**part one**: Intro Complexity & ABM
**part two**: A case study and a tool for ABM
**part three**: Synthetic Populations & FOSSR service

Rocco Paolillo
CNR-IRPPS

rocco.paolillo@cnr.it
@roccopaolillo.bsky.social
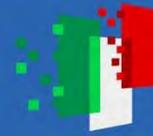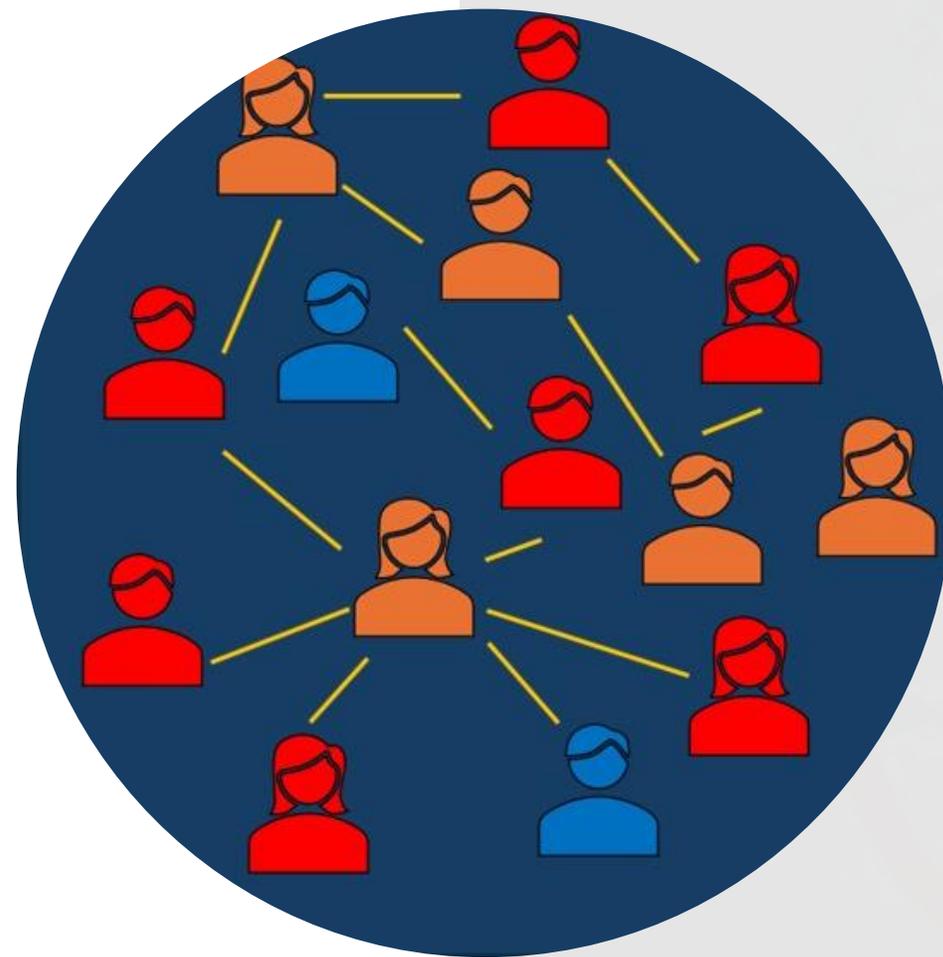
Missione 4 • **Istruzione e Ricerca**

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

Consiglio Nazionale
delle Ricerche

# Social complexity and social mechanisms

- Most of social phenomena are inherently collective phenomena
- They assume the definition of a system made of interacting components (e.g. market place, urban landscape, welfare state) where the phenomenon unfolds
- Common scope of different disciplines is to unfold the mechanisms that unfold the dynamics of the phenomenon

If the concepts of system and the goal to identify mechanisms to disclose the phenomenon is common to many disciplines, different perspectives apply

**Different perspectives on the concept of *emergence* and *mechanisms***

**Macro**

## Macro to Macro

- Markov Chain
- Evolution of processes depending on the previous state of the system as unit of analysis

**Meso**

## Meso to Macro

- Random graph models
- Evolution of networks as units of analysis

**Micro**

- Interactions of individuals and institution as the interacting components of the system

# Micro to Macro

- But sometimes, we might be interested in those mechanisms that move from the micro level, e.g. citizens/institutions with their attitudes, motivations and course of action, but get **outside of the individual agency and inglobe the interaction** of individuals as explicative mechanism of emergence

- The phenomenon is an aggregated, mutual adaptation of individuals rather than the sum of individual action

v
**Analytical Sociology
Agent-based Modeling**



Coleman Boat (1994), additions by Hedström and Ylikoski (2010), adapted

# Example of spatial segregation

- **actors:** households evaluating their neighborhood
- **individual behavior:** preference for a percentage (threshold) of similar ones in their neighborhood (homophily)
- **mechanism of emergence:** cascade effects where the behavior of one household influences the composition of neighborhood and preference of others
- **unexpected outcome:** high levels of spatial segregation, also for mild threshold preferences

**Schelling Model**
(1969,1971)



**macro constraints**
population density

**individual representation**
threshold: percentage of similar ones wanted in the neighborhood

**macro outcome**
spatial segregation

**individual action**
random relocation to an empty spot

# Agent-based Modeling

- **Simulation method** tailored to model the interacting components that constitute the system, e.g. agents representing citizens in an **artificial society**
- We can manipulate both attributes and plan of actions of agents and observe the consequences of interaction of agents executing their plans.
- By manipulating plans and conditions where the agents interact and adapt, we can experiment on and formalize the dynamics of emergence of the collective phenomenon

# Social computing with agent-based modeling

Design of the conditions, actors and initial mechanism we want to test to study the phenomenon

Formalization into rules and functions

$$if \; \vartheta < \theta : \text{leave} \; ;$$
$$if \; \vartheta \geq \theta : stay$$

Translation into code to translate the theoretical model we want to test and investigate setting-up what-if scenarios

```
set happy? similar-nearby >=
(%-similar-wanted*total-nearby / 100)

to move-unhappy-turtles
  ask turtles with [ not happy? ]
    [ find-new-spot ]
end
```

Define what-if scenarios to create experimental conditions

["density population" 70 95]
["%-similar-wanted" 0 30 60]

Collection of data as measurement of changes in the system

mean exposure to similars in the neighborhood of agents when no one relocates anymore

**KISS**
Keep it short simple, stupid
**KIDS**
Keep it descriptively simple

- A society where people differentiate by some traits
- They stay in a neighborhood if certain threshold of similarity are satisfied
- Can segregation emerge even if the threshold is not that high?
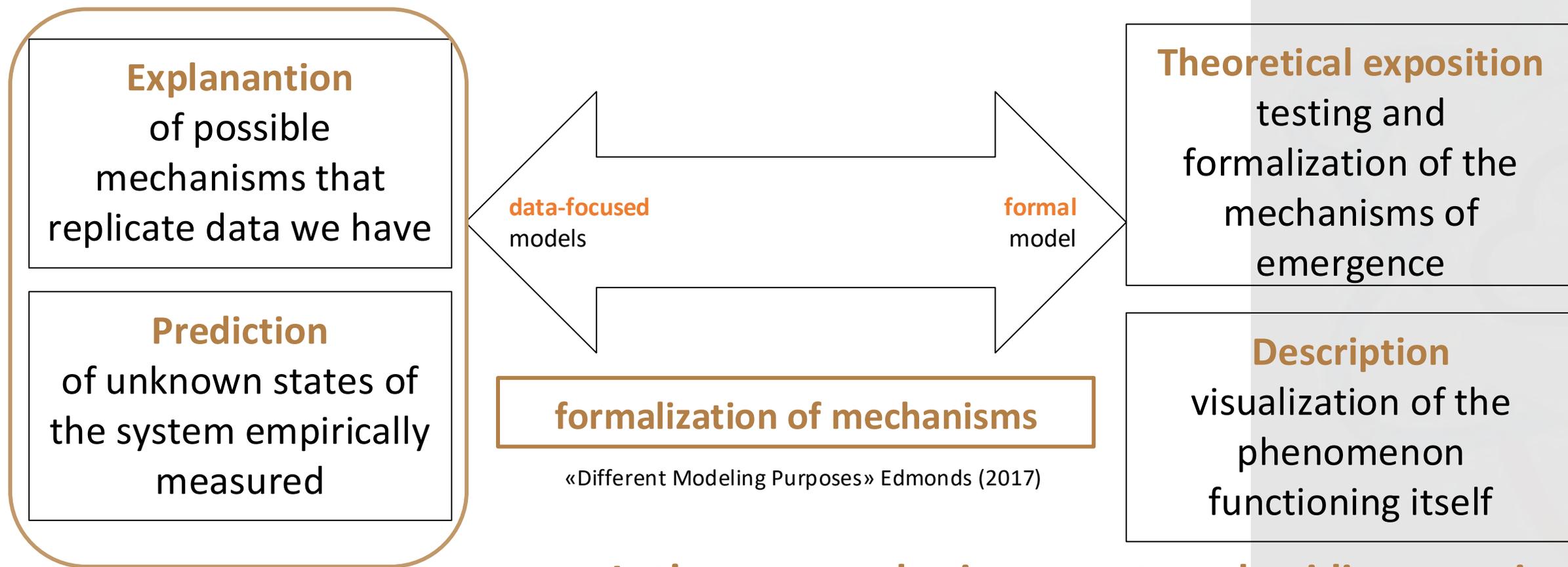
'If you can grow it, you have explained it' (Epstein, 2006)

'If you don't know how you grew it, you didn't explain it.' (Macy & Flache, 2009, p.263)

# What is agent-based modeling useful for?

**Explanantion**
of possible mechanisms that replicate data we have

**Prediction**
of unknown states of the system empirically measured

**data-focused** models

**formal** model

**formalization of mechanisms**

«Different Modeling Purposes» Edmonds (2017)

**Theoretical exposition**
testing and formalization of the mechanisms of emergence

**Description**
visualization of the phenomenon functioning itself

**Let's see some basic concepts and guiding questions**

# Agents – Who are the actors involved in the phenomenon?

**Agents**: a virtual object capable of elaborating information and able to execute an action (individuals, institutions, households...)

•**Intentionality**: acting based on goals or plans
•**Proactivity**: initiating actions rather than waiting
•**Reactivity**: responding to external stimuli or changes
•**Prosociality**: acting in coordination with others (social agents)



**State variable:** what characteristics cannot change through time? Ethnicity, gender
**Dynamic variable**: what characteristics can change through time? Opinions, preferences
**Global variable**: shared by all agents
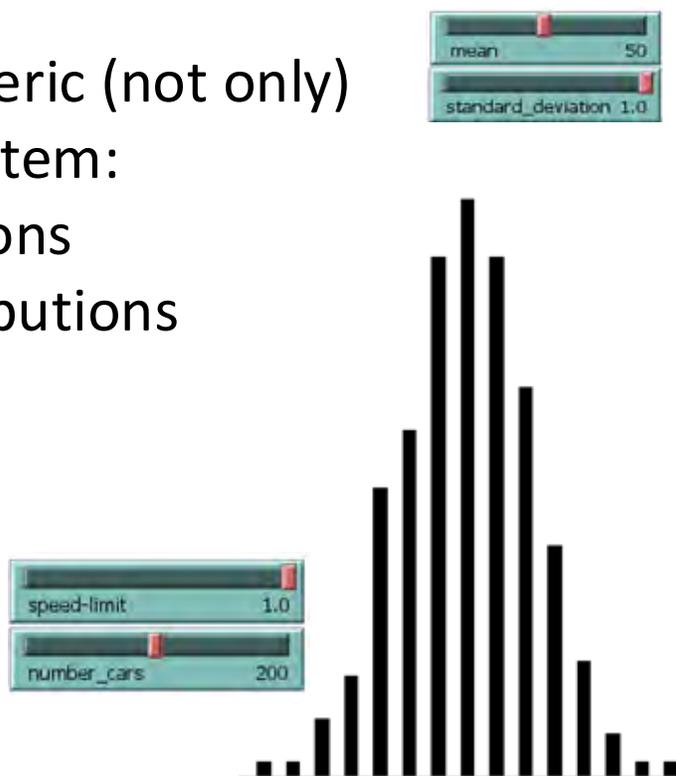**Local variable**: shared by specific agents or class of agents
**Heterogeneous** vs **Homogeneous** (attributes distribution)

**Attributes, Beliefs, Desire, Intentions (A+BDI)**

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

Consiglio Nazionale
delle Ricerche

# Parameters – What are the conditions affecting the phenomenon?

**Parameters**: tunable numeric (not only) variables to modify the system:
- calibrate initial conditions
- agents' variables distributions

mean                50

standard_deviation  1.0

speed-limit         1.0

number_cars         200

**Global parameter:** variable affecting all agents, and every agent can interact with
Belief shared by all agents

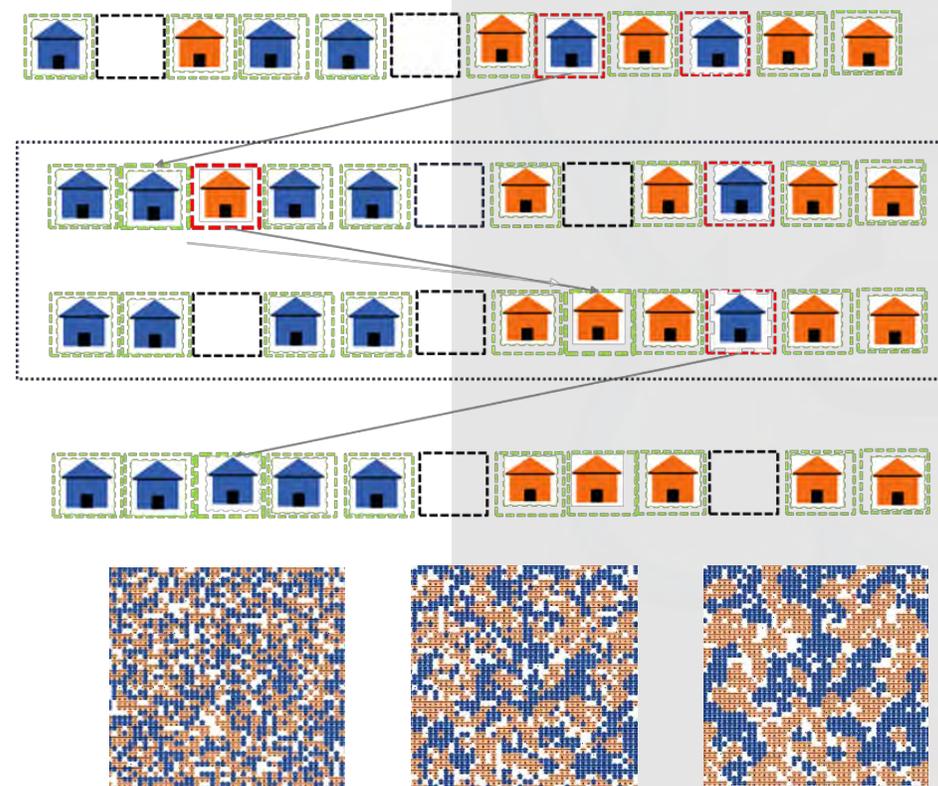**Local variable**: accessible only to some specific agents
Norm specific to a class of agents

# Evolution (I) – How the phenomen emerges through agents' interaction?

It is not much **time** as a continuous variable, rather the evolution of the system along two interconnected concepts:

- **micro level: schedule** of activation of agents' behavior
- **macro level transition phase** of the system changing due to mutual adaptation of the agents

**Cascade effect** of the behavior of one agent on the neighborhood composition to other agents, affecting segregation at macro level

# Evolution (II) – How does the order of agents' action influence each other?
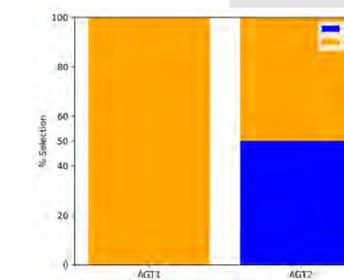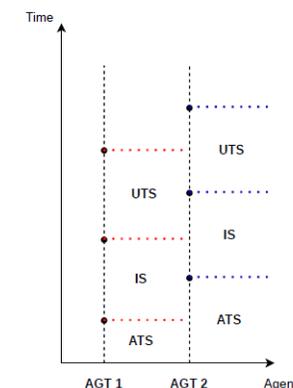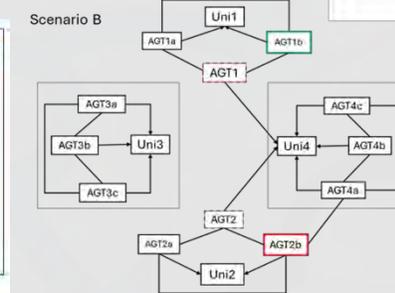
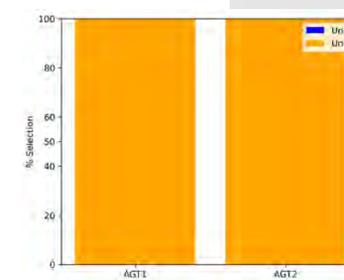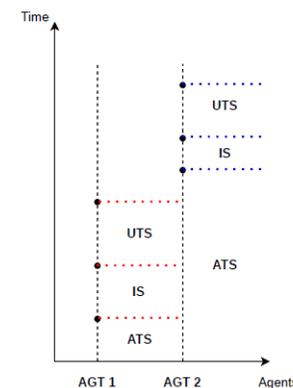**Parallelization:** how the behavior is executed

**Synchronous behavior:** agents act together in parallel

**Asynchronous behavior:** agents act sequentially
(physical threads)

**Synchrony:** when the behavior is executed

**Synchronization**: agents decide based on the same knowledge of the world, including effects of actions of others
(they act «in parallel»)

**Example**: scholar agents choose between two universities based on the chance to be introduced to elective authors based on potential shared connections. The sequential order of agents can affect the decision of those who select after
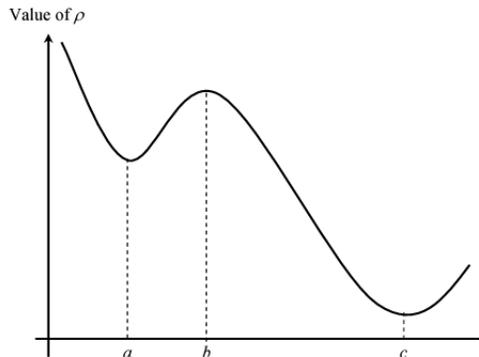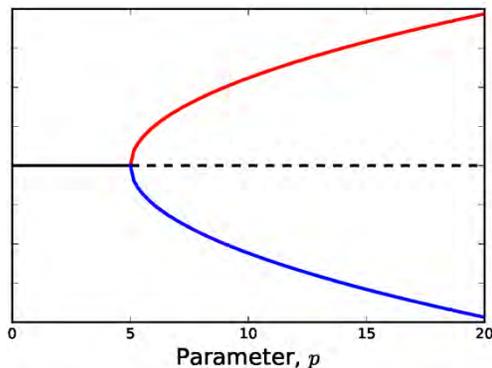
UTS > ATS: update shared knowledge

Longo, Paolillo, Ceriani (2026, forthcoming)

# Outcome (I) – How can I read the evolution of the phenomenon?
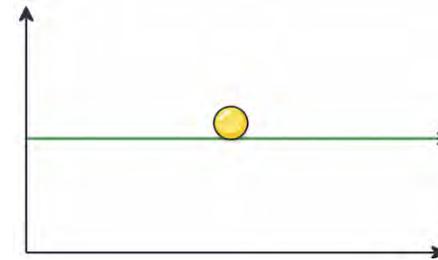
## Tipping points

Value of $\rho$



A sudden transition of the system is narrowed to one direction
e.g. in Schelling model the local level of segregation triggers relocation so that segregation becomes steady

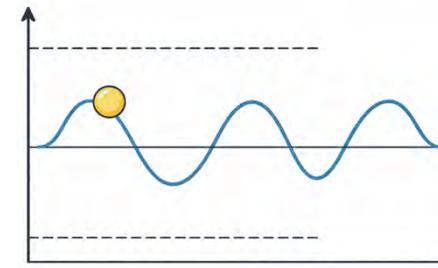**stable equilibrium**

not mutable system
(Schelling, consensu)

**cyclic equilibrium**

system follows a trend of sequential cycles
(grass & sheeps, gentrification)

## Bifurcation



A moment where the phenomenon can diverge in two opposite directions with equal probability
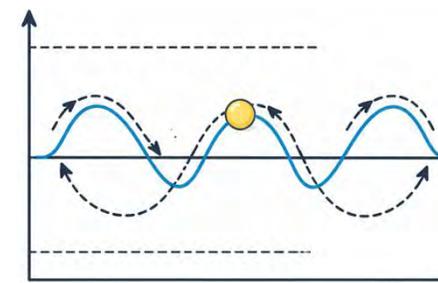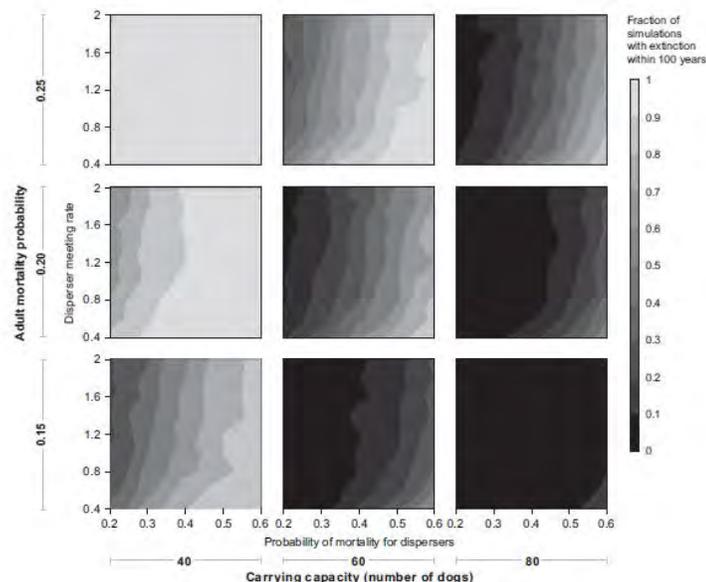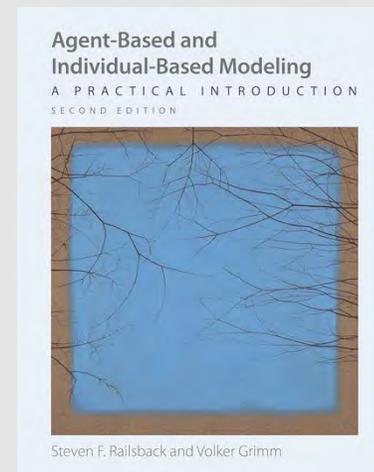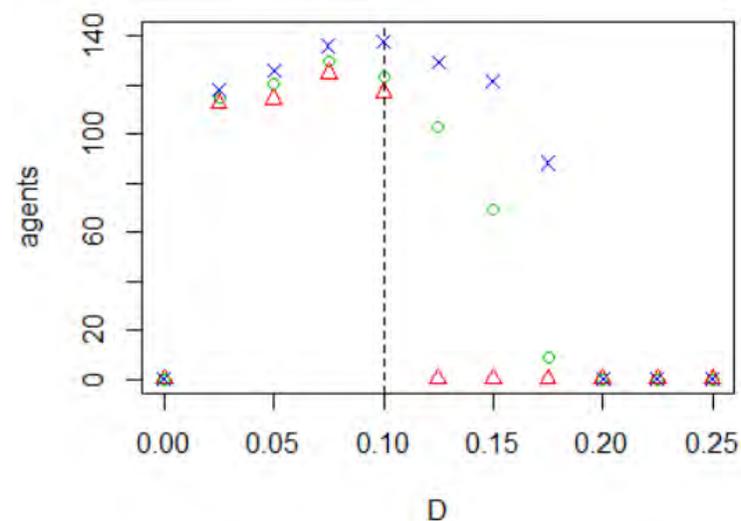e.g. in case of political polarization

**dynamic equilibrium**

oscillations/inflows/outflows causing the system to apparently remain in balance
(e.g. supply & demand)

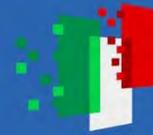# Outcome (II) – How can I identify how the conditions affect the phenomenon?

