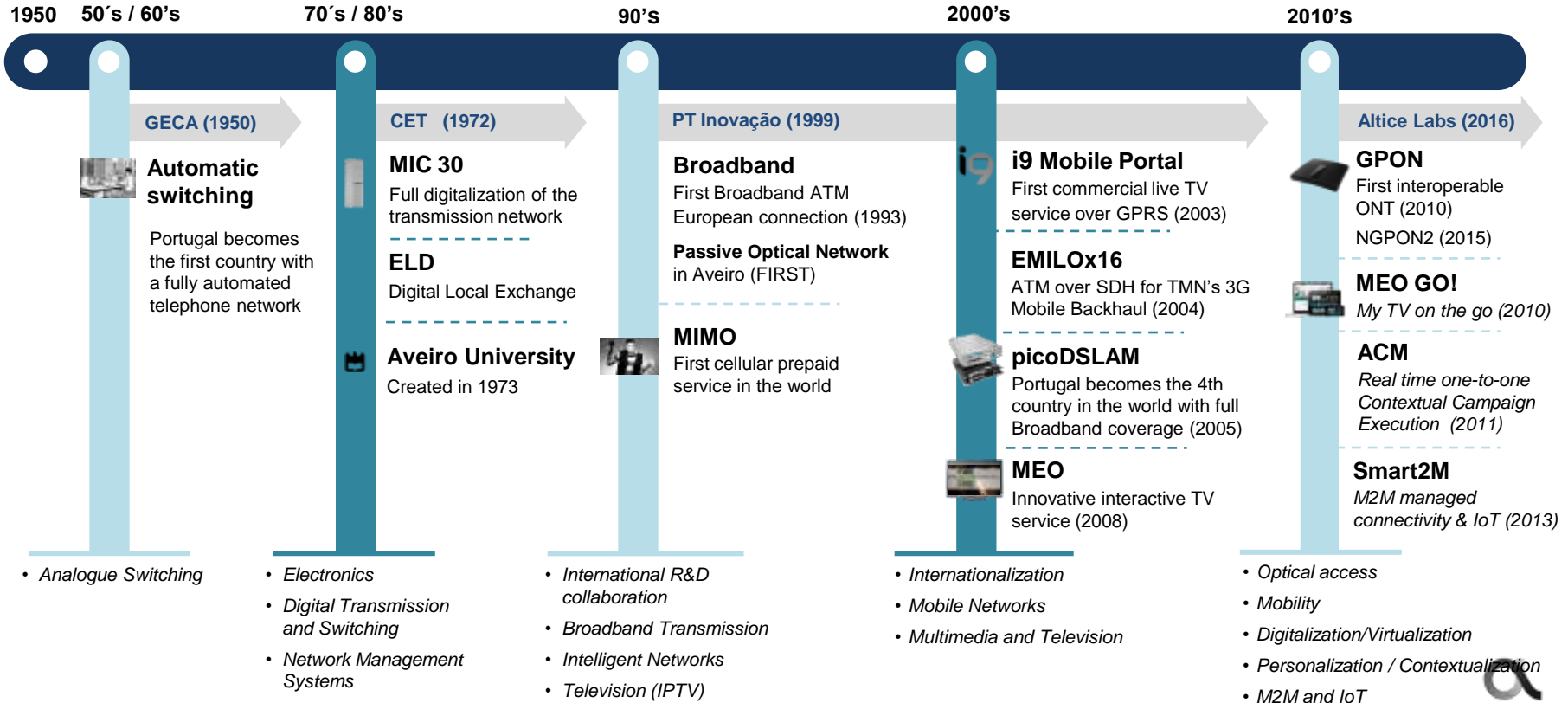# Big Data platform architecture at Altice Labs

