# Rapid Analytics and Model Prototyping

## Akin Kazakçi

MdC / Mines ParisTech

CGS

## Alexandre Gramfort

MdC / Telecom ParisTech

LTCI

## Balázs Kégl

DR / CNRS - UPSaclay

LAL & LRI

# OUTLINE

- Paris-Saclay Center for Data Science

  - the **data science ecosystem**

- Analytics tools

  - **data challenges**

  - **rapid analytics and model prototyping**

# DATA SCIENCE

Design of **automated methods**

to analyze **massive** and **complex** data

to extract useful **information**

# DATA SCIENCE

# ≠

# BIG DATA

We are focusing on **inference**:

**data** ➝ **knowledge**

**Interfacing with infrastructure, security, production**

# UNIVERSITÉ PARIS-SACLAY

## 19 founding partners

# UNIVERSITÉ PARIS-SACLAY

**19** *fondateurs*

**60 000** *étudiants*

**6 000** *doctorants*

**15 000** *étudiants en master*

**8** *Schools*

**11 000** *chercheurs et enseignants-chercheurs*

**300** *laboratoires*

**8 000** *publications /an*

**15 %** *de la recherche publique française*

**10** *départements*

**+ horizontal multi-disciplinary and multi-partner initiatives ("lidex") to create cohesion**

# Paris-Saclay Center for Data Science

A multi-disciplinary initiative to **define, structure, and manage** the **data science ecosystem** at the Université Paris-Saclay

http://www.datascience-paris-saclay.fr/

**250** researchers in **35** laboratories

**Biology & bioinformatics**
IBISC/UEvry
LRI/UPSud
Hepatinov
CESP/UPSud-UVSQ-Inserm
IGM-I2BC/UPSud
MIA/Agro
MIAj-MIG/INRA
LMAS/Centrale

**Chemistry**
EA4041/UPSud

**Earth sciences**
LATMOS/UVSQ
GEOPS/UPSud
IPSL/UVSQ
LSCE/UVSQ
LMD/Polytechnique

**Economy**
LM/ENSAE
RITM/UPSud
LFA/ENSAE

**Neuroscience**
UNICOG/Inserm
U1000/Inserm
NeuroSpin/CEA

**Particle physics astrophysics & cosmology**
LPP/Polytechnique
DMPH/ONERA
CosmoStat/CEA
IAS/UPSud
AIM/CEA
LAL/UPSud

**Machine learning**
LRI/UPSud
LTCI/Telecom
CMLA/Cachan
LS/ENSAE
LIX/Polytechnique
MIA/Agro
CMA/Polytechnique
LSS/Supélec
CVN/Centrale
LMAS/Centrale
DTIM/ONERA
IBISC/UEvry
LIST/CEA

