES

STEP 1





# Cut criteria: maximise global quality score

Optimum iteration to end the algorithm is suggested by the first local maximum of S

incoming flows in  
Municipality  $j$

$$S = \sum_{i,j} F(i,j) - \sum_{i,j} F(i \rightarrow) * \frac{F(\rightarrow j)}{K}$$

Flows of vehicles  
by the municipality  $i$   
to municipality  $j$

Total outgoing flows starting  
from Municipality  $i$

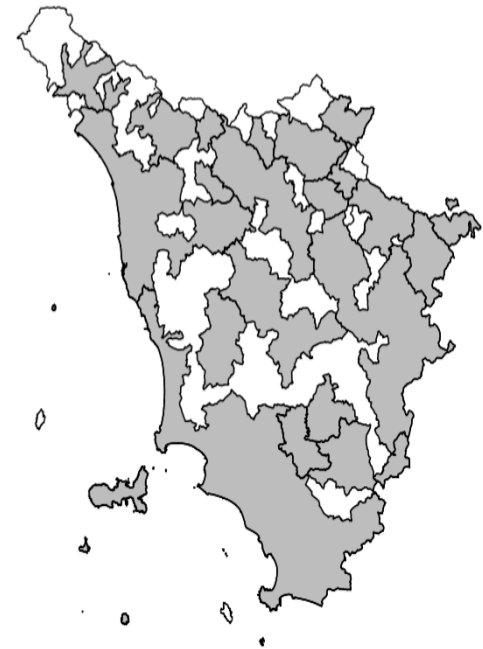
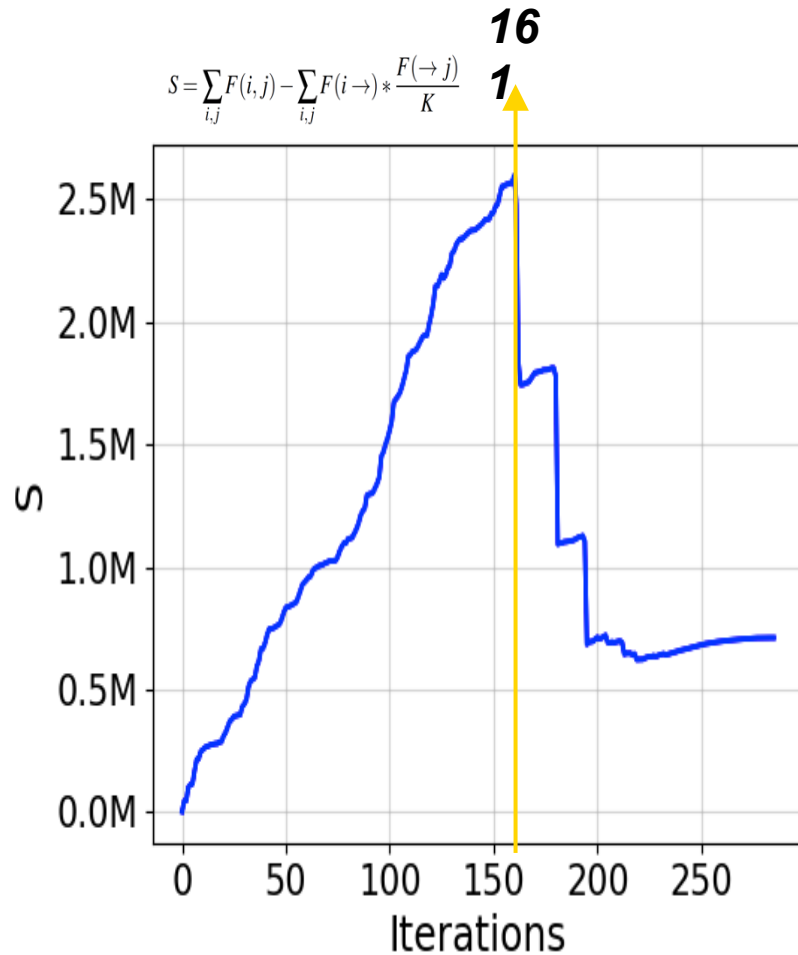
incoming flows in  
Municipality  $j$