**Mário Moreira**
Head of Experimentation and Technology Coordination
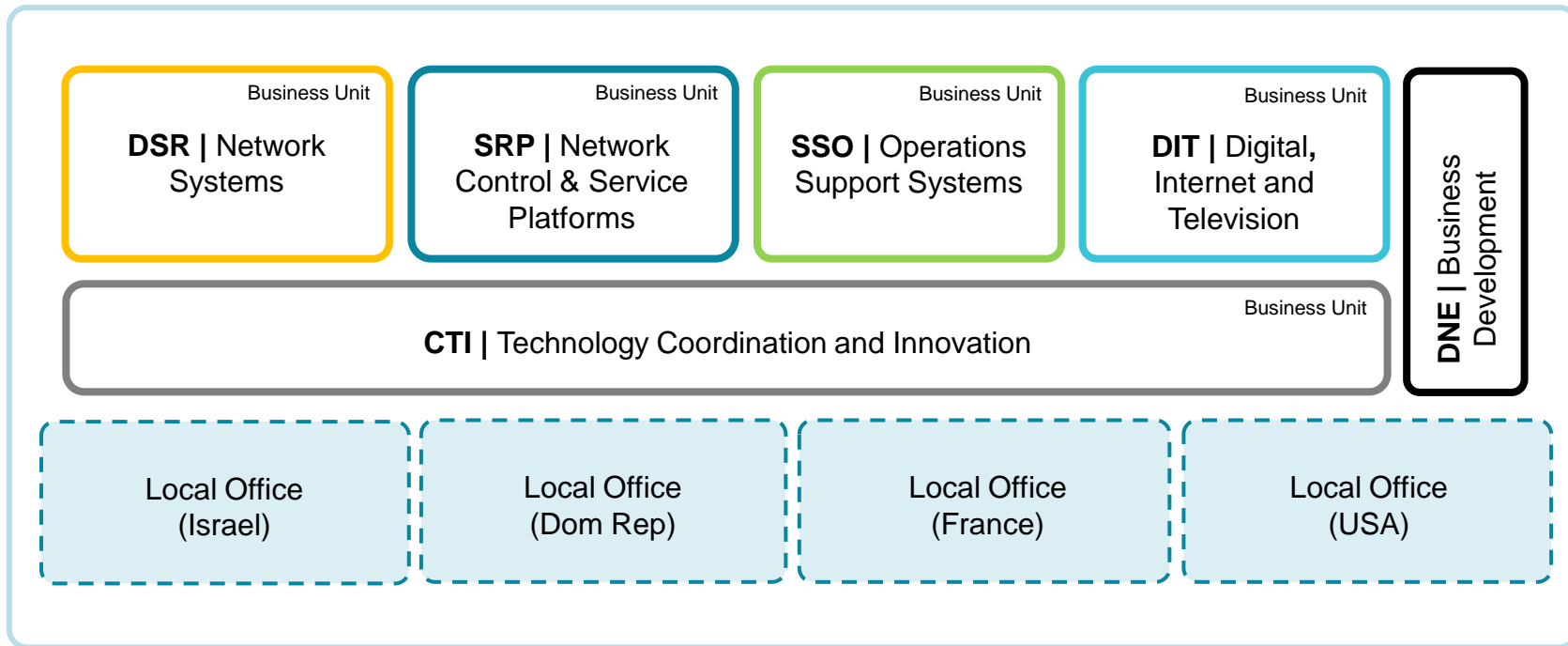
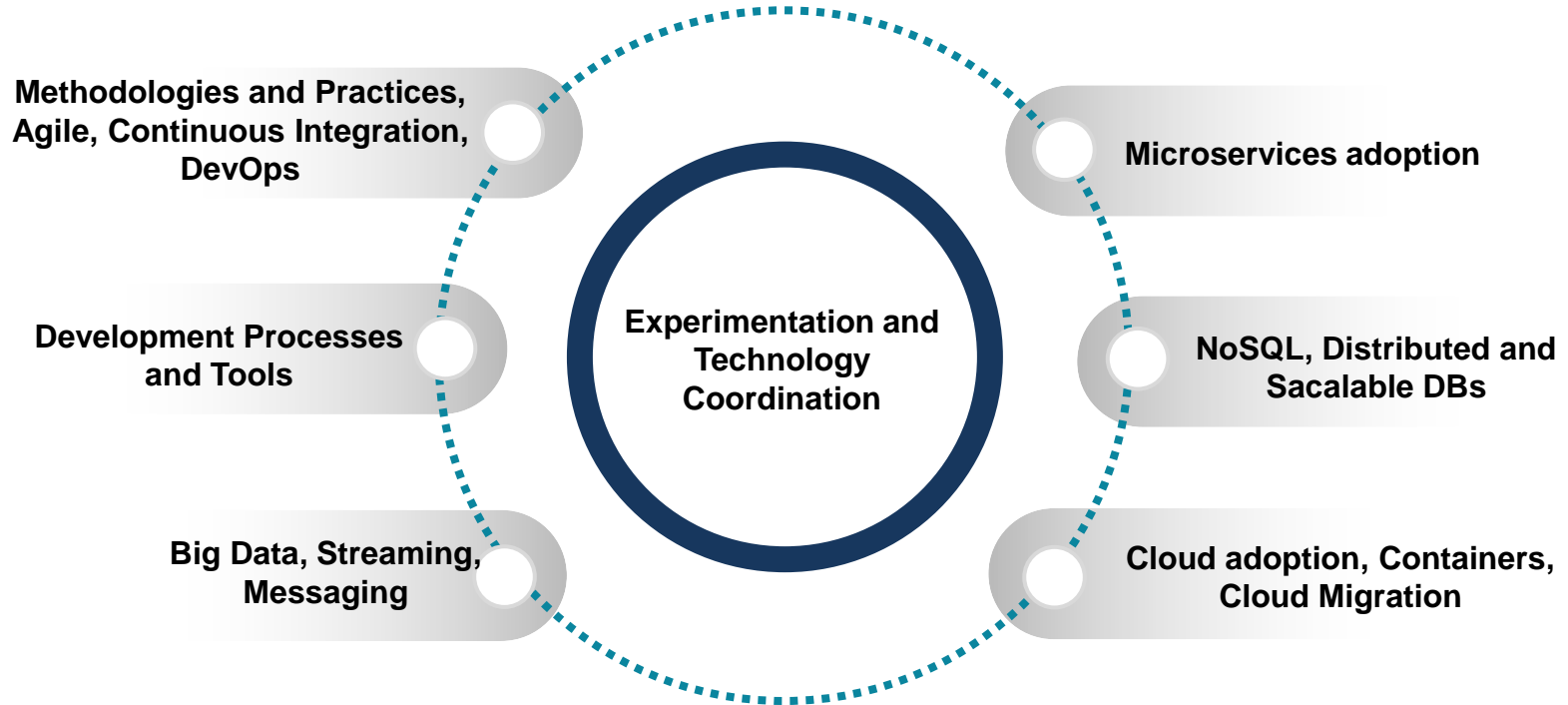# The history of Altice Labs is linked to the Portuguese telecommunications sector evolution

