



# Data Engineering interview case

*A glimpse into the world of Big Data at OLX*

Version 15

**Note:** Throughout this presentation we use the term **NNL** = 'net new listings'

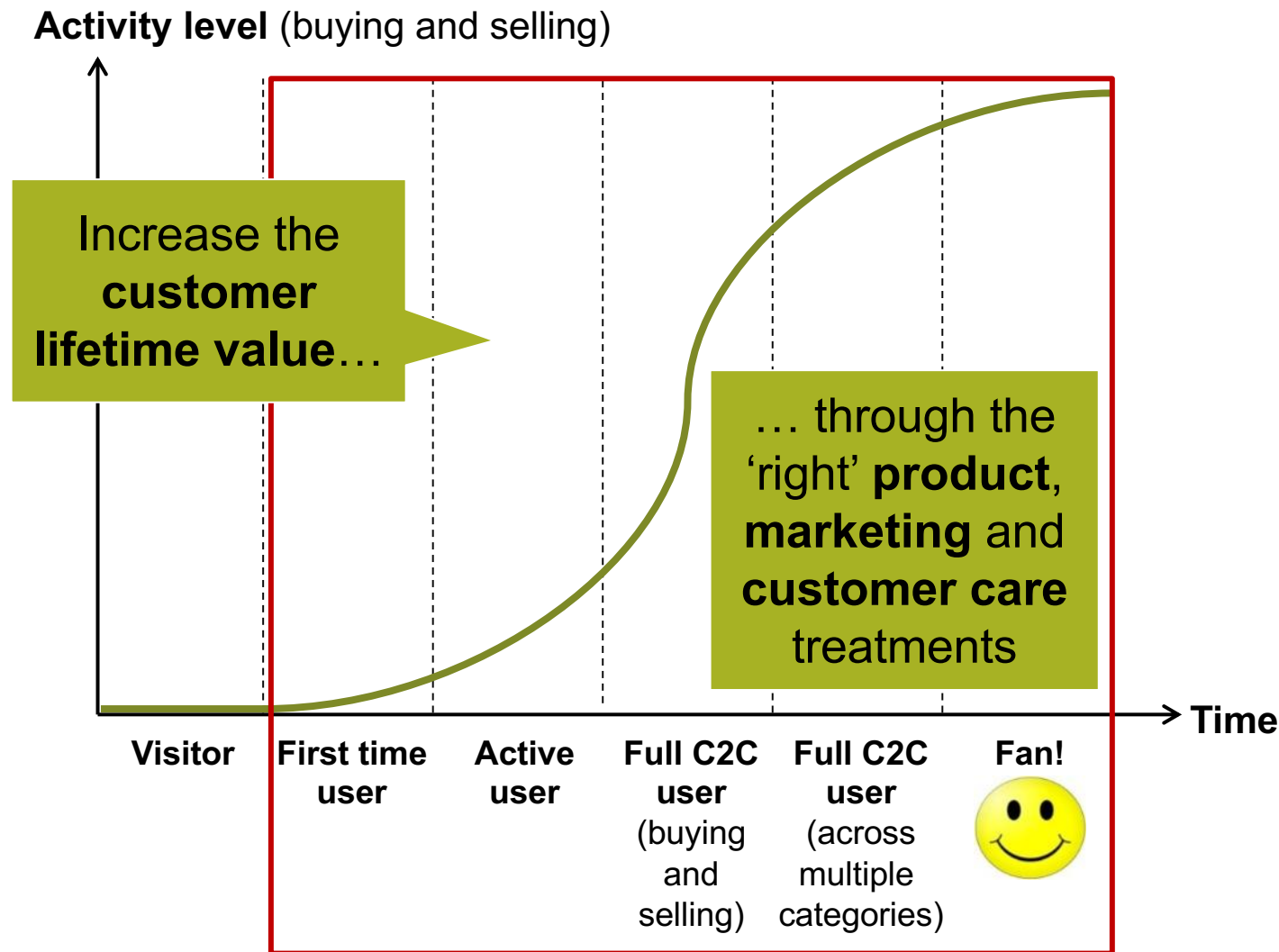
**Net** = Approved, not rejected, not spam  
**New** = First version of listing (not renewed)

# What is Customer Lifecycle Management?

---



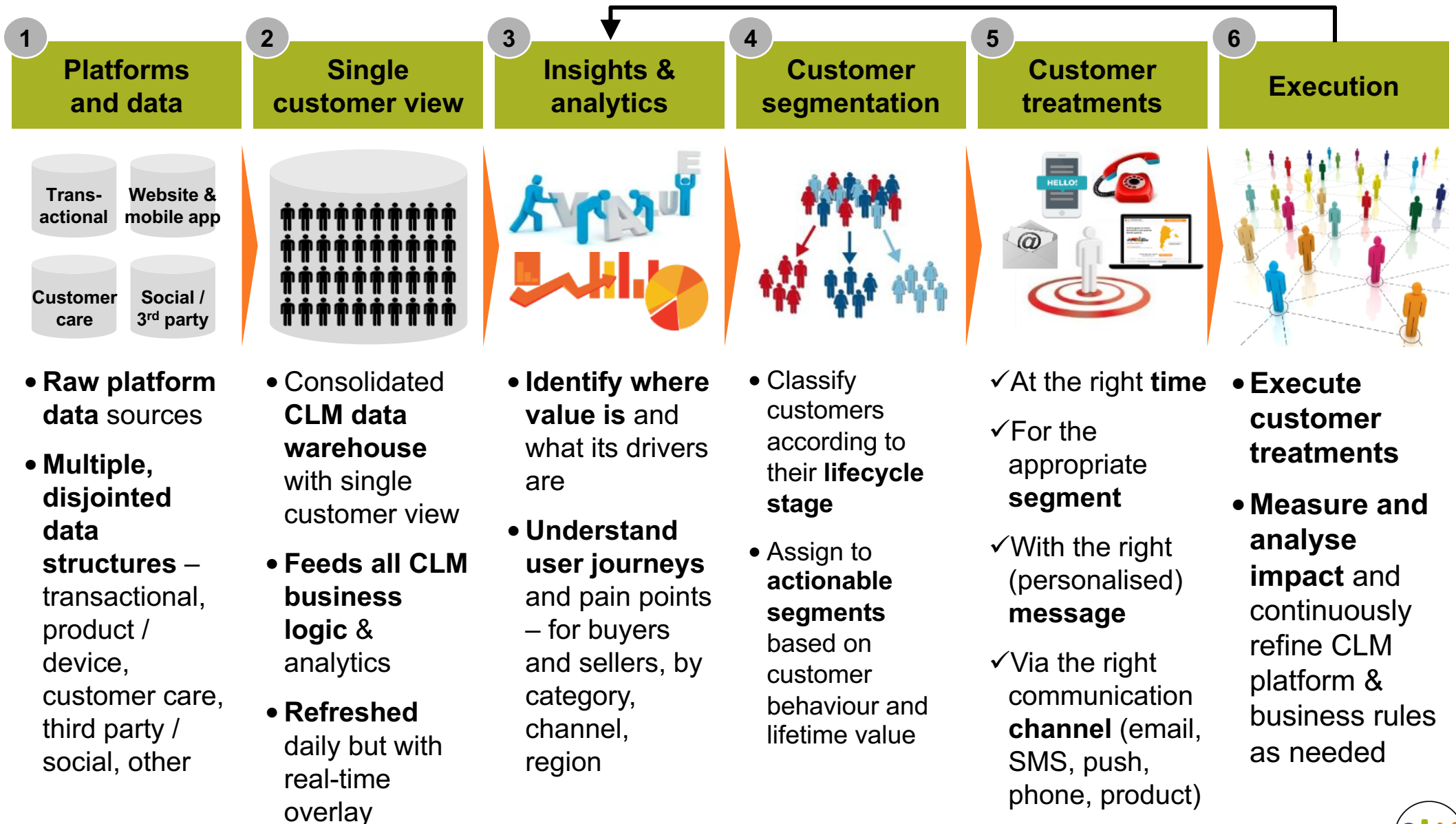
CLM is all about **treating our customers the best way possible** in the lifecycle stage they are in, turning them into 'fans'



