



# The OLX data theory of everything

Caspar Schönau

Head of Global BI

Jakub Orłowski

Data engineering manager



*'The biggest internet  
company that you have  
never heard of'*

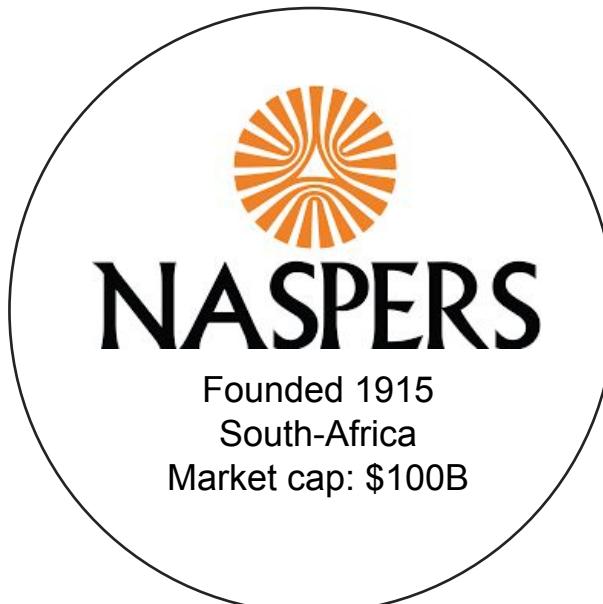
*Tencent* 腾讯



OLX GROUP



*Delivery Hero*



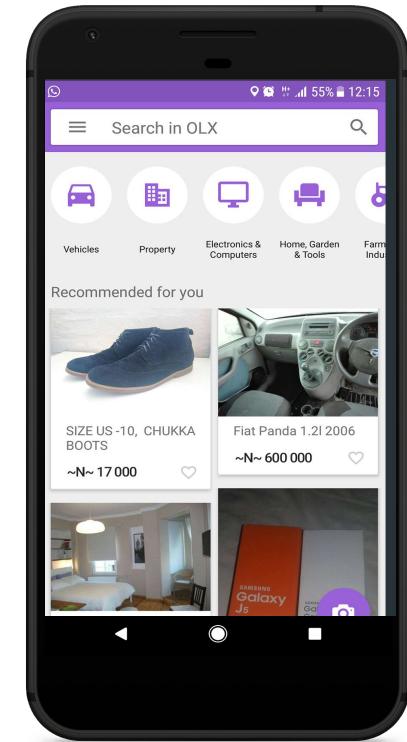
eMAG

BRAINLY

# SELLER WINS

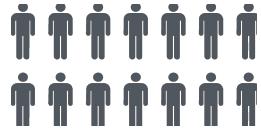


# BUYER WINS



**43**  
Countries

**35**  
Offices



**+5,000** Employees

 **+350M** MAU

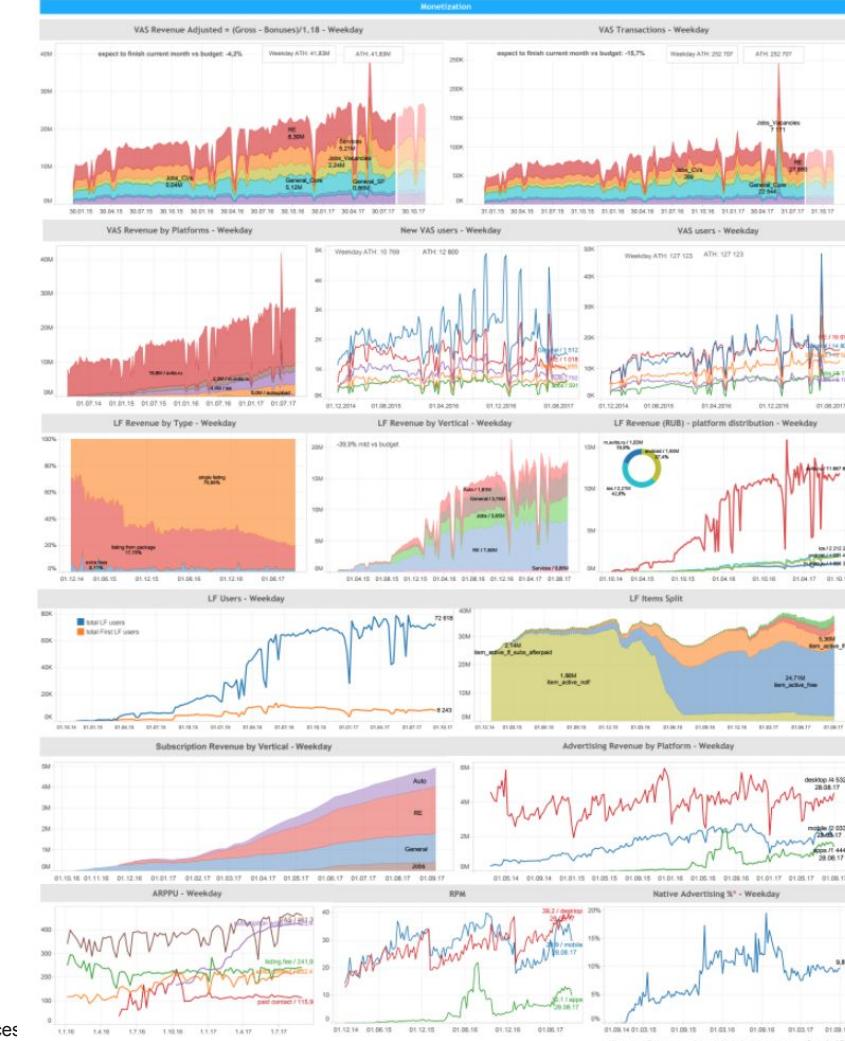
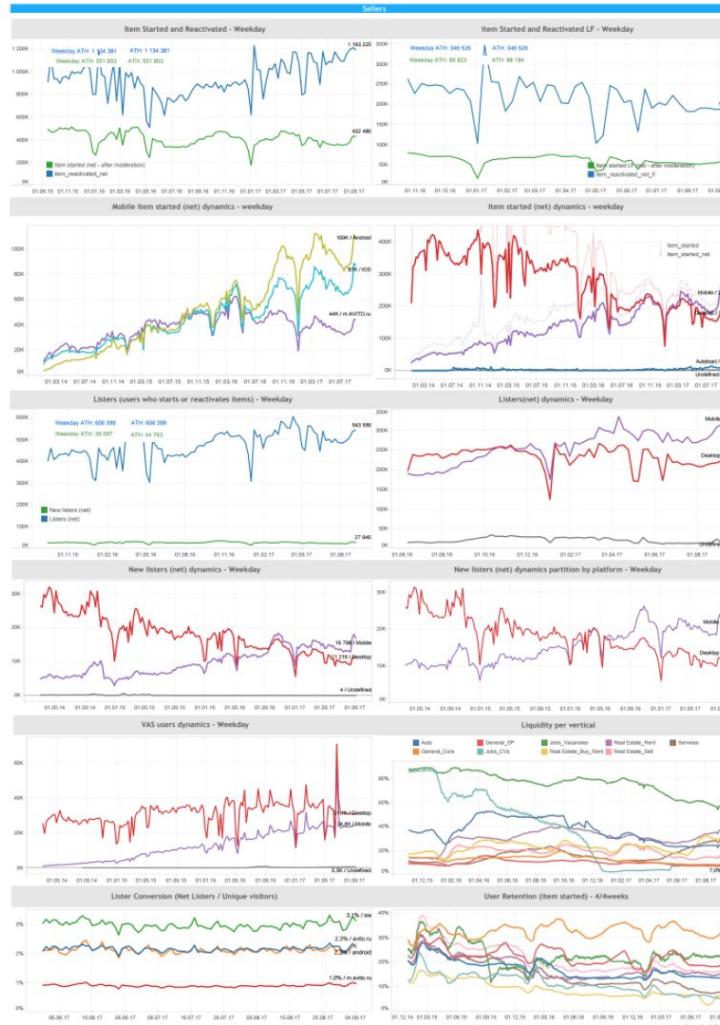
 **+4B** events / day



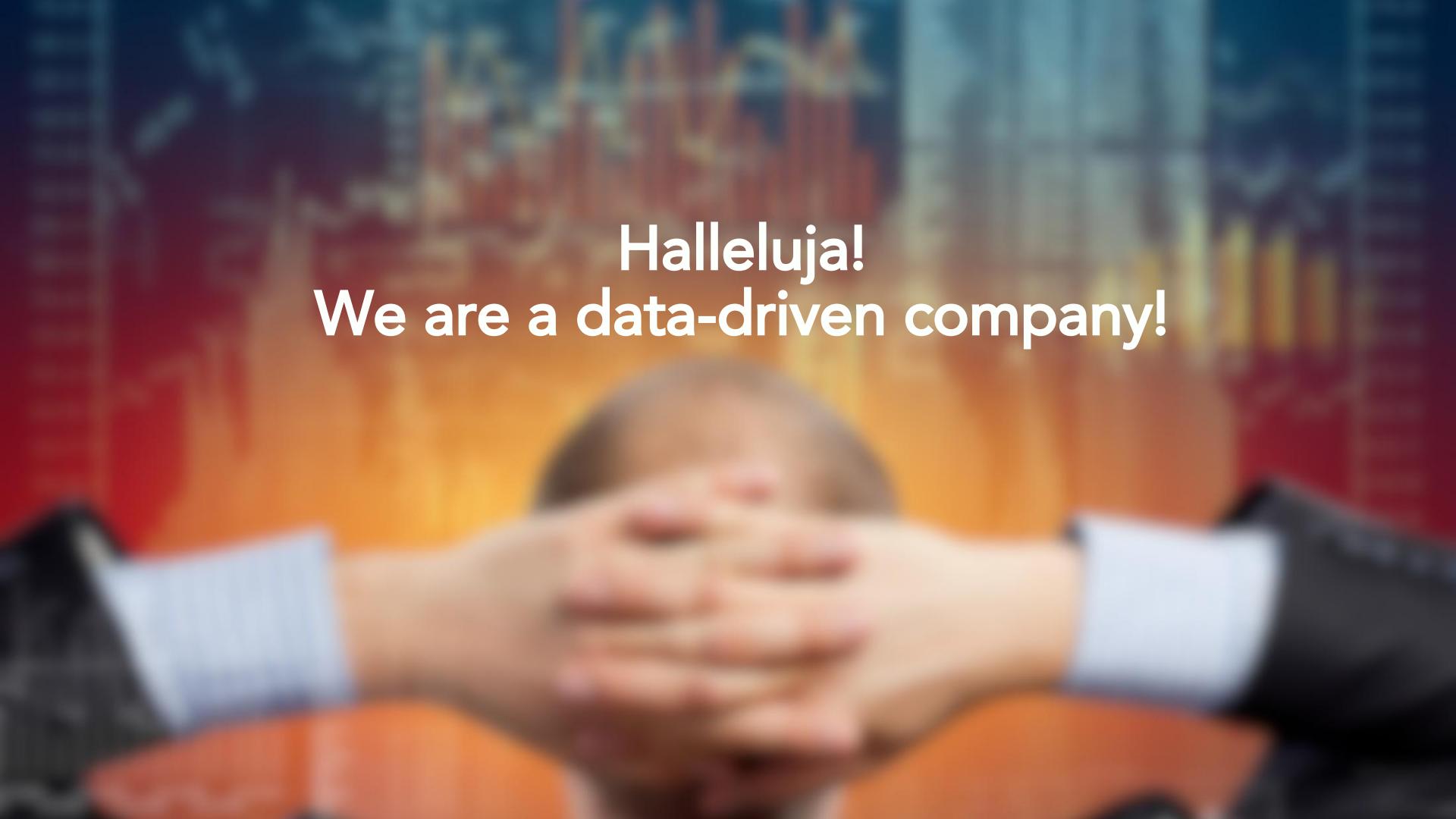
The OLX challenge:

