



# Towards a cloud-based computing and analysis framework to process environmental science big data

Eleonora Luppi, Sebastiano Fabio Schifano, Luca Tomassetti  
University of Ferrara, Italy

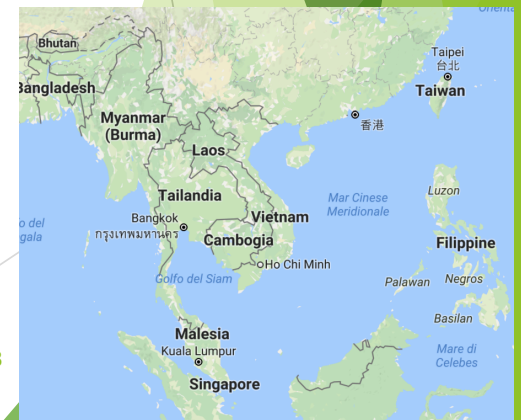
# Introduction

- ▶ Environmental sciences use data coming from several sources:
  - ▶ satellites
  - ▶ large network of sensors installed on the ground or sea-floating stations
  - ▶ devices installed on balloons or aircrafts
- ▶ These networks produce a big amount of data that needs to be appropriately processed and analyzed to extract information useful for scientists to investigate natural phenomenas
- ▶ Needs:
  - ▶ to collect and store huge amount of data together with space and time information
  - ▶ large and powerful computing resources to run analysis and visualization codes

# TORUS Project

Toward Open Resources Using Services

- ▶ Interdisciplinary EU - ERASMUS+ Capacity Building - **TORUS** project, which includes Europe's and South East Asia's partners with a strong expertise in distributed computing and earth and environmental sciences.
- ▶ **TORUS** project aims at making available to environmental scientists a cloud based computing and analysis framework to manage and process big-data:
  - ▶ ability to access clouds to virtualize the computing resources, and knowledge to use software tools to process and analyze data coming from the different sources
  - ▶ data correlation with time and space meta-data information and data storage
  - ▶ high-level data presentation to facilitate management and analysis by user scientists
  - ▶ investigation of high-performance computing integration to boost tasks, also using recent accelerators like GP-GPUs or many-core processors



# TORUS Project

## ► Partners:



## ► Regular Workshops:

- Hanoi (Jan, 2016)
- Ferrara (Jun, 2016)
- Pathumthani (Nov, 2016)
- Brussel (Mar, 2017)
- Ho Chi Min (Sep, 2017)
- Wailalak Univ. (2018)
- Pau (2018)

# TORUS Project Goals

- ▶ Develop research on cloud computing in the environmental sciences and promote its education in the countries of South East Asian partners.
- ▶ Installation of two computation mini-clusters with private cloud:
  - ▶ VNU - Hanoi
  - ▶ AIT - Pathumthani
  
  - ▶ Dual-socket CPUs (>10 cores each)
  - ▶ 64GB of RAM per socket
  - ▶ 2x10Gbits network
  - ▶ ~100TB storage server with SSD cache
  - ▶ Linux based (Debian) OS
- ▶ Setup will be finalized in H2 2017

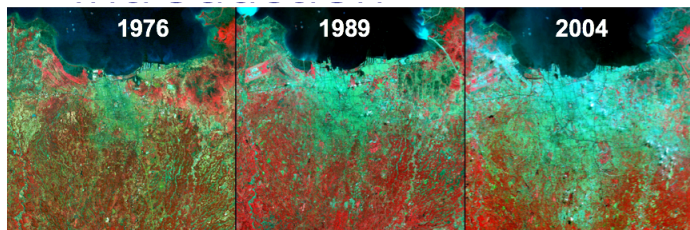
# TORUS Project

- ▶ Several applications in Earth and environmental sciences, geography, satellite image processing are the main focus of the project partners:
  - ▶ AIT: Air Pollution Modeling Applications in Thailand
  - ▶ VNU: Air Pollution Mapping from Space in Vietnam
  - ▶ VUB: Water Resources Management
  - ▶ Toulouse: Statistical approach to geographic applications

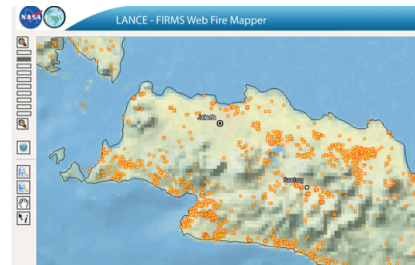
# AIT - Air Pollution Modeling Applications

- ▶ Dr. D. A. Permadi, Prof. N. T. Kim Oanh
- ▶ Asian Institute of Technology, Pathumthani, Thailand

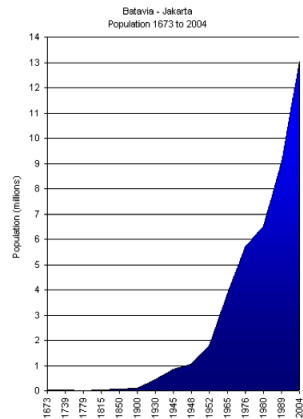
# AIT - Air Pollution Modeling Applications



