

# From (big) data mess to data as an innovation enabler

Olaf Zschiedrich – CTO OLX Group – GOTO Berlin 2018



**BRAINLY**



**OLX GROUP**



**NASPERS**

Founded 1915  
South-Africa  
Market cap: \$100B

**Tencent 腾讯**

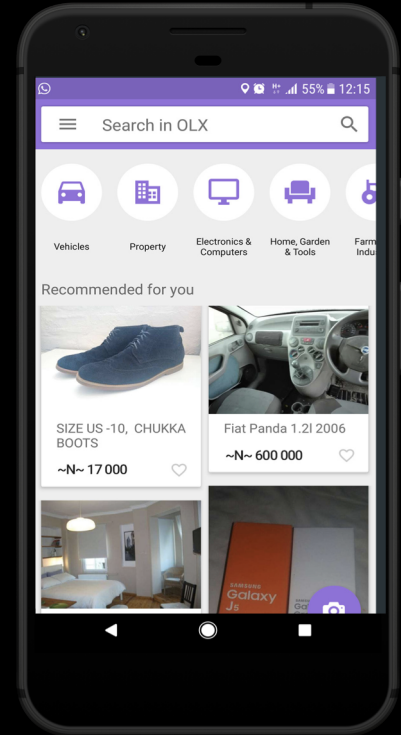


***Delivery Hero***

**eMAG**

# SELLER WINS

# BUYER WINS





OLX GROUP

## Horizontals



letgo



dubizzle

### Real Estate

property24

Domofond.ru

imovirtual

otodom.

storia

### Other verticals

shedd

TRADUS

fixly

### Cars

OTOMOTO

Auto  
Trader  
.co.za

FCG FRONTIER CAR GROUP

AUTOVIT.RO

STANDVIRTUAL

STRADIA





**43** Countries



**35** Offices



**+5,000**  
Employees



**+350M**  
MAU



**+4B**  
Events/Day



The background of the image is a blurred city skyline at sunset. The sky is a mix of orange, yellow, and blue, with the sun's glow creating a warm, hazy atmosphere. In the foreground, two hands are clasped together, belonging to people wearing dark suits and light blue shirts. The hands are positioned in the center of the frame, below the text.