*'Give everybody the  
data that he or she  
needs'*

(but also not much more)





A blurred background image of a person's hands holding a smartphone. The screen of the phone displays a vibrant, abstract painting with reds, blues, and yellows. The hands are positioned in the center, with fingers wrapped around the device.

Halleluja!  
We are a data-driven company!

- Kg of food in the valley
- # of apes in the valley
- Kg of food / ape needed per day
- # days of left this winter

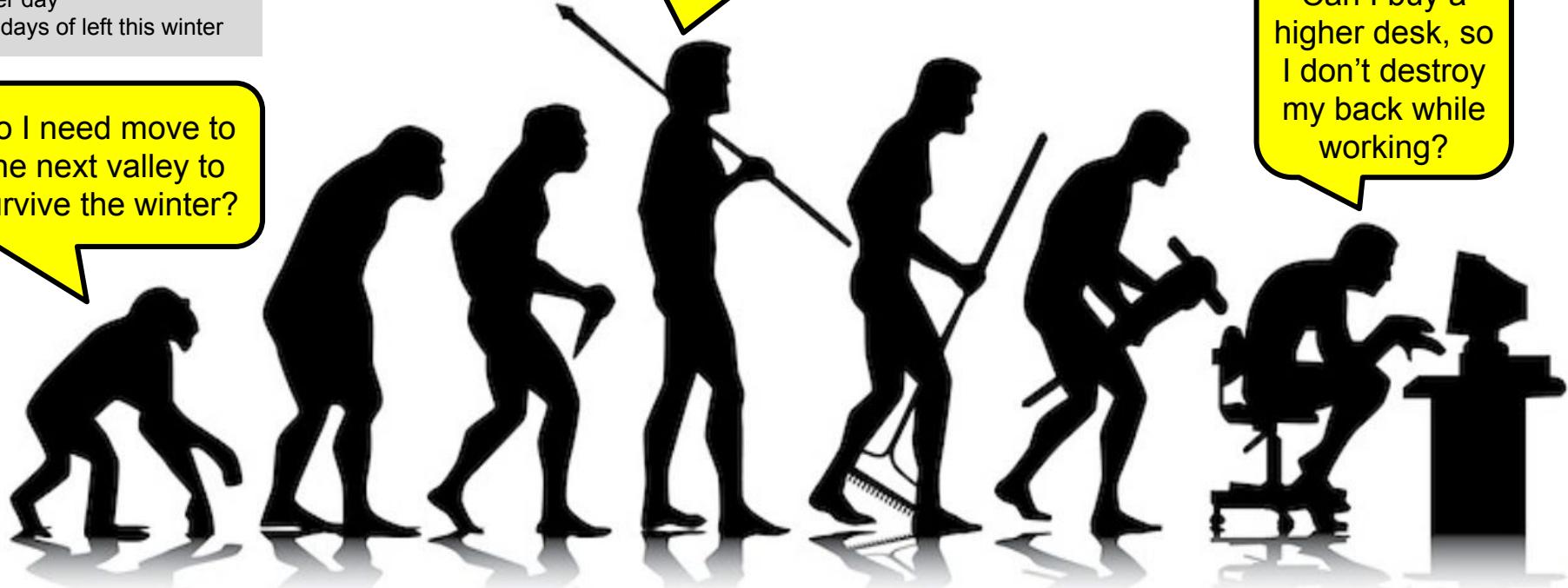
Do I need move to the next valley to survive the winter?

Can I win the javelin throw competition at the next Olympics?

- WR javelin throw in m
- PB javelin throw in m
- Time between date of PB and next Olympics

- \$ in the bank
- Price of a decent desk in \$

Can I buy a higher desk, so I don't destroy my back while working?



**What decisions are you really taking on a daily basis?  
And how does data play a role?**

- Size of the prize of online cars market in Mexico
  - Cost of success
  - Chance of success
  - Available war chest

- ROI of a typical offline marketing campaign
  - ROI of continuous online marketing
  - Expected reach of both offline and online marketing

- Average # of ads moderated by agent #6 in the last month
  - Average # of ads moderated by the team in the last month
  - Error rate agent #6
  - Error rate of the team

- Detailed properties of all listings (# pictures, attributes, length of the title, etc)
  - All individual replies, including related buyers and seller data

- Requests per second
  - Post per second

# Do I launch a new car portal in Mexico?



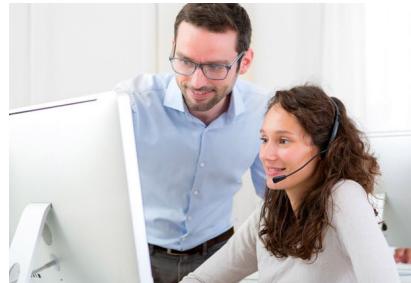
CEO

Shall I invest  
more in online  
or in offline  
marketing?



GM

Should I fire  
CS agent #253?



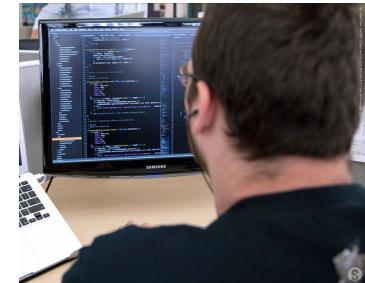
CS Manager

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



# Business analyst

Is the platform still online?



# Infra engineer

The same goes for an organization like OLX.  
Which data points are really influencing your decisions?