**Visualization**
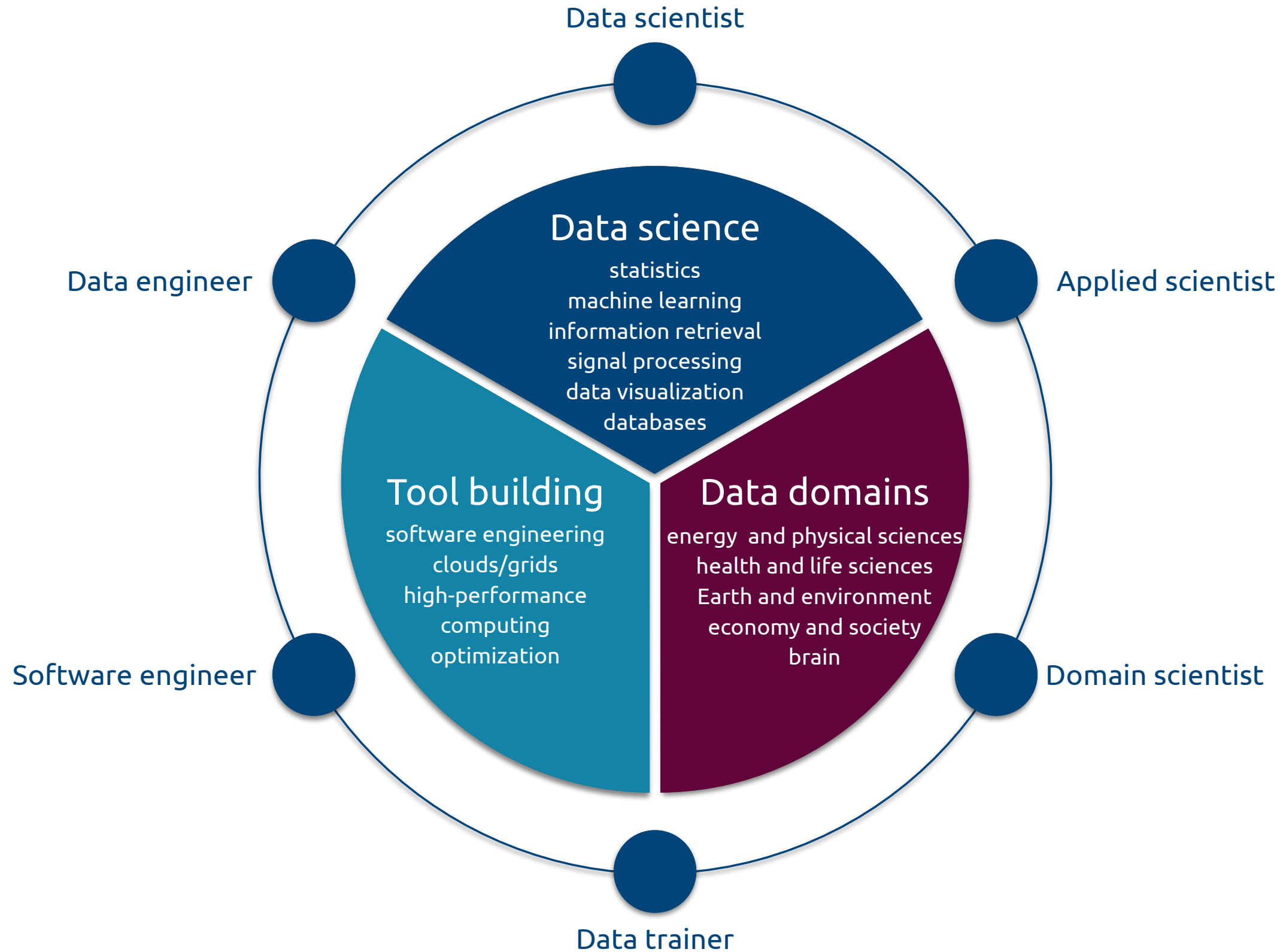INRIA
LIMSI

**Signal processing**
LTCI/Telecom
CMA/Polytechnique
CVN/Centrale
LSS/Supélec
CMLA/Cachan
LIMSI
DTIM/ONERA

**Statistics**
LMO/UPSud
LS/ENSAE
LSS/Supélec
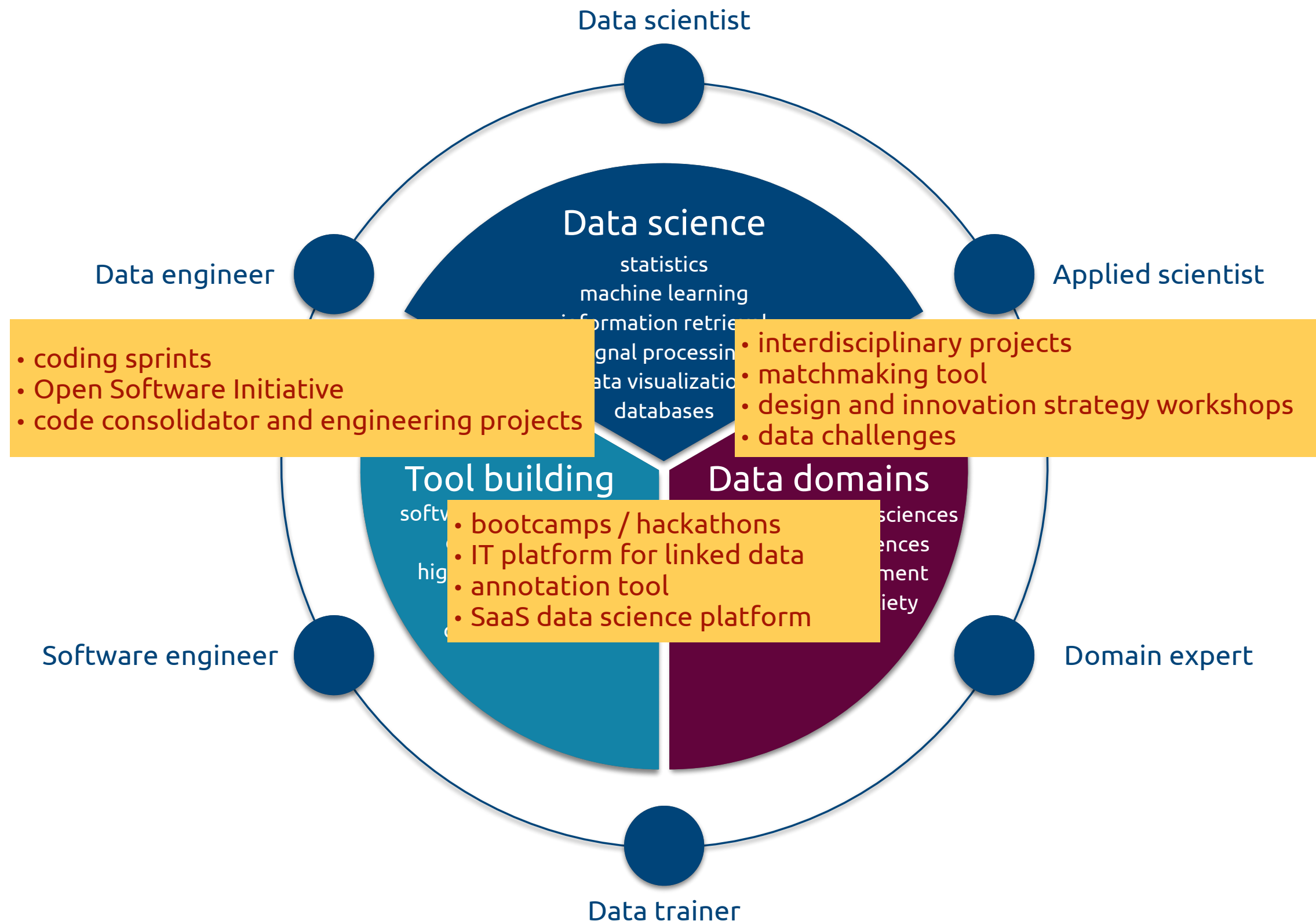CMA/Polytechnique
LMAS/Centrale
MIA/AgroParisTech

cnrs  TELECOM ParisTech  MINES ParisTech  université PARIS-SACLAY  Paris-Saclay Center for Data Science

# THE DATA SCIENCE ECOSYSTEM



Data scientist

Data engineer

Applied scientist

## Data science

statistics
machine learning
information retrieval
signal processing
data visualization
databases

## Tool building

software engineering
clouds/grids
high-performance
computing
optimization

## Data domains

energy and physical sciences
health and life sciences
Earth and environment
economy and society
brain

Software engineer

Domain scientist

Data trainer

8

Paris-Saclay
Center for Data Science

# TOOLS

We are **designing** and **learning to manage**

**tools**

to **accompany** data science projects

with **different needs**

# Tools: landscape to ecosystem



Data scientist

Data engineer

Applied scientist

**Data science**
statistics
machine learning
information retrieval
signal processing
data visualization
databases

- coding sprints
- Open Software Initiative
- code consolidator and engineering projects

- interdisciplinary projects
- matchmaking tool
- design and innovation strategy workshops
- data challenges

**Tool building**

**Data domains**

- bootcamps / hackathons
- IT platform for linked data
- annotation tool
- SaaS data science platform

Software engineer

Domain expert

Data trainer

# TWO ANALYTICS TOOLS

## DATA CHALLENGES

## RAPID ANALYTICS AND MODEL PROTOTYPING

# DATA CHALLENGES

- A **data challenge** is a recently developed unconventional **dissemination** and **communication** tool