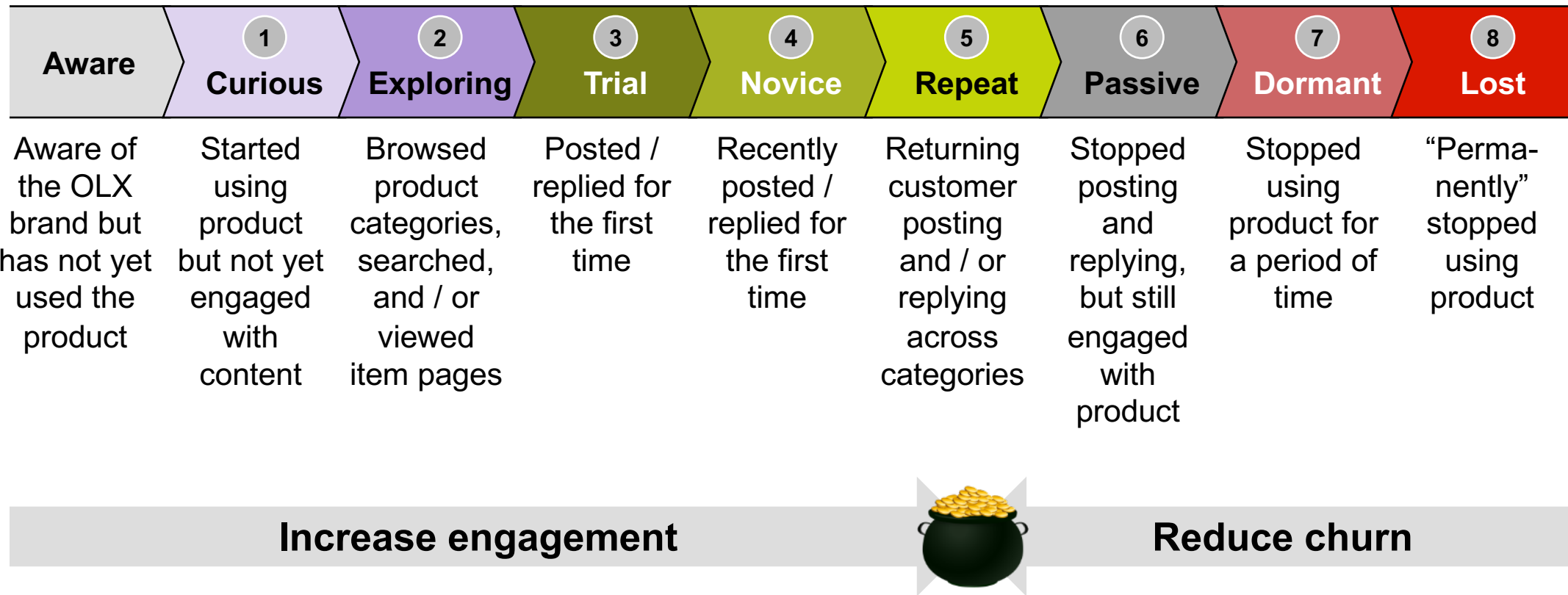
### CLM enables:

- ✓ Dramatic improvement in marketing effectiveness
- ✓ Improved retention
- ✓ Higher engagement
- ✓ Happier customers
- ✓ \$xB value in the endgame