**Rapid development of the city (Parasati, 2011)**



**Biomass open burning practices**

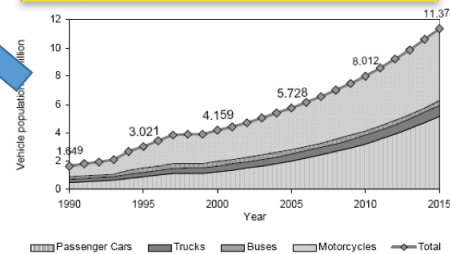


**Air Pollution Problems:**

- Health impacts
- Ecosystem impacts

**Urgent calls for mitigation**

**High rate of motorization**



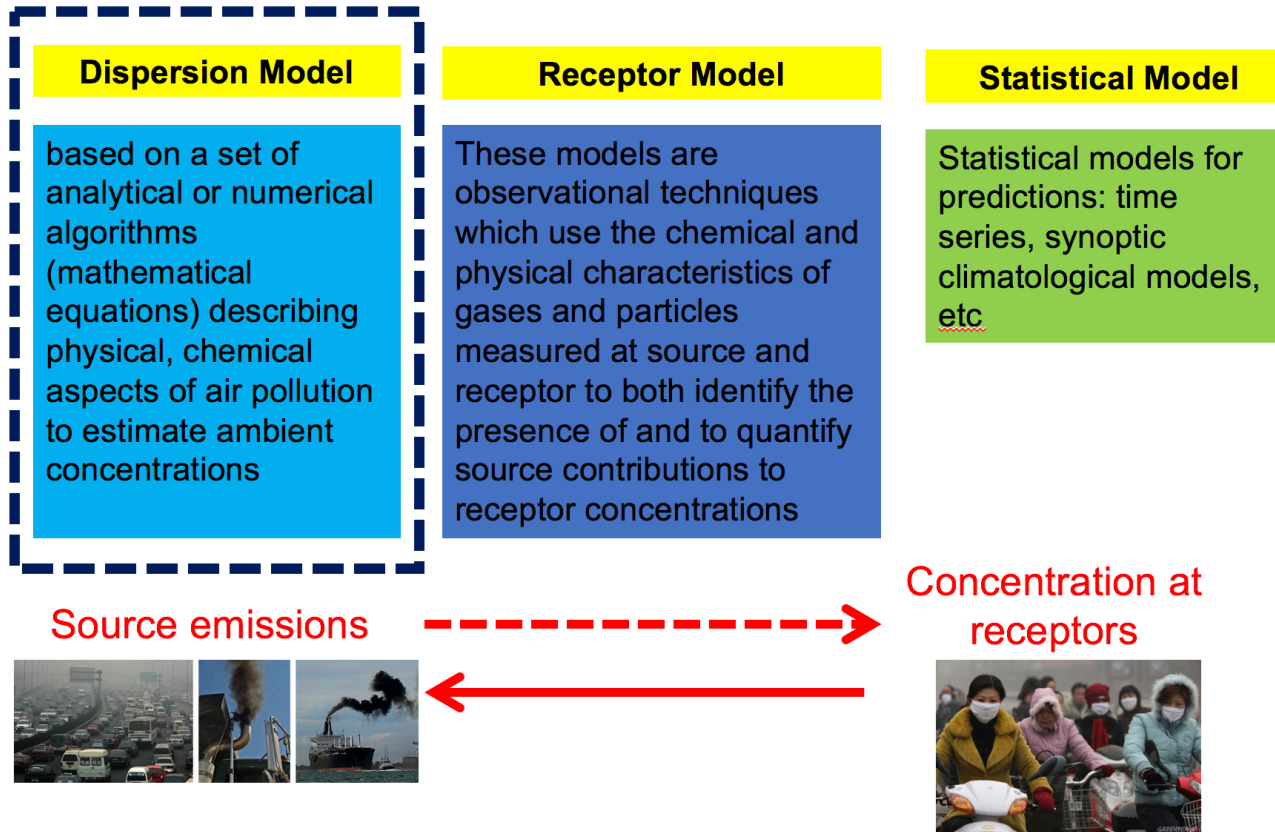
**High population growth**



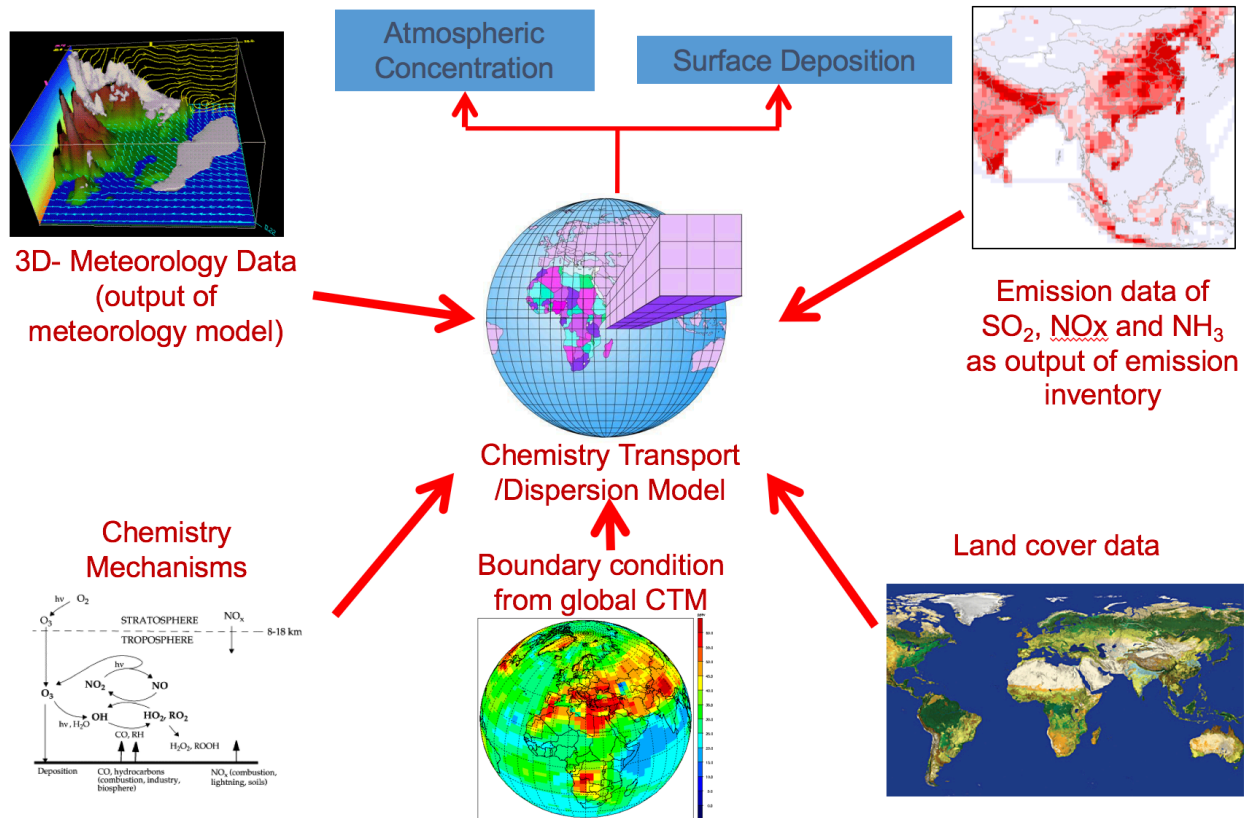
# AIT - Air Pollution Modeling Applications

