

# Data-centric Machine Learning Applied to Sustainability

**Laure Berti-Equille**

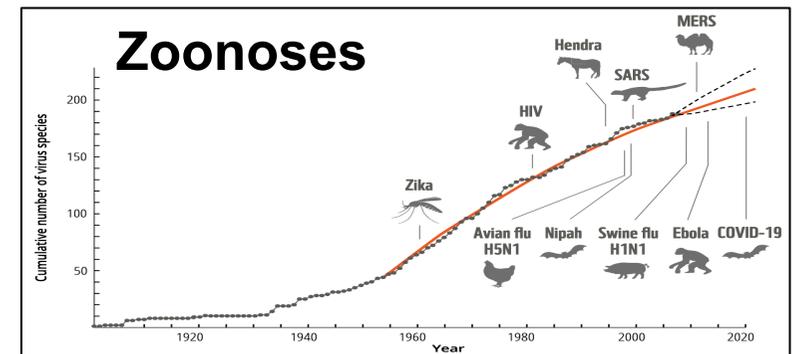
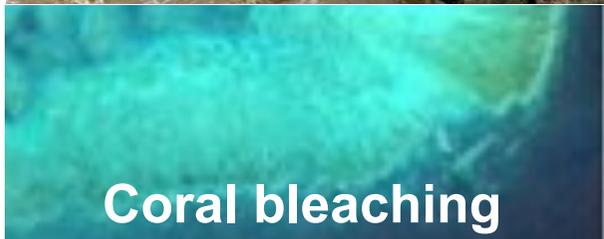
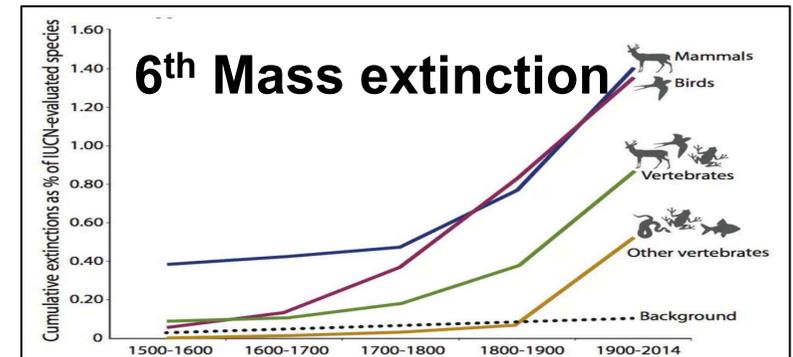
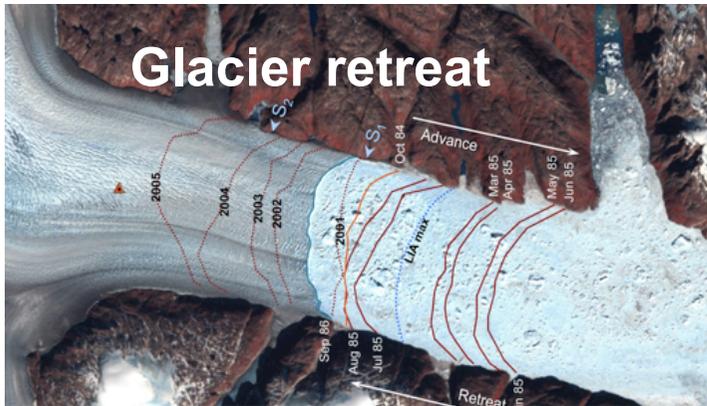
IRD ESPACE-DEV  
Montpellier, France

[laure.berti@ird.fr](mailto:laure.berti@ird.fr)

<https://laureberti.github.io/website/>



# Today's View of Our World



# 2030 Sustainable Development Goals



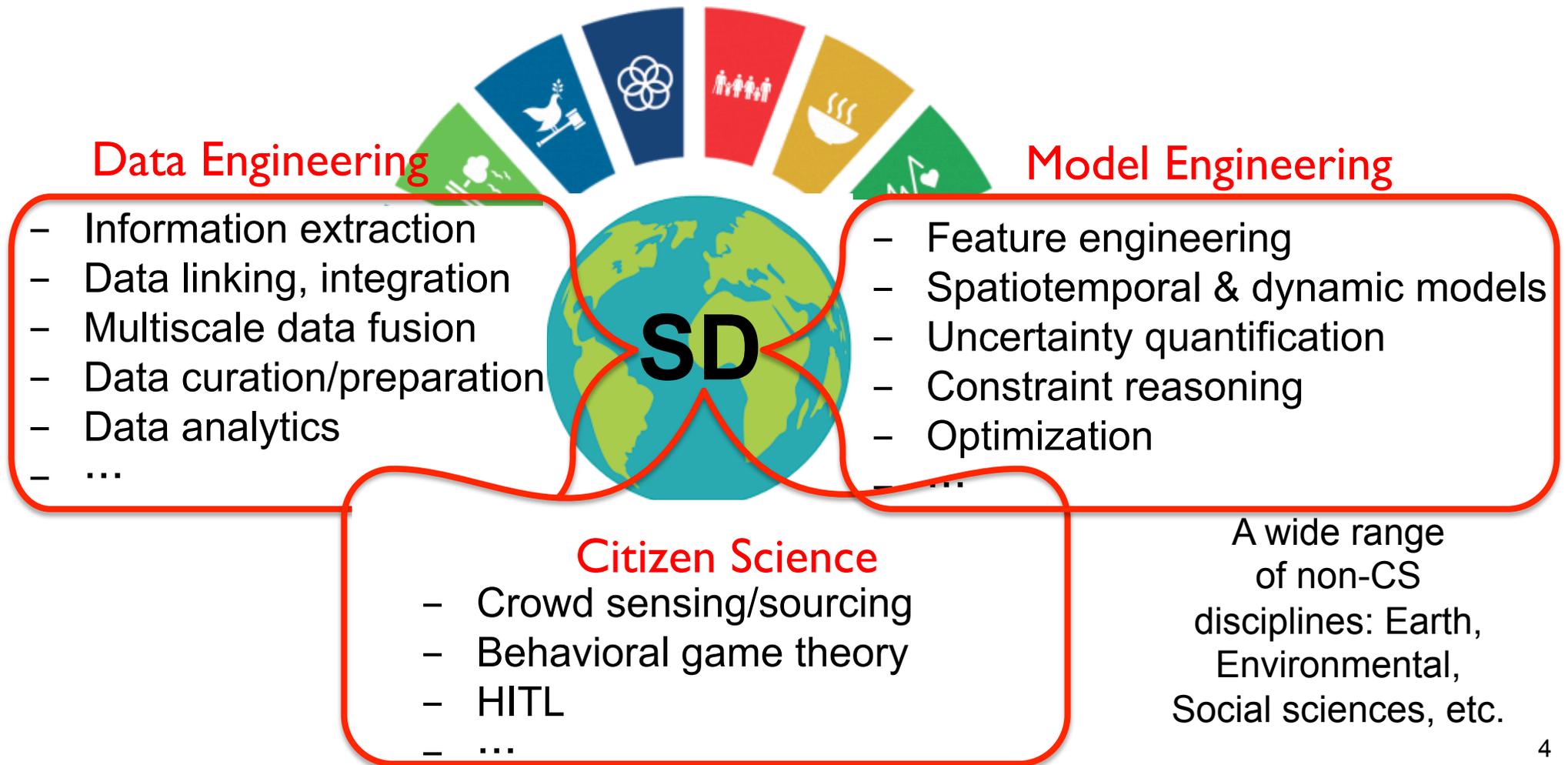
2030 Agenda for Sustainable Development: 17 goals, 169 targets, 232 Indicators  
*New norms to integrate the principles of sustainable development into country policies and programs*

<https://sdgs.un.org/2030agenda>

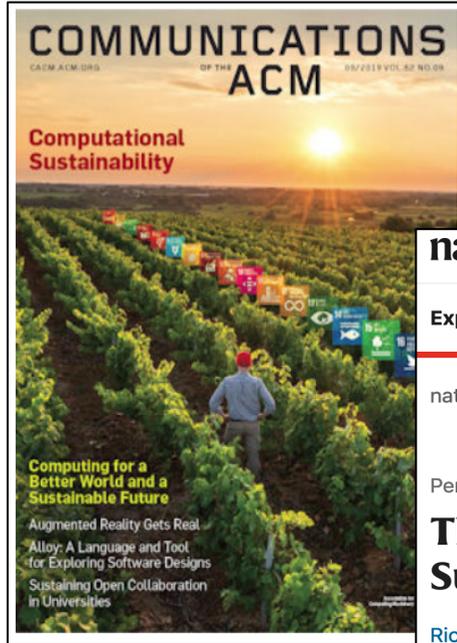
<https://sdgs.un.org/goals>

# Applying Data Science & AI to SDGs

Transdisciplinary research, reasoning, and discovery from interconnected information contents



