# Big Data Architectures & Real World Experiments

Yann Vernaz

# Outline

# Introduction

## ERODS Team (27 pers.)

**E**fficient and **RO**bust **D**istributed **S**ystems

**Responsable** : Noël De Palma

Equipe de recherche commune CNRS, Grenoble INP, UJF, UPMF

**Site Web** : http://erods.liglab.fr

# Processing Models

## Batch Processing

- Familiar concept of processing data en masse.
- Generally incurs a high-latency.

## (Event-) Stream Processing

- A one-at-a-time processing model.
- A datum is processed as it arrives.
- Sub-second latency.
- Difficult to process state data efficiently.

## Micro-Batching

- A special case of batch processing with very small batch sizes.
- A enjoyable mix between batching and streaming.
- At cost of latency.
- Gives stateful computation, making windowing an easy task.

# Big Data Architectures

# Requirements dictate the choice

### Latency
Is performance of streaming application paramount.

### Development Cost
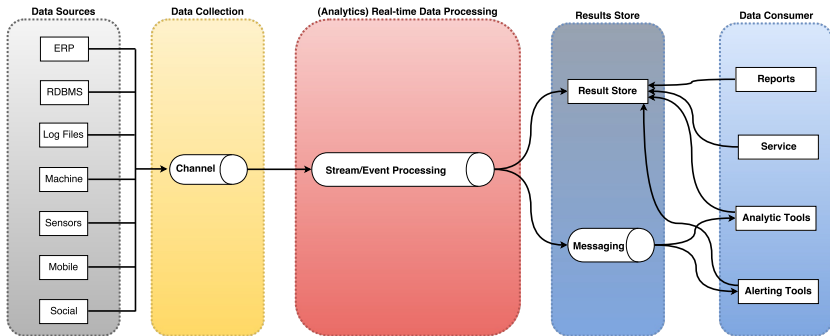Is it desired to have similar code bases for batch and stream processing.

### Message Delivery Guarantees
Is there high importance on processing every single record, or is some normal amount of data loss acceptable.

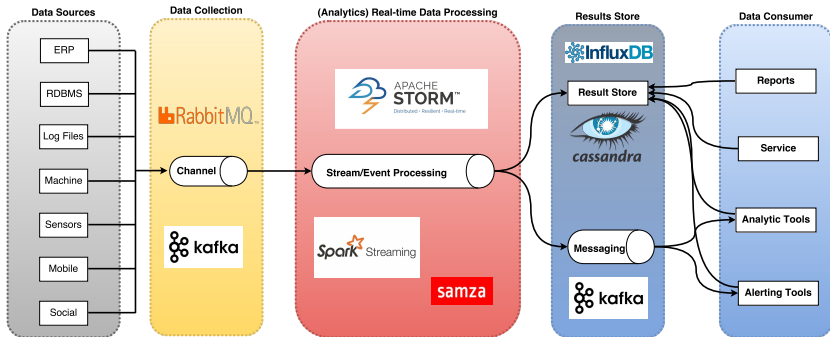### Process Fault Tolerance
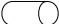Is high-availability of primary concern.

