#STechDay2020

A Data Journey: from the Data Warehouse to the Data Lakehouse





ORGANIZATION



SPONSORS





#### THANK YOU!

# TECH DAY

#### #STechDay2020



## Pablo A. Doval

Principal Data Architect – Plain Concepts UK

"I work with code and data, but don't tell my mom; she thinks I'm a piano player in a whorehouse."

@PabloDoval

palvarez@plainconcepts.com





### Innovation through Data and AI...

Smart	Video Content	Real Time	Knowledge	Lead Cases
Maintenance	classifier	Churn Analysis	Extraction	Selection
Anomaly	Audience	Automatic	Document	Semantic
Detection	Segmentation	Case Classifier	Classification	Search
Behavioural	Smart Pricing	Case Outcome	Clause Outlier	Settlement
Analysis	Models	Forecasting	Detection	Forecast



## Where are we?

Unstructured, chaotic data estate

Lack of skills and productivity

#### Solutions not Enterprise-ready







### We have evolved and transformed...



## ... but what about our data systems?

#### Inventory Data Entry SBEM0412 FEC There is something different about how or where this Xero account was accessed. You won't be able to see any contact inform SISTEMA DE PAGOS MISMO DIA но SBE0412/ Main Status History Order G CONSULTA DETALLE DE LOTES RECIBIDOS Type: BOIL STA LOTE 2 RECHAZADO STA RESPUESTA Account: 05-2852 # INV-001 The Hub King' JuneFee F Is: LOTE 12091 SECUENCIA 00001 FECHA APLICACIO TPO REGISTRO 01 ORDENANTE BANCO 00021 DESCONOCIDO TPO CTA 01 CHQ CUENTA 000000 Class: NOMBRE OTROS BANCO RFC BC00306 Item Quantity PROCESO 00003 TRAS COMB PRODUCTO 00002 AC CTAS EXT Cost center: CARP :: Office Space Rental 1 INSTRUCCION 01 CARGO 1:1 MONEDA 01 MN CANAL ACCES LEY. CARGO CARGO 1 DETALLE **1**/Å Account Rent -Tax Rate 20% (VAT on Expenses) -ENVIADO RECIBIDO Location MOVTOS 00005 5 Add another item. Bet wall IMPORTE 00003000001500000 COD. RECH. 00500 DESCRIPCION FONDOS INSUF SE INTENTARA EN S Belt wall peg 1 CTA ORD REAL 4000000018 TP CODE TEST FOLIO 56 CC Boiler room cab HORA RECH HORA ALTA HORA CANC HORA REP Bet rack 17:07:52 17:07:56 <PF5> MENU ANTERIOR Item #: ASBY-27553 То Ke Switch to classic invoicing 01/00

## TECH DAY

## We have lots of data, but it is in silos

Mobile/Web Transactional Social ΙoΤ Ads Marketplace J محمح R  $\langle \rangle$ 

#### SINGULARITY TECH DAY





### **Reactive System Example**



#### #STechDay2019

### **Proactive System Example**



## Some challenges ahead...

Multi-modal support (tabular, video, audio...)

Support advanced analytics through the data lifetime

Promote rapid prototyping, and rapid productionalization

Data dispersal due to proprietary systems and formats

Lack of common semantic model



TECH DAY

## A bit of data history...





### We need semantic models





### Once upon a time...











Ads

## Marketplace





### Once upon a time...







IT







#### #STechDay2019

### **Lessons Learned**

End users / departments require certain technical skills

Ensure veracity and data quality is practically impossible

Data lineage management is also nearly impossible

Balance between **cost** and **data governance** 



### The mighty Data Warehouse...





#### #STechDay2019

### **Lessons Learned**

High cost in trying to create a single version on truth

Difficult to change without impacting the rest of the org

Lack of multi-modal support (Video, Audio...) and ML

Limited Support for Streaming







#### #STechDay2019

### **Lessons Learned**

Manage master data is a very complex challenge

Requires extra capacity and skills on IT team

Lack of multi-modal support (Video, Audio...) and ML

Limited Support for Streaming

## **The New World**

Real Time Data

**Highly Volatile Data Structures** 

Hybrid and Multi-vendor Ecosystems

**AI/ML** Capabilities



## **Data Lakes**

#STechDay2019

## The promise of the Data Lake







## What a Data Lake is \*not\*

SINGULARITY

**TECH DAY** 

Just storage / Azure Data Lake Storage 😳

Just a Hadoop/Spark/HPC cluster

An more modern kind of Data Warehouse



## **Disclaimer!**

sidra Data Platform

> The next slides will be heavily based on our approach to data lakes, as implemented by Sidra Data Platform. Far from a sales pitch, I will use them to explain our approach to some technical challenges; what has worked and what needed to be improved over time.





C

 $\bigcirc$ 

sidra Data Platform





Q









Q



#510





C

 $\bigcirc$ 







C

 $\bigcirc$ 



sidra 🕎

 $\bigcirc$ 

R

 $\bigcirc$ 



sidra 🕎

R

 $\bigcirc$ 





 $\bigotimes$ 

 $\bigcirc$ 





## Enabling Business Cases via Client Apps





 $\bigcirc$ 





 $\bigcirc$ 





 $\bigcirc$ 







## The road ahead...

#### #STechDay2019

## **Challenges with Data Lakes**

**Deployment Complexity** 

Real-Time (lambda architectures)

Reprocessing (Validations, etc...)

Append, Update and Merge (Right to be forgotten, etc...)

Keeping Historical Versions of Data



#### #STechDay2019

### **Enter the Delta Engine**

Full ACID Transactions

Based on Parquet – Open Standard, Open Source

Runs on Spark





### **Enter the Delta Engine**





## TECH DAY

#### #STechDay2019

## **Delta Lake and the Data Lakehouse**



Source: What is a Lakehouse? - The Databricks Blog (30-Jan-2020)

TECH DAY

## **Our lessons learnt...**

## What have we learnt?

Integrated Data Quality processes

Integrated Data Mastering processes

3-tier Approach using Transactional Tables



## TECH DAY



## THANKS AND...

## **SEE YOU SOON!**

ORGANIZATION

SPONSORS

SUPPORT