# A Decade of AI Contributions & Initiatives



### Tackling Climate Change with Machine Learning

David Rolnick<sup>1\*</sup>, Priya L. Donti<sup>2</sup>, Lynn H. Kaack<sup>3</sup>, Kelly Kochanski<sup>4</sup>, Alexandre Lacoste<sup>5</sup>, Kris Sankaran<sup>6,7</sup>, Andrew Slavin Ross<sup>9</sup>, Nikola Milojevic-Dupont<sup>10,11</sup>, Natasha Jaques<sup>12</sup>, Anna Waldman-Brown<sup>12</sup>, Alexandra Luccioni<sup>6,7</sup>, Tegan Maharaj<sup>6,8</sup>, Evan D. Sherwin<sup>2</sup>, S. Karthik Mukkavilli<sup>6,7</sup>, Konrad P. Körding<sup>1</sup>, Carla Gomes<sup>13</sup>, Andrew Y. Ng<sup>14</sup>, Demis Hassabis<sup>15</sup>, John C. Platt<sup>16</sup>, Felix Creutzig<sup>10,11</sup>, Jennifer Chayes<sup>17</sup>, Yoshua Bengio<sup>6,7</sup>

nature communications

Explore content ▾ Journal information ▾ Publish with us ▾

nature > nature communications > perspectives > article

Perspective | Open Access | Published: 13 January 2020

## The role of artificial intelligence in achieving the Sustainable Development Goals

Ricardo Vinuesa ✉, Hossein Azizpour, Iolanda Leite, Madeline Balaam, Virginia Dignum, S. Domsch, Anna Felländer, Simone Daniela Leachman, Max Tegmark & Francesco Ferrero

## CompSustNet

Computational Sustainability Network

[www.compsust.net](http://www.compsust.net)



[www.climatechange.ai](http://www.climatechange.ai)

## SDGacademy

<https://sdgacademy.org/>

KDD 2017 Applied Invited Talk KDD'17, August 13–17, 2017, Halifax, NS, Canada

### Big Data in Climate: Opportunities and Challenges for Machine Learning

Anuj Karpatne  
 Department of Computer Science and Engineering,  
 University of Minnesota  
 karpa009@umn.edu

Vipin Kumar  
 Department of Computer Science and Engineering,  
 University of Minnesota  
 kumar001@umn.edu

**ABSTRACT**  
 The climate and Earth sciences have recently undergone a rapid transformation from a data-poor to a data-rich environment. In particular, massive amount of data about Earth and its environment is now continuously being generated by a large number of Earth observing satellites as well as physics-based earth system models running on large-scale computational platforms. These massive and information-rich datasets offer huge potential for understanding how the Earth's climate and ecosystem have been changing and how they are being impacted by humans actions. We discuss the challenges involved in analyzing these massive data sets as well as opportunities they present for both advancing machine learning as well as the science of climate change.



## Climate Informatics

[www.climateinformatics.org](http://www.climateinformatics.org)

ENERGY

## How A.I. Will Revolutionize Climate Tech

This week on The Interchange: we look deeper at artificial intelligence as a climate

STEPHEN LACEY | JUNE 17, 2021

APPLIED ECONOMIC PERSPECTIVES AND POLICY

Submitted Article

### Big Data in Agriculture: A Challenge for the Future

Keith H Coble ✉, Ashok K Mishra, Shannon Ferrell, Terry Griffin

First published: 16 February 2018 | <https://doi.org/10.1093/aep/pxx056> | Citations: 10

# SD as Optimization Problems

- Maximizing the probability of achieving an SD target
- Minimizing the degradations of environmental and human conditions

## Complexity:

- Multi-objective: improve the quality of human life, preserve the Earth's diversity, minimize the depletion of non-renewable resources
- Multi-disciplinary: environmental, social sciences, etc.
- Multi-scale in time and space: global, national, regional, local with various horizons
- Multi-actor: civil society, private companies, government

# ML for SD: Main Challenges

- Data**
  1. Lack of representative ground truth data
  2. Dirty and imperfect data with data quality issues, uncertainty, label noise, bias
- Models**
  3. Inappropriate feature engineering and inadequate data preprocessing
  4. Energy efficiency and carbon footprint reduction of ML models
- Output**
  5. Non actionable and misleading outputs
  6. Difficult to measure the impact of ML-based solutions
  7. Interpretability and trust issues for local actors

# Progress on the SDGs



The screenshot shows a web browser window with the URL <https://sdg-tracker.org>. The page has a dark blue header with "SDG Tracker" on the left and "About Our World in Data" on the right. The main content area features a large heading "Measuring progress towards the Sustainable Development Goals" followed by three paragraphs of introductory text.

SDG Tracker About Our World in Data

## Measuring progress towards the Sustainable Development Goals

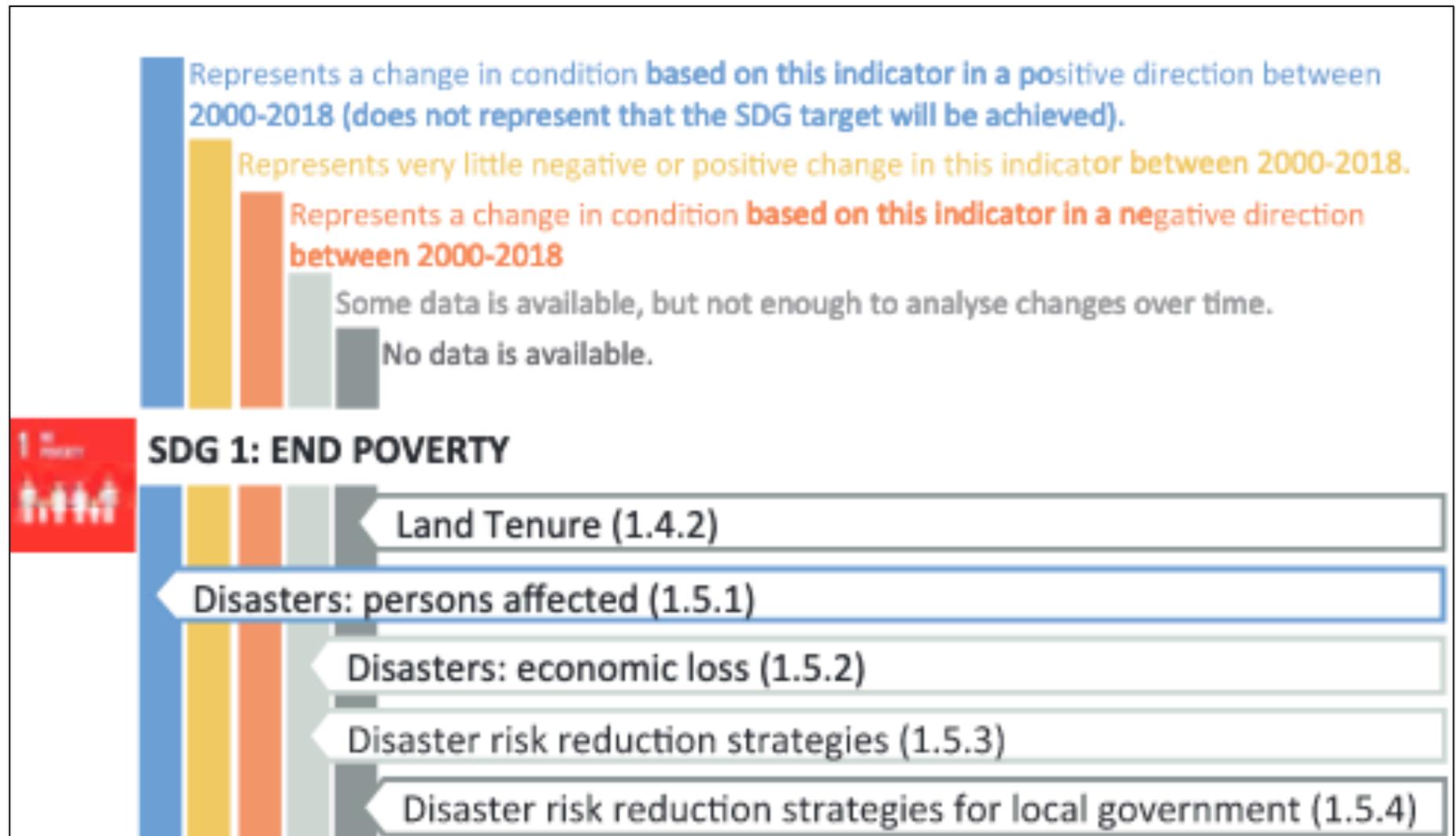
The United Nations [Sustainable Development Goals](#) (SDGs) are targets for global development adopted in September 2015, set to be achieved by 2030. All countries of the world have agreed to work towards achieving these goals.

Our SDG Tracker presents data across all available indicators from the [Our World in Data](#) database, using official statistics from the UN and other international organizations. It is a free, open-access publication that tracks global progress towards the SDGs and allows people around the world to hold their governments accountable to achieving the agreed goals.