- ▶ Environment effects are product of complex dynamic system driven by multiple processes (e.g. main processes determining air pollutant dispersion)
  - ▶ Atmospheric transport by mean wind field
  - ▶ Atmospheric turbulent diffusion
  - ▶ Atmospheric chemical and photochemical reactions
  - ▶ Interactions between surface (sea, land) and atmosphere
  - ▶ Wet and dry removal process
- ▶ Modeling tool used to integrate these processes in a systematic approach to assess impacts of different scenarios on environment (causal links)
- ▶ Hindcast, nowcast, and forecast are possible

# AIT - Air Pollution Modeling Applications

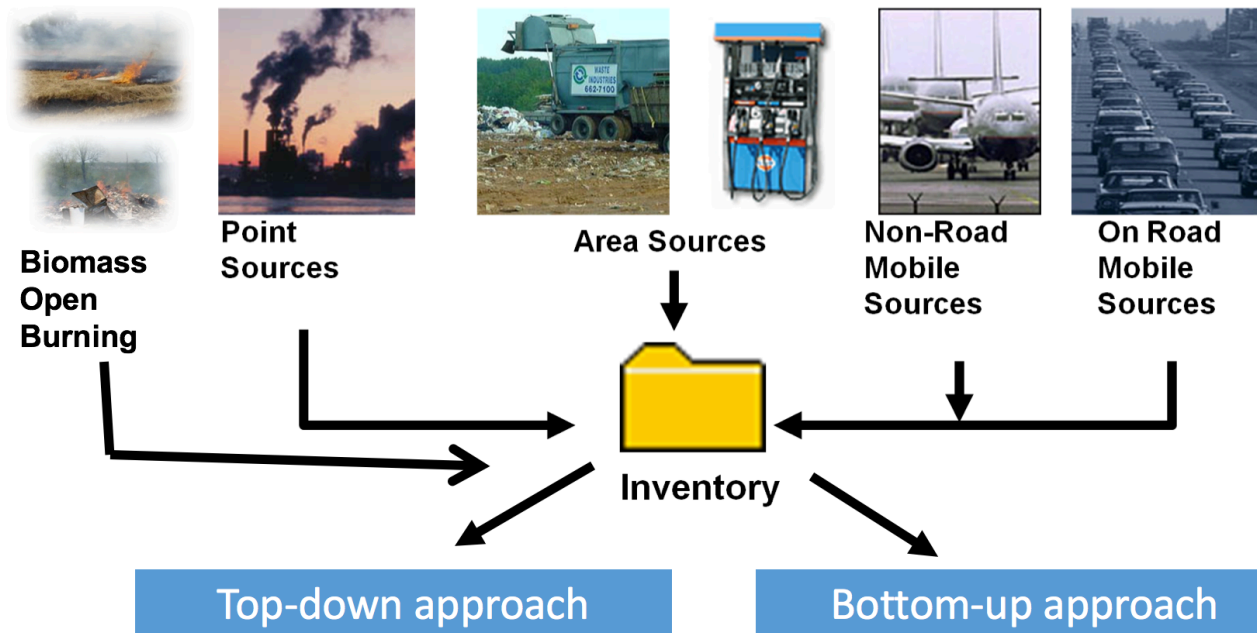


# AIT - Air Pollution Modeling Applications



# AIT - Air Pollution Modeling Applications

**Emission Inventory** - a comprehensive listing by sources of air pollutant emissions in a geographic area during a specific time period



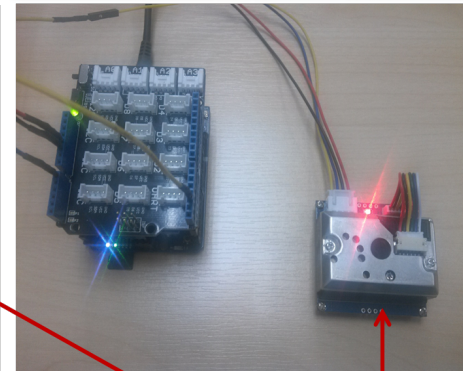
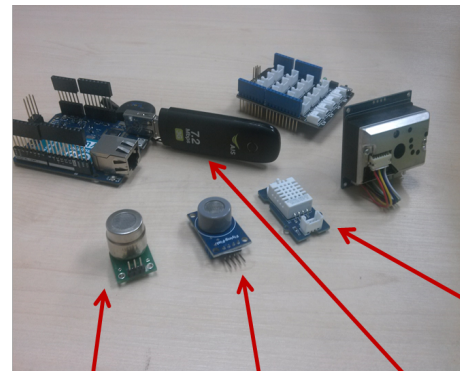
# AIT - Air Pollution Modeling Applications

- ▶ Air quality models require extensive data transfer and storage (input - output of meteorology and chemistry)
  - ▶ Satellite images and metadata from MODIS/VIIR S/LandSat/etc..., albedo, green fraction, land-use, USGS landcover, orography, soil type, and topography
  - ▶ The Emission Database for Global Atmospheric Research (EDGAR),
  - ▶ The Atmospheric Composition Change by the European Network of Excellence (ACCENT),
  - ▶ The Regional Emission inventory in ASia (REAS),
  - ▶ Global Fire Emission Database (GFED)
  - ▶ Inventory for: Ozone, NO<sub>x</sub>, CO<sub>2</sub>, SO<sub>2</sub>, CO, N<sub>2</sub>O, NH<sub>3</sub>, Black-Carbon, Organic-Carbon, CH<sub>4</sub>, PM<sub>2.5</sub>, Total Particulate Matter, and Non-Methan Volatile Organic Compounds
- ▶ High performance computing is important for model simulations
- ▶ Integrated application for data visualization/dissemination through web-based interface can be developed using Cloud services



# AIT - Air Pollution Modeling Applications

- ▶ Network of connected ground sensors
- ▶ PaaS for retrieval & visualization of collected data