**Global sensitivity analysis**: pattern-oriented-modeling (POM) interaction between parameters (space of the model) to understand the overall mechanisms of the model

**Local sensitivity analysis**: one-factor-at-time (OFAT), focus on the effect of one specific parameter (nominal value) over the others

**Compare what-if scenarios and measures**

**Parameters**
What are the conditions affecting the phenomenon?

**Outcome**
How can I read the evolution of the phenomenon?

**Design**

**Agents**
Who are the actors involved in the phenomenon?

**Evolution**
How the phenomenon emerges through agents' interaction?

**Formalization**

**Agents**
Attributes,
Beliefs,
Desires,
Intentions

**Outcome**
How can I identify how the conditions affect the phenomenon?

**Code**

**What-if scenarios**

**Evolution**
How does the order of agents' action influence each other

macro constrains

macro-level association

macro outcome

situational mechanisms

transformational mechanisms

individual representation

individual actions

Adapt your model to the grammars of a programming language

**Let's implement to a case study and one programming tool (NetLogo)**

**ODD+D protocol to clarify ideas**

Müller et al., 2013

# References, suggested readings and tools

Macy, M. W., & Willer, R. (2002). From factors to actors: Computational sociology and agent-based modeling. Annual review of sociology, 28(1), 143-166.

Edmonds, B. (2017). Different modelling purposes. Simulating social complexity: A handbook, 39-58.

Edmonds, B., & Moss, S. (2004, July). From KISS to KIDS–an 'anti-simplistic' modelling approach. In International workshop on multi-agent systems and agent-based simulation (pp. 130-144). Berlin, Heidelberg: Springer Berlin Heidelberg.

Epstein, J. M., & Axtell, R. (1996). Growing artificial societies: social science from the bottom up. Brookings Institution Press.

Hedström, P., & Ylikoski, P. (2010). Causal mechanisms in the social sciences. Annual review of sociology, 36, 49-67.

Macy, Michael W., and Andreas Flache. 2009. "Social Dynamics from the Bottom Up: Agent-based Models of Social Interaction." In Hedström, P. and Bearman, P. (Eds.) The Oxford Handbook of Analytical Sociology. Oxford, UK: Oxford University Press.

Schelling, T. C. (1969). Models of segregation. The American economic review, 59(2), 488-493.

Schelling, T. C. (1971). Dynamic models of segregation. Journal of mathematical sociology, 1(2), 143-186.

Longo, C.F., Paolillo, R., & Ceriani, M. (Forthcoming). T2B2T: The Ontology for adaptive Agent-driven seamless integration with the Semantic Web. In Proceedings of the Joint Ontology Workshops (JOWO) - Episode XI: The Sicilian Summer under the Etna, co-located with the 15th International Conference on Formal Ontology in Information Systems (FOIS 2025), To appear on CEUR

Müller, B., Bohn, F., Dreßler, G., Groeneveld, J., Klassert, C., Martin, R., ... & Schwarz, N. (2013). Describing human decisions in agent-based models–ODD+ D, an extension of the ODD protocol. Environmental Modelling & Software, 48, 37-48.

## Softwares free

NetLogo
https://www.netlogo.org/

MESA python
https://mesa.readthedocs.io/latest/

GAMA Platform
https://gama-platform.org/

COMPLEXITY EXPLORER
SANTA FE INSTITUTE

JASSS
THE JOURNAL OF ARTIFICIAL SOCIETIES AND SOCIAL SIMULATION

essa
https://essa.eu.org/

MAILING LIST
SIMSOC@jiscmail.ac.uk

li.u
Institute for Analytical Sociology
Linköping University

rocco.paolillo@cnr.it

## Thank you! Questions?

@roccopaolillo.bsky.social

**Parameters**
What are the conditions affecting the phenomenon?

**Design**

**Agents**
Who are the actors involved in the phenomenon?

**Formalization**

**Agents**
Attributes,
Beliefs,
Desires,
Intentions

**Code**

**Outcome**
How can I read the evolution of the phenomenon?

**Evolution**
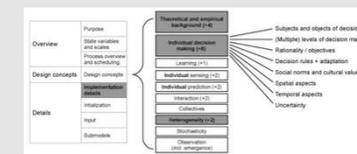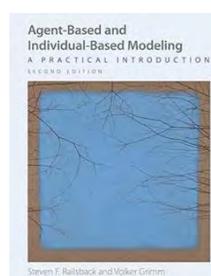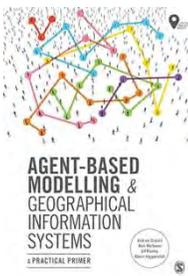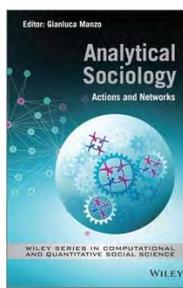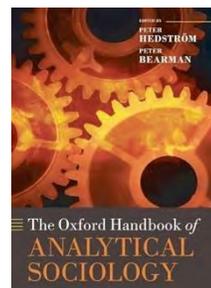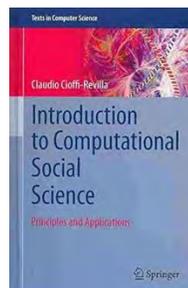How the phenomenon emerges through agents' interaction?

**Outcome**
How can I identify how the conditions affect the phenomenon?

**What-if scenarios**

| macro constrains | macro-level association | macro outcome |

situational mechanisms — transformational mechanisms

| individual representation | individual actions |

**Evolution**
How does the order of agents' action influence each other

Adapt your model to the grammars of a programming language

**Let's implement to a case study and one programming tool (NetLogo)**

**ODD+D protocol to clarify ideas**

Müller et al., 2013

# NetLogo

- Platform IDE (interfaccia) e programming language (java + starlogo) specific to agent-based modeling and experiments
- Open Source & User-friendly
- Allows many extensions (shapefile, csv import, random-wheel selection...)
- Widely used in the social science community and continuously maintained (7+ version)
- Programming language tailored to be as intuitive as possible and ready functions
- Supported by documentation (and Chat-GPT)
  - https://docs.netlogo.org/dictionary



https://www.netlogo.org/

# NetLogo



**Interface tab** to interact
**Info tab tab** to document
**Code tab** to build the model
(also in conjunction with interface)



**A command line** to interact *on the fly*

**world** where things happen

**Sliders, buttons, chooser** to facilitate interaction with parameters in the interface to explore conditions



**Plots, monitors** to detect how the phenomenon is emerging

# NetLogo

0-indexed

`list [a b c]`

show item 1 (list 1 2 4) > 2
first agent to appear has who 0

grid space world

Agents are called **turtles** identified
by **who** ID

Grid cells are called **patches** and can interact as agents

# NetLogo

**simulated synchrony:** agents execute commands asynchronously, but every agent acts knowing update in the model, based on coding

```
to ask dosomething
  ask turtles [do_A]
  ask turtles [do_B]
  ask turtles [do_C]
end
```

An agent in random order does A, then another agent does A. When all have done A, one random agent does B, then another does B. When all have done B, dosomething is executed
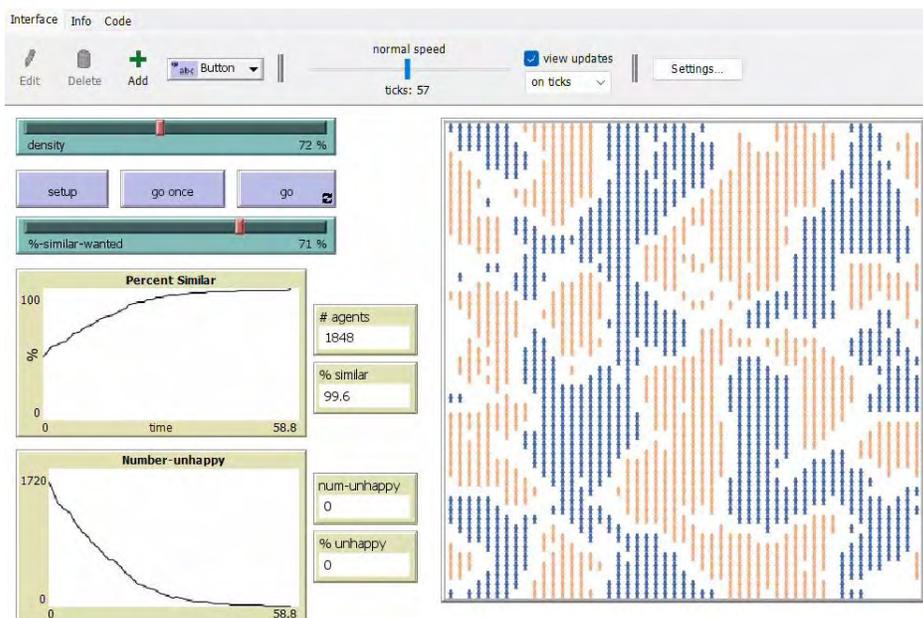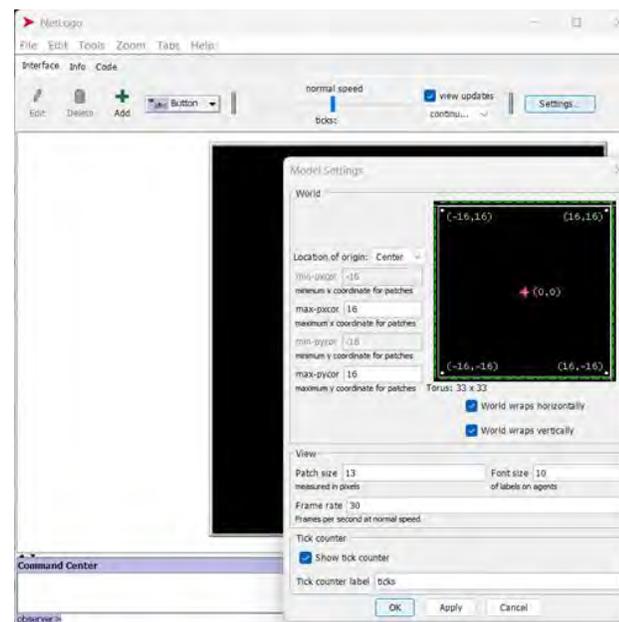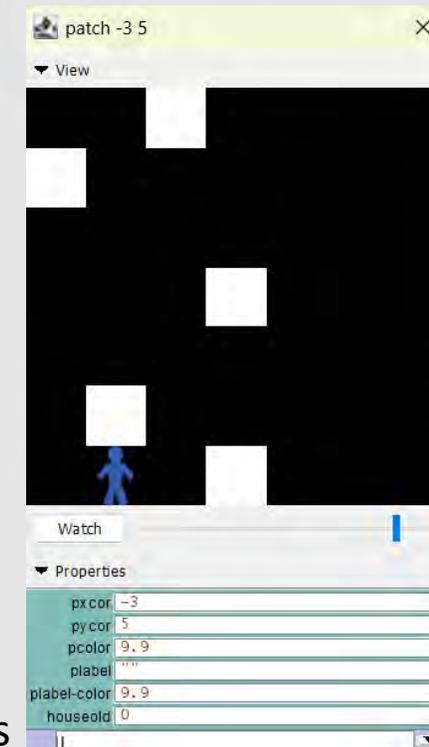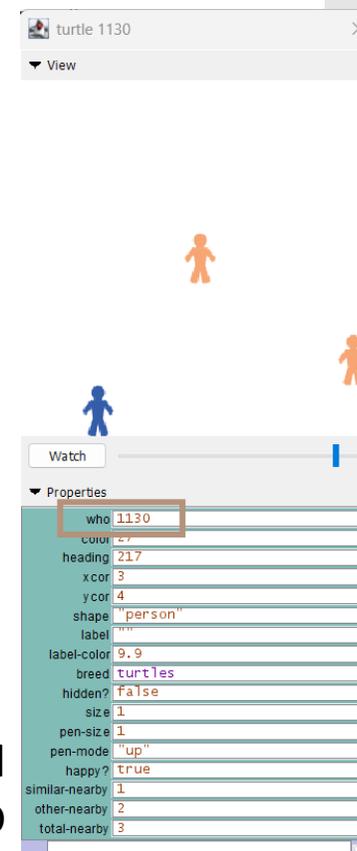
```
to ask dosomething
  ask turtles [
    do_A
    do_B]
end
```

An agent in random order does A then B, then another agent does A then B. When all have done A then B, dosomething is executed

```
to-report sumall [a b]    report 5 6 > 11
  report a + b
end
```

activate native extensions →

agent-class (breed) →

global variable →

agent-class level variable →

```
extensions [gis table csv rnd profiler]
turtles-own [PRO_COM]
breed [hospital hospitals]
breed [women womens]
breed [counselcenter counselcenters]
globals [tuscany distservices distservicesnorm]
counselcenter-own [ID capacity utility womencounsel]
hospital-own [ID hospitalizations utility capacity womenhospital mob
women-own [pregnant givenbirth selcounsel counselstay rankinglist di
```

**command block** →
that translates model components to be run

```
to setup
;  random-seed 10
  clear-all
  ask patches [set pcolor white]
  gis:load-coordinate-system "C:/Users/LENOVO/Documents/GitHub/childl
  set tuscany gis:load-dataset "C:/Users/LENOVO/Documents/GitHub/chi
  gis:set-world-envelope (gis:envelope-union-of (gis:envelope-of tus
  displaymap
```
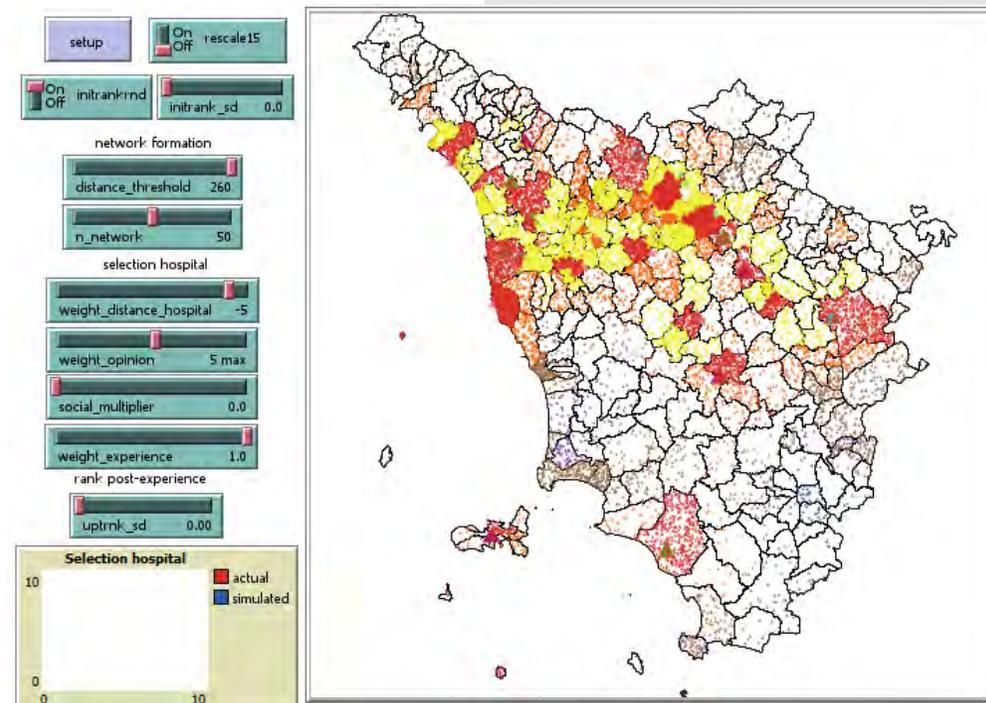
`mean`, `count`, `sort`  primitive reporters

`ask`, `set`, `forward`  primitive commands

`let` h who  local variable existing within a command block (to alleviate memory)

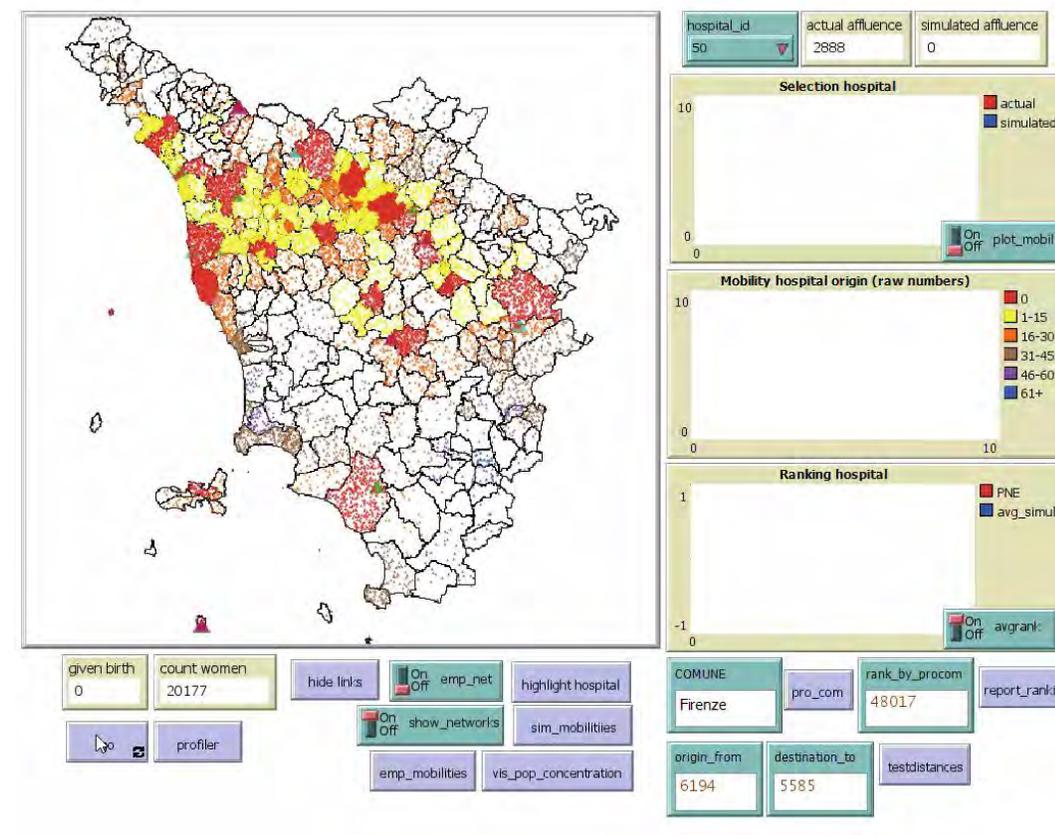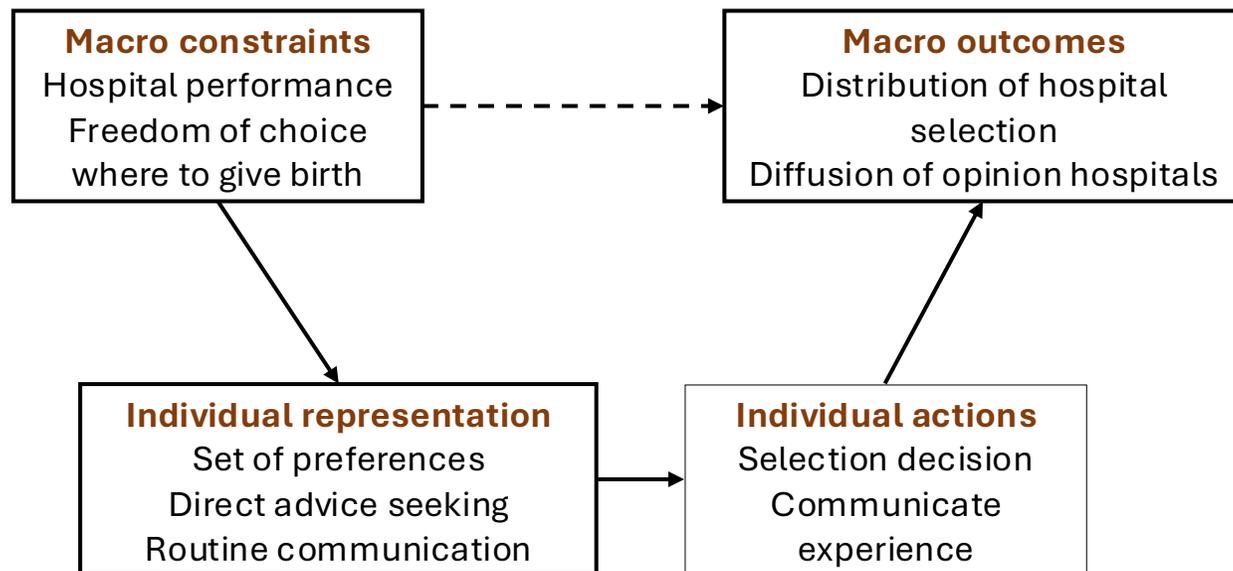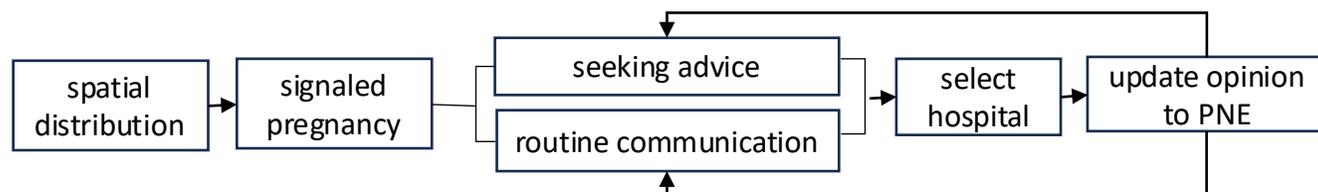**Let's implement in a case study...**

# Childbirth Mobilities: a Geo-Spatial Simulation Approach

- **Context**: While some determinants of hospital maternity selection are identified in the literature, the individual decisional processes, and social influence processes underlying the choice are unknown, neither specific data available.

