# ScienceMesh + JupyterLab
## Collaborative Data Science services in scientific use cases and in business across different fields of study

Marcin Sieprawski
Head of Big Data Lab
Software Mind

Software Mind
part of Ailleron Group

**Software Mind**

**Software house** focused on building **dedicated teams** to extend **product engineering** and **digital transformation** capabilities

**4**

R&D Labs in Poland

HQ in Cracow
Branches: Warsaw,
Rzeszow Bielsko-Biala

**20**

Years of experience
established in 1999

**Branches**
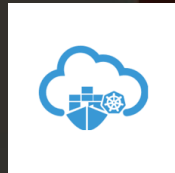and representations

USA, Australia, Singapore,

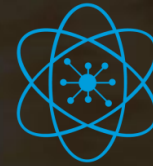# Big enough to scale, small enough to care

# ScienceMesh

CS3MESH4EOSC project

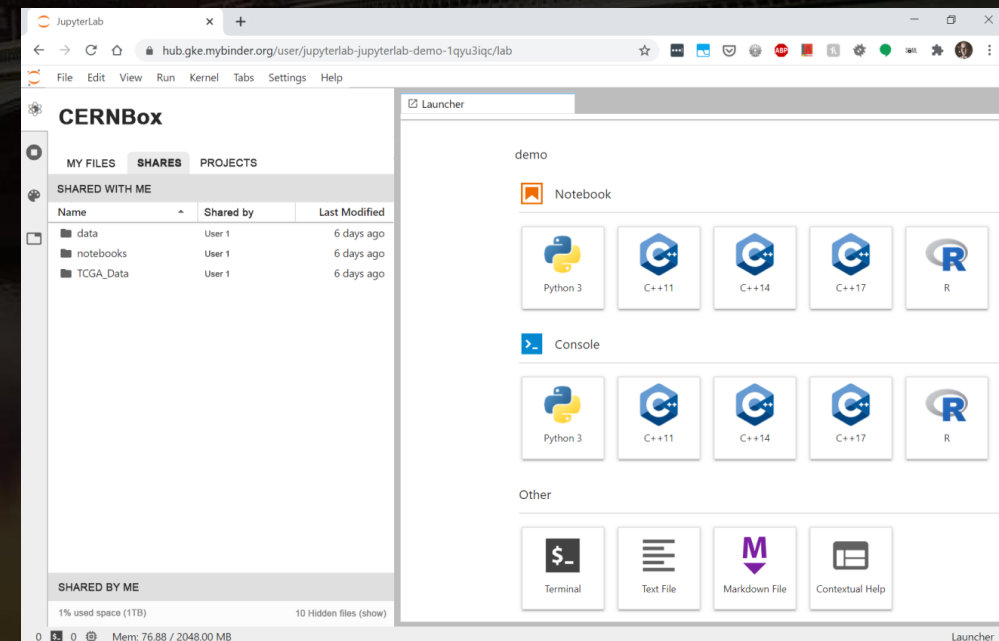 Cloud inter-operability platform

 cloud-native
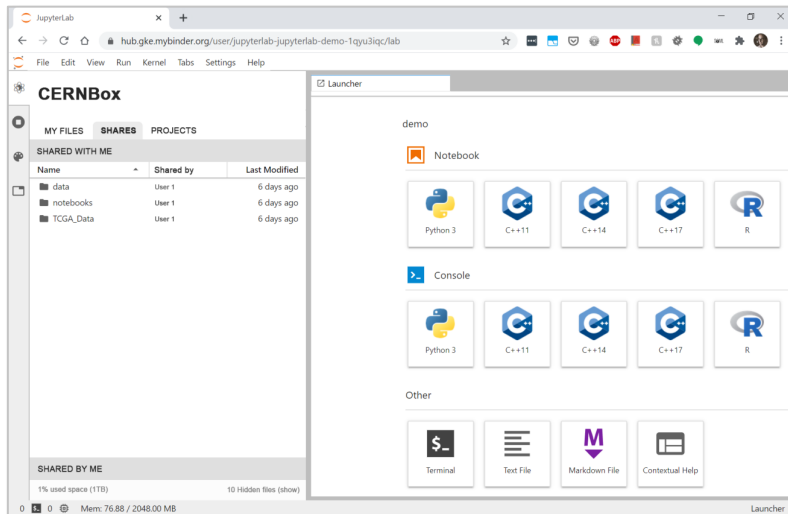
 Distributed Data Science environments

 JupyterLab extension (CS3 APIS)

- Leading tasks on
  - Reference inter-operability platform
  - Distributed Data Science environments

- ScienceMesh Inter-operability platform
  - make cloud storage and application providers inter-operable, via the CS3 APIS