MG811 – CO2 sensor

MQ 9 – CO sensor

3G USB

DHT22 – Temp and Humidity sensor

Particulate matter sensor

# AIT - Air Pollution Modeling Applications

## Main Components

- ▶ **Atmospheric modeling system**
- ▶ **Meteorological model (WRF: Weather Research and Forecasting)**
  - ▶ Developed by National Center for Atmospheric Research (NCAR) and National Oceanic and Atmospheric Administration (NOAA): it's a supported *community model* with free and shared resources and distributed development.
  - ▶ 2 dynamical cores:
    - ▶ NMM (Nonhydrostatic Mesoscale Model) for atmospheric physics, real-time and forecast.
    - ▶ **ARW** (Advanced Research WRF) for global and regional climate, coupled-chemistry applications, and idealized simulations.
- ▶ **Chemistry Transport Models (Chimere and CAMx)**
  - ▶ Chimere is a multi-scale model primarily designed to produce **daily forecasts of ozone, aerosols and other pollutants** and make **long-term simulations** for emission control scenarios
  - ▶ Comprehensive Air quality Model with eXtensions (CAMx) is an open-source modeling system for multi-scale integrated assessment of gaseous and particulate air pollution.

# Test and prototyping

- ▶ Collaboration between Unife and AIT to early prototyping and optimization of WRF / air pollution modeling applications in a HPC cluster
- ▶ Use of the Ferrara's cluster
  - ▶ 5 nodes with 2 CPUs, 8 cores per CPU
  - ▶ 2 Infiniband FDR per node
  - ▶ 8 dual GPU Nvidia K80 per node
- ▶ Goal:
  - ▶ optimized run @AIT and @VNU clusters
  - ▶ future exploitation of GPU computing

