

# *Data science per lo studio di ecosistemi complessi*



**Edoardo Pasoli**

*Dipartimento di Agraria*

*Università degli Studi di Napoli Federico II*



Caffè scientifico di Agraria – Portici, Sala Cinese – 26 giugno 2019

[edoardo.pasoli@unina.it](mailto:edoardo.pasoli@unina.it)

 [@epasoli](https://twitter.com/epasoli)



# Agenda

- Biografia in breve
- Introduzione alla *data science*
- Le mie linee di ricerca



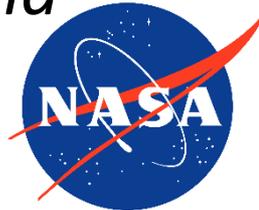
---

edoardo.pasolli@unina.it

 @epasolli

# Biografia in breve

- **Laurea in Ingegneria delle Telecomunicazioni @ Univ. di Trento, 2008**
- **Dottorato in *Information and Communication Technology* @ Univ. di Trento, 2011**  
Tesi: "*Active learning methods for classification and regression problems*"
- **Postdoc @ NASA Goddard Space Flight Center, 2012-13**
- **Postdoc @ Purdue University, 2013-14**



edoardo.pasolli@unina.it

 @epasolli

# Biografia in breve

- **Assegnista di ricerca & Marie Skłodowska-Curie Individual Fellow @ Univ. di Trento, 2014-18**



- **RTD-B @ Agraria, Da Dicembre 2018**  
S.S.D.: ING-INF/03 (Telecomunicazioni)

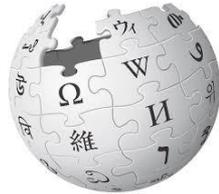


edoardo.pasolli@unina.it

 @epasolli

# Introduzione alla *data science*

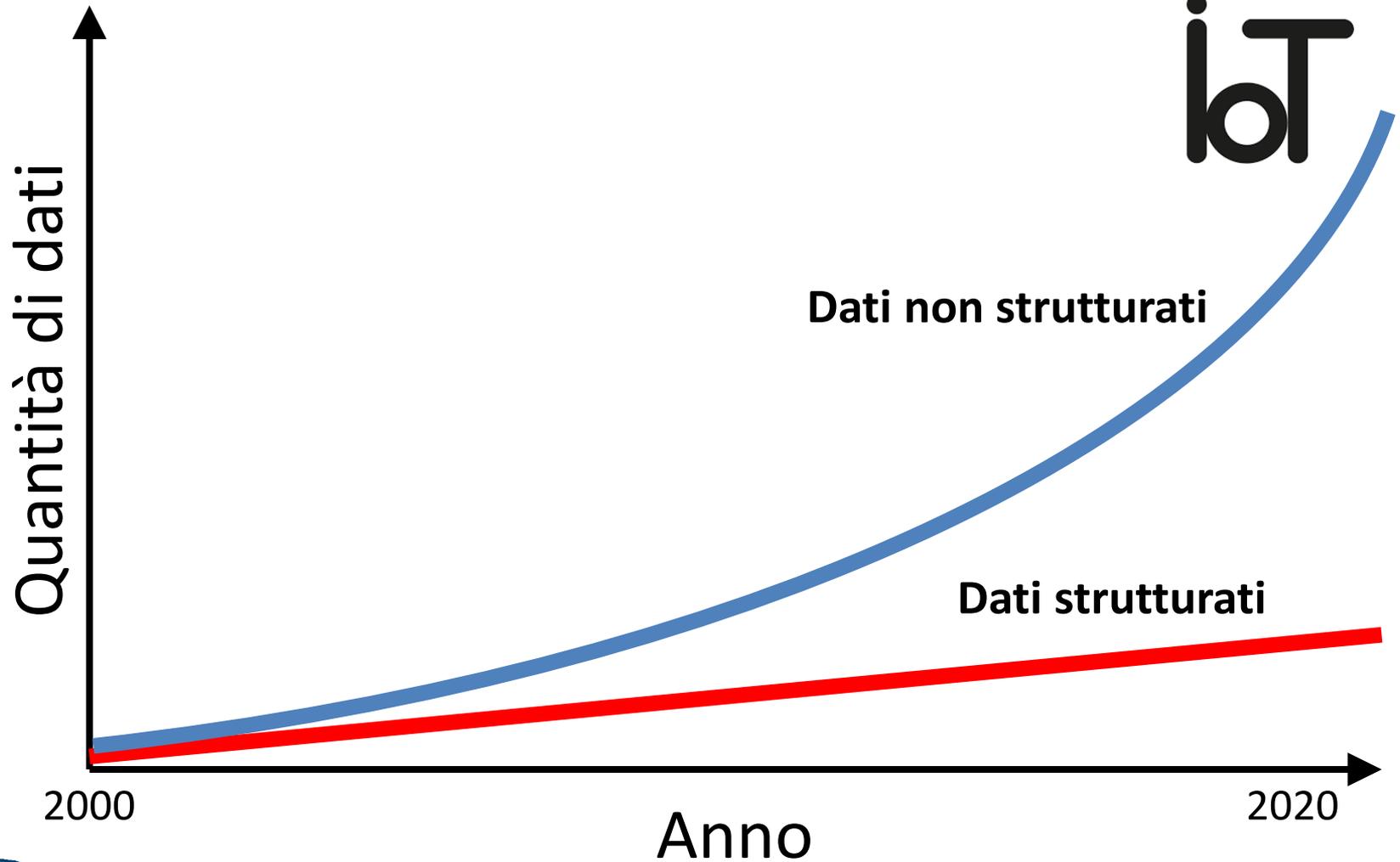
- Una possibile definizione: "*Data science is an **interdisciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from data in various forms***"



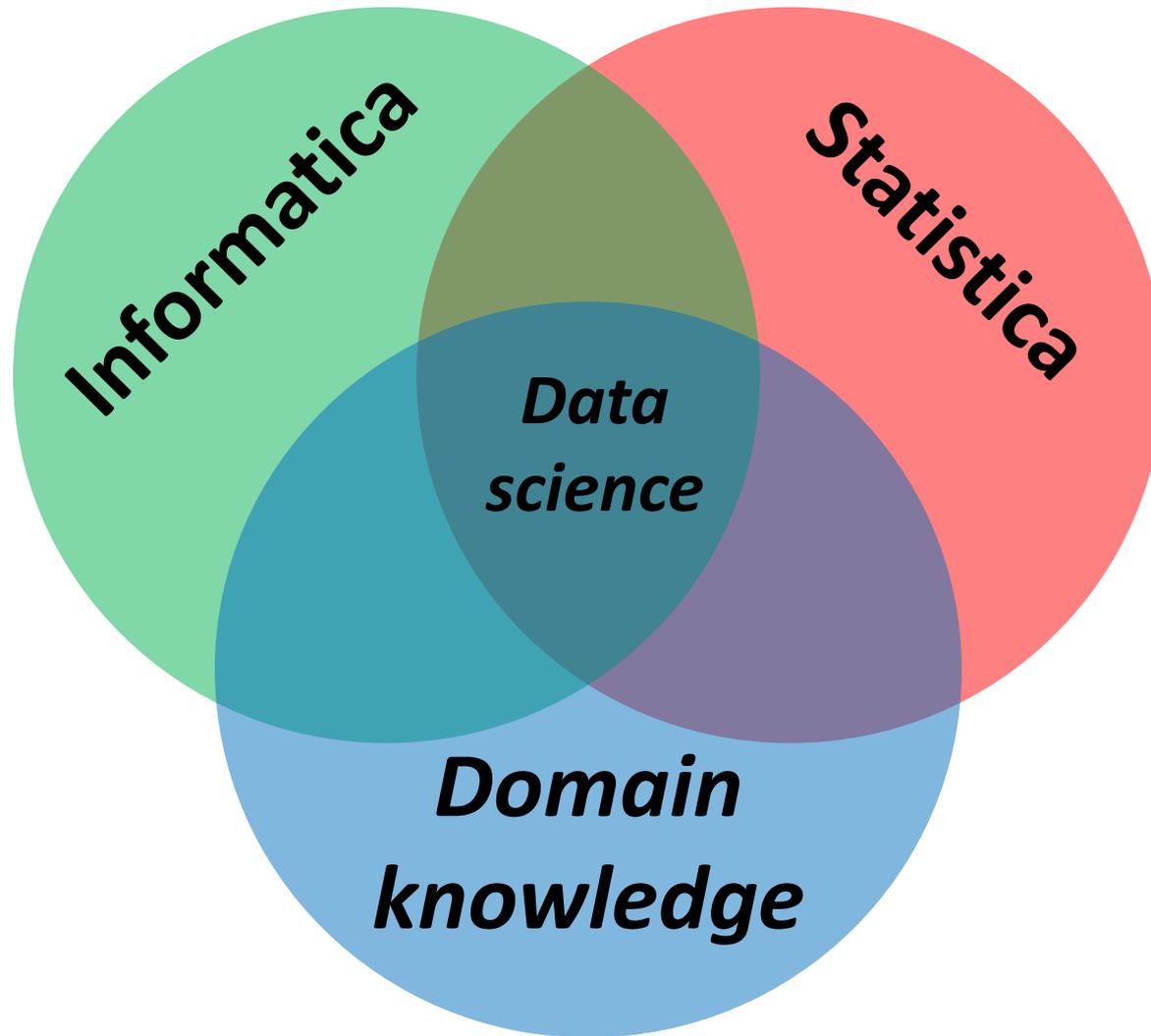
- In breve: "*The processes used to extract insights from data*"



