# Big Data,
## between science and society

## Mario Rasetti
### ISI Foundation, Torino – New York

# COMPLEX SYSTEMS
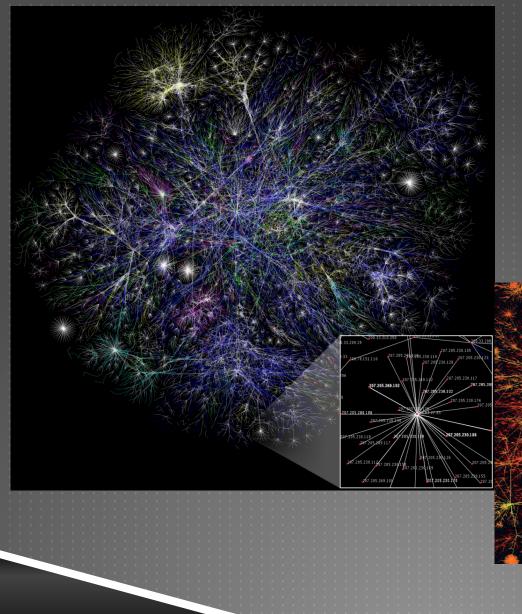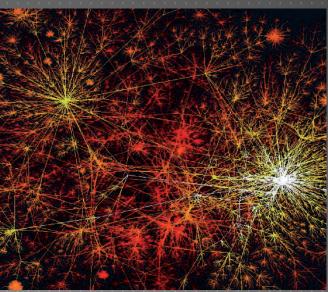
## Bevys of Starlings

## Multi-agent – Multi-scale – Emergent effects

# Network representable

## The Internet
## Clustering (handshake)

✓ **large number of components**

✓ **interactions between components**

✓ **multi-scale hierarchical structures**

✓ **coupling between scales**

✓ **self-organization (*no blueprint*)**

✓ **emergent properties**

✓ **"complex" is <u>more</u> than "complicated"**

✓ **<u>more</u> is different (P.W. Anderson)**

★ **the end of linear thinking**

★ **interdependence and systemic approach**

★ **the problem of causal inference**

★ **<u>universality</u>**

★ **(BIG) <u>DATA</u> representable**

# BIG DATA : how many ?

We live in a world where:

➢ <u>every day</u> 320 billions e-messages and 25 billions SMS are exchanged ; and over 500 millions photographs are posted and shared on Facebook ;

➢ more than 4.5 billion people (> ½ world population!) have a mobile phone (or similar device)

➢ the total quantity of information created and exchanged has reached, in 2013, 4 *zettabyte*s [a zettabyte means $10^{21}$ bytes (characters) : the 1250 pages of Tolstoy's *War and Peace* could be stored 323 billion times in a zettabyte, and the whole content of the Library of Congress 4 million times] and grows 40% per year (will reach in 4 years a *yottabyte*, $10^{24}$)

- A world in which <u>every year</u> 1 billion cars take the road, over 2 billions people fly on airplanes; population and urbanization, commercial exchanges and migrations grow tumultuosly - entangled with a parallel universe of more and more sophisticated technologies, giving life to a unique, interconnected **socio-technical system**.

- The growing complexity of this scenario hides enormous opportunities as well as potential **risks**, generated just by the vulnerability that lies in the interdipendencies between different systems.

- For this society requires an increased capacity of making **predictions**, to anticipate, evaluate and correlate risks and understand the **complexity** of the novel, different, world that new technologies are giving life to.

- Big Data 2.0 will be <u>sensors</u>

**Thanks to Data,** digital imaging is tracking this world more and more closely

- this allows us to use computation to **extract patterns and establish causal inferences** using tools, e.g., from A.I. : **data mining**, **machine learning**, **statistics** → <u>**MINING ALGORITHMS**</u> **(Google)**

- mathematical modeling and forecast may then take place on a **data-rich landscape,** fed by **data streams** from multiple sources : (BI)<u>SIMULATION</u>

- thus we can **assess** at unprecedented speed and scale **our worldview against reality,** and feed it back to models, and help – through virtual scenarios – decision makers & policy makers to behave in a more rational way : <u>PREDICTIONS</u>