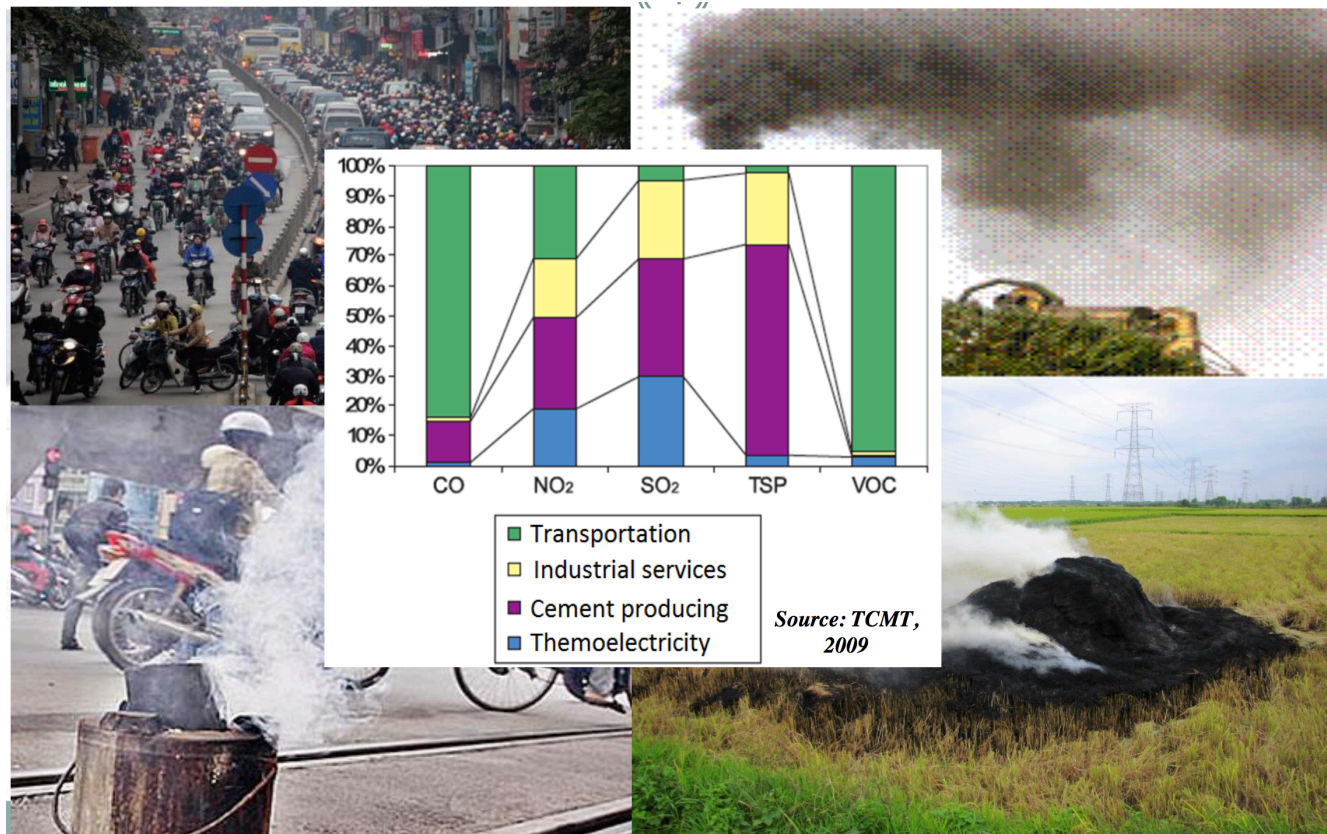




# VNU - Air Pollution Mapping from Space

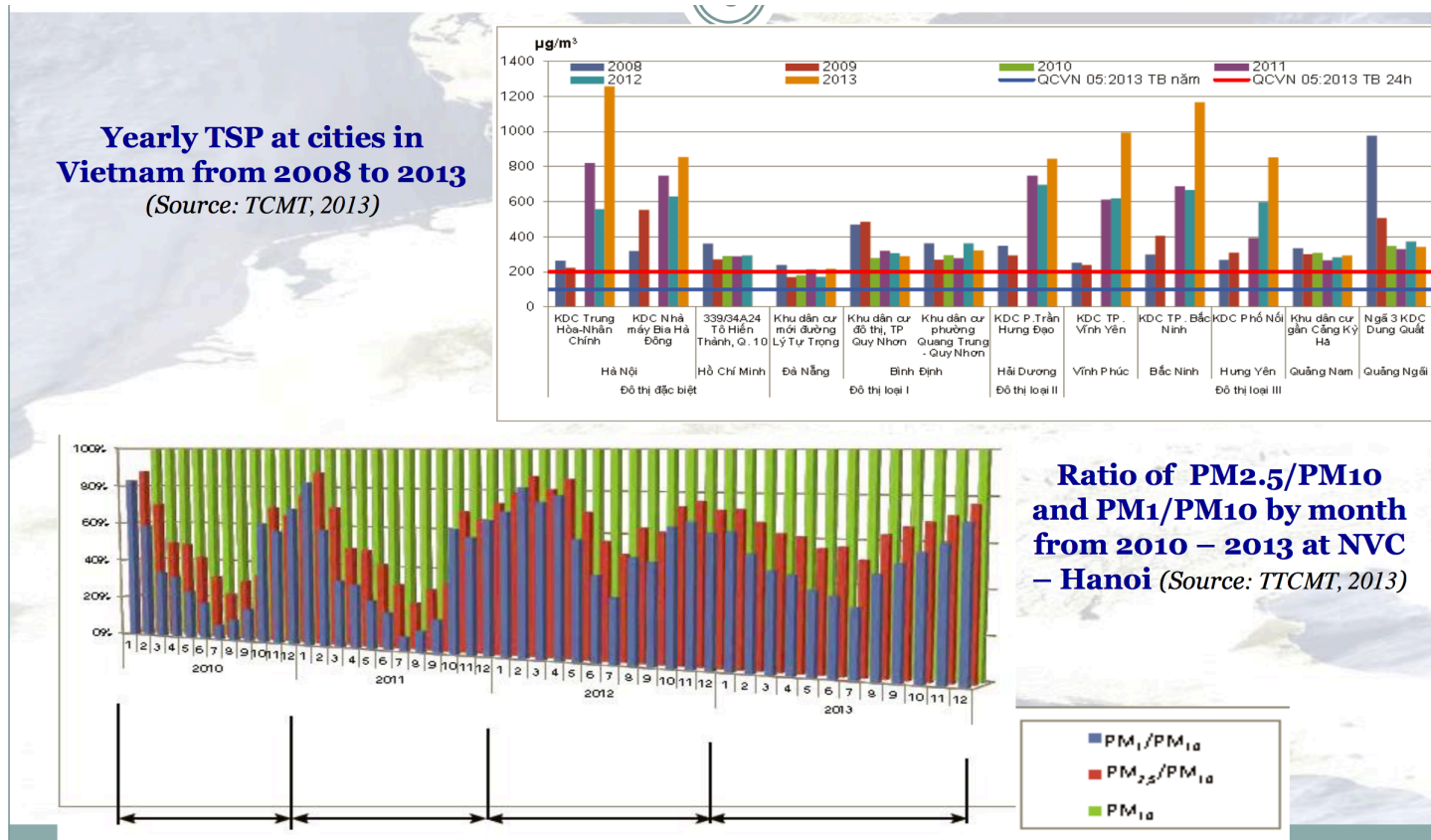
- ▶ NGUYEN THI NHAT THANH, BUI QUANG HUNG, LE THANH HA, NGUYEN NAM HOANG, NGUYEN HAI CHAU, NGUYEN THANH THUY, PHAM VAN HA, LUU VIET HUNG, MAN DUC CHUC, PHAM NGOC HAI, PHAM HUU BANG, LE XUAN THANH PHAN VAN THANH, DO XUAN TU
- ▶ CENTER OF MULTIDISCIPLINARY INTEGRATED TECHNOLOGIES FOR FIELD MONITORING  
UNIVERSITY OF ENGINEERING AND TECHNOLOGY, VIETNAM NATIONAL UNIVERSITY HANOI

# VNU - Air Pollution Mapping from Space



TSP: Total Suspended Particles  
VOC: Volatile Organic Compounds

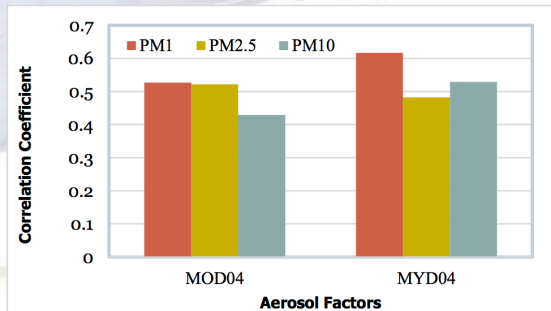
# VNU - Air Pollution Mapping from Space



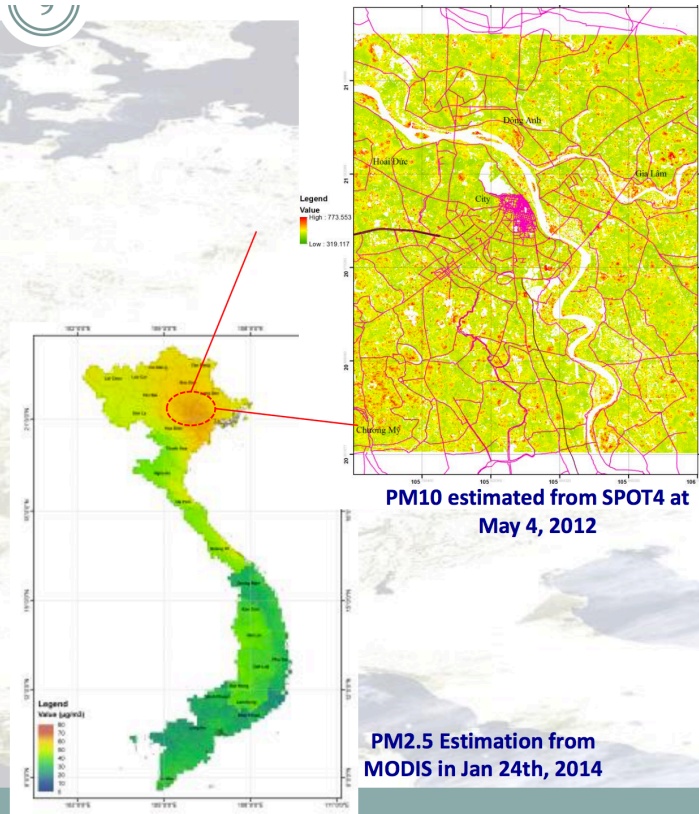
# VNU - Air Pollution Mapping from Space

