

Big Data e Linked Open Data per la statistica ufficiale: verso nuove forme di conoscenza generata dai dati



Stefano De Francisci

25 maggio 2017

1. Demistificare i Big Data
2. Il contesto dei Big Data nella statistica ufficiale
3. Altri punti di vista
4. Problemi aperti
5. Esperienze correnti in Istat

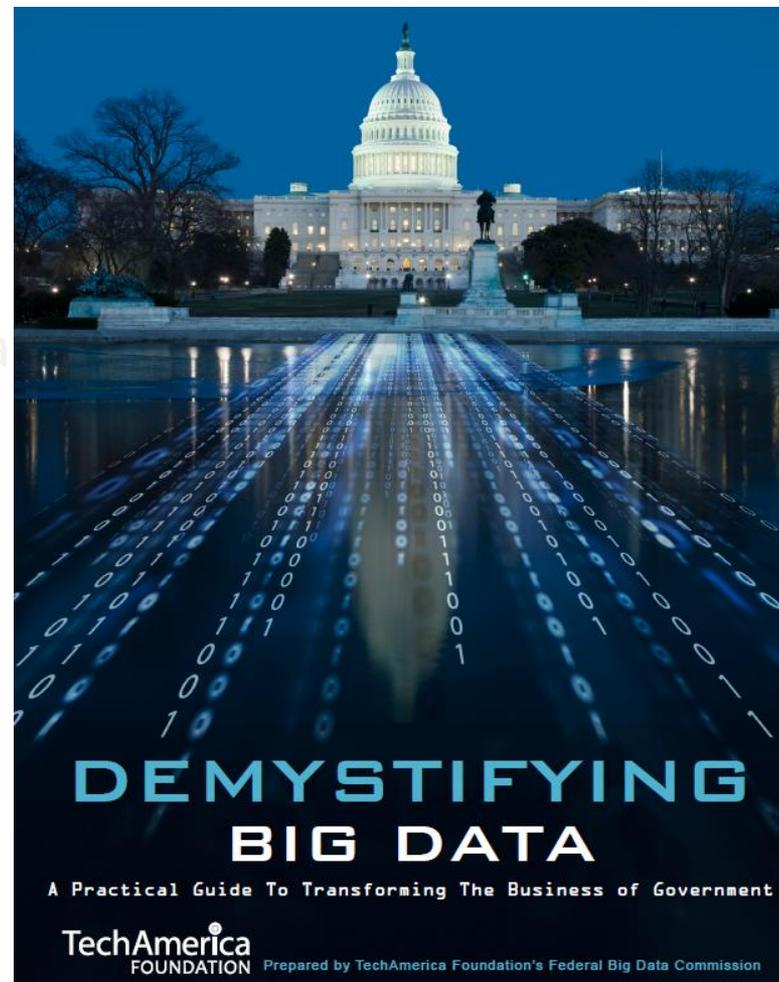
1. Demistificare i Big Data

2. Il contesto dei Big Data nella statistica

3. Altri punti di vista

4. Problemi aperti

5. Esperienze correnti in Istat





Big Data ... and the Next Wave of *InfraStress*

John R. Mashey
Chief Scientist, SGI

Technology Waves:
NOT technology for technology's sake
IT'S WHAT YOU DO WITH IT
But if you don't understand the trends
IT'S WHAT IT WILL DO TO YOU



Breve storia dei Big Data secondo Gil Press

 Volume  Velocità  Nascita dei Big Data moderni

1 simbolo per inch³ → 500 caratteri per inch³ → 1.25X10¹¹ byte per inch³

Densità dei dati: Sumeri → Gutenberg → 2000
(Data Communications)

“Tracking the Flow of Information” 1983
(Science)

“Volume of Information” 1981
(NSI Ungheria)

1980 **“Data expands to fill the space available”** (4° IEEE Symposium)

1975 **“La produzione di informazione cresce più rapidamente del suo consumo”**
(Ministero Poste giapponese)

1971 **“Un uomo si misura da quanti bytes occupano il suo dossier”**
(A. Miller)

1967 **Information explosion** (B.A. Marron, P. A. D. de Maine)

1961 **Crescita conoscenza scientifica attraverso aumento esponenziale di giornali e riviste** (D. Price)

1944 **2040: 200 Milioni volumi... 6000 miglia di scaffali.... 6000 catalogatori** (F. Rider)

<https://www.forbes.com/sites/gilpress/2013/05/09/a-very-short-history-of-big-data/#533cdbfc65a1>

(J.R. Mashey)

1990

1994

1997

1998

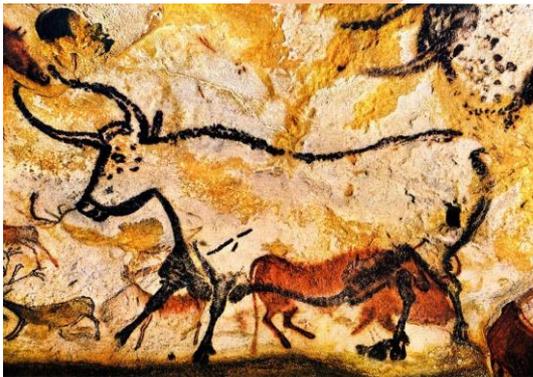
2001

(M. Cox, D. (J.R. Mashey)
Hellsworth)

Invenzione del termine Big Data?

(Gartner)
Scoperta delle 3v

Breve storia dei Big Data secondo Gil Press



Quadro di contesto dei Big Data. Voci critiche

«È difficile dare un senso a grandi quantità di dati in modo significativo»

The Big Data Deception

What was I supposed to do - call him for cheating better than me, in front of the others?

1 conference, read a blog (ahem) or open a [tech mag](#) without someone talking about [Big Data](#) these days. Now the next person whenever new techniques, approaches, tools, frameworks, whatever come along, but equally, my penchant for hype, it's important to keep one eye out for denuded emperors keen to show off their new duds with sales targets to hit.

nds after it was announced that Barack Obama had won the US e had, along with only a handful of others, predicted the result with He hadn't used a [psychic octopus](#), a [gifted cat](#), [halloween masks](#) or ss, he used Big Data.



its and along with everyone else admired Nate Silver's analytical skills. ter all that data, adjust for various biases and survey limitations, was rely if you needed a business case for Big Data that was it. This is a dict who's going to win an election and he's 100% right for every side, it also crossed my mind that I wouldn't want to be him at the cause the media will whip up such high expectations that, even if he result right, you're going to feel like

«È difficile che le imprese IT tradizionali adottino i Big Data»

WIRED

PARTNER CONTENT STEVE DODSON, FREELER
BIG DATA, BIG HYPE?



«I Big Data sono solo una grande quantità di dati»



«As with so many buzzwords, "big data" is a vague term, often thrown around by people with something to sell»

abcNEWS U.S. International Politics Lifestyle Entertainment Virtual Reality

Growing Doubts About Big Data

April 8, 2014
By GARY LANGER

There's quite a kerfuffle going on in the world of big data, with a range of prominent articles in the past month suggesting it's not the analytical holy grail it's been made out to be. Taken together, these pieces suggest the start of a serious rethink of what big data can and can't actually do.

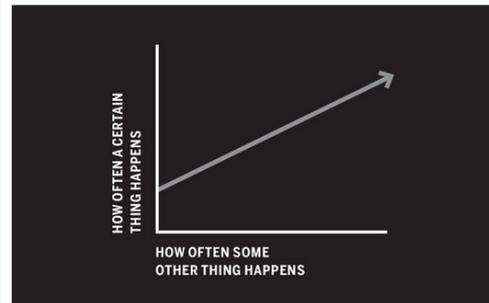
«The core challenge is that most big data that have received popular attention are not the output of instruments designed to produce valid and reliable data amenable for scientific analysis»

The New York Times

The Opinion Pages OP-ED CONTRIBUTORS

Eight (No, Nine!) Problems With Big Data

By GARY MARCUS and ERNEST DAVIS APRIL 6, 2014



«Non è vero che grandi moli di dati comportino cambiamenti nei modi in cui interagiamo con essi per esplorarli e dare loro un senso»

Big Data, Big Ruse

Stephen Few, Perceptual Edge
Visual Business Intelligence Newsletter
July/August/September 2012

If you're like me, the mere mention of Big Data now turns your stomach. Nearly every business intelligence (BI) vendor, publication, and event has Big Data flashing in neon colors in Times Square dimensions. Never before have I seen an idea in the BI space elicit this much obsession. Why all the fuss? Why, indeed. Essentially, Big Data is a marketing campaign, pure and simple.

«Non è vero che le nuovi fonti di dati siano davvero nuove»

INTERNET

Google Flu Trends' Failure Shows Good Data > Big Data

by Kaiser Fung

CUSTOMERS

Big Data's Dangerous New Era of Discrimination

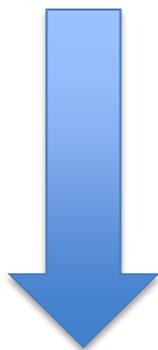
by Michael Schrage

http://www.perceptualedge.com/articles/visual_business_intelligence/big_data_big_ruse.pdf; <http://www.julianbrowne.com/article/viewer/big-data-deception>
<http://www.wired.com/2014/04/big-data-big-hype/>; <http://www.ft.com/cms/s/2/21a6e7d8-b479-11e3-a09a-00144feabdc0.html#axzz2yQ2QQfQX>
http://www.nytimes.com/2014/04/07/opinion/eight-no-nine-problems-with-big-data.html?_r=1; <http://abcnews.go.com/blogs/politics/2014/04/growing-doubts-about-big-data/>; <https://hbr.org/2014/03/google-flu-trends-failure-shows-good-data-big-data>