Disc

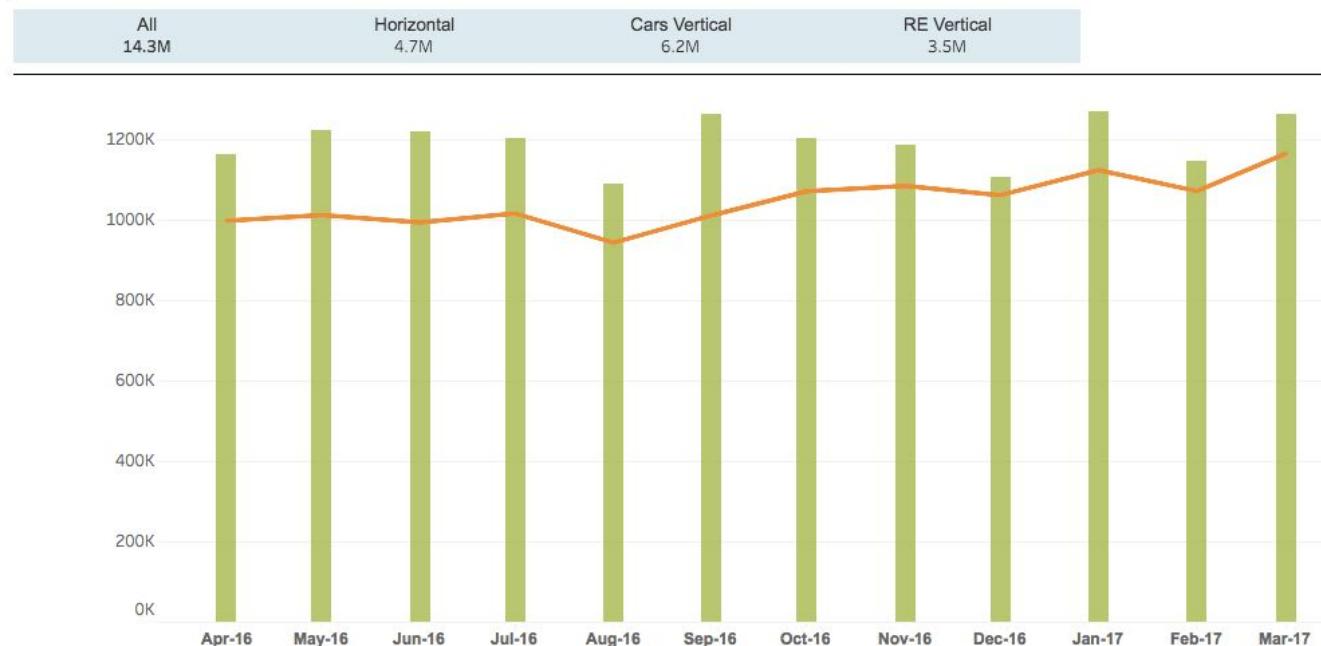
42

is broken

System:	PK_Sales																											
Package:	IS Management																											
Created:	10/14/2016 3:43:53 PM																											
Modified:	4/21/2017 12:53:49 PM																											
Notes:	Nugent Management Structure: Nugent Management Structure (Nugent), IS-Account: ISG1000008																											
Columns:	Scenario, Actual, Scenario, ISG100008, Period: FY 2017																											
Cursors:	Currency: Currency, Transaction Type: Transaction Type, Level: Level, Consolidation Percent: Consolidation Percent, PR Rates: PR Rates, YTD: YTD, Measure: Monthly Value																											
Other rows:	Filter columns:																											
		Actual												Budget														
		FY 2017	2016-04	2016-05	2016-06	2016-07	2016-08	2016-09	2016-10	2016-11	2016-12	2017-01	2017-02	2017-03	FY 2017	2016-04	2016-05	2016-06	2016-07	2016-08	2016-10	2016-11	2016-12	2017-01	2017-02	2017-03		
Classifieds managed	TOTAL REVENUE	388,991,353	29,752,207	29,236,519	30,889,393	30,527,947	33,812,124	33,716,726	34,441,518	36,553,599	36,747,988	36,517,727	39,965,728	38,058,078	26,167,703	26,750,418	27,146,963	27,187,764	28,348,653	30,926,589	34,044,987	37,361,365	33,671,680	32,021,275	34,974,765	37,089,339	37,099,411	38,428,460
	e-Commerce revenue	317,137,880	25,372,880	28,287,091	26,585,849	26,507,463	27,957,180	26,718,292	27,238,029	27,764,764	26,460,807	27,542,452	34,911,648	320,931,060	22,402,735	22,598,162	23,218,711	23,382,779	24,369,520	26,429,215	29,077,246	30,116,279	28,605,869	26,990,219	29,832,082	33,488,81		
	Classifieds revenue	315,477,251	25,318,026	28,286,550	26,596,580	26,398,848	28,845,180	25,586,012	27,118,100	27,064,983	26,286,110	27,380,544	26,980,888	34,780,658	320,082,893	22,234,792	23,525,209	23,248,434	23,870,450	24,805,240	26,429,215	29,077,246	30,116,279	28,605,869	26,990,219	29,832,082	33,488,81	
	Classifieds - Subscriptions	28,040,793	3,460,034	2,062,808	3,755,162	3,742,680	3,434,463	3,194,016	3,630,774	3,139,937	3,112,899	3,745,045	3,441,464	3,140,694	17,792,428	3,760,266	3,777,089	3,854,173	3,677,234	3,818,059	3,569,444	3,341,614	3,651,674	3,644,742	3,900,877	3,276,119	4,414,054	
	Classifieds - Listing fees	157,404,314	11,600,667	12,101,602	12,566,647	12,808,938	13,161,409	13,682,596	18,586,578	18,699,840	12,707,344	18,757,652	18,074,452	17,018,514	26,218,298	18,869,014	10,040,622	20,808,025	29,942,479	20,799,701	20,531,030	24,625,444	25,480,109	28,956,425	21,938,100	24,633,789	27,638,585	
	Classifieds - Promotional fees	128,580,437	9,146,943	9,067,584	9,739,092	9,292,669	10,246,196	10,625,726	10,915,024	11,202,618	11,460,150	10,864,996	11,468,494	14,570,042	22,943,763	1,699,429	1,702,409	1,624,861	1,644,371	1,641,164	2,114,415	2,136,907	2,009,890	2,145,414	1,998,017	2,145,414		
	Classifieds - Success fees	66,508	5,450	5,057	4,538	4,363	4,019	4,672	5,792	8,589	5,976	6,076	6,219	73,069	6,082	6,265	6,937	6,052	5,992	5,812	5,928	6,047	6,264	6,264	6,264			
	Payments revenue	2,040	2,054	2,314	2,314	2,314	2,314	2,314	2,314	2,314	2,314	2,314	2,314	2,314	0	0	0	0	0	0	0	0	0	0	0	0		
	E-shop platforms revenue	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Total sales of goods	102,820	21,094	36,033	73,330	23,866	23,604	146,152	18,001	142,696	19,931	18,033	18,086	16,295	0	0	0	0	0	0	0	0	0	0	0	0		
Africa	Sales of financial products	722	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Agency revenue on sales of courier services	2,112,205	85,502	84,261	92,462	92,479	145,504	178,452	131,304	355,446	214,578	220,278	226,054	237,219	795,644	83,319	63,335	60,683	57,803	59,568	64,381	68,092	89,285	74,523	89,440	75,870		
	Advertising revenue	74,489,339	5,312,367	5,427,310	7,796,051	5,810,348	2,722,972	6,357,208	6,799,778	6,314,229	6,688,528	5,584,481	6,052,484	6,081,274	44,018,491	3,073,303	3,179,516	3,088,519	3,085,073	3,158,429	3,594,687	4,022,408	4,198,398	3,895,522	4,126,544	4,562,423		
	Advertising - Barter	85,682	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Advertising - AdSense	22,794,233	2,595,359	2,680,302	2,704,467	2,738,252	2,422,279	2,913,303	2,193,479	2,892,537	2,015,589	2,718,376	1,799,507	2,169,587	8,768,634	636,639	672,250	629,897	397,489	593,021	745,962	786,388	821,872	792,200	791,263	794,733	902,987	
	Advertising - Banners	47,715,012	3,367,422	3,345,640	3,858,938	3,836,870	4,824,208	4,296,187	4,274,134	4,085,248	4,274,793	4,342,540	3,696,647	4,412,312	5,341,390	411,899	409,387	412,348	363,592	450,760	492,932	534,537	415,397	462,619	419,429	4,216,733		
	Advertising - Other advertising	4,079,776	218,299	206,394	225,698	221,218	199,189	347,818	372,128	306,163	402,249	423,365	365,314	382,862	293,980,860	20,010,717	20,978,880	20,647,274	21,234,467	23,995,145	24,711,518	28,816,399	2,662,925	2,776,863	2,899,978	3,112,214		
	Total other revenues	7,364,385	1,067,357	510,116	314,493	210,177	372,017	446,603	472,870	504,804	547,638	297,976	2,072,141	2,011,793	2,011,793	2,011,793	2,011,793	2,011,793	2,011,793	2,011,793	2,011,793	2,011,793	2,011,793	2,011,793	2,011,793			
	TOTAL REVENUE	5,462,388	576,636	472,396	236,677	268,676	357,244	424,276	347,017	543,463	567,494	406,582	453,279	454,254	1,138,112	44,390	51,250	61,135	45,771	76,405	85,734	104,739	218,522	201,179	122,382	140,388		
Corporate - Africa	e-Commerce revenue	200,170	577	(10,845)	70,158	14,424	3,507	18,524	38,158	37,423	32,470	32,228	30,672	38,549	705,703	25,595	20,636	29,771	23,597	41,681	51,377	68,354	82,391	64,266	82,277	90,234	126,295	
	Classifieds revenue	200,178	577	(10,845)	70,124	20,222	7,749	18,524	38,158	37,423	32,470	32,228	30,672	38,549	705,703	25,595	20,636	29,771	23,597	41,681	51,377	68,354	82,391	64,266	82,277	90,234	126,295	
	Classifieds - Subscriptions	(116,750)	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Classifieds - Listing fees	316,931	571	5,910	70,124	20,071	7,902	18,524	38,158	37,423	32,470	32,228	30,672	38,549	705,703	25,595	20,636	29,771	23,597	41,681	51,377	68,354	82,391	64,266	82,277	90,234	126,295	
	Classifieds - Success fees	83	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Payments revenue	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Total sales of goods	60	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Advertising revenue	639,651	4,448	5,548	6,773	6,209	6,244	10,058	184,016	102,879	30,933	62,560	95,044	95,037	42,7409	28,795	30,614	31,364	32,844	33,784	34,359	36,385	37,131	36,912	40,706	41,151	43,424	
	Advertising - AdSense	333,572	4,448	5,548	6,773	6,209	6,244	10,058	184,016	102,879	30,933	62,560	95,044	95,037	42,7409	28,795	30,614	31,364	32,844	33,784	34,359	36,385	37,131	36,912	40,706	41,151	43,424	
	Advertising - Other advertising	86,079	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
Kenya	Total other revenues	4,622,166	372,243	323,692	296,734	211,937	1,243	8,706	27,033	30,698	88,807	25,565	25,413	28,059	26,420	28,553	3,045	3,488	3,815	3,505	7,782	24,361	28,424	35,672	37,064	39,500		
	e-Commerce revenue	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Classifieds revenue	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Classifieds - Listing fees	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Classifieds - Success fees	10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Payments revenue	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Total sales of goods	16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
	Advertising revenue	193,888	414	368	682	652	656	3,475	66,279	25,409	21,089	28,039	26,219	26,493	46,289	2,950	3,353	4,204	3,156	3,360	3,392	4,250	4,084	4,697	4,879	5,212		
	Advertising - AdSense	203,809	414	368	682	652	656	3,475	66,279	25,409	21,089	28,039	26,219	26,493	46,289	2,950	3,353	4,204	3,156	3,360	3,392	4,250	4,084	4,697	4,879	5,212		
	Advertising - Other advertising	86,079	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
Nigeria	Total other revenues	217,603	216,523	231	0	3,273	31,976	36,223	22,538	157	4,324	0	302	2,092	0	0	0	0	0	0	0	0	0	0	0	0		
	TOTAL REVENUE	142,540	291,764	298,114	168,558	9,945	9,444	6,940	15,442	26,241	23,131	23,799	15,402	304,500	28,500	26,500	29,000	22,000	22,000	27,000	30,000	35,000	39,500	39,500	36,220			
	e-Commerce revenue	92,215	0	5,519	15,876	9,643	9,087	7,298	7,220	11,136	6,482	7,055	6,003	8,416	246,000	10,500	13,500	15,000	18,000	20,000	22,000	25,000	27,500	30,000	32,000	32,000		
	Classifieds revenue	92,215	0	5,519	15,876	9,643	9,087	7,298	7,220	11,136	6,482	7,055	6,003	8,416	246,000	10,500	13,500	15,000	18,000	20,000</td								

# Financial revenue dashboard - summary

## Business Model



## Filters

### Fiscal Year

- FY17
- FY18

### Region

Europe 

### Country

- (All)
- Angola
- Belarus
- Bosnia
- Bulgaria
- Kazakhstan
- Mozambique
- Other
- Poland
- Portugal
- Romania
- Ukraine
- Uzbekistan

