

---

*Customer Retention Strategy*

---

## ***1. Business Challenge/Requirement***

***Customer retention and acquisition strategies are on top of every organization's agenda. To offer better customer service and boost loyalty, a company has to invest in a state-of-the-art CRM tool. In pursuit of these goals, every organization implements CRM as a strategy that integrates the concepts of data mining and data warehousing. The data collected through the CRM helps the leadership team make actionable decisions in real time. It helps them build and retain long-term and profitable relationships with customers.***

***FutureCart Inc. is a hypothetical leading retail company with an omnipresence in India with more than 5000 retail stores and hypermarkets across and e-commerce in the country.***

***The company has formed a dedicated team to handle after-sales services. The team is entrusted with the responsibility to address customer complaints and delight them - and eventually increase brand loyalty***

***Below is an abstract of end to end process:***

- The company has multiple contact centers across India to provide support service to their customers***
- Customers can reach out to the care team over different communication channels depending on their preference and convenience: Calls, Chat, or Email.***
- CCR (Customer Care Representative) registers the complaint by collecting all the necessary details - - which is called a case***
- A case can have a status -- open or closed***
- Each case can belong to a category and sub-category. This category and sub-category will determine***

**case priority. Depending on the priority key, CCR has an SLA (in hours) to close the case within the SLA hours**

- **Once a case is closed, the customer is sent a survey link to rate the overall experience of interacting with the contact center representative**
- **The customer can take a survey or leave it unattended. The customer can rate the experience on a scale of 1-10 on various questions**
- **Survey response is captured for that particular case. The data collected through complete CRM process is used by the company for analysis. The analytics team working on this data captures the below KPIs to further enhance and optimize the CRM process. KPIs (Both on real-time data and batch-processed data)**
- **Total numbers of cases**
- **Total open cases in the last 1 hour**
- **Total closed cases in the last 1 hour**
- **Total priority cases**
- **Total positive/negative responses in the last 1 hour**
- **Total number of surveys in the last 1 hour**
- **Total open cases in a day/week/month**
- **Total closed cases in a day/week/month**
- **Total positive/negative responses in a day/week/month**
- **Total number of surveys in a day/week/month Real-time KPIs**
- **Total numbers of cases that are open and closed out of the number of cases received**

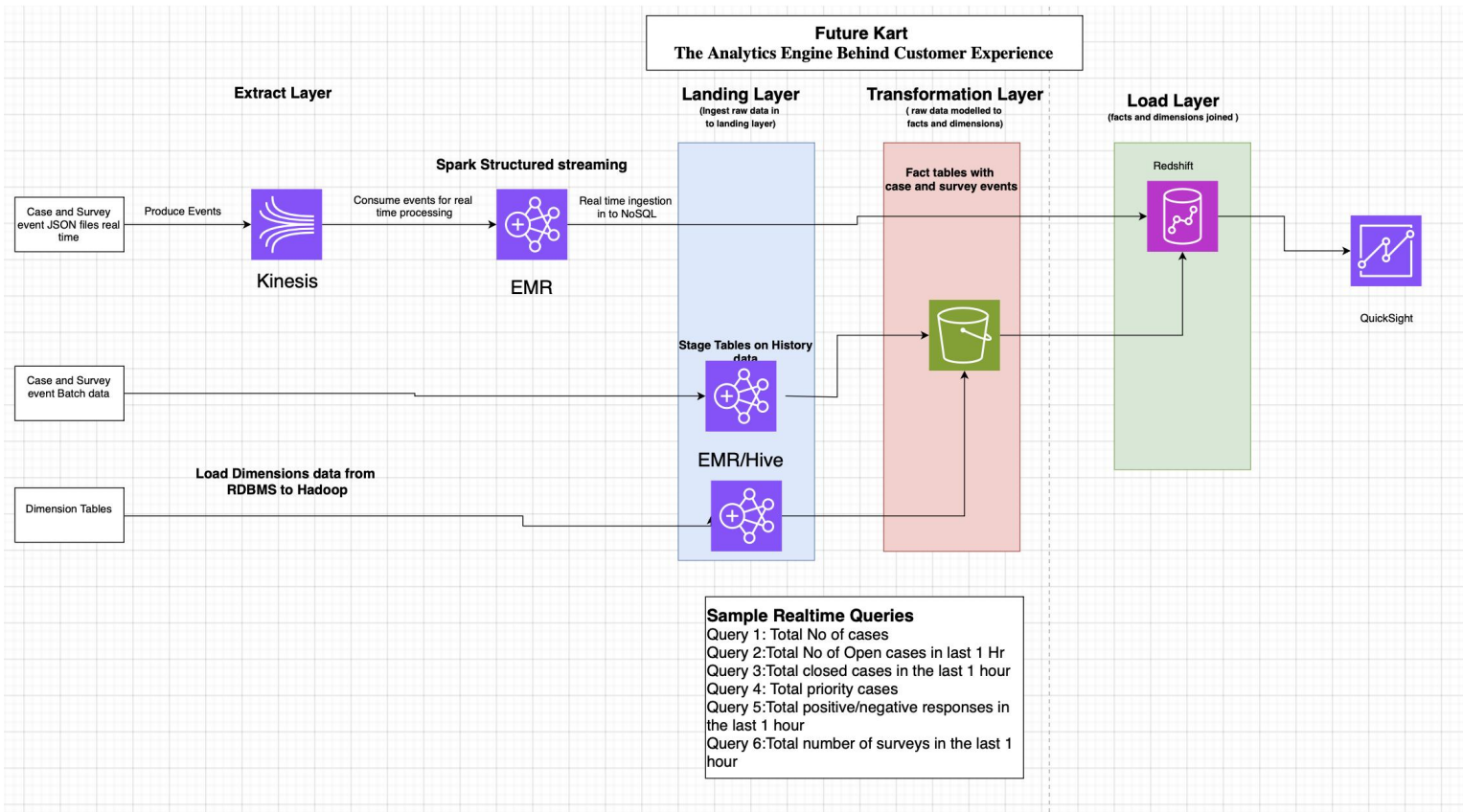
- **Total number of cases received based on priority and severity**