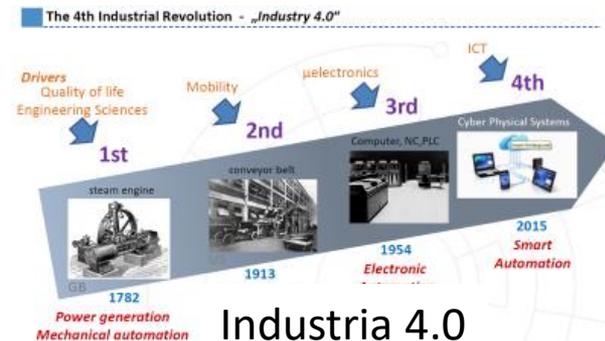
Statistica ufficiale



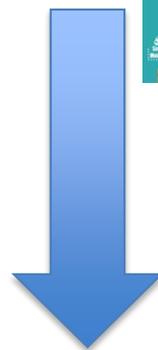
Scienza



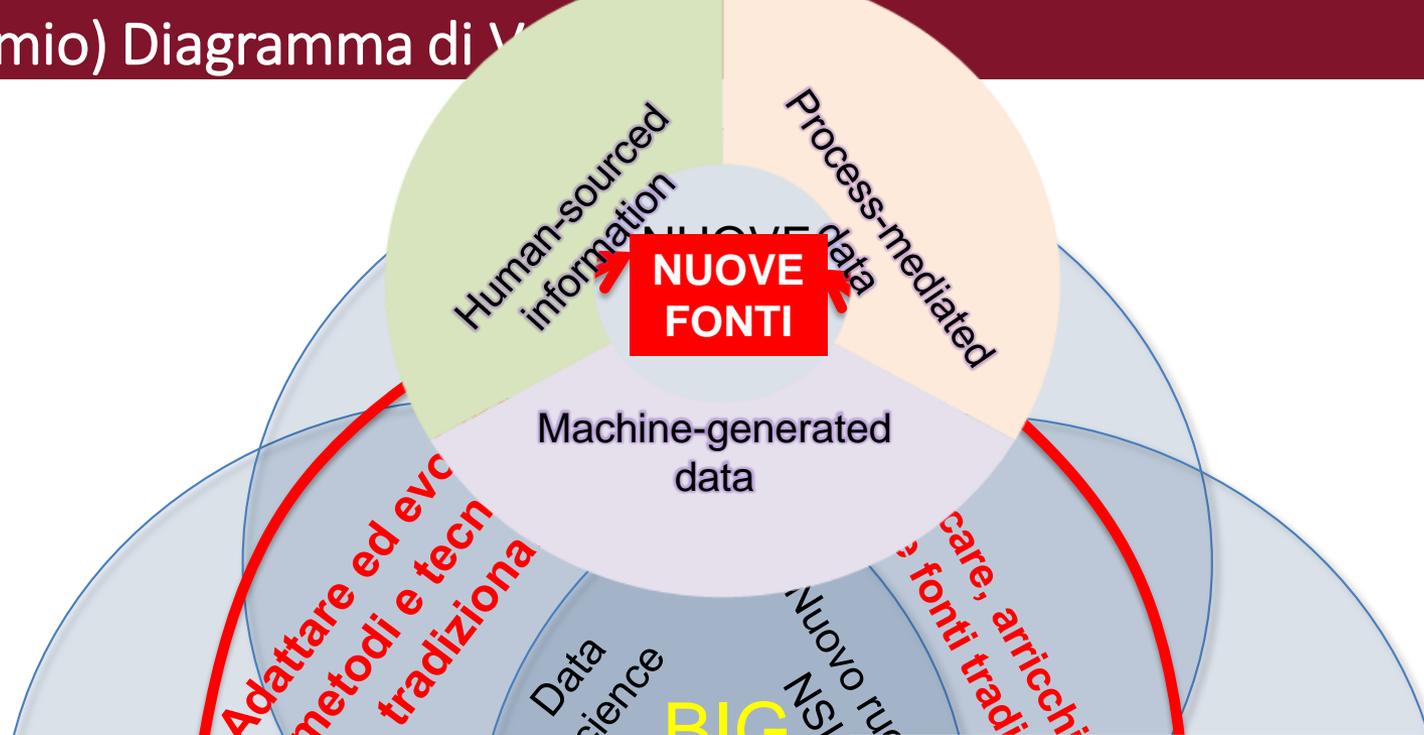
Business



Smartness



- **NUOVE FONTI** che affiancano, arricchiscono e sostituiscono le fonti tradizionali
- **NUOVI METODI E TECNOLOGIE** per adattare ed evolvere metodi e tecniche tradizionali
- Sono suggeriti, sostenuti e favoriti **NUOVI USI** dell'informazione



BIG DATA LANDSCAPE, VERSION 3.0

Export: Acquisition of IPO

NUOVI METODI E TECNOLOGIE

© Matt Turck (@matsturck), Sution Dong (@sutiondong) & FirstMark Capital (@firstmarkcap)

BIG DATA

Smart World

re, suggerire,

Libelium Smart World

NUOVI USI

- Air Pollution:** Control of CO₂, emissions of factories, pollution emitted by cars and boats, gaseous pollutants in forests.
- Forest Fire Detection:** Monitoring of environmental parameters for the conditions to define alert areas.
- Wine Quality Enhancing:** Monitoring and measure and track diameter in response to control the amount of sugar in grapes and grapevine health.
- Outspiring Care:** Risk signs monitoring in high pressure centers and hotels.
- Sportsmen Care:** Control of operating conditions of the playing ground to ensure its survival and health.
- Structural Health:** Monitoring vibrations and wear on buildings, bridge, embankments.
- Smartphones Detection:** Detect Phone and Android devices and to generate log device which works with WiFi or Bluetooth interfaces.
- Perimeter Access Control:** Access control to restricted areas and selection of people in non-authorized areas.
- Radiation Levels:** Distributed measurement of radiation levels in nuclear power stations, nuclear power plants, waste storage sites.
- Smartphones Detection:** Detect Phone and Android devices and to generate log device which works with WiFi or Bluetooth interfaces.
- Perimeter Access Control:** Access control to restricted areas and selection of people in non-authorized areas.
- Radiation Levels:** Distributed measurement of radiation levels in nuclear power stations, nuclear power plants, waste storage sites.
- Smart Roads:** Warning messages and diversions according to climate, weather, and compromised events like accidents or traffic jams.
- Smart Lighting:** Intelligent and weather adaptive lighting systems.
- Intelligent Shopping:** Getting advice on the spot of sale according to customer habits, preferences, presence of change components for them in shopping lists.
- Noise Urban Maps:** Sound monitoring in bar areas and central areas in the town.
- Water Leakages:** Detection of water pressure, pressure and pressure variations along pipes.
- Vehicle Auto-diagnosis:** Information collection from CAN bus to send real time across the smartphone to provide advice to drivers.
- Item Location:** Control of movement items in big surfaces like warehouses or harbours.
- Waste Management:** Intelligent rubbish bins, autonomous to optimize the trash collection routes.
- Smart Parking:** Monitoring of parking spaces availability in the city.
- Golf Courses:** Selective irrigation in dry zones to reduce the water resources required in the green.
- Water Quality:** Stock of water turbidity, pH and the use for future and ecology for products use.

libelium
www.libelium.com

BIG DATA LANDSCAPE, VERSION 3.0

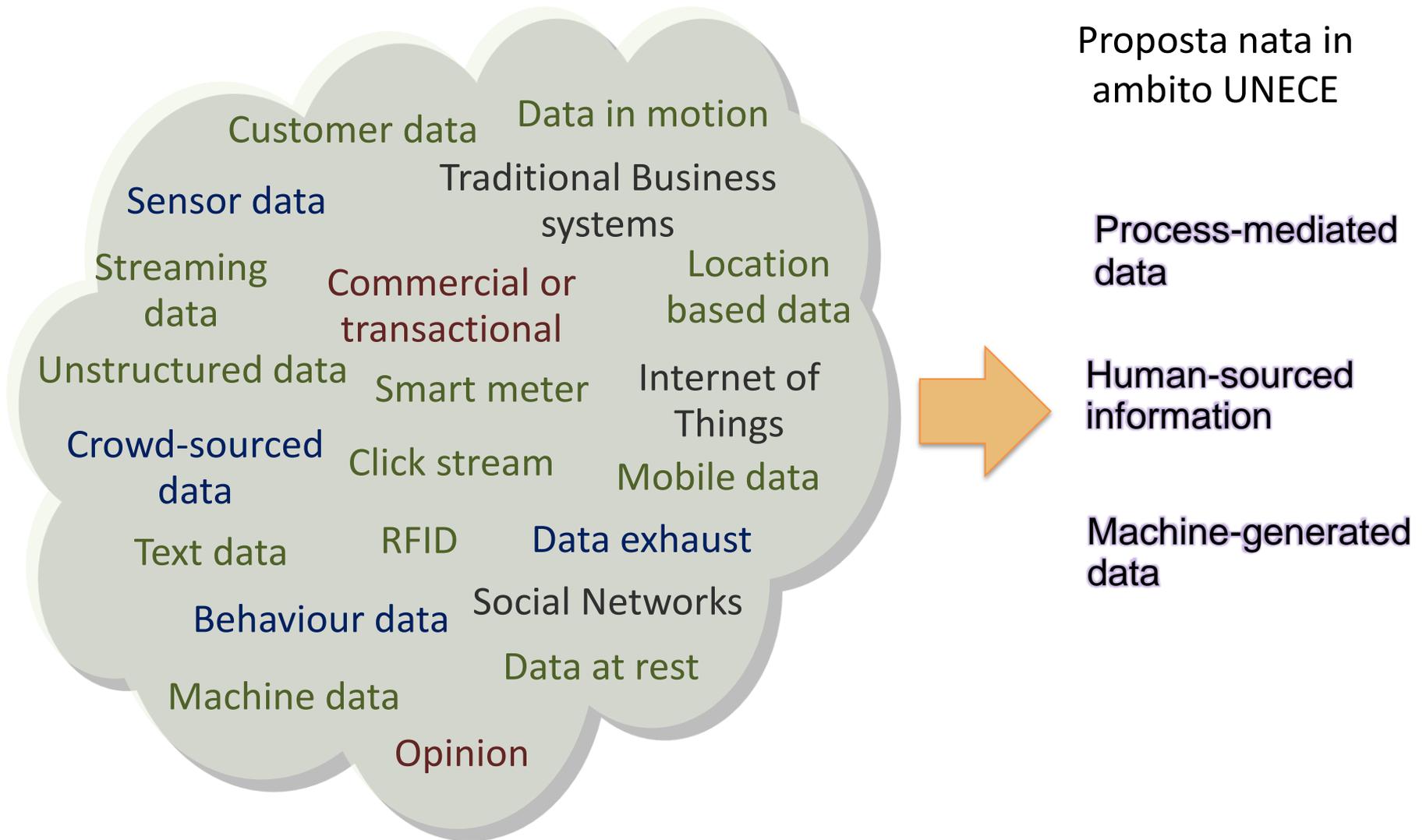
Exit: Acquisition or IPO



Categoria	2014	2016
Analytics	120	169
Applications	67	167
Cross Infrastructure/Analytics	12	13
Data Sources	29	56
Incubators & Schools	0	7
Infrastructure	99	136
Open Source	31	62
Totale complessivo	358	610