# Perchè adesso?



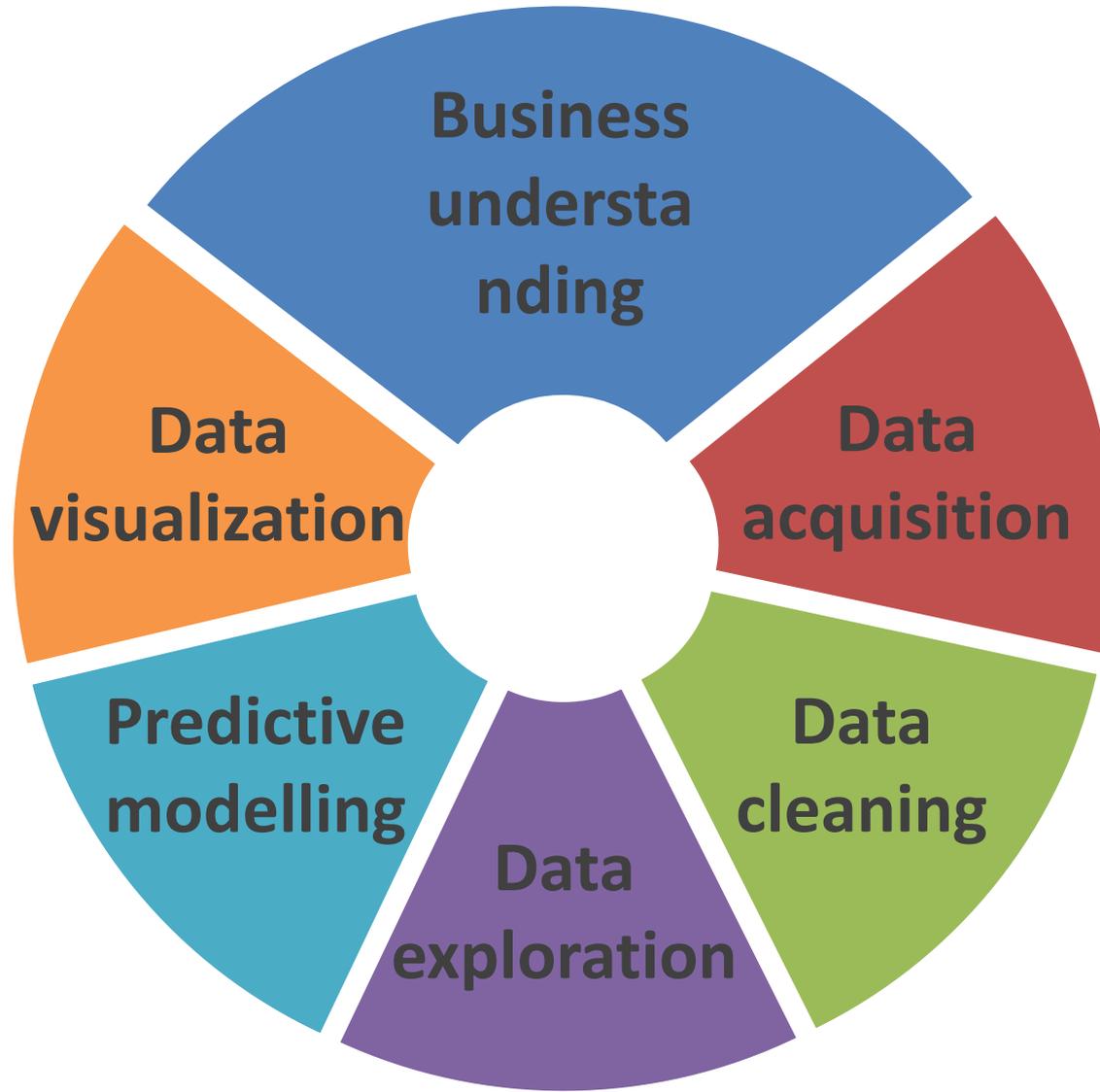
# Un ambito fortemente multidisciplinare



[edoardo.pasolli@unina.it](mailto:edoardo.pasolli@unina.it)

 @epasolli

# *Il data science lifecycle*



edoardo.pasolli@unina.it

 @epasolli



# *Garbage in = Garbage out*



**Dati**



**Analisi**



**Risultato**



**Dati**



**Analisi**



**Risultato**



# Le mie linee di ricerca

- Sviluppo di **algoritmi di *machine learning*** per **dati telerilevati**
  - Metodologie di *active learning*
  - Fusione di dati multi-sorgente
- Sviluppo ed applicazione di ***tool* computazionali** per il **microbioma da dati di metagenomica**
  - *Tools* per l'analisi a livello di ceppo
  - Microbioma umano... e da alimenti



# Introduzione al telerilevamento

- Definizione: *“Remote sensing is the acquisition of information about an object or phenomenon without making physical contact with it”*