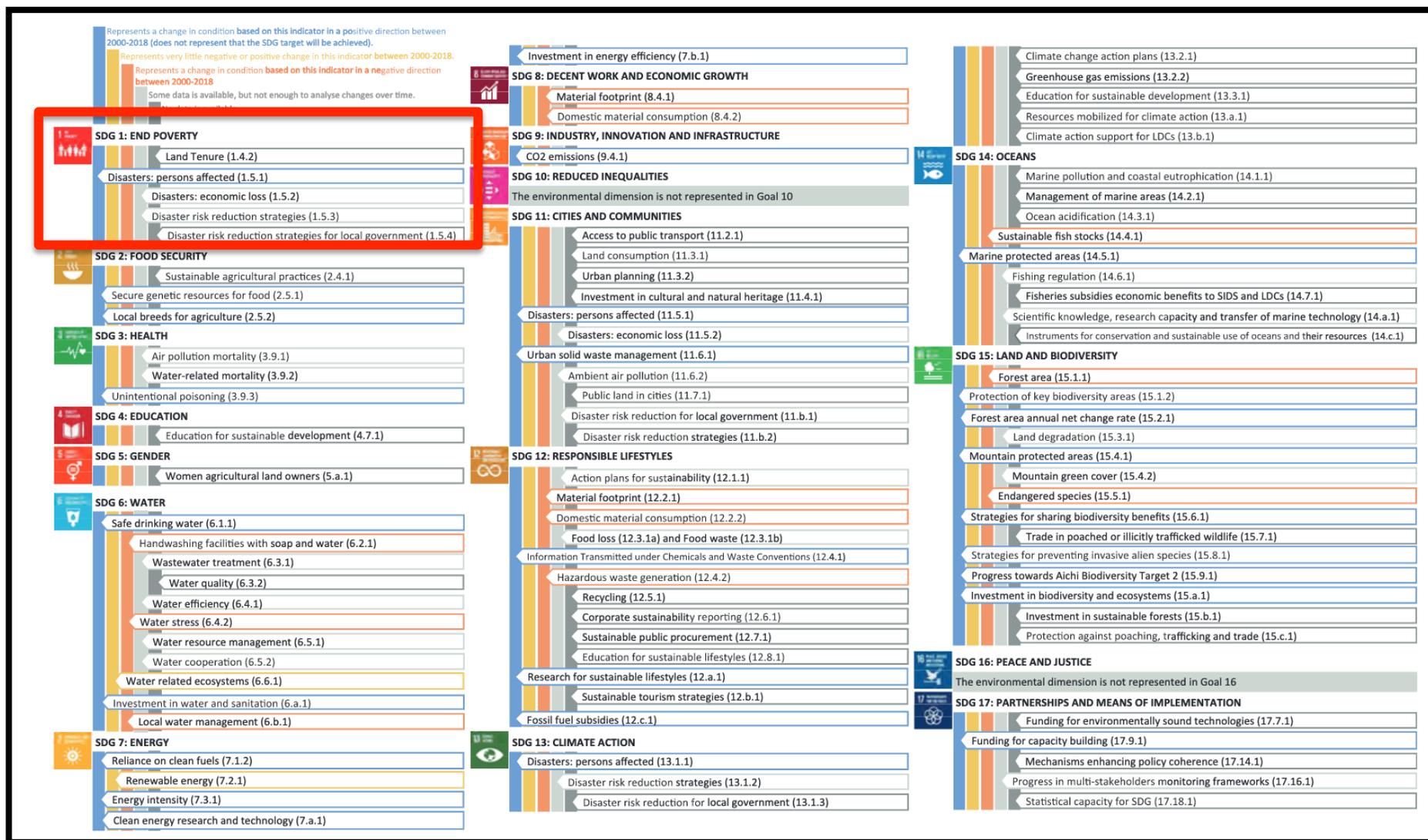
The 17 Sustainable Development Goals are defined in a list of 169 SDG Targets. Progress towards these Targets is agreed to be tracked by 232 unique Indicators. Here is the [full list of definitions](#).

<https://sdg-tracker.org/>

# Example: The Lack of Data



# Progress on the SDGs



# Outline

## I- Challenges in ML & data science pipelines

- Building the pipelines
- Preparing the data

## II- Current Projects

- Combining satellite imagery and socio-economic indicators to estimate poverty evolution in Africa
- Finding sustainable transition pathways in Nordeste



# Illustrative Example I

1. SD Question

2. Actors

3. Time/space of Interest

4. Data & Knowledge

## Multi-Scale

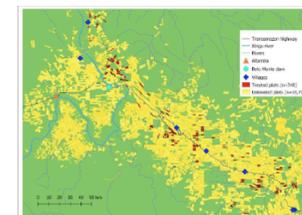
Brazil



Minas Gerais



Area of study



## Multi-Actor

*at country-scale*

- National authorities
- Policymakers
- Scientific community
- General public

*at state-scale*

- Federal authorities
- Regional actors

*at territory-scale*

- Administrator representatives of the population (stakeholders private and public)
- Site managers
- Local administrations
- Elected representatives of the local authorities
- Agroforestry exploitation owners
- Farmers
- Indigenous representatives

# Illustrative Example I

## 1. SD Question

What are the fines effective for reducing deforestation?

## 2. Actors

Coffee farmers, policy makers, local administrators

## 3. Time/space of Interest

2018-2022, Minas Gerais, Brazil

## 4. Data & Knowledge

Satellite images, surveys, agroecology datasets, etc.

### Multi-Scale

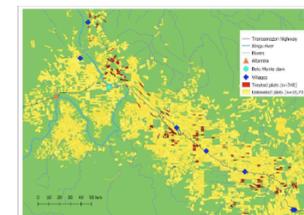
Brazil



Minas Gerais



Area of study



### Multi-Actor

*at country-scale*

- National authorities
- Policymakers
- Scientific community
- General public

*at state-scale*

- Federal authorities
- Regional actors

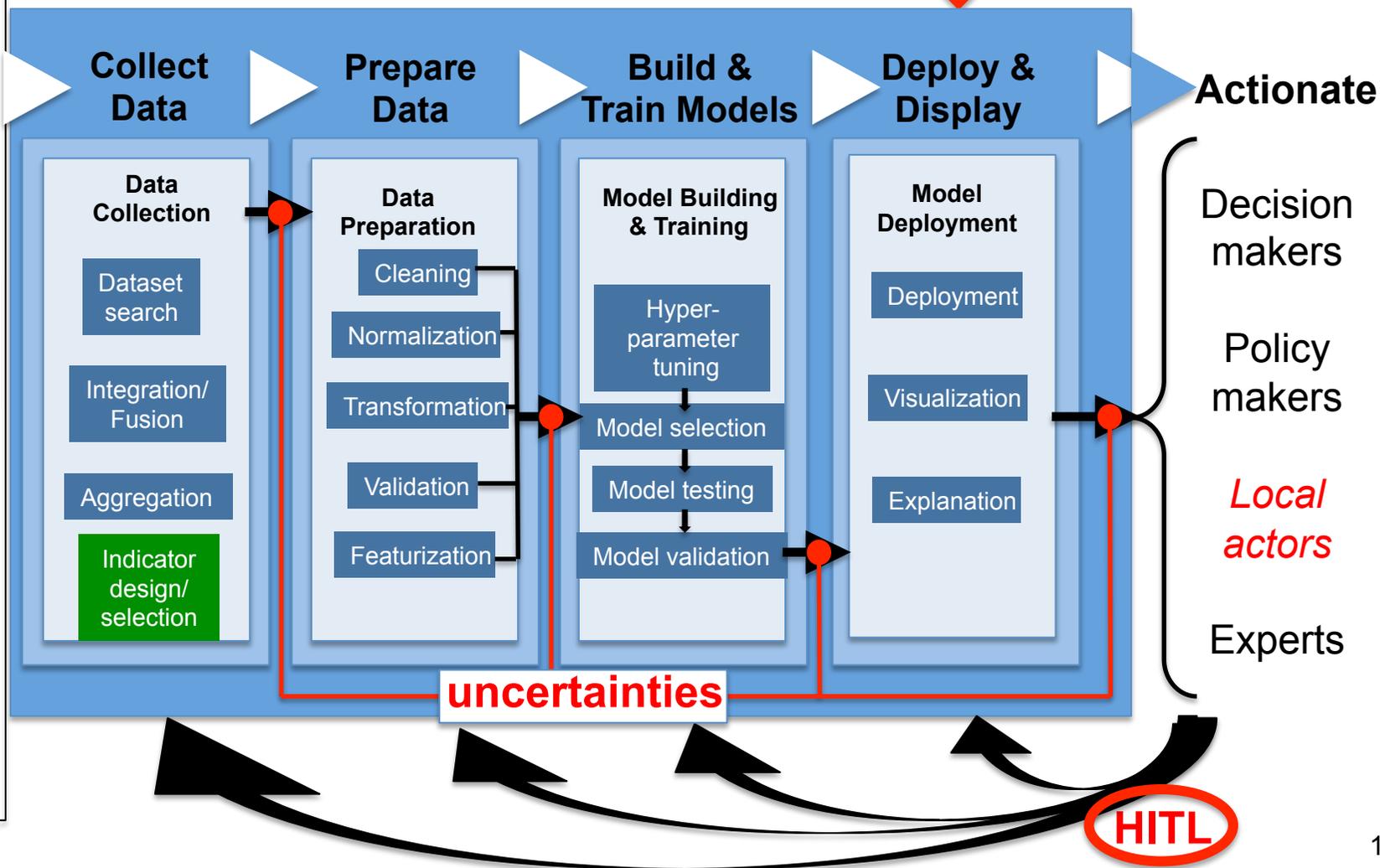
*at territory-scale*

- Administrator representatives of the population (stakeholders private and public)
- Site managers
- Local administrations
- Elected representatives of the local authorities
- Agroforestry exploitation owners
- Farmers
- Indigenous representatives

# Example 2. SD Data Science Pipeline

Multiple datasets

Are the fines effective for reducing deforestation?