© Matt Turck (@matturck), Sujan Dong (@sujan Dong) & FirstMark Capital (@firstmarkcap)

Classificazione dei Big Data



Classificazione delle nuove fonti Big Data

Social Networks



Dati prodotti dall'interazione con mezzi di informazione e social media o tramite dispositivi (anche mobili)

Blog, Twitter, Facebook
User-generated maps

Human-sourced information

Traditional Business systems



Dati prodotti da sistemi transazionali tradizionali e in modo passivo:

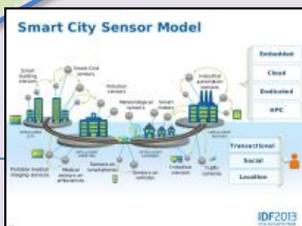
Scanner data
Log ricerca,
Record medici,
Transazioni commerciali e bancarie

Process-mediated data



Machine-generated data

Internet of Things



Dati prodotti da sensori e macchinari utilizzati per misurare e registrare eventi e situazioni nel mondo fisico: immagini satellitari, sensori stradali e di traffico, sensori climatici e ambientali, ecc

1. Demistificare i Big Data

2. Il contesto dei Big Data nella statistica ufficiale

3. Altri punti di vi



« ... Big Data is also potentially very interesting as an input for Official Statistics; either for use on its own, or in combination with more traditional data sources such as sample surveys and administrative registers»

Gennaio 2013

High Level Group for Modernization of Statistical Production and Services

...della statistica ufficiale

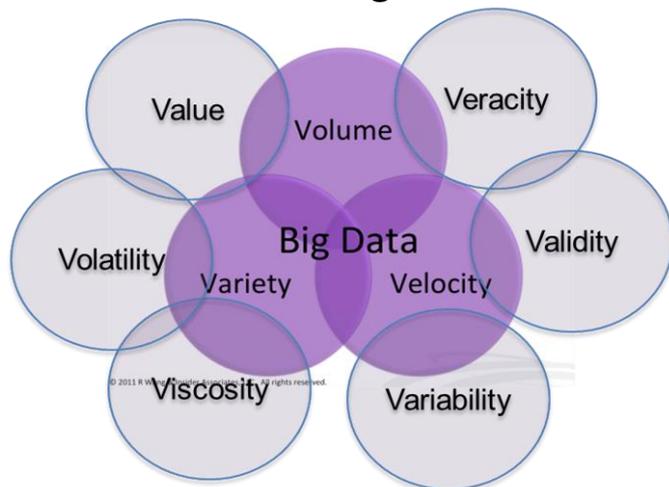
- imparzialità
- affidabilità
- obiettività
- indipendenza scientifica
- efficienza economica
- riservatezza statistica
- non comporta oneri eccessivi per gli operatori economici

CARTA DEI DIRITTI FONDAMENTALI

Art. 338 del trattato sul funzionamento dell'UE

http://europa.eu/pol/pdf/consolidated-treaties_it.pdf

...dei Big Data



La sfida

Modernisation Committee on Products and Sources

*"We must move from a paradigm of producing the best estimates possible from a **survey** to that of producing the best possible estimates to meet user needs from **multiple data sources**" (Conny Citro)*

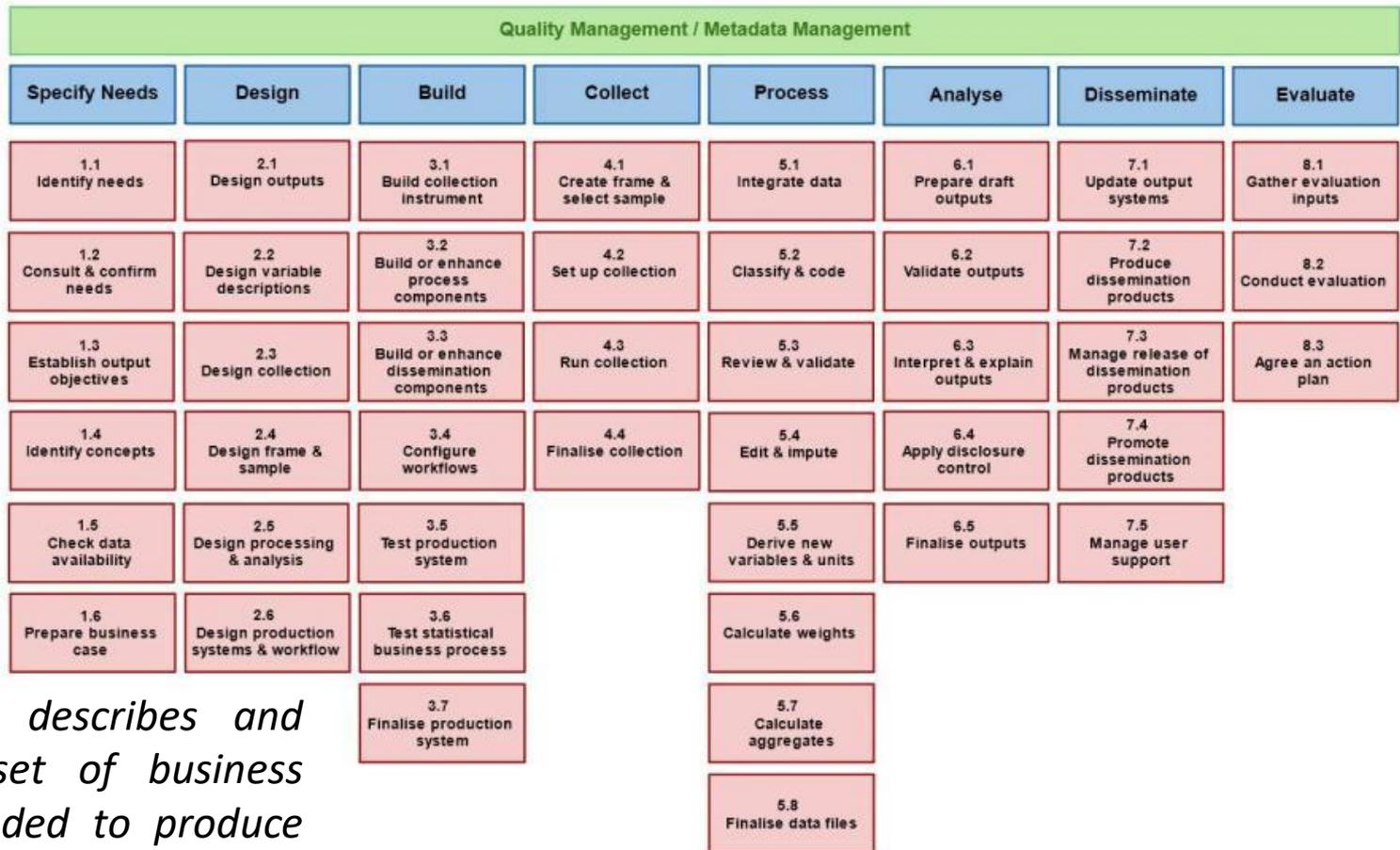
Produce stable output with unstable ever changing inputs

Opportunities
Big data and administrative data
New technologies

Challenges

Produce stable output with unstable ever changing inputs

(HLG-MOS)

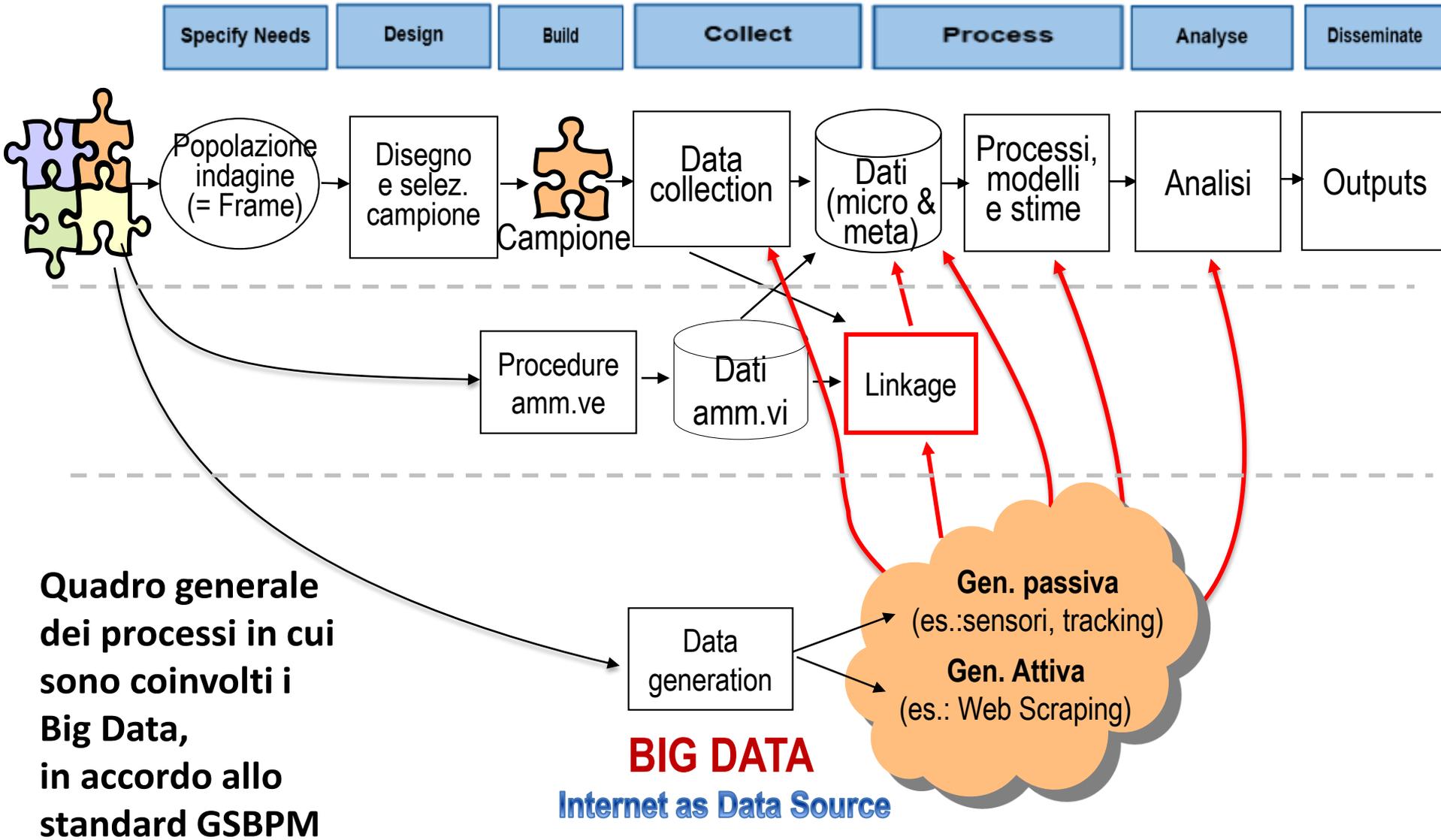


“The GSBPM describes and defines the set of business processes needed to produce official statistics.