# Classificazione & Mappe di copertura del suolo



**Input:**  
Dato telerilevato

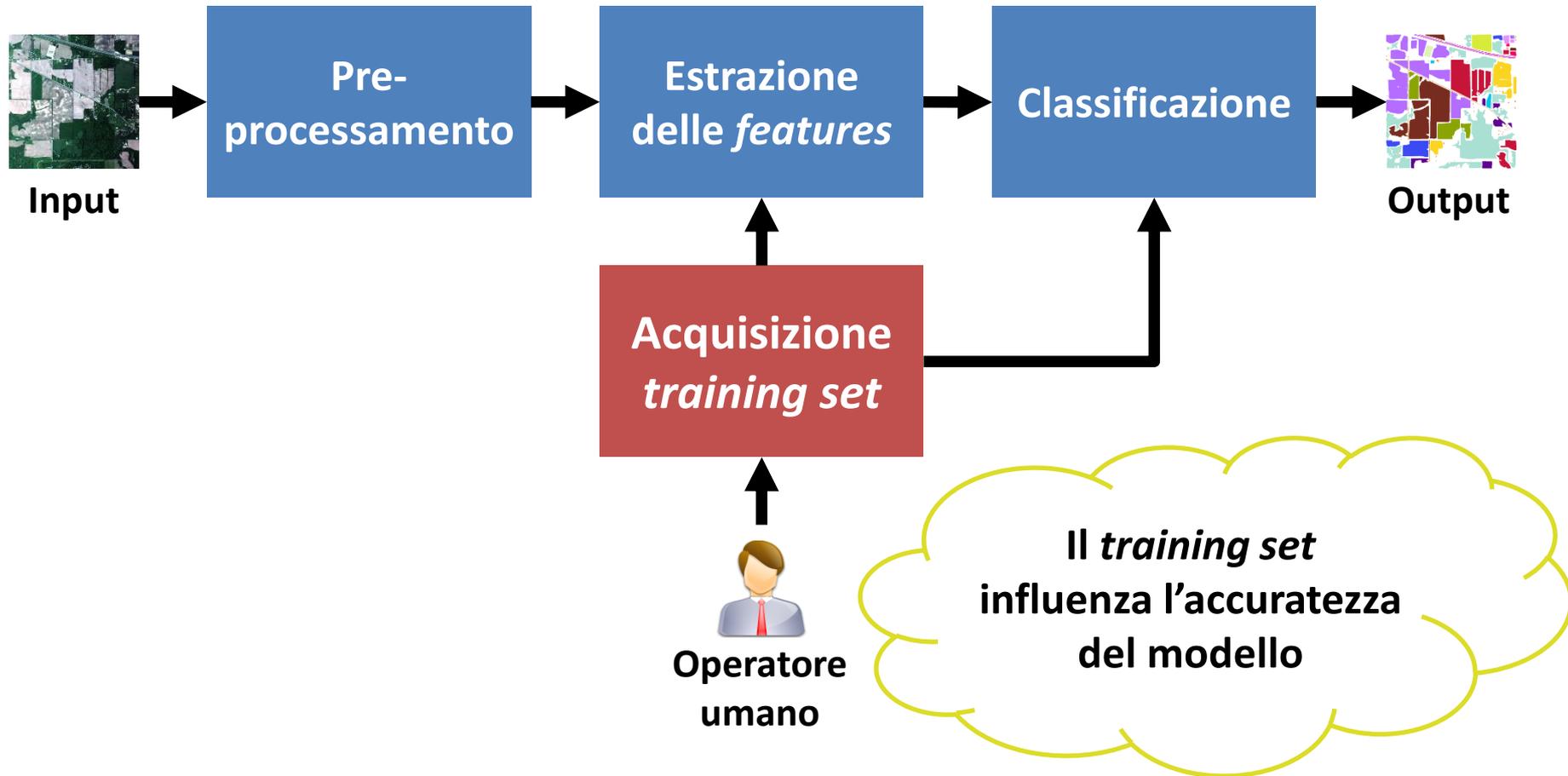


**Output:**  
Mappa di copertura del suolo  
(*Semantic segmentation*)

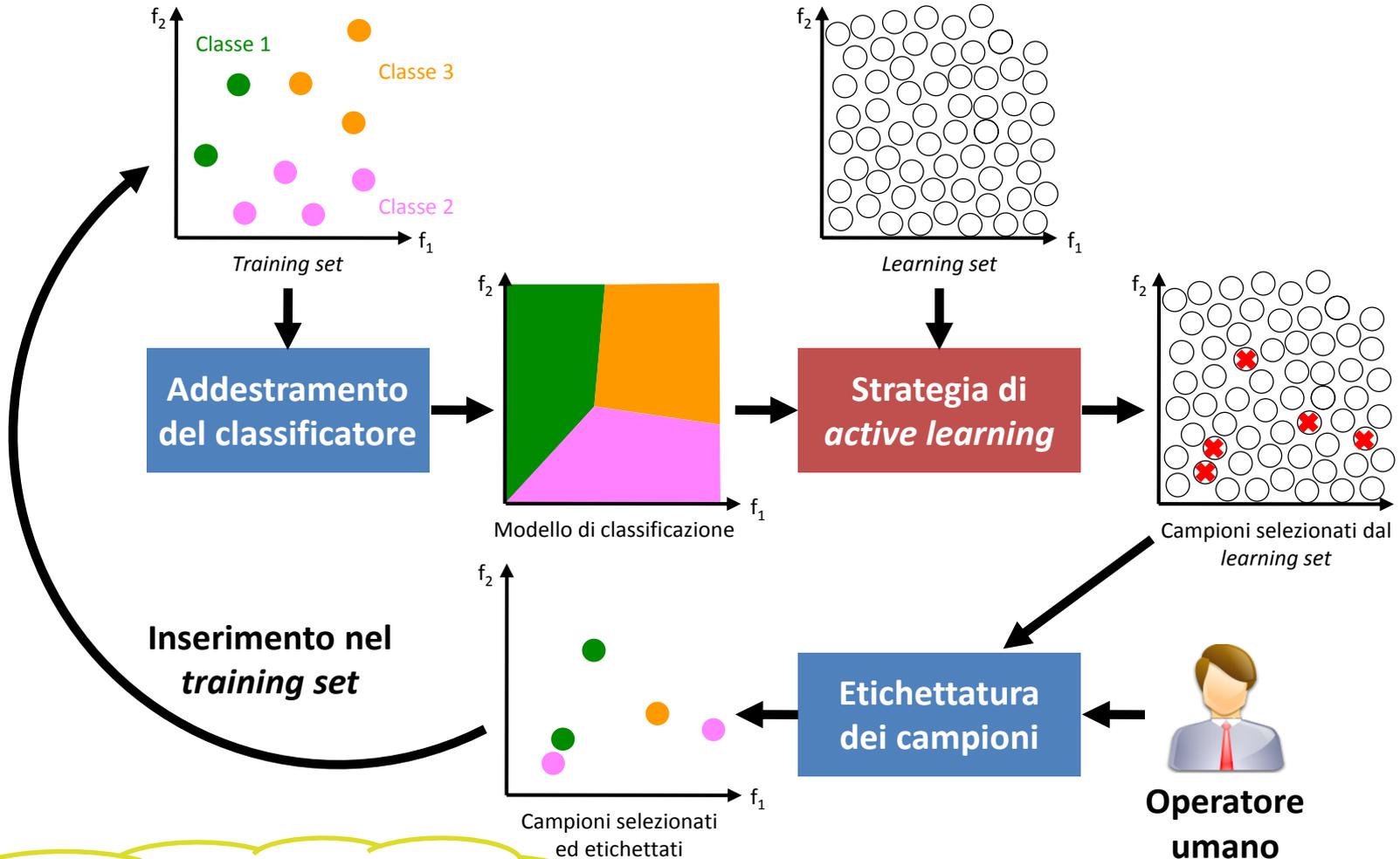
- Lavorato principalmente con sensori ottici + fusione con dati LiDAR



# Approccio di classificazione supervisionata



# Active learning & Acquisizione training set



**Obiettivo: Selezionare il training set migliore**

*Human-in-the-loop*



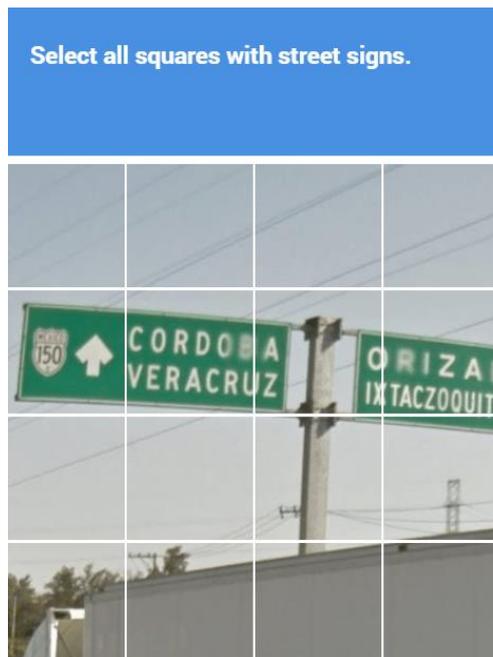
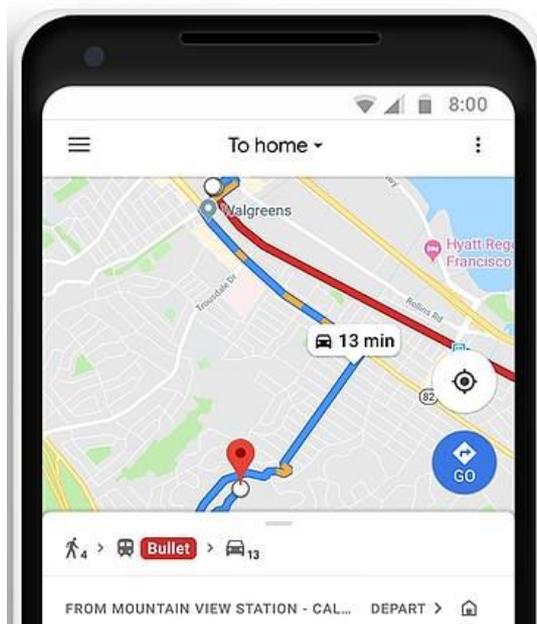
# Human-in-the-loop, qualche esempio



