

Real time Spend intent prediction...

For each transaction.

- Collate historic features
- **Compute Dynamic Features**
- Combine Static and Dynamic Features related to CM, SE and geo-location
- For each category execute the GBDT model In-parallel and compute the predictions
- If the predicted probability is greater than the threshold for that category, then, suggest the category as a possible category that the CM might

visit next. Peak traffic of 30 TPS equating to 30K lookups per

second and 5K complex calculations per second; 900 GBDT executions in-parallel at peak

Through training rank order the features that predict better and assign each of these features some weights such that better predicting feats.

Have more weights. Eg. Phase of day might have lesser weights as compared to the feature indicating the Periodicity of the CM visiting a particular SE.

Calculate the prediction thresholds for each Category.

Eg. If prediction probability for Restaurant Category is greater than 0.1 then consider it true

~1000 feature collation and lookup; ~100 dynamic feature computation and running 30+ models

Predict if a CM will transact at a particular category in the next two hours Define

Objective

Predict

Learn

Model building and Train

Data points for Modeling. (i) Transaction data (ii) SE Characteristics

Static Features

periodically? How much time does

Does the CM visit SE

Entertainment, etc.)?

the CM spend in a category (Lifestyle,

CM's home location CM's FICO score

CM's share of wallet

Features.

Dynamic Features

Time of the day, day of the week, phase of day

Trends in the neighboring merchants

Distance of the CM from the home location

Number of times the CM has transacted in a session

Feature Engineering search for all CMs and SEs

Build several decision on different subsets of the

~60 GB of data size; ~1000 features,~100 dynamic

3MM CMs, 150K SEs. Facilitating sub-second

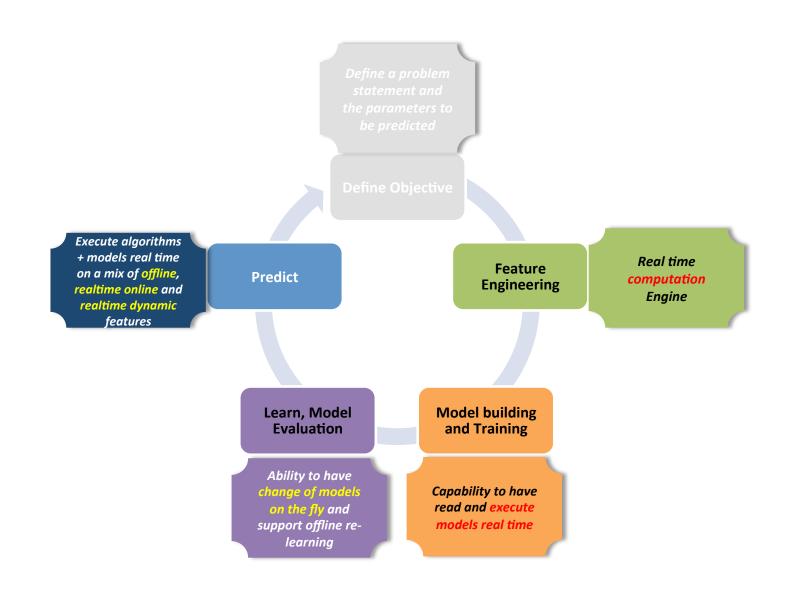
Eg. if (phase of day="morning") and (CM has not transacted at "restaurant" in this session) and (CM periodically visits "restaurants" before Noon); then – there is a 95% chance that the CM

Will transact in the "RESTAURANT" category.