- JupyterLab extension (Cs3Api4Lab)
  - Integration with ScienceMesh IOP (CS3 APIS)
  - replaces the default file manager
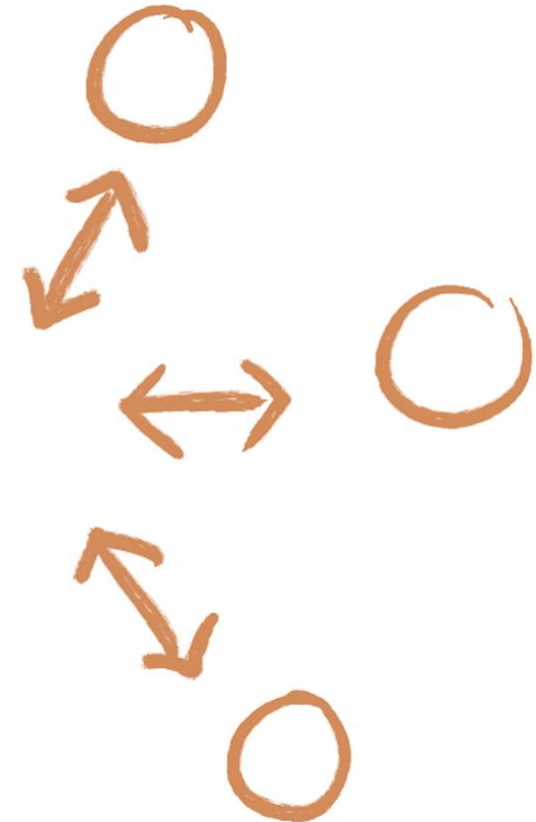  - new UI elements for share functionalities

# JupyterLab extension (**Cs3Api4Lab**)
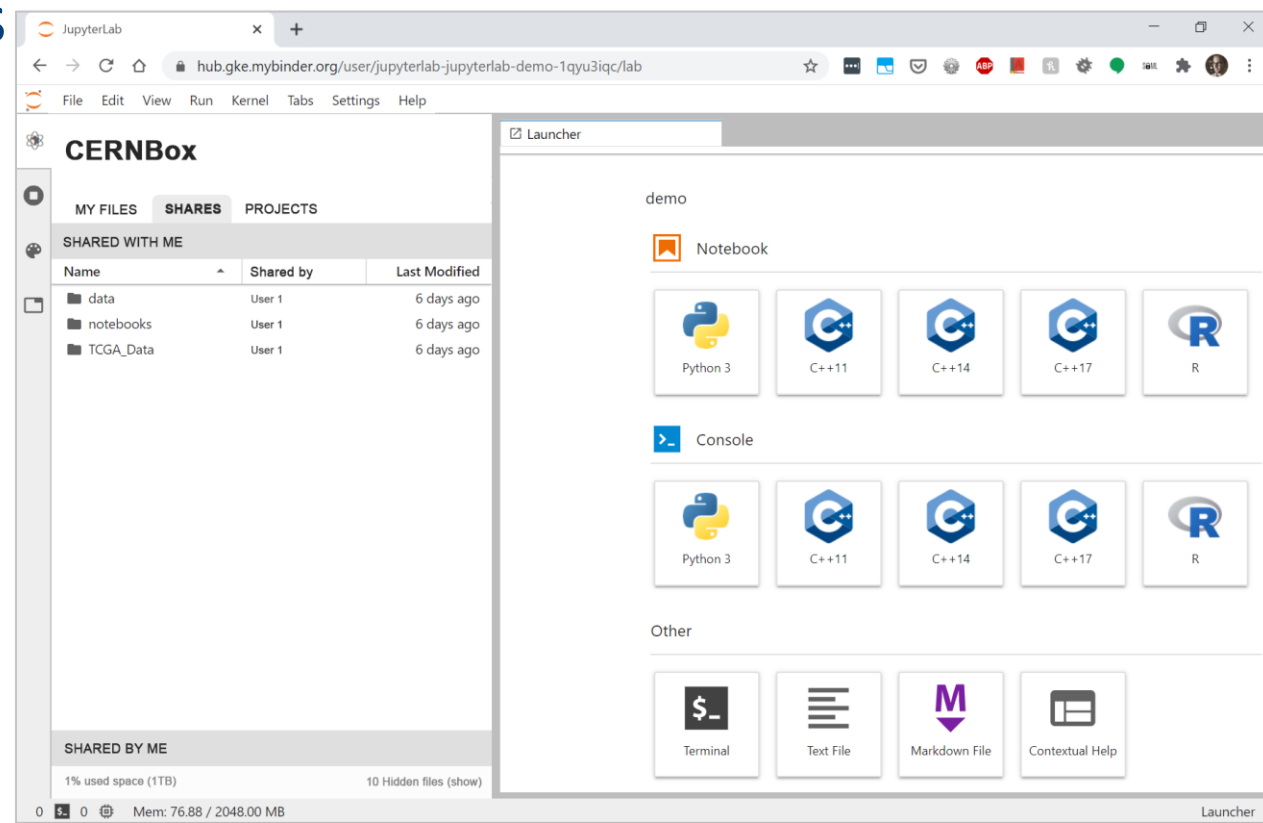
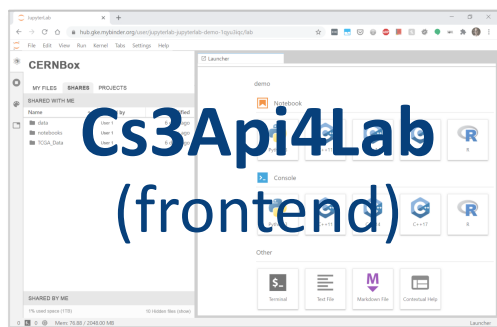\# Integration with ScienceMesh IOP (CS3 APIS)