Consigliati in base alla tua Lista dei Desideri [Visualizza altro](#)



amazon.it



edoardo.pasoli@unina.it

 @epasoli

# Active learning & Mappe di copertura del suolo

## Remote Sensing of Environment

Using active learning to adapt remote sensing image classifiers

D. Tuia <sup>a,\*</sup>, E. Pasolli <sup>b</sup>, W.J. Emery <sup>c</sup>

IEEE GEOSCIENCE AND REMOTE SENSING LETTERS, VOL. 8, NO. 3, MAY 2011

### Support Vector Machine Active Learning Through Significance Space Construction

Edoardo Pasolli, *Student Member, IEEE*, Farid Melgani, *Senior Member, IEEE*, and Yakoub Bazi, *Senior Member, IEEE*

IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING, VOL. 50, NO. 10, OCTOBER 2012

### Active Learning Methods for Biophysical Parameter Estimation

Edoardo Pasolli, *Student Member, IEEE*, Farid Melgani, *Senior Member, IEEE*, Naif Alajlan, *Member, IEEE*, and Yakoub Bazi, *Senior Member, IEEE*

IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING, VOL. 51, NO. 6, JUNE 2013

### Optical Image Classification: A Ground-Truth Design Framework

Edoardo Pasolli, *Member, IEEE*, Farid Melgani, *Senior Member, IEEE*, Naif Alajlan, *Member, IEEE*, and Nicola Conci, *Member, IEEE*

IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING, VOL. 52, NO. 4, APRIL 2014

### SVM Active Learning Approach for Image Classification Using Spatial Information

Edoardo Pasolli, *Member, IEEE*, Farid Melgani, *Senior Member, IEEE*, Devis Tuia, *Member, IEEE*, Fabio Pacifici, *Senior Member, IEEE*, and William J. Emery, *Fellow, IEEE*

edoardo.pasolli@unina.it

 @epasolli



# Active learning & Mappe di copertura del suolo

IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING, VOL. 54, NO. 4, APRIL 2016

## Active-Metric Learning for Classification of Remotely Sensed Hyperspectral Images

Edoardo Pasoli, *Member, IEEE*, Hsiuhan Lexie Yang, *Member, IEEE*, and Melba M. Crawford, *Fellow, IEEE*



Dato telerilevato



Verità a terra

**Dataset:** New Indian Pine  
**Data di acquisizione:** Maggio 2010  
**Sensore:** SpecTIR  
**Risoluzione spaziale:** 2 m  
**Risoluzione spettrale:** 5 nm  
**# campioni etichettati:** 1,094,132  
**# classi tematiche:** 19

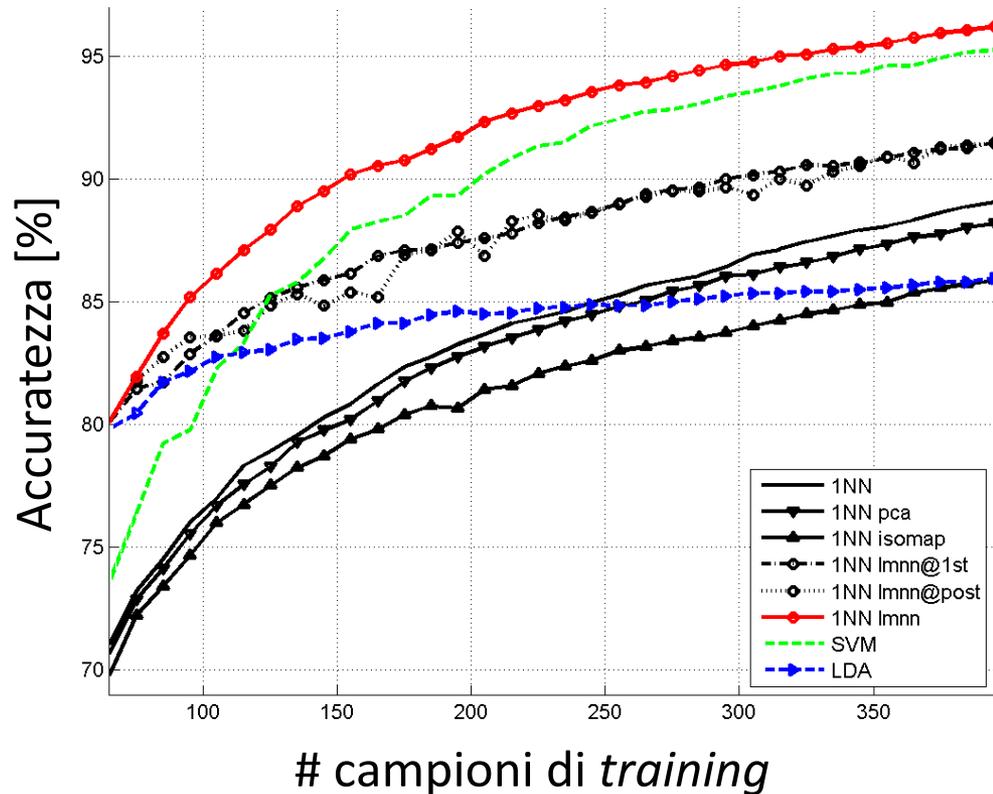


edoardo.pasoli@unina.it

 @epasoli

# Active learning & Mappe di copertura del suolo

## Esempio di risultati di classificazione

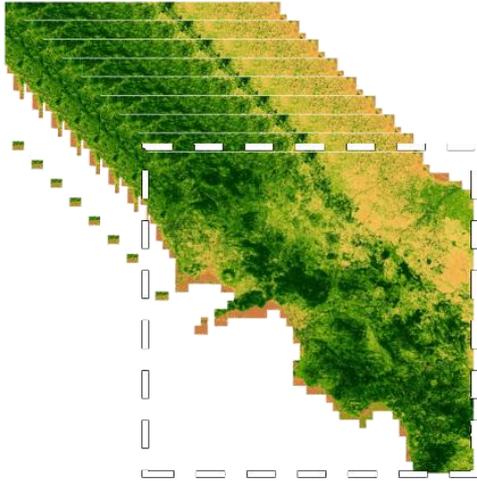


Sviluppi futuri: **Active + deep + transfer learning**  
In collaborazione con prof. Antonia Maria Tulino @ DIETI

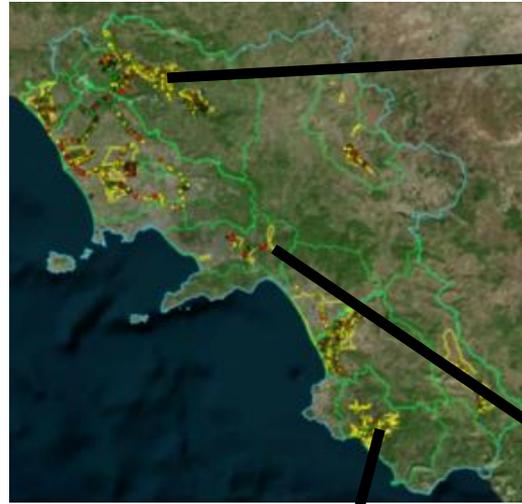


# Identificazione di aree irrigate (in corso)

Con prof. Guido D'Urso



Dato multi temporale  
Landsat-8 & Sentinel-2  
(2018)



**Classe 1:**  
*Bare soil*



**Classe 3:**  
*Tree crop*



**Classe 2:**  
*Herbaceous crop*

*Verità a terra*  
2.992 campioni

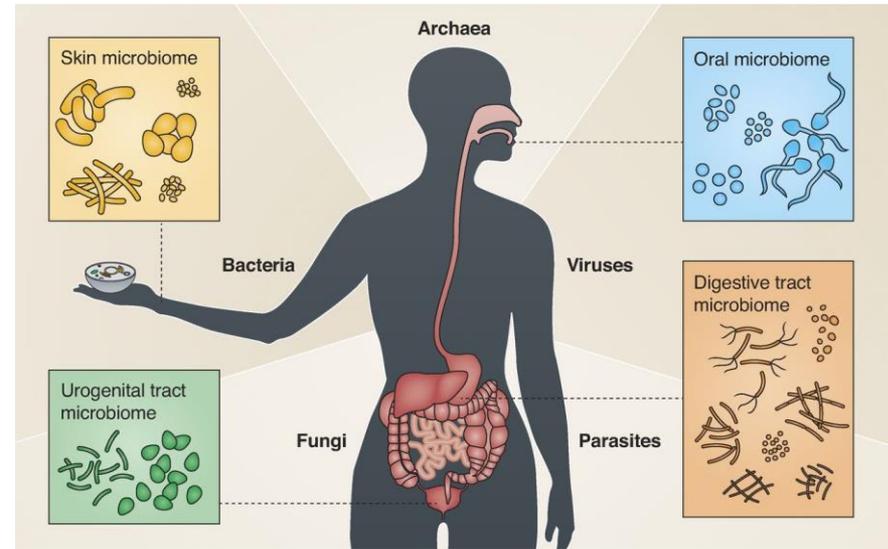
edoardo.pasolli@unina.it

 @epasolli



# Introduzione al microbioma

*“The ecological community of commensal, symbiotic, and pathogenic microorganisms living in the human body”*