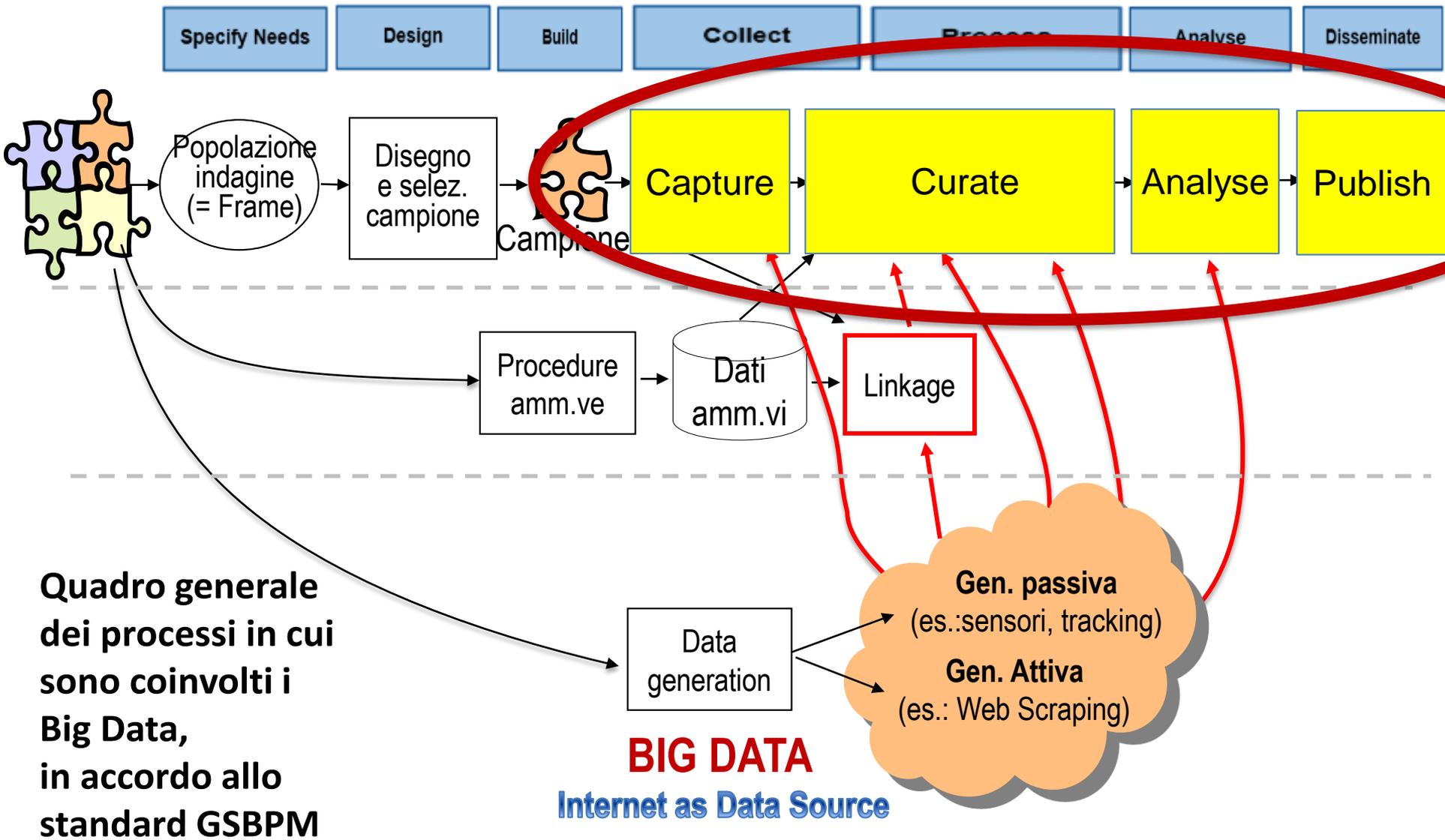
It provides a standard framework and harmonised terminology to help statistical organisations to modernise their statistical production processes, as well as to share methods and components.”

<http://www1.unece.org/stat/platform/display/GSBPM/I.Introduction#I.Introduction-Toc375051192>

Big Data: possibile uso nelle fasi del processo statistico



Big Data: possibile uso nelle fasi del processo statistico



1. Demistificare i Big Data

2. Il contesto dei Big Data nella statistica ufficiale

3. Altri punti di vista

4. Problemi aperti

5. Esperienze correnti in statistica

**La ricerca
sociale**

**Il legame col
territorio**



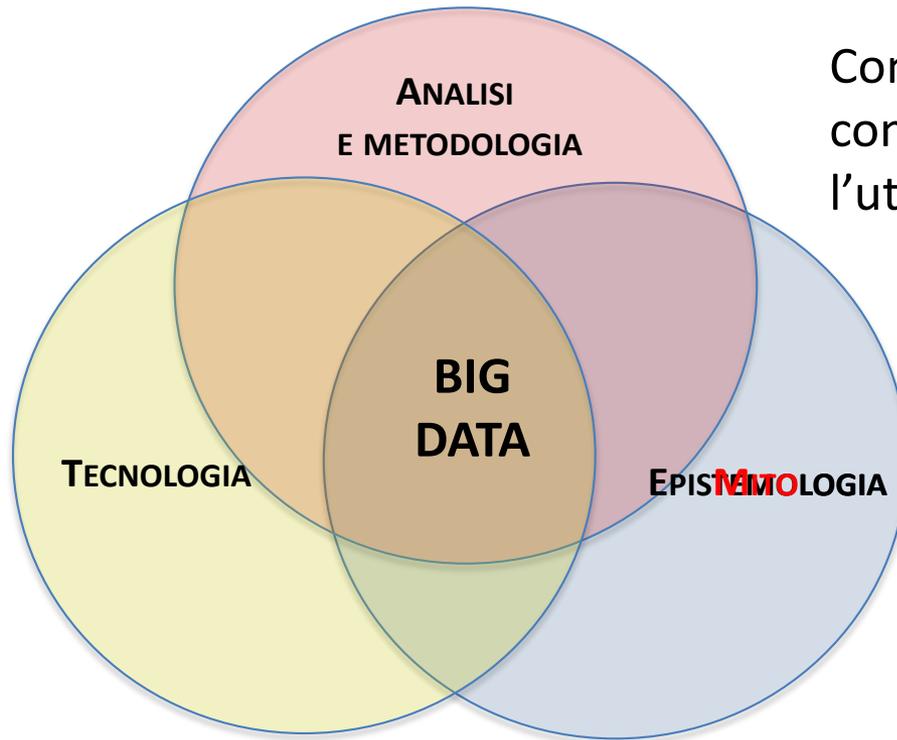
28 aprile 2017

Big Data, Big Challenges.
Convegno di metà mandato della Sezione AIS di Metodologia

StatCities

L'ordito e la trama.
I fili della statistica:
un tessuto
per il governo
delle città.

PRATO
13 • 14 OTTOBRE 2016
MUSEO DEL TESSUTO
Via Puccetti, 3



Come metodi e condizioni della conoscenza scientifica cambino con l'utilizzo di nuove fonti informative

La diffusa convinzione che i grandi set di dati offrono una forma più alta di Intelligenza e conoscenza che possono generare intuizioni in precedenza impossibili, con l'aura della verità, dell'obiettività e dell'accuratezza
[BOYD, CRAWFORD]

<https://www.danah.org/papers/2012/BigData-ICS-Draft.pdf>

La tematica dei Big Data nella ricerca sociale

«I dati non richiedono sforzi specifici per essere raccolti, essendo il sottoprodotto digitale di operazioni di routine svolte entro il sistema»

- Gli algoritmi come principi ordinatori dei più diversi ambiti della vita sociale
- I dati come scatola nera da aprire

- Trasformare molti aspetti della vita delle persone in dati digitali
- Trasformare queste informazioni in nuove forme di valore

IL POTERE DEI GRANDI NUMERI E
LA GOVERNANCE BY NUMBERS

IL PROBLEMA DELLA
DISINTERMEDIAZIONE

**Big Data nella
ricerca sociale**

LA DATIFICAZIONE
DEL MONDO

Mise en données du monde

RAW DATA VS.
COOKED DATA

CAUSALITÀ VS.
CORRELAZIONE

«Raw data is both an oxymoron and a bad idea; to the contrary, data should be cooked with care»

«Correlation doesn't mean causation»

La tematica dei Big Data nella ricerca sociale

Use of more than one method of data collection or research.

Mixed methods research is more specific in that it includes the mixing of qualitative and quantitative data, methods, methodologies, and/or paradigms.

MULTIMETHODOLOGY E
MIXED METHODS

GROUND
THEORY

Osservazione ed elaborazione teorica procedono di pari passo, in un'interazione continua. Il ricercatore scopre la teoria nel corso della ricerca empirica

**Big Data nella
ricerca sociale**

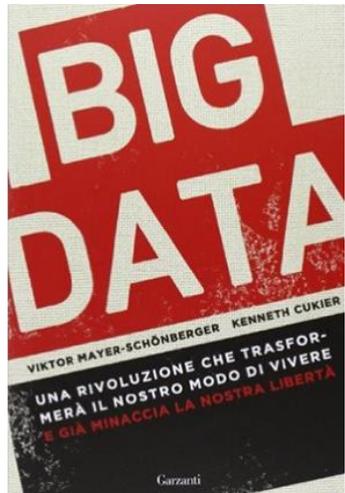
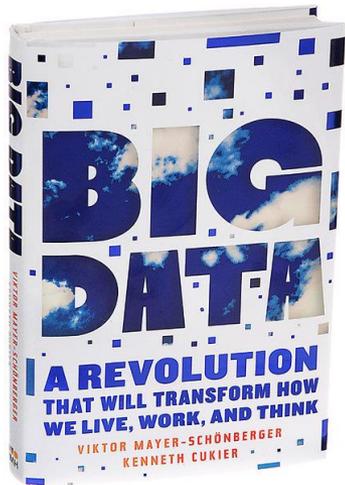
PARADIGMA
INTERPRETATIVO

SCIENCE AND
TECHNOLOGY STUDIES

La realtà sociale non può essere semplicemente osservata ma necessita di interpretazione. Comprendere significa cogliere l'intenzionalità dell'agire umano, attraverso il senso soggettivo attribuito dall'individuo al proprio comportamento.