- Monitor PM based on satellite images

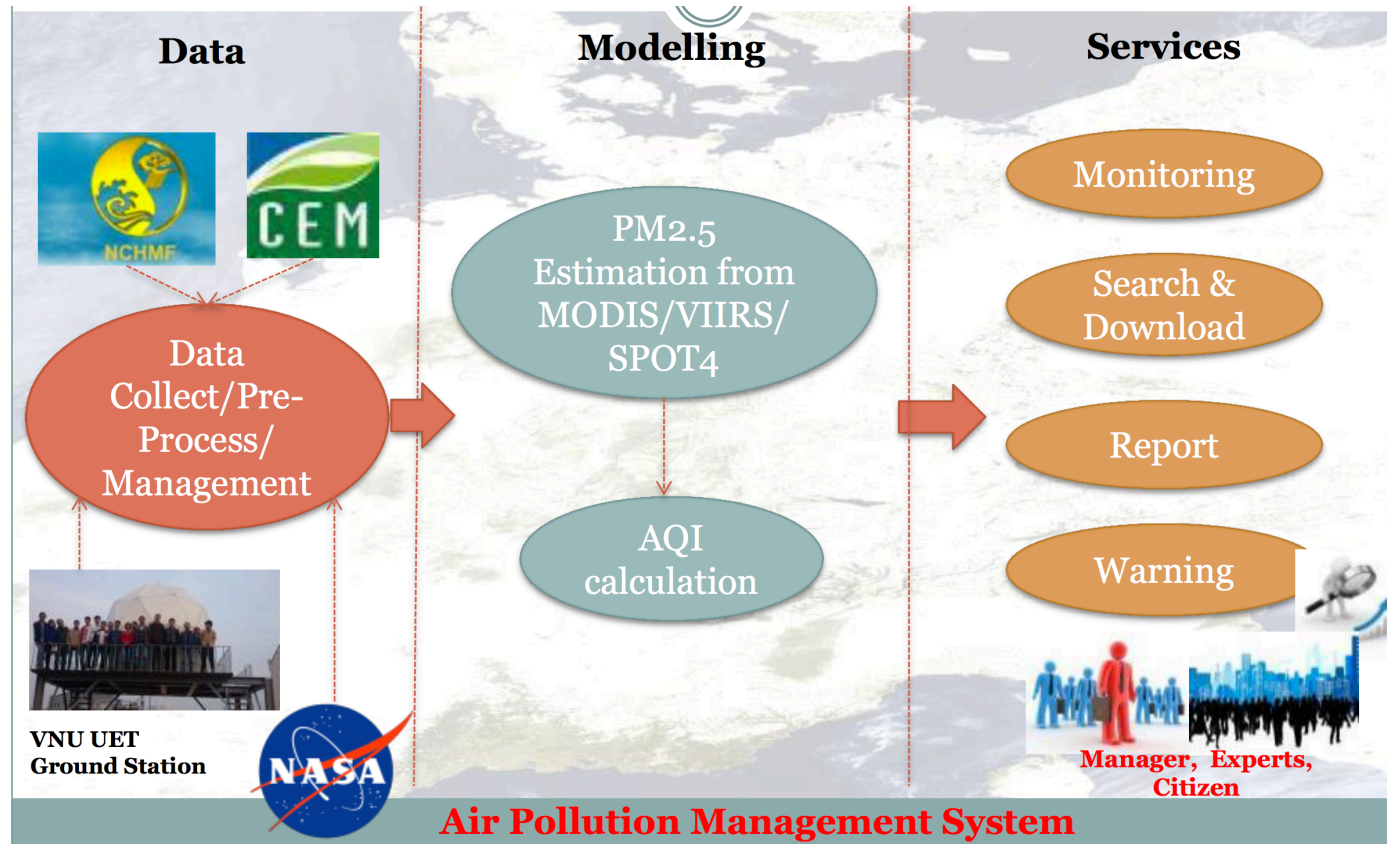
MODIS AOT vs. PM Correlation at PhuTho, Hanoi, Hue, DaNang, 2010-2014



- Provide products at different spatial and temporal scales



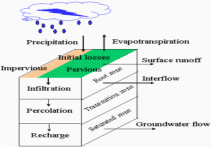

# VNU - Air Pollution Mapping from Space



# VBU - Water resources management


- ▶ Ann van Griensven, Hichem Sahli, Imeshi Weerasinghe
- ▶ Vrije Universiteit Brussel