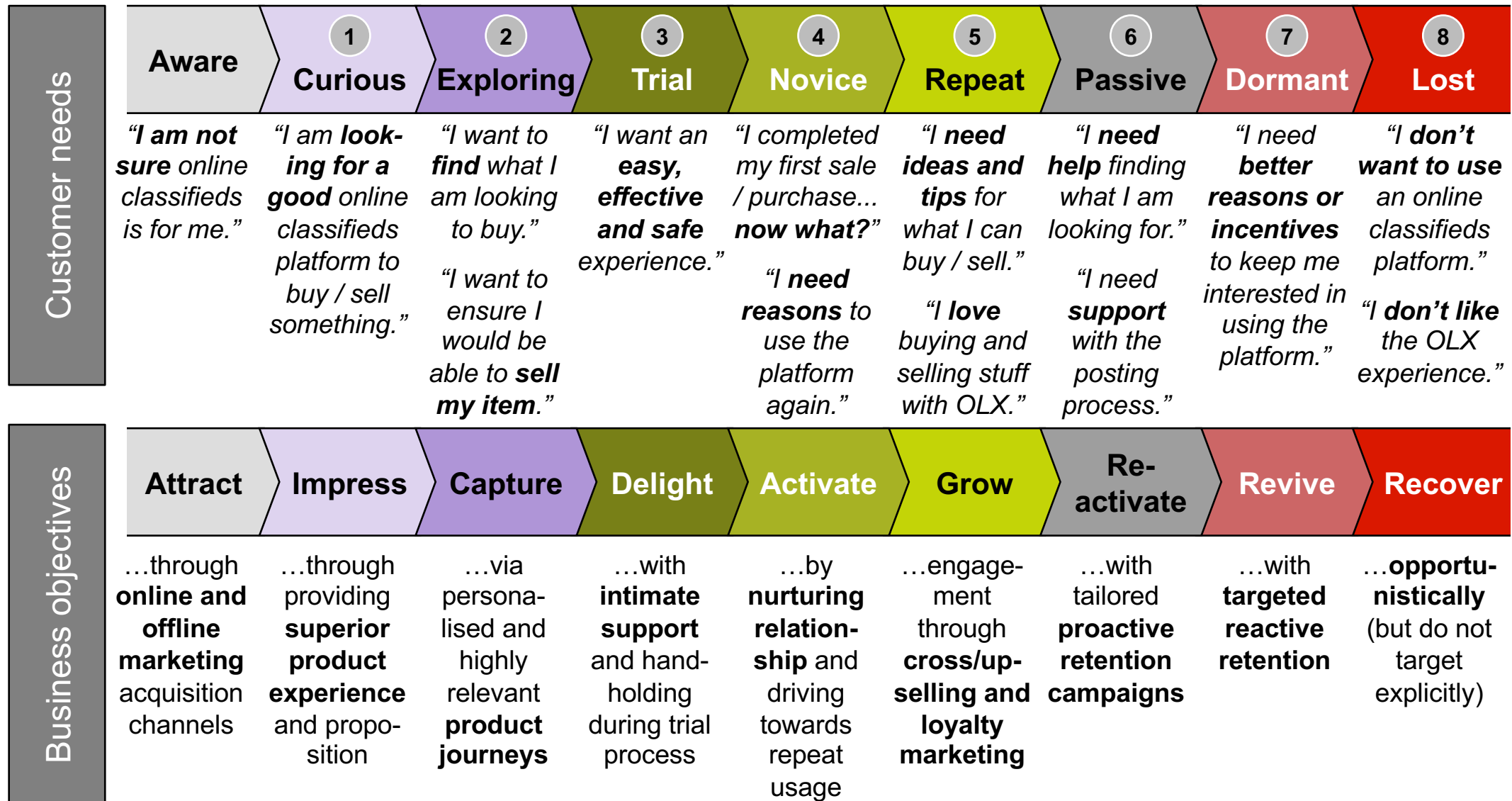
# Sweeping overview of the CLM machine



# Overview of the **customer lifecycle**, from ‘*cradle to grave*’

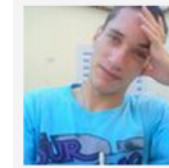


# Customer needs and business objectives evolve over the lifecycle

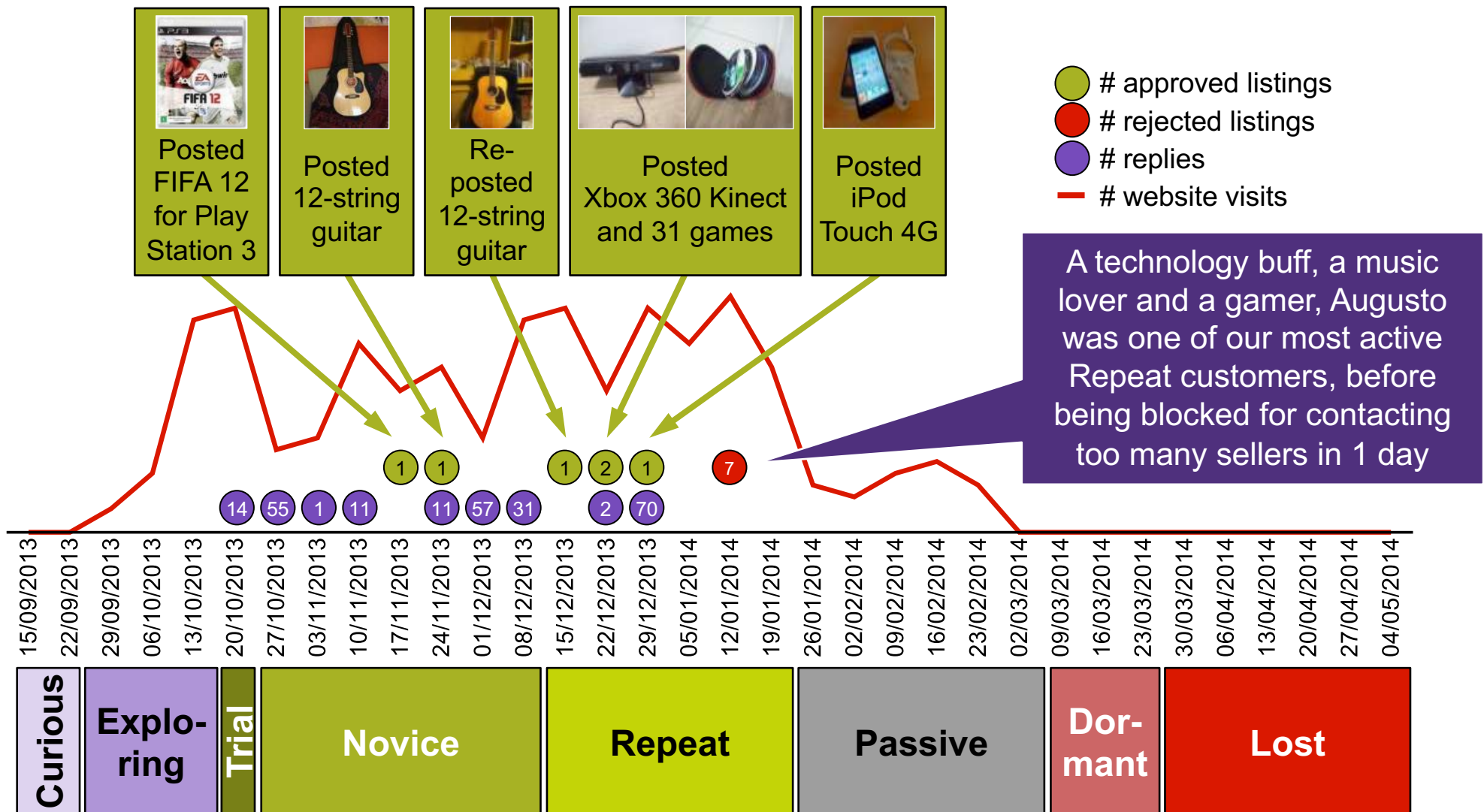


# Illustration of a (real) customer lifecycle

ID: 16342062; email: [djguto00@gmail.com](mailto:djguto00@gmail.com)