Structure, function and diversity of the healthy human microbiome **nature**

**Role of the Microbiota** **Cell**  
**in Immunity and Inflammation**

Potential of fecal microbiota for early-stage detection of **colorectal cancer**

molecular  
systems  
biology

**Richness of human gut microbiome** **nature**  
correlates with **metabolic markers**

A metagenome-wide association study of gut microbiota in **type 2 diabetes** **nature**

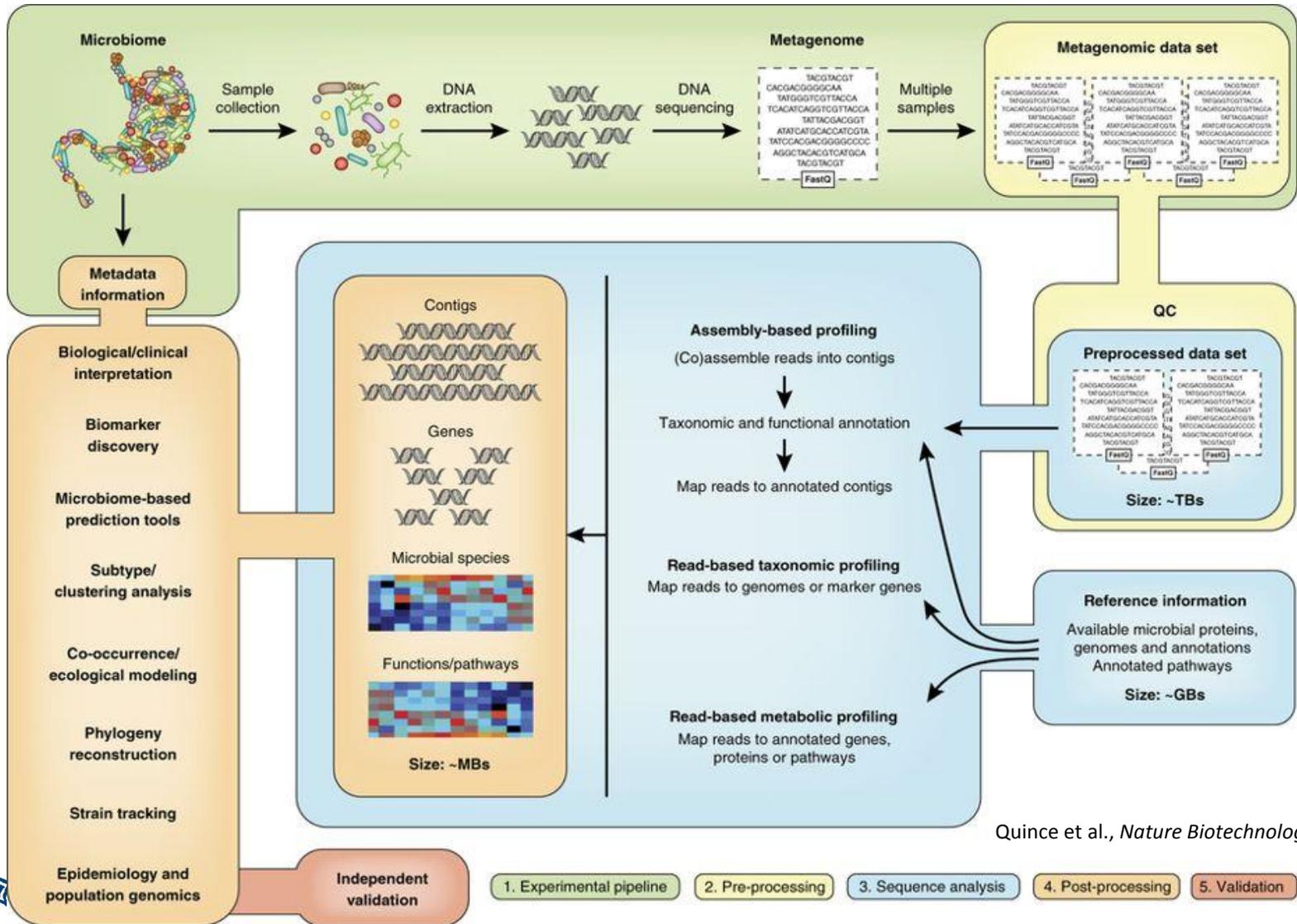
edoardo.pasolli@unina.it

 @epasolli

 **DIPARTIMENTO DI AGRARIA**



# Microbiome & Metagenomics



edoardo.pasoli@unina.it

@epasoli

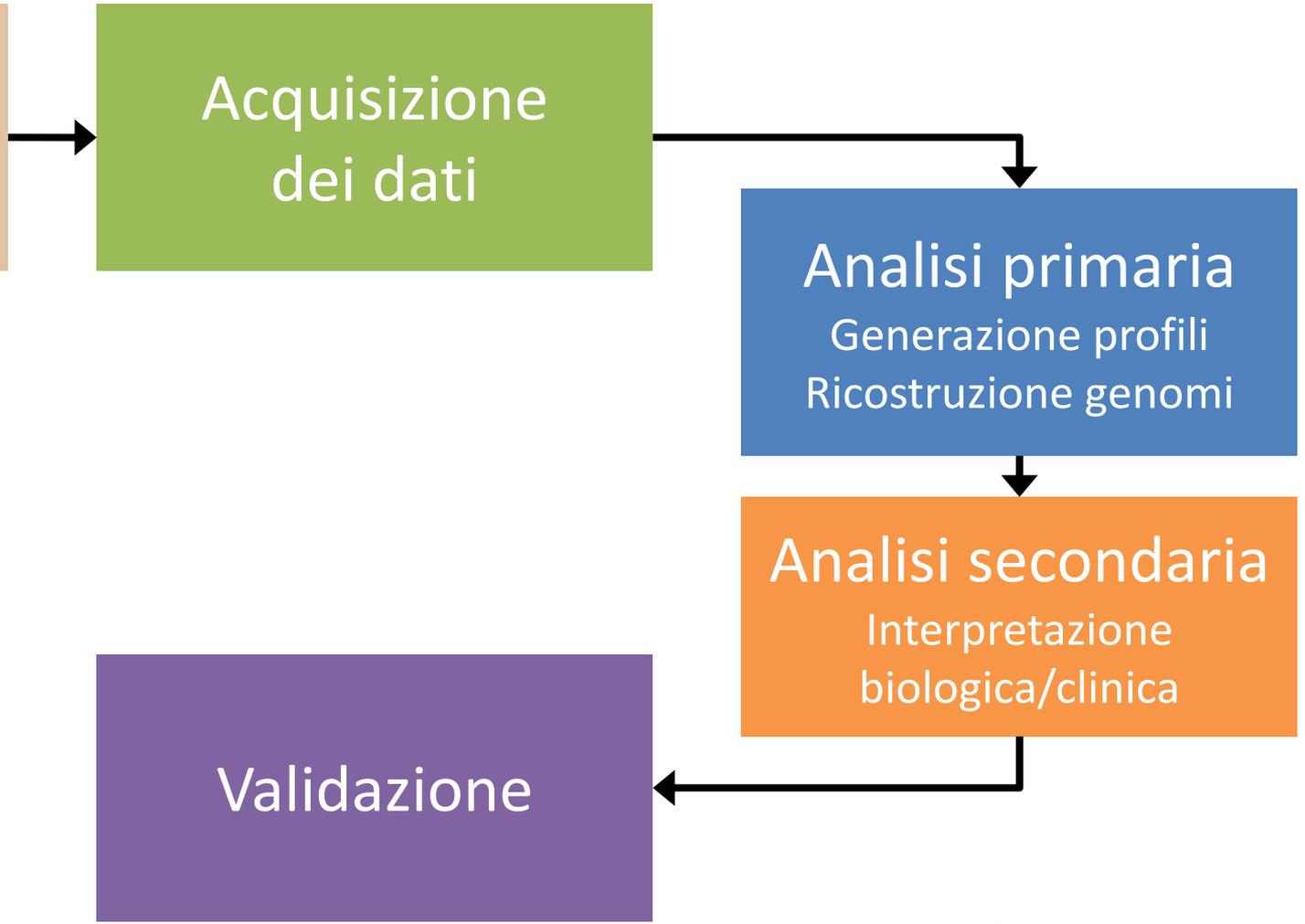


# Microbioma & Metagenomica

Da una prospettiva di *Data science*



Campioni biologici



# Microbioma & Metagenomica

Le mie linee di ricerca

*Machine  
learning per il  
microbioma*

*Tools per  
l'analisi a  
livello di ceppo*



edoardo.pasolli@unina.it

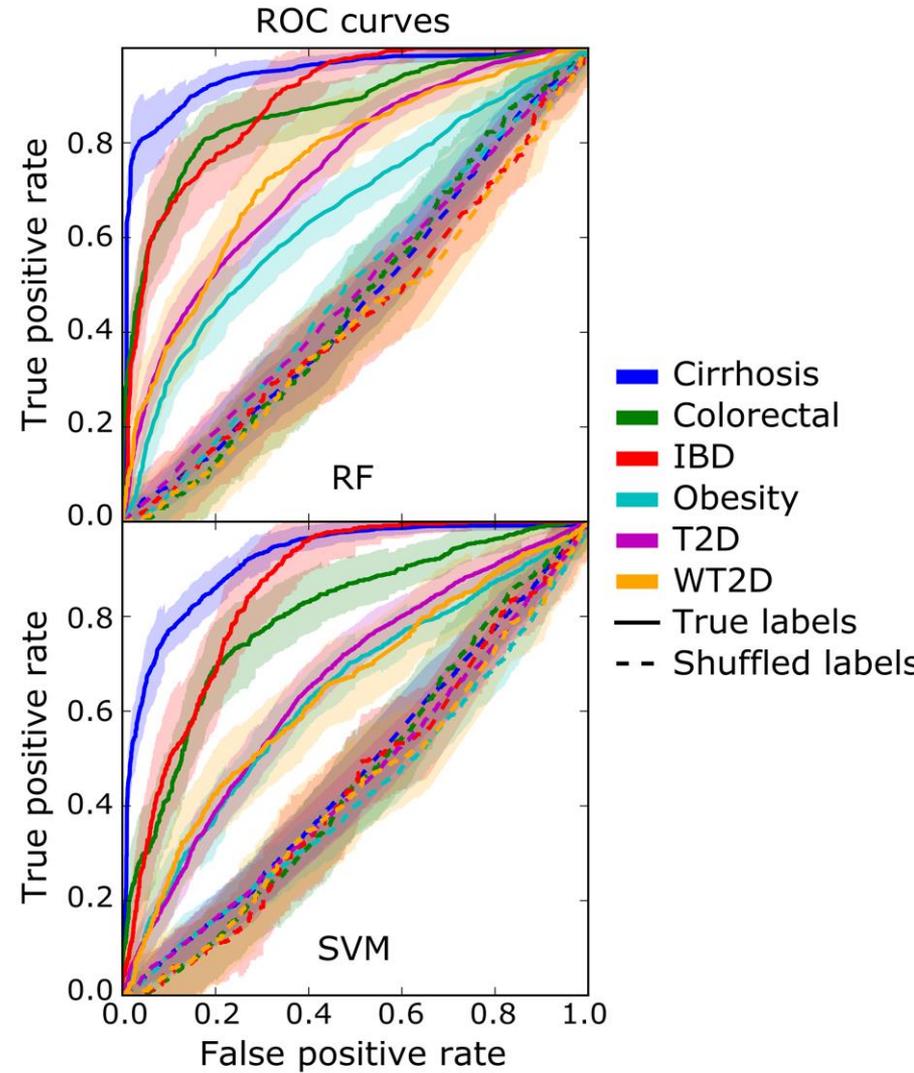
 @epasolli