# The great challenge :
# understand, govern and control the process

## Data → Information → Knowledge → Wisdom

Data ➔ Information    ➔ Data Mining

Information ➔ Knowledge ➔ Patterns & Correlations

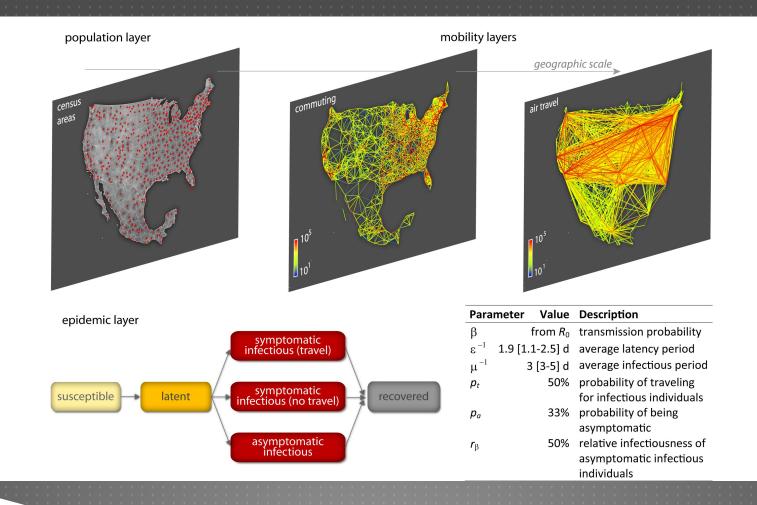Knowledge ➔ Wisdom    ➔ A 'Field Theory of Data'

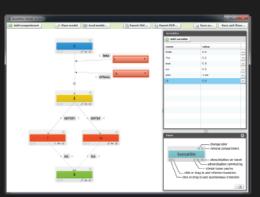## The methodology

**RISK: decisions & policies - scenarios**

**complex systems modeling**

**machine learning**
**data mining**
**natural language processing**
**modern mathematics**

**IT infrastructures**

**Data → Virtual Reality → (Simulation) → Predictions**

# Predicting Pandemics



population layer

mobility layers

*geographic scale*

census areas

commuting

air travel

epidemic layer

| Parameter | Value | Description |
|---|---|---|
| $\beta$ | from $R_0$ | transmission probability |
| $\varepsilon^{-1}$ | 1.9 [1.1-2.5] d | average latency period |
| $\mu^{-1}$ | 3 [3-5] d | average infectious period |
| $p_t$ | 50% | probability of traveling for infectious individuals |
| $p_a$ | 33% | probability of being asymptomatic |
| $r_\beta$ | 50% | relative infectiousness of asymptomatic infectious individuals |

susceptible → latent → symptomatic infectious (travel) / symptomatic infectious (no travel) / asymptomatic infectious → recovered

GLEaMviz is a multiplatform application that allows to interactively program and generate data with the GLEaM computational model.
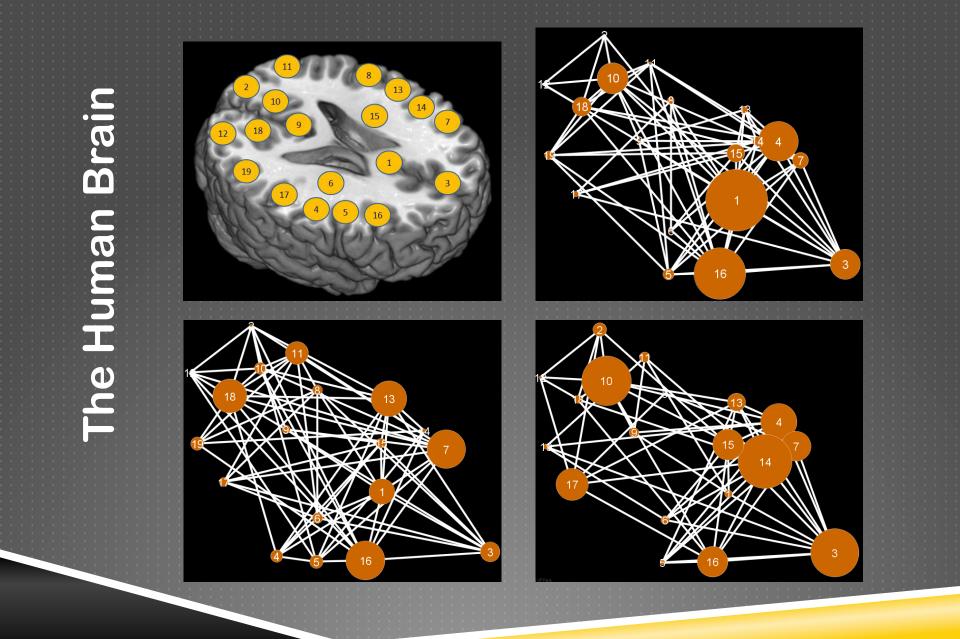
With the GLEaMviz application you can:
· Configure the epidemic model by setting the disease natural history.
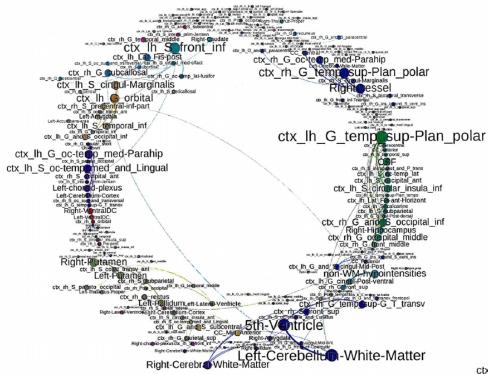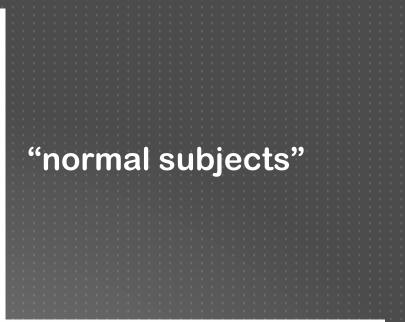· Define the simulation scenario by including environmental effects and the initial conditions of the outbreak.
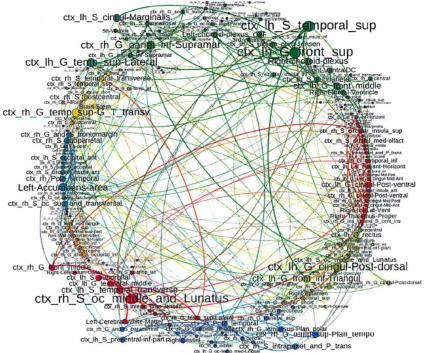· Define mitigation and containment policies and analyze their effectiveness.
· Explore simulation results through dynamic plots and maps.
· Download and share the generated data through a user-friendly interface.