## 2. The Goal of the Project

**Below are some of the high-level technical and non-technical goals for this project:**

- **Get an overall understanding of the CRM domain**
- **Learn the fundamentals & standards of ETL and data warehousing**
- **Real-time and batch ingestion of data from multiple sources to Big Data storage like Hive/ DynamoDB /HDFS using Kinesis and Spark**
- **Data cleansing/wrangling/transformation using Hive and Spark**
- **Lambda architecture where data can be processed in both batch and real-time**
- **Reporting KPIs(Key Performance Indicators)**

## 3. Data Flow Architecture/Process Flow



#### 4. Dataset Explanation and Schema

*We have three types of data sources:*

- ❖ *Data for which static/dimension tables to be created in MySQL*
- ❖ *Historical data of 10 days for cases and survey events to created in Hive(JSON)*
- ❖ *Real-time data for the current date for cases and survey events in JSON format*

##### 4.1 Data for which static/dimension tables to be created in MySQL

*We have the below datasets -*

*futurecart\_calendar\_details.txt - Calendar details for the company*

<i>column Name</i>	<i>Data type</i>	<i>Column description</i>	<i>sample value</i>
<i>calendar_date</i>	<i>date,</i>	<i>Calendar date in yyyy-mm-dd format</i>	<i>2011-02-20</i>
<i>date_desc</i>	<i>varchar(50)</i>	<i>Calendar date in words</i>	<i>Sunday, February 20, 2011</i>
<i>week_day_nbr</i>	<i>smallint</i>	<i>Number of days in a week</i>	<i>2</i>

<i><b>week_number</b></i>	<i><b>smallint</b></i>	<i><b>Week number of the year</b></i>	<i><b>_____</b></i> <i><b>4</b></i>
<i><b>week_name</b></i>	<i><b>varchar(50)</b></i>	<i><b>Week name</b></i>	<i><b>Week 04</b></i>
<i><b>year_week_number</b></i>	<i><b>int</b></i>	<i><b>Week number with year</b></i>	<i><b>_____</b></i> <i><b>201104</b></i>
<i><b>month_number</b></i>	<i><b>smallint</b></i>	<i><b>Month number in the year</b></i>	<i><b>1</b></i>
<i><b>month_name</b></i>	<i><b>varchar(50)</b></i>	<i><b>Month name</b></i>	<i><b>february</b></i>
<i><b>quarter_number</b></i>	<i><b>smallint</b></i>	<i><b>Quarter number in the year</b></i>	<i><b>1</b></i>
<i><b>quarter_name</b></i>	<i><b>varchar(50)</b></i>	<i><b>Quarter name</b></i>	<i><b>Q1</b></i>
<i><b>half_year_number</b></i>	<i><b>smallint</b></i>	<i><b>Half-year number in the year</b></i>	<i><b>1</b></i>
<i><b>half_year_name</b></i>	<i><b>varchar(50)</b></i>	<i><b>Half-year name</b></i>	<i><b>1st Half</b></i>
<i><b>geo_region_cd</b></i>	<i><b>char(2)</b></i>	<i><b>Geographic region code</b></i>	<i><b>US</b></i>