Total flows in the network

Real

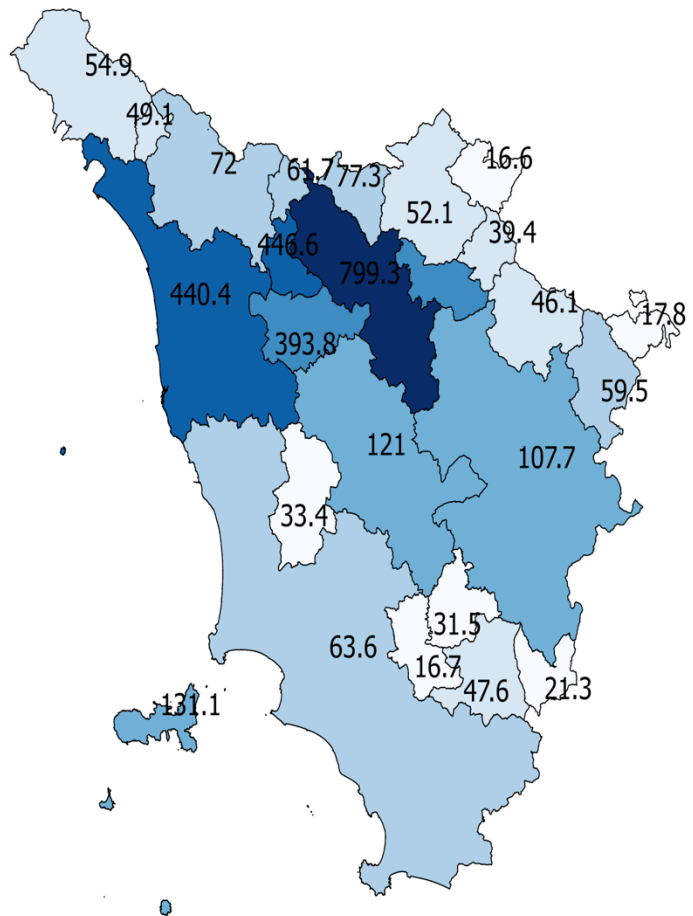
Expected

# Policentric cities RESULTS



**24** Communities discovered

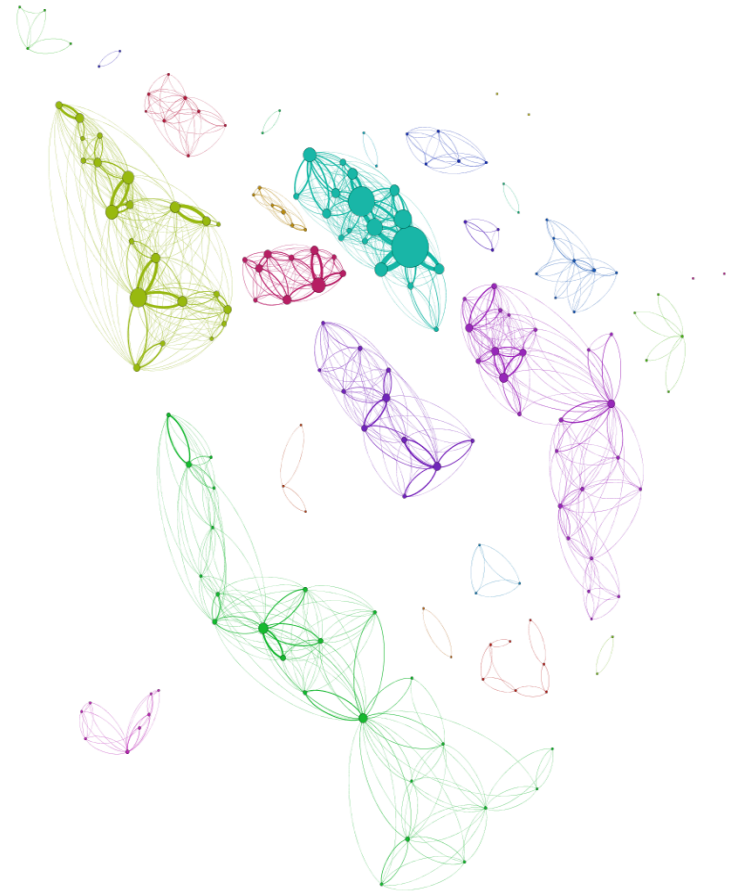