CS3 APIs

# JupyterLab extension (**Cs3Api4Lab**): Frontend

\# Full client in Lab

\# File browser – share functionalities
   \# Shared by/with tab
   \# Sharing buttons
   \# Entries in the context menus
   \# Pop-up windows: file information and sharing status
   \# Account info

\# File browsing

# JupyterLab extension  (**Cs3Api4Lab**): Backend

- Replaces ContentsManager and Checkpoints
- REST endpoints for integration with the frontend:
  - API for content operations
  - API for checkpoints operations  (todo)
  - API for share operations
- Connecting IOP: gRPC (CS3 APIs)



**Cs3Api4Lab**
(frontend)

REST

**Cs3Api4Lab**
(backend)

CS3 APIs
(gRPC)

Inter-operability
Platform

# Cloud interoperability

Main factors

- Hybrid / multi- cloud
  - Preventing vendor lock-in
  - Cost optimization (private cloud)
  - Managing sensitive data (Privacy by Design)
  - Supporting digital transformation (a process: multiple environments)
- Distributed cloud computing
  - location of cloud-delivered services - part of its definition
  - Important in distributed data science environments
- Main factors of cloud adoption
  - Integration skills (hybrid cloud -> connections and integration points)
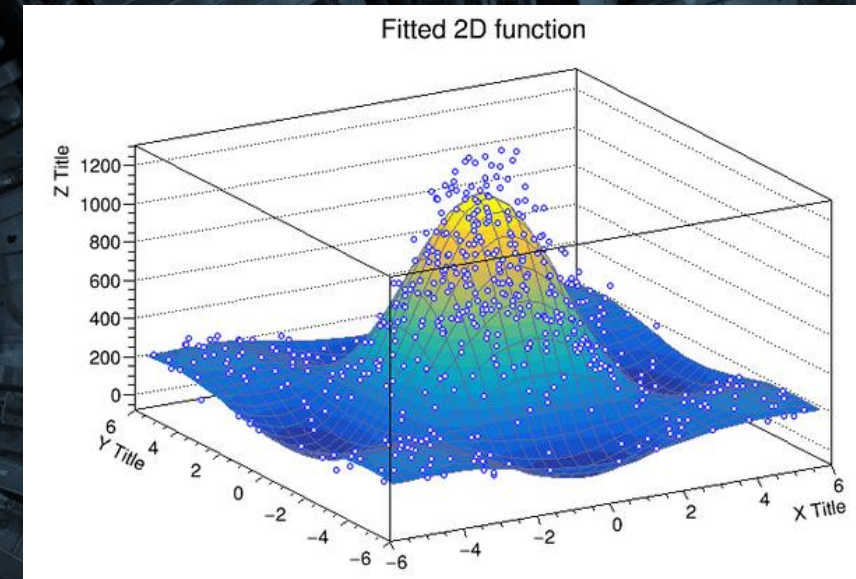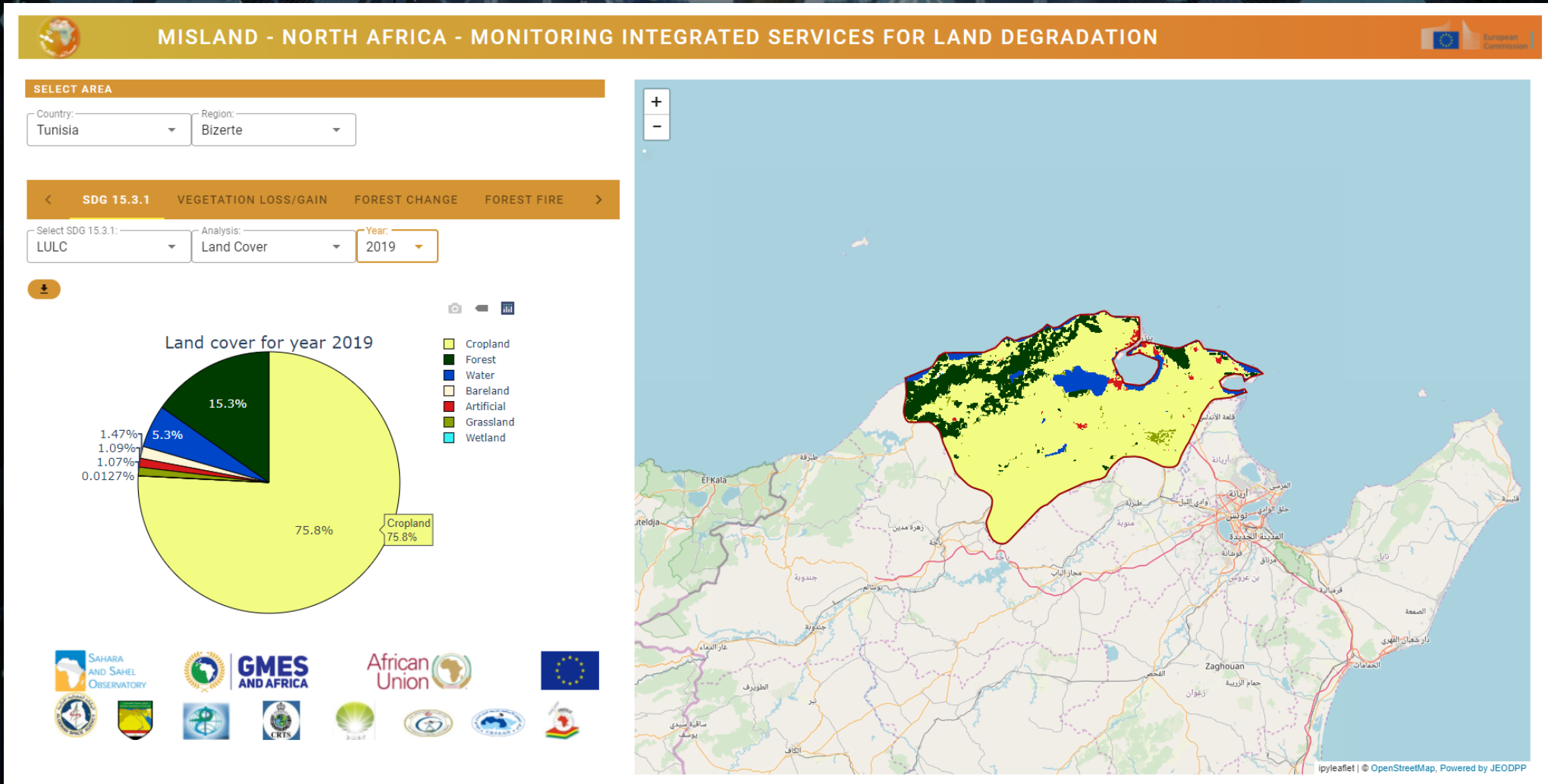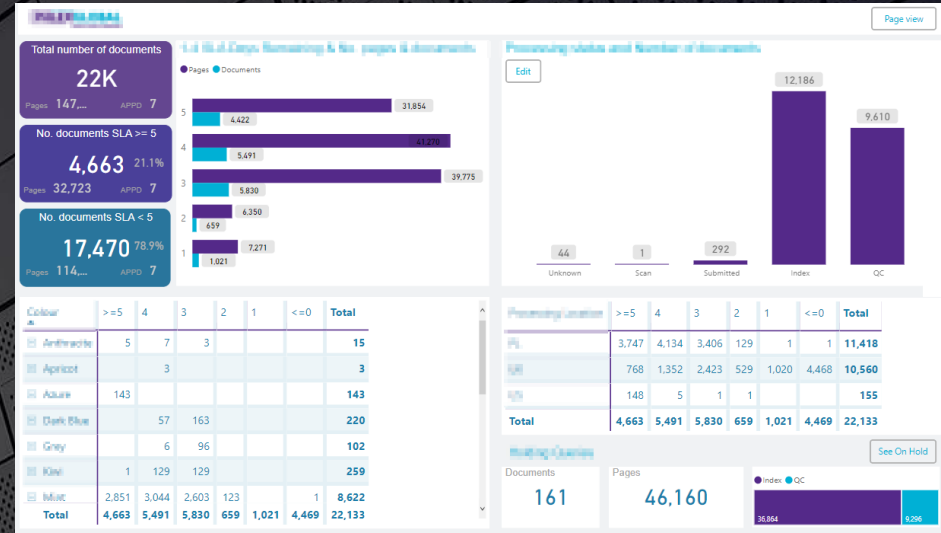  - native-cloud skills
  - Interoperability tools