# Streaming Analytics Architecture

# Streaming Analytics Architecture (Technologies)

# Lambda Architecture

# Kappa Architecture

# Unified Architecture
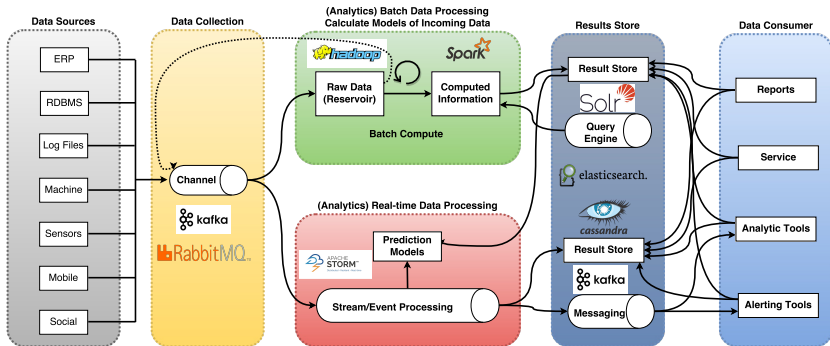
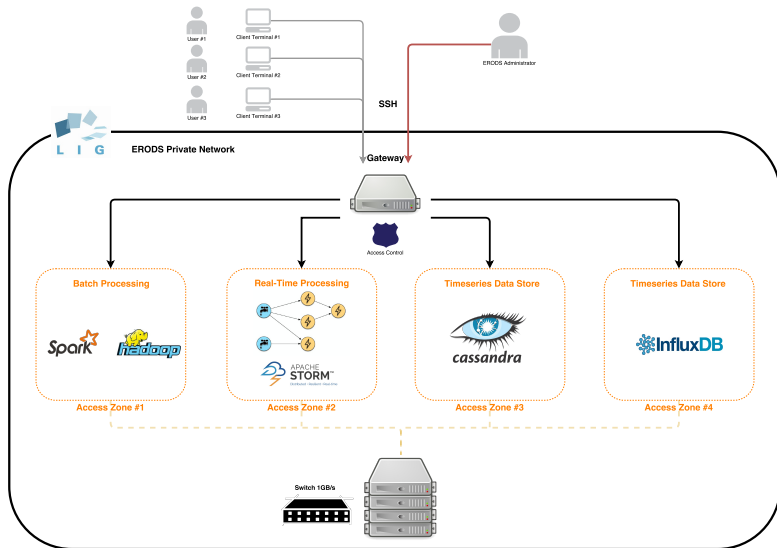# Unified Architecture (Technologies)

# Real World Experiments
# Proof-Of-Concept
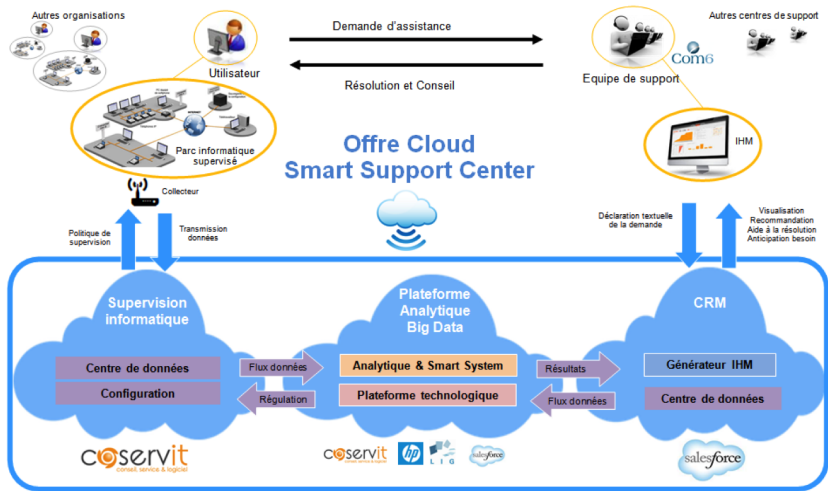
# Virtualization: OpenStack + KVM



- Our experimental platform is Openstack-based.
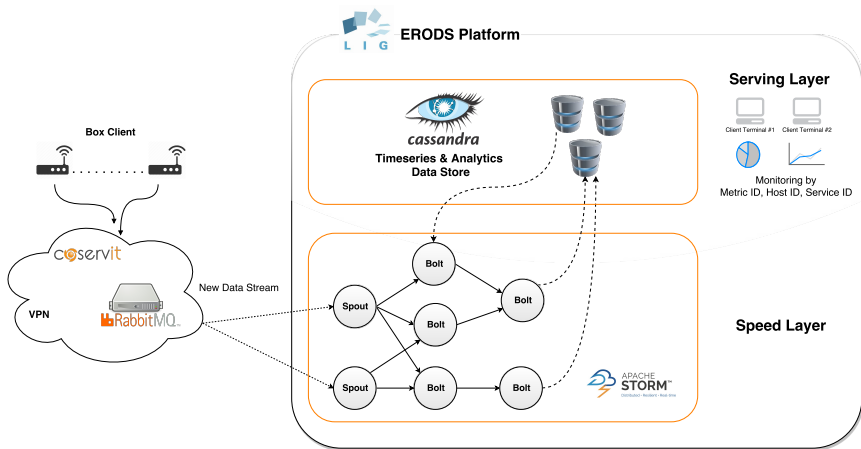- OpenStack has been deployed using Fuel.
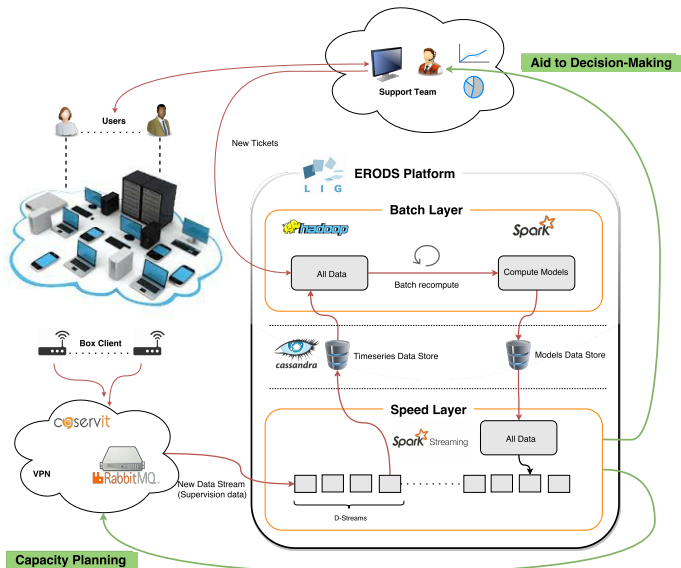
# Services Architecture