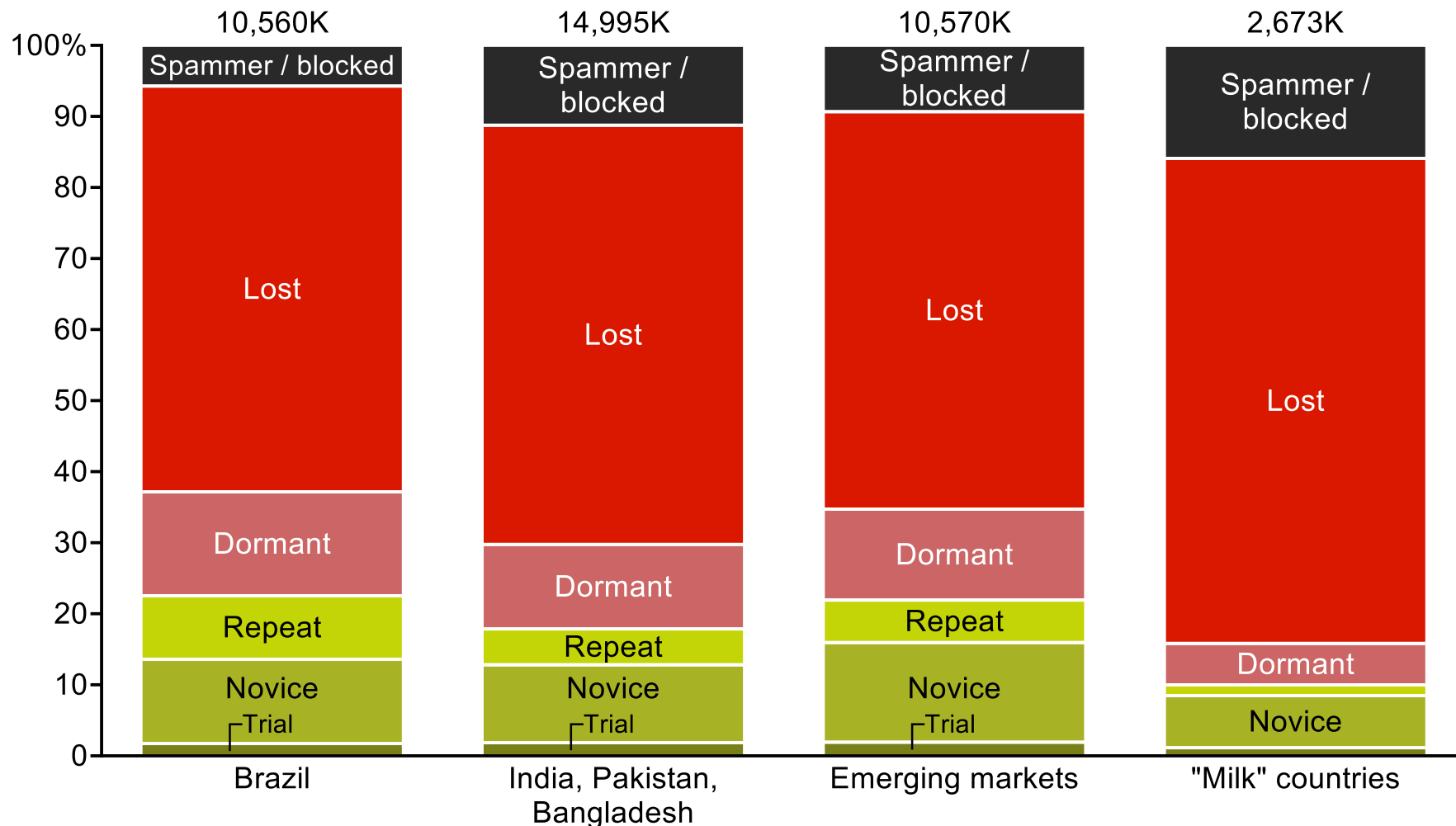
augusto cesar simplicio carlos  
@djguto00  
Sao paulo  
<http://Djguto>



Source: OLX BI data warehouse  
Note: Visits is dummy data

Today's snapshot: **65-70%** of our customers are inactive (**Dormant + Lost**) and only **5-10%** are coming back regularly (**Repeat**)