***Halleluja!***  
***We are a data-driven company!***

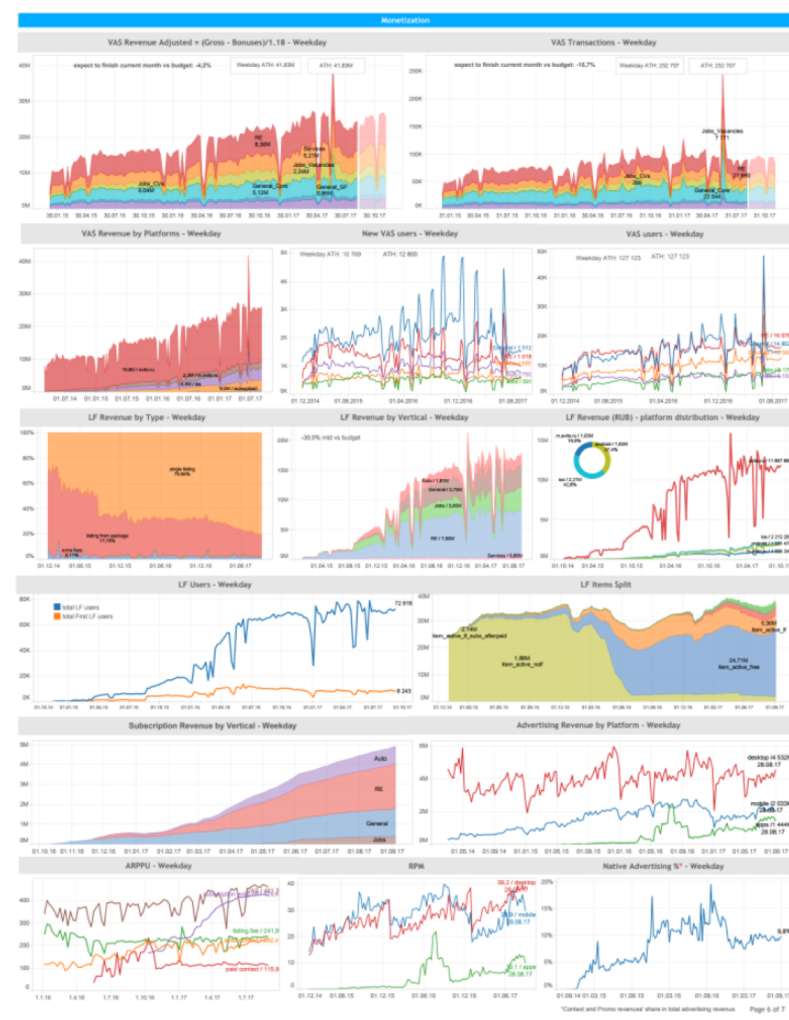
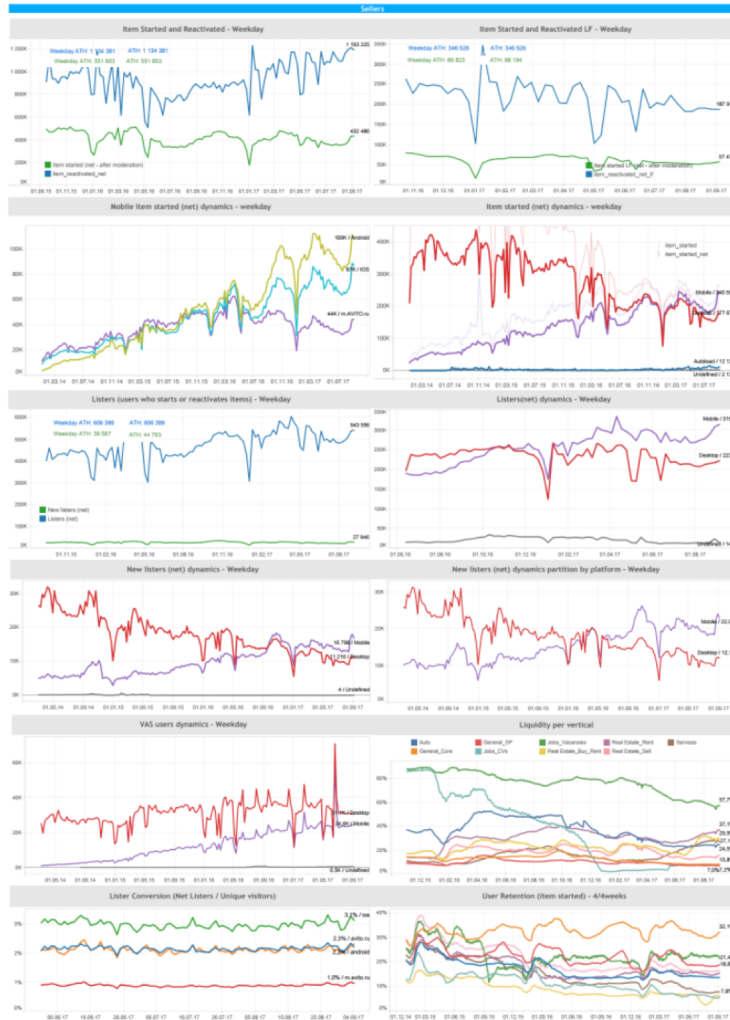
















Disc

42

is broken



Do I launch a new car portal in Mexico?



**CEO**

Shall I invest more in online or in offline marketing?



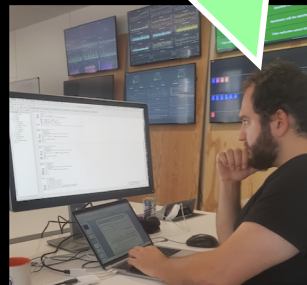
**GM**

How is CS agent #253 performing?



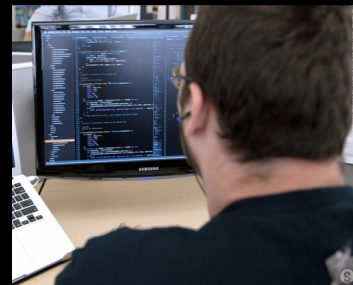
**CS Manager**

Can I predict which listings have the highest probability to sell?



**Business Analyst**

Are we still online?