- Relazioni fra innovazioni scientifiche e tecnologiche, partendo dal presupposto che ambedue sono socialmente costruite e che la società è essa stessa un aggregato sociotecnico.
- Analisi degli effetti, rischi, ridefinizione dei parametri sociali

Le due facce dei Big Data

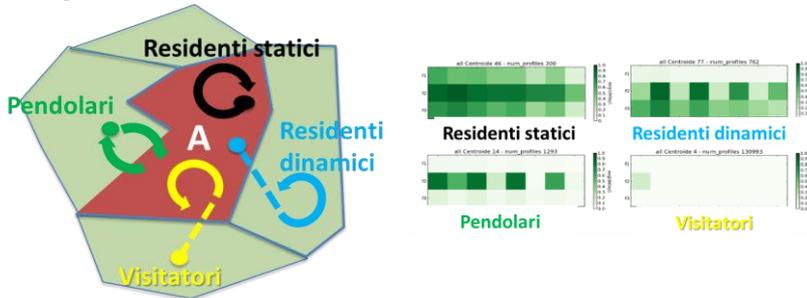


- Le persone come “somma” delle relazioni sociali, delle interazioni online e delle connessioni con i *contenuti* delle azioni che svolgono
- La conoscenza di una persona passa attraverso la penombra (il più larga possibile) di dati che la circonda
- **BIG BROTHER:**
La privacy diventa più difficile da gestire
- **MINORITY REPORT:**
“[...] predictions seem so accurate that people can be arrested for crimes before they are committed”

Big Data e territorio

Persons & Places

Popolazione che insiste su un territorio



Utilizzo di dati GSM e applicazione di modelli distinguere tra residenti e pendolari dinamici (non possibile con i soli dati amministrativi)

Disegno di nuovi «territori» in base alla mobilità della popolazione

Analisi degli spostamenti tramite GPS per delimitare le aree prevalenti entro cui si svolgono le attività. La densità di traffico permette di costruire confini geografici



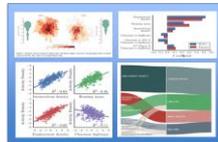
[Rinzivillo et al. KI-Künstliche Intelligenz, 26 (2012)]

Fonti nuove e tradizionali per analizzare le condizioni della vita urbana sulla base del grado di *vita pedonale*

Death VS. Life (Jacobs, 1961)

La struttura urbana «spiega» la vita urbana (77%)

- 1) *Mixed land uses*
- 2) *Small blocks*
- 3) Diversificazione edilizia
- 4) Concentrazione equilibrata di persone ed edifici

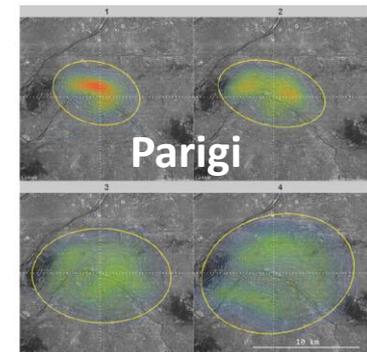


[Fondazione Bruno Kessler, University of Trento, Sorbonne Universités, Telecom Italia, Bell Labs Cambridge]

Dati di sensori di telefonia mobile per monitorare l'attività umana nelle città

Nuovi indicatori riferiti alla vita nelle città ➔ **Tempo sociale** delle attività
Differenze spaziali per misurare il «giorno attivo»

Mappe di densità calcolate per quartili di durata della vita sociale (approssimata dalla *vitalità* delle celle di telefonia mobile)



1. Demistificare i Big Data
2. Il contesto dei Big Data nella statistica ufficiale
3. Altri punti di vista
- 4. Problemi aperti**
5. Esperienze correnti in Istat

Six Provocations for Big Data

1. AUTOMATING
RESEARCH CHANGES
THE DEFINITION OF
KNOWLEDGE

6. LIMITED ACCESS TO BIG DATA
CREATES NEW DIGITAL DIVIDES

2. CLAIMS TO OBJECTIVITY AND
ACCURACY ARE MISLEADING

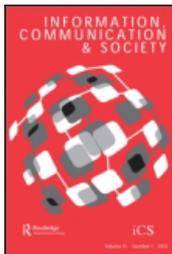
5. JUST BECAUSE IT IS ACCESSIBLE
DOESN'T MAKE IT ETHICAL

Six provocations for Big Data

— Danah Boyd, Kate Crawford —
Samira Shaikh, Veena Ravishankar

3. BIGGER DATA ARE NOT
ALWAYS BETTER DATA

4. NOT ALL DATA ARE
EQUIVALENT



Big data is at its best when analyzing things that are extremely common, but often falls short when analyzing things that are less common

Although big data is very good at detecting correlations, especially subtle correlations that an analysis of smaller data sets might miss, it never tells us which correlations are meaningful

Big data can work well as an adjunct to scientific inquiry but rarely succeeds as a wholesale replacement

UTILI SOLO QUANDO IL PROBLEMA È SEMPLICE

COMPRESIONE CAUSE

COMPLEMENTO /SOSTITUZIONE

CRITICITÀ SU DOMANDE IMPRECISE

The New York Times

MANIPOLAZIONE DELLE TECNICHE

Many tools that are based on big data can be easily gamed

Big data is prone to giving scientific-sounding solutions to hopelessly imprecise questions.

TROPPE CORRELAZIONI

POCA ROBUSTEZZA

EFFETTO ECO-CAMERA

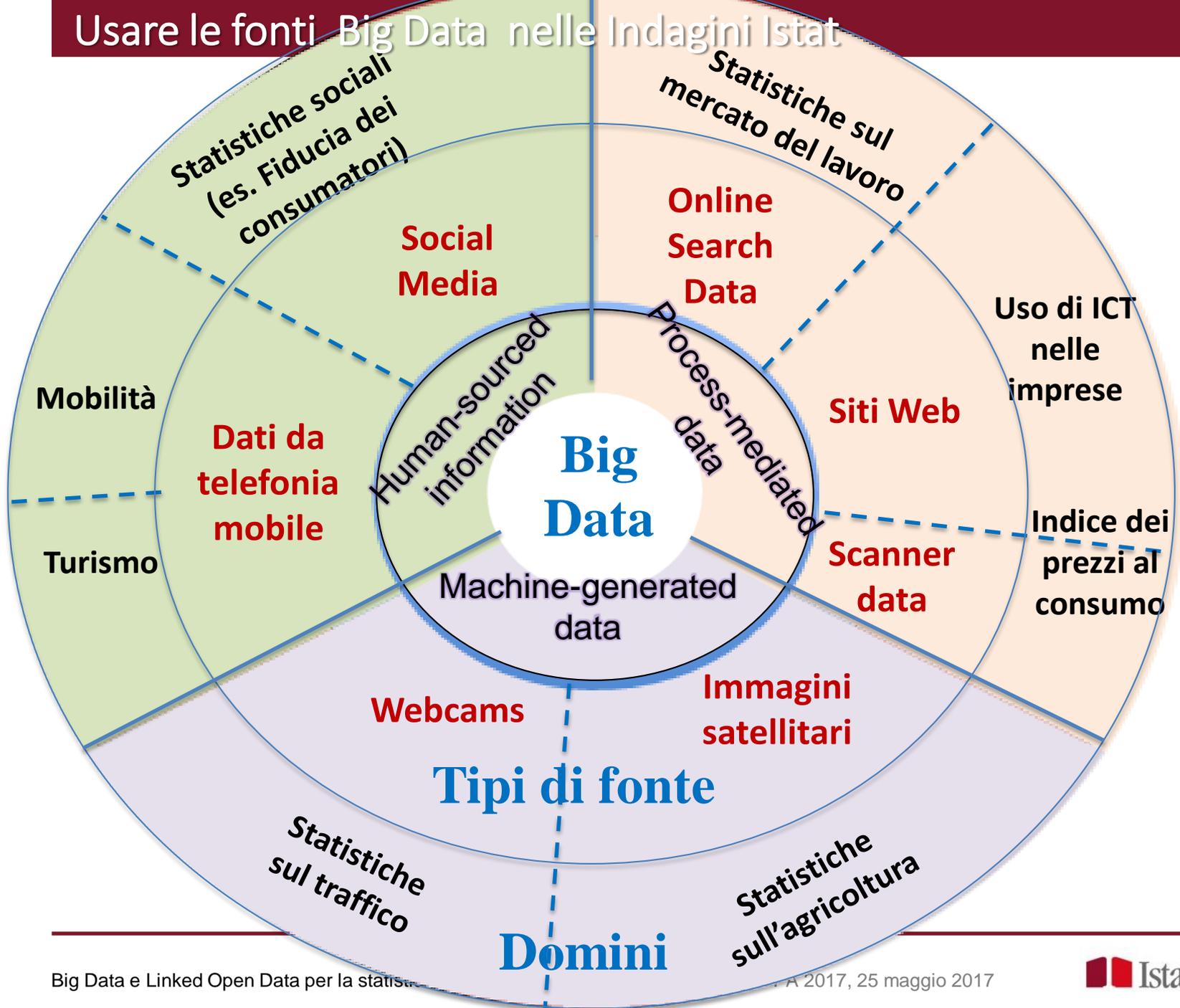
Absent careful supervision, the magnitudes of big data can greatly amplify such errors.