**1950** **50´s / 60's** **70´s / 80's** **90's** **2000's** **2010's**

**GECA (1950)**

**Automatic switching**

Portugal becomes the first country with a fully automated telephone network

**CET (1972)**

**MIC 30**
Full digitalization of the transmission network

**ELD**
Digital Local Exchange

**Aveiro University**
Created in 1973

**PT Inovação (1999)**

**Broadband**
First Broadband ATM European connection (1993)

**Passive Optical Network**
in Aveiro (FIRST)

**MIMO**
First cellular prepaid service in the world

**i9 Mobile Portal**
First commercial live TV service over GPRS (2003)

**EMILOx16**
ATM over SDH for TMN's 3G Mobile Backhaul (2004)

**picoDSLAM**
Portugal becomes the 4th country in the world with full Broadband coverage (2005)

**MEO**
Innovative interactive TV service (2008)

**Altice Labs (2016)**

**GPON**
First interoperable ONT (2010)
NGPON2 (2015)

**MEO GO!**
*My TV on the go (2010)*

**ACM**
*Real time one-to-one Contextual Campaign Execution (2011)*

**Smart2M**
*M2M managed connectivity & IoT (2013)*

- *Analogue Switching*

- *Electronics*
- *Digital Transmission and Switching*
- *Network Management Systems*

- *International R&D collaboration*
- *Broadband Transmission*
- *Intelligent Networks*
- *Television (IPTV)*

- *Internationalization*
- *Mobile Networks*
- *Multimedia and Television*

- *Optical access*
- *Mobility*
- *Digitalization/Virtualization*
- *Personalization / Contextualization*
- *M2M and IoT*