Software Mind

# ScienceMesh
## CS3MESH4EOSC project

# Collaborative Data Science: Beyond High Energy Physics



Fitted 2D function

- **All scientific disciplines nowadays are data-driven**
  - Data analytics play an increasing role in all types of research
  - Distributed data science environments => all fields of study
  - A more effective collaboration between scientific institutions

- **Business: develop new products in all sectors**
  - Finance, IoT, SmartCities, energy and many others

- **Gartner - Critical Capabilities for Data Science and Machine Learning Platforms**
  - (13 March 2021)
  - **By 2023, 30% of organizations will harness the collective intelligence** of their analytics communities, outperforming competitors that rely solely on centralized analytics or self-service.
  - **By 2024, 70% of enterprises will use cloud and cloud-based AI infrastructure** to operationalize AI, thereby significantly alleviating concerns about integration and upscaling.
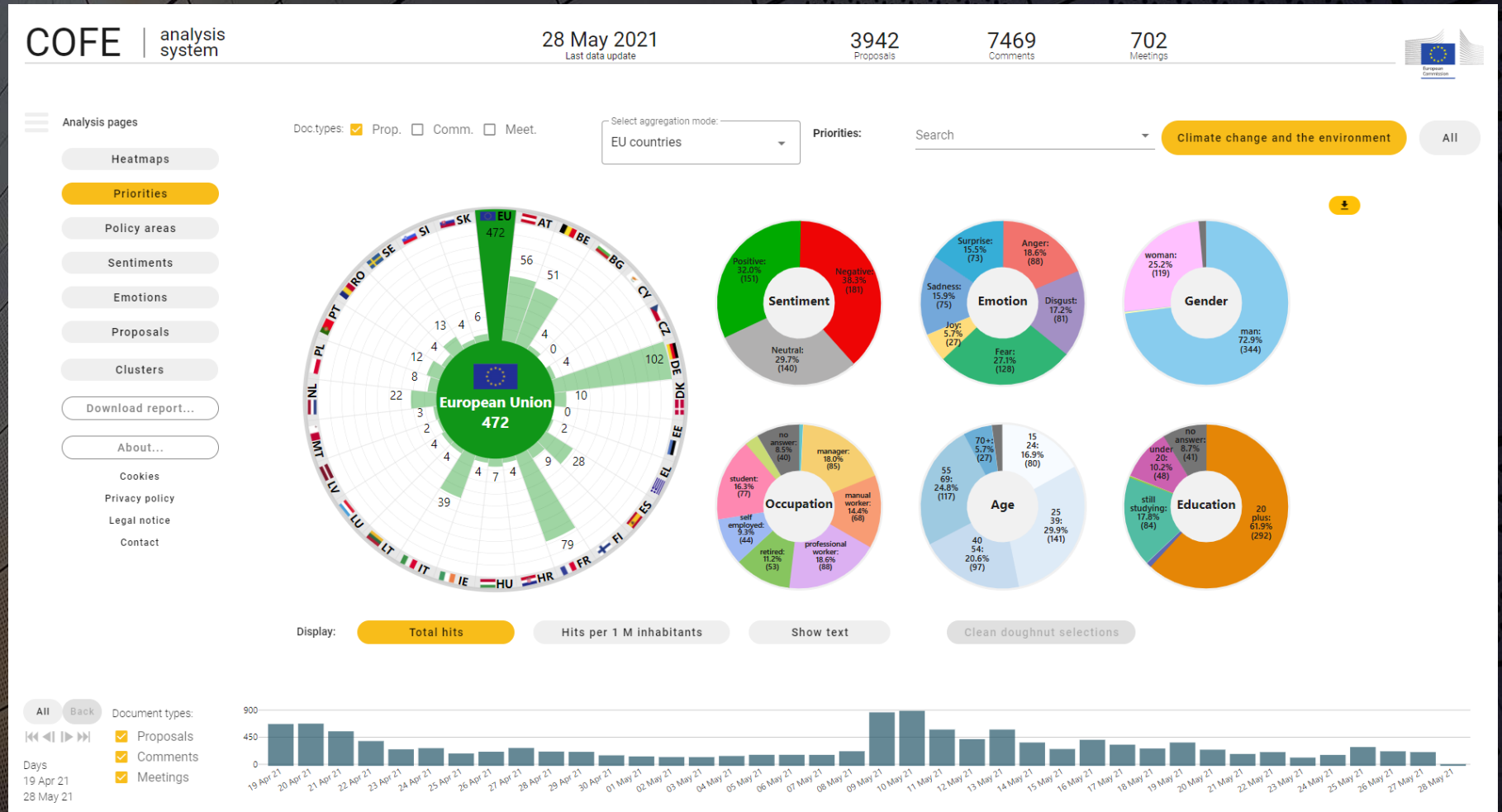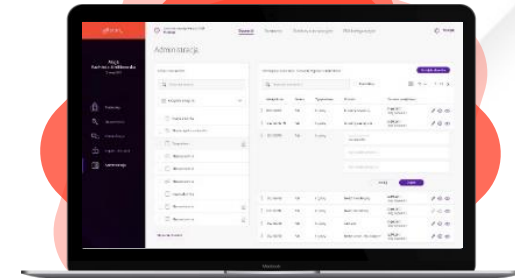
# Ailleron offers a comprehensive
## DIGITAL BANKING PLATFORM