***futurecart\_call\_center\_details.txt – Contact/Call center details for the company***

<b><i>column Name</i></b>	<b><i>Data type</i></b>	<b><i>Column description</i></b>	<b><i>sample value</i></b>
<b><i>call_center_id</i></b>	<b><i>varchar(10)</i></b>	<b><i>Unique identifier for a call center</i></b>	<b><i>C-101</i></b>
<b><i>call_center_vendor</i></b>	<b><i>varchar(50)</i></b>	<b><i>Vendor company name which is handling the call center</i></b>	<b><i>Concentrix</i></b>
<b><i>location</i></b>	<b><i>varchar(50)</i></b>	<b><i>Call center location</i></b>	<b><i>New york</i></b>
<b><i>country</i></b>	<b><i>varchar(50)</i></b>	<b><i>Call center country</i></b>	<b><i>US</i></b>

***futurecart\_case\_category\_details.txt - Category details of a case event***

<b><i>column Name</i></b>	<b><i>Data type</i></b>	<b><i>Column description</i></b>	<b><i>sample value</i></b>
<b><i>category_key</i></b>	<b><i>varchar(10)</i></b>	<b><i>Unique identifier for a case category</i></b>	<b><i>CAT1</i></b>
<b><i>sub_category_key</i></b>	<b><i>varchar(10)</i></b>	<b><i>Unique identifier for a case sub category</i></b>	<b><i>SCAT1</i></b>
<b><i>category_description</i></b>	<b><i>varchar(50)</i></b>	<b><i>Category description</i></b>	<b><i>Subscription</i></b>
<b><i>sub_category_description</i></b>	<b><i>varchar(50)</i></b>	<b><i>Subcategory description</i></b>	<b><i>Renewal</i></b>
<b><i>priority</i></b>	<b><i>varchar(10)</i></b>	<b><i>Priority key</i></b>	<b><i>P1</i></b>

*futurecart\_case\_country\_details.txt - Country details*

<i>column Name</i>	<i>Data type</i>	<i>Column description</i>	<i>sample value</i>
<i>id</i>	<i>int</i>	<i>Unique identifier for a country</i>	<i>4</i>
<i>Name</i>	<i>varchar(75)</i>	<i>Country name</i>	<i>India</i>
<i>Alpha_2</i>	<i>varchar(2)</i>	<i>Country short name 2 chars</i>	<i>IN</i>
<i>Alpha_3</i>	<i>varchar(2)</i>	<i>Country short name 3 chars</i>	<i>IND</i>

*futurecart\_case\_priority\_details.txt - Priority details of a case*

<i>column Name</i>	<i>Data type</i>	<i>Column description</i>	<i>sample value</i>
<i>Priority_key</i>	<i>varchar(5)</i>	<i>Unique identifier for a case priority</i>	<i>P1</i>
<i>priority</i>	<i>varchar (20)</i>	<i>Priority level</i>	<i>Highest</i>
<i>severity</i>	<i>varchar (100)</i>	<i>Severity level</i>	<i>critical</i>
<i>SLA</i>	<i>varchar (100)</i>	<i>SLA in HOURS for the priority and severity combination</i>	<i>1</i>