2

altice labs

**200 million people communicate everyday through technology developed by Altice Labs**

# Governance



| Business Unit | Business Unit | Business Unit | Business Unit | |
|---|---|---|---|---|
| **DSR \|** Network Systems | **SRP \|** Network Control & Service Platforms | **SSO \|** Operations Support Systems | **DIT \|** Digital, Internet and Television | **DNE \|** Business Development |

**CTI \|** Technology Coordination and Innovation — Business Unit

| Local Office (Israel) | Local Office (Dom Rep) | Local Office (France) | Local Office (USA) |
|---|---|---|---|

altice labs

# Current Innovation Areas

**Methodologies and Practices, Agile, Continuous Integration, DevOps**

**Development Processes and Tools**

**Big Data, Streaming, Messaging**

**Experimentation and Technology Coordination**

**Microservices adoption**

**NoSQL, Distributed and Sacalable DBs**

**Cloud adoption, Containers, Cloud Migration**

altice
labs

# Big Data @ Altice Labs

# Big Data @ Altice Labs

**4T 2012**             **2013 - 2014**             **2015**             **2016 - Today**

## Phase 1

**To experiment and to obtain know how in Big Data Technologies**

Technologies used Cloudera Hadoop, Map Reduce, Pig, Impala, Hbase, Flume, Sqoop, Oozie, Zookeeper, …

Benchmarks against the traditional approach (Oracle DB + Java code)

Learn what works and what doesn't work

## Phase 2

**To "sell" and disseminate Big Data technologies to product teams**

Create a reference Big Data technology stack and a technology adoption strategy for ALB products

Comparative studies on NoSQL Databases (MongoDB, Cassandra, Redis, Neo4j, VoltDB, HBase)

## Phase 3

**Experiment with Real-time Big Data technologies**

Select and validate the technologies better suited to implement high performance and high scalability versions of some ALB products

Comparative analysis and benchmark of Storm and Spark (we didn't consider Flink and Samza mature enough at that time)

## Phase 4

**Management of MEO Big Data platform (ex-Sapo)**

Most data is from Web Portal, mobile apps and IPTV

Participation on Altice Big Data Group project to standardize Big Data technology stack, reference architecture and data models and governance

altice labs

# About the "Golias" platform

## NUMBERS

**Total of 98 servers holding 191 TB (compressed data) and 8 TB (uncompressed aggregated data) divided by 10 different clusters**

**Data is received from +40 different distinct systems and is processed by +240 distinct processes**

## DATA INGESTION

**By batch processes that import data from external databases**

**By Sapo Broker, and in house developed pub/sub messaging platform**

**By service listeners that receive events from webpages, mobile apps, set top boxes apps, etc**

## SOFTWARE USED

**Cloudera for the Hadoop clusters, 3 PostgreSQL clusters, 2 Cassandra clusters, MonetDB, and 49 instances of Solr**

**Data is handled by Pig, Java or Python processes**

## DATA ACCESS

**In house developed web applications for analytics and dashboards**

**Using in house developed API interfaces**

**Ad hoc reports as needed by the business areas**

altice labs

# Altice Labs Big Data reference architecture

# Reference Architecture



External Databases  BI Tools  Data Visualization  Others …

**Data Service Layer (Public Interfaces)**
- KPIs DBs
- Recent Data
- **Power User Interface**
- **AdHoc Reporting**
- **Service APIs**

**Data Processing Layer**
- Internal DBs
- **Altice Common Libs (Aggregation, Indexing, Correlation, …)**
- **BigData Distribution Base Components (Hadoop, Yarn, HDFS, Pig, …)**
- CEP
- Data Mining and Machine Learning
- Streaming Processing

**Data Ingestion Layer**
- **Batch Ingestion**
- **Event / Streaming Ingestion**

Process Orchestration | Big Data Cluster Management | Data Governance Policy and Tools

External Databases  Enterprise Applications  Network Elements  Service Platforms  Others …