# Smart Support Center - Description
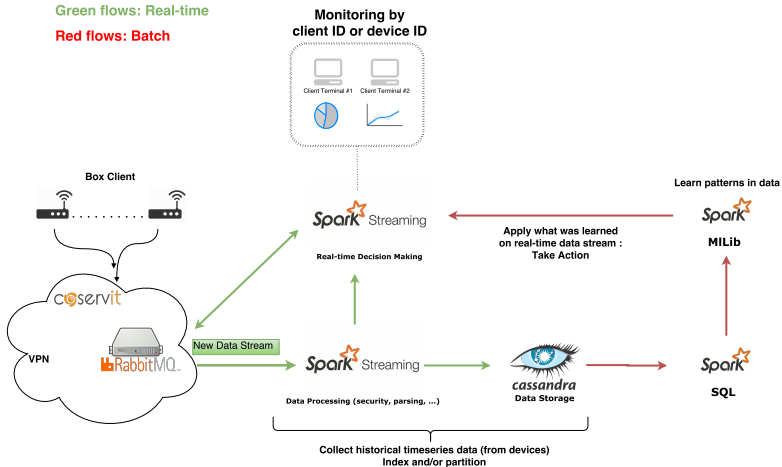
# Smart Support Center - Streaming Analytics

# Smart Support Center - Analytics System

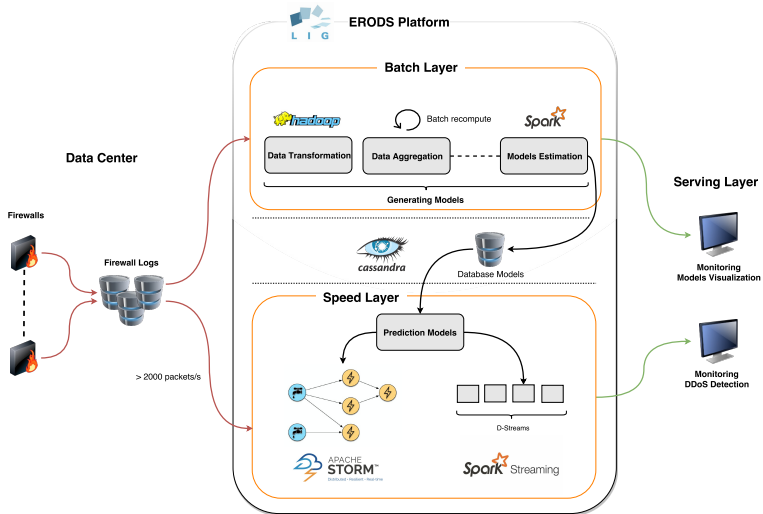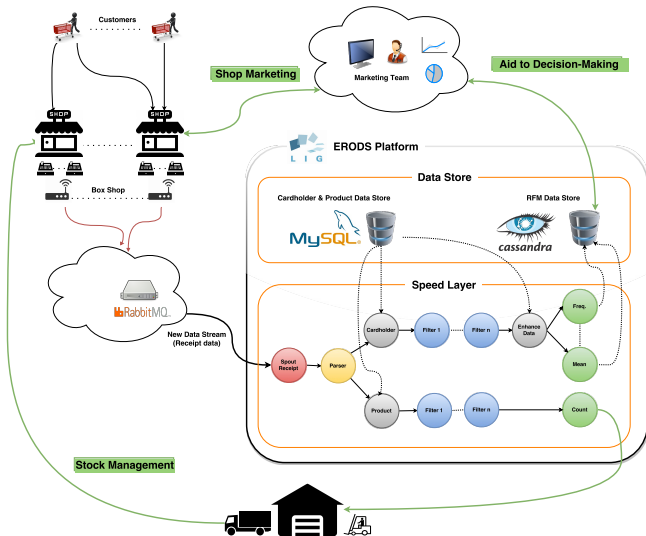# Smart Support Center - Data Flow