# Data prep is crucial (2/2)

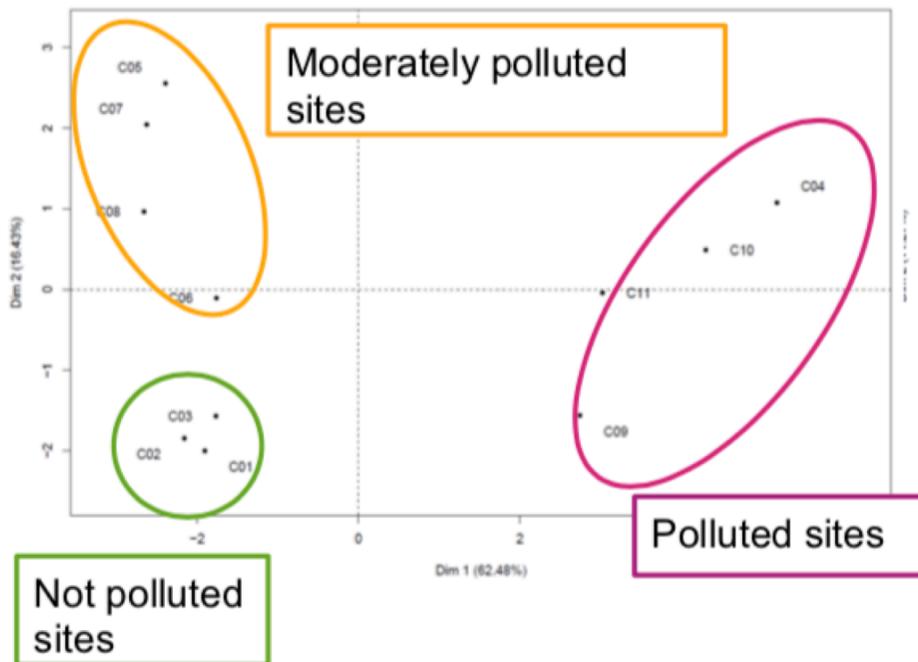
**Various data cleaning/preparation strategies lead to different and misleading conclusions**

• Duplicates   
 • Inconsistencies   
 • Outliers   
 • Incomplete and missing data

C
3.2
6.3
NA

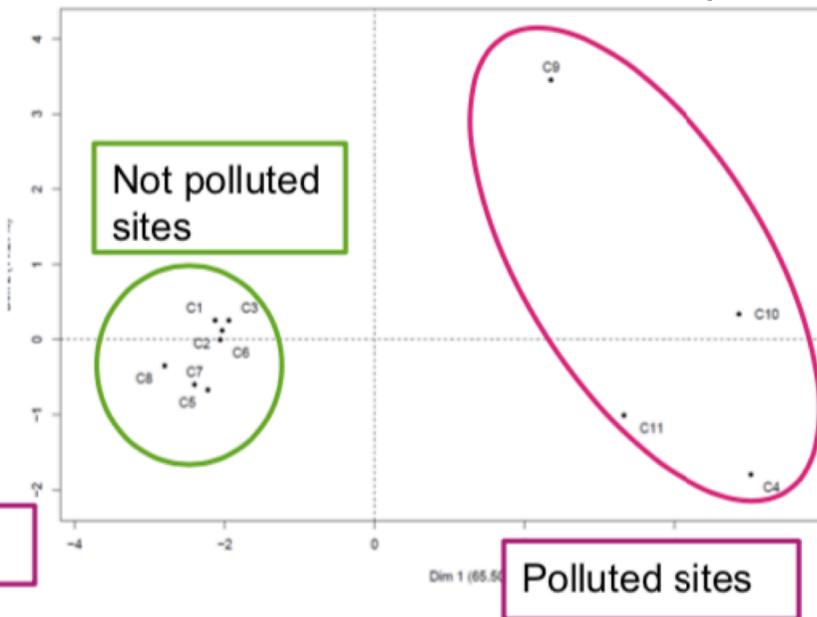
*Preprocessing 1*

Individuals factor map (PCA)

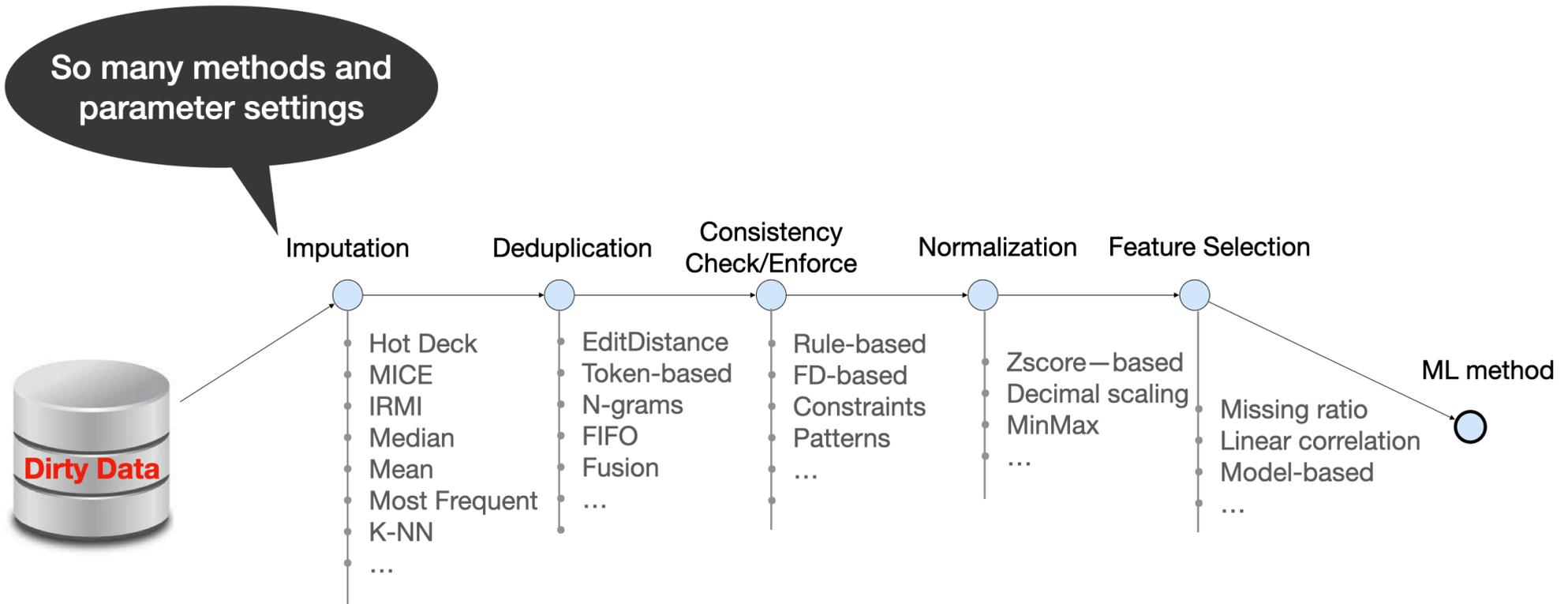


*Preprocessing 2*

Individuals factor map (PCA)