# Microbioma & Modelli predittivi



## Machine Learning Meta-analysis of Large Metagenomic Datasets: Tools and Biological Insights

Edoardo Pasolli<sup>1</sup>, Duy Tin Truong<sup>1</sup>, Faizan Malik<sup>2</sup>, Levi Waldron<sup>2</sup>, Nicola Segata<sup>1\*</sup>



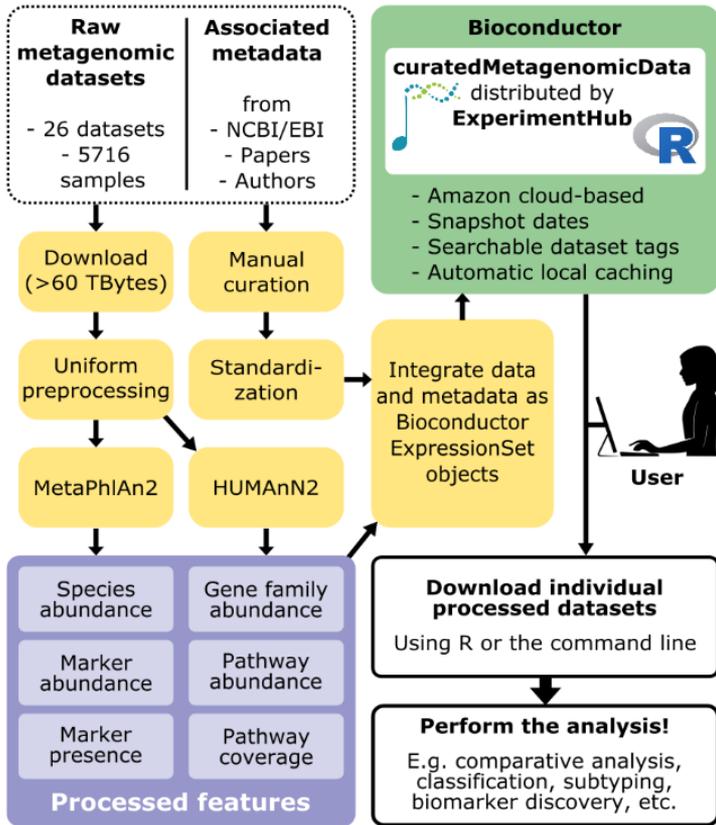
edoardo.pasolli@unina.it

@epasolli

# curatedMetagenomicData

## Offline high computational load pipeline

(incrementally performed on new data)

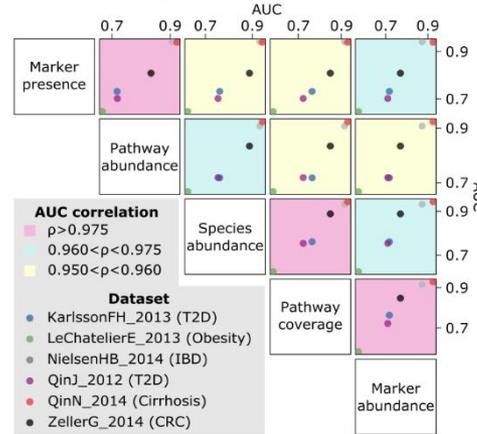


## nature **methods**

### Accessible, curated metagenomic data through ExperimentHub

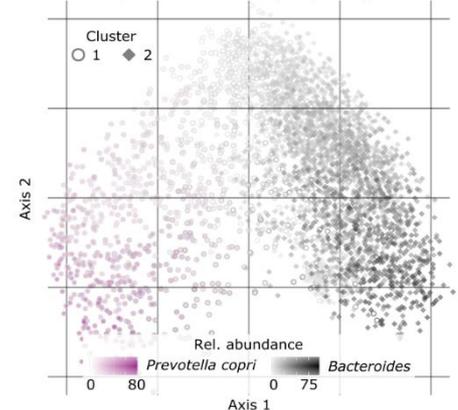
Edoardo Pasolli<sup>1,8</sup>, Lucas Schiffer<sup>2,3,8</sup>, Paolo Manghi<sup>1,8</sup>, Audrey Renson<sup>2,3</sup>, Valerie Obenchain<sup>3</sup>, Duy Tin Truong<sup>1</sup>, Francesco Beghini<sup>1</sup>, Faizan Malik<sup>2</sup>, Marcel Ramos<sup>2-4</sup>, Jennifer B Dowd<sup>2,5</sup>, Curtis Huttenhower<sup>6,7</sup>, Martin Morgan<sup>4</sup>, Nicola Segata<sup>1</sup> & Levi Waldron<sup>2,3</sup>