- **Why ABM**: we can model the weights of preferences for hospital attributes at agents' micro-level, compare different social influence processes and compare how they replicate the data

- **Data available**: Mobility patterns in Tuscany 2023:
  - municipality residencies of women who gave birth
  (aggregated and anonymous)
  - municipality hospital where they gave birth
  - ranking of hospital (PNE performance indicator)
  - matrix of ditances
  - shapefile to map geographies to data



**Goal**

- Model a combination of individual decisional processes and social influence processes that can underline the selection of maternity hospital

- Which condition can best replicate the mobilities we observe?

Paolillo, Accordino, Pecoraro, https://github.com/RoccoPaolillo/childbirthnet/tree/MIE

spatial distribution → signaled pregnancy → seeking advice / routine communication → select hospital → update opinion to PNE

**Macro constraints**
Hospital performance
Freedom of choice
where to give birth

**Macro outcomes**
Distribution of hospital selection
Diffusion of opinion hospitals

**Individual representation**
Set of preferences
Direct advice seeking
Routine communication

**Individual actions**
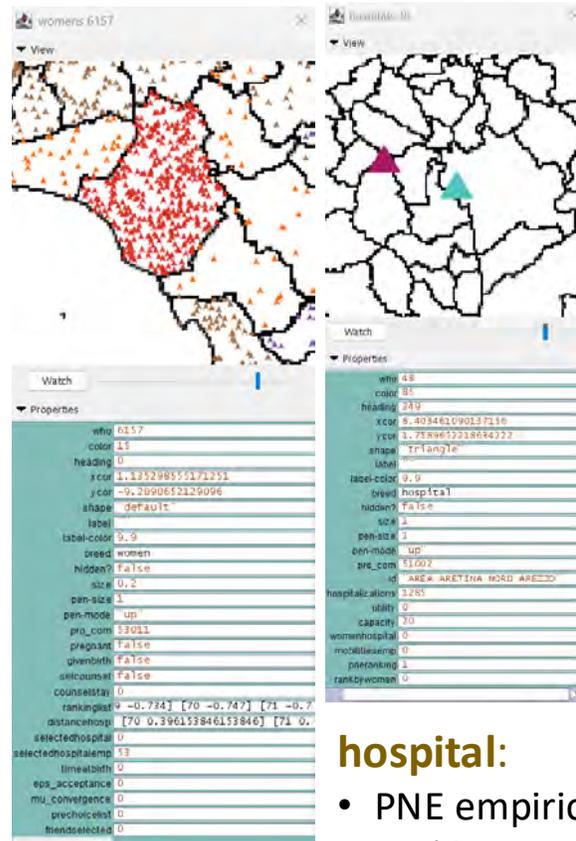Selection decision
Communicate experience

**Let's compare with concepts presented in part one...**

# Agents – Who are the actors involved in the phenomenon?

**women**:

- they hold an initial random distribution of ranking opinion for each hospital in the region, when they become pregnant activate for choosing one hospital.

- they can expressely ask advice to friends in their municipality or base on common opinion of hospitals from routine communication

- after selecting one hospital, they can vehiculate the opinion of actual performance of hospital (PNE)



**hospital**:

- PNE empirical ranking

# Parameters – What are the conditions affecting the phenomenon?

```
to setup

  clear-all
  ask patches [set pcolor white]
  gis:load-coordinate-system "C:/../comuni_consultori_2019.prj"
  set tuscany gis:load-dataset "C:/../comuni_consultori_2019.shp"
  gis:set-world-envelope (gis:envelope-union-of (gis:envelope-of tuscany))
  displaymap

  set distservices csv:from-file "C:/../matrice_distanze_consultori.csv"
  set distservicesnorm csv:from-file "C:/../normalized_distance.csv«

  create-counselcenters
  create-hospitals
  create-womens

  let sorted-hospitals sort-by [[a b] -> [hospitalizations] of a >
  [hospitalizations] of b] hospital

ask women [options_hospital]
  plot-hospitals

  reset-timer
  reset-ticks
end
```
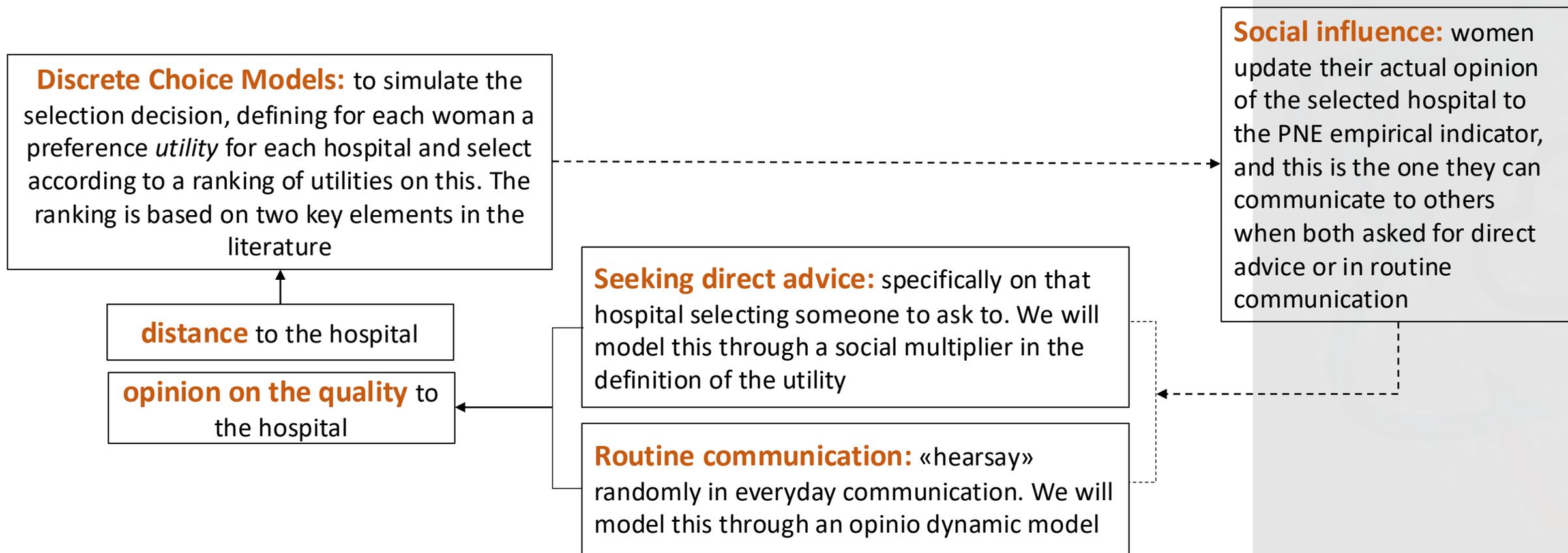
initialize the model with shapefiles

fetch data available

create agent classes

setup time schedule

# Evolution - How does the order of agents' action influence each other?

**Discrete Choice Models:** to simulate the selection decision, defining for each woman a preference *utility* for each hospital and select according to a ranking of utilities on this. The ranking is based on two key elements in the literature

**distance** to the hospital

**opinion on the quality** to the hospital

**Seeking direct advice:** specifically on that hospital selecting someone to ask to. We will model this through a social multiplier in the definition of the utility

**Routine communication:** «hearsay» randomly in everyday communication. We will model this through an opinio dynamic model

**Social influence:** women update their actual opinion of the selected hospital to the PNE empirical indicator, and this is the one they can communicate to others when both asked for direct advice or in routine communication

# Parameters – What are the conditions affecting the phenomenon?

**Discrete choice modeling:** modeling the selection decision of agents, defining a utility (U) for each hospital h, based on a weight (parameter $\beta$) of how two characteristics of each hospital are relevant to the agent:
$D_h$: distance from the agent to the hospital
$O_h$: opinion on quality

The utility is used to define a probability to select that hospital $h$ over the others hospital $k$

The higher is $\beta$, the more deterministic the selection is based on differences for attribute, the closer $\beta$ is to 0, the more the selection is random $\varepsilon$

stable softmax to avoid numerical overflow
* 10 to harmonize ß of ranking and distance due to different scales

weight_distance_hospital    -5

weight_opinion    5 max

We can input and manipulate the weight of each characteristic in the mind of agents
- equivalent to coefficients from regressions (clogit), not available
- test the consequences of combining different weights

$$U_h = -\beta(D_h) + \beta(O_h) + \varepsilon$$

```
set utility ((weight_distance_hospital *
(distancefrom*10)) + (weight_opinion * opinionquality)  )
```

$$P_h = \frac{e^{(U_h - \max(U_k))}}{\sum e^{(U_k - \max(U_k))}}$$

```
set selectedhospital [who] of rnd:weighted-one-of hospital
[exp(utility - max [utility] of hospital)]
```

**rnd:weighted-one-of** *agentset reporter*

# Parameters – What are the conditions affecting the phenomenon?

**Social multiplier:** A weight $\theta[0,1]$ in the definition of opinion quality $O_h$ at the moment of decision.

$\theta = 1$: the opinion quality completly aligns to that of people to whom asked for advice

$\theta = 0$: the advice of others is not taken into consideration

We also included a weighted average to allocate different weights to friend who gave birth to that hospital ($p$), and whose opinion $o$ is based on actual experience, and those who speak for hearsay ($a$)

$a = 1 - w$

$w = 1$, only those who gave birth influence



We can manipulate how influenced people will be by those they seek advice to

We can manipulate how many people advice is searched for advice and how far from hometown

$$\left( \frac{o_w + o_a + o_w + \cdots}{w + a + w + \cdots} \right)$$

```
foreach sort friends  [ z ->
  let weightfriend ifelse-value ([selectedhospital] of z = [who] of self)
  [weight_experience][(1 - weight_experience)]
  set totweightfriend lput weightfriend totweightfriend
  set ranking_othweight lput (table:get [rankinglist] of z [who] of self * weightfriend)
ranking_othweight]
```

$$O_h = OwnOpinion_h + \theta \left( \left( \frac{o_w + o_a + o_w + \cdots}{w + a + w + \cdots} \right) - OwnOpinion_h \right)$$

```
set opinionquality [( opinionquality + social_multiplier *
((reduce + ranking_othweight / reduce + totweightfriend) - opinionquality ) )]
```

# Parameters – What are the conditions affecting the phenomenon?

**Opinion dynamics:** a method to model the routine communication from women who gave birth to others in their municipality, spreading the own (updated) opinion $a$ of that hospital. The receiver agent $i$ accepts to listen if the distance between the own opinion of the hospital and that of the sender falls below a latitude of acceptace $|o_t^i - o_t^a| \leq \varepsilon$.
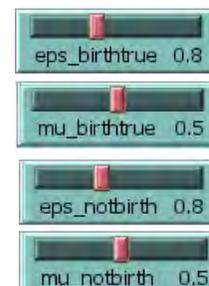If so, the receiver aligns to the sender with convergence $\mu$

$\varepsilon$ = 0, not communication occurs
$\varepsilon$ = 1, everyone is listened
$\mu$ = 0, not influence occurs
$\mu$ = 1, complete alignment occurs

We set to communicate every 80 time steps and to random 10% of women in municipality

eps_birthtrue   0.8
mu_birthtrue   0.5
eps_notbirth   0.8
mu_notbirth   0.5

We can manipulate how available to listen to those who gave birth and to what degree they will be influenced already in the *hearsay* routine communication

$$if \ |o_t^i - o_t^a| \leq \ \varepsilon$$
$$o_t^i = o_{t-1}^i + \mu(o_{t-1}^a - o_{t-1}^i)$$

```
ask alter [
if abs(table:get rankinglist topic –
table:get [rankinglist] of myself topic) <= eps_acceptance
[table:put rankinglist topic
( table:get rankinglist topic +
(mu_convergence * (table:get [rankinglist] of myself topic –
table:get rankinglist topic)))]
```

# Outcome – How can I identify how the conditions affect the phenomenon?

**BehaviorSpace:** a tool provided by NetLogo to set many experiments to run independently, setting the conditions for each parameter, define specific report measures, how many repetition wanted, and collect data in csv file

- Tools > BehaviorSpace
- Supports batch mode (headless)
- Better with a Server! For computational power, can run on laptop anyway



['eps_birthtrue' 0 2]
will run the conditions with the variable set 0 and 2
['eps_birthtrue' [0 0.1 2]]
will run all the conditions with the variable set from 0 to 2 in increments 0.1
(e.g. 0 0.1 0.2...1.9 2)

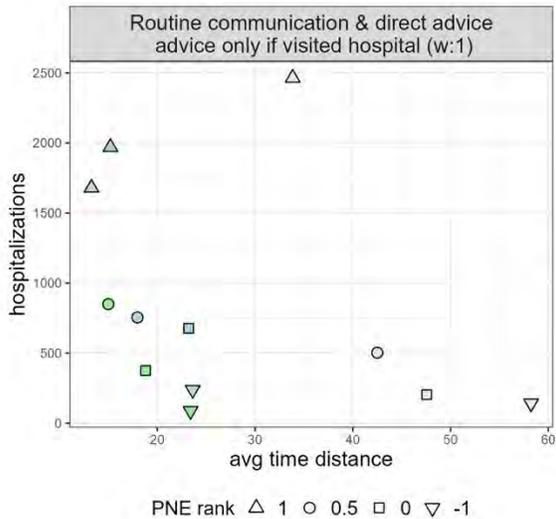# Outcome – How can I identify how the conditions affect the phenomenon?



**What-if scenarios**

- Women communicate with everyone in routine communication ($\varepsilon$ = 2), updating their opinion to the actual experience of those who gave birth, but they do not seek for advice (social multiplier $\theta$ = 0).
- We manipulate the weights for distance [0 -1,-5] and opinion quality updated via routine communication [0 1 5]
  - **condition A**: only distance matters (opinion weight 0)
    - with minimal weight of ditance ($\beta$ = -1), women select hospitals more sparsely and difference between rankings do not emerge. Increasing weight distance ($\beta$ = -5), the simulated distribution overestimates proximity of selected hospitals, and still sort by ranking not appearing
  - **condition B**: when we include also high opinion weight ($\beta$ = 5), simulation results better approximate empirical data when coupled with high weight of distance (green condition $\beta$ = -5, $\beta$ = 5), both distance-wise and ranking-wise

**What would be the effect of seeking advice then?**

# Outcome – How can I identify how the conditions affect the phenomenon?



**What-if scenarios**

- Now women only seek for advice at the moment of selection (latitude opinion dynamic $\varepsilon = 0$), they can be influenced only by those who actually experience the hospital ($w$ = 1) or by everyone equally ($w$ = 0.5)

- Being influenced by those who gave birth ($w$ = 1), hospitals with high PNE are overestimated, and more when the selection is more random by distance ($\beta$ = -1)
- Being influenced with equal weight by everyone, those with actual experience and those with random opinion, underestimates the match with empirical data instead

**What if we combine the two types of social influence?**

# Outcome – How can I identify how the conditions affect the phenomenon?



**What-if scenarios**
- Now women undergo both types of influence. When seeking advice, they can be influenced only by those who actually experience the hospital ($w$ = 1) or by everyone equally ($w$ = 0.5)

- Being influenced by those who gave birth ($w$ = 1), still shows overestimation of hospitals with high, and more when the selection is more random by distance ($β$ = -1)
- Being influenced with equal weight by everyone, those with actual experience and those who updated opinion by *hearsay* in common routine, better approximates the empirical data and reduces overestimation (sligthly)

# Outcome – How can I read the evolution of the phenomenon?

## Conclusions

- The best approximation to empirical mobilities is due to a combination of preference for shorter distance and high opinion quality. But high opinion quality with different weight of distance doesn't produce the same effect. So, distance seems more relevant, and conditioning the diffusion of opinion updates

- Concerning the two modalities of social influence, seeking for advice would overestimate the effect of PNE ranking of hospitals, since the difference in opinion quality becomes more salient. The effect is higher if agents relocate randomly in space, probably because more likely to find high PNE hospitals, that are more and in more populated areas.