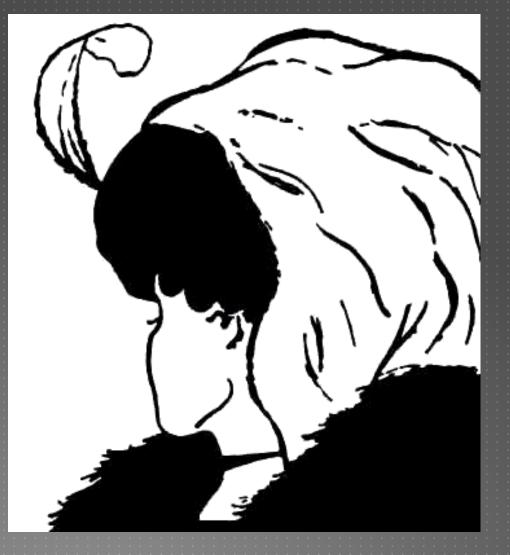"normal subjects"

Psychoactive drug treated subjects

13

**Francisco Varela**





**Ambiguous pictures**

# MacKinsey BD Report (2012)



*Big data—capturing its value*

**$300 billion**
potential annual value to US health care—more than double the total annual health care spending in Spain

**€250 billion**
potential annual value to Europe's public sector administration—more than GDP of Greece

**$600 billion**
potential annual consumer surplus from using personal location data globally

**60%** potential increase in retailers' operating margins possible with big data

**140,000–190,000** (USA)
more deep analytical talent positions, and

**1.5 million**
more data-savvy managers needed to take full advantage of big data in the United States

# IBM estimate (2014) - worldwide



**4.4MILLION** data scientists needed by 2015

ONLY 1/3 FILLED

DATA SCIENTIST

## The hardware side

**What do we need to achieve?**

Data comes today from different sources, in different formats and at different times.

Easy access requires:

- Scalability
- Fast & efficient data storage and access
- Tiering and compression
- Friendly visualization

# What infrastructure features do we need?

Big data analytics are useful only if the insights are easily available to people and processes that need them.

This requires an infrastructure with:

- Self-healing capability
- High-grade continuous management
- Single- and multiple-site solutions

## What speed do we need?

Also very efficient big data analytics tools are useless if insights require infinite time.

This requires an infrastructure capable of extracting information from big data and performing data analytics in real time:

- Low latency capability
- Ability to scale up, in and out fast
- Analytics customized for the system of record

- Over 50% business leaders don't have access to the data they need
- Only 38% of organizations are prepared to deal with the onslaught of big data
- Infrastructures are needed able to handle up to 20 Petabytes of data
- Presently available storage systems can be improved to store in the same space up to 5 times the amount of data now stored

- In 2020 over 30 billion devices will be in use
- By 2017, 60% of cars will have connected services
- Infrastructures providing data analytics with essentially > 95% fewer outages and performance problems than present commodity platforms will be necessary
- Big data and data analytics infrastructures should deliver up to 99.999% uptime

- ➤ **40% of executives need real-time information accessible within the process time**
- ➤ **Estimated minimum need is capability of processing complex queries thousands of times faster**
- ➤ **Infrastructures providing data analytics results with essentially > 95% higher speed than present commodity platforms will be necessary**
- ➤ **Capacity of storing data to be processed enabling delivery insights at least 20 times larger**