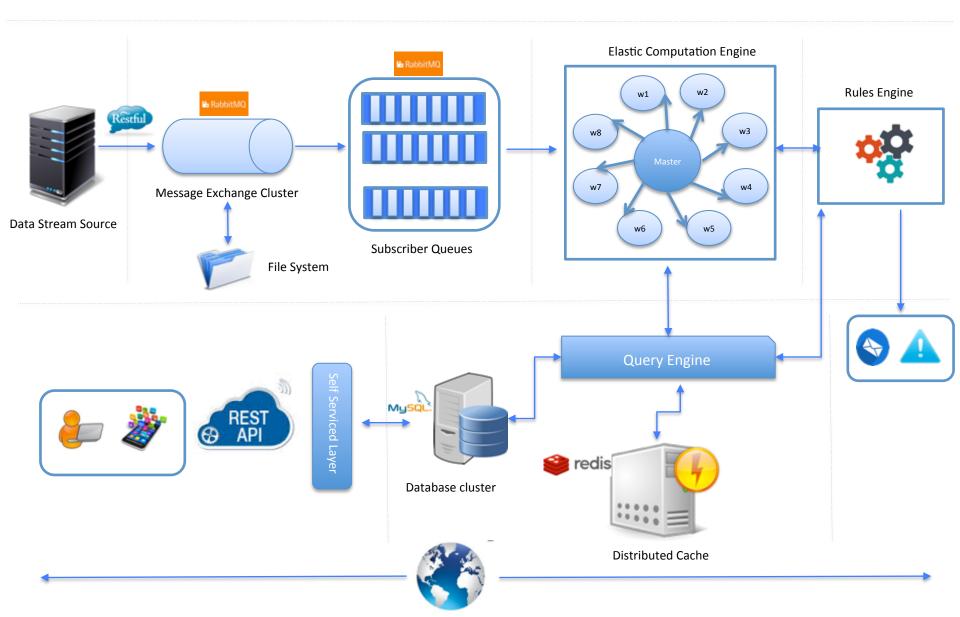
Create a Gradient Boosting Decision Trees(GBDT) model of the training set for each category of merchants

30 GBDT models with 200 Decision Trees each leading 30 parallel executions of GBDT ~900ms execution time

Real time Prediction at 30k feet...



Big Data – Real Time Analytics Platform (CAPSTONE) : End To End Logical Architecture



Big Data – Real Time Analytics Platform (CAPSTONE): Physical Architecture

Database (MySQL)

Database Cache (Memcache) Distributed Cache (Redis)

Message Queue (ActiveMQ/ Solace) Logging/ Monitoring (Splunk) Compute Engine (Custom)

Virtual IP Required

Data Nodes (2 Nodes)

Database (MySQL)

Database
Cache
(Memcache) Database (MySQL)

Database Cache (Memcache)

Virtual IP Required Dist Cache Nodes (2 Nodes)

Distributed Cache (Redis)

Distributed
Cache
(Redis)

Virtual IP Required

Msg + Compute Nodes (2 Nodes)

Message Queue (ActiveMQ/ Solace) Compute Engine (Custom)

Message Queue (ActiveMQ/ Solace)

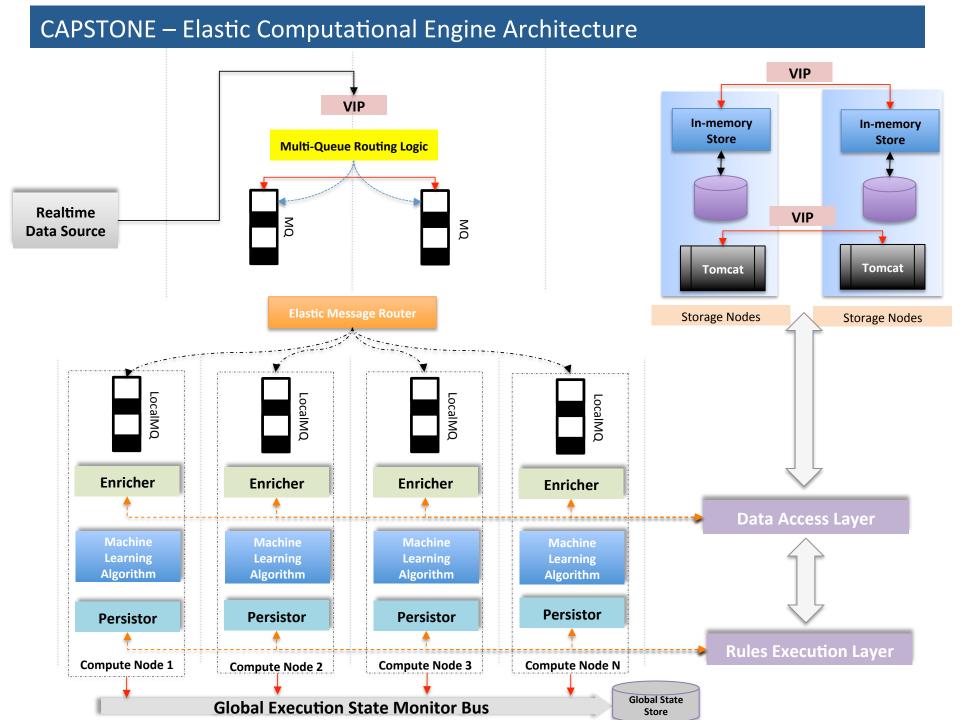
> Compute Engine (Custom)