- Being exposed to different opinions when seeking advice and in combination with routine communication ameliorates the polarization effect of ranking coming closer to the empirical data, and routine communication seems to suffice

# Outcome - How the phenomenon emerges through agents' interaction?

**Limits and Next Steps**

- To better understand the actual differentiation between seeking advice modality and routine communication, looking at the evolution through time and wider parameter space
- To differentiate action of women and measures considering the actual microspace within the region

## BUT

- Overall, we had quite amount of data here
- Sometimes information on sociodemographic population might be missing, how could we do?
- **Synthetic Populations**

**Thank you! Questions?**

rocco.paolillo@cnr.it

@roccopaolillo.bsky.social

**Let's see what synthetic populations are...**

# A service for synthetic populations extraction...

Creation of an *Italian Open Science Cloud for the Social Sciences* guided by *Open Science* principles

which shall provide **innovative tools and services** to investigate issues related to the **economic and societal change of contemporary societies** through the enhancement of **research infrastructures**

https://www.fossr.eu/

**Ex-ante** policy analysis evaluation
**Post-ante** policy analysis evaluation
**Conterfactual** policy scenarios



Luciana Taddei
Mario Paolucci *Editors*

**Longitudinal Data Infrastructures in Europe**

Tools for Open Science in Social Science Research

OPEN ACCESS    Springer

FOSSR
Fostering Open Science in Social Science Research
Innovative tools and services to investigate economic and societal change

Chapter 12
**Synthetic Populations in Research Infrastructures**

Rocco Paolillo, Nicholas Roxburgh, Alice Sbrana, Gary Polhill, Evelina Carmen Sabatella, and Mario Paolucci

**12.1  Collective Phenomena and Social Complexity**

- An artificial society is a stylized social system where to study the mechanisms of the phenomenon

- Especially when agent-based modeling is used for policy purposes, the mechanisms observed need to be bounded to the conditions of the system they want to operate in, e.g. Digital Twin systems
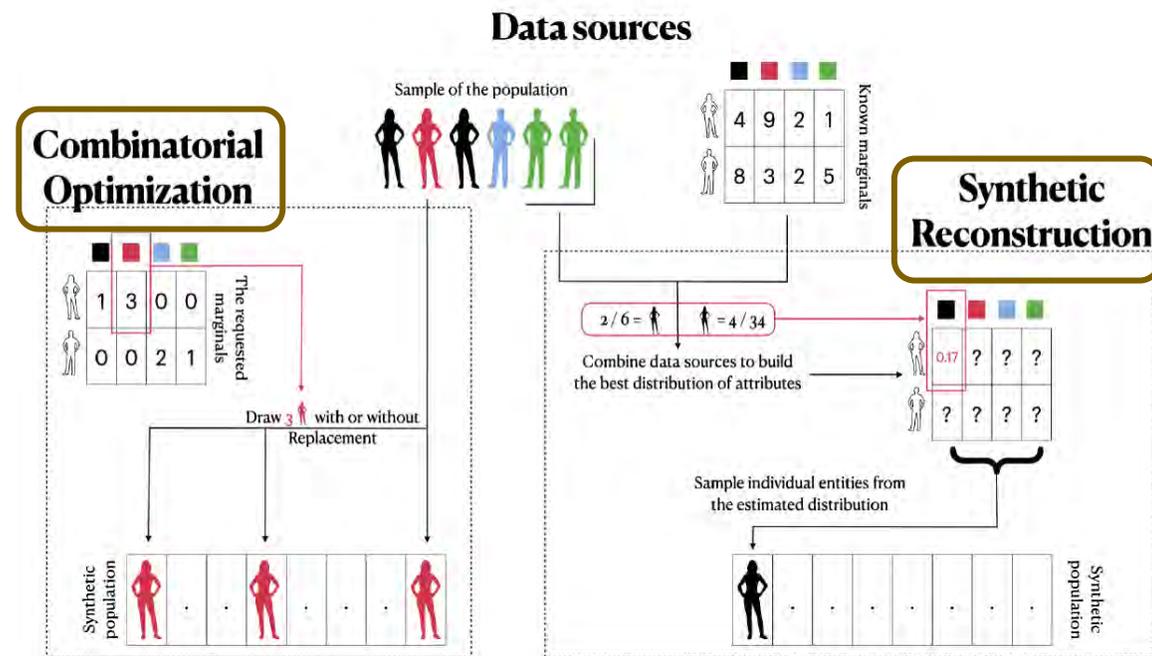- The system needs to be a synthesis of the information available of the target society

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

Consiglio Nazionale
delle Ricerche

- Challenge to identify attributes at micro-level:
  - data not collected
  - separate datasets
  - privacy issues

**Synthetic populations**: a series of techniques to handle available data and replicate attributes of the target population

'While a synthetic population is implicitly an artificial population, an artificial population is not necessarily a synthetic population'

1111112234
1111112234
1111112234
442rrr34333

*synthesis*
of information

*generation*
of synthetic data
Bigi et al. (2022)

*comparison*
with available data



Chapuis et al. (2022)

Machine Learning

Every method has its peculiarities and boundaries not set in stone

Finanziato dall'Unione europea
NextGenerationEU

Ministero dell'Università e della Ricerca

Italiadomani
PIANO NAZIONALE DI RIPRESA E RESILIENZA

Consiglio Nazionale delle Ricerche

# Synthetic Reconstruction

- Focus on estimate of unknown joint distributions from data available from marginal distribution
- macro to micro

**Extensions**

- **Multiple Iterative Proportional Fitting (MIPF)** taking the estimated joints as marginal to next step dimension
- **Hierarchical Iterative Proportional Fitting (HIPF)** nested data fitting marginals from one level (e.g. household) to the narrow (e.g. citizens) (Yamaego et al., 2021)
- **Iterative Proportional Updating (IPU)**: from micro data finds weight for cross-category multiplied by marginals and correct backwards -> combinatorial optmization

- Archetype in the **Iterative Proportional Fitting** (IPF, raking)



weight for each cell
update **cell row**
marginal **estimates**
update **cell columns**
update **marginals**
stop ← target → no

$$weight = \frac{observed\ marginal}{fitted\ marginal}$$

✅ mathematically transparent, robust, focus on weights

🟥 require ad hoc setting of algorithms, zero-cell problem

FOSSR:SPG

R::synthpop (MIPF)

Finanziato dall'Unione europea
NextGenerationEU

Ministero dell'Università e della Ricerca

Italiadomani
PIANO NAZIONALE DI RIPRESA E RESILIENZA

Consiglio Nazionale delle Ricerche

# Combinatorial Optimization

- Scaling of the synthetic population
- Draw random data and optimize against marginals observed

- **Simulated Annealing (SA)** from the micro-data identifies some seed numbers and compares synthetic marginas to empirical marginals to correct backwards
- **Markov Chain Monte Carlo (MCMC)** random number generations for all intersection and correct backwards

Seed data observations

Random noise

Initial solution

Evaluate

target

Rearrange

Optimal solution met

☑ overcome some setbacks of IPF family

☑ better with complex intersections (handled *at once*)

◼ focus on the outcome rather than weights & inner conditions

◼ computationally more demanding

py::simmaneal

R::MCMCpack   py::PyMC

# Machine Learning Approach

- Most recent in time
- Statistical learning
- micro to macro

- **Generative Adversarial Networks (GAN)**
  - originated from images AI, applied to data
  - two competing (neural) networks:
    - **Generator** who produces random data
    - **Discriminator** that discriminate realistic data from not realistic
    - Goal of Generator is to get better to *deceive* the Discriminator who gets better in discriminate, meaning that synthetic data are very realistic



✅ promising because they integrate the performance of combinatorial methods with transparency, multidimensionally compared to IPF family

🟥 they require a training dataset from which the learning process occurs, where underrepresented groups are likely to be ignored in estimates

- Pre-adjust the training dataset with ad hoc weights to marginals (Falck, 2025)
- Post-adjust the synthetic outcome with weights to marginals

py::PyTorch

# Synthetic Populations Generator (SPG)

**FOSSR**
Fostering Open Science in Social Science Research

Service to enable researchers and policy makers to extract synthetic populations at desired level of information from input dataset

- agent-based modeling
- spatial analysis
- conditional model
- ...

## Open Science & Source Software

https://github.com/RoccoPaolillo/IPF_multidim.git > synthpopgen.py



executability

VRE
(container)

modifiability

demo
(standalone program)

GitHub
(open code)

(Jimenez et al., 2017; Hong et al., 2022)

# Iterative Proportional Fitting (IPF)

weight for each cell

↓

update cell row

↓

marginal estimates

↓

update cell columns

↓

update marginals

↓

stop ← target → no

$$weight = \frac{observed\ marginal}{fitted\ marginal}$$

## Multiple Iterative Proportional Fitting (MIPF)

marginal varible 1     marginal varible 2

↓

joint variable 1,2     marginal variable 3

↓

joint variables 1,2,3     marginal variable n...

↓

...

Selected for higher transparency, robustness, light computing, core mechanism common to other method (somehow), but provided more service-oriented direction is guaranteed

First version released (ISTAT gender X age for validation)     10.5281/zenodo.10638800

# Aims of Synthetic Populations Generator (SPG)

- Include **multidimensionality**
- Increase **generalizability** of variable handling
- Enable **automation** input-execution-output
- Customize **filtering** selection

  - leverage estimate of joint and conditional probability over in-cell weight iteration

Tested with **opensalute Lazio** health data:
- gender
- age
- hyptertension (hpt)
- heart failure (hf)

known joints:
- hypertension over age
- heart failure over age

goal of service: identify joint distribution for all combinations

| gender | age | hpt | hf | value |
|--------|------|-----|-----|---------|
| male | | | | 3073047 |
| female | | | | 3259977 |
| | 30 | | | 1745215 |
| | 3060 | | | 2832088 |
| | 60100 | | | 1755721 |
| | | yes | | 1193445 |
| | | no | | 5139579 |
| | | | yes | 93926 |
| | | | no | 6239098 |
| | 30 | yes | | 3547 |
| | 3060 | yes | | 252543 |
| | 60100 | yes | | 937355 |
| | 30 | no | | 1741668 |
| | 3060 | no | | 2579545 |
| | 60100 | no | | 818366 |
| | 30 | | yes | 424 |
| | 3060 | | yes | 8459 |
| | 60100 | | yes | 85043 |
| | 30 | | no | 1744791 |
| | 3060 | | no | 2823629 |
| | 60100 | | no | 1670678 |

# The algorithm

input csv

identify total population

identify marginals

identify combinations (one value for variable)

identify joint distribution

**known joint categories?**

no
**conditional product = 1**

yes
**conditional_product = known joint category / basevar**

**estimate percentage** : **(basevar/ population) * conditional_product**

**estimate percentage** : **(indvar / population)** * **(basevar/ population) * conditional_product**

**size synthetic population differs from input?**

no
**estimated counts = estimate percentage * size population**

yes
**estimated counts = estimate percentage * size scaled population**

display

filter

output csv

**Estimate percentageM30HPTHF**
gender male, age 30, hptyes, hf yes

$$\frac{male}{population} * \left( \frac{age30}{population} * \frac{hptyes, age30}{age30} * \frac{hfno, age30}{age30} \right)$$

**Estimate countM30HPTHF**
**Estimate percentageM30HPTHF * size synthetic population**

# Open code

https://github.com/RoccoPaolillo/IPF_multidim.git  >

synthpopgen.py

## cmd line

python synthpopgen.py -i input_file_tuples.csv \
  -f (filter)
   'all'
   'gender:female,age:3060' \
-d (display)
  'split'
  'aggregate'\
- v (validate)*
--synth-total 20303*
-o results.csv



→ **pro**: high modifiability for users(rewrite, retest...)

→ **vs**: knowledge coding, dependencies

measure of validation
**Average percentage error** between input marginals
(and joints) data and stimated marginals, **RMSE 0.6**

| constraint | observed | predicted | avg_percentage_err |
|---|---|---|---|
| age=30 | 1745215 | 1745215 | 0.0 |
| age=30,hf=no | 1744791 | 1744791 | 0.0 |
| age=30,hf=yes | 424 | 424 | 0.0 |
| age=30,hpt=no | 1741668 | 1741668 | 0.0 |
| age=30,hpt=yes | 3547 | 3547 | 0.0 |
| age=3060 | 2832088 | 2832087 | 3,53E-02 |
| age=3060,hf=no | 2823629 | 2823628 | 3,54E-02 |
| age=3060,hf=yes | 8459 | 8459 | 0.0 |
| age=3060,hpt=no | 2579545 | 2579545 | 0.0 |
| age=3060,hpt=yes | 252543 | 252542 | 0.00039597 |
| age=60100 | 1755721 | 1755722 | 5,70E-02 |
| age=60100,hf=no | 1670678 | 1670678 | 0.0 |
| age=60100,hf=yes | 85043 | 85044 | 0.00117588 |
| age=60100,hpt=no | 818366 | 818366 | 0.0 |
| age=60100,hpt=yes | 937355 | 937356 | 0.00010668 |
| gender=female | 3259977 | 3259977 | 0.0 |
| gender=male | 3073047 | 3073047 | 0.0 |
| hf=no | 6239098 | 6239097 | 1,60E-02 |
| hf=yes | 93926 | 93927 | 0.00106467 |
| hpt=no | 5139579 | 5139579 | 0.0 |
| hpt=yes | 1193445 | 1193445 | 0.0 |

* only if whole population combinations are stimated (-f all)

# Standalone program

.exe local file
py:: tkinter

→ **pro**: not coding needed, no dependencies
→ **vs**: local CPU, no modifiability

# FOSSR VRE Container

enhance collaboration and tools of researchers through digital platform
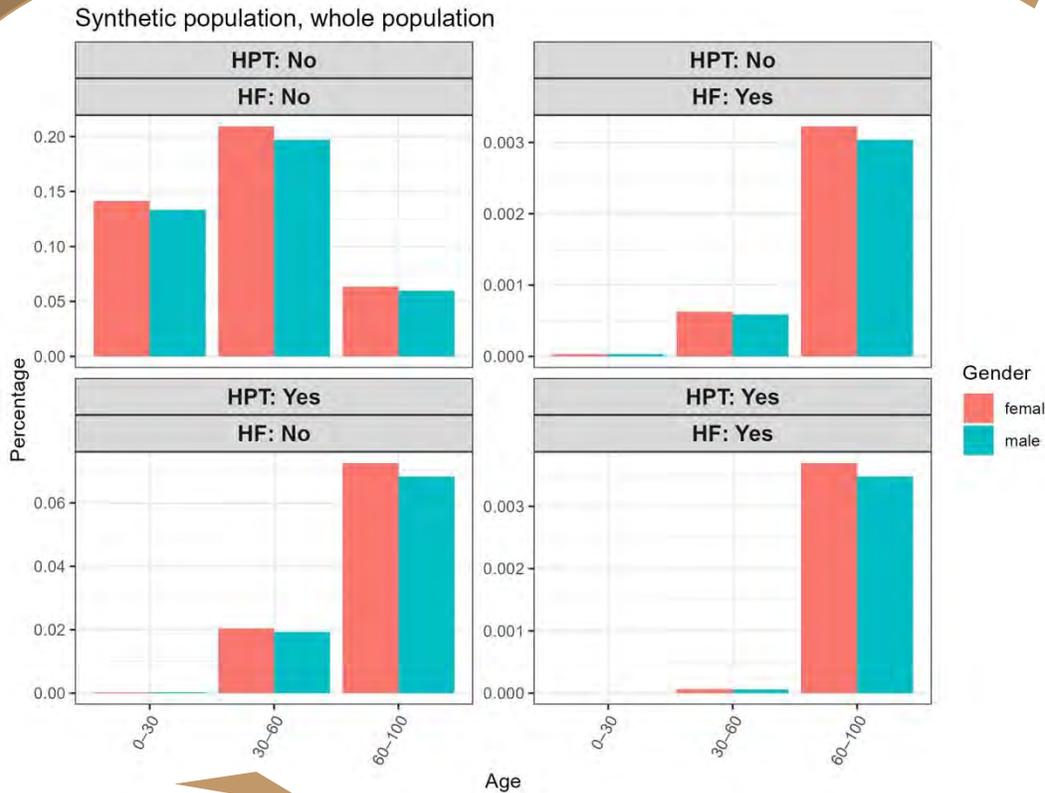
https://fossr.d4science.org

➔ **pro**: uses D4Science servers, web-app
➔ **vs**: no modifiable, internet-dependent

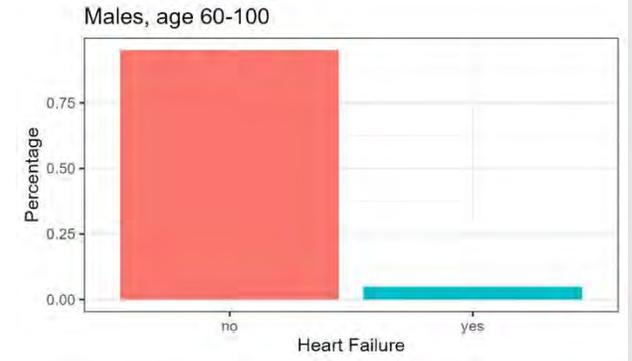VRE > CCP (Cloud Computing Platform)> Synthetic Populations Generator (SPG)

# Output & applications

gender:male,age:60100, hpt:yes -d aggregate: 454844

gender:male,age:60100, hpt:yes -d split

- individual-level attributes as per request
- preserving original data and restrictions
- identify intersections and heterogeneity in the population
- actor-based models initialization (agent-based models, social network analysis,...)

# Current and future steps

- Assumes marginals from the same total population
- looking at machine learning as micro-to-macro approach
- Automated iteration over spatial units (census tract)
- Automated multi-source input
- Enhance UX experience and assitance
  - **LLM**

https://github.com/RoccoPaolillo/IPF_multidim.git

Deployment into the **FOSSR market place**

**References**

Paolillo, R., Roxburgh, N., Sbrana, A., Polhill, G., Sabatella, E. C., & Paolucci, M. (2025). Synthetic Populations in Research Infrastructures. In Longitudinal Data Infrastructures in Europe: Tools for Open Science in Social Science Research (pp. 153-164). Cham: Springer Nature Switzerland.

Chapuis, K., Taillandier, P., & Drogoul, A. (2022). Generation of synthetic populations in social simulations: a review of methods and practices. Journal of Artificial Societies and Social Simulation, 25(2).

Bigi, F., Rashidi, T. H., & Viti, F. (2024). Synthetic population: A reliable framework for analysis for agent-based modeling inmobility.Transportation Research Record, 2678(11),1–15.

Yameogo, B. F., Vandanjon, P. O., Gastineau, P., & Hankach, P. (2021). Generating a two-layered synthetic population for French municipalities: Results and evaluation of four synthetic reconstruction methods. Journal of Artificial Societies and Social Simulation, 24,27.

Falck, V. (2025). Generating spatial synthetic populations using Wasserstein generative adver-sarial network: A case study withEU-SILC data for Helsinki and Thessaloniki. Preprint.arXiv:2501.16080.

Jiménez, R. C., Kuzak, M., Alhamdoosh, M., Barker, M., Batut, B., Borg, M., ... & Crouch, S. (2017). Four simple recommendations to encourage best practices in research software. F1000Research, 6, ELIXIR-876.

Chue Hong, N. P., Katz, D. S., Barker, M., Lamprecht, A. L., Martinez, C., Psomopoulos, F. E., ... & Yehudi, Y. (2022). FAIR principles for research software (FAIR4RS principles). Zenodo.

# THANK YOU!

✉ **fossr.dissemination@ircres.cnr.it**

🐦 **@fossrproject**

📘 **fossr.eu**

in **fossr-eu**

▶ **@fossr**

zenodo **zenodo.org/communities/fossr**

**Thank you! Questions?**

**FOSSR DAYS 2026, 4-5-6 February**



l.cnr.it/fossr-days-2026-registration-form

rocco.paolillo@cnr.it
@roccopaolillo.bsky.social

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA
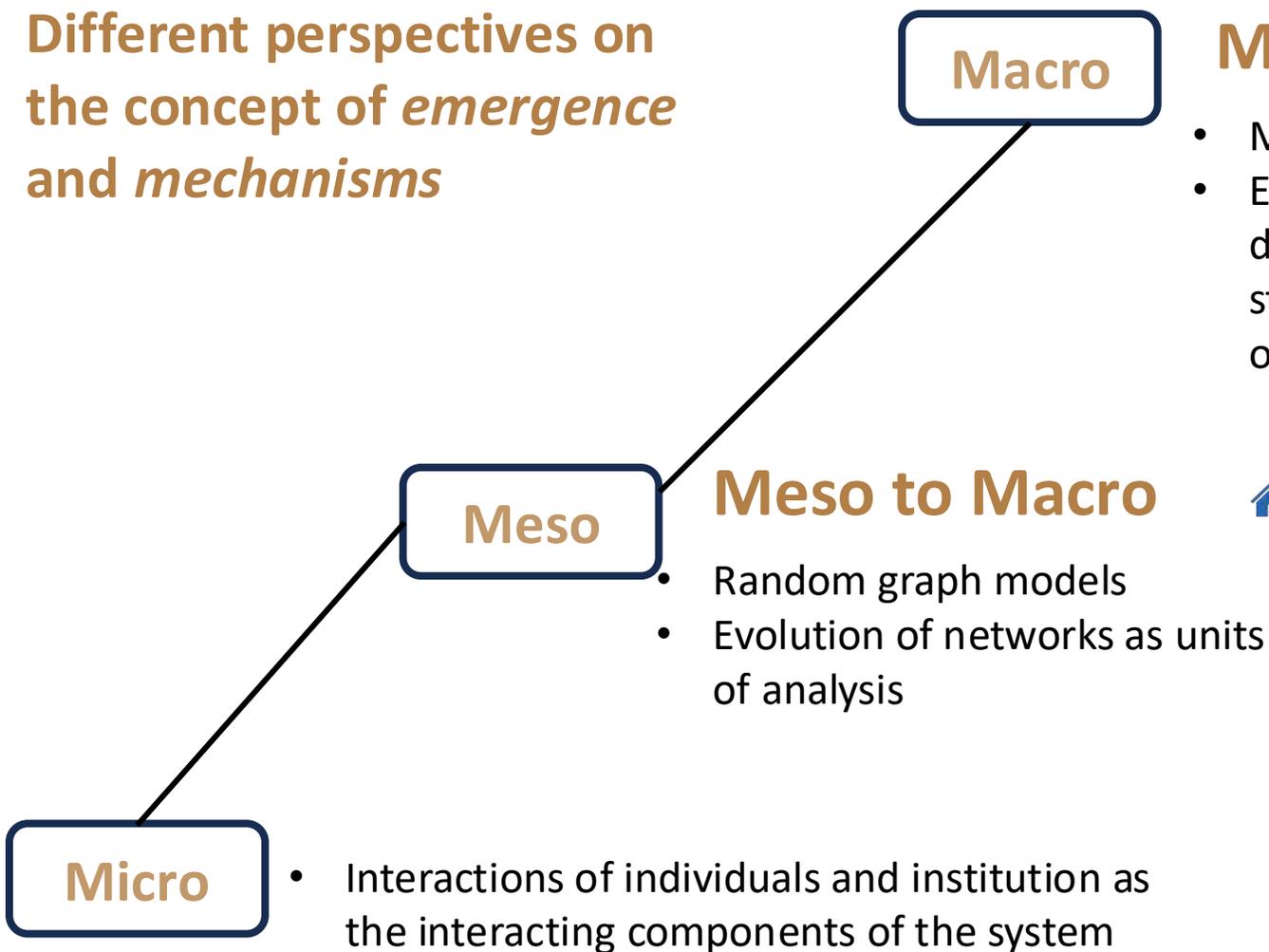
Consiglio Nazionale
delle Ricerche

# Social complexity and social mechanisms



- Most of social phenomena are inherently collective phenomena
- They assume the definition of a system made of interacting components (e.g. market place, urban landscape, welfare state) where the phenomenon unfolds
- Common scope of different disciplines is to unfold the mechanisms that unfold the dynamics of the phenomenon