### Example 1: Classification

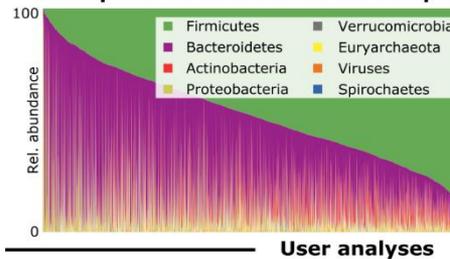


### Example 2: Clustering

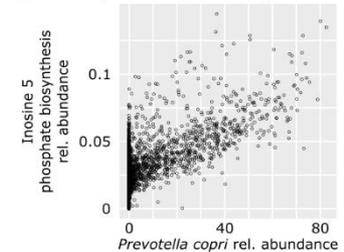
PCoA on species abundance, displaying 2 clusters



### Example 3: Abundance across samples



### Example 4: Species-pathway correlation



**Bioconductor**  
OPEN SOURCE SOFTWARE FOR BIOINFORMATICS

edoardo.pasolli@unina.it

@epasolli



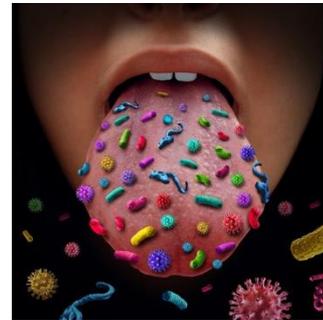
# Analisi a larga scala

Thomas\*, Manghi\*, et al.

nature  
medicine

**Metagenomic analysis of colorectal cancer datasets identifies cross-cohort microbial diagnostic signatures and a link with choline degradation**

**Progetto di  
dottorato industriale**



**Modelli predittivi basati su dieta e microbioma  
per lo studio di malattie del cavo orale**



edoardo.pasolli@unina.it

 @epasolli

# Tools per l'analisi a livello di ceppo

Assunzione: Ciascun **ceppo** ha una combinazione unica di...

... geni dal pangenoma della specie

**PanPhlAn**

**Strain-level microbial epidemiology and population genomics from shotgun metagenomics**

Matthias Scholz<sup>1,4</sup>, Doyle V Ward<sup>2,4</sup>, Edoardo Pasolli<sup>1,4</sup>, Thomas Tolio<sup>1</sup>, Moreno Zolfo<sup>1</sup>, Francesco Asnicar<sup>1</sup>, Duy Tin Truong<sup>1</sup>, Adrian Tett<sup>1</sup>, Ardythe L Morrow<sup>3</sup> & Nicola Segata<sup>1</sup>

nature|**methods**

... SNVs nel *core genome*

**StrainPhlAn**

**Microbial strain-level population structure and genetic diversity from metagenomes**

Duy Tin Truong,<sup>1</sup> Adrian Tett,<sup>1</sup> Edoardo Pasolli,<sup>1</sup> Curtis Huttenhower,<sup>2,3</sup> and Nicola Segata<sup>1</sup>

GENOME  
RESEARCH



edoardo.pasolli@unina.it

 @epasolli

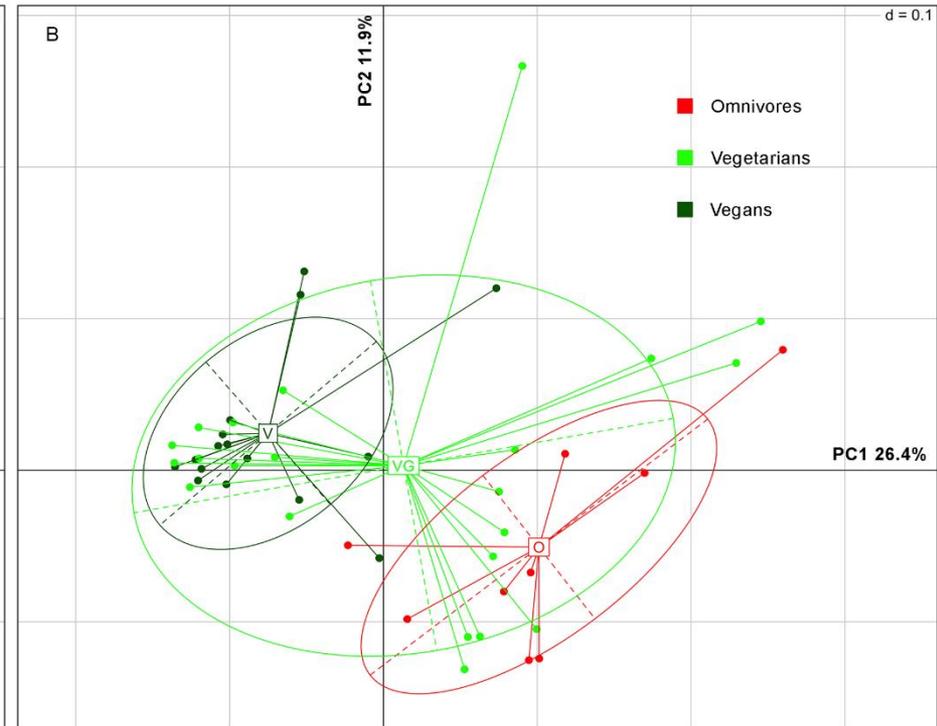
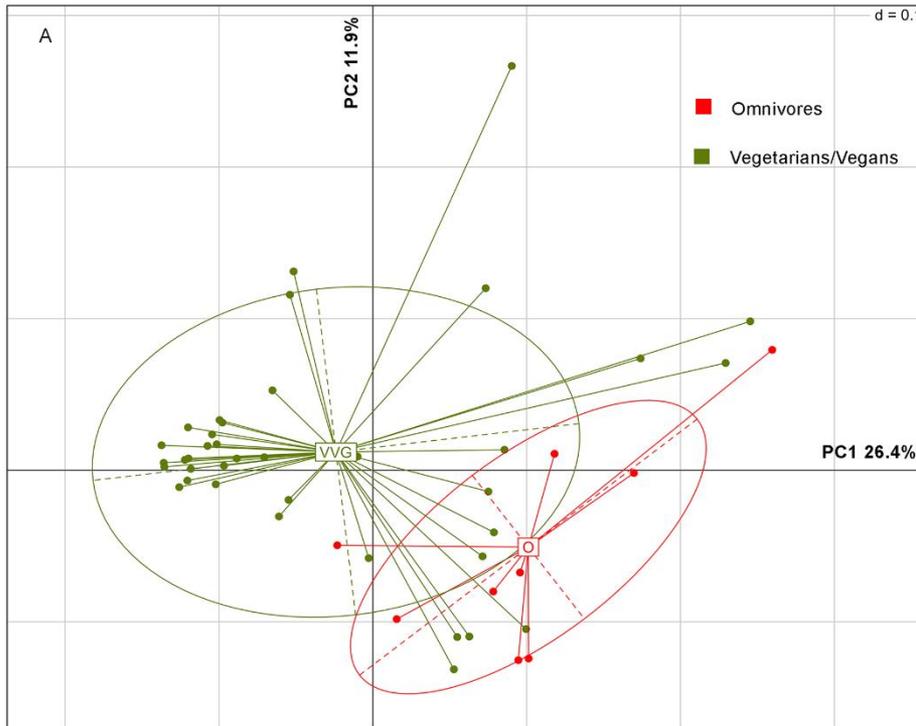
# Prevotella copri & Dieta



## Cell Host & Microbe

### Distinct Genetic and Functional Traits of Human Intestinal *Prevotella copri* Strains Are Associated with Different Habitual Diets