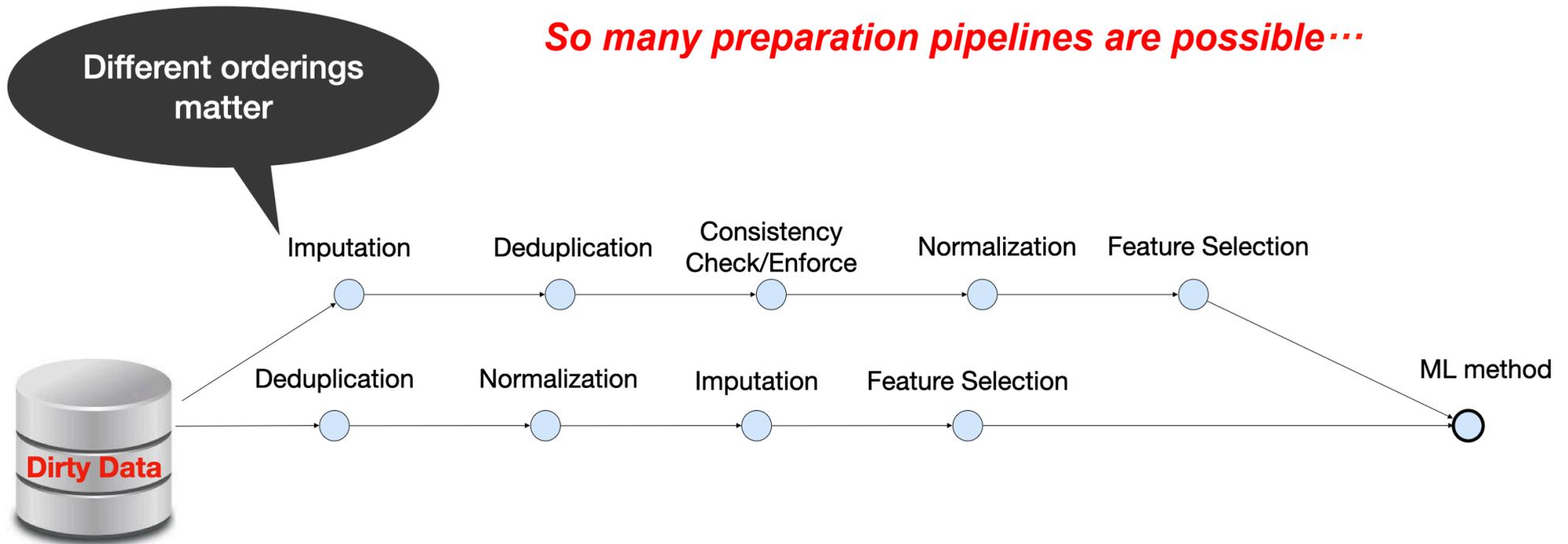


# Data prep is challenging



# Data prep is challenging

*So many preparation pipelines are possible...*

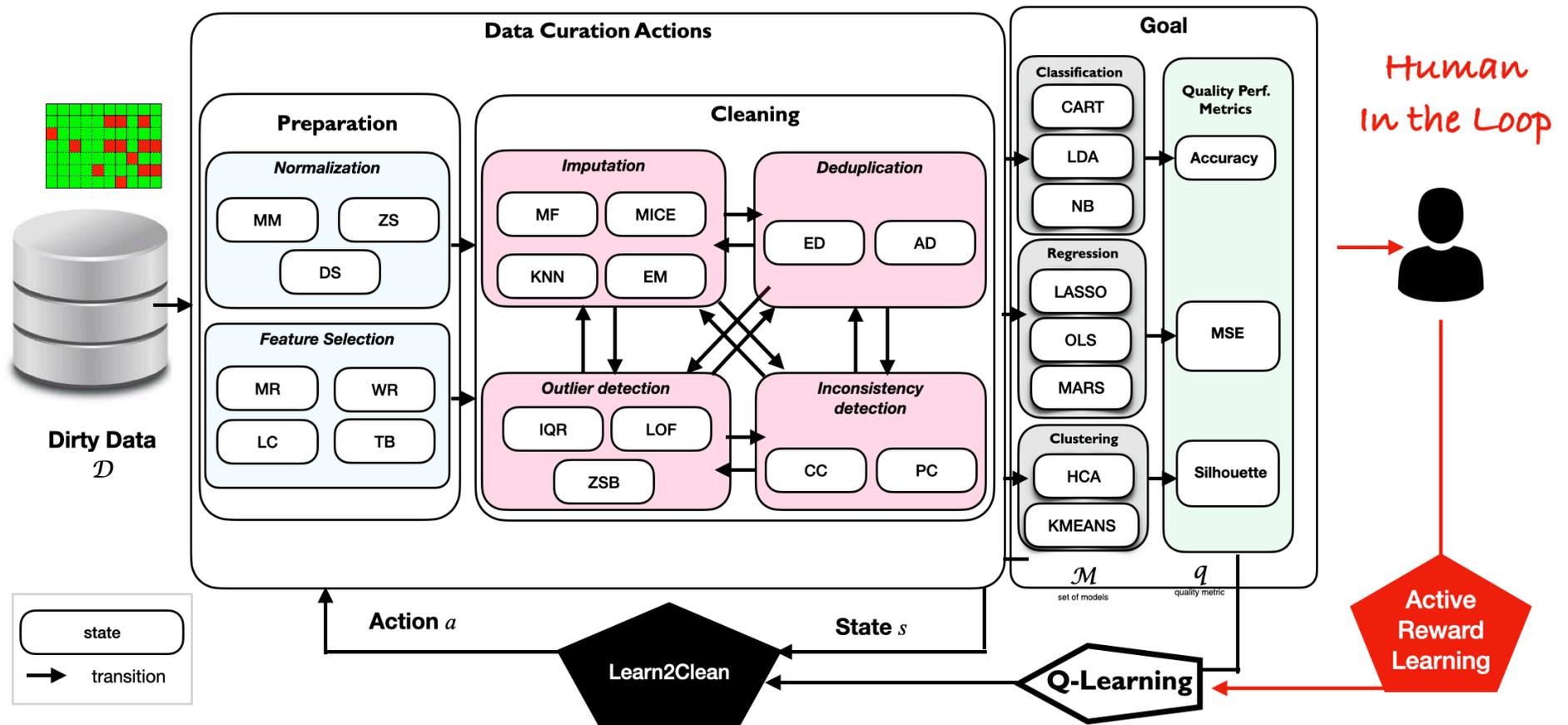


*Reiterate some steps is possible...*

*Selective prep is also possible  
... and even recommended*

# Reinforcement Learning can help data prep

[Berti-Equille, 2019, 2020]



<https://github.com/LaureBerti/Learn2Clean>

# Outline

## I- Challenges in ML & data science pipelines

- Building the pipelines
- Preparing the data

## II- Current Projects

- Combining satellite imagery and socio-economic indicators to estimate poverty evolution in Africa
- Finding sustainable transition pathways in Brazil



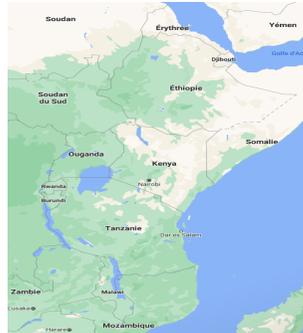
# Predict Poverty Evolution from Series of Satellite Images



MPA Poverty  
funded by ANR

➤ Incomplete surveys from 2009 to 2016 for 4 African Countries

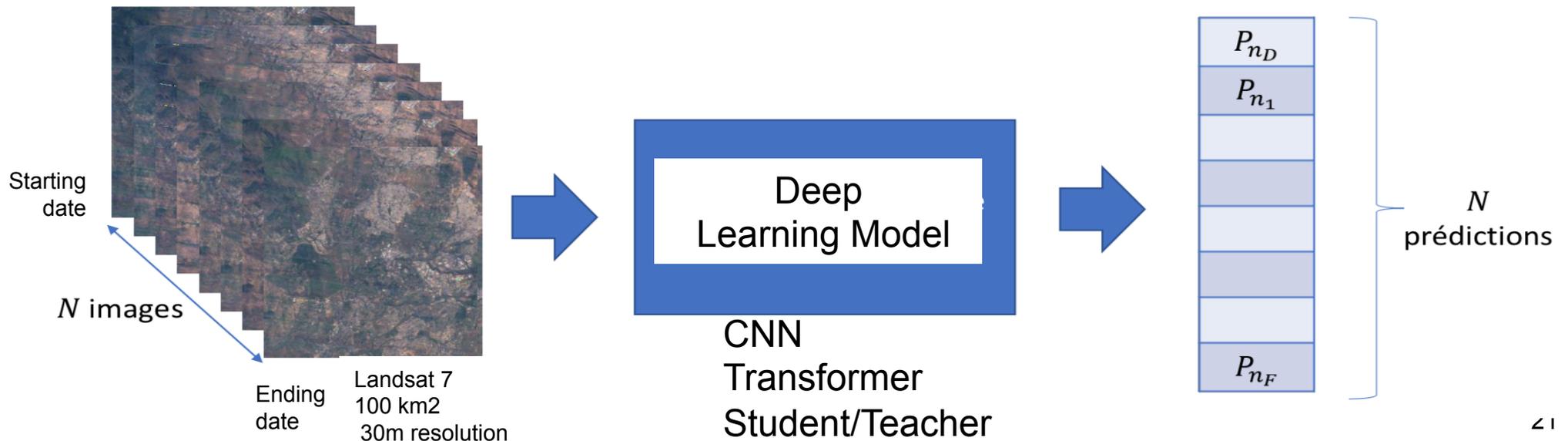
Economic  
wealth  
indicators



- Malawi : 3 surveys
- Ethiopia : 3 surveys
- Ouganda : 2 surveys
- Tanzania : 3 suveys

5800 ground truth  
data points as  
indicators

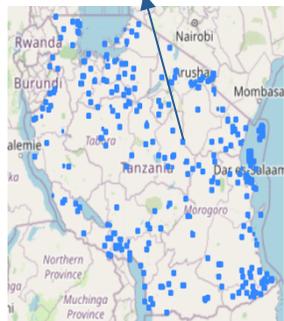
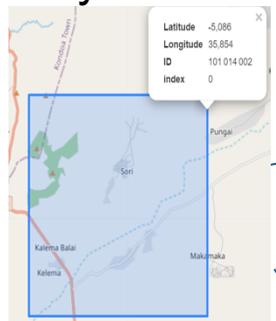
➤ Sequences of images at different resolutions





# Poverty Evolution Prediction: A Naive Spatial Approach

Area covered  
by SITS



Wealth  
indicators

2000



$R^2 = 0.71$

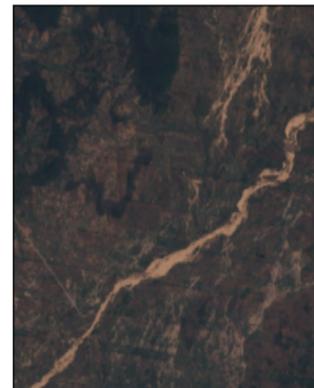
Spatial Model  
(Yeh, 2020)

2000

1.98 \$/d/p

·  
·  
·

2020



$R^2 = 0.71$

Spatial Model  
(Yeh, 2020)

2020

2.74 \$/d/p

<sup>1</sup>C. Yeh, A. Perez, A. Driscoll, *et al.* "Using publicly available satellite imagery and deep learning to understand economic well-being in Africa." *Nat Commun* 11, 2583 (2020).