If the concepts of system and the goal to identify mechanisms to disclose the phenomenon is common to many disciplines, different perspectives apply
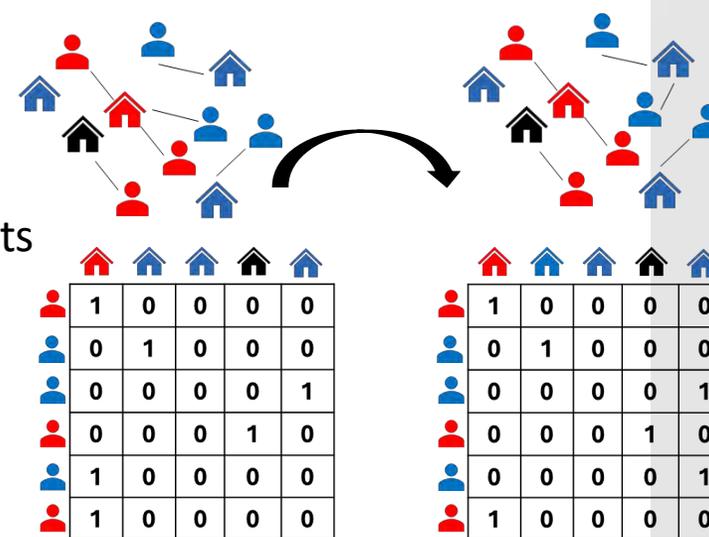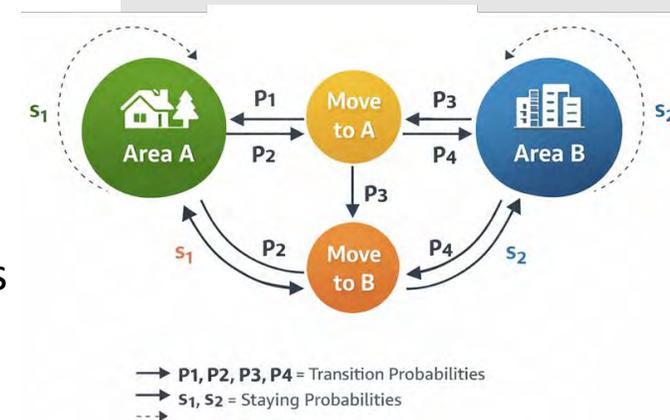
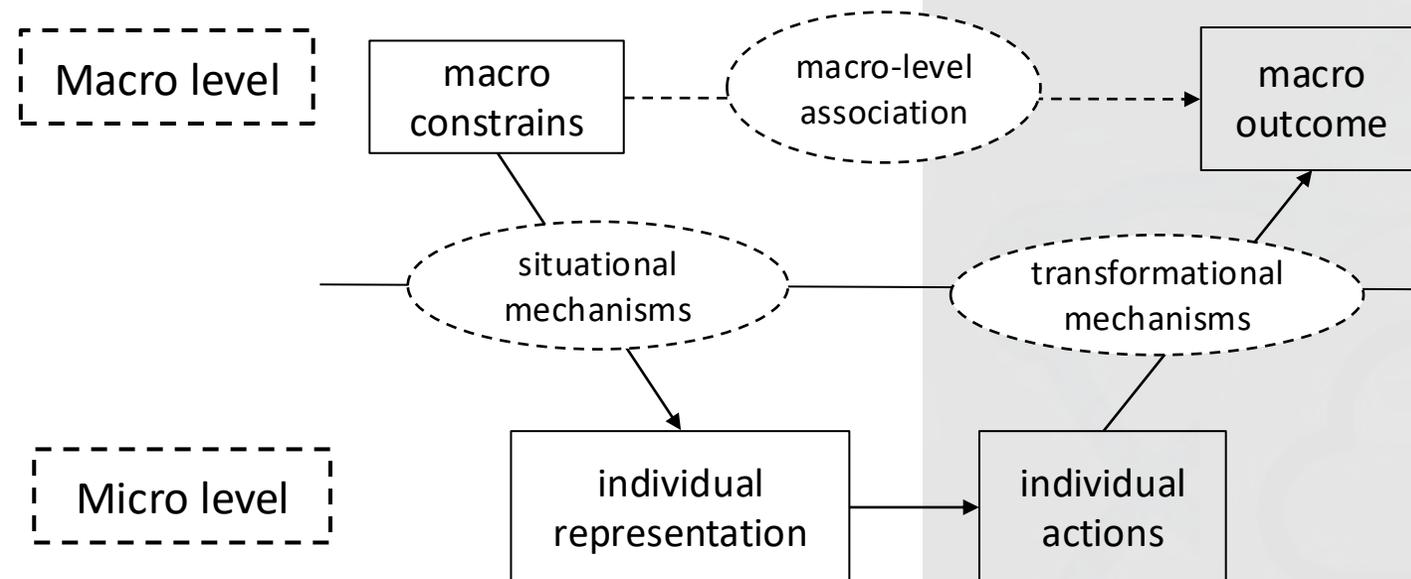**Different perspectives on the concept of *emergence* and *mechanisms***

**Macro**

## Macro to Macro

- Markov Chain
- Evolution of processes depending on the previous state of the system as unit of analysis

**Meso**

## Meso to Macro

- Random graph models
- Evolution of networks as units of analysis

**Micro**

- Interactions of individuals and institution as the interacting components of the system
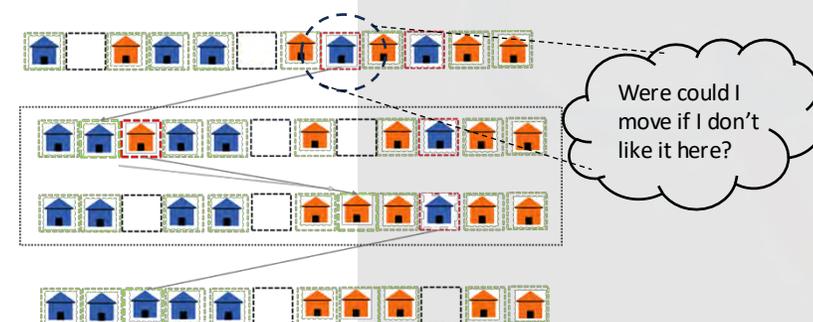
# Micro to Macro

- But sometimes, we might be interested in those mechanisms that move from the micro level, e.g. citizens/institutions with their attitudes, motivations and course of action, but get **outside of the individual agency and inglobe the interaction** of individuals as explicative mechanism of emergence

- The phenomenon is an aggregated, mutual adaptation of individuals rather than the sum of individual action
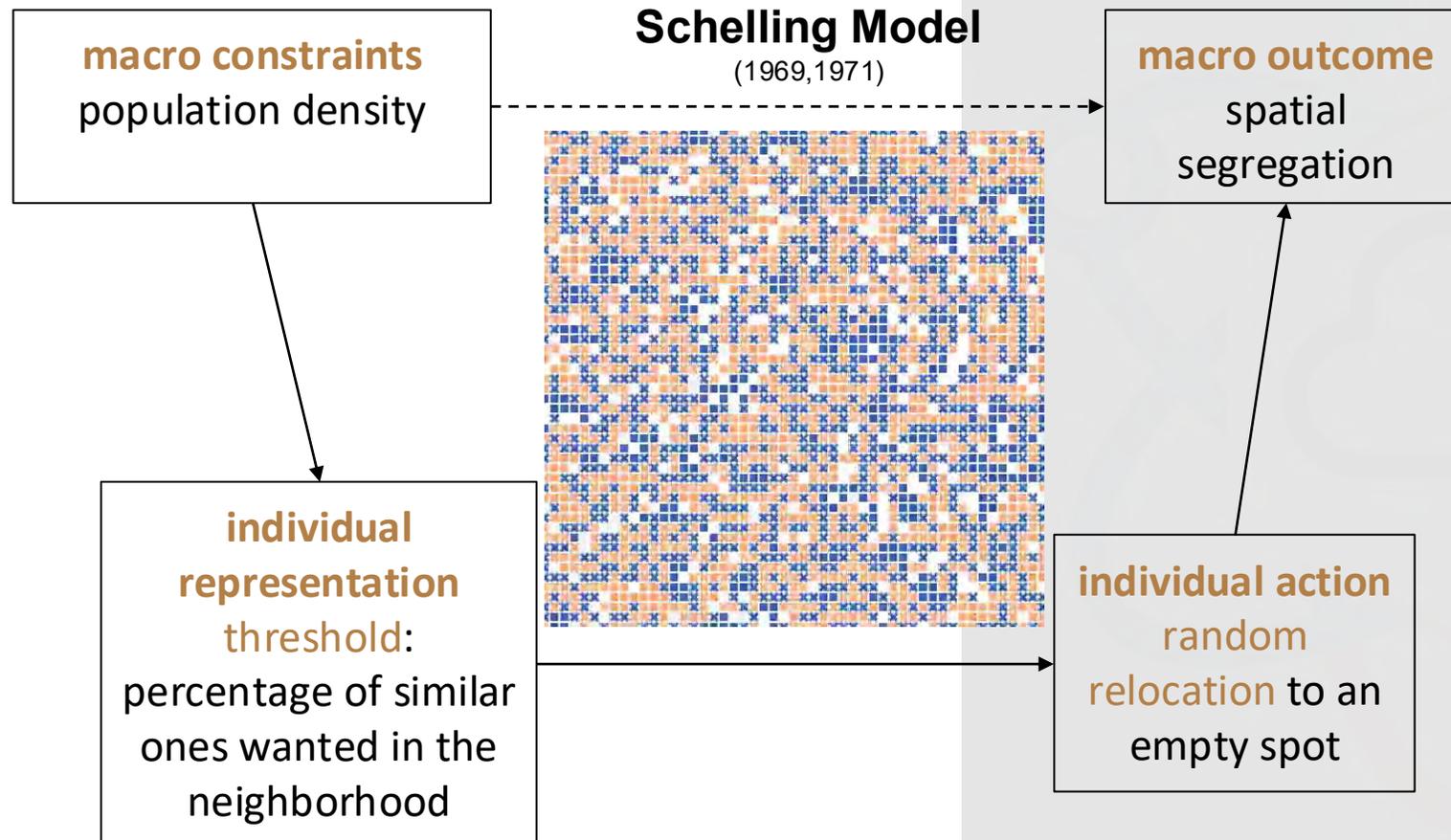
v
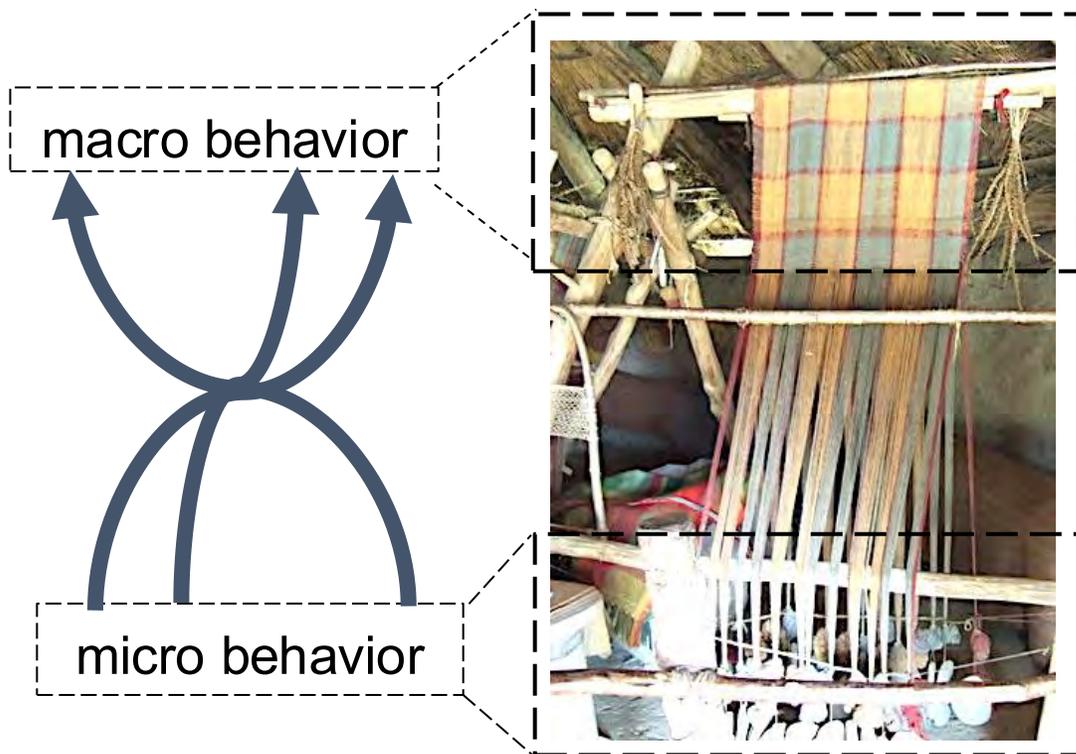**Analytical Sociology
Agent-based Modeling**



Coleman Boat (1994), additions by Hedström and Ylikoski (2010), adapted
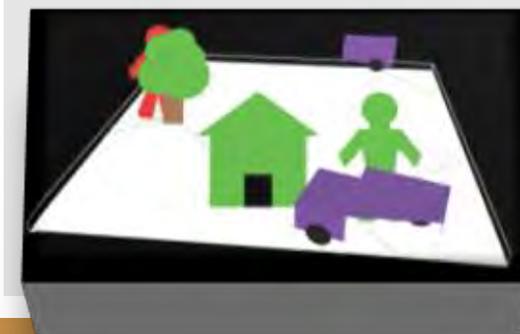
# Example of spatial segregation

- **actors:** households evaluating their neighborhood
- **individual behavior:** preference for a percentage (threshold) of similar ones in their neighborhood (homophily)
- **mechanism of emergence:** cascade effects where the behavior of one household influences the composition of neighborhood and preference of others
- **unexpected outcome:** high levels of spatial segregation, also for mild threshold preferences

**macro constraints** population density

**Schelling Model**
(1969,1971)

**macro outcome** spatial segregation

**individual representation** threshold: percentage of similar ones wanted in the neighborhood

**individual action** random relocation to an empty spot

## Agent-based Modeling



macro behavior

micro behavior

- **Simulation method** tailored to model the interacting components that constitute the system, e.g. agents representing citizens in an **artificial society**
- We can manipulate both attributes and plan of actions of agents and observe the consequences of interaction of agents executing their plans.
- By manipulating plans and conditions where the agents interact and adapt, we can experiment on and formalize the dynamics of emergence of the collective phenomenon

# Social computing with agent-based modeling

KISS
Keep it short simple, stupid
KIDS
Keep it descriptively simple

**Design of the conditions, actors and initial mechanism we want to test to study the phenomenon**

- A society where people differentiate by some traits
- They stay in a neighborhood if certain threshold of similarity are satisfied
- Can segregation emerge even if the threshold is not that high?

**Formalization into rules and functions**

$$if \ \vartheta < \theta : \text{leave} \ ;$$
$$if \ \vartheta \geq \theta : stay$$

**Translation into code to translate the theoretical model we want to test and investigate setting-up what-if scenarios**

```
set happy? similar-nearby >=
(%-similar-wanted*total-nearby / 100)

to move-unhappy-turtles
  ask turtles with [ not happy? ]
    [ find-new-spot ]
end
```

**Define what-if scenarios to create experimental conditions**

["density population" 70 95]
["%-similar-wanted" 0 30 60]

**Collection of data as measurement of changes in the system**

mean exposure to similars in the neighborhood of agents when no one relocates anymore
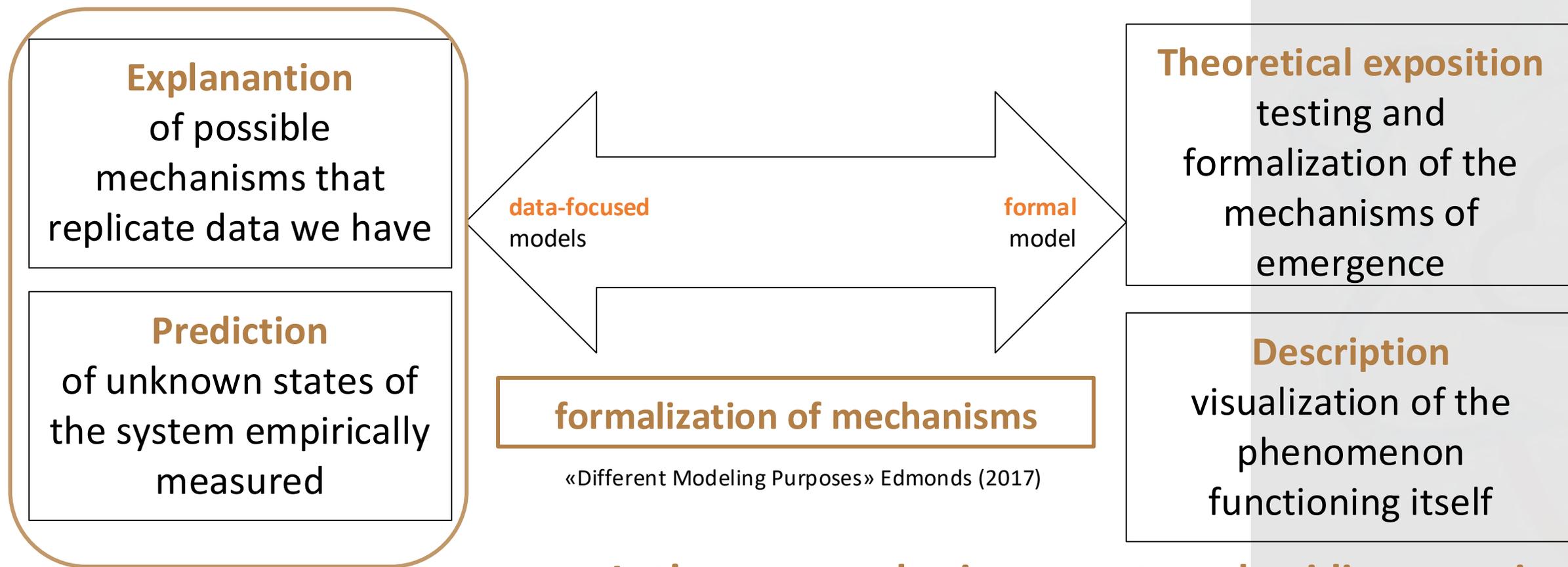
'If you can grow it, you have explained it' (Epstein, 2006)

'If you don't know how you grew it, you didn't explain it.' (Macy & Flache, 2009, p.263)

# What is agent-based modeling useful for?

**Explanantion**
of possible mechanisms that replicate data we have

**Prediction**
of unknown states of the system empirically measured

**data-focused** models

**formal** model

**formalization of mechanisms**

«Different Modeling Purposes» Edmonds (2017)

**Theoretical exposition**
testing and formalization of the mechanisms of emergence

**Description**
visualization of the phenomenon functioning itself

**Let's see some basic concepts and guiding questions**

# Agents – Who are the actors involved in the phenomenon?

**Agents**: a virtual object capable of elaborating information and able to execute an action (individuals, institutions, households...)

•**Intentionality**: acting based on goals or plans
•**Proactivity**: initiating actions rather than waiting
•**Reactivity**: responding to external stimuli or changes
•**Prosociality**: acting in coordination with others (social agents)



**State variable:** what characteristics cannot change through time? Ethnicity, gender
**Dynamic variable**: what characteristics can change through time? Opinions, preferences
**Global variable**: shared by all agents
**Local variable**: shared by specific agents or class of agents
**Heterogeneous** vs **Homogeneous** (attributes distribution)

## Attributes, Beliefs, Desire, Intentions (A+BDI)

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

Consiglio Nazionale
delle Ricerche

# Parameters – What are the conditions affecting the phenomenon?

**Parameters**: tunable numeric (not only) variables to modify the system:
- calibrate initial conditions
- agents' variables distributions

mean          50
standard_deviation 1.0

speed-limit      1.0
number_cars      200

**Global parameter:** variable affecting all agents, and every agent can interact with
Belief shared by all agents

**Local variable**: accessible only to some specific agents
Norm specific to a class of agents

# Evolution (I) – How the phenomen emerges through agents' interaction?

It is not much **time** as a continuous variable, rather the evolution of the system along two interconnected concepts:

- **micro level: schedule** of activation of agents' behavior
- **macro level transition phase** of the system changing due to mutual adaptation of the agents

**Cascade effect** of the behavior of one agent on the neighborhood composition to other agents, affecting segregation at macro level

# Evolution (II) – How does the order of agents' action influence each other?

**Parallelization:** how the behavior is executed

**Synchronous behavior:** agents act together in parallel

**Asynchronous behavior:** agents act sequentially
(physical threads)

**Synchrony:** when the behavior is executed

**Synchronization**: agents decide based on the same knowledge of the world, including effects of actions of others
(they act «in parallel»)

**Example**: scholar agents choose between two universities based on the chance to be introduced to elective authors based on potential shared connections. The sequential order of agents can affect the decision of those who select after

UTS > ATS: update shared knowledge

Longo, Paolillo, Ceriani (2026, forthcoming)

# Outcome (I) – How can I read the evolution of the phenomenon?

## Tipping points

A sudden transition of the system is narrowed to one direction
e.g. in Schelling model the local level of segregation triggers relocation so that segregation becomes steady

## Bifurcation

A moment where the phenomenon can diverge in two opposite directions with equal probability
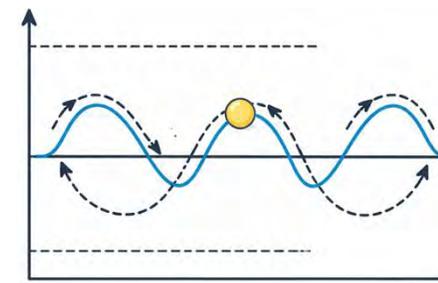e.g. in case of political polarization

### stable equilibrium

not mutable system
(Schelling, consensu)

### cyclic equilibrium

system follows a trend of sequential cycles
(grass & sheeps, gentrification)

### dynamic equilibrium

oscillations/inflows/outflows causing the system to apparently remain in balance
(e.g. supply & demand)

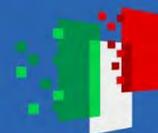# Outcome (II) – How can I identify how the conditions affect the phenomenon?