  - a scientific or industrial **data producer** arrives with a **well-defined problem** and a corresponding **annotated data set**

  - defines a **quantitative goal**

  - makes the **problem** and part of the data set (the **training set**) **public** on a **dedicated site**

  - **data science experts** then take the public training data and **submit solutions (predictions)** for a **test set** with hidden annotations

  - submissions are **evaluated numerically** using the **quantitative measure**

  - contestants are listed on a **leaderboard**

  - after a **predefined time**, typically a couple of months, the **final results** are revealed and the **winners are awarded**

# DATA CHALLENGES



- The **HiggsML** challenge on **Kaggle**

  - https://www.kaggle.com/c/higgs-boson

# Huge publicity

# CLASSIFICATION FOR DISCOVERY

| # | Δ1w | Team Name ‡ model uploaded * in the money | | Entries | Last Submission UTC (Best − Last Submission) |
|---|-----|-------------------------------------------|----|---------|---------------------------------------------|
| 1 | ↑4 | Gábor Melis ‡ * | 3.80581 | 1?0 | Sun, 14 Sep 2014 09:10:04 (-0h) |
| 2 | ↓1 | Tim Salimans ‡ * | | 57 | Mon, 15 Sep 2014 23:49:02 (-40.6d) |
| 3 | — | nhlx5haze ‡ * | 3.78682 | 254 | Mon, 15 Sep 2014 16:50:01 (-76.3d) |
| 4 | ↑55 | ChoKo Team | 3.77526 | 216 | Mon, 15 Sep 2014 15:21:36 (-42.1h) |
| 5 | ↑23 | cheng chen | 3.77384 | 21 | Mon, 15 Sep 2014 23:29:29 (-0h) |
| 6 | ↓2 | quantify | 3.77086 | 8 | Mon, 15 Sep 2014 16:12:48 (-7.3h) |
| 7 | ↑73 | Stanislav Semenov & Co (HSE Yandex) | 3.76211 | 68 | Mon, 15 Sep 2014 20:19:03 |
| 8 | ↓1 | Luboš Motl's team | 3.76050 | 589 | Mon, 15 Sep 2014 08:38:49 (-1.6h) |
| 9 | ↓1 | Roberto-UCIIIM | 3.75864 | 292 | Mon, 15 Sep 2014 23:44:42 (-44d) |
| 10 | ↑5 | Davut & Josef | 3.75838 | 161 | Mon, 15 Sep 2014 23:24:32 (-4.5d) |
| 990 | ↓65 | sandy | 3.20546 | 5 | Fri, 29 Aug 2014 18:14:30 (-0.7h) |
| 991 | ↓65 | Rem. | | 2 | Mon, 16 Jun 2014 21:53:43 (-30.4h) |
| 📍 | | simple TMVA boosted trees | 3.19956 | | |
| 992 | ↓65 | Xiaohu SUN | | 3 | Tue, 03 Jun 2014 13:14:47 |
| 993 | ↓65 | Pierre Boutaud | 3.19956 | 10 | Fri, 25 Jul 2014 15:25:07 (-30d) |

cnrs · TELECOM ParisTech · MINES ParisTech · université PARIS-SACLAY · Paris-Saclay Center for Data Science

# HUGE PUBLICITY

# SIGNIFICANT IMPROVEMENT OVER THE BASELINE

**yet partially missing the objectives**

# DATA CHALLENGES

- Challenges are useful for

  - generating **visibility** in the **data science community** about **novel application domains**

  - **benchmarking** in a fair way **state-of-the-art techniques** on **well-defined problems**

  - **finding** talented **data scientists**

- Limitations

  - **not** necessary **adapted** to solving **complex** and **open-ended** data science problems in **realistic environments**