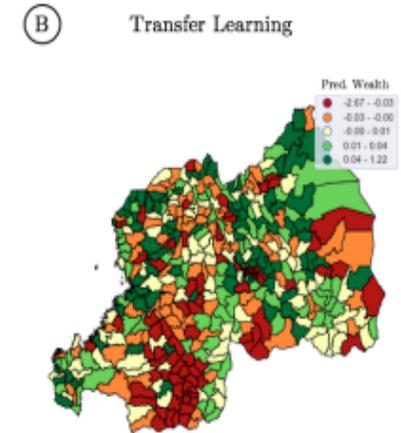
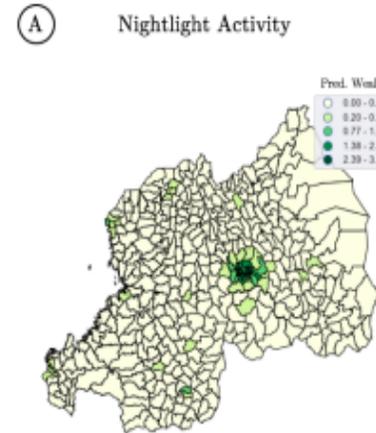
# Poverty Evolution Prediction: Limited Results of the Naive Approach



[Kondmann et al., 2020]<sup>1</sup>

	Mean Training $R^2$	Test $R^2$	Test $R^2$
	2005	2010	2015
Transfer Learning & ResNet50	0.47	0.69	0.57
Mean Nightlights per cluster	0.52	0.74	0.61
ResNet50	0.55	0.43	0.36

Time	Change TL	Change NTL	Poverty Reduction	GDP Growth
2005 - 2010	0.00	0.03***	0.11	0.49
2010 - 2015	0.00	0.02***	0.07	0.44
2005 - 2015	0.00	0.05***	0.18	1.15



- Very heterogeneous evolution map
- **Consistency in time and space may not be preserved**

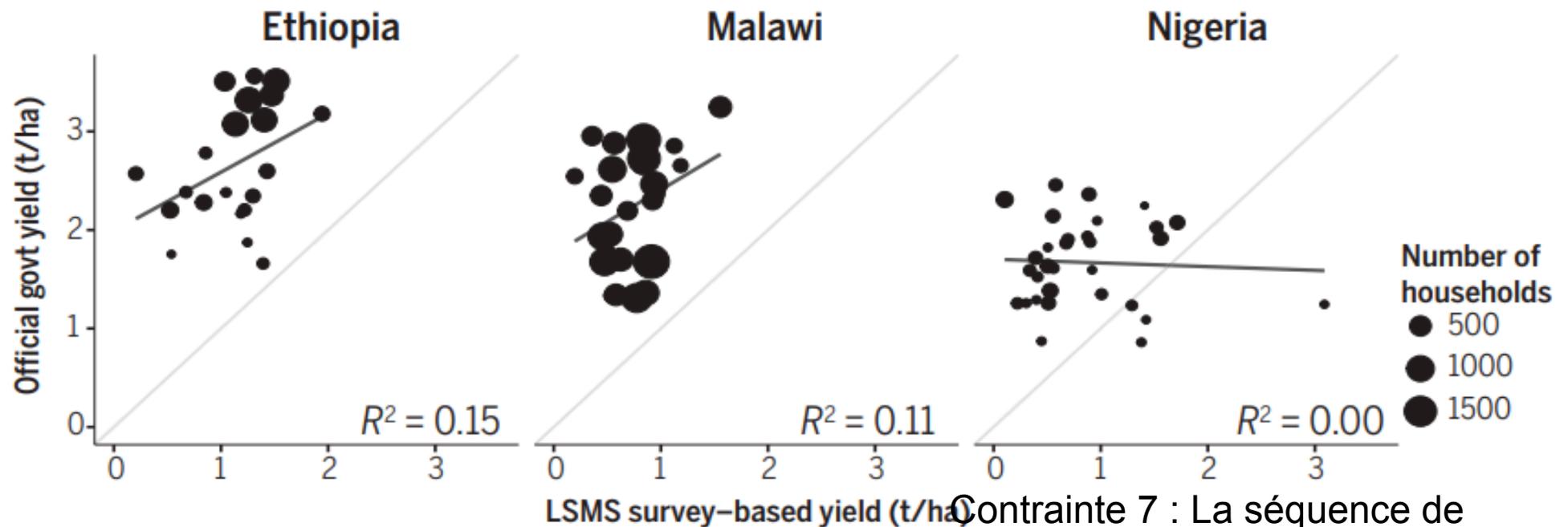
<sup>1</sup>L. Kondmann, et Zhu, X. X. "Measuring Changes in Poverty with Deep Learning and Satellite Images", 2020 ICLR Practical ML for Developing Countries Workshop

<sup>2</sup>C. Yeh, A. Perez, A. Driscoll et al. "Using publicly available satellite imagery and deep learning to understand economic well-being in Africa." 2020 Nat Commun 11.



# Poverty Evolution Prediction: Naive Approach : Limited Results

Burke et al, 2021<sup>1</sup>



Contrainte 7 : La séquence de pauvreté est altérée

<sup>1</sup>Marshall Burke, Anne Driscoll, David B. Lobell, Stefano Ermon. "Using satellite imagery to understand and promote sustainable development". *2021 Science*

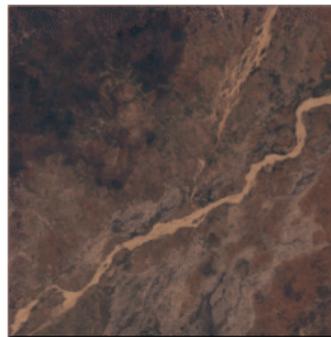
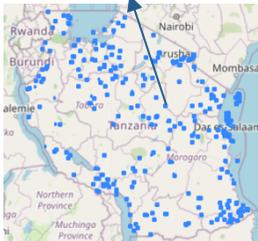
# Poverty Evolution Prediction using Transformers



Input: Series of images

2000-2020

2000



Spatio-Temporal Model

Output 1 : series of poverty indicators

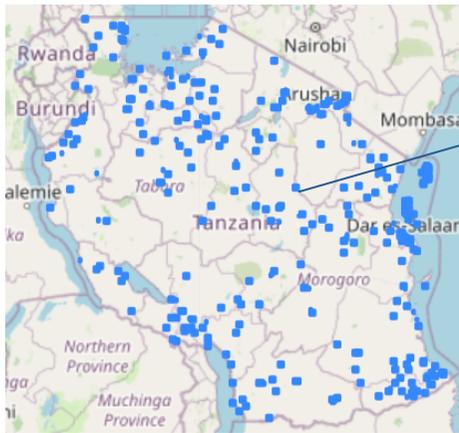
2000  
2.359 \$/d/p

Increase /  
Decrease /  
Stable

Output 2 : classes of evolution trend

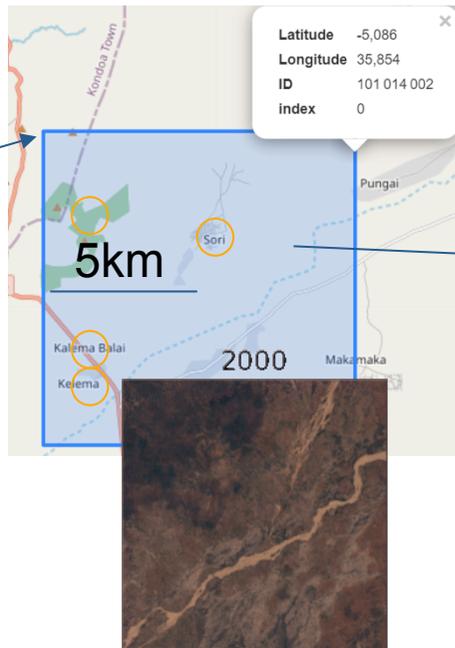
- Constraint 1** : Avg Spatial Resolution (30m)
- Constraint 2** : Avg Temporal Resolution (1 an)
- Constraint 3** : Sensor Noise in the signal

# Poverty Evolution Prediction: Constraints

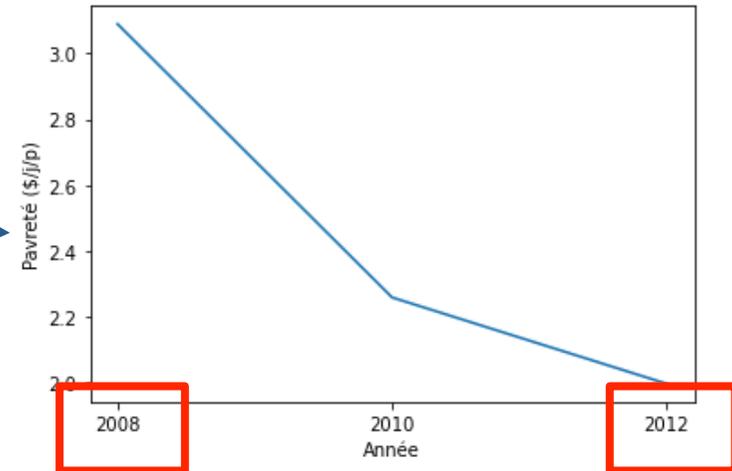


409 points in  
Tanzania  
(LSMS ISA)

**Constraint 4 :**  
**Data Sparsity**



**Constraint 5 :**  
**Randomization of the  
poverty indicator for  
anonymization**



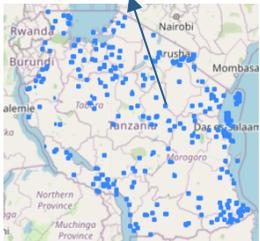
**Constraint 6 :**  
**Relatively Small time  
window (few number  
of points)**



# Poverty Evolution Prediction: Need for intermediate training

Input: Series of images

2000-2020



Spatio-Temporal Model

Output 1 : series of poverty indicators

2000

2.359 \$/d/p

Increase /Decrease /Stable

Output 2 : classes of evolution trend

- Much more data available
- More accurate locations
- Longer time series
- *Presumably correlated with poverty evolution*

Intermediate Training on NLI  
(proxy) Night Light Intensity

# Challenges in Poverty Evolution Prediction from Satellite Images



- 1. Neural architecture search and tuning is difficult, especially when training data is scarce in time/space**
- 2. Finding the right proxy for pre-training is crucial**
- 3. Finding the right resolution and aggregation level in time/space is challenging**
- 4. Uncertainties are at every stage of the process**

# Outline

## I- Challenges in ML & data science pipelines

- Building the pipelines
- Preparing the data

## II- Current Projects

- Combining satellite imagery and socio-economic indicators to estimate poverty evolution in Africa
- **Finding sustainable transition pathways in Nordeste**