***futurecart\_employee\_details.txt - Employee details of the company***

<b><i>column Name</i></b>	<b><i>Data type</i></b>	<b><i>Column description</i></b>	<b><i>sample value</i></b>
<b><i>emp_key</i></b>	<b><i>Int</i></b>	<b><i>Unique ID of an employee</i></b>	<b><i>10001</i></b>
<b><i>first_name</i></b>	<b><i>varchar</i></b>	<b><i>First name</i></b>	<b><i>Georgi</i></b>
<b><i>last_name</i></b>	<b><i>varchar</i></b>	<b><i>Last name</i></b>	<b><i>Facello</i></b>
<b><i>email</i></b>	<b><i>varchar</i></b>	<b><i>email</i></b>	<b><i>Georgi.Facello01@testmail.com</i></b>
<b><i>gender</i></b>	<b><i>varchar</i></b>	<b><i>gender</i></b>	<b><i>M</i></b>
<b><i>ldap</i></b>	<b><i>varchar</i></b>	<b><i>User id</i></b>	<b><i>5941CF7D</i></b>
<b><i>hire_date</i></b>	<b><i>Date</i></b>	<b><i>Hire date</i></b>	<b><i>2014-04-06</i></b>
<b><i>manager</i></b>	<b><i>varchar</i></b>	<b><i>Manager key</i></b>	<b><i>455246</i></b>

***futurecart\_product\_details.txt - Product details of the company***

<b><i>column Name</i></b>	<b><i>Data type</i></b>	<b><i>Column description</i></b>	<b><i>sample value</i></b>
<b><i>product_id</i></b>	<b><i>varchar</i></b>	<b><i>Unique id for a product</i></b>	<b><i>26355</i></b>
<b><i>department</i></b>	<b><i>varchar</i></b>	<b><i>Department description</i></b>	<b><i>GROCERY</i></b>
<b><i>brand</i></b>	<b><i>varchar</i></b>	<b><i>Brand description</i></b>	<b><i>Private</i></b>
<b><i>commodity_desc</i></b>	<b><i>varchar</i></b>	<b><i>Commodity description</i></b>	<b><i>COOKIES/CONES</i></b>
<b><i>sub_commodity_desc</i></b>	<b><i>varchar</i></b>	<b><i>Subcommodity description</i></b>	<b><i>SPECIALTY COOKIES</i></b>

***futurecart\_survey\_question\_details.txt - Question details for the survey***

<b><i>column Name</i></b>	<b><i>Data type</i></b>	<b><i>Column description</i></b>	<b><i>sample value</i></b>
<b><i>question_id</i></b>	<b><i>varchar</i></b>	<b><i>Unique id for a survey question</i></b>	<b><i>Q1</i></b>
<b><i>question_desc</i></b>	<b><i>varchar</i></b>	<b><i>Question text</i></b>	<b><i>How would you rate your overall e customer support process?</i></b>
<b><i>response_type</i></b>	<b><i>varchar</i></b>	<b><i>Response type (scale or options)</i></b>	<b><i>Scale</i></b>
<b><i>range</i></b>	<b><i>varchar</i></b>	<b><i>Scale range if the response type is scale else NA</i></b>	<b><i>1-10</i></b>

<i><b>negative_response_range</b></i>	<i><b>varchar</b></i>	<i><b>Scale range to qualify a survey response as negative</b></i>	<i><b>1-4</b></i>
<i><b>neutral_response_range</b></i>	<i><b>varchar</b></i>	<i><b>Scale range to qualify a survey response as neutral</b></i>	<i><b>5-7</b></i>
<i><b>positive_response_range</b></i>	<i><b>varchar</b></i>	<i><b>Scale range to qualify a survey response as positive</b></i>	<i><b>8-10</b></i>

*4.2 Historical data of 10 days for cases and survey events to be created in hive*

*futurecart\_case\_details.txt – Case details of the company*

<i>column Name</i>	<i>Data type</i>	<i>Column description</i>	<i>sample value</i>
<i>case_no</i>	<i>varchar</i>	<i>Unique ID of a case</i>	<i>2024</i>
<i>create_timestamp</i>	<i>varchar</i>	<i>Case create timestamp</i>	<i>2020-04-20 01:01:29</i>
<i>last_modified_timestamp</i>	<i>varchar</i>	<i>Case last modified timestamp</i>	<i>2020-04-20 01:01:29</i>
<i>created_employee_key</i>	<i>varchar</i>	<i>Employee key who created the case</i>	<i>274649</i> _____
<i>call_center_id</i>	<i>varchar</i>	<i>Call center id where case is logged and handled</i>	<i>C-104</i>
<i>status</i>	<i>varchar</i>	<i>Current status of the case</i>	<i>Open</i>
<i>category</i>	<i>varchar</i>	<i>Category key of the case</i>	_____ <i>CAT1</i>
<i>sub_category</i>	<i>varchar</i>	<i>Subcategory key of the case</i>	<i>S CAT1</i>
<i>communication_mode</i>	<i>varchar</i>	<i>Mode of communication</i>	<i>Email</i>
<i>country_cd</i>	<i>varchar</i>	<i>Country code</i>	<i>PY</i>
<i>product_code</i>	<i>varchar</i>	<i>Product code</i>	<i>997719</i>