Whenever the source of information for a big data analysis is itself a product of big data, opportunities for vicious cycles abound

Even when the results of a big data analysis aren't intentionally gamed, they often turn out to be less robust than they initially seem

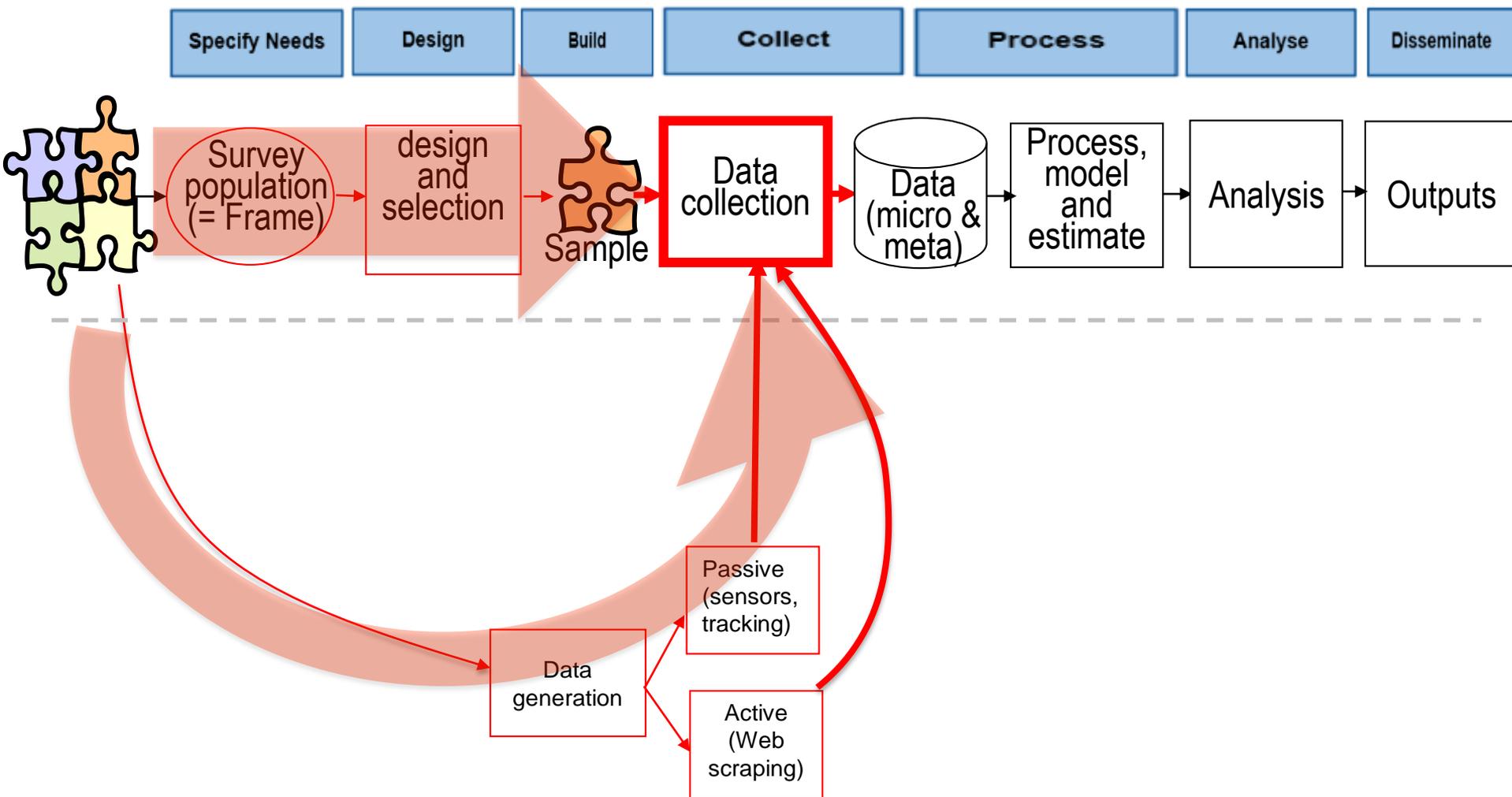
1. Demistificando i Big Data
2. Il contesto dei Big Data nella statistica ufficiale
3. Altri punti di vista
4. Problemi aperti
- 5. Esperienze correnti in Istat**

Usare le fonti Big Data nelle Indagini Istat

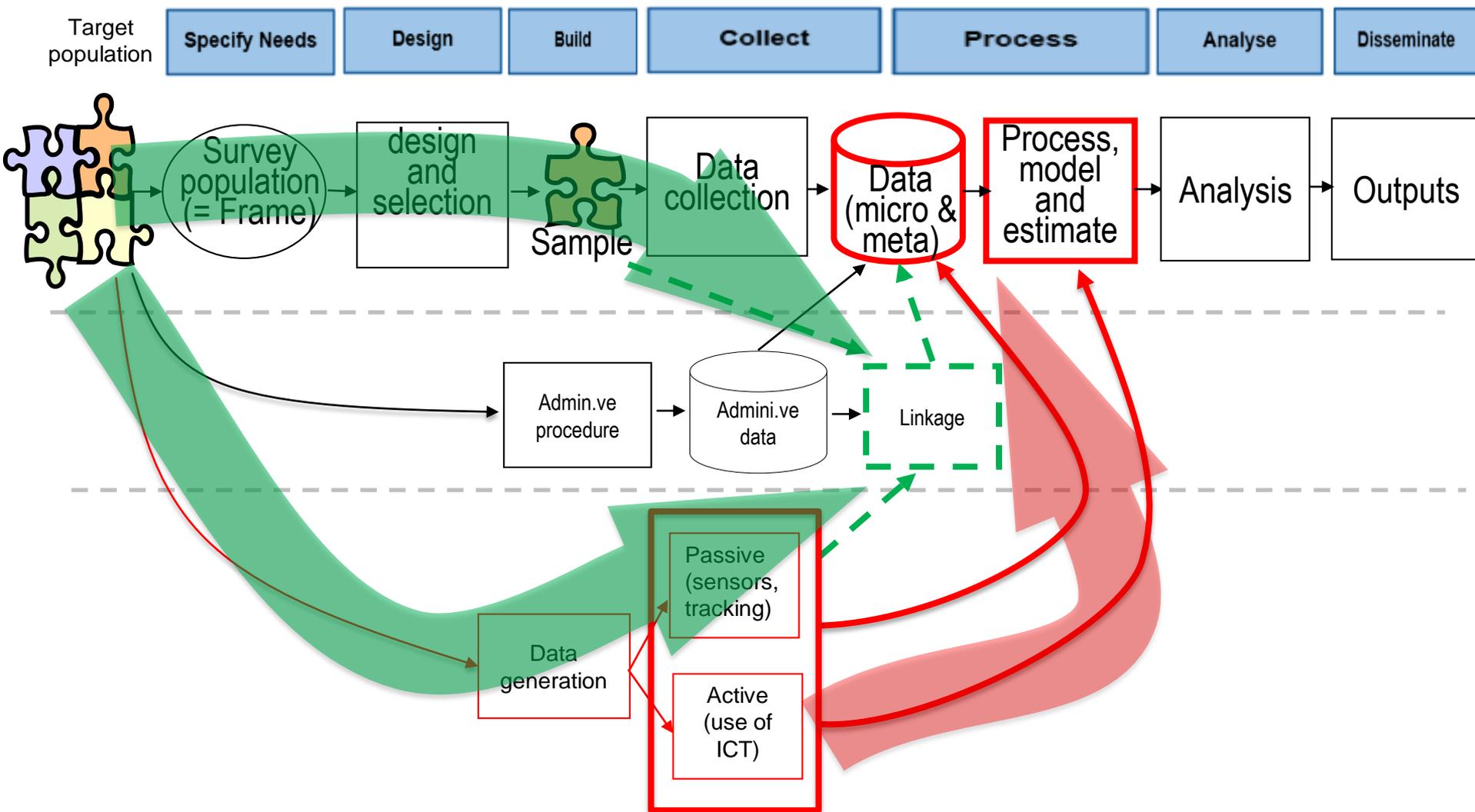


- Scenario 1: **Uso delle fonti Big limitato alla raccolta dati**
Vengono raccolte le stesse informazioni utilizzate nel processo statistico tradizionale, ma accedendo direttamente alla fonte Big e senza interventi significativi su approcci/ tecniche di analisi
- Scenario 2: **Uso delle fonti Big in combinazione o integrato con le altre fonti di dati (da indagine e amministrativi)**
Per le stime si utilizzano sia i dati da fonte tradizionale sia Big Data, dopo opportuno e specifico passo di integrazione (RL) e/o trattamento (NLP, Text Mining, ML, ecc.)
- Scenario 3: **Uso delle fonti Big in sostituzione delle (o alternativo rispetto alle) fonti tradizionali**
Per le stime si usano solo Big Data e non dati di indagine, con individuazione di specifiche tecniche e nuovi metodi di analisi/trattamento lungo tutto il processo statistico (da RD in poi)

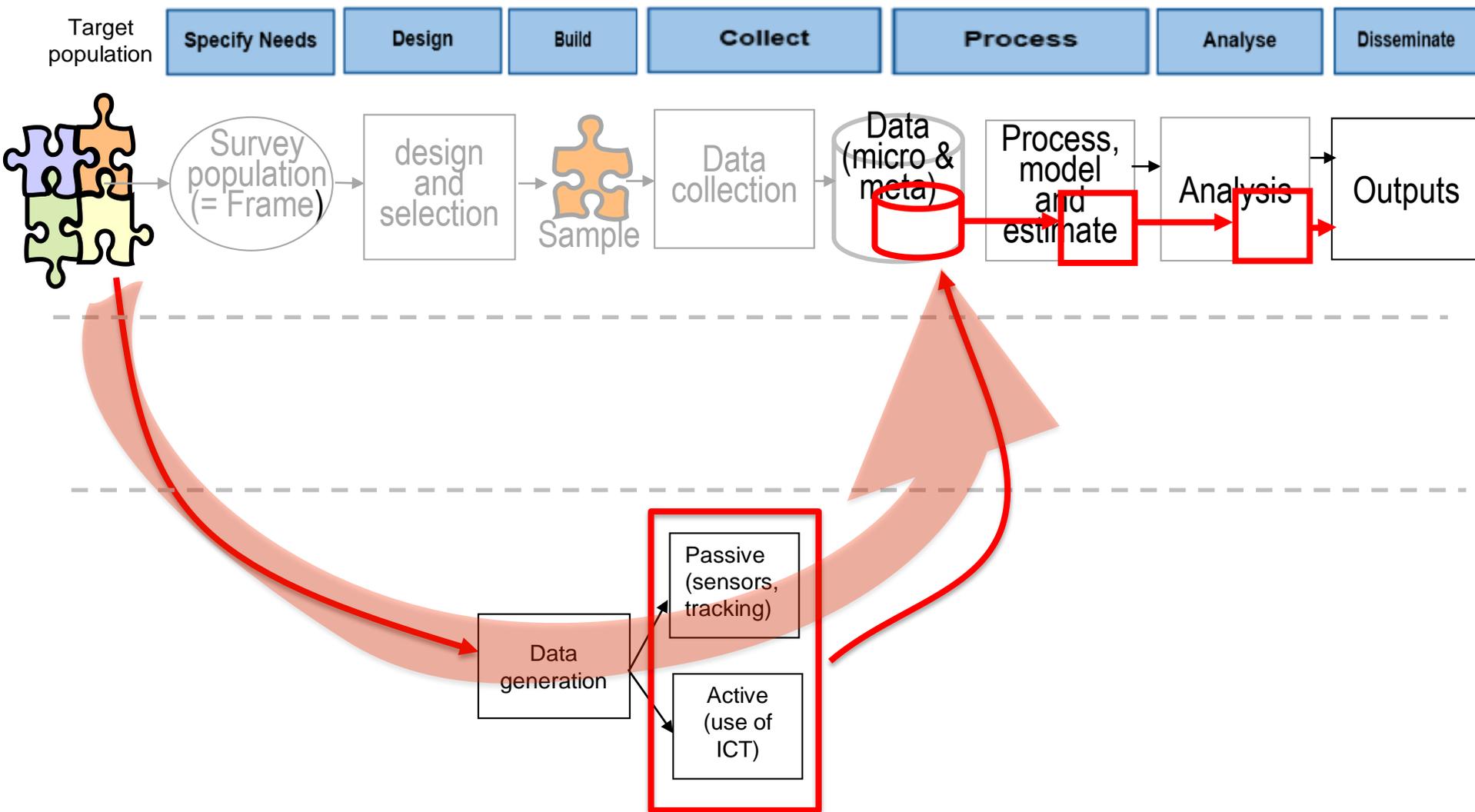
Scenario 1: tecniche alternative di data collection