# Finding Sustainable Pathways

<https://ideal.ufpb.br/>



IDEAL International Joint Lab

## Area of study



## Actors

- Administrators
- Site Managers
- Elected Representatives
- Exploitation owners
- Farmers
- Indigenous representatives



# Finding Sustainable Pathways

<https://ideal.ufpb.br>



IDEAL International Joint Lab

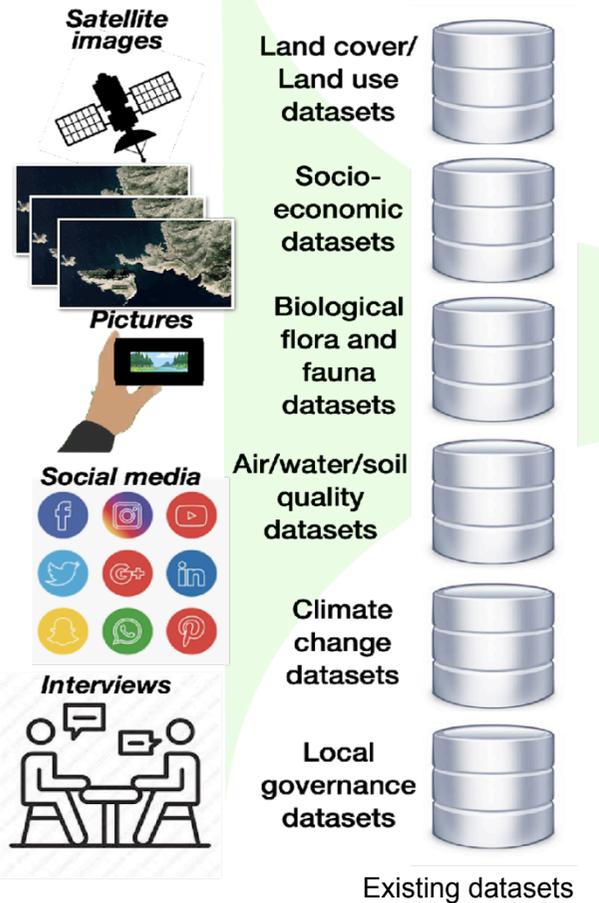
## Area of study



## Actors

- Administrators
- Site Managers
- Elected Representatives
- Exploitation owners
- Farmers
- Indigenous representatives

## Data Science & Engineering

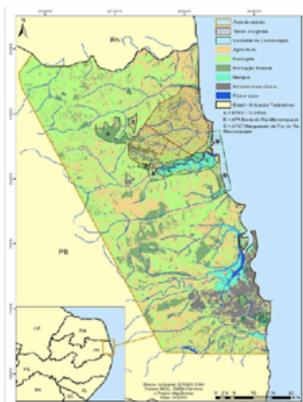


# Finding Sustainable Pathways

<https://ideal.ufpb.br>



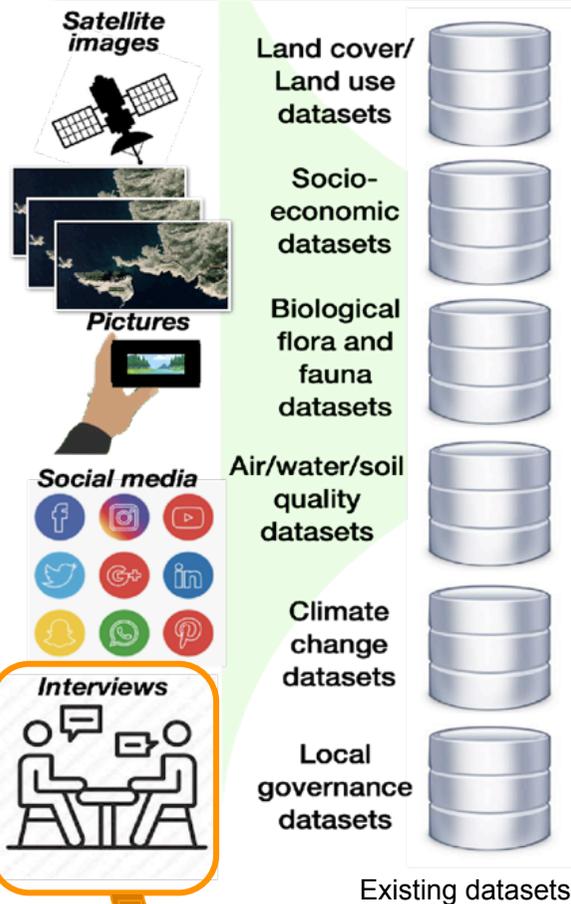
## Area of study



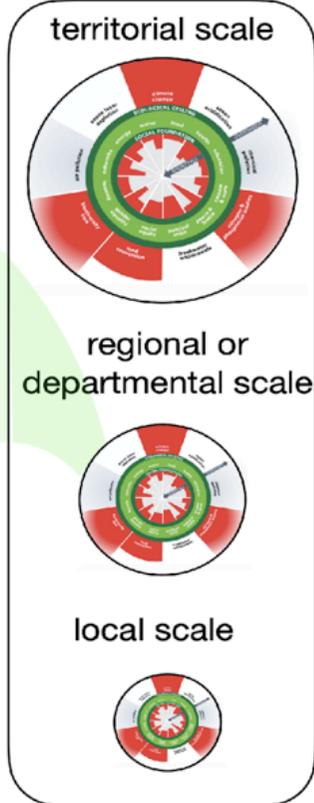
## Actors

- Administrators
- Site Managers
- Elected Representatives
- Exploitation owners
- Farmers
- Indigenous representatives

## Data Science & Engineering

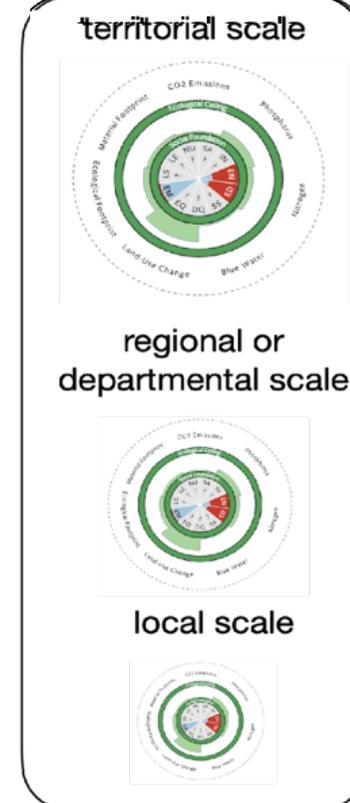


## Coviability Diagnosis



Multiple views as many as actors

## Ideal Target



Multiple co-specified targets

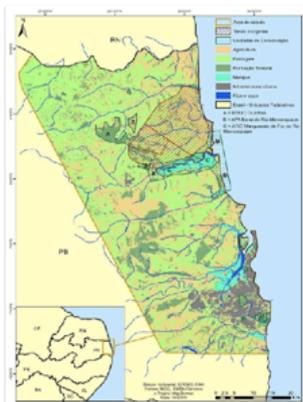
SHS Fieldwork with local populations

# Finding Sustainable Pathways

<https://ideal.ufpb.br/>



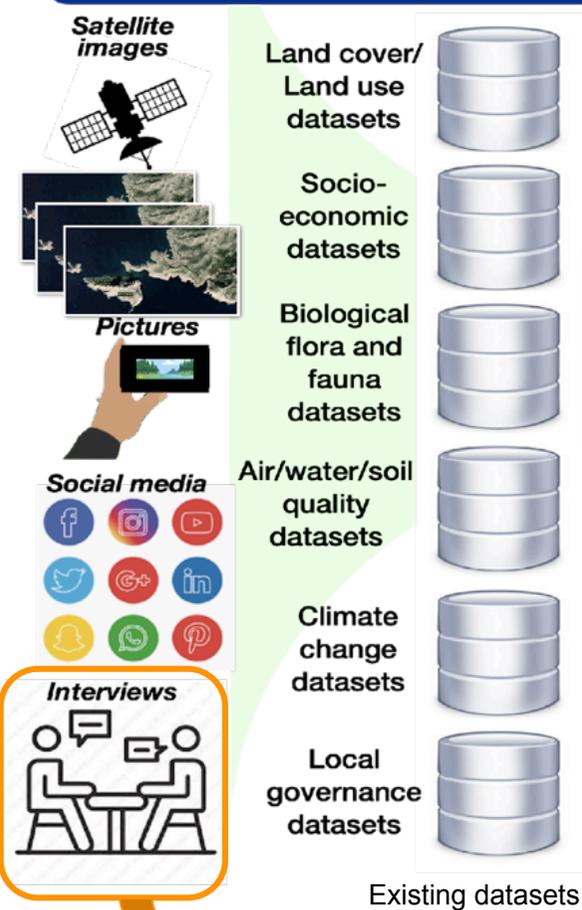
## Area of study



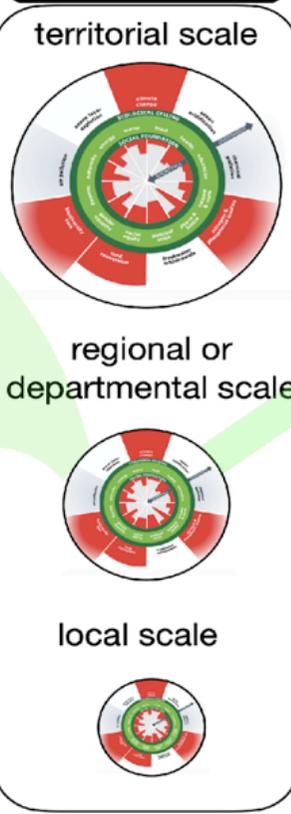
## Actors

- Administrators
- Site Managers
- Elected Representatives
- Exploitation owners
- Farmers
- Indigenous representatives