Virtual IP Required

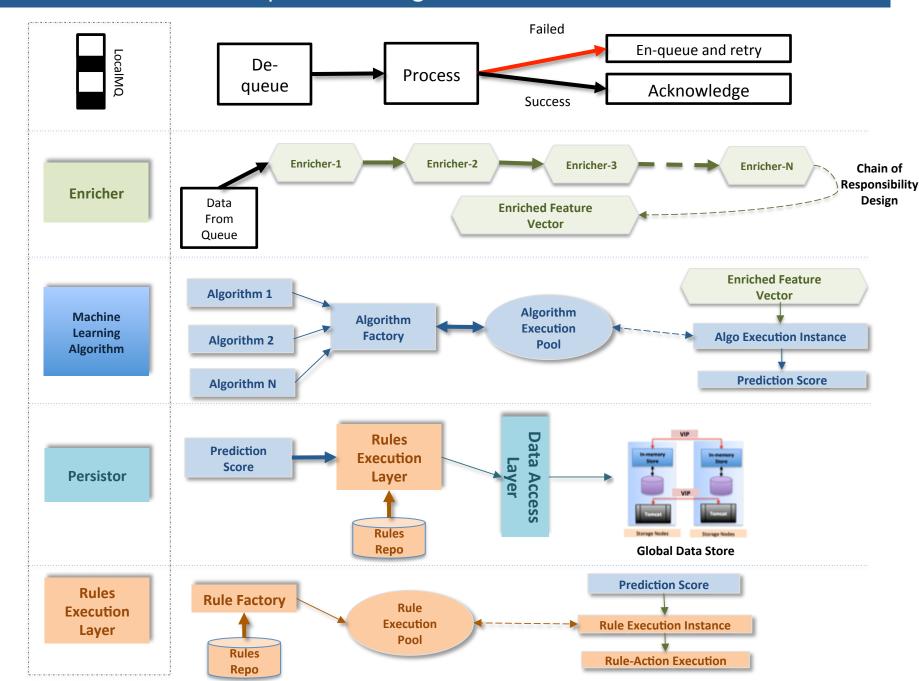
Misc Nodes – Logging, API, Self-service (2)

Logging/ Monitoring/ API (Splunk/ Tomcat)

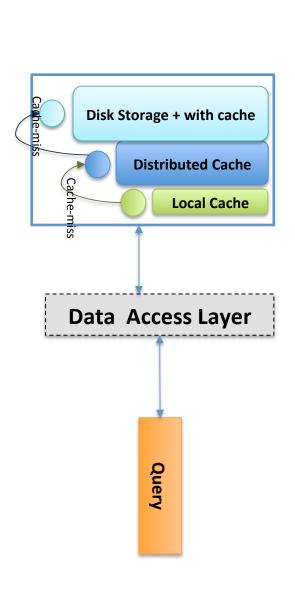
Logging/
Monitoring/
API
(Splunk/
Tomcat)

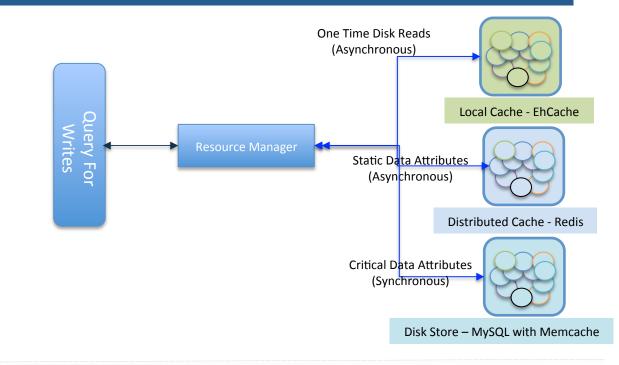


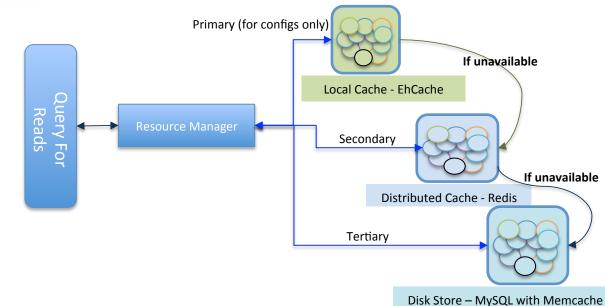
CAPSTONE – Elastic Computational Engine Architecture