### Revenue Type

- (All)
- Advertising revenue
- Classifieds revenue
- Total other revenues

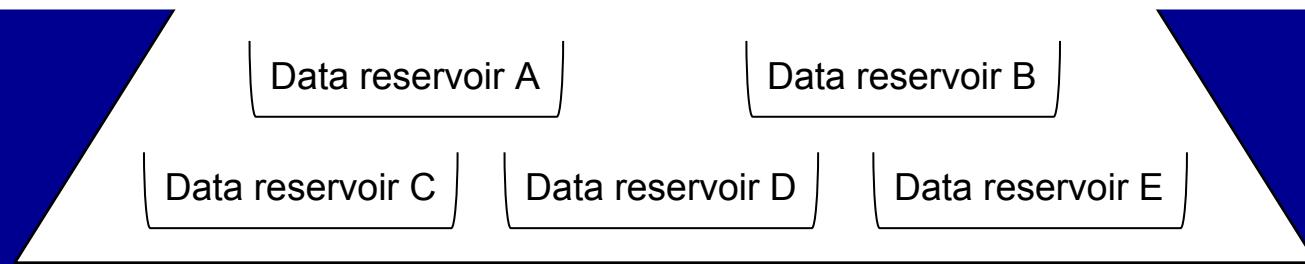
Monthly actual (USD)	1.2M	1.2M	1.2M	1.2M	1.1M	1.3M	1.2M	1.2M	1.1M	1.3M	1.1M	1.3M
Monthly budget (USD)	1.0M	1.0M	1.0M	1.0M	0.9M	1.0M	1.1M	1.1M	1.1M	1.1M	1.1M	1.2M
% vs. budget	16.5%	20.8%	22.8%	18.2%	15.4%	24.6%	12.1%	9.4%	4.1%	12.9%	6.9%	8.5%
Cumulative actual (USD)	1.2M	2.4M	3.6M	4.8M	5.9M	7.2M	8.4M	9.6M	10.7M	11.9M	13.1M	14.3M
Cumulative budget (US..)	1.0M	2.0M	3.0M	4.0M	5.0M	6.0M	7.1M	8.1M	9.2M	10.3M	11.4M	12.6M
% vs. budget cumulative	16.5%	18.6%	20.0%	19.5%	18.8%	19.8%	18.6%	17.4%	15.8%	15.5%	14.7%	14.1%

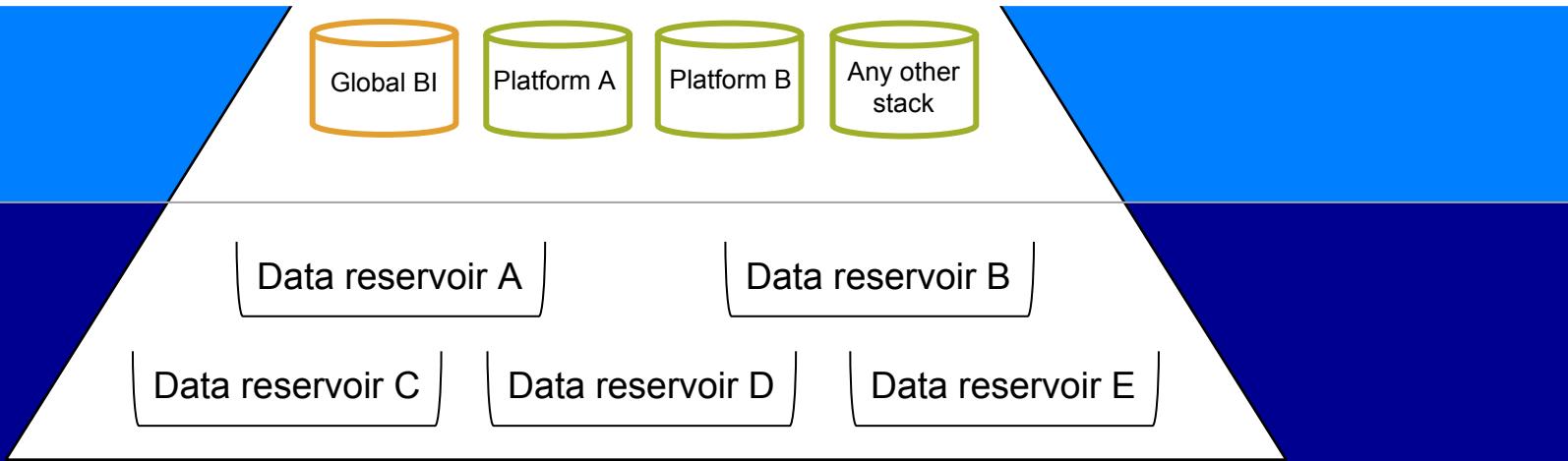
# The OLX data iceberg model

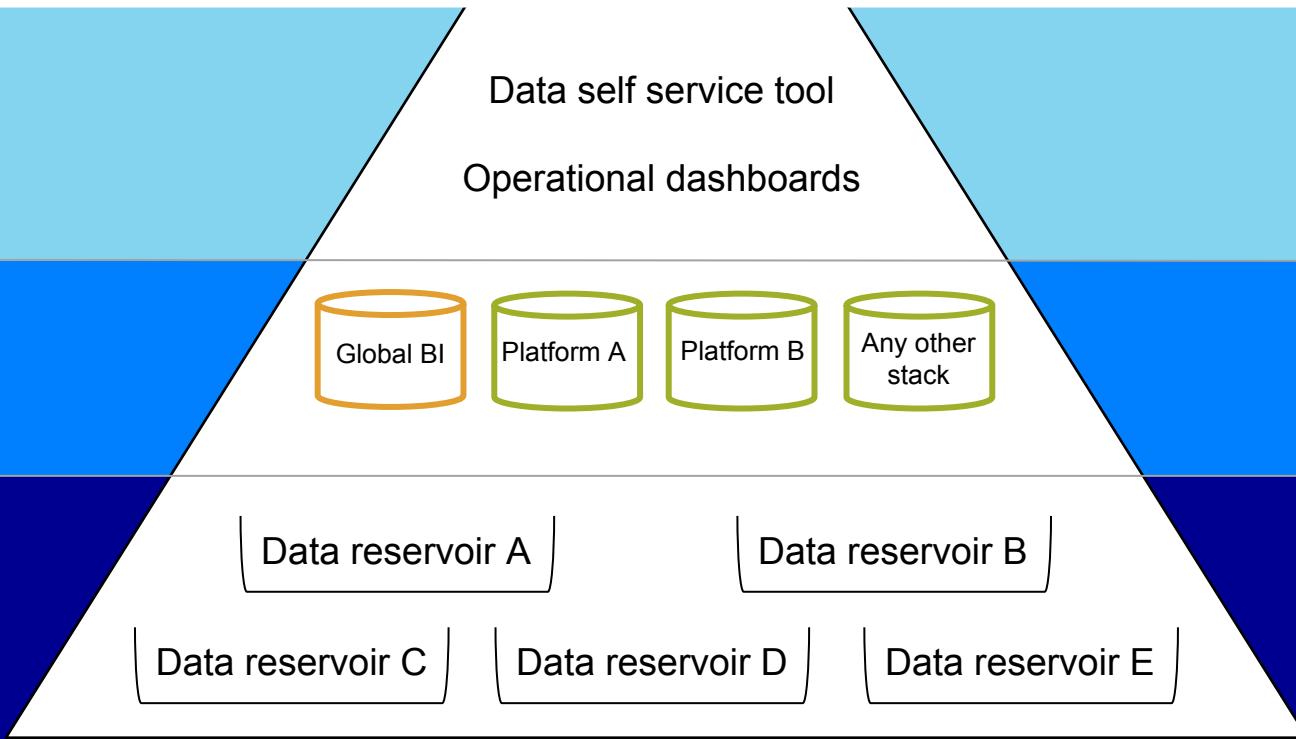
*The old school way of providing data...*