**Global sensitivity analysis**: pattern-oriented-modeling (POM) interaction between parameters (space of the model) to understand the overall mechanisms of the model

**Local sensitivity analysis**: one-factor-at-time (OFAT), focus on the effect of one specific parameter (nominal value) over the others



Agent-Based and Individual-Based Modeling
A PRACTICAL INTRODUCTION
SECOND EDITION

Steven F. Railsback and Volker Grimm

**Compare what-if scenarios and measures**

**Parameters**
What are the conditions
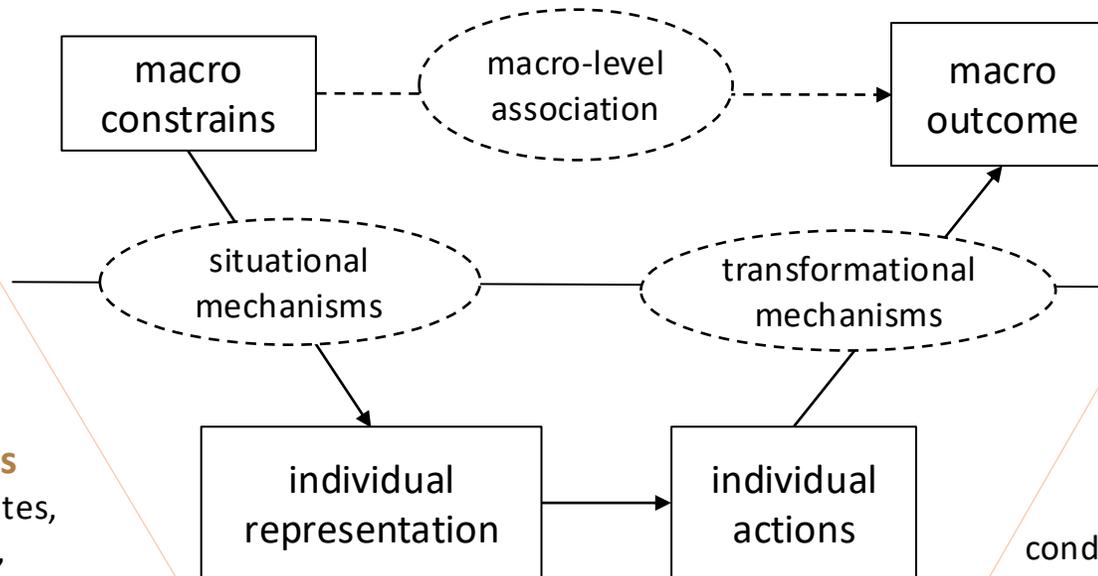affecting the phenomenon?

**Design**

**Agents**
Who are the actors involved
in the phenomenon?

**Formalization**

**Agents**
Attributes,
Beliefs,
Desires,
Intentions

**Code**

**Outcome**
How can I read the evolution
of the phenomenon?

**Evolution**
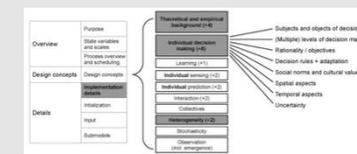How the phenomenon emerges
through agents' interaction?

**Outcome**
How can I identify how the
conditions affect the phenomenon?

**What-if scenarios**

macro
constrains

macro-level
association

macro
outcome

situational
mechanisms

transformational
mechanisms

individual
representation

individual
actions

**Evolution**
How does the order of agents' action influence each other

Adapt your model to the grammars of a programming language

**Let's implement to a case study and one programming tool (NetLogo)**

**ODD+D protocol to clarify ideas**

Müller et al., 2013

# References, suggested readings and tools

Macy, M. W., & Willer, R. (2002). From factors to actors: Computational sociology and agent-based modeling. Annual review of sociology, 28(1), 143-166.

Edmonds, B. (2017). Different modelling purposes. Simulating social complexity: A handbook, 39-58.

Edmonds, B., & Moss, S. (2004, July). From KISS to KIDS–an 'anti-simplistic' modelling approach. In International workshop on multi-agent systems and agent-based simulation (pp. 130-144). Berlin, Heidelberg: Springer Berlin Heidelberg.

Epstein, J. M., & Axtell, R. (1996). Growing artificial societies: social science from the bottom up. Brookings Institution Press.

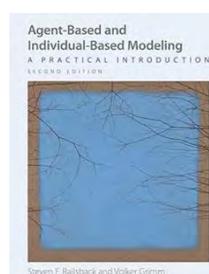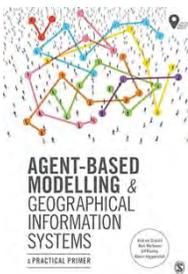Hedström, P., & Ylikoski, P. (2010). Causal mechanisms in the social sciences. Annual review of sociology, 36, 49-67.
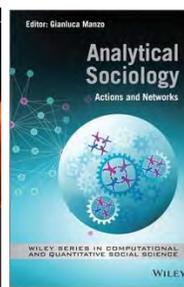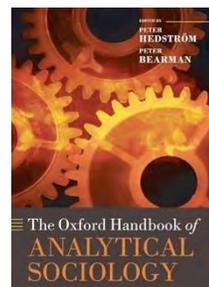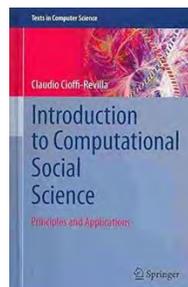
Macy, Michael W., and Andreas Flache. 2009. "Social Dynamics from the Bottom Up: Agent-based Models of Social Interaction." In Hedström, P. and Bearman, P. (Eds.) The Oxford Handbook of Analytical Sociology. Oxford, UK: Oxford University Press.

Schelling, T. C. (1969). Models of segregation. The American economic review, 59(2), 488-493.

Schelling, T. C. (1971). Dynamic models of segregation. Journal of mathematical sociology, 1(2), 143-186.

Longo, C.F., Paolillo, R., & Ceriani, M. (Forthcoming). T2B2T: The Ontology for adaptive Agent-driven seamless integration with the Semantic Web. In Proceedings of the Joint Ontology Workshops (JOWO) - Episode XI: The Sicilian Summer under the Etna, co-located with the 15th International Conference on Formal Ontology in Information Systems (FOIS 2025), To appear on CEUR

Müller, B., Bohn, F., Dreßler, G., Groeneveld, J., Klassert, C., Martin, R., ... & Schwarz, N. (2013). Describing human decisions in agent-based models–ODD+ D, an extension of the ODD protocol. Environmental Modelling & Software, 48, 37-48.

## Softwares free

**NetLogo**
https://www.netlogo.org/

**MESA python**
https://mesa.readthedocs.io/latest/

**GAMA Platform**
https://gama-platform.org/

https://essa.eu.org/    SIMSOC@jiscmail.ac.uk

Institute for Analytical Sociology
Linköping University

rocco.paolillo@cnr.it

## Thank you! Questions?

@roccopaolillo.bsky.social

**Parameters**
What are the conditions
affecting the phenomenon?

**Outcome**
How can I read the evolution
of the phenomenon?

**Design**

**Agents**
Who are the actors involved
in the phenomenon?

**Evolution**
How the phenomenon emerges
through agents' interaction?

**Formalization**

**Agents**
Attributes,
Beliefs,
Desires,
Intentions

**Outcome**
How can I identify how the
conditions affect the phenomenon?

**Code**

**What-if scenarios**

**Evolution**
How does the order of agents' action influence each other

| macro constrains | --- | macro-level association | ---> | macro outcome |

situational mechanisms → transformational mechanisms

individual representation → individual actions

Adapt your model to the grammars of a programming language

**Let's implement to a case study and one programming tool (NetLogo)**
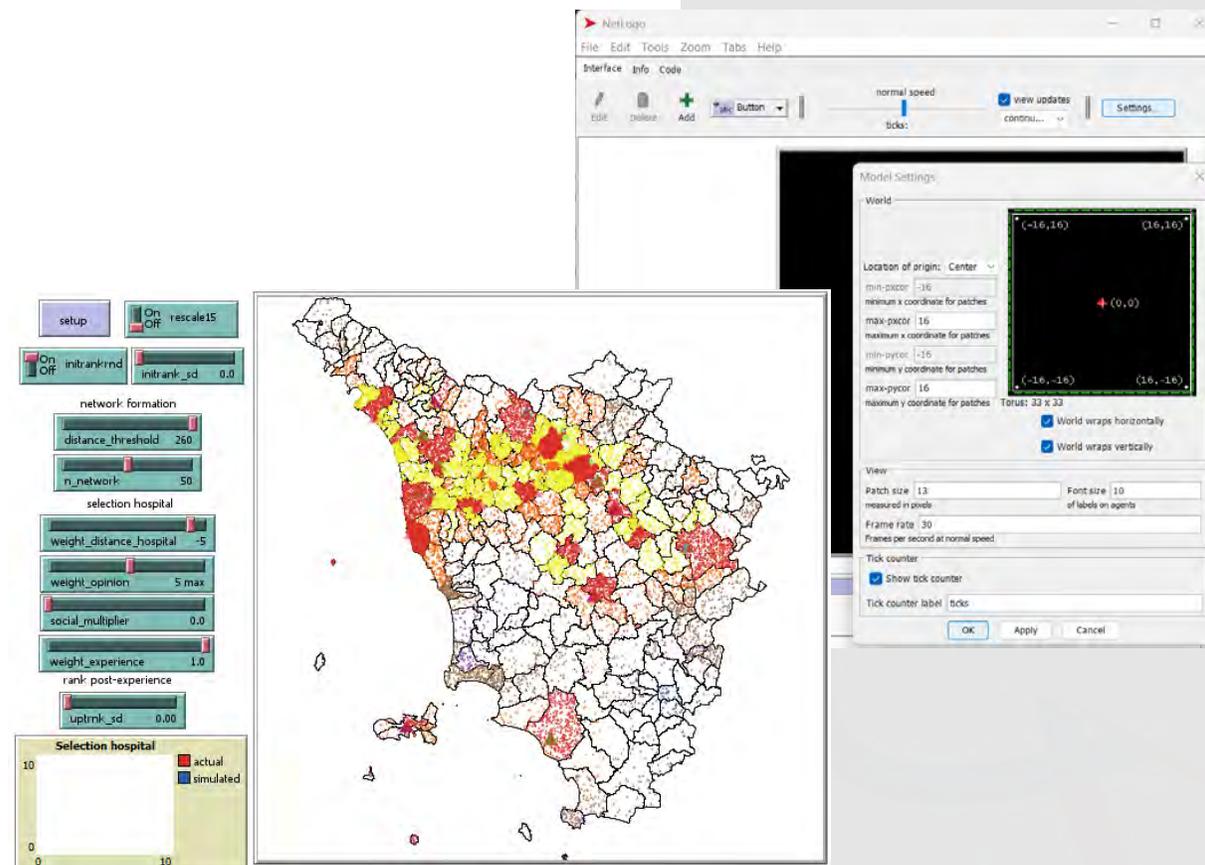
**ODD+D protocol to clarify ideas**
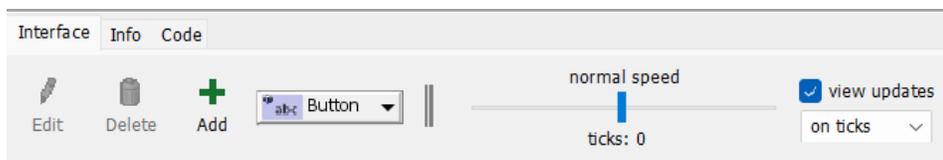
Müller et al., 2013

# NetLogo

- Platform IDE (interfaccia) e programming language (java + starlogo) specific to agent-based modeling and experiments
- Open Source & User-friendly
- Allows many extensions (shapefile, csv import, random-wheel selection…)
- Widely used in the social science community and continuously maintained (7+ version)
- Programming language tailored to be as intuitive as possible and ready functions
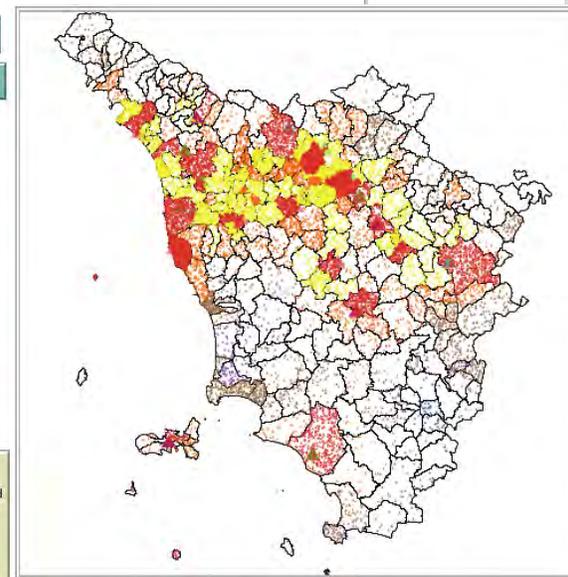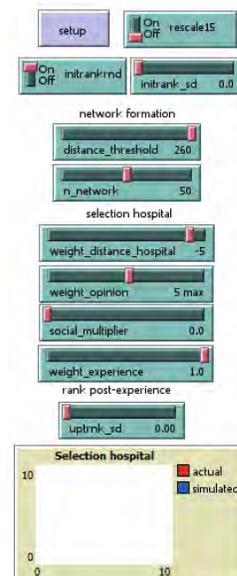- Supported by documentation (and Chat-GPT)
  - https://docs.netlogo.org/dictionary

https://www.netlogo.org/
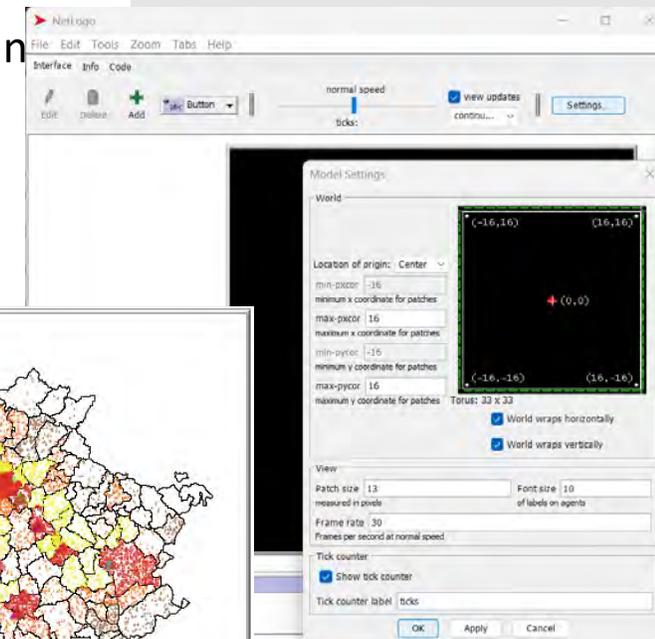
# NetLogo

**world** where things happen

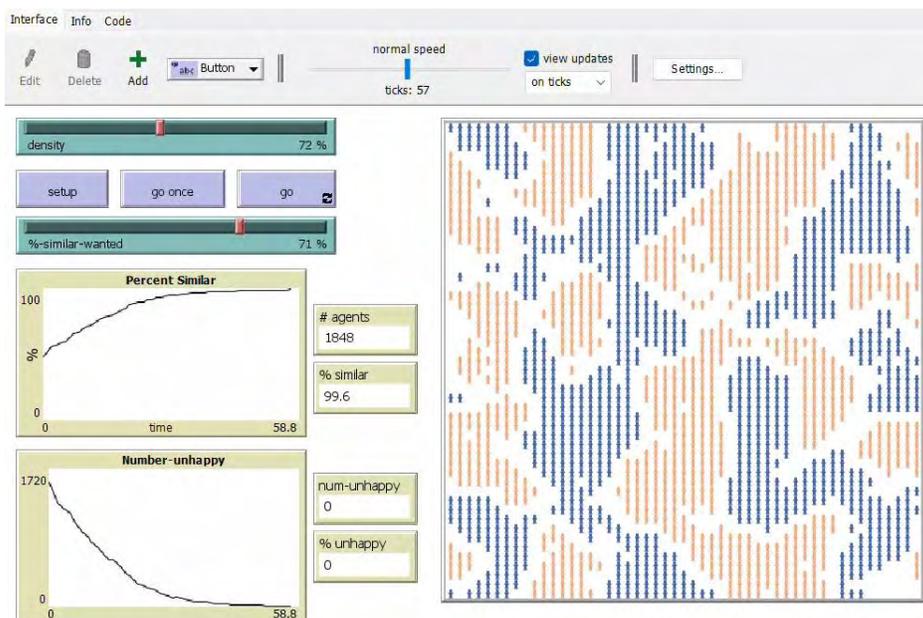**Sliders, buttons, chooser** to facilitate interaction with parameters in the interface to explore conditions

**Interface tab** to interact
**Info tab tab** to document
**Code tab** to build the model
(also in conjunction with interface)

**A command line** to interact *on the fly*

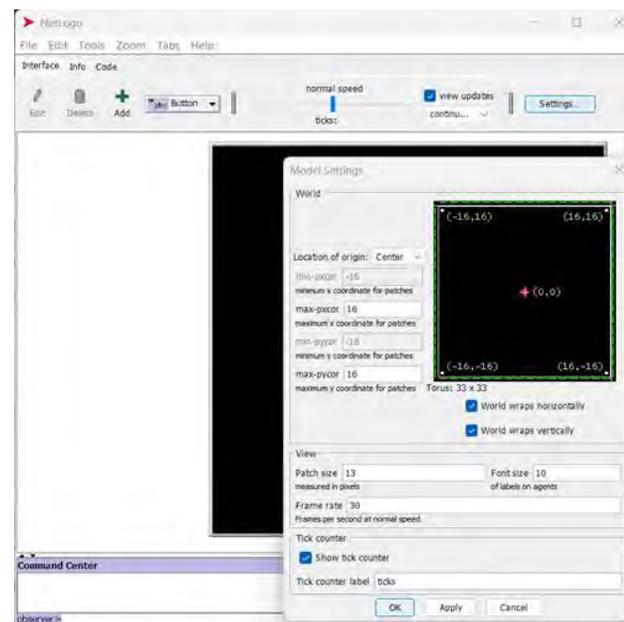**Plots, monitors** to detect how the phenomenon is emerging
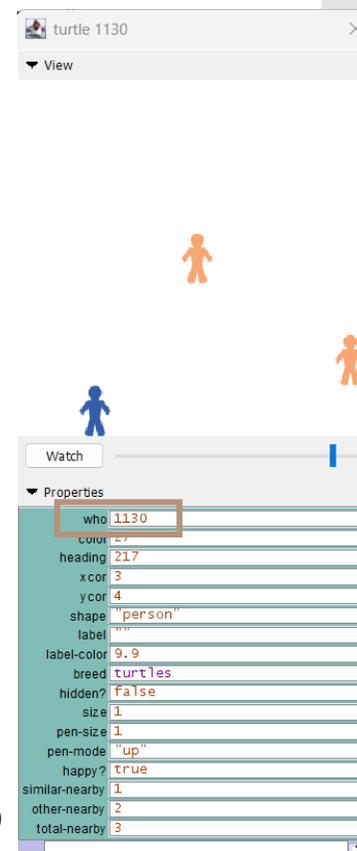
# NetLogo

grid space world

0-indexed

`list` [a b c]
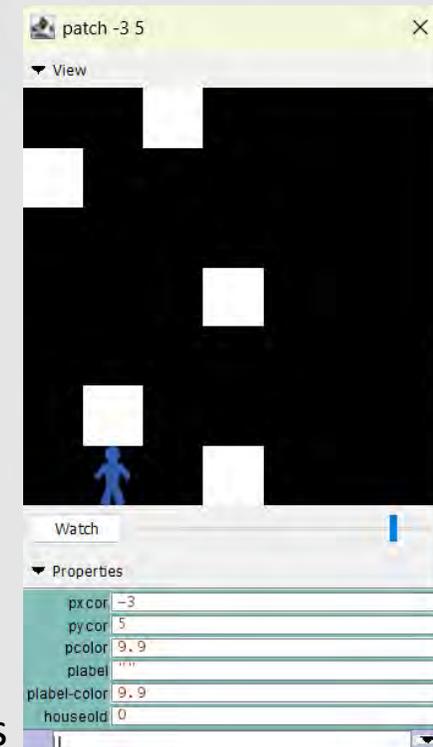
show item 1 (list 1 2 4) > 2
first agent to appear has who 0

Agents are called **turtles** identified by **who** ID

Grid cells are called **patches** and can interact as agents

# NetLogo

**simulated synchrony:** agents execute commands asynchronously, but every agent acts knowing update in the model, based on coding

```
to ask dosomething
  ask turtles [do_A]
  ask turtles [do_B]
  ask turtles [do_C]
end
```

An agent in random order does A, then another agent does A. When all have done A, one random agent does B, then another does B. When all have done B, dosomething is executed

```
to ask dosomething
  ask turtles [
    do_A
    do_B]
end
```