# customers by region (top 24 countries), as of 29 June 2014

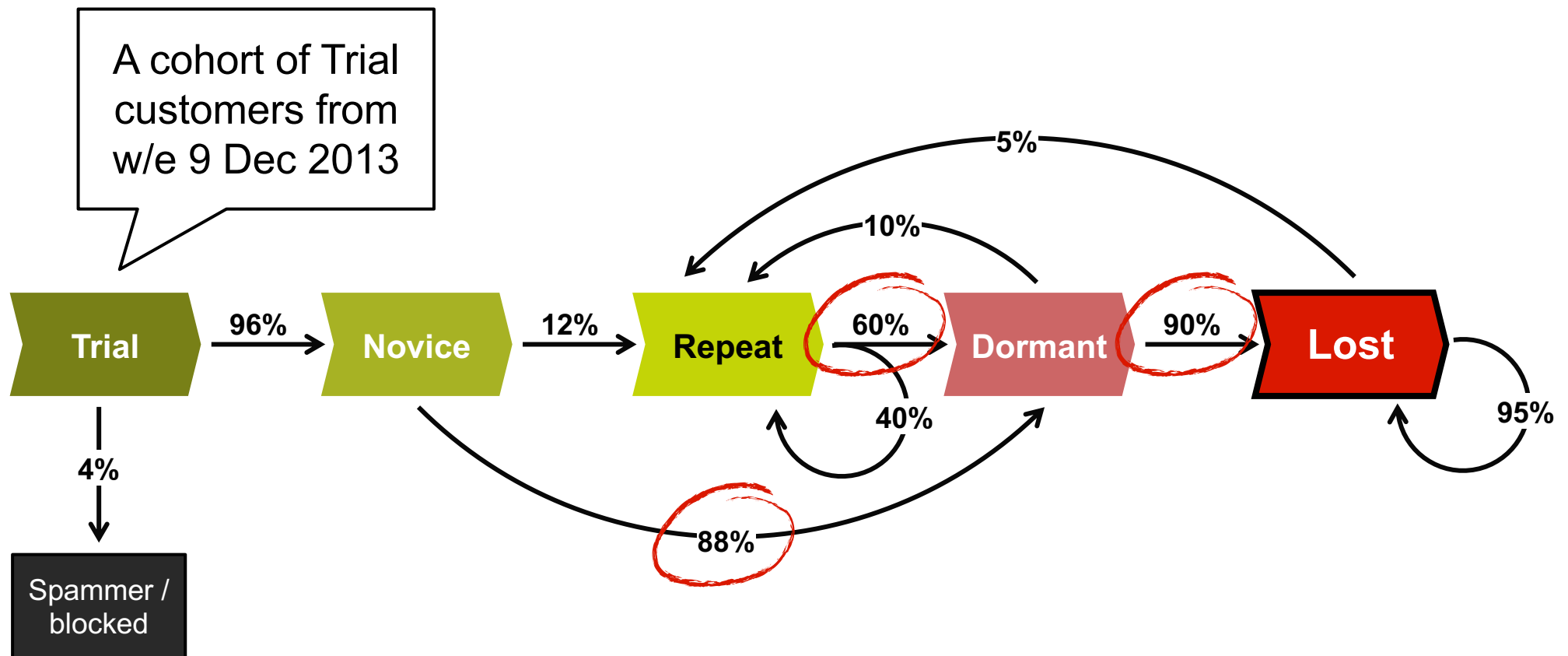




# There is significant customer 'leakage' across all lifecycle stages



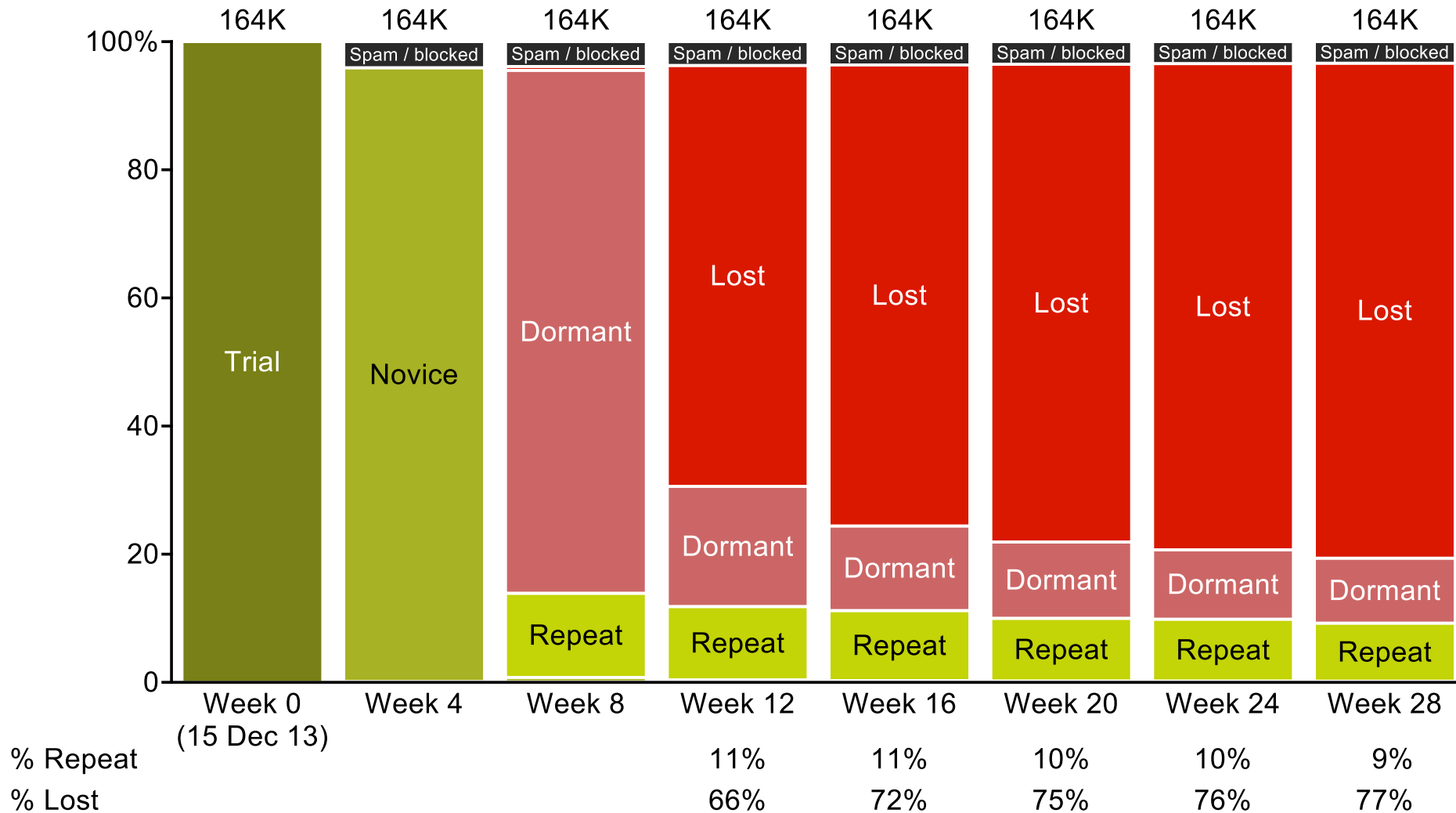
## Behavioural propensities by lifecycle stage