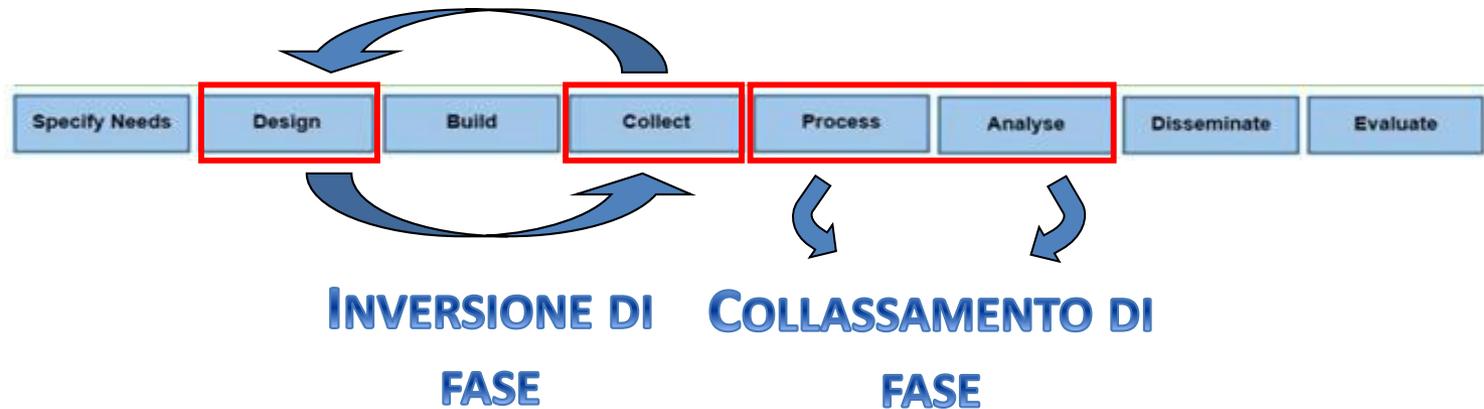
Scenario 2: uso integrato



Scenario 3: uso in sostituzione delle fonti tradizionali

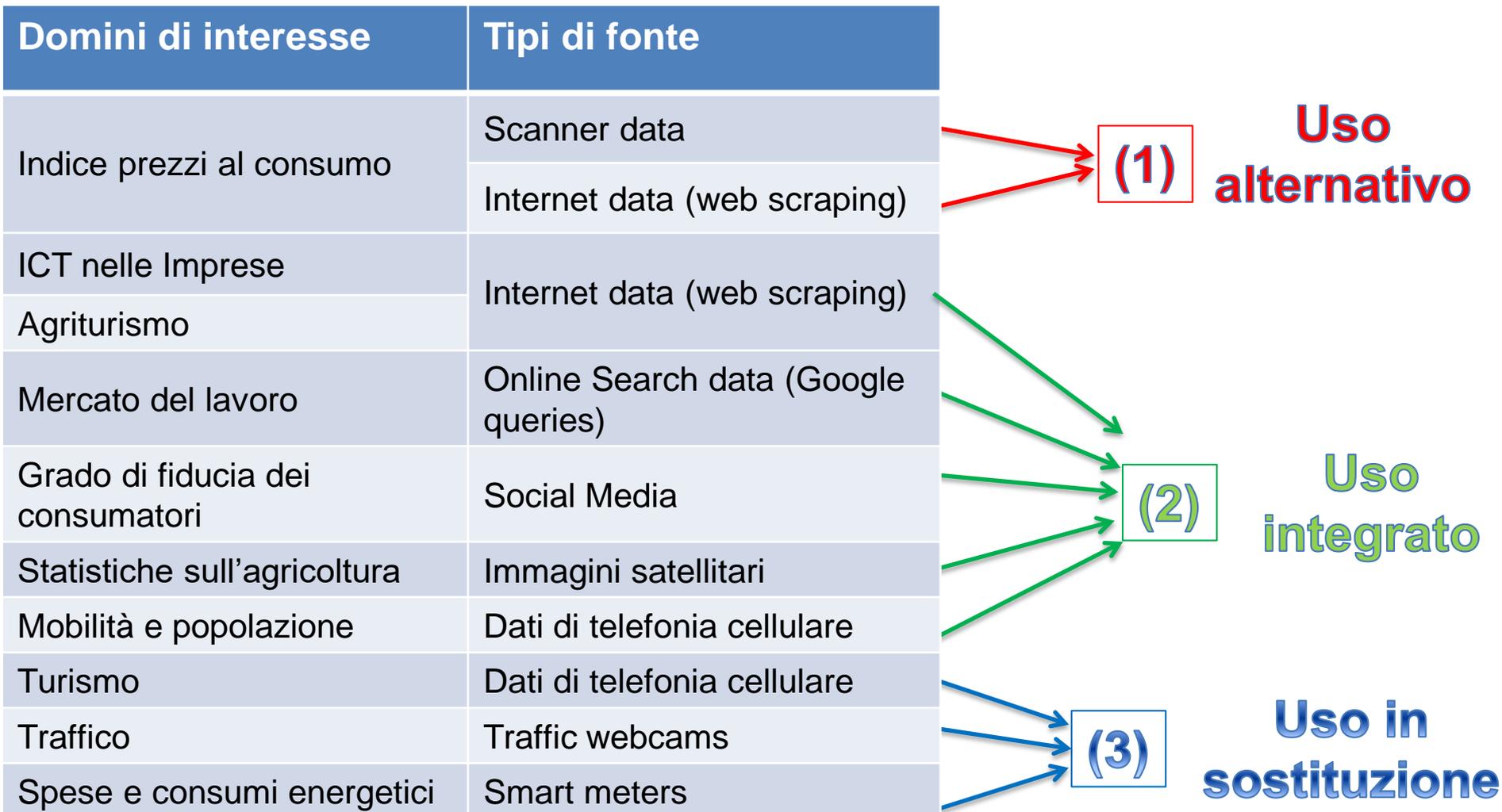


Impatto dei Big Data nelle fasi del processo statistico

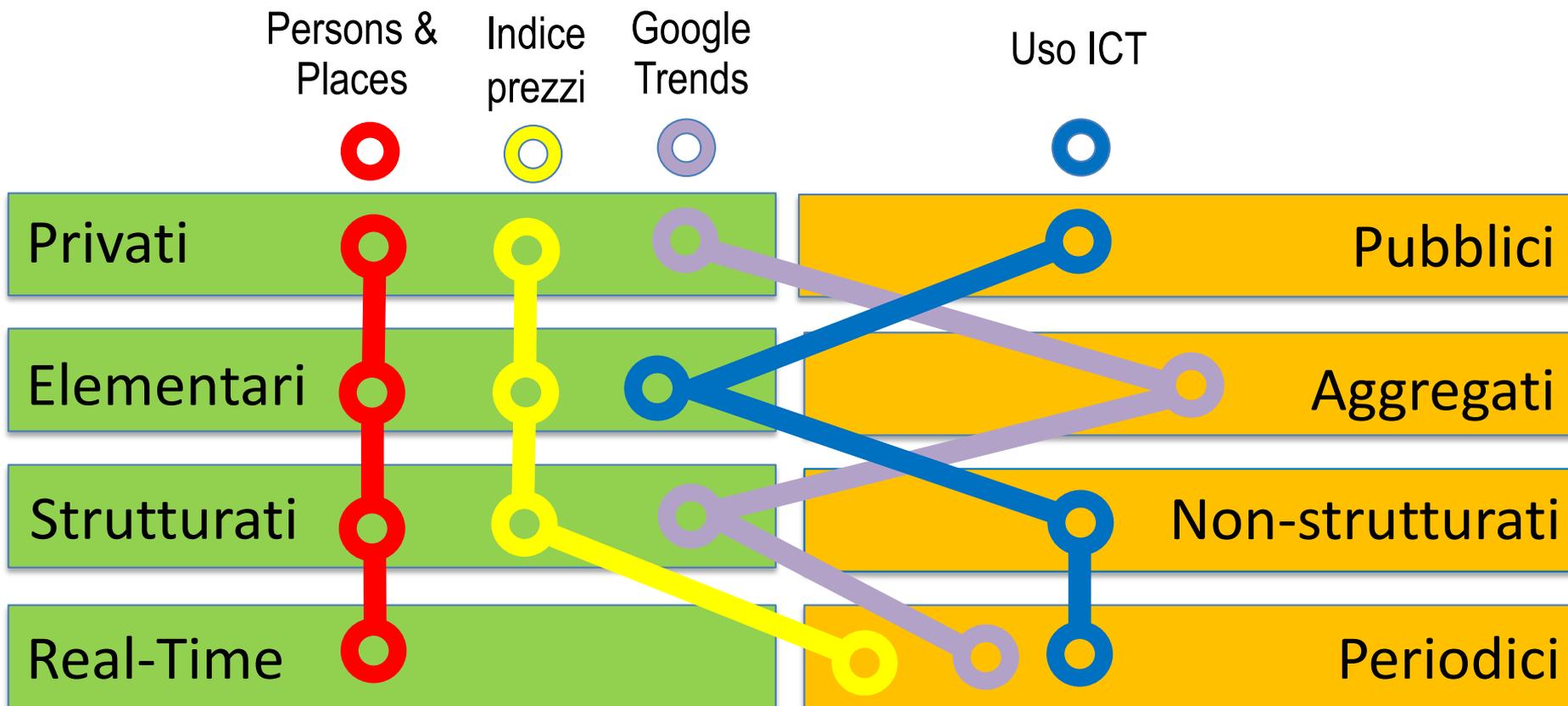


- Possibile inversione di alcune fasi (*Design* e *Collect*)
- La fase di *collezione dati* può a volte essere sostituita da quella di *generazione dati*
- Possibile collapsamento delle fasi di *Process* e *Analyse* (possono avvalersi degli stessi metodi)
- Altre fasi (ad es. *Dissemination*) non sono ancora coinvolte