  - no direct access to **solutions** and **data scientist**

  - emphasizes **competition**

# We decided to design something better

# Rapid analytics and model prototyping

- Single-day **coding sessions**

  - **20-30** participants

  - **preparation** is similar to challenges

- Goals

  - **focusing** and **motivating** top talents

  - promoting **collaboration**, **speed**, and **efficiency**

  - **solving** (prototyping) **real** problems

# Rapid analytics and model prototyping

# ANALYTICS TOOLS TO PROMOTE COLLABORATION AND INNOVATION

# ANALYTICS TOOLS TO PROMOTE COLLABORATION AND INNOVATION

# RESEARCH
## (BEYOND SOLVING PROBLEMS)

- **Algorithm selection** and **hyperparameter optimization**

  - studying **human problem-solving**

  - combining **human solutions** with **automatic tools**

  - comparing and **tuning hyperparameter optimizers**

  - meta-learning: **embedding** data sets and models, **collaborative optimization**
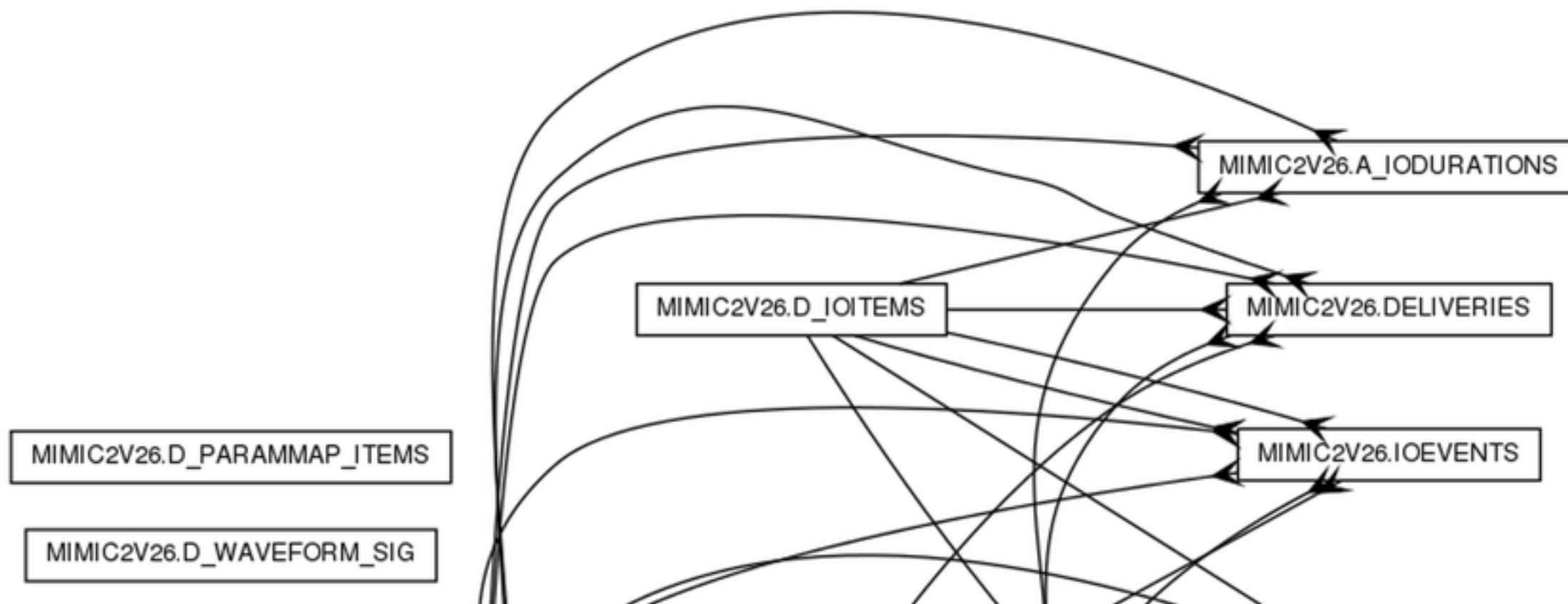
2015 Jan 15 replaying the **HiggsML** challenge