10

altice labs

# Data Ingestion - Batch



**① Explore**

- Apache NiFi
- StreamSets

**② Experiment**

- Sqoop
- Flume

**③ Adopt**

- Shell scripts
- Python

altice labs

# Data Ingestion - Streaming



① **Explore**

② **Experiment**

③ **Adopt**

- Kafka
- REST APIs

# Data Ingestion – Shared Storage



**①  Explore**

- CephFS
- BeeGFS (old FGFS from Fraunhofer)

**②  Experiment**

**③  Adopt**
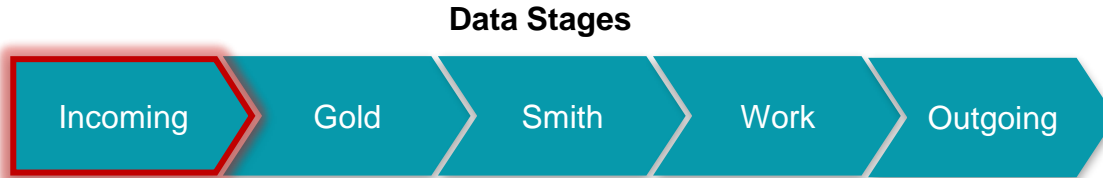
- GlusterFS

# Hadoop Base Distribution



**1** **Explore**

**2** **Experiment**

**3** **Adopt**

- Cloudera

# Data Layer – Incoming stage



**Data Stages**

| Incoming | Gold | Smith | Work | Outgoing |

① **Explore**

② **Experiment**

③ **Adopt**

- Kudu
- HDFS

# Data Layer – Gold and other stages



**Data Stages**



| Incoming | Gold | Smith | Work | Outgoing |

**1** **Explore**

**2** **Experiment**

**3** **Adopt**

- Impala
- HDFS

altice labs

# Data Processing - Batch



**1** **Explore**

**2** **Experiment**

**3** **Adopt**

- PIG
- Python

# Data Processing - Streaming



**(1) Explore**

- Spark

**(2) Experiment**

- Storm
- Flink

**(3) Adopt**

altice
labs

# Working In-Memory Storage



**①** **Explore**

**②** **Experiment**

**③** **Adopt**

- Redis

# Data Warehouse



**1** **Explore**

**2** **Experiment**

- Citus DB
- Apache Drill

**3** **Adopt**

- PostgreSQL

# Data Science Playground - Storage



**1** **Explore**

**2** **Experiment**

- Citus DB
- Apache Drill

**3** **Adopt**

- HDFS
- Impala
- PostgreSQL

altice labs

# Data Access APIs



**1** **Explore**

**2** **Experiment**

- Apache Drill
- REST APIs

**3** **Adopt**

# Cluster Management and Monitoring



**1** **Explore**

**2** **Experiment**

**3** **Adopt**

- Cloudera Manager

# Process Orchestration



**1** **Explore**

**2** **Experiment**

- Apache Airflow

**3** **Adopt**

# Data Governance, Lineage and Audit



**(1) Explore**

- Cloudera Navigator
- Apache Atlas

**(2) Experiment**

**(3) Adopt**

altice labs

# Power User Interface and Ad Hoc Reporting



① **Explore**

② **Experiment**

- Apache Zepplin
- Jupyter

③ **Adopt**

altice labs

# CEP and Machine Learning



**①** **Explore**

- Esper (CEP) (within storm)
- Apache Mahout
- Python – numpy, pandas, scikit-learn
- Yet no experience with Tensorflow, Mxnet, H2O, Theano, Torch, …

**②** **Experiment**

**③** **Adopt**

# About the Technologies Presented

Data visualization and exploration tools are out of this scope (assuming reuse of already installed tools like Tableau and others).

Bussiness users need to be able to take advantage of the Data Lake. Traditional BI and Analytical tools may not play well with Hadoop ecosystem of technologies (proprietary formats, informaton silos, designed assuming relational model, …). We will need new tools that integrate well with these technologies for bussiness users (ex: Dataiku, others). This might imply user retraining.

It's still a Work in Progress – as we experiment and learn more we might change a few things.

Your comments and sugestions are wellcome.

altice
labs

# Big Data platform architecture at Altice Labs

**Mário Moreira**
Head of Experimentation and Technology Coordination

Any questions?