Big Data: possibili scenari e applicazioni



Caratteristiche dei Big Data utilizzati nelle sperimentazioni



Internet access

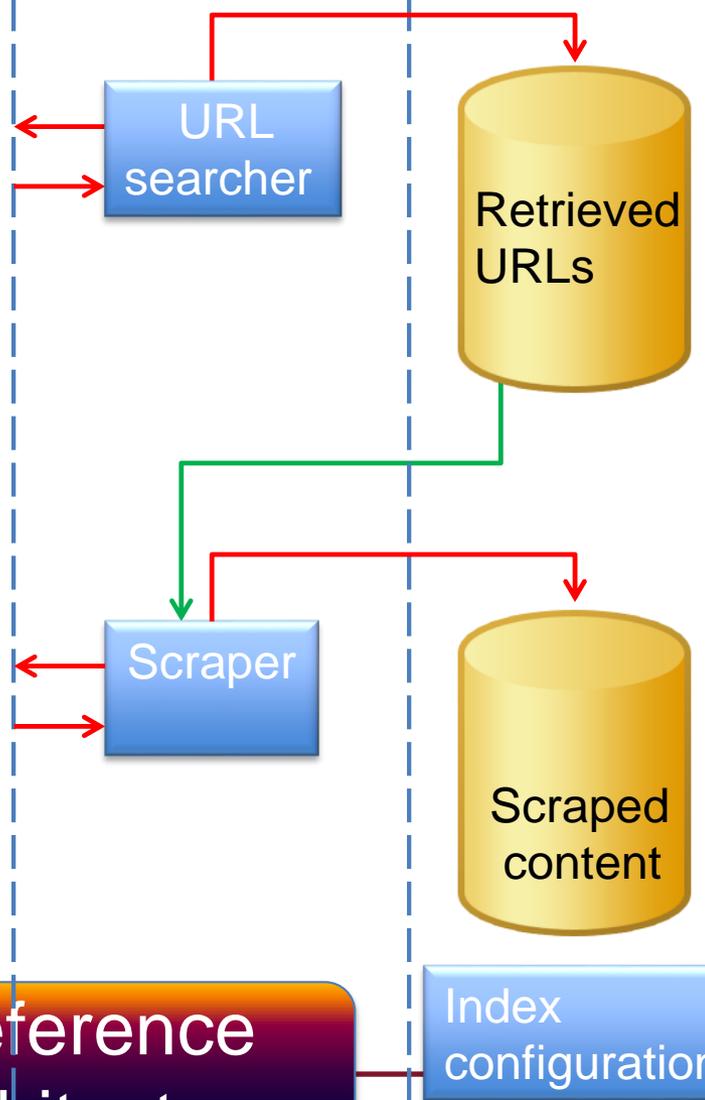
Storage

Data preparation

Analysis



W
E
B



Feature extraction

URL scorer

Tokenization

Data Parsing

Word filters (eg. stopwords)

Language specific lemmatization

Term document matrix generation

Machine Learning

Build training & test sets

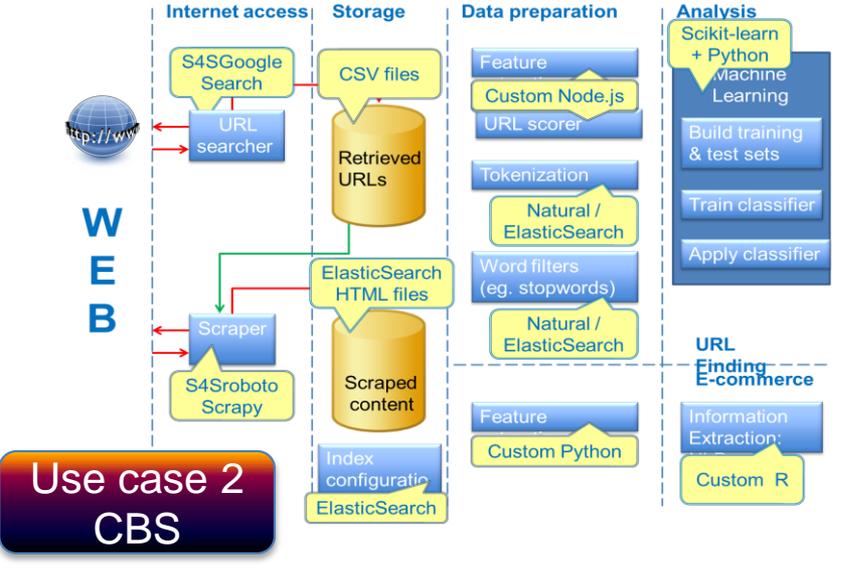
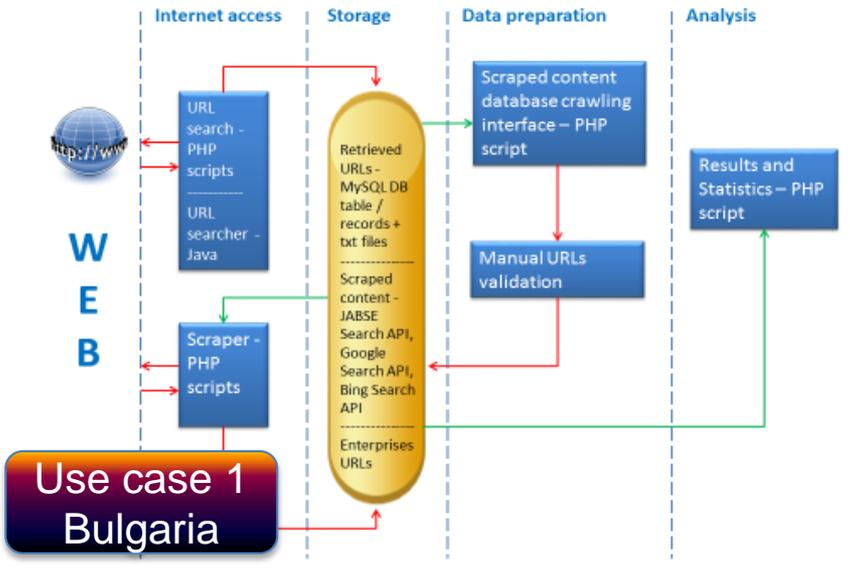
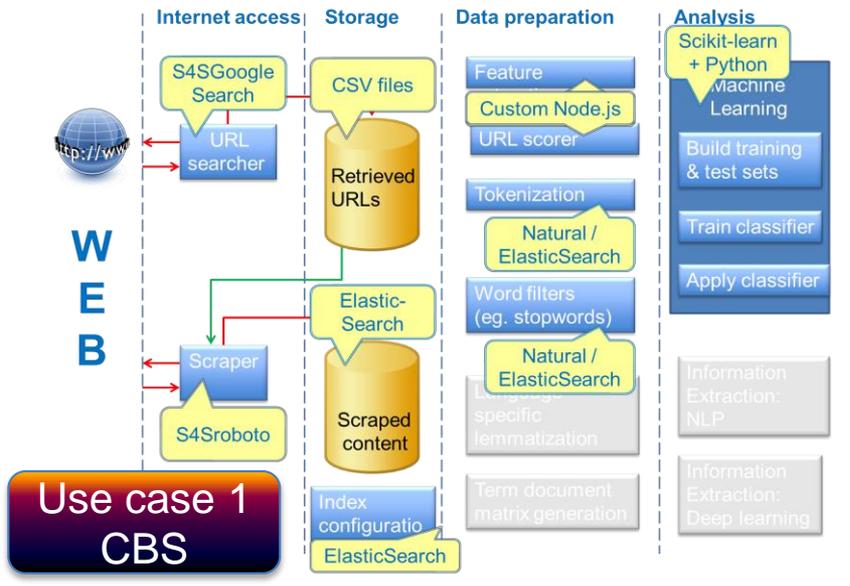
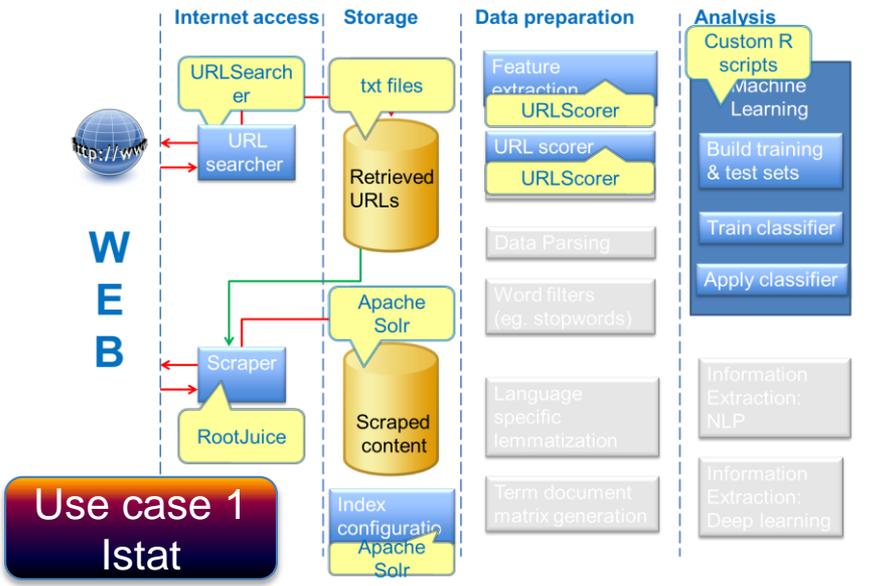
Train classifier

Apply classifier

Information Extraction: NLP

Information Extraction: Deep learning

Reference Architecture





*“There is no reason
anyone would want a
computer in their
home.” -*

*Ken Olson, president, chairman
and founder of Digital
Equipment Corp. (DEC), maker
of big business mainframe
computers, arguing against the
PC in 1977*

FINE

Grazie

stefano.defrancisci@istat.it