**DevOps Engineer**

months ← ————— urgency ————— → seconds

***Which data points are really influencing your decisions?***

*‘Give everybody the data  
that he or she needs’*

*(but also not much more)*

**Product innovation (AI/ML)**



**BI and reporting**

Data marts

**Data democratization**

Data reservoirs



**Data collection**

Data lake

# #1 Data Collection



S3

400 TB  
after Compression



**4B events / day**  
**~ 2T events total**

# LIVESYNC

**1.3M listings / day**  
**~ 1B listings total**

## #2 Data Democratization





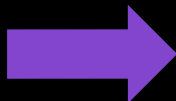
Amazon S3



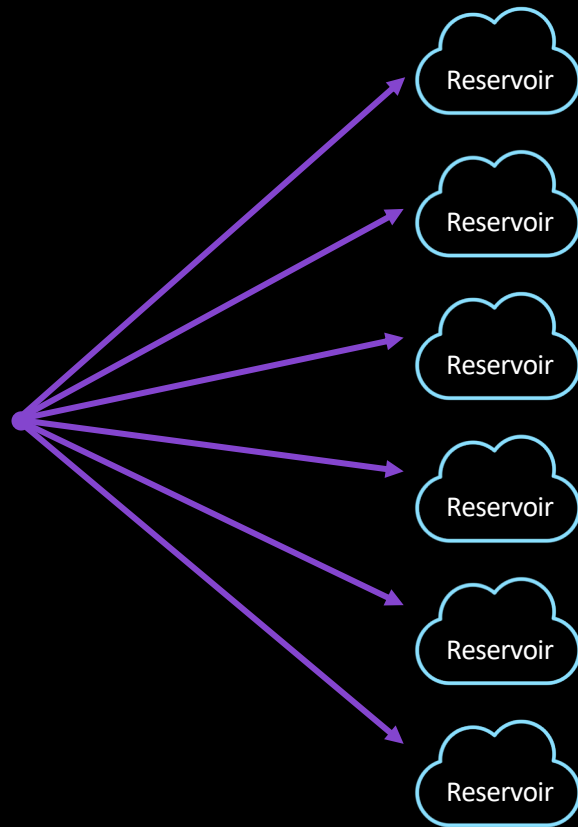
Amazon EMR



AWS Glue



Catalog





# Welcome to Catalog

What do you need to do?



## Browse available data sources

There are 59,233 fields in 3,539 sources and 26 reservoirs.

100% of fields are classified

 All sources

 All reservoirs



## Manage your data sources

You are not producing any data.

→ Start publishing data



## Overall health

- ✓ No owners with bad sources
- ✓ No subscribers consuming unclassified fields
- ✓ No sources have unclassified fields
- ⚠ 16 of 26 reservoirs have risky subscriptions  
(subscriptions using personal fields not anonymized)



## Subscription requests

You have 387 requests to review

 Review



## My reservoirs

Fully anonymised reservoir 2 subscriptions

- ✓ All subscriptions use only classified fields!
- ⚠ 1 subscription are using personal fields not anonymized



## Help

Introduction to Catalog  
How to classify data  
How to subscribe to data  
How to approve subscriptions