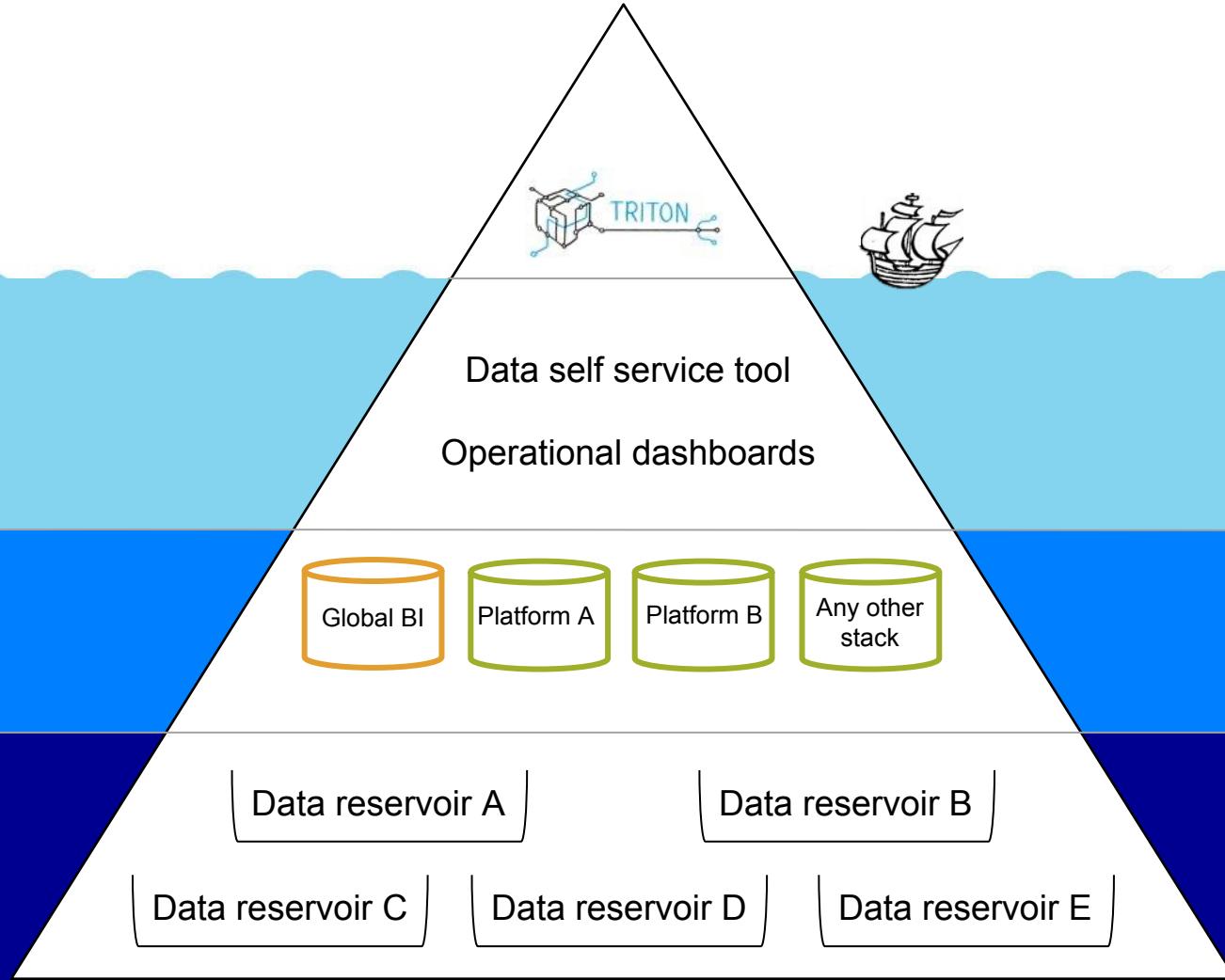


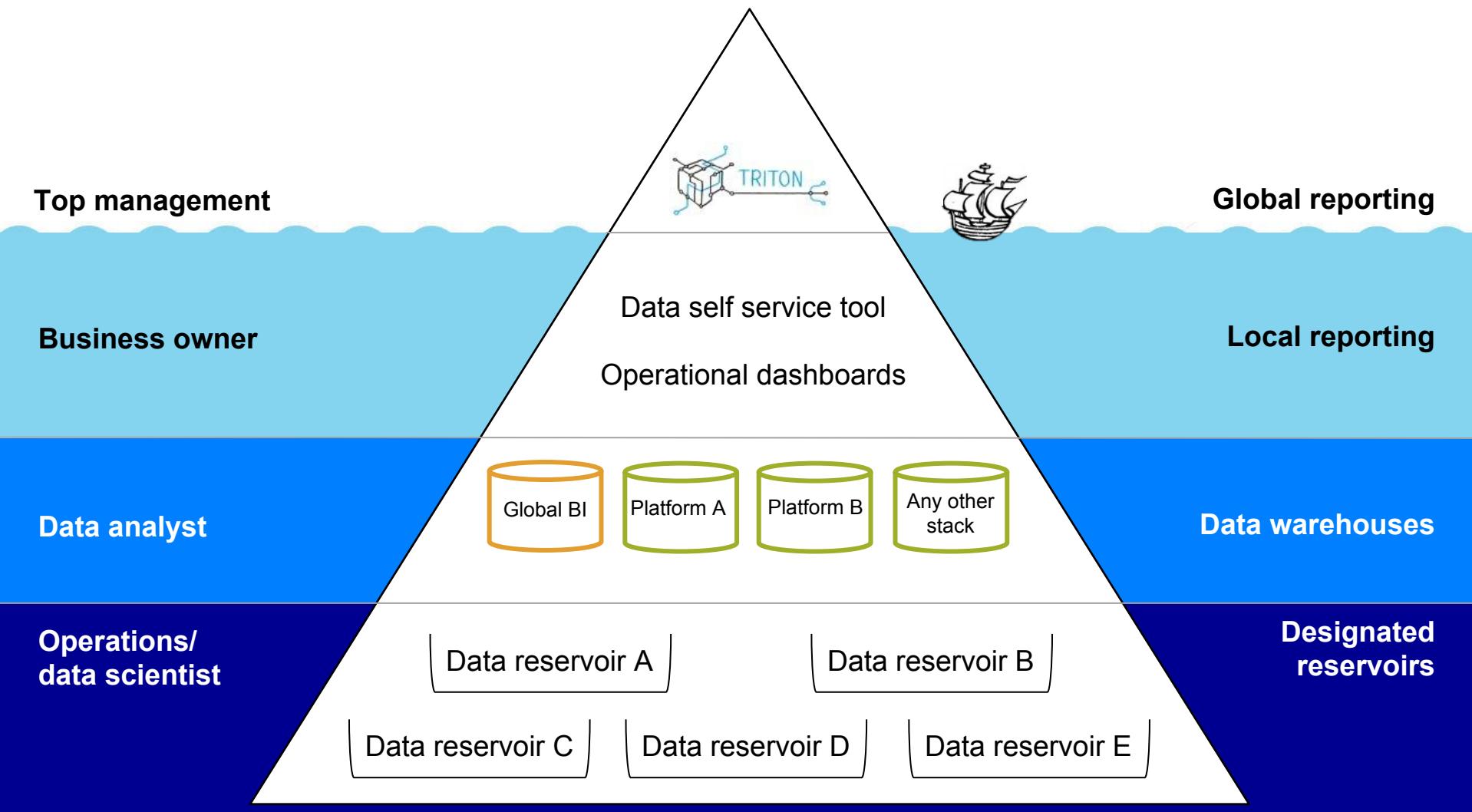
OLX data lake

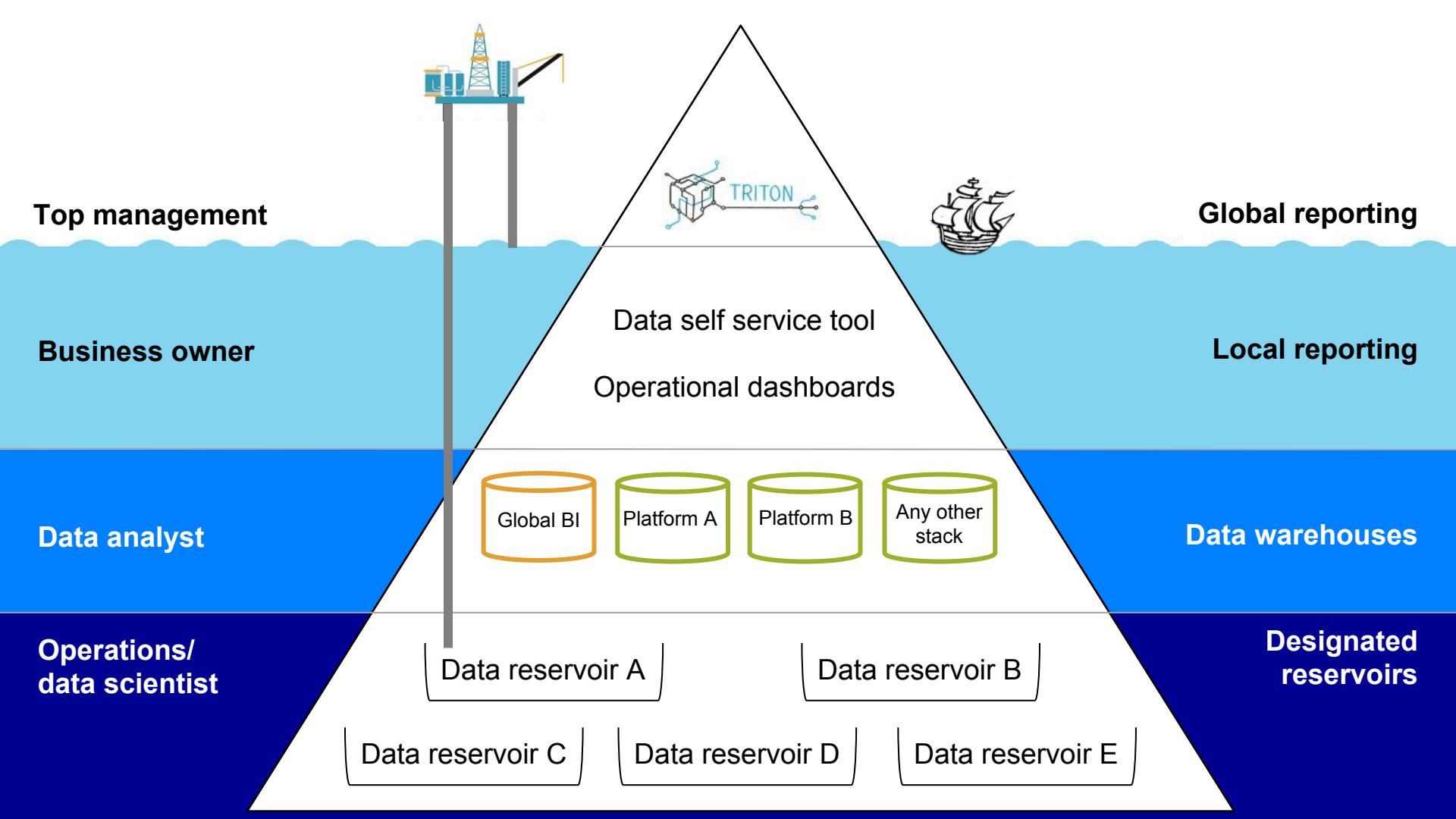












# 4 “V”s of Big Data

top challenges of data processing

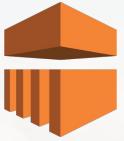




Amazon  
S3



Amazon  
Redshift



Amazon  
EMR

# Volume

data at rest





Amazon  
Kinesis

Amazon  
SQS



# Velocity

## data in flight

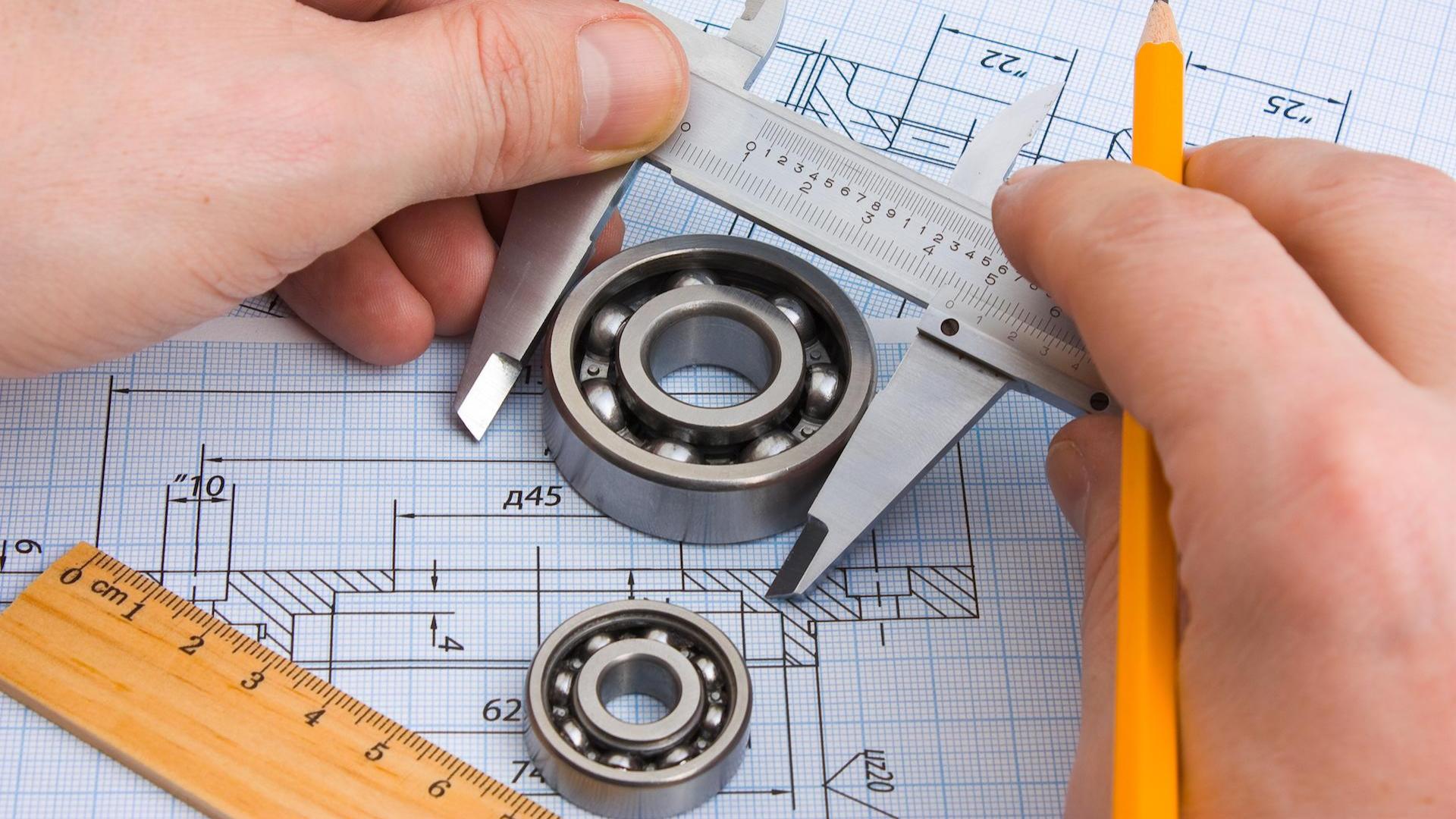




AWS Glue

# Variety

data in many shapes



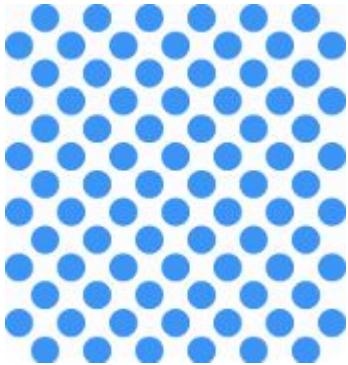


Amazon  
Athena

# Veracity

## data quality

## Volume



### Data at rest

Terabytes of existing, historical data that needs to be stored for extended period

## Velocity



### Data in flight

Streaming, near real-time data sources, short reaction and computation time required