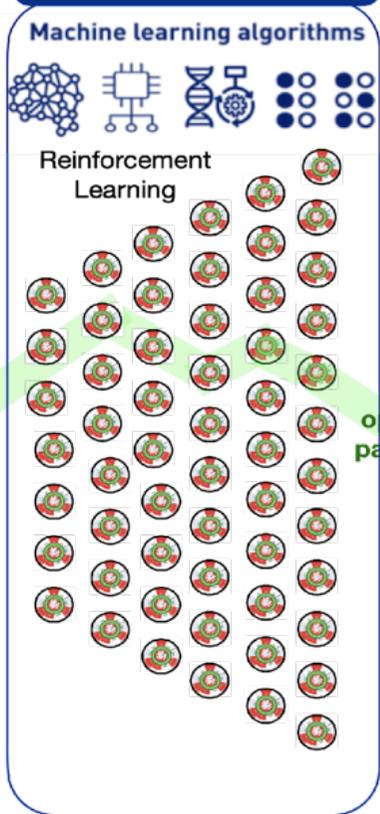
## Data Science & Engineering



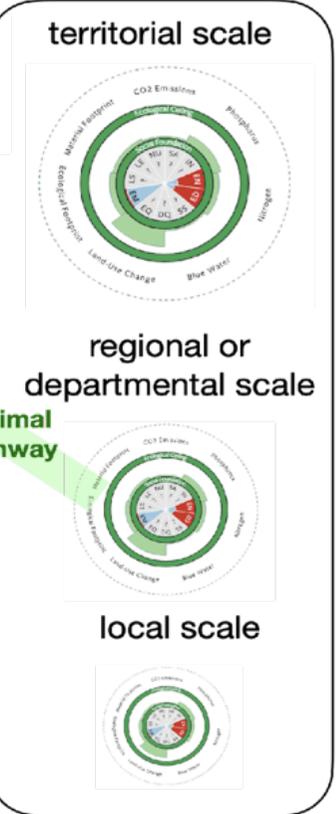
## Coviability Diagnosis



## Artificial Intelligence

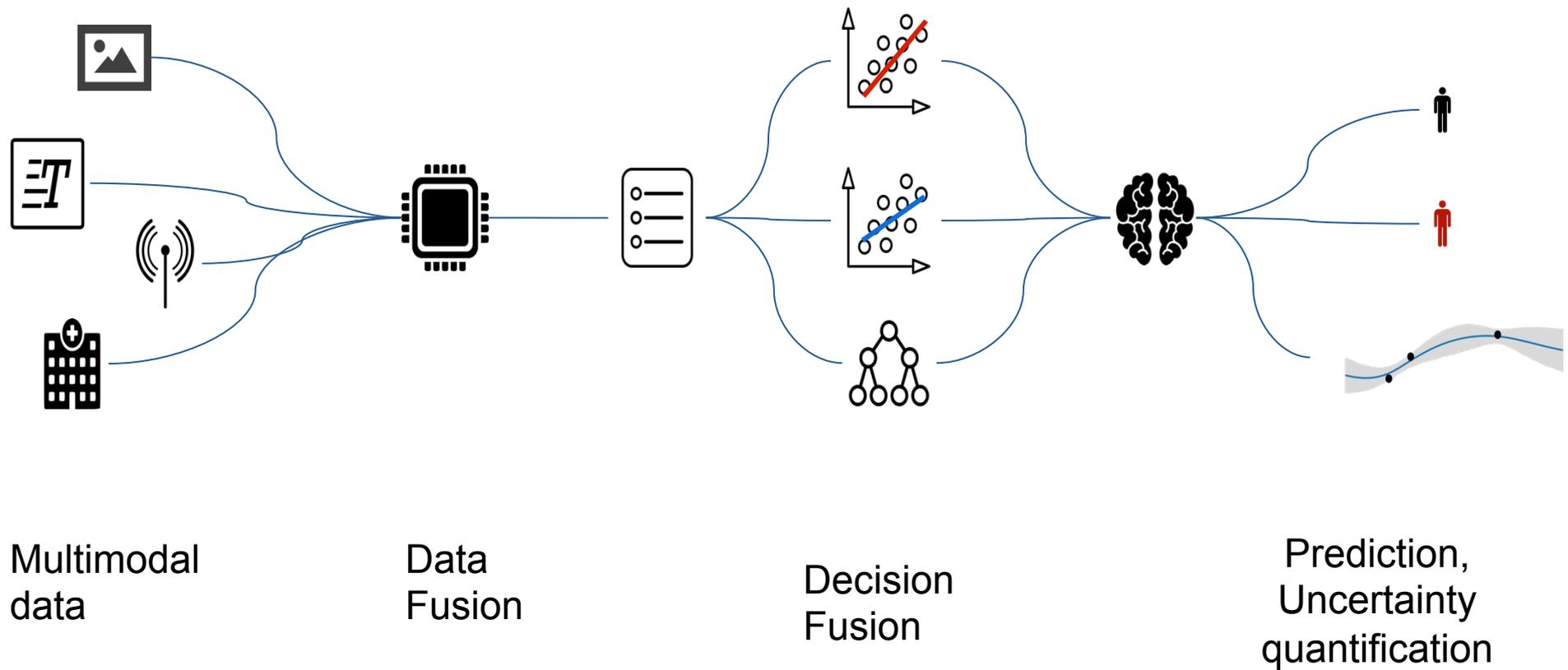


## Ideal Target

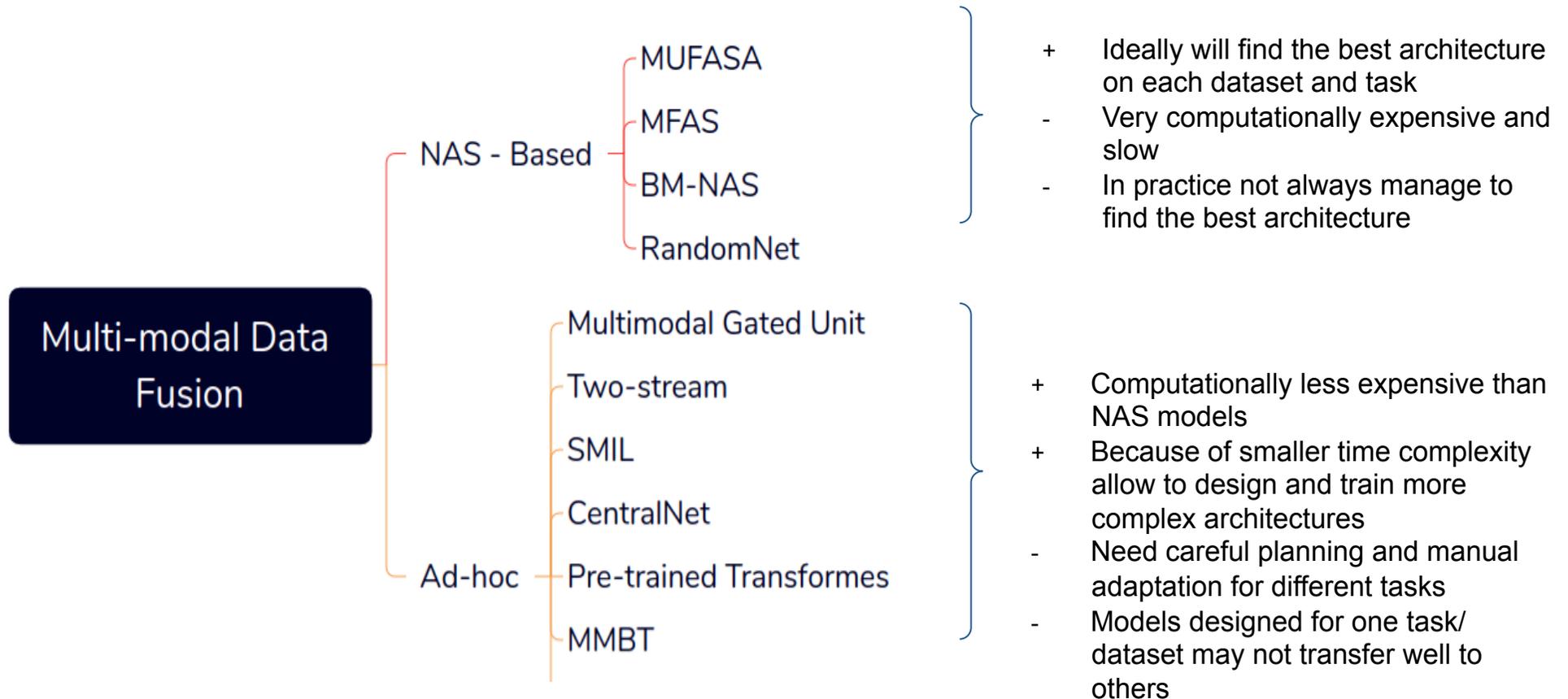


SHS Fieldwork with local populations

# Challenges in Data and Decision Fusion



# Architectures for Multimodal Data Fusion



# Concluding Remarks

- ML and data analytics provide efficient tools to help us answering many SD questions
- But crucially need principled data curation and prep
- ML for SD require **humans in the loop** and **UQ to provide actionable outputs**
- There are many opportunities for:
  - Managing and orchestrating human/machine resources
  - Proposing impactful ML applications for SDGs
  - Revisiting our methods & technologies **to serve SD**

**Thanks!**