### Top Owners

- By health
- By amount of sources

### Top Reservoirs

- By health
- By amount of subscriptions

### Top Subscribers

- By health
- By amount of subscriptions

## Reservoir Structure – S3 Bucket

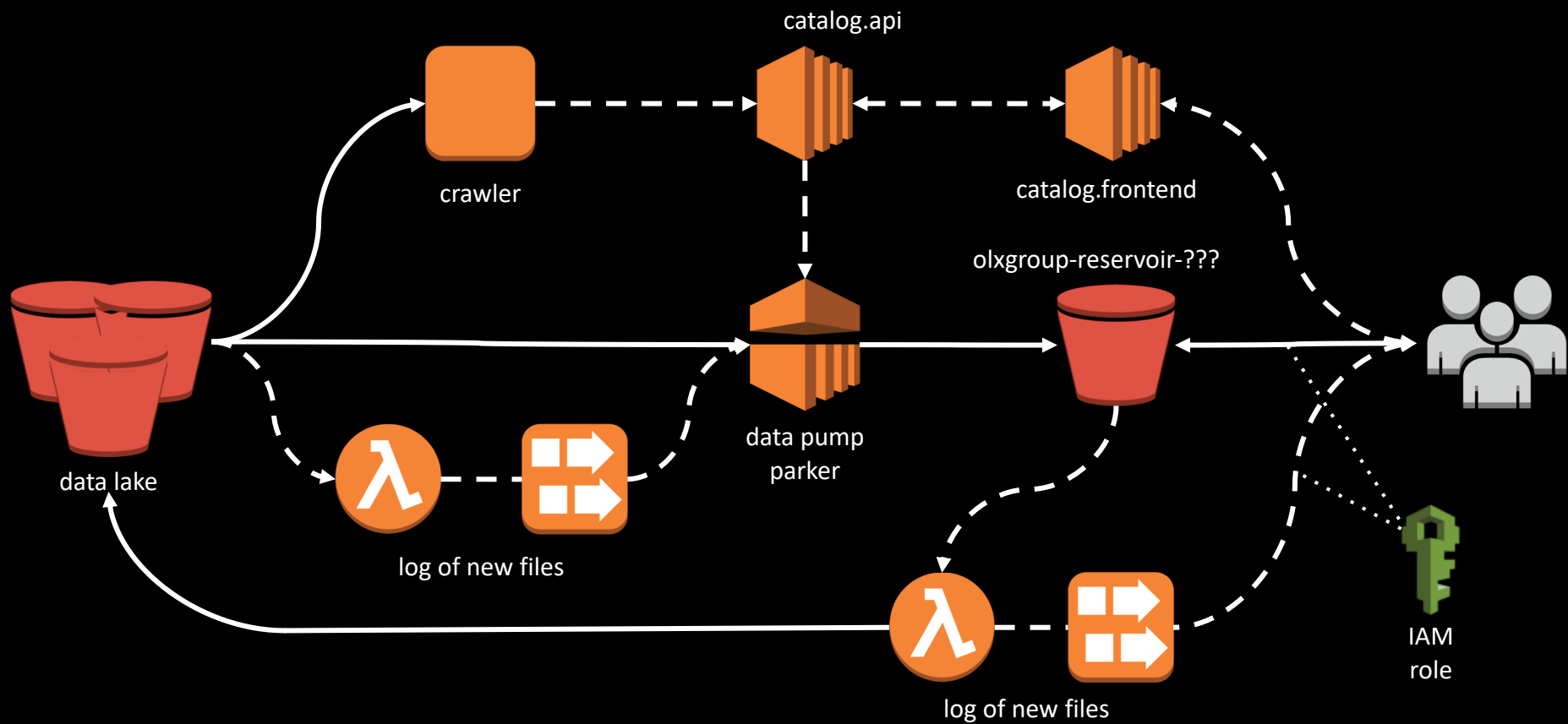


***/in***

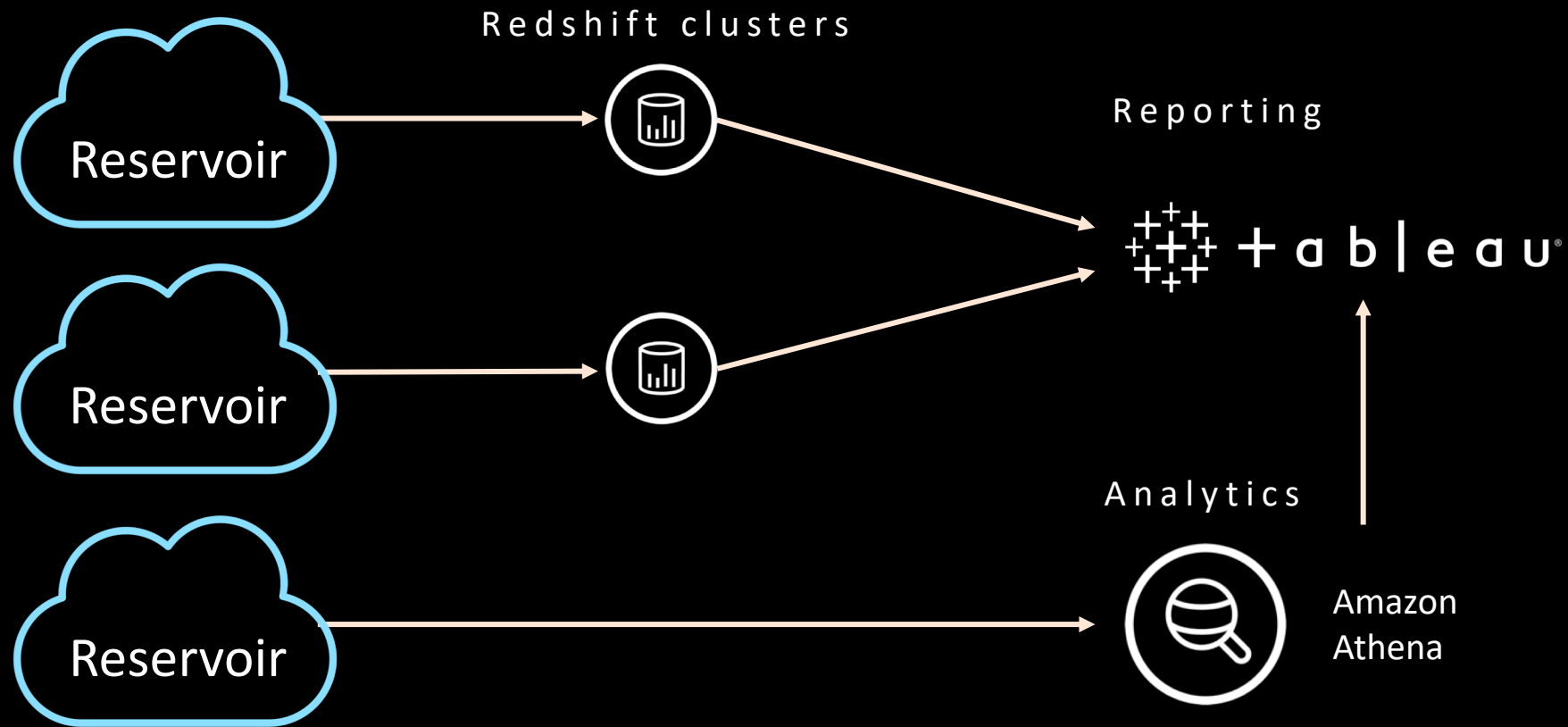