# Internal structures



## Density

population/square km

- 17 - 33
- 33 - 55
- 55 - 77
- 77 - 131
- 131 - 394
- 394 - 447
- 447 - 799

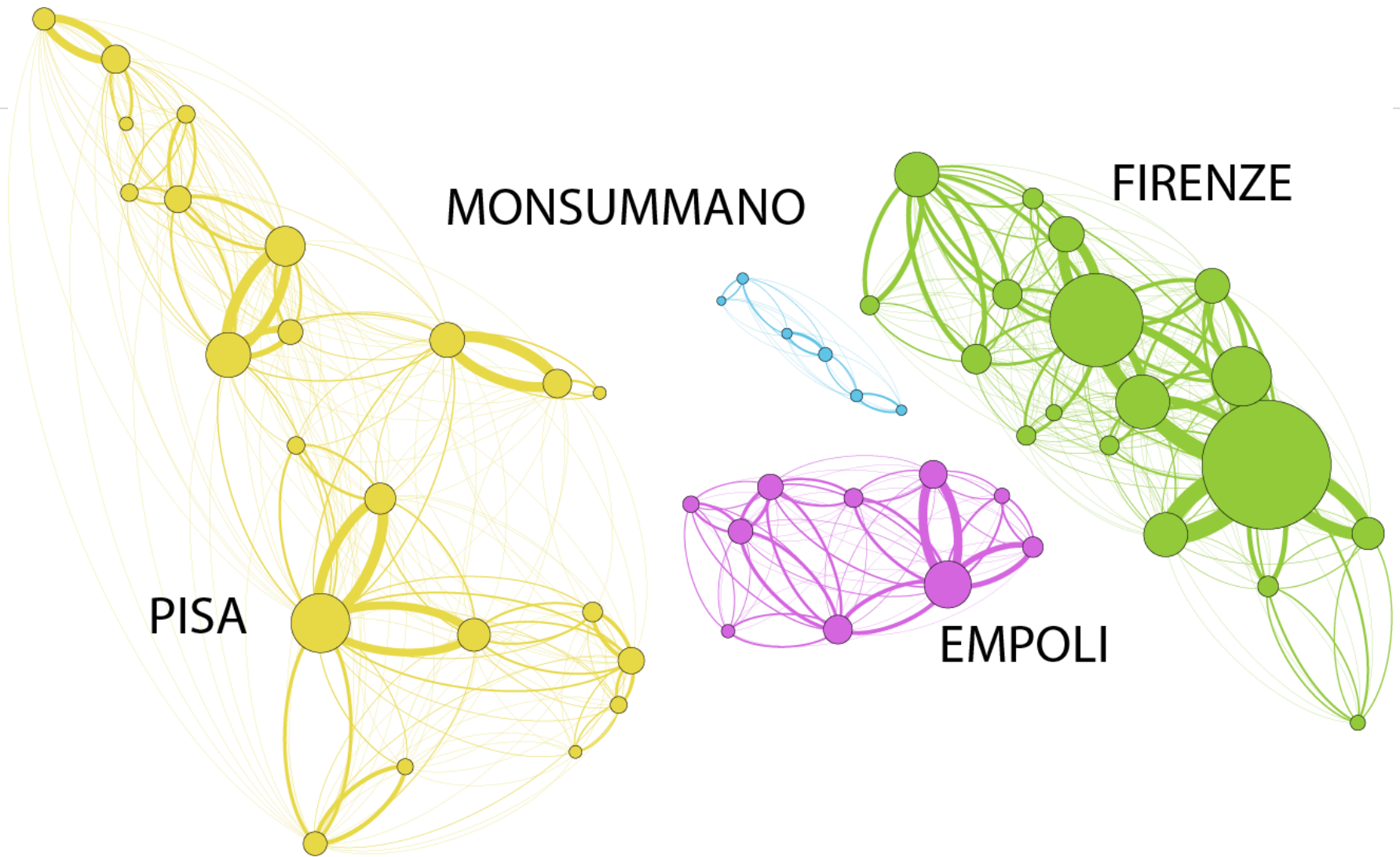