Francesca De Filippis,<sup>1,2</sup> Edoardo Pasoli,<sup>1,3</sup> Adrian Tett,<sup>3</sup> Sonia Tarallo,<sup>4</sup> Alessio Naccarati,<sup>4</sup> Maria De Angelis,<sup>5</sup> Erasmo Neviani,<sup>6</sup> Luca Cocolin,<sup>7</sup> Marco Gobetti,<sup>8</sup> Nicola Segata,<sup>3</sup> and Danilo Ercolini<sup>1,2,9,\*</sup>



edoardo.pasoli@unina.it

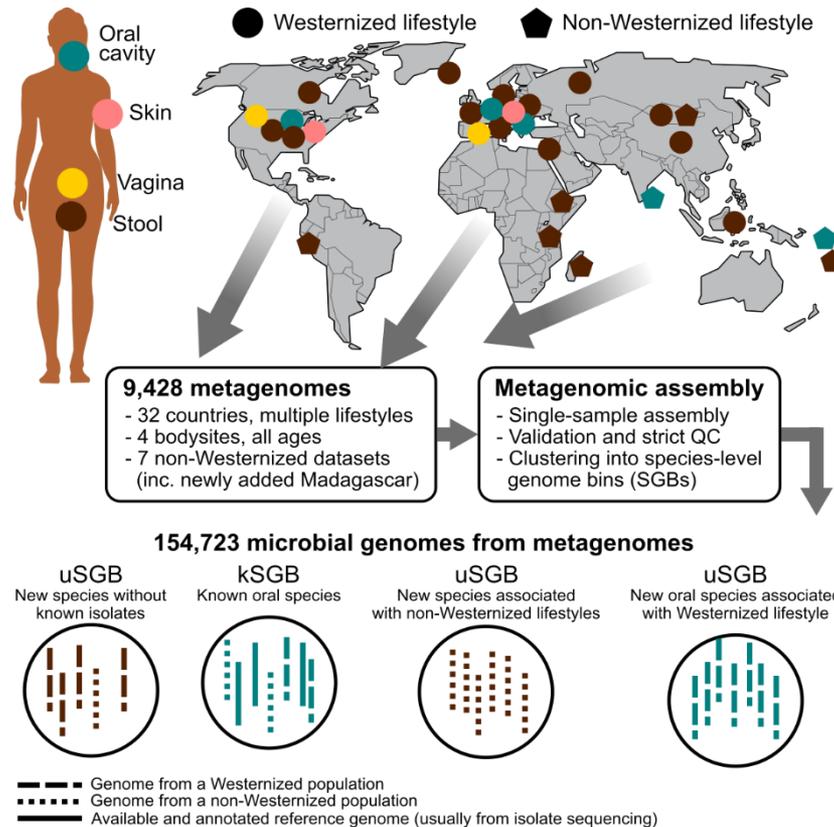
 @epasoli

# Il più ricco catalogo del microbioma umano

## Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle



Edoardo Pasoli,<sup>1</sup> Francesco Asnicar,<sup>1,8</sup> Serena Manara,<sup>1,8</sup> Moreno Zolfo,<sup>1,8</sup> Nicolai Karcher,<sup>1</sup> Federica Armanini,<sup>1</sup> Francesco Beghini,<sup>1</sup> Paolo Manghi,<sup>1</sup> Adrian Tett,<sup>1</sup> Paolo Ghensi,<sup>1</sup> Maria Carmen Collado,<sup>2</sup> Benjamin L. Rice,<sup>3</sup> Casey DuLong,<sup>4</sup> Xochitl C. Morgan,<sup>5</sup> Christopher D. Golden,<sup>4</sup> Christopher Quince,<sup>6</sup> Curtis Huttenhower,<sup>4,7</sup> and Nicola Segata<sup>1,9,\*</sup>

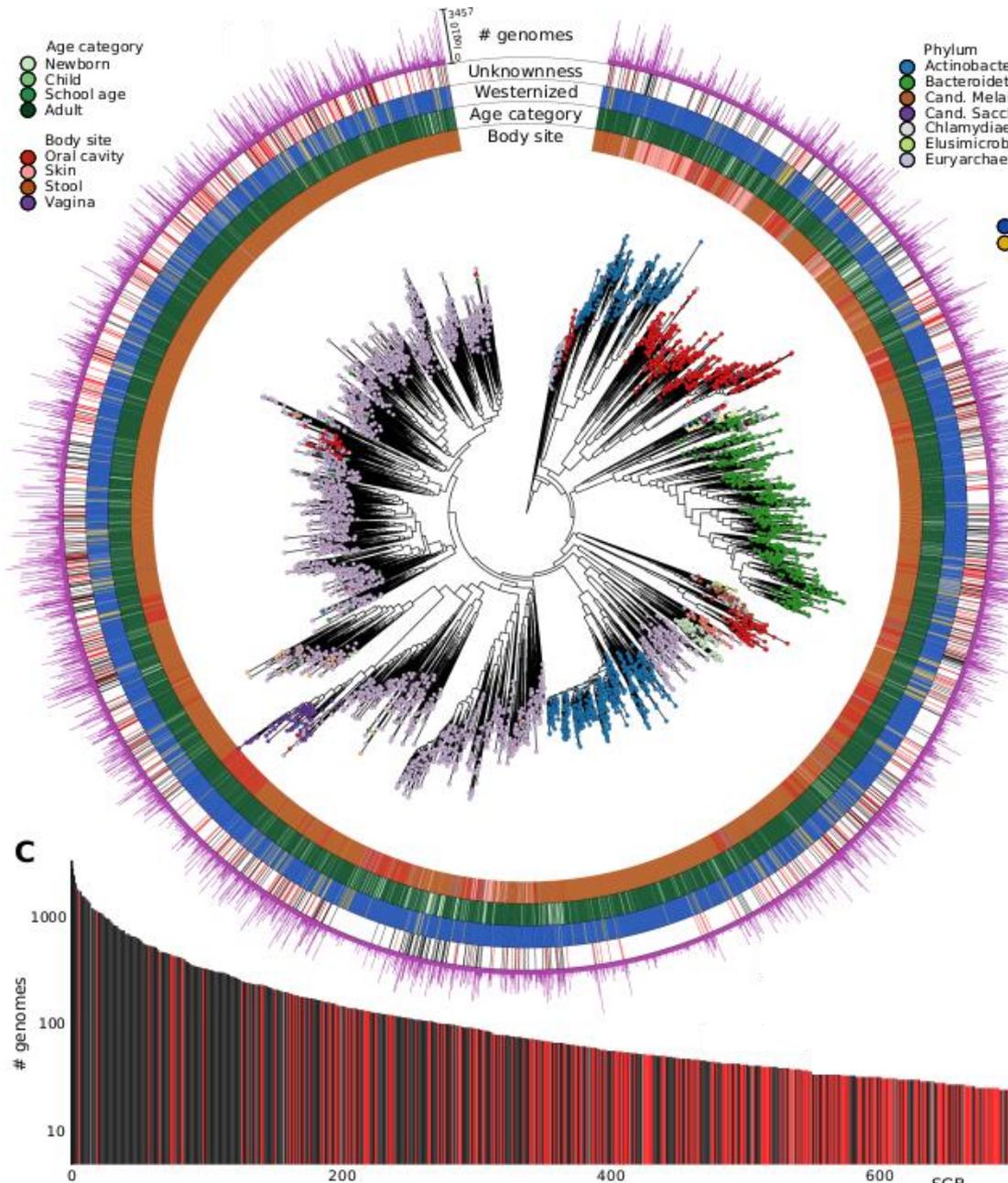


edoardo.pasoli@unina.it

@epasoli



# Il più ricco catalogo del microbioma umano



- 154.723 genomi ricostruiti da 9.428 metagenomi
- Raddoppiato il numero di genomi disponibili
- Identificate 4,930 specie (77% nuove)
- Nuove specie prevalenti in popolazioni non-occidentalizzate
- *Mappability* da 68% a 88% (mediana 94%) nell'intestino
- Da 65% a 82% nella cavità orale

# Lactic acid bacteria, bridging food and gut health

**Batteri lattici sono  
altamente prevalenti  
negli alimenti fermentati.  
Ma quanto lo sono  
nell'intestino?**

- Analisi condotta su 9.000 metagenomi umani
- *S. thermophilus* e *Lactobacillus* in popolazioni occidentalizzate
- *Leuconostoc* e *Weissella* in popolazioni non-occidentalizzate

**Batteri lattici negli  
alimenti e nell'intestino,  
come differiscono  
funzionalmente?**

- Ricostruiti 650 nuovi genomi da 250 metagenomi da cibo
- Integrazione con il catalogo di genomi del microbioma umano
- **Analisi funzionale in corso**



In collaborazione con prof. Danilo Ercolini e dott. Francesca De Filippis

edoardo.pasolli@unina.it

 @epasolli



**Prossimo appuntamento**

**18 Settembre 2019**

***Il valore della diversità***

**Francesco Caracciolo di Torchiarolo**