# Impact over time: only ~9% of Trial customers remain active after ~6 months, 77% end up Lost



Evolution of Trial customers, cohort w/e 15 Dec 2013, 000's



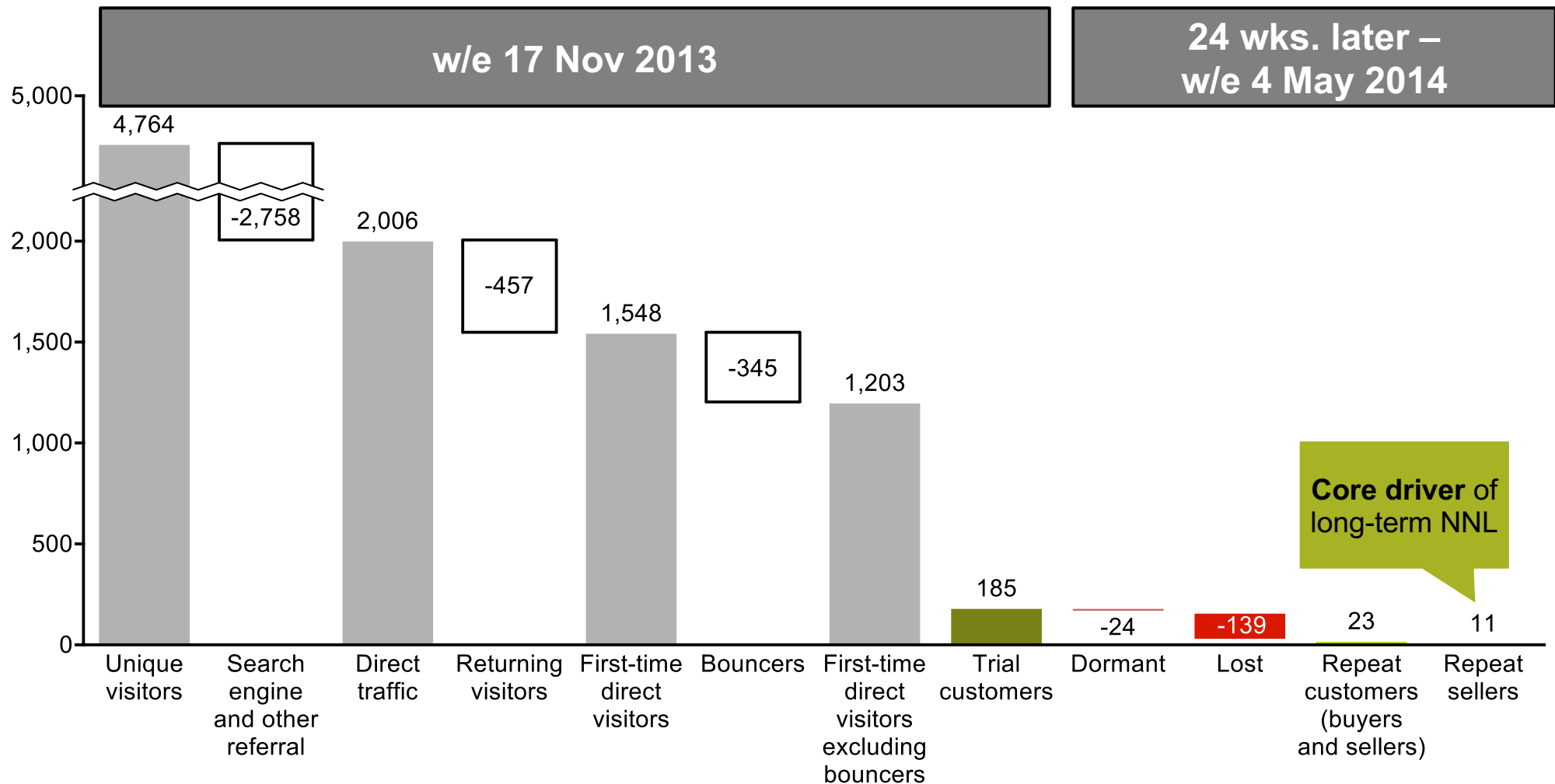
Source: OLX BI data warehouse

Note: Based on extrapolated 10% Brazil sample. \*Active = New + Trial + Repeat

# Ultimately very few customers become Repeat sellers – huge opportunity to improve conversion across the lifecycle



Lifecycle funnel, # customers, 000's (excludes mobile app usage)