## Variety



### Data types

Platform databases, user behaviour, infrastructure monitoring, pictures

## Veracity



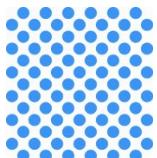
### Data quality

Tracking coverage, downtime, implementation errors, schema changes

## Data democratization

### Data access

Volume



Velocity



### Data understanding

Variety



Veracity



producers



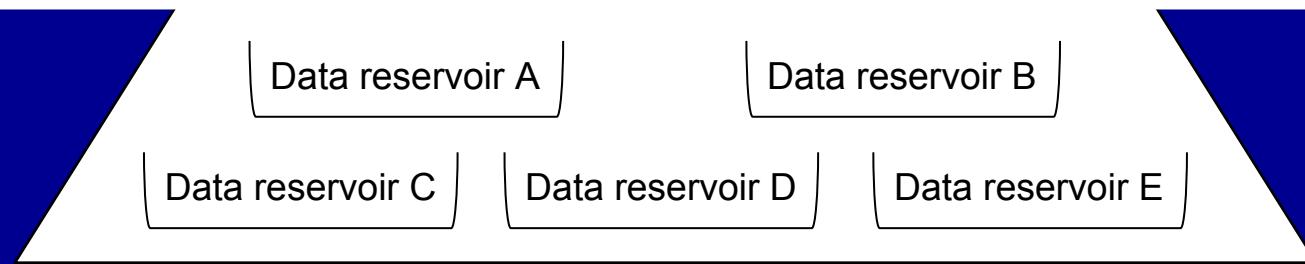
consumers



compliance

# Data democratization at OLX

architecture overview



producers



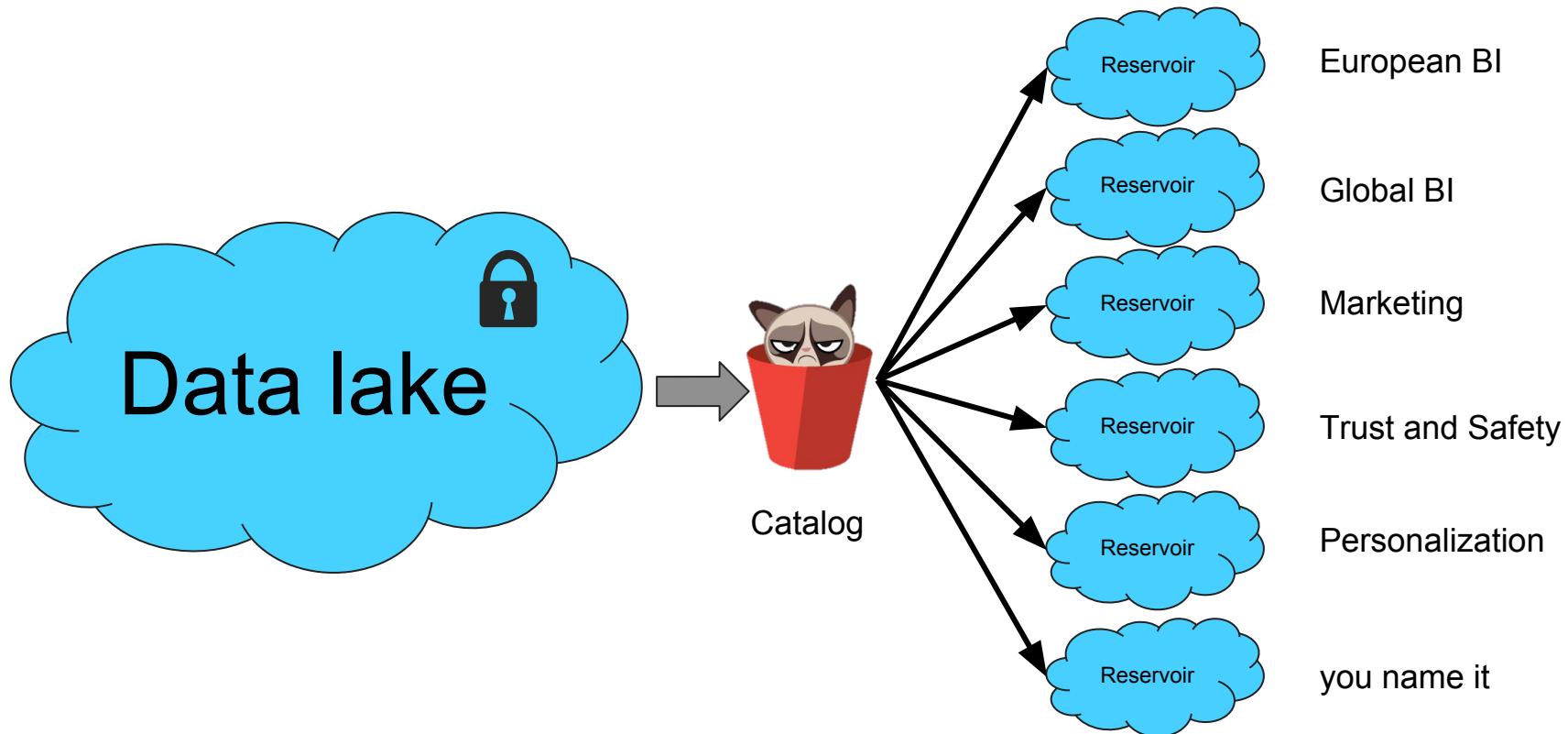
consumers

compliance

The OLX challenge:

*'Give everybody the  
data that he or she  
needs'*

(but also not much more)



# Reservoir structure - S3 bucket



**/in**

Incoming raw files in JSON format.  
Files appear right after saving in data lake.

**/out**

Outgoing raw files in JSON format.  
Files will be copied to data lake in real time.