# 2015 Feb 9

# Mortality prediction in septic patients

## MIMIC II V2.6

**Description:**



ERD of the schema

# 2015 Apr 10
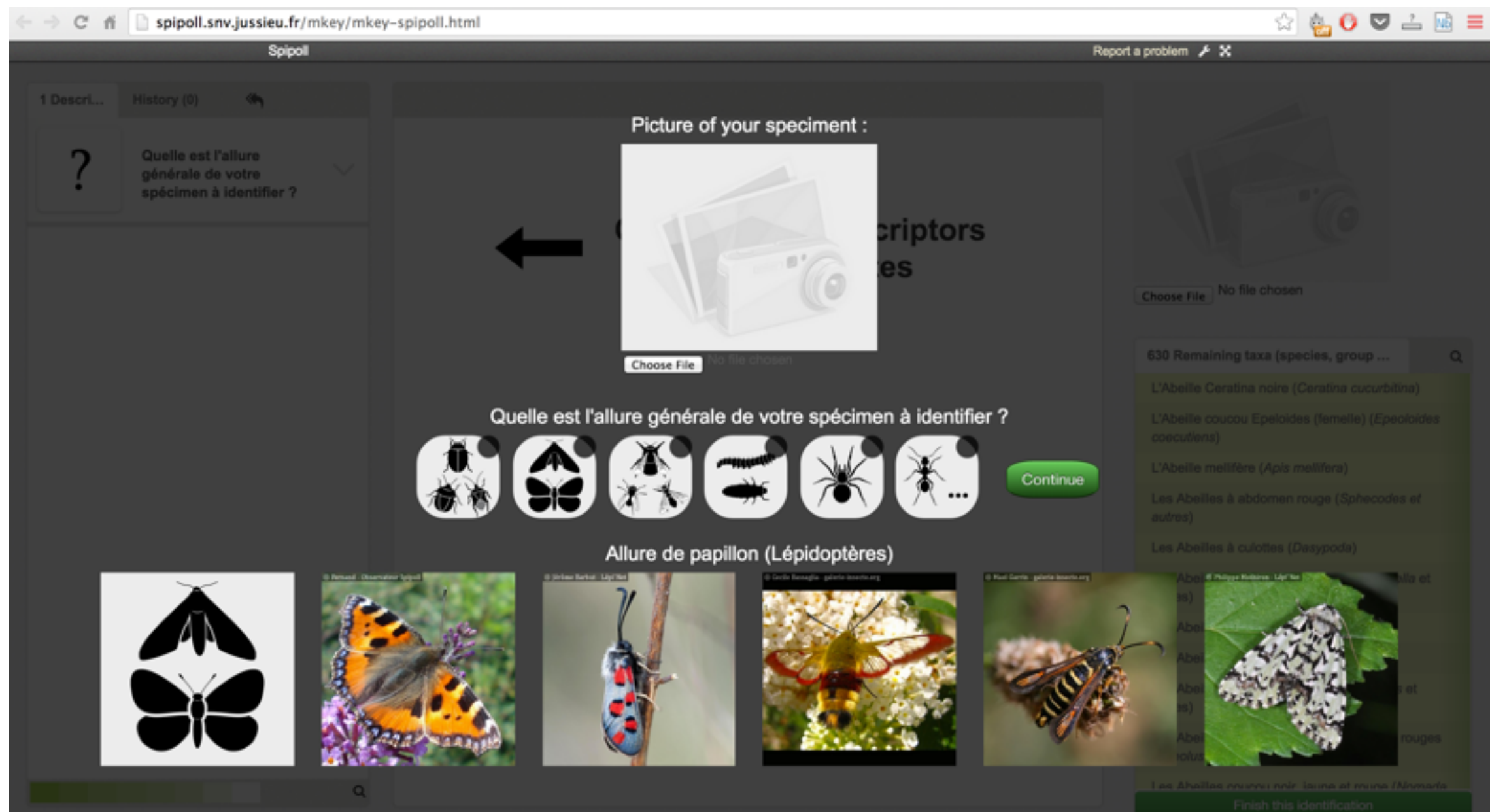# Classifying variable stars

# 2015 May

# Drug identification from spectra

# RAPID ANALYTICS AND MODEL PROTOTYPING

## 2015 June
## Insect classification

# IMPLEMENTING RAMPS IN AN INDUSTRIAL CONTEXT

- Short-term, **ad-hoc teams** assembled for a given task

- **Low-engagement** consulting job

- **Efficient** use of scarce data scientists (i.e., your time)

- Developing and practicing **marketable skills**

  - fast-feedback experimentation is **also useful in research**

- Networking

- Management meta-tools to **track** your performance, to guide your **training**

# Thank you!