# DDoS Attack Detection - Architecture

# Churn Risk Analysis

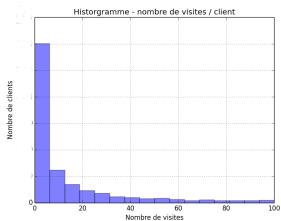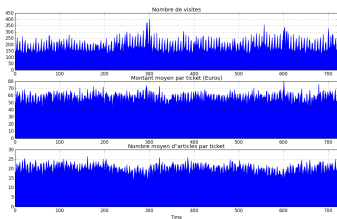# Churn Risk Analysis - Streaming Analytics



Compute real-time statistics by client, shop or country.

# Churn Risk Analysis - Real Time Reporting

# Churn Risk Analysis - Customer Churn Probability



Compute real-time churn risk probability by client.

# Churn Risk Analysis - Results.

# Projet STREAM



Objet de l'étude : écoulements turbulents et/ou multiphasique.



L : 200-0,5 mm; f : 0,1-10 kHz

# Projet STREAM



Plateforme Coriolis

Canal à houle

**Grande soufflerie**

**Soufflerie climatique**

**Cameras phantom (x2)**

**Instrumentation acoustique**

Transducteur ultrasonore

Halles expérimentales

Fibre optique

Zone tertiaire

**Postes de travail**

**Stockage massif**

**Cluster BIG DATA**

**Adaptation des installations**

**Nouvelle instrumentation**

**Infrastructure fibre optique**

**Nouvel environnement Big Data**

# Projet STREAM



Dispositif expérimental

L : 5 cm-50 μm; t : 1-100 ms (eau)
**Mesure d'accélération :
100 000 i/s**

**Suivi lagrangien**

Suivi dans le temps de chaque particule

- 3x24 GB toutes les 20 s soit **300 TB/jour** d'images brutes...
- Extraction des positions en temps réel.
- Gestion des tailles, trajectoires et vitesses des particules.

# Projet STREAM - Topology Storm

# Summary

# Big Data is still "work in progress"

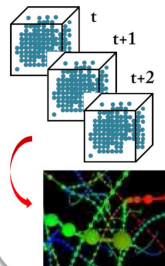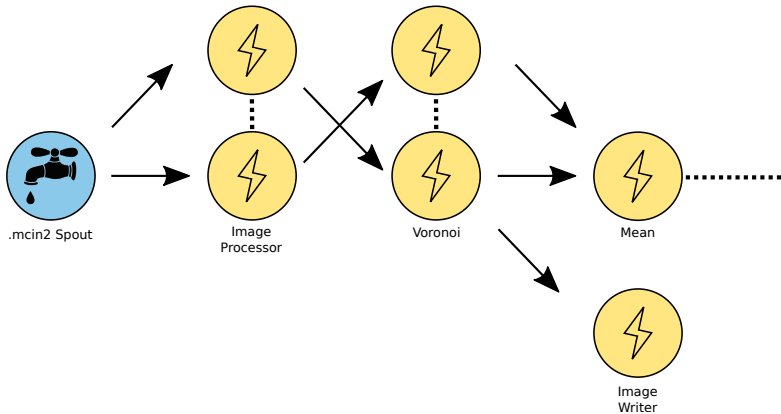Choosing the right architecture is key for any (big data) project.

Know the use cases before choosing your architecture.

There are no standard architectures available which have been used for years (young field).

In the past few years, a few architectures have evolved and have been discussed.

To have one/a few reference architectures can help in choosing the right components.

# From Offline to Online

Recently, Big Data analytics has been transitioning from **offline** (or batch) to primarily **online** (or streaming).

## Forrester remarked the following in Q3 2014

*"The high velocity, white-water flow of data from innumerable real-time data sources such as market data, Internet of Things, mobile, sensors, click-stream, and even transactions remain largely unnavigated by most firms. The opportunity to leverage streaming analytics has never been greater."*

**Velocity** is one of the 7Vs commonly used to characterize Big Data.