MONSUMMANO

FIRENZE

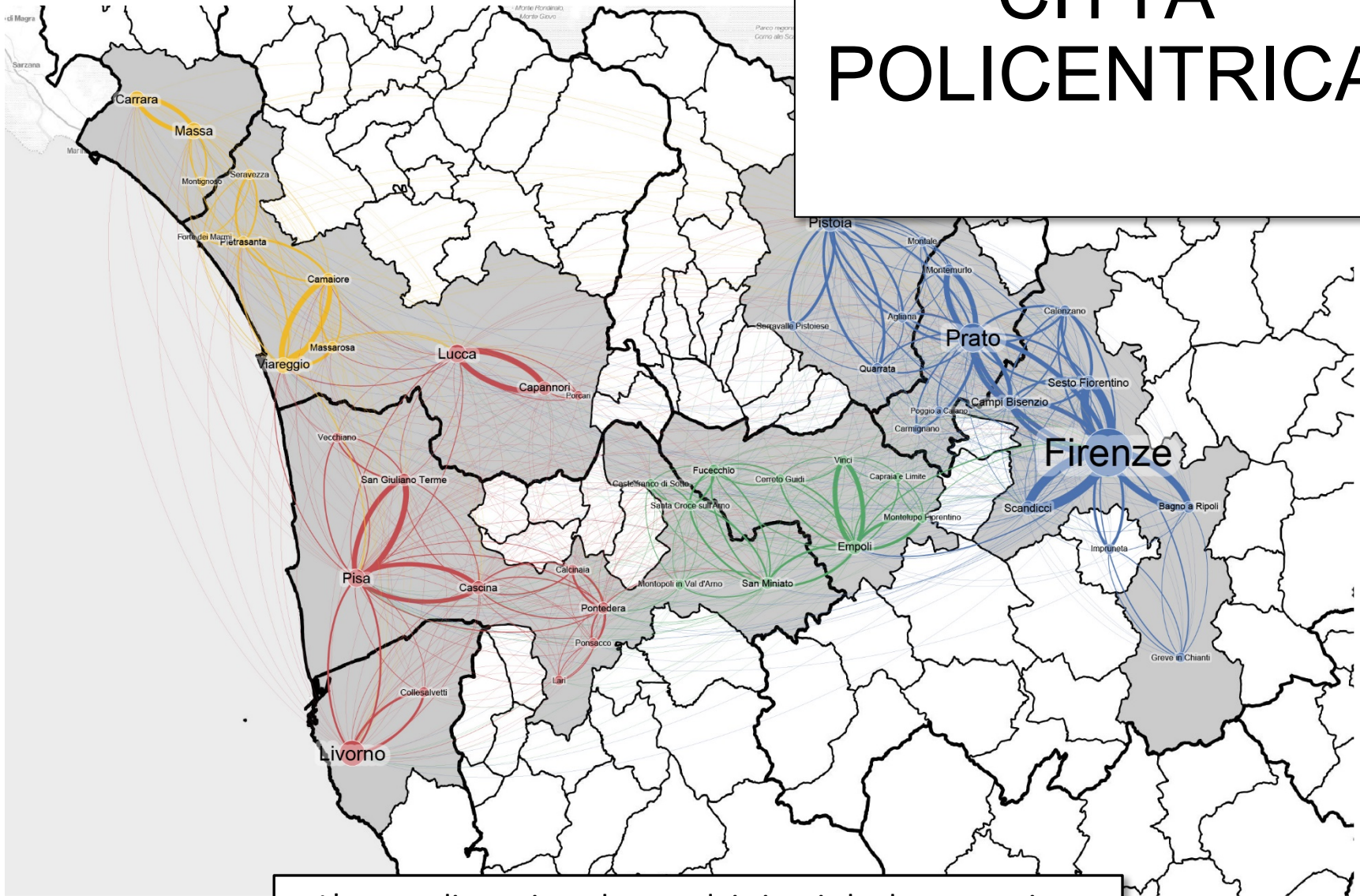
PISA

EMPOLI





# CITTÀ POLICENTRICA



L'area policentrica e la rete dei viaggi che la caratterizza