An agent in random order does A then B, then another agent does A then B. When all have done A then B, dosomething is executed

```
to-report sumall [a b]     report 5 6 > 11
  report a + b
end
```

activate native extensions →
```
extensions [gis table csv rnd profiler]
```
```
turtles-own [PRO_COM]
```
agent-class (breed) →
```
breed [hospital hospitals]
breed [women womens]
breed [counselcenter counselcenters]
```
global variable →
```
globals [tuscany distservices distservicesnorm]
counselcenter-own [ID capacity utility womencounsel]
hospital-own [ID hospitalizations utility capacity womenhospital mob
```
agent-class level variable →
```
women-own [pregnant givenbirth selcounsel counselstay rankinglist di
```

command block → that translates model components to be run
```
to setup
;  random-seed 10
  clear-all
  ask patches [set pcolor white]
  gis:load-coordinate-system "C:/Users/LENOVO/Documents/GitHub/childl
  set tuscany gis:load-dataset "C:/Users/LENOVO/Documents/GitHub/chi
  gis:set-world-envelope (gis:envelope-union-of (gis:envelope-of tus
  displaymap
```
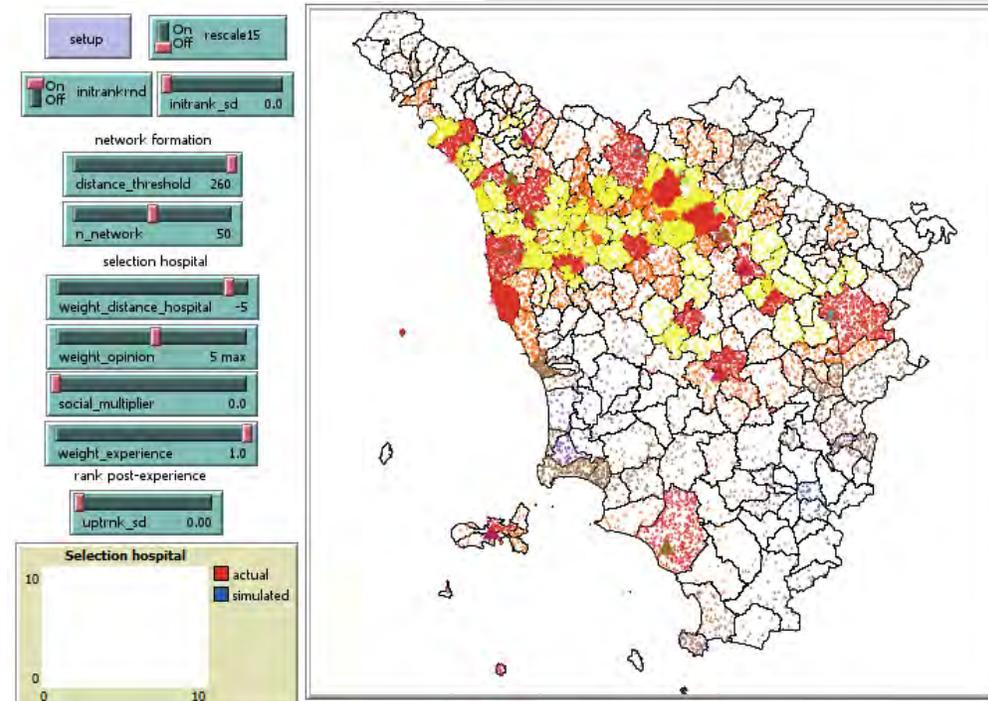
`mean`, `count`, `sort`  primitive reporters

`ask`, `set`, `forward`  primitive commands

`let h who`  local variable existing within a command block (to alleviate memory)
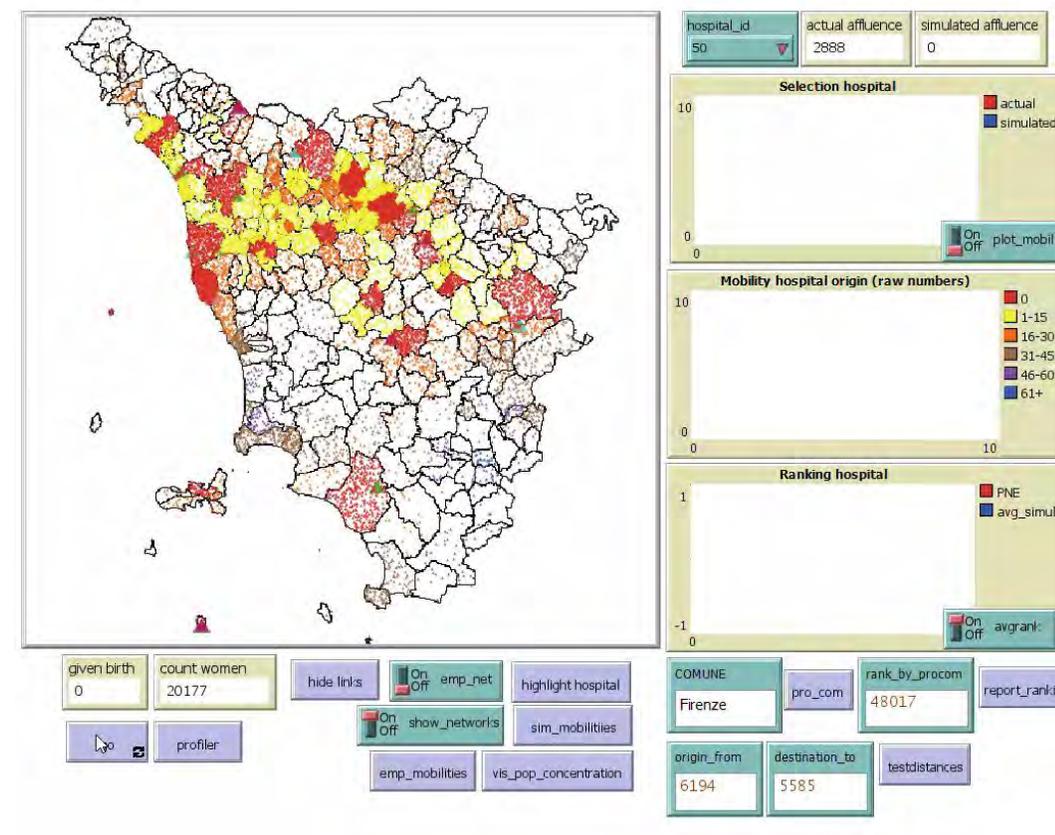
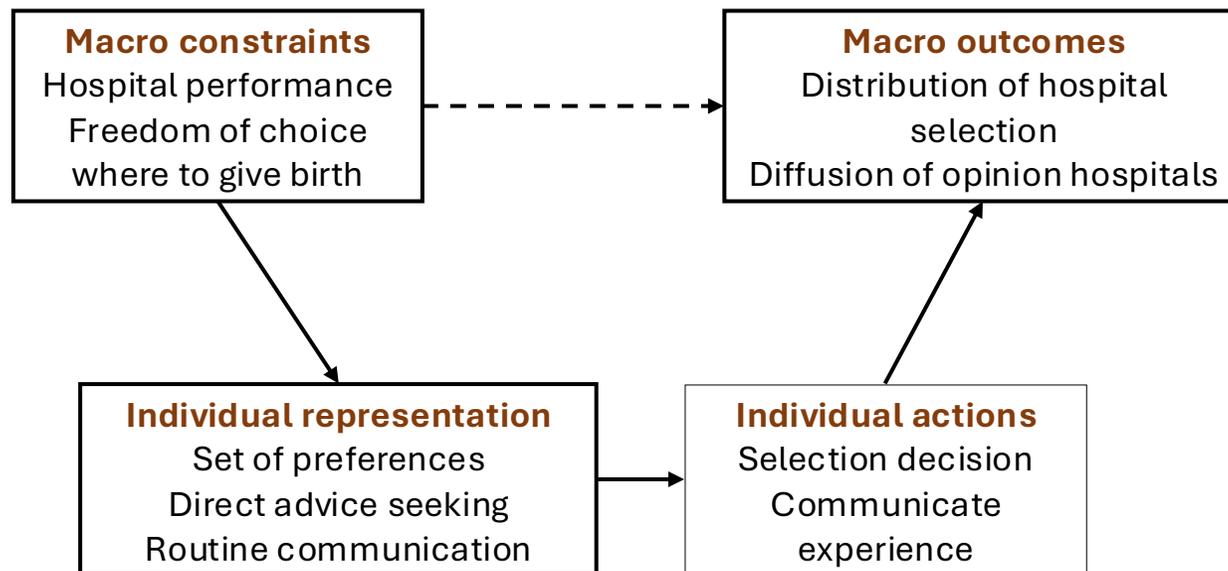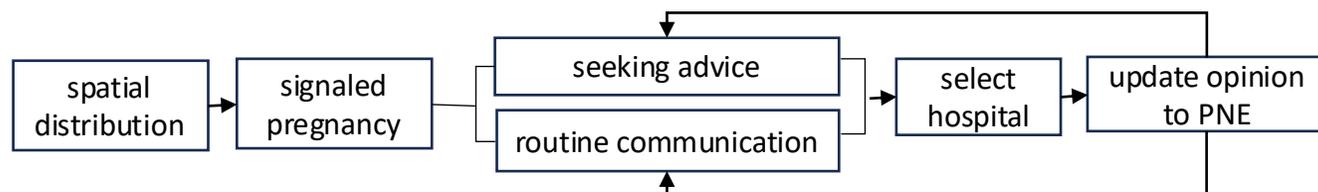**Let's implement in a case study...**

# Childbirth Mobilities: a Geo-Spatial Simulation Approach

- **Context**: While some determinants of hospital maternity selection are identified in the literature, the individual decisional processes, and social influence processes underlying the choice are unknown, neither specific data available.

- **Why ABM**: we can model the weights of preferences for hospital attributes at agents' micro-level, compare different social influence processes and compare how they replicate the data

- **Data available**: Mobility patterns in Tuscany 2023:
  - municipality residencies of women who gave birth
  
  (aggregated and anonymous)
  - municipality hospital where they gave birth
  - ranking of hospital (PNE performance indicator)
  - matrix of ditances
  - shapefile to map geographies to data

Paolillo, Accordino, Pecoraro, https://github.com/RoccoPaolillo/childbirthnet/tree/MIE



**Goal**

- Model a combination of individual decisional processes and social influence processes that can underline the selection of maternity hospital

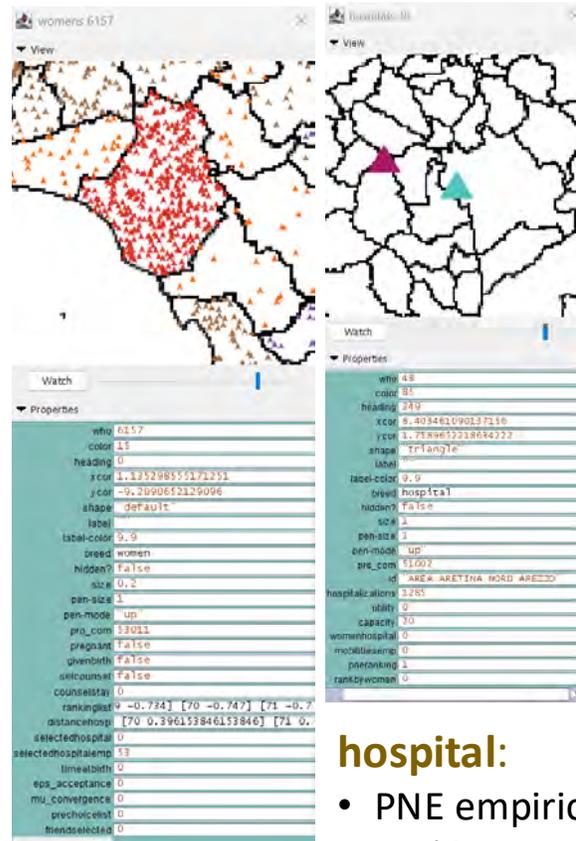- Which condition can best replicate the mobilities we observe?

Let's compare with concepts presented in part one…

# Agents – Who are the actors involved in the phenomenon?

# Parameters – What are the conditions affecting the phenomenon?

**women**:

- they hold an initial random distribution of ranking opinion for each hospital in the region, when they become pregnant activate for choosing one hospital.

- they can expressely ask advice to friends in their municipality or base on common opinion of hospitals from routine communication

- after selecting one hospital, they can vehiculate the opinion of actual performance of hospital (PNE)



**hospital**:

- PNE empirical ranking

```
to setup
  clear-all
  ask patches [set pcolor white]
  gis:load-coordinate-system "C:/../comuni_consultori_2019.prj"
  set tuscany gis:load-dataset "C:/../comuni_consultori_2019.shp"
  gis:set-world-envelope (gis:envelope-union-of (gis:envelope-of tuscany))
  displaymap

  set distservices csv:from-file "C:/../matrice_distanze_consultori.csv"
  set distservicesnorm csv:from-file "C:/../normalized_distance.csv«

  create-counselcenters
  create-hospitals
  create-womens

  let sorted-hospitals sort-by [[a b] -> [hospitalizations] of a >
[hospitalizations] of b] hospital

ask women [options_hospital]
  plot-hospitals

  reset-timer
  reset-ticks
end
```
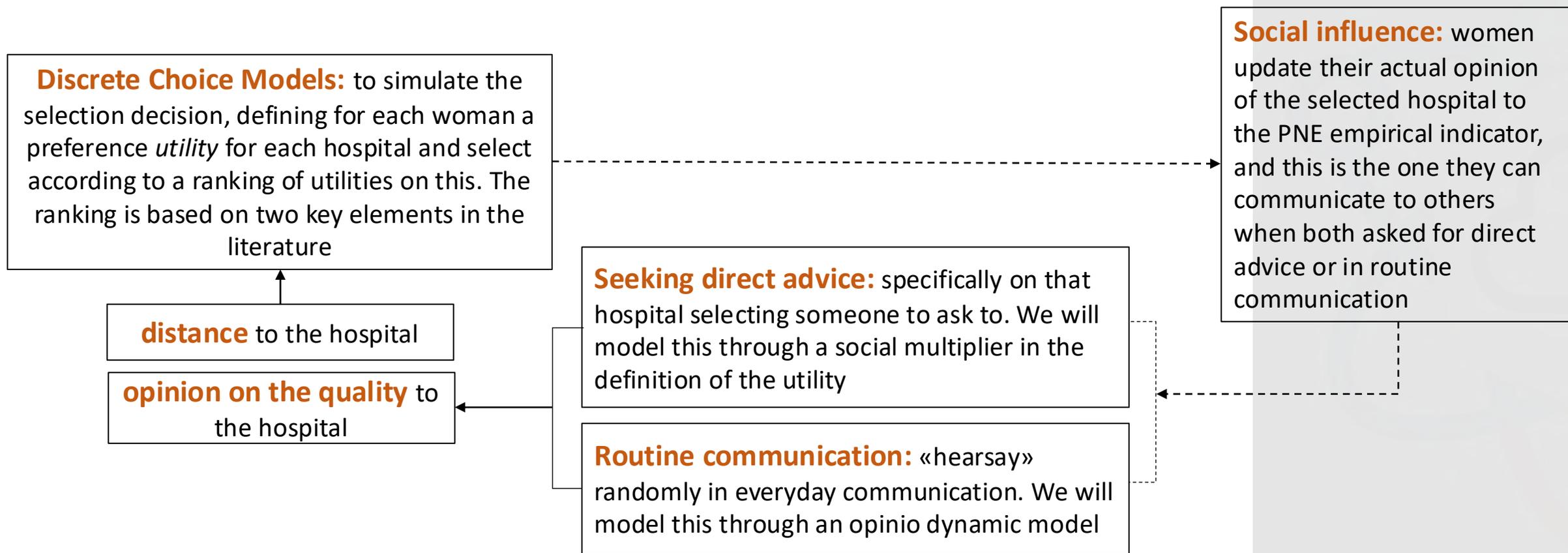
initialize the model with shapefiles

fetch data available

create agent classes

setup time schedule

# Evolution - How does the order of agents' action influence each other?

**Discrete Choice Models:** to simulate the selection decision, defining for each woman a preference *utility* for each hospital and select according to a ranking of utilities on this. The ranking is based on two key elements in the literature

**distance** to the hospital

**opinion on the quality** to the hospital

**Seeking direct advice:** specifically on that hospital selecting someone to ask to. We will model this through a social multiplier in the definition of the utility

**Routine communication:** «hearsay» randomly in everyday communication. We will model this through an opinio dynamic model

**Social influence:** women update their actual opinion of the selected hospital to the PNE empirical indicator, and this is the one they can communicate to others when both asked for direct advice or in routine communication

# Parameters – What are the conditions affecting the phenomenon?

**Discrete choice modeling:** modeling the selection decision of agents, defining a utility (U) for each hospital h, based on a weight (parameter $\beta$) of how two characteristics of each hospital are relevant to the agent:

$D_h$: distance from the agent to the hospital

$O_h$: opinion on quality

The utility is used to define a probability to select that hospital *h* over the others hospital *k*

The higher is $\beta$, the more deterministic the selection is based on differences for attribute, the closer $\beta$ is to 0, the more the selection is random $\varepsilon$

stable softmax to avoid numerical overflow

* 10 to harmonize ß of ranking and distance due to different scales

We can input and manipulate the weight of each characteristic in the mind of agents
- equivalent to coefficients from regressions (clogit), not available
- test the consequences of combining different weights

$$U_h = -\beta(D_h) + \beta(O_h) + \varepsilon$$

```
set utility ((weight_distance_hospital *
(distancefrom*10)) + (weight_opinion * opinionquality) )
```

$$P_h = \frac{e^{(Uh - \max(Uk))}}{\sum e^{(U_k - \max(U_k))}}$$

```
set selectedhospital [who] of rnd:weighted-one-of hospital
[exp(utility - max [utility] of hospital)]
```

**rnd:weighted-one-of** *agentset reporter*

# Parameters – What are the conditions affecting the phenomenon?

**Social multiplier:** A weight $\theta\,[0,1]$ in the definition of opinion quality $O_h$ at the moment of decision.

$\theta = 1$: the opinion quality completly aligns to that of people to whom asked for advice

$\theta = 0$: the advice of others is not taken into consideration

We also included a weighted average to allocate different weights to friend who gave birth to that hospital ($p$), and whose opinion $o$ is based on actual experience, and those who speak for hearsay ($a$)

$a = 1 - w$

$w = 1$, only those who gave birth influence



We can manipulate how influenced people will be by those they seek advice to

We can manipulate how many people advice is searched for advice and how far from hometown

$$\left(\frac{o_w + o_a + o_w + \cdots}{w + a + w + \cdots}\right)$$

```
foreach sort friends  [ z ->
   let weightfriend ifelse-value ([selectedhospital] of z = [who] of self)
   [weight_experience][(1 - weight_experience)]
   set totweightfriend lput weightfriend totweightfriend
   set ranking_othweight lput (table:get [rankinglist] of z [who] of self * weightfriend)
ranking_othweight]
```

$$O_h = OwnOpinion_h + \theta\left(\left(\frac{o_w + o_a + o_w + \cdots}{w + a + w + \cdots}\right) - OwnOpinion_h\right)$$

```
set opinionquality [( opinionquality + social_multiplier *
((reduce + ranking_othweight / reduce + totweightfriend) - opinionquality ) )]
```

# Parameters – What are the conditions affecting the phenomenon?

**Opinion dynamics:** a method to model the routine communication from women who gave birth to others in their municipality, spreading the own (updated) opinion $a$ of that hospital. The receiver agent $i$ accepts to listen if the distance between the own opinion of the hospital and that of the sender falls below a latitude of acceptace $|o_t^i - o_t^a| \leq \varepsilon$.
If so, the receiver aligns to the sender with convergence $\mu$

$\varepsilon = 0$, not communication occurs
$\varepsilon = 1$, everyone is listened
$\mu = 0$, not influence occurs
$\mu = 1$, complete alignment occurs

We set to communicate every 80 time steps and to random 10% of women in municipality

eps_birthtrue  0.8
mu_birthtrue  0.5
eps_notbirth  0.8
mu_notbirth  0.5

We can manipulate how available to listen to those who gave birth and to what degree they will be influenced already in the *hearsay* routine communication

$$if \ |o_t^i - o_t^a| \leq \ \varepsilon$$
$$o_t^i = o_{t-1}^i + \mu(o_{t-1}^a - o_{t-1}^i)$$

```
ask alter [
if abs(table:get rankinglist topic –
table:get [rankinglist] of myself topic) <= eps_acceptance
[table:put rankinglist topic
( table:get rankinglist topic +
(mu_convergence * (table:get [rankinglist] of myself topic –
table:get rankinglist topic)))]
```

# Outcome – How can I identify how the conditions affect the phenomenon?

**BehaviorSpace:** a tool provided by NetLogo to set many experiments to run independently, setting the conditions for each parameter, define specific report measures, how many repetition wanted, and collect data in csv file

- Tools > BehaviorSpace
- Supports batch mode (headless)
- Better with a Server! For computational power, can run on laptop anyway



['eps_birthtrue' 0 2]
will run the conditions with the variable set 0 and 2
['eps_birthtrue' [0 0.1 2]]
will run all the conditions with the variable set from 0 to 2 in increments 0.1
(e.g. 0 0.1 0.2…1.9 2)

# Outcome – How can I identify how the conditions affect the phenomenon?



**What-if scenarios**

- Women communicate with everyone in routine communication ($\varepsilon$ = 2), updating their opinion to the actual experience of those who gave birth, but they do not seek for advice (social multiplier $\theta$ = 0).
- We manipulate the weights for distance [0 -1,-5] and opinion quality updated via routine communication [0 1 5]
  - **condition A**: only distance matters (opinion weight 0)
    - with minimal weight of ditance ($\beta$ = -1), women select hospitals more sparsely and difference between rankings do not emerge. Increasing weight distance ($\beta$ = -5), the simulated distribution overestimates proximity of selected hospitals, and still sort by ranking not appearing
  - **condition B**: when we include also high opinion weight ($\beta$ = 5), simulation results better approximate empirical data when coupled with high weight of distance (green condition $\beta$ = -5, $\beta$ = 5), both distance-wise and ranking-wise

**What would be the effect of seeking advice then?**

# Outcome – How can I identify how the conditions affect the phenomenon?



**What-if scenarios**
- Now women only seek for advice at the moment of selection (latitude opinion dynamic $\varepsilon = 0$), they can be influenced only by those who actually experience the hospital ($w = 1$) or by everyone equally ($w = 0.5$)

- Being influenced by those who gave birth ($w = 1$), hospitals with high PNE are overestimated, and more when the selection is more random by distance ($\beta = -1$)
- Being influenced with equal weight by everyone, those with actual experience and those with random opinion, underestimates the match with empirical data instead

**What if we combine the two types of social influence?**

# Outcome – How can I identify how the conditions affect the phenomenon?



**What-if scenarios**
- Now women undergo both types of influence. When seeking advice, they can be influenced only by those who actually experience the hospital ($w = 1$) or by everyone equally ($w = 0.5$)

- Being influenced by those who gave birth ($w = 1$), still shows overestimation of hospitals with high, and more when the selection is more random by distance ($\beta = -1$)
- Being influenced with equal weight by everyone, those with actual experience and those who updated opinion by *hearsay* in common routine, better approximates the empirical data and reduces overestimation (sligthly)

# Outcome – How can I read the evolution of the phenomenon?

## Conclusions

- The best approximation to empirical mobilities is due to a combination of preference for shorter distance and high opinion quality. But high opinion quality with different weight of distance doesn't produce the same effect. So, distance seems more relevant, and conditioning the diffusion of opinion updates

- Concerning the two modalities of social influence, seeking for advice would overestimate the effect of PNE ranking of hospitals, since the difference in opinion quality becomes more salient. The effect is higher if agents relocate randomly in space, probably because more likely to find high PNE hospitals, that are more and in more populated areas.