**/parquet**

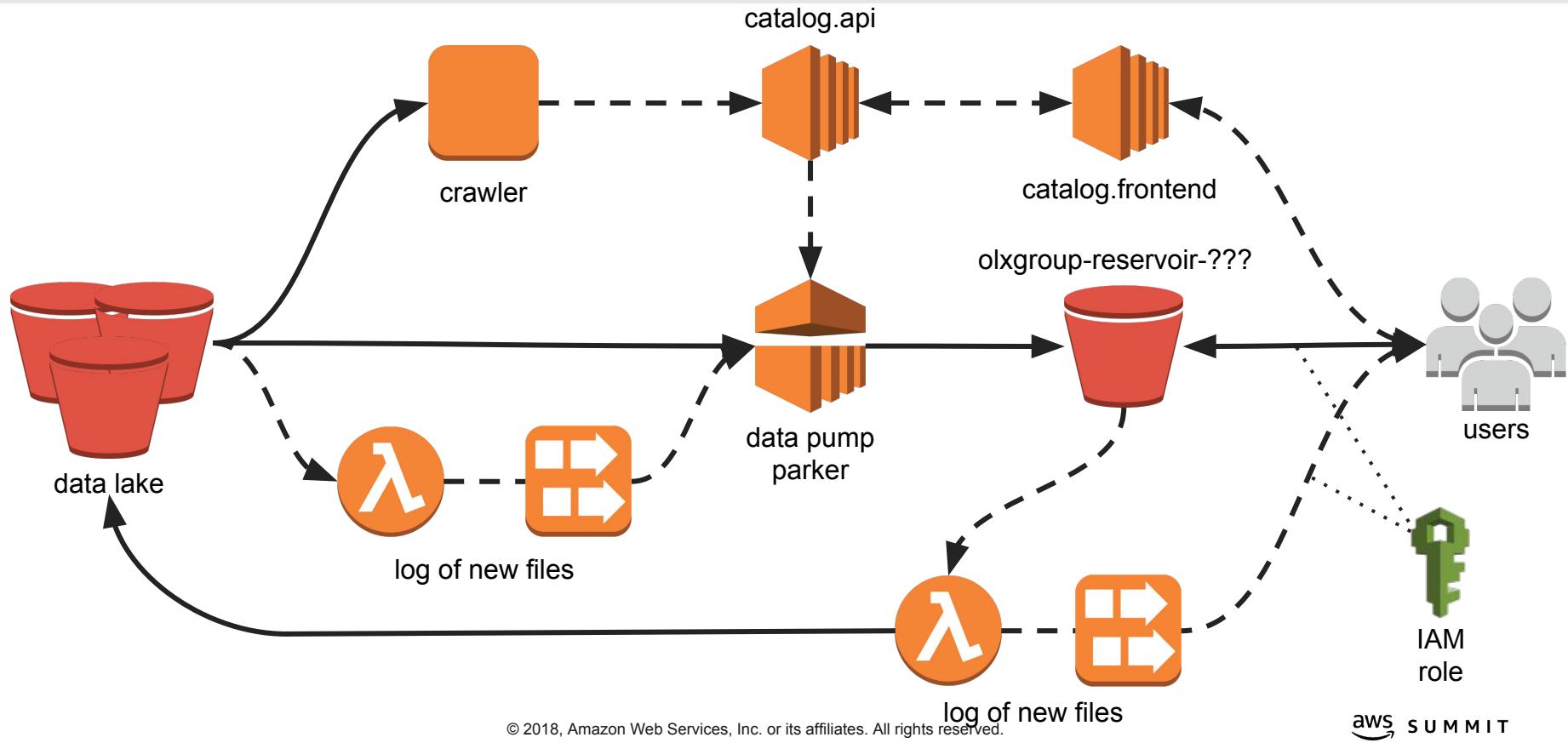
Incoming files in Parquet format.  
Files are partitioned hourly.

**/tmp**

Folder for temporary files, can be used for  
higher-level data processing apps.

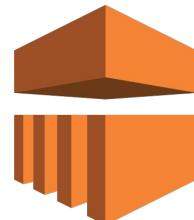
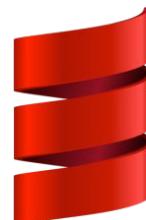
Reservoirs have individual data retention policies attached.

# Architecture overview



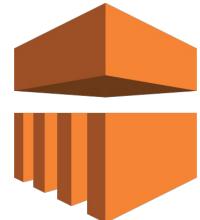
# Data Pump – raw data pre-processor

- **Technology:** Scala, Spark Streaming, EMR
- **Type:** CPU intensive
- **Cluster:**
  - Master: c4.xlarge
  - Core: 15 \* c4.xlarge
  - Spot
- **Throughput:**
  - 220K files / hour
  - 280GB (compressed) / hour
- **Price:** \$1500 / month



# Parker – raw-to-parquet converter

- **Technology:** Python3, PySpark, EMR
- **Type:** Memory intensive
- **Cluster:**
  - Master: r4.xlarge
  - Core: 10 \* r4.xlarge
  - Auto scaling
  - Spot
- **Price:** \$1100 / month





# 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

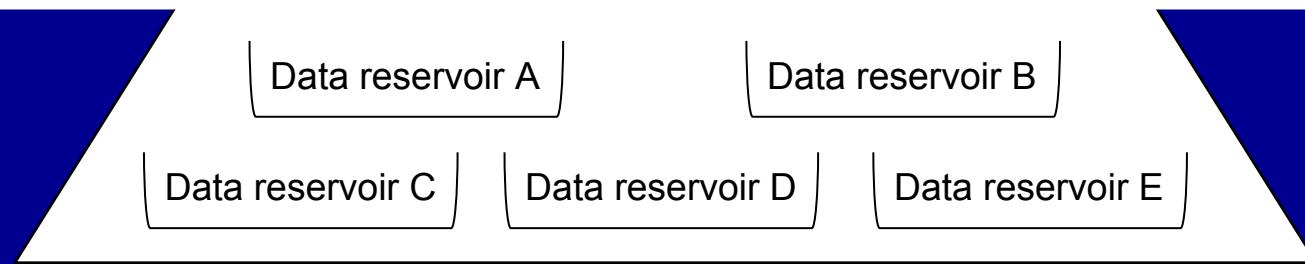
- By health
- By amount of sources

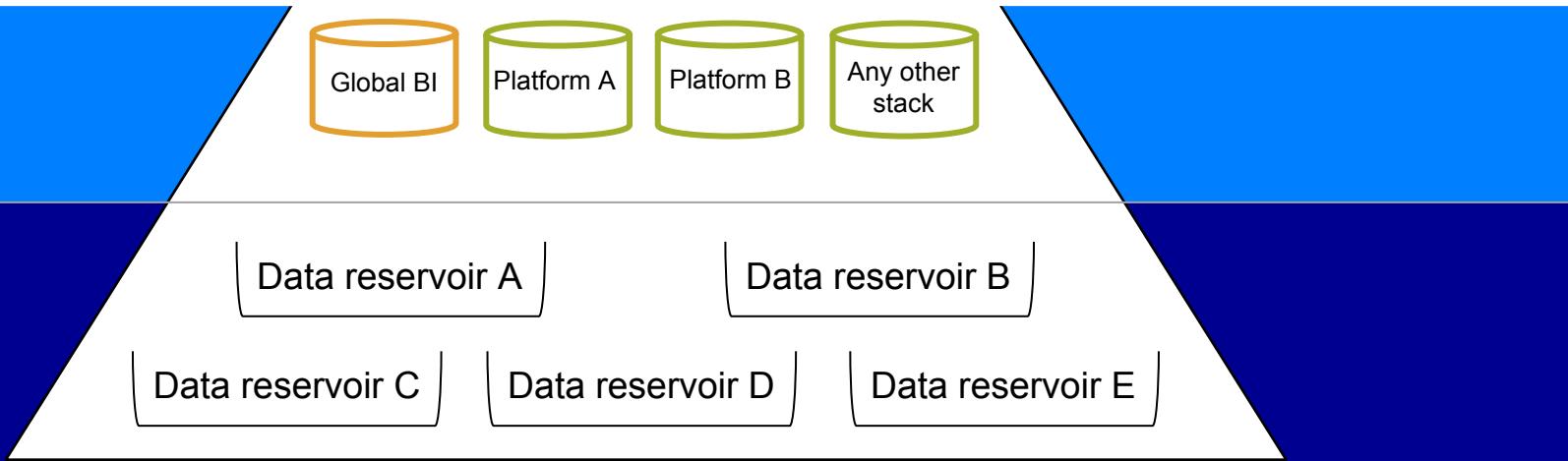
## Top Reservoirs

- By health
- By amount of subscriptions

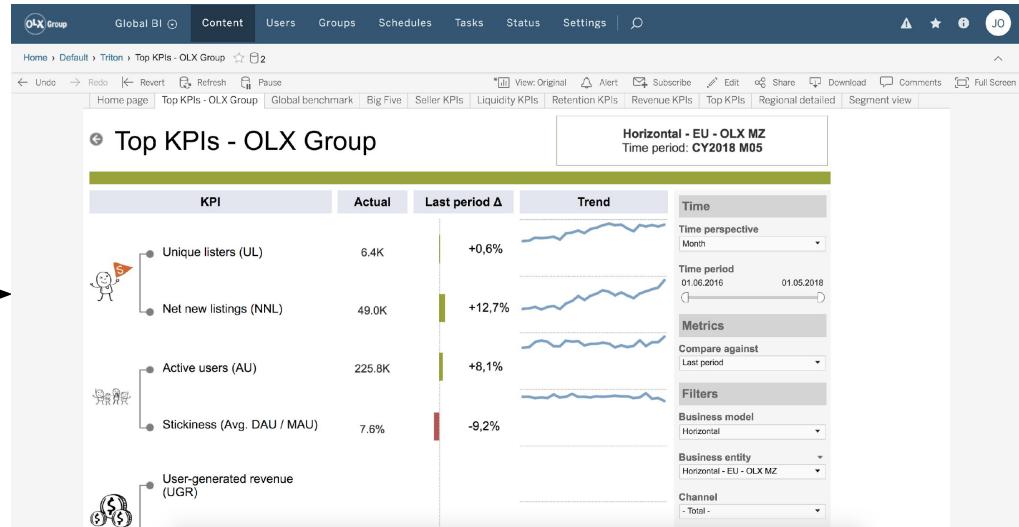
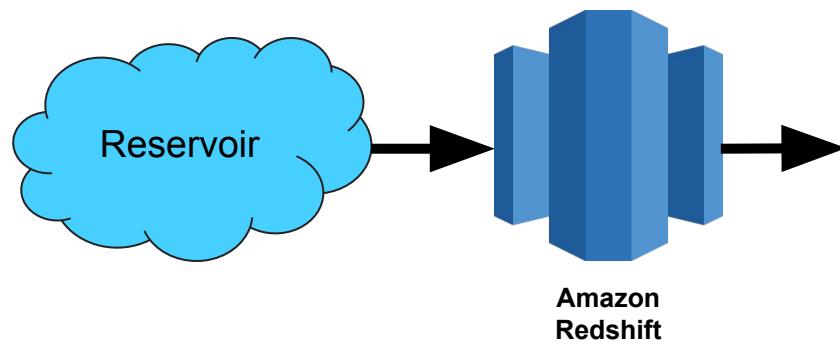
## Top Subscribers