### *futurecart\_case\_survey\_details.txt – survey details of the cases closed*

<i>column Name</i>	<i>Data type</i>	<i>Column description</i>	<i>sample value</i>
<i>survey_id</i>	<i>varchar</i>	<i>Unique ID of a survey</i>	<i>S-1000</i>
<i>Case_no</i>	<i>varchar</i>	<i>Case number for which survey has been filled</i>	<i>130114</i>
<i>survey_timestamp</i>	<i>varchar</i>	<i>Survey taken timestamp</i>	<i>2020-04-20 01:01:2</i>
<i>Q1</i>	<i>varchar</i>	<i>Q1 response</i>	<i>2</i>
<i>Q2</i>	<i>varchar</i>	<i>Q2 response</i>	<i>7</i>
<i>Q3</i>	<i>varchar</i>	<i>Q3 response</i>	<i>3</i>
<i>Q4</i>	<i>varchar</i>	<i>Q4 response</i>	<i>N</i>
<i>Q5</i>	<i>varchar</i>	<i>Q5 response</i>	<i>7</i>

### *4.3 Real-time data for the current date for cases and survey events in JSON files*

*Copy the file - stream\_to\_kinesis.py to your VM and generate real-time data. This real-time simulator script which will generate JSON data and writes it to the kinesis stream of your choice. Change the script to point it to your stream.*

- *futurecart\_case\_event*
- *futurecart\_survey\_event*

**JSON formats:**

**Sample data :**

**case\_data:**

```
[{  
  "status": "Open", "category": "CAT3", "sub_category": "SCAT14",  
  "last_modified_timestamp": "2020- 06-17 18:42:19", "case_no": "600999",  
  "create_timestamp": "2020-06-17 18:42:19",  
  "created_employee_key": "240604", "call_center_id": "C-116", "product_code": "9829787",  
  "country_cd": "PR", "communication_mode": "Chat"  
}, {  
  "status": "Open", "category": "CAT3", "sub_category": "SCAT14",  
  "last_modified_timestamp": "2020- 06-17 18:42:19", "case_no": "601000",  
  "create_timestamp": "2020-06-17 18:42:19",  
  "created_employee_key": "215285", "call_center_id": "C-114", "product_code":  
  "12457101", "country_cd": "EE", "communication_mode": "Call"  
}]
```

**Survey Data:**

```
[{  
  "Q1": 9, "Q3": 1, "Q2": 8, "Q5": 3, "Q4": "N", "case_no": "600991", "survey_timestamp":  
  "2020-06-17  
  19:42:04", "survey_id": "S-500014"  
}, {  
  "Q1": 8, "Q3": 9, "Q2": 1, "Q5": 1, "Q4": "N", "case_no": "600992", "survey_timestamp":  
  "2020-06-17  
  19:42:04", "survey_id": "S-500015"  
}]
```

## **5. Problem Statements/Tasks**

**The high-level task is to create a Data Mart on CRM data with a lambda architecture where we will ingest and process data in both batch and real time. We also want to enable reporting of KPIs in both batch and real time.**

**Technical tasks in details :**

**Refer data flow and architecture for additional reference:**

- 1. Companies generally store transactional data in RDBMS because they provide faster read and write operations and support ACID properties. Hence, create MySQL tables for the dimension datasets shared. Create a database in an EC2 instance with MySQL 5.6 server and create all the required dimension tables.**
- 2. Perform batch ingestion from MySQL to Hive tables for static dimensions.**
- 3. Historical data for the last 10 days is generated using the script (generate\_historical\_data.py) which generates json data for cases and surveys. Load this data in HDFS and create historical tables for cases and surveys.**
- 4. Generate new cases and survey events in json format from the python script and send them to kinesis stream**

***5. Create an application that will consume and process real-time data.***

***6. Load the realtime data coming from kinesis in to redshift as fact tables. If the incoming record from stream is a “Case” data load it into case table else load it in to survey table***

***7. Load the historical and dimension tables from hive in to s3 using spark with dataframes. Load these s3 data in to redshift tables***

***8. Create the above queries on the redshift tables.***

***9. Create dashboard for some of the queries in quicksight with redshift as the source***