***/out***

***/parquet***

***/tmp***



# #3 BI & Reporting



# #4 Applied Intelligence

REVO<sup>LUTION</sup>





**Image recognition**

**Fraud detection**

**Moderation**

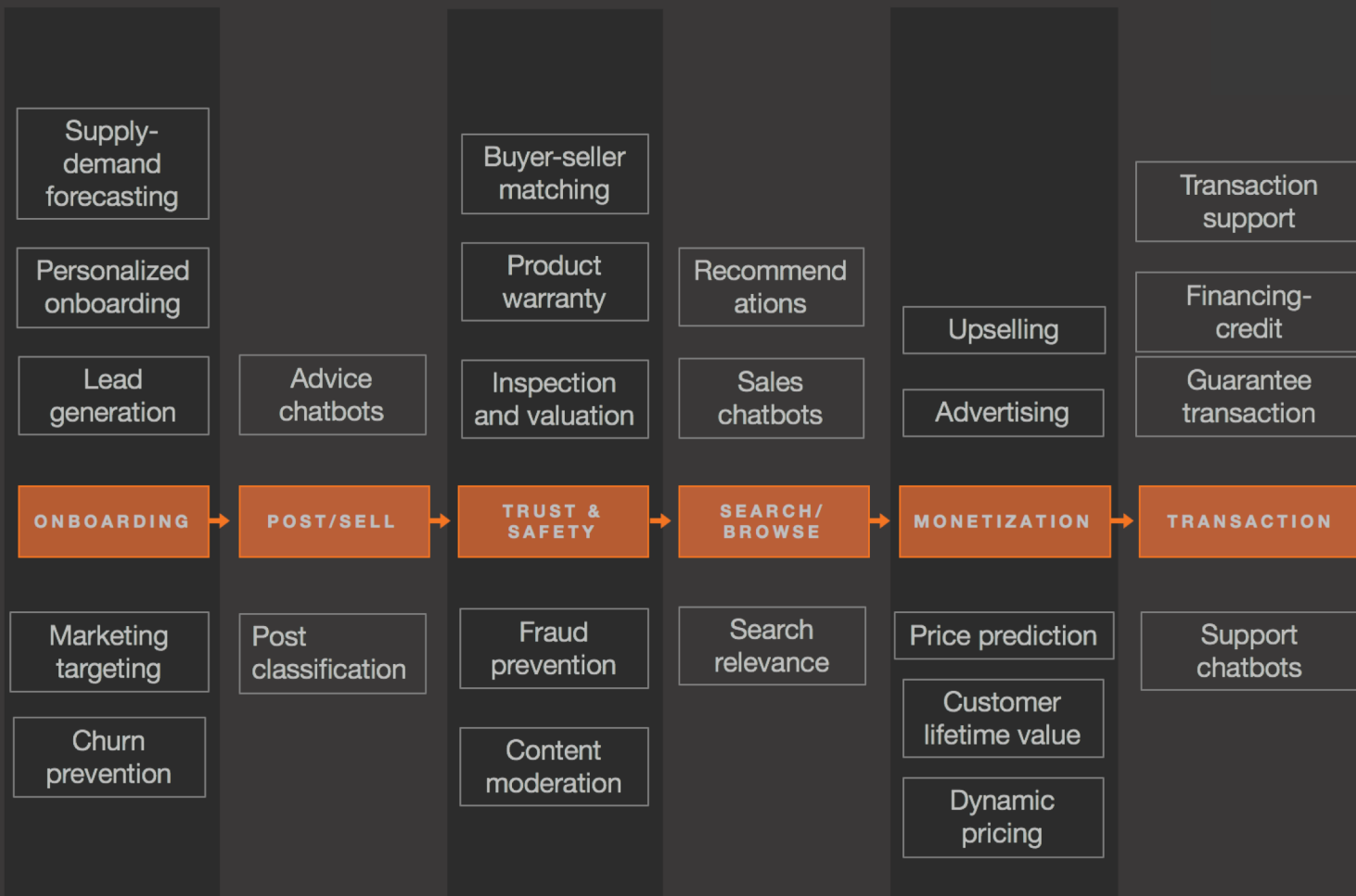
**Recommendations**

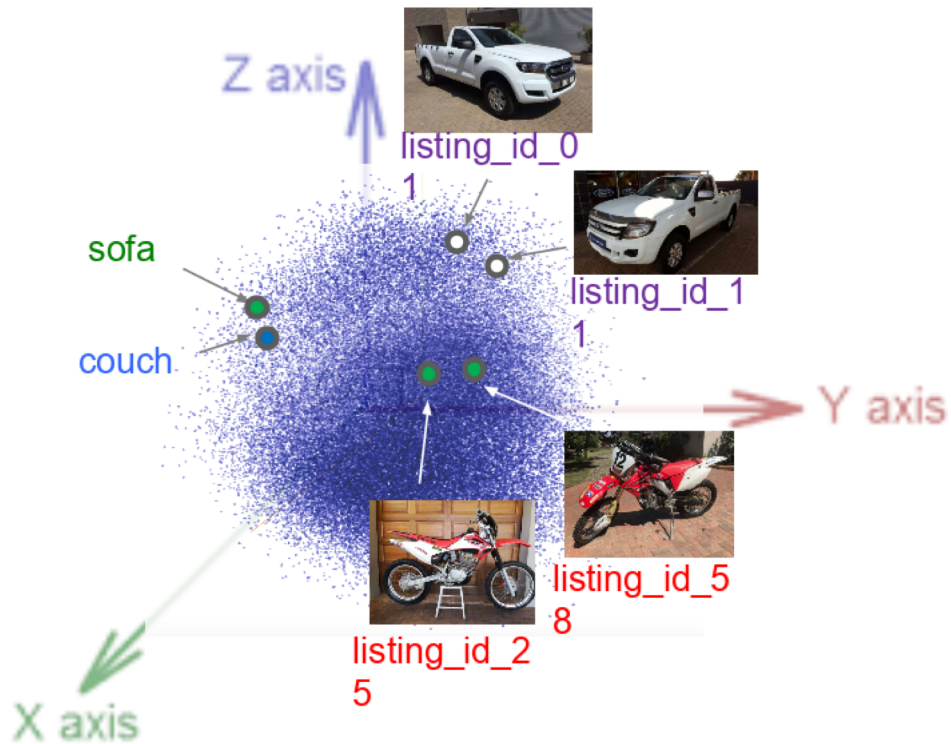
**And more...**

TOP-LINE DRIVERS  
NEW MARKETS  
INCREASE MARKET SHARE



COST REDUCTION  
CAPITAL ALLOCATION  
COST EFFICIENCY



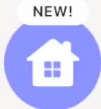


Invite

letgo



CARS



HOUSING



TECH

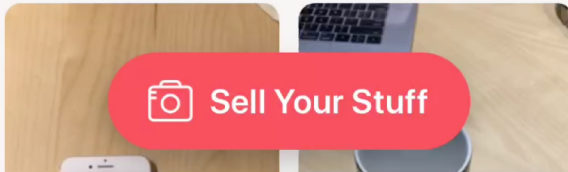
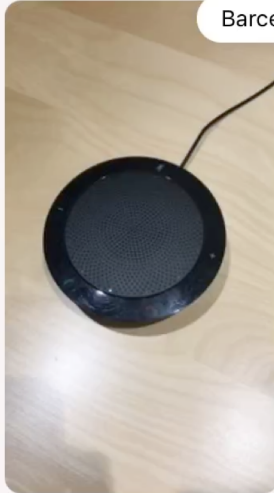


HOME



LEISURE

Barcelona



**WHAT'S  
NEXT ?**

**THANK YOU**