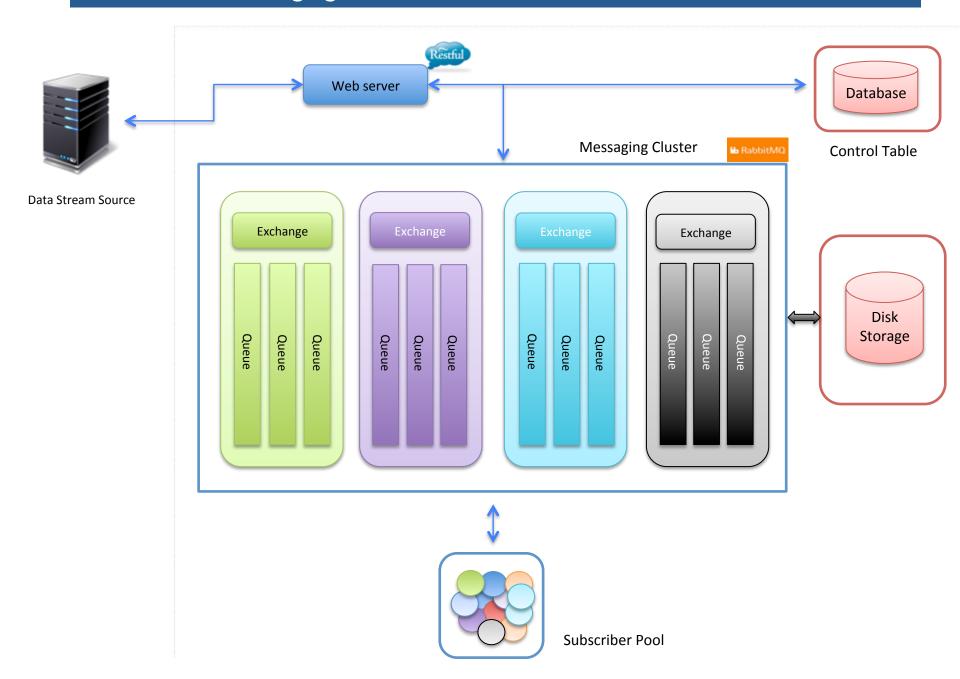
CAPSTONE – Architecting Caching Solution for Realtime







CAPSTONE – Messaging Cluster Architecture



CAPSTONE – Self Service Layer Architecture

Databases



In-Memory Store





Query Engine

Data Access Layer

Security, Authentication and Authorization Layer

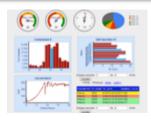
API Layer

Visualization and Reporting Layer

Performance and System Monitoring Layer

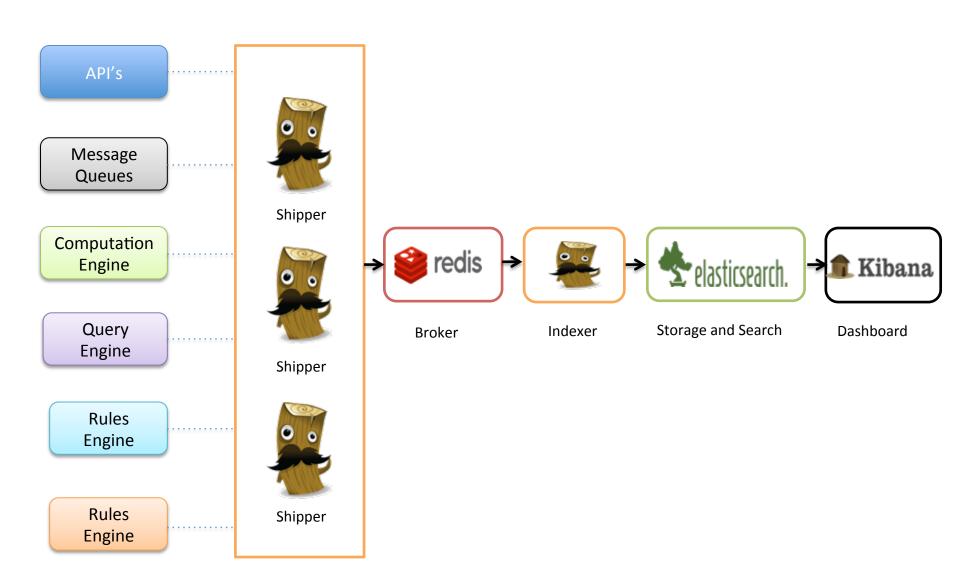
External Application Interface







CAPSTONE – Logging reporting layer



CAPSTONE – Summary of prototype