- By health
- By amount of subscriptions

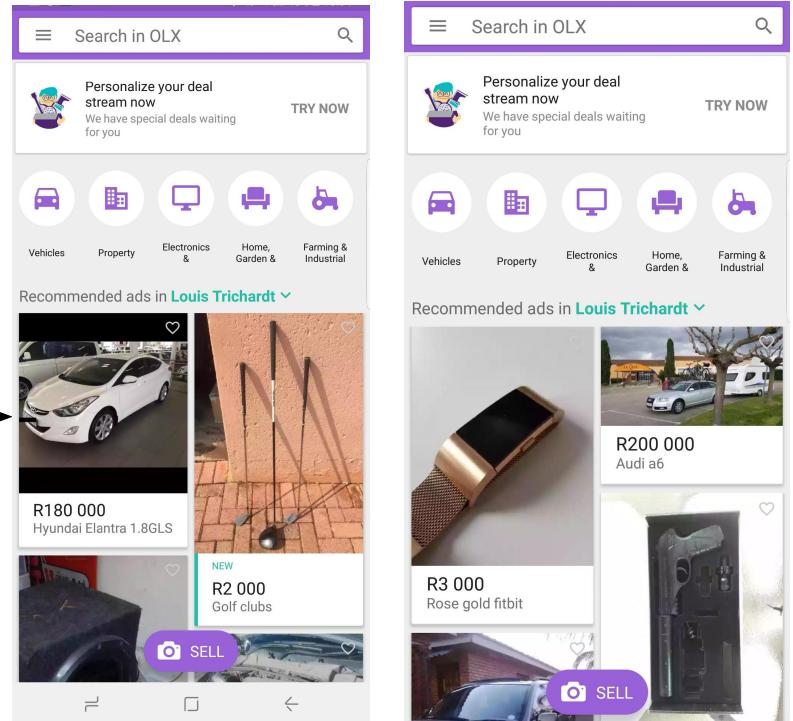
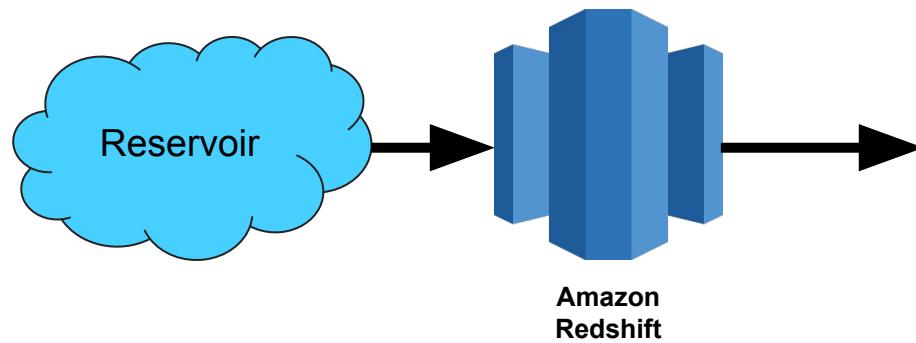




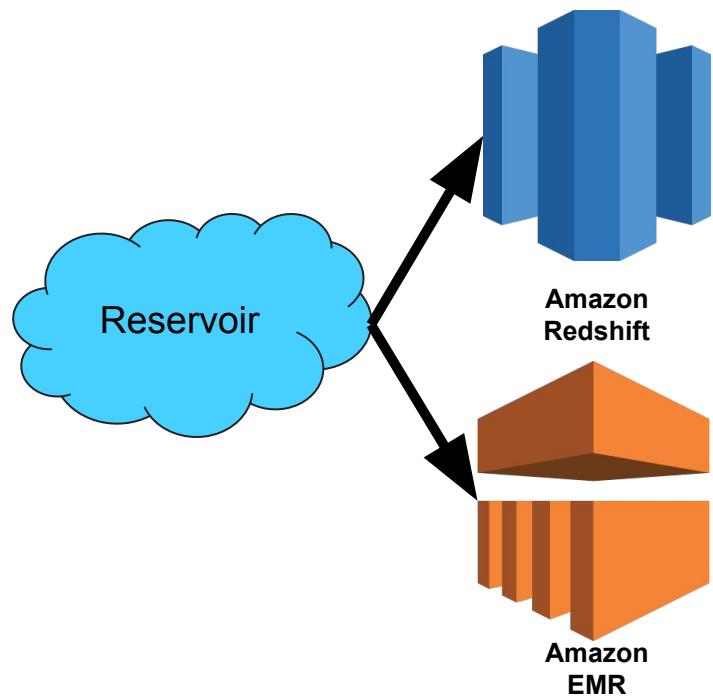
# Usage examples – Business intelligence



# Usage examples – Personalization & Relevance



# Usage examples – User communication



Edit treatment 8040436

1 Basics > 2 Targeting > 3 Timing > 4 Content > 5 Review and save

Basic Advanced

Prev Completed

**Basics**

- ✓ Name  
L2 buying ideas | L2 suggestions
- ✓ Description  
Indicate the number of deals in 2nd & 3rd recommended L2 categories
- ✓ Next best action  
Make a reply
- ✓ Priority  
Medium

**Targeting**

- ✓ Legacy who are novice and are all buyers

**Content**

- ✓ Subject
- ✓ Message
- ✓ Landing page  
rec(b2b,category\_l2,category\_l2,listing\_url\_app)
- ✓ Push type  
Push with 2 blocks

**Timing**

- ✓ Trigger  
first activity of user
- ✓ Delay  
39 day(s)
- ✓ Repetitions  
32767 time(s)  
every 42 day(s)

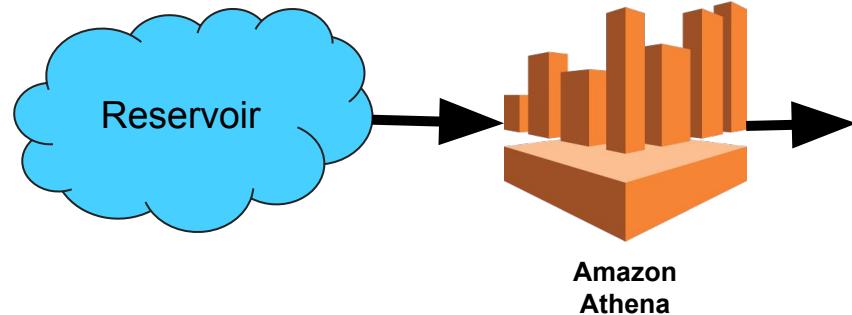
OLX Ещё не нашли, что искали?  
Рассмотрите свежие предложения в рубриках  
strollers и tablets • Возможно, покупка  
вашей мечты именно там!

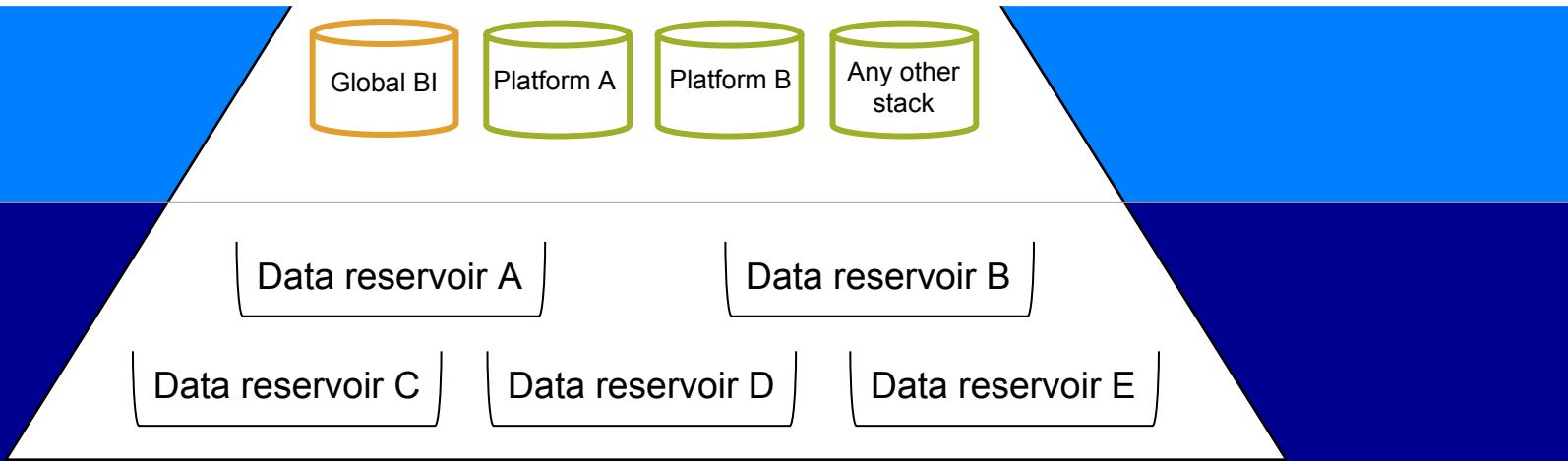
strollers tablets

Apply changes

This screenshot shows the "Edit treatment" interface for a marketing campaign. The navigation bar at the top includes steps: Basics, Targeting, Timing, Content, and Review and save. Below the steps are tabs for Basic and Advanced settings. The Content section is currently active, displaying fields for Subject, Message, Landing page (with a placeholder for a specific URL), and Push type (set to "Push with 2 blocks"). The Timing section shows a trigger for the first user activity, a delay of 39 days, and 32,767 repetitions every 42 days. A preview window on the right shows a mobile device displaying a promotional message from OLX about finding new items, with links to strollers and tablets. The bottom right corner has a green "Apply changes" button.

# Usage examples – Exploration and monitoring





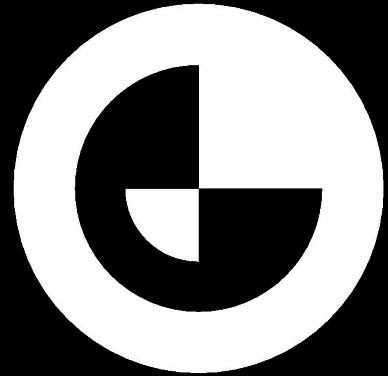
# The (data) Theory of Everything:

~ '*Everything*' is over-rated, nobody needs *everything*

# Key takeaways

---

- Not everybody needs all data
- Different stakeholders need different data solutions
- When it comes to user data, go with privacy by design and default
- Make sure to follow AWS Well-Architected framework
- Use spot instances and auto scaling where possible – it will help you focus on fault tolerance, and you will save money in return



**OLX GROUP**

twitter

blog

open roles

@olxtechberlin

[tech.olx.com](http://tech.olx.com)

[olxgroup.com](http://olxgroup.com)