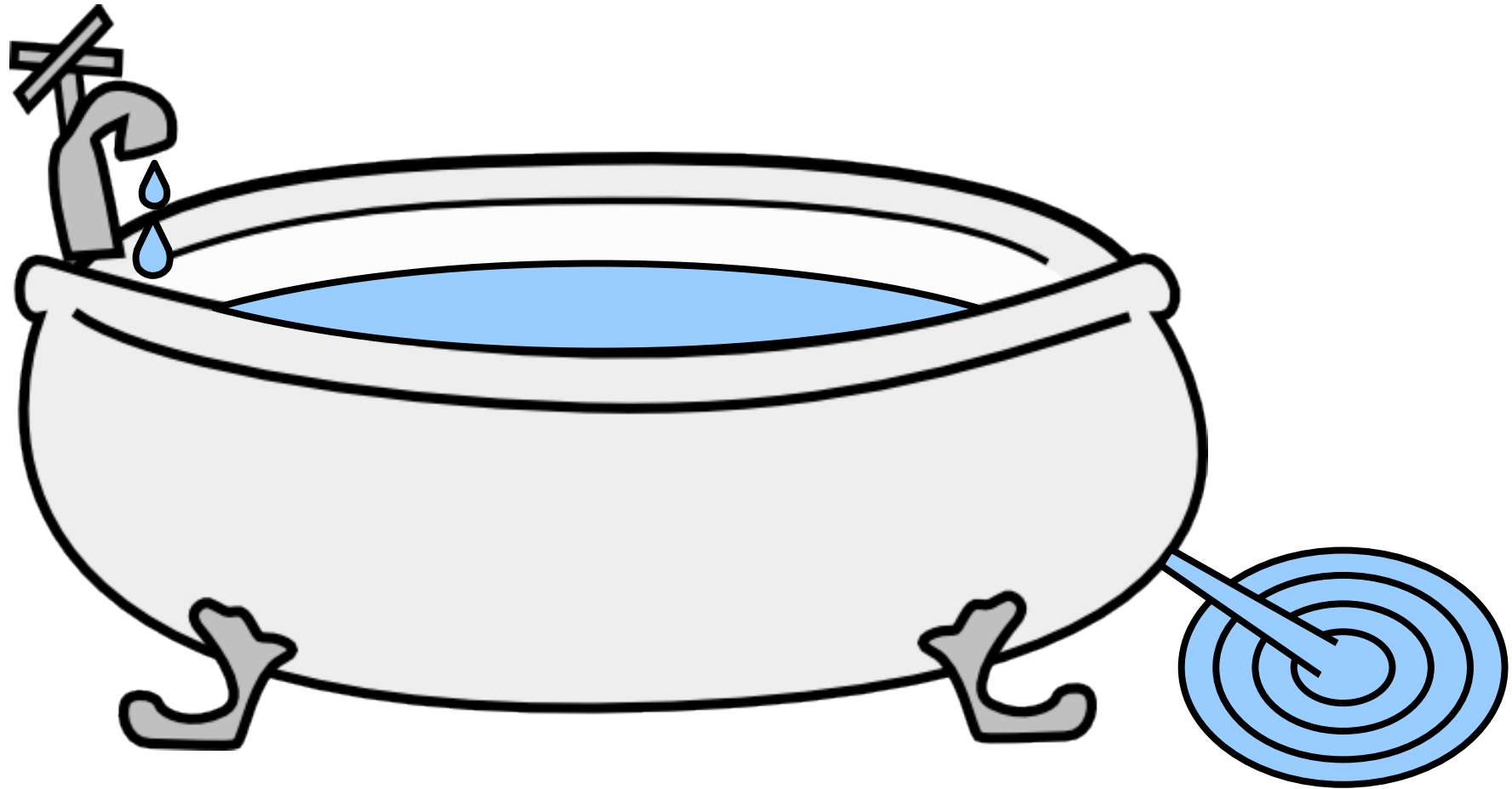


Source: AT Internet; OLX BI data warehouse; internal marketing spend data; analysis


Note: Usage data based on extrapolated 10% Brazil sample and excludes usage from mobile app users. Marketing spend in Nov 2013 in Brazil was \$6.3M.

So... how do we fix the leakage? Let's look at 5 opportunities...

---



**5 strategic CLM objectives** to fix this, leveraging (email, push and phone) marketing, product changes and customer care to get there

Objective	Current state	Target	Impact
1 Capture <u>Exploring</u> customers	<ul style="list-style-type: none"> <li>Only <b>15%</b> of Exploring customers become new Trial users (to be validated)</li> </ul>	<b>17.5%</b>	 <p><b><u>Double</u> daily NNL rate</b></p>
2 Activate <u>Novice</u> customers	<ul style="list-style-type: none"> <li>Only <b>12%</b> of Novice customers become active Repeat users</li> </ul>	<b>17%</b>	
3 Retain <u>Repeat</u> customers	<ul style="list-style-type: none"> <li><b>40%</b> of Repeat customers remain active after 6 months</li> </ul>	<b>50%</b>	
4 Grow <u>Repeat</u> customers	<ul style="list-style-type: none"> <li><b>47%</b> of Repeat customers are sellers (53% are buyers only)</li> <li><b>3 NNL</b> on average per active Repeat seller per month</li> </ul>	<b>50%</b> <b>3.3 NNL</b>	
5 Recover <u>Dormant / Lost</u> customers	<ul style="list-style-type: none"> <li><b>10%</b> of Dormant customers (5% of Lost) re-activate and become Repeat customers again over 6 months</li> </ul>	<b>20%</b>	

# Fun exercises for potential new Data Engineers



<b>Objectives</b>	<ul style="list-style-type: none"><li>• Load data to Redshift</li><li>• Implement CLM segmentation on sample customer data</li><li>• Calculate listing liquidity</li></ul>
<b>Deliverable</b>	<ul style="list-style-type: none"><li>• Table in Redshift cluster</li><li>• SQL code</li></ul>
<b>Format</b>	<ul style="list-style-type: none"><li>• Redshift SQL</li></ul>

## Some tips

- Read the Redshift documentation 😊

<http://aws.amazon.com/documentation/redshift/>

- Recommended SQL clients:

- **Datagrip**

<https://www.jetbrains.com/datagrip/>

- The given data is only a small sample of the real dataset (which is 500x bigger), so think carefully about performance of your query as it is supposed to still run fast on a bigger dataset.
- Depending on the size of your test cluster, your queries should usually finish in less than a couple of minutes