- Being exposed to different opinions when seeking advice and in combination with routine communication ameliorates the polarization effect of ranking coming closer to the empirical data, and routine communication seems to suffice

# Outcome - How the phenomenon emerges through agents' interaction?

## Limits and Next Steps

- To better understand the actual differentiation between seeking advice modality and routine communication, looking at the evolution through time and wider parameter space
- To differentiate action of women and measures considering the actual microspace within the region

## BUT

- Overall, we had quite amount of data here
- Sometimes information on sociodemographic population might be missing, how could we do?
- **Synthetic Populations**

**Let's see what synthetic populations are...**

## Thank you! Questions?

rocco.paolillo@cnr.it

@roccopaolillo.bsky.social

# A service for synthetic populations extraction...

Creation of an *Italian Open Science Cloud for the Social Sciences* guided by *Open Science* principles

which shall provide **innovative tools and services** to investigate issues related to the **economic and societal change of contemporary societies** through the enhancement of **research infrastructures**

https://www.fossr.eu/

**Ex-ante** policy analysis evaluation
**Post-ante** policy analysis evaluation
**Conterfactual** policy scenarios

Luciana Taddei
Mario Paolucci *Editors*

**Longitudinal Data Infrastructures in Europe**

Tools for Open Science in Social Science Research

OPEN ACCESS

Springer

FOSSR
Fostering Open Science in Social Science Research
Innovative tools and services to investigate economic and societal change



Chapter 12
**Synthetic Populations in Research Infrastructures**

Rocco Paolillo, Nicholas Roxburgh, Alice Sbrana, Gary Polhill, Evelina Carmen Sabatella, and Mario Paolucci

**12.1 Collective Phenomena and Social Complexity**



- An artificial society is a stylized social system where to study the mechanisms of the phenomenon

- Especially when agent-based modeling is used for policy purposes, the mechanisms observed need to be bounded to the conditions of the system they want to operate in, e.g. Digital Twin systems
- The system needs to be a synthesis of the information available of the target society

- Challenge to identify attributes at micro-level:
  - data not collected
  - separate datasets
  - privacy issues

**Synthetic populations**: a series of techniques to handle available data and replicate attributes of the target population

*synthesis* of information

*generation* of synthetic data
Bigi et al. (2022)

*comparison* with available data

'While a synthetic population is implicitly an artificial population, an artificial population is not necessarily a synthetic population'



Chapuis et al. (2022)

Machine Learning

Every method has its peculiarities and boundaries not set in stone

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

Consiglio Nazionale
delle Ricerche

# Synthetic Reconstruction

- Focus on estimate of unknown joint distributions from data available from marginal distribution
- macro to micro

**Extensions**

- **Multiple Iterative Proportional Fitting (MIPF)** taking the estimated joints as marginal to next step dimension
- **Hierarchical Iterative Proportional Fitting (HIPF)** nested data fitting marginals from one level (e.g. household) to the narrow (e.g. citizens) (Yamaego et al., 2021)
- **Iterative Proportional Updating (IPU)**: from micro data finds weight for cross-category multiplied by marginals and correct backwards -> combinatorial optmization

- Archetype in the **Iterative Proportional Fitting** (IPF, raking)





✅ mathematically transparent, robust, focus on weights

🟥 require ad hoc setting of algorithms, zero-cell problem

FOSSR:SPG

R::synthpop (MIPF)

# Combinatorial Optimization

- Scaling of the synthetic population
- Draw random data and optimize against marginals observed

- **Simulated Annealing (SA)** from the micro-data identifies some seed numbers and compares synthetic marginas to empirical marginals to correct backwards
- **Markov Chain Monte Carlo (MCMC)** random number generations for all intersection and correct backwards

Seed data observations → Initial solution

Random noise → Initial solution

Initial solution → Evaluate

Evaluate → target

target → Rearrange

Rearrange → Evaluate

target → Optimal solution met

✅ overcome some setbacks of IPF family

✅ better with complex intersections (handled *at once*)

🟥 focus on the outcome rather than weights & inner conditions

🟥 computationally more demanding

py::simmaneal

R::MCMCpack   py::PyMC

# Machine Learning Approach

- Most recent in time
- Statistical learning
- micro to macro

- **Generative Adversarial Networks (GAN)**
  - originated from images AI, applied to data
  - two competing (neural) networks:
    - **Generator** who produces random data
    - **Discriminator** that discriminate realistic data from not realistic
    - Goal of Generator is to get better to *deceive* the Discriminator who gets better in discriminate, meaning that synthetic data are very realistic

Training dataset → Sample → Discriminator → ✅ / 🟥

Random noise → Generator → Sample → Discriminator

✅ promising because they integrate the performance of combinatorial methods with transparency, multidimensionally compared to IPF family

🟥 they require a training dataset from which the learning process occurs, where underrepresented groups are likely to be ignored in estimates

- Pre-adjust the training dataset with ad hoc weights to marginals (Falck, 2025)
- Post-adjust the synthetic outcome with weights to marginals

py::PyTorch

# Synthetic Populations Generator (SPG)

Service to enable researchers and policy makers to extract synthetic populations at desired level of information from input dataset

- agent-based modeling
- spatial analysis
- conditional model
- …

**Open Science & Source Software**

https://github.com/RoccoPaolillo/IPF_multidim.git > synthpopgen.py

executability

VRE
(container)

modifiability

demo
(standalone program)

GitHub
(open code)

(Jimenez et al., 2017; Hong et al., 2022)

# Iterative Proportional Fitting (IPF)

**weight** for each cell

update **cell row**

marginal **estimates**

update **cell columns**

update **marginals**

target → stop / no

$$\text{weight} = \frac{\text{observed marginal}}{\text{fitted marginal}}$$

## Multiple Iterative Proportional Fitting (MIPF)

marginal varible 1

marginal varible 2

joint variable 1,2

marginal variable 3

joint variables 1,2,3

marginal variable n...

...

Selected for higher transparency, robustness, light computing, core mechanism common to other method (somehow), but provided more service-oriented direction is guaranteed

First version released (ISTAT gender X age for validation)   10.5281/zenodo.10638800

# Aims of Synthetic Populations Generator (SPG)

- Include **multidimensionality**
- Increase **generalizability** of variable handling
- Enable **automation** input-execution-output
- Customize **filtering** selection

  - leverage estimate of joint and conditional probability over in-cell weight iteration

Tested with **opensalute Lazio** health data:
- gender
- age
- hyptertension (hpt)
- heart failure (hf)

known joints:
- hypertension over age
- heart failure over age

goal of service: identify joint distribution for all combinations

| gender | age | hpt | hf | value |
|--------|-----|-----|-----|-------|
| male | | | | 3073047 |
| female | | | | 3259977 |
| | 30 | | | 1745215 |
| | 3060 | | | 2832088 |
| | 60100 | | | 1755721 |
| | | yes | | 1193445 |
| | | no | | 5139579 |
| | | | yes | 93926 |
| | | | no | 6239098 |
| | 30 | yes | | 3547 |
| | 3060 | yes | | 252543 |
| | 60100 | yes | | 937355 |
| | 30 | no | | 1741668 |
| | 3060 | no | | 2579545 |
| | 60100 | no | | 818366 |
| | 30 | | yes | 424 |
| | 3060 | | yes | 8459 |
| | 60100 | | yes | 85043 |
| | 30 | | no | 1744791 |
| | 3060 | | no | 2823629 |
| | 60100 | | no | 1670678 |

The algorithm

# Open code

https://github.com/RoccoPaolillo/IPF_multidim.git >
synthpopgen.py

**cmd line**



python synthpopgen.py -i input_file_tuples.csv \
  -f (filter)
   'all'
   'gender:female,age:3060' \
-d (display)
   'split'
   'aggregate'\
- v (validate)*
--synth-total 20303*
-o results.csv

→ **pro**: high modifiability for users(rewrite, retest...)
→ **vs**: knowledge coding, dependencies

measure of validation
**Average percentage error** between input marginals (and joints) data and stimated marginals, **RMSE 0.6**

| constraint | observed | predicted | avg_percentage_err |
|---|---|---|---|
| age=30 | 1745215 | 1745215 | 0.0 |
| age=30,hf=no | 1744791 | 1744791 | 0.0 |
| age=30,hf=yes | 424 | 424 | 0.0 |
| age=30,hpt=no | 1741668 | 1741668 | 0.0 |
| age=30,hpt=yes | 3547 | 3547 | 0.0 |
| age=3060 | 2832088 | 2832087 | 3,53E-02 |
| age=3060,hf=no | 2823629 | 2823628 | 3,54E-02 |
| age=3060,hf=yes | 8459 | 8459 | 0.0 |
| age=3060,hpt=no | 2579545 | 2579545 | 0.0 |
| age=3060,hpt=yes | 252543 | 252542 | 0.00039597 |
| age=60100 | 1755721 | 1755722 | 5,70E-02 |
| age=60100,hf=no | 1670678 | 1670678 | 0.0 |
| age=60100,hf=yes | 85043 | 85044 | 0.00117588 |
| age=60100,hpt=no | 818366 | 818366 | 0.0 |
| age=60100,hpt=yes | 937355 | 937356 | 0.00010668 |
| gender=female | 3259977 | 3259977 | 0.0 |
| gender=male | 3073047 | 3073047 | 0.0 |
| hf=no | 6239098 | 6239097 | 1,60E-02 |
| hf=yes | 93926 | 93927 | 0.00106467 |
| hpt=no | 5139579 | 5139579 | 0.0 |
| hpt=yes | 1193445 | 1193445 | 0.0 |

* only if whole population combinations are stimated (-f all)

# Standalone program

.exe local file
py:: tkinter

➔ **pro**: not coding needed, no dependencies
  ➔ **vs**: local CPU, no modifiability

# FOSSR VRE Container

enhance collaboration and tools of researchers through digital platform
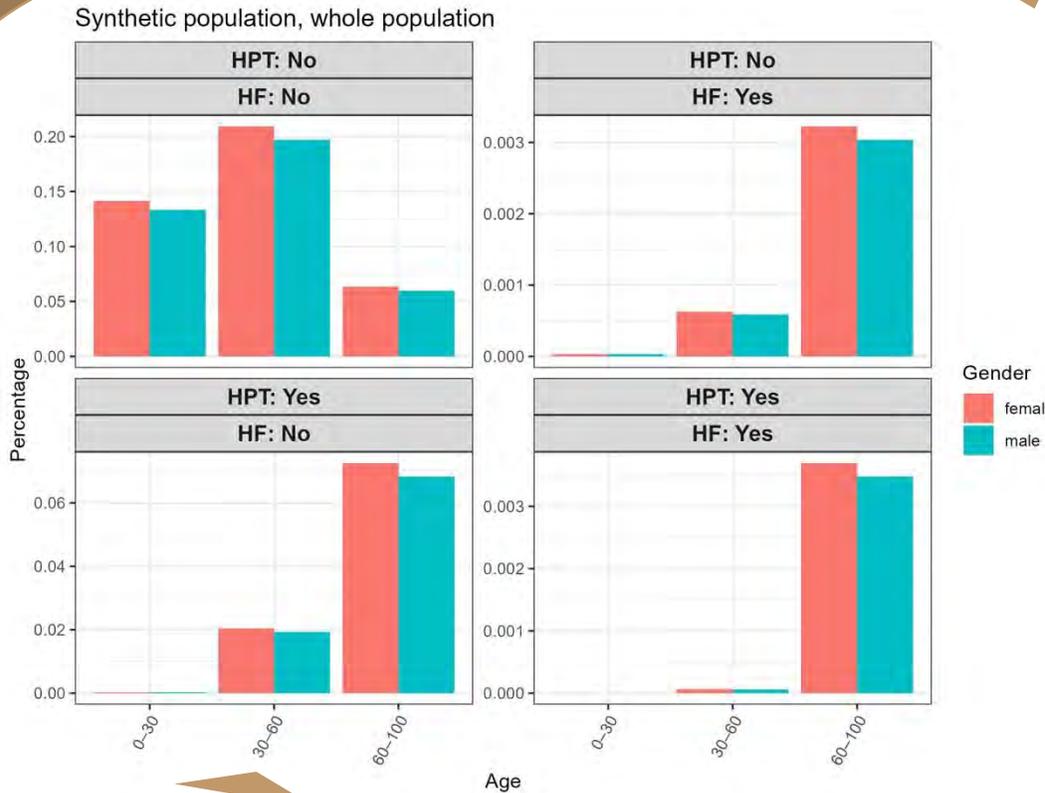
https://fossr.d4science.org

→ **pro**: uses D4Science servers, web-app
→ **vs**: no modifiable, internet-dependent

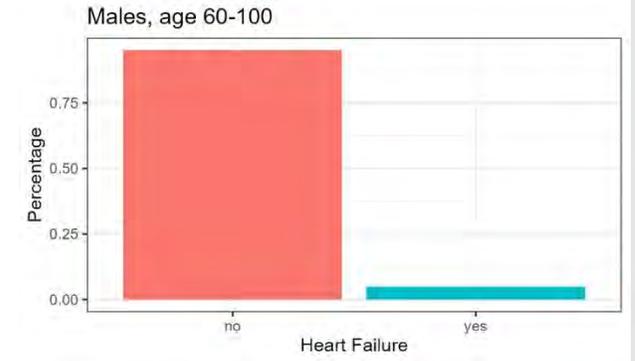VRE > CCP (Cloud Computing Platform)> Synthetic Populations Generator (SPG)

# Output & applications



gender:male,age:60100, hpt:yes -d aggregate: 454844

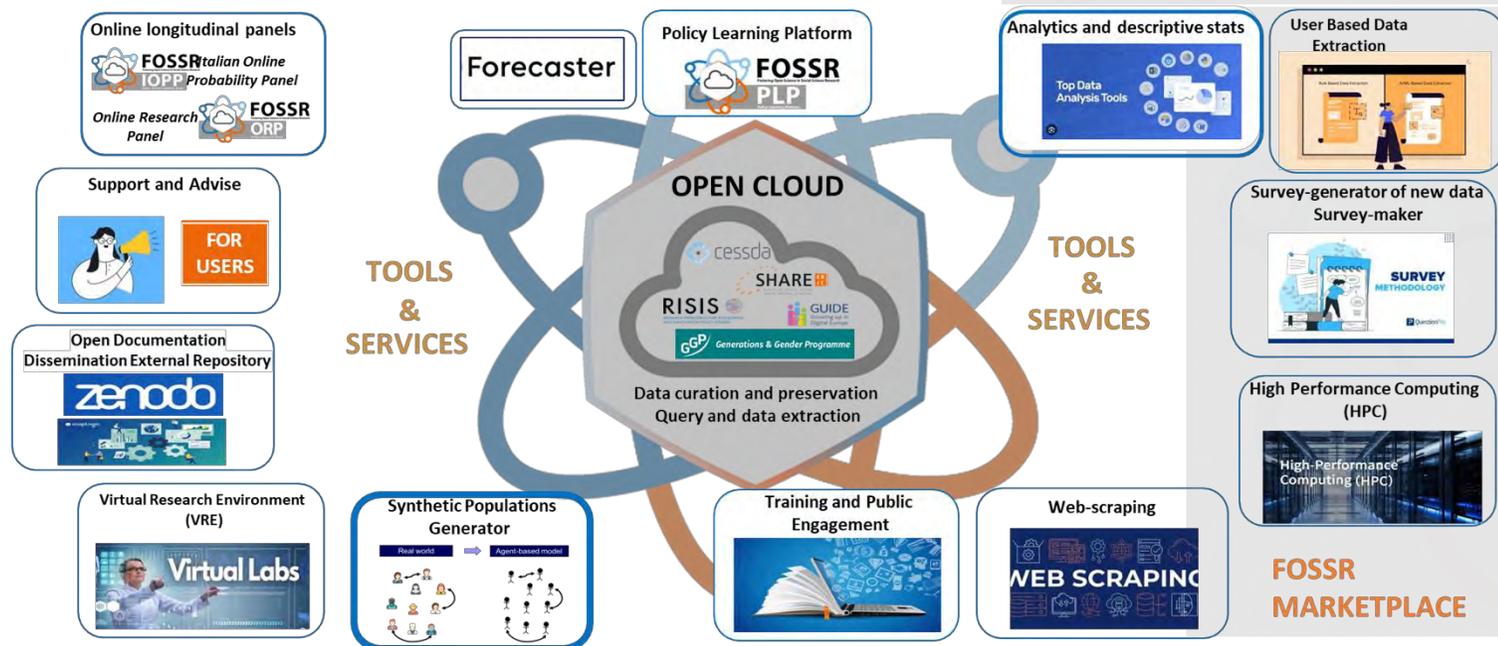gender:male,age:60100, hpt:yes -d split

- individual-level attributes as per request
- preserving original data and restrictions
- identify intersections and heterogeneity in the population
- actor-based models initialization (agent-based models, social network analysis,…)

# Current and future steps

- Assumes marginals from the same total population
- looking at machine learning as micro-to-macro approach
- Automated iteration over spatial units (census tract)
- Automated multi-source input
- Enhance UX experience and assitance
  - **LLM**

https://github.com/RoccoPaolillo/IPF_multidim.git

Deployment into the **FOSSR market place**

**References**

Paolillo, R., Roxburgh, N., Sbrana, A., Polhill, G., Sabatella, E. C., & Paolucci, M. (2025). Synthetic Populations in Research Infrastructures. In Longitudinal Data Infrastructures in Europe: Tools for Open Science in Social Science Research (pp. 153-164). Cham: Springer Nature Switzerland.

Chapuis, K., Taillandier, P., & Drogoul, A. (2022). Generation of synthetic populations in social simulations: a review of methods and practices. Journal of Artificial Societies and Social Simulation, 25(2).

Bigi, F., Rashidi, T. H., & Viti, F. (2024). Synthetic population: A reliable framework for analysis for agent-based modeling inmobility.Transportation Research Record, 2678(11),1–15.

Yameogo, B. F., Vandanjon, P. O., Gastineau, P., & Hankach, P. (2021). Generating a two-layered synthetic population for French municipalities: Results and evaluation of four synthetic reconstruction methods. Journal of Artificial Societies and Social Simulation, 24,27.

Falck, V. (2025). Generating spatial synthetic populations using Wasserstein generative adver-sarial network: A case study withEU-SILC data for Helsinki and Thessaloniki. Preprint.arXiv:2501.16080.

Jiménez, R. C., Kuzak, M., Alhamdoosh, M., Barker, M., Batut, B., Borg, M., ... & Crouch, S. (2017). Four simple recommendations to encourage best practices in research software. F1000Research, 6, ELIXIR-876.

Chue Hong, N. P., Katz, D. S., Barker, M., Lamprecht, A. L., Martinez, C., Psomopoulos, F. E., ... & Yehudi, Y. (2022). FAIR principles for research software (FAIR4RS principles). Zenodo.

# THANK YOU!

**FOSSR DAYS 2026, 4-5-6 February**

✉ **fossr.dissemination@ircres.cnr.it**

🐦 **@fossrproject**

f **fossr.eu**

in **fossr-eu**

▶ **@fossr**

zenodo **zenodo.org/communities/fossr**

**Thank you! Questions?**

[l.cnr.it/fossr-days-2026-registration-form](l.cnr.it/fossr-days-2026-registration-form)

rocco.paolillo@cnr.it
@roccopaolillo.bsky.social