Hydrological modelling using field and remote sensing data



Model developments using open source software

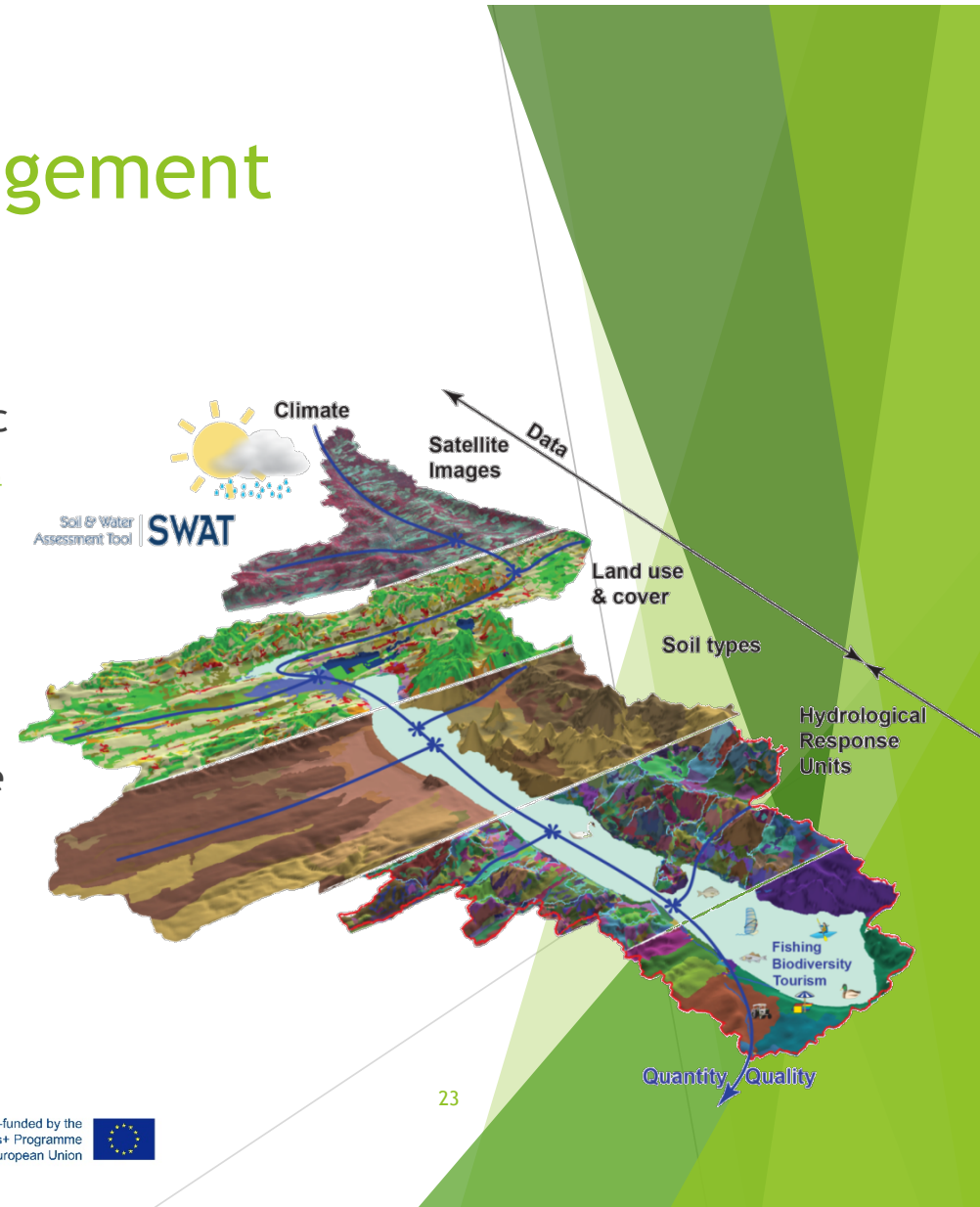
- WETSPA-python
- SWAT
- ...



Vrije Universiteit Brussel

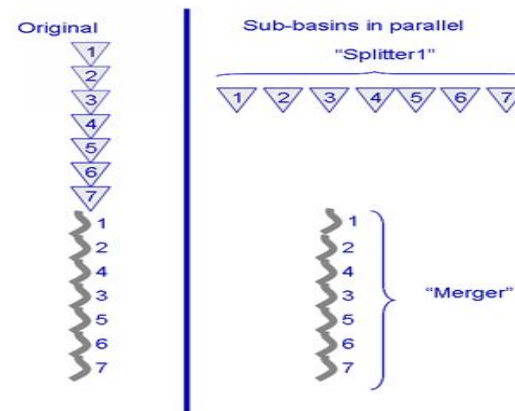
# VBU - Water resources management

- ▶ The Soil and Water Assessment Tool (SWAT) is a public domain model jointly developed by [USDA Agricultural Research Service \(USDA-ARS\)](#) and [Texas A&M AgriLife Research](#), part of The Texas A&M University System.
- ▶ SWAT is a small watershed to river basin-scale model to simulate the quality and quantity of surface and ground water and predict the environmental impact of land use, land management practices, and climate change.
- ▶ SWAT is widely used in assessing soil erosion prevention and control, non-point source pollution control and regional management in watersheds.



# VBU - Water resources management

- ▶ GRID Computing of SWAT
- ▶ SWAT Model Parallelization:
  - ▶ Split large SWAT models at sub-basin level
  - ▶ Compute them separately as independent tasks
  - ▶ Merge individual outputs from each sub-basin and route the outputs through the river network



7 sub-basins, 7 HRU's:	Computation time (seconds)		Number of CPUs	Speedup
<b>Full model ("sequence")</b>	<b>32</b>			
<b>Parallelisation Experiment</b>	<b>Approach I</b>	<b>Approach II</b>		
Splitting	1.2	1.4		
Sub-basin	3.3	5		
Merging	6.3	4.4		
<b>Parallel computing</b>	<b>10.8</b>	<b>10.8</b>	7	2.96

S. Yalew, A. van Griensven, N. Ray, L. Kokoszkiwicz, G.D. Betrie, Distributed computation of large scale SWAT models on the Grid, Environmental Modelling & Software 41 (2013) 223-230



# VBU - Water resources management

- ▶ **Future developments** (community/demand driven)
- ▶ **STANDARDISATION** for
  - ▶ Data exchange, model exchange and data-model exchange
  - ▶ Interoperability
- ▶ **QUALITY CONTROL**
  - ▶ Data models and metadata for observed data and model results
  - ▶ User rating
- ▶ **LIBRARIES & PORTALS**
  - ▶ Repositories for data, models and model applications
  - ▶ Open access