## Description of data structure of sample table **actions**

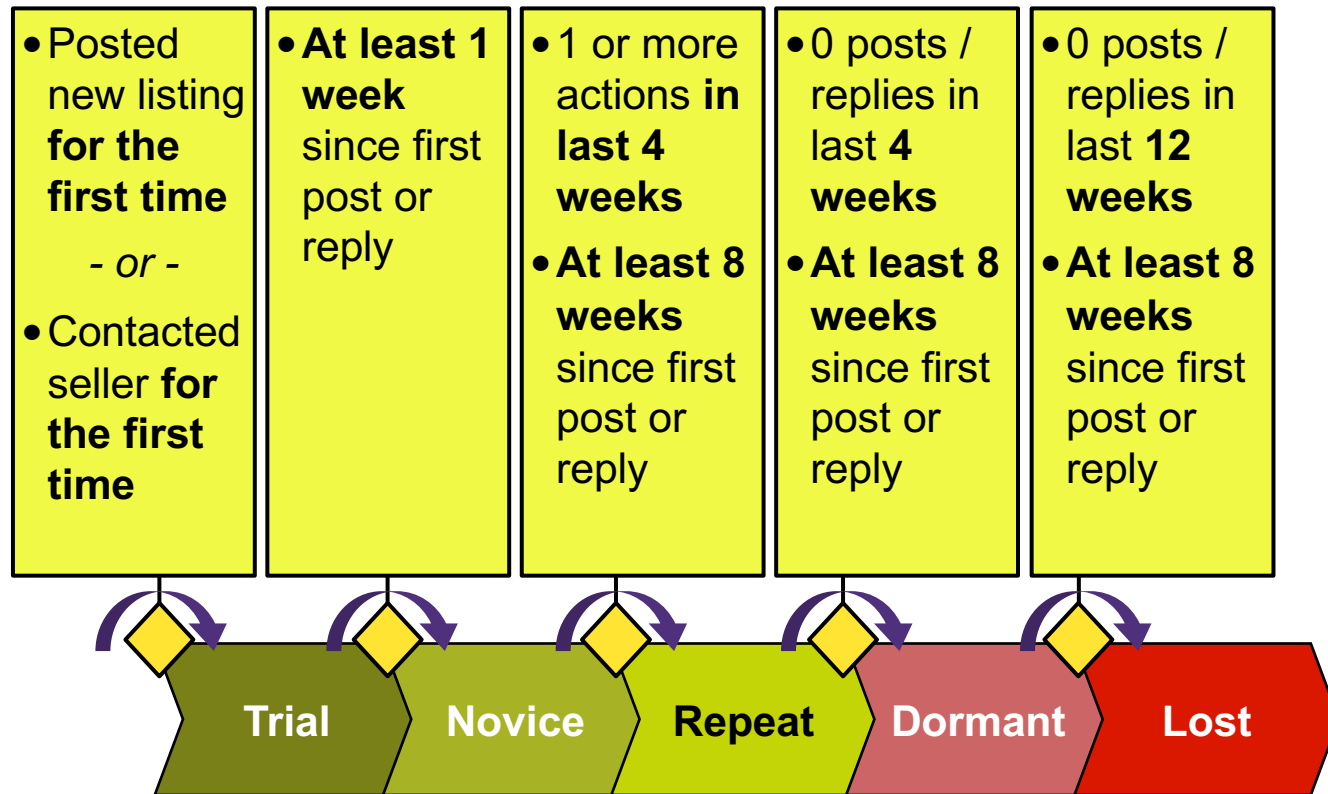
Name	Type	Values
<b>user_id</b>	<b>VARCHAR(48)</b>	Unique id of each user that executed the action
<b>action_type</b>	<b>CHAR(1)</b>	One of two possible actions: 'P' – post (create new ad); 'R' – reply (contact seller)
<b>action_ts</b>	<b>TIMESTAMP</b>	Date-time timestamp of the action
<b>item_id</b>	<b>VARCHAR(48)</b>	The unique id of the item that was posted or replied to
<b>device</b>	<b>VARCHAR(2)</b>	Channel / device used for action: 'DW' – desktop web; 'MW' – mobile web; 'MA' – mobile app; 'U' – Unknown
<b>b2c</b>	<b>BOOLEAN</b>	Whether the item belongs to a B2C product category (Jobs, Real Estate, Vehicles) 1 – item is in B2C category NULL – item is not in B2C category



## Exercise 1: Load sample data from S3 to a new Redshift cluster

- Read Redshift documentation
- Set up free Redshift cluster
- Load sample data (stored in publicly accessible AWS S3 directory)
  - Bucket=s3://tradus-bi-recruiting/actions/
  - Size=5 GB compressed (gzip)
  - aws\_access\_key\_id=AKIAIQGLDQCPBKH6VQCA
  - aws\_secret\_access\_key=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
  - region: eu-west-1
  - You can list the files using AWS CLI using:  
**aws s3 ls s3://tradus-bi-recruiting/actions/**
  - Files were unloaded from Redshift using:  
**UNLOAD ... PARALLEL ON GZIP**
- Tip: Think about distribution, sort keys and encoding!

## Exercise 2: Segment customer base for below 5 lifecycle stages



Note: **action** = post -or- reply

### Outputs:

- A table with the designated segmentation for each user as of 1 July 2018 + SQL code  
***TIP:** Use daily precision (not weekly precision)*
- A query that shows the distribution (relative size) of each lifecycle stage based on this table

## Exercise 3: Calculate liquidity

### Definitions:

- **Replies within X days** → Number of replies received within X days of an item's posting date
  - Example: **'Replies within 7 days'** is number of replies received (calculated for each item) between the item's posting date and (posting date + 7 days)
- **Liquid items X replies within Y days** → Number of items posted on a given date that receive X or more replies within Y days
  - Example: **'Liquid items 3 replies within 7 days'** for a given date is the number of items with 3 or more replies received within 7 days since posting date
- **% liquidity X replies within Y days** → % of all items posted on a given date that receive X or more replies within Y days
  - Example: **'% liquidity 3 replies within 7 days'** for a given date is the number of items with 3 or more replies received within 7 days since posting date ÷ All ads posted on that date

### Exercise:

- Write a query that creates a fact table **fact\_item\_liquidity** with the following information:
  - **Dimensions:**
    - **date** (all dates available in dataset)
    - **item\_id** (all items available in dataset)
  - **Measures:**
    - # replies received within 1 day
    - # replies received within 7 days
- Write a query that creates a fact table **fact\_liquidity** with the following information (Tip: Use previous table as input):
  - **Dimensions:**
    - **date** (all dates available in dataset)
  - **Measures:**
    - # items posted on date
    - Liquid items 1 reply within 1 day
    - Liquid items 3 replies within 7 days
    - Liquid items 5 replies within 7 days
    - % liquidity 1 reply within 1 day
    - % liquidity 3 replies within 7 days
    - % liquidity 5 replies within 7 days

# Final notes

- Sample table is highly simplified for this case. As simple as it might look, it will give us a good insight in your code style and coding skills.
- Comment your code and make it readable . Use 4 spaces or tabs that represent 4 spaces. See sample formatting below.
- Send your output in a zip file including access details to your test cluster and console
- If in doubt make sure to clarify exercises with your interviewer to avoid risk of incorrect output

```
-- Example code formatting
SELECT  x,
        y,
        MAX(z)
  FROM  a
  JOIN  t ON a.x = t.x
  LEFT  JOIN t2 ON a.x = t2.x
 GROUP BY x,y
HAVING  MAX(z) > 1
 ORDER BY x,y
 LIMIT 100;
```