Omni-channel platform providing flawless user experience

Ready to go capabilities – value added services

Back-office modules which support retail services administration

## with Value Added Modules

**01 Digital Onboarding (eKYC) > 02 Meeting Scheduler > 03 Subscription Manager > 04 Notifications & Campaign Manager > 05 Robo Assistant**

# AI Bank
## SOLUTION

**AI First Banking**

Profitability | Personalisation | Automation | Innovation | Omnichannel Experience

| Smart alerts and automations in the area of recurrent transactions | Daily Banking Support (FAQ) | Smart product recommendations | Budget analysis & forecast | Customer behavior predictions (4 scenarios) | Transaction anomalies (security) |

**Ailleron AI Assistant**

Data Analysis| Monitoring | Recommendations | Alerts | Integration

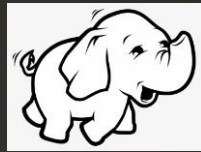| Card transactions | Internet & mobile banking | CRM | Transfers & payments | DWH | OTHERS |

# Selected case studies

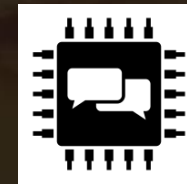# Startup: Electronic Identity Protection Platform

Web-scale Semantic Web startup (2005-2008) commercializing new technology

Bare-metal cluster

Hadoop

Natural Language Processing
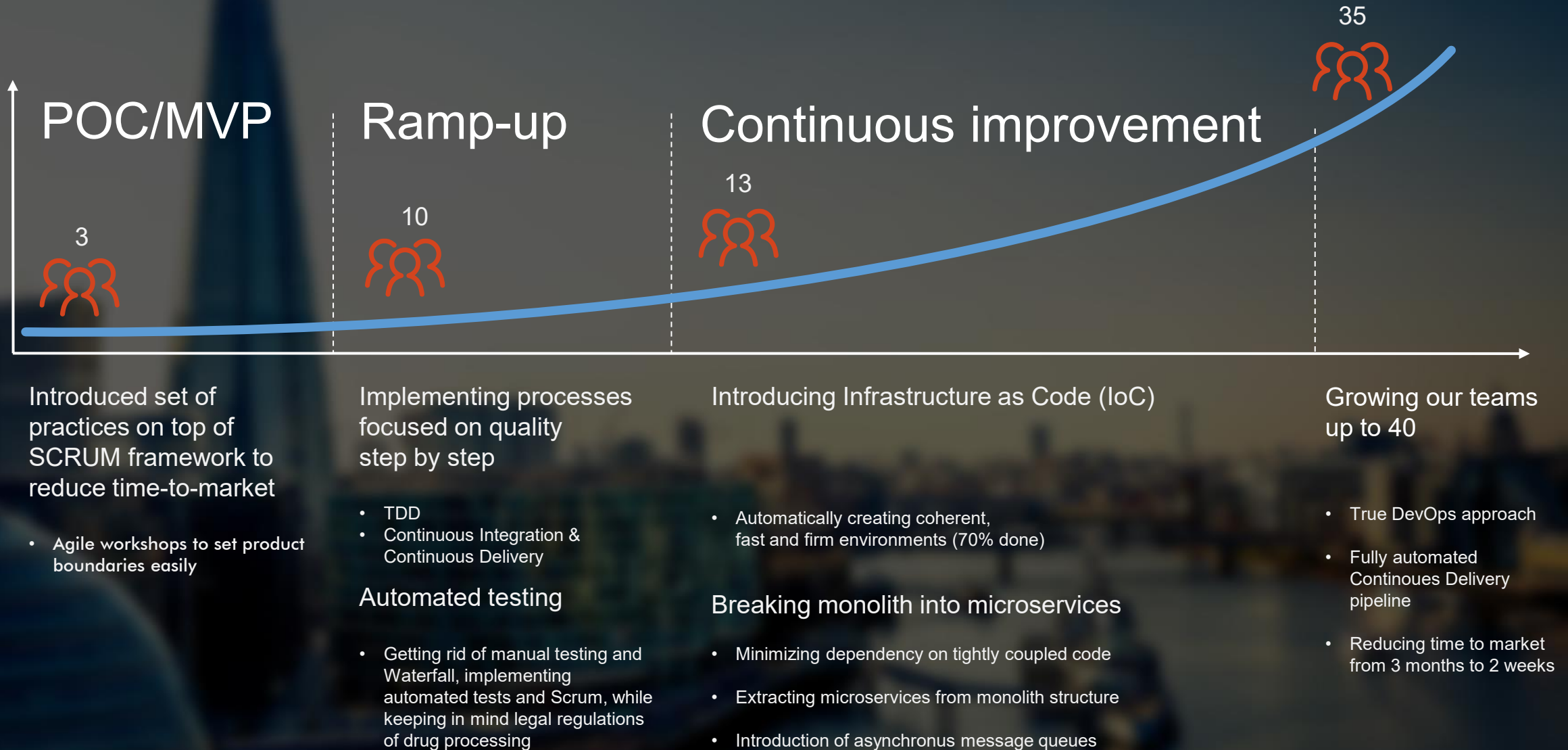
## Main Facts:

**Customer**
UK Startup

**Project**
Building new product
from scratch – electronic
identity protection platform

**Operating Model**
Dedicated Teams,
MVP approach, R&D
with universities (University
of Southampton, University
of Sheffield)

## Challenges:

- Transfer of innovative technologies (Semantic Web, NLP, Graph Databases) from universities to real business

- Scaling the niche technology to fully blown commercial solution (e.g. browsing 4 billion of web pages)

- Building Big Data solution even before this term was defined

# Software Mind Digital Transformation Services



**POC/MVP**  ·  3

**Ramp-up**  ·  10

**Continuous improvement**  ·  13  ·  35

**Introduced set of practices on top of SCRUM framework to reduce time-to-market**

- Agile workshops to set product boundaries easily

**Implementing processes focused on quality step by step**

- TDD
- Continuous Integration & Continuous Delivery

## Automated testing

- Getting rid of manual testing and Waterfall, implementing automated tests and Scrum, while keeping in mind legal regulations of drug processing

**Introducing Infrastructure as Code (IoC)**

- Automatically creating coherent, fast and firm environments (70% done)

## Breaking monolith into microservices

- Minimizing dependency on tightly coupled code
- Extracting microservices from monolith structure
- Introduction of asynchronus message queues

**Growing our teams up to 40**

- True DevOps approach
- Fully automated Continoues Delivery pipeline
- Reducing time to market from 3 months to 2 weeks

# SETA

Big Data - large scale smart mobility managent platform (2017-2019)
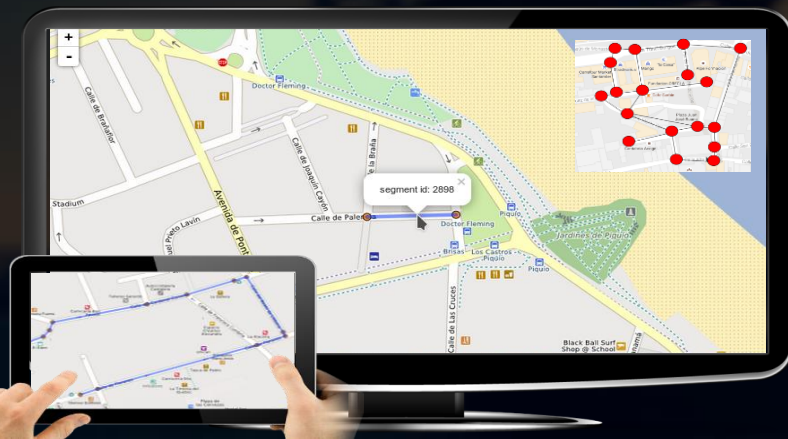
Private Cloud | Cloud - native | Apache Spark | Data Science environment

- Big Data technologies for monitoring and managing mobility in large metropolitan areas.
- Solution based on data from millions of citizens, thousands of connected cars, thousands of city sensors and hundreds of distributed databases