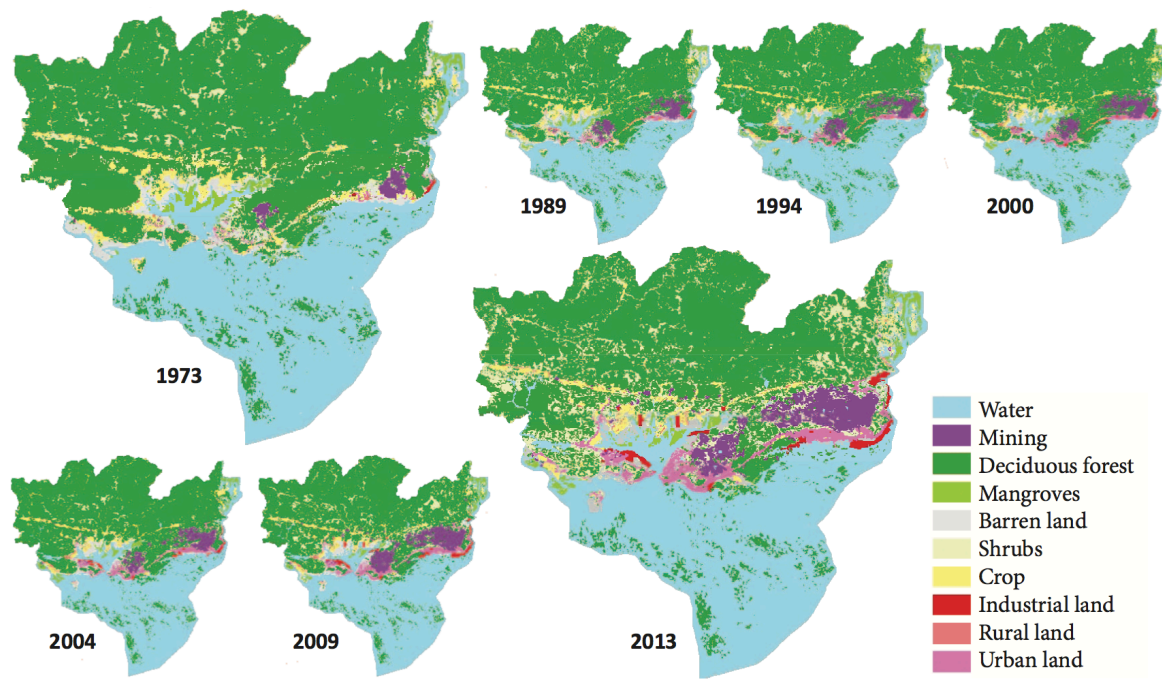
# JJT2 - Statistical approach to Geography

- ▶ Dominique Laffly, Nathalie Hernandez, Florent Devin, Astrid Jourdan, Yannik Le Nir
- ▶ Toulouse University 2 and EISTI Pau

# JJT2 - Statistical approach to Geography

- ▶ Understanding environment changes using statistical analysis of several datasets: satellite images, in-situ measurements, online databases, etc...
  - ▶ Land use change over time (e.g. Vietnam rural to urban areas)
  - ▶ Urban management (e.g. Predict effects of urban changes on quality of life in the city)
  - ▶ East Loven glacier mass balance in Spitsbergen - 78°N, 12°E, Svalbard, Norway (e.g. Predict evolution of glacier size/mass/etc...)
- ▶ Use of Multiple Correspondence Analysis (MCA), Agglomerative Hierarchical Clustering (AHC), Supervised Classification, ...

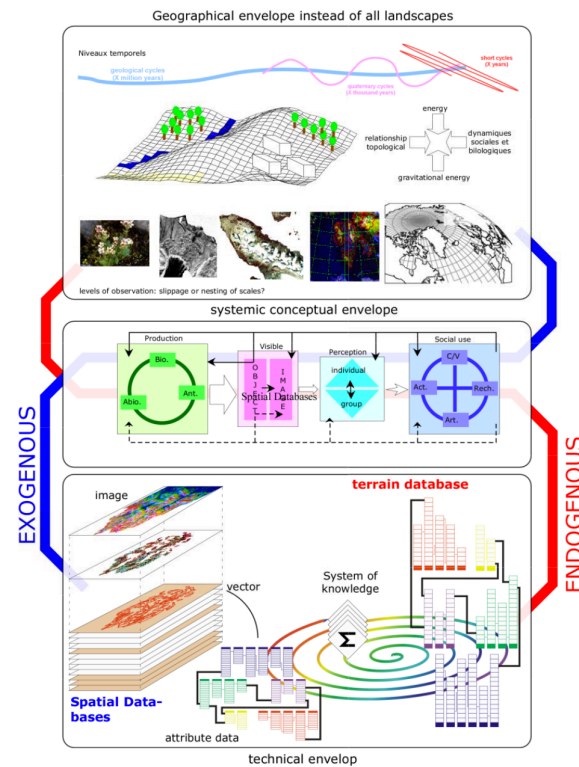
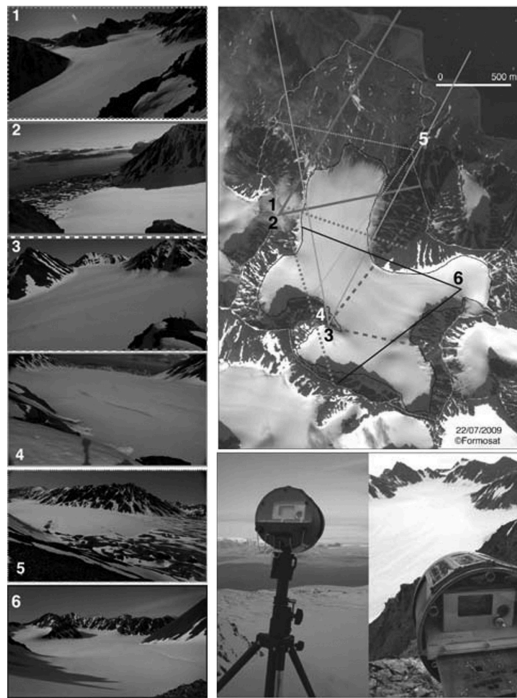
# JJT2 - Statistical approach to Geography



Land use land cover patterns in the Ha Long bay area from 1973 to 2013

# JJT2 - Statistical approach to Geography

Evolution of glacier correlating data coming from satellite images and in-situ monitoring



# JJT2 - Statistical approach to Geography

- ▶ Tools and frameworks used for data collection, storage and statistical analysis:
  - ▶ Spark
  - ▶ R
  - ▶ Scala
  - ▶ MongoDB
  - ▶ Hupi

# Conclusions

- ▶ Collection of requirements is almost finalized
- ▶ Some applications are ready to be “easily” run on private or commercial clouds (standard software on SaaS/PaaS)
- ▶ For other applications (more HPC-oriented) studies are in progress
  - ▶ WRF + Chimere/CAMx
    - ▶ Optimization of use-cases workflow in private cluster and exploration of existing solutions (e.g. WRF4G and its evolution to Clouds)
  - ▶ SWAT
    - ▶ Further develop current solutions and asses performances of runs on grid
    - ▶ Evaluating existing solutions (e.g. SWAT watershed calibration on Azure, ...)