## Main Facts

- Management of huge geospatial and spaciotemporal data
- GPU-based acceleration of geospatial indexes
- Privacy by Design
- High performance geo-located event processing engine
- Scalable backend for mobile aps
- Integrated machine learning components
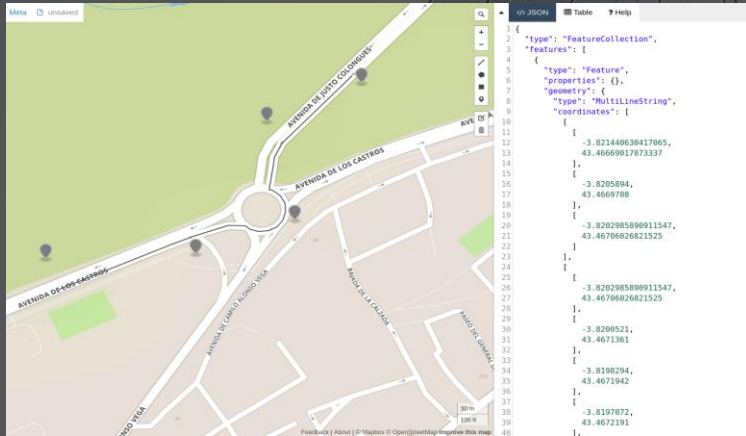- Data Science environment built into Agile software development process

segment id: 2898
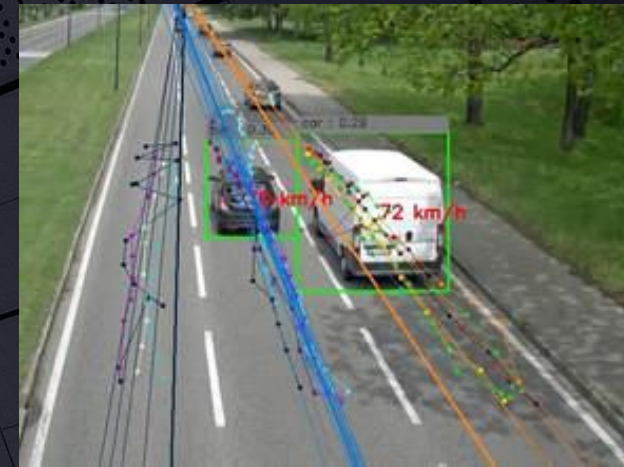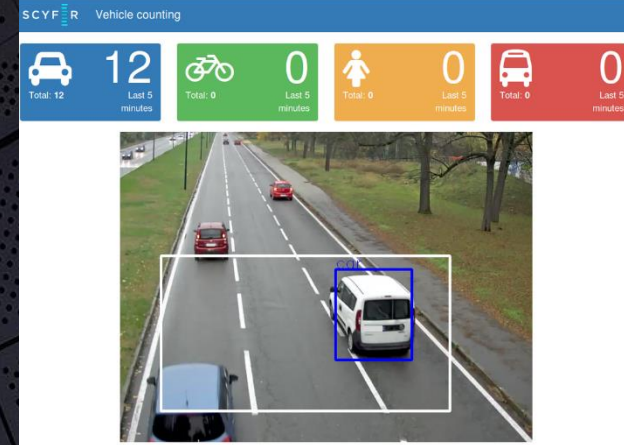
# Seta

Data Science Environment - examples

**Geospatial visualisation
+ map matching**
- low veracity of GPS tracks
- map matching: Hidden Markov models

**Graph analysis/algorithms**
- Analysis and visualisation
  in the context of road infrastructure graph
- Jupyter Notebook + Spark/GraphX

# Data Science environment - FinTech



- Collecting blog posts
- Natural Language Processing
  - Named Entity regocnition
  - Concepts (topic extraction)
  - Insight Engine,
    - eg. patterns of named entities: distance between entities, entities in the same sentences, statistical patterns, etc.)
- Combining data
- Visualisation and data mining
- Text Classification
  - Binary, Multiple class, Multilabel, Complex taxonomy
  - common use cases
    - IPTC Subject Codes
    - EuroVoc
    - Business Reputation
    - IAB Taxonomy
    - Social Media

# Software Mind

# Conversational AI

Chatbot for financial institutions

- Ailleron product

- Dialogue-based machine learning (advanced conversations not just question and answer)

- Intent classification, Named Entity Recognition state-of-the-art NLP methods implemented and used

- An Advanced administration panel means adding and managing bot knowledge with ease

# Software Mind

# Thank you for your attention!

Marcin Sieprawski
Head of Big Data Lab
marcin.sieprawski@softwaremind.com

**Headquarters**
ul. Zyczkowskiego 20
31-864 Krakow

Tel.: + 48 12 252 34 00
web: www.softwaremind.com

**Branches in Poland**
Warsaw, Rzeszow, Bielsko-Biala

**Worldwide**
Australia / Sydney, Brisbane
Singapore / Singapore
USA / Boise, Dallas

Krakow, Poland