



When Distributed Database Meets Cloud Lessons Learned

Yanqing Weng

lweng@pivotal.io

About me

- Principal Software Engineer in Pivotal
- Apache HAWQ Committer
- Apache HAWQ PMC Member

Agenda

- Introduction to Distributed Database
- Distributed Database on Cloud
- Lessons Learned
- Q & A



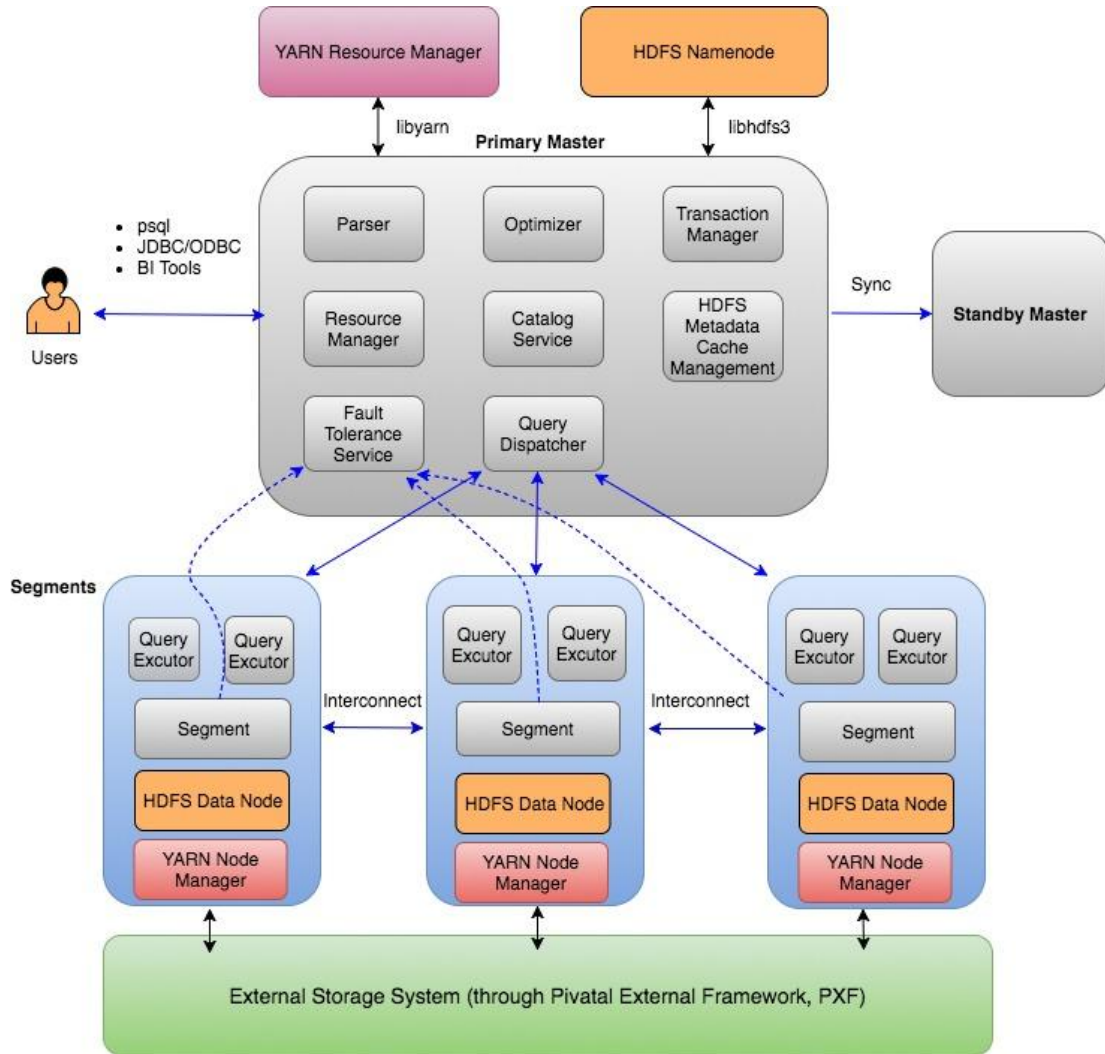


Distributed Database

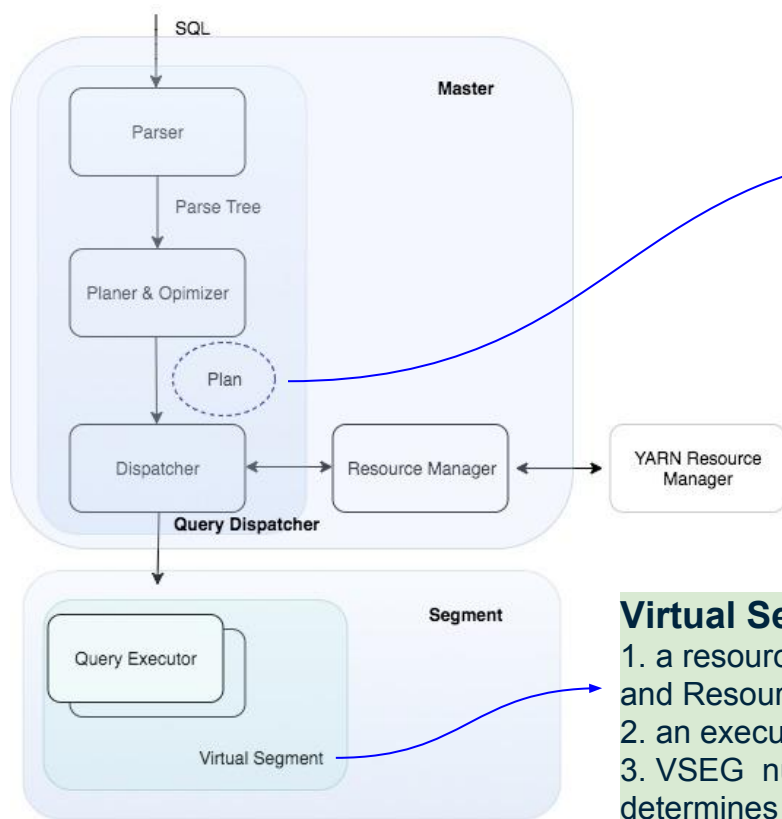
Apache HAWQ

- Apache Hadoop Native SQL, Advanced, MPP, Elastic Query Engine.
- Apache Top Level Project in 2018.8

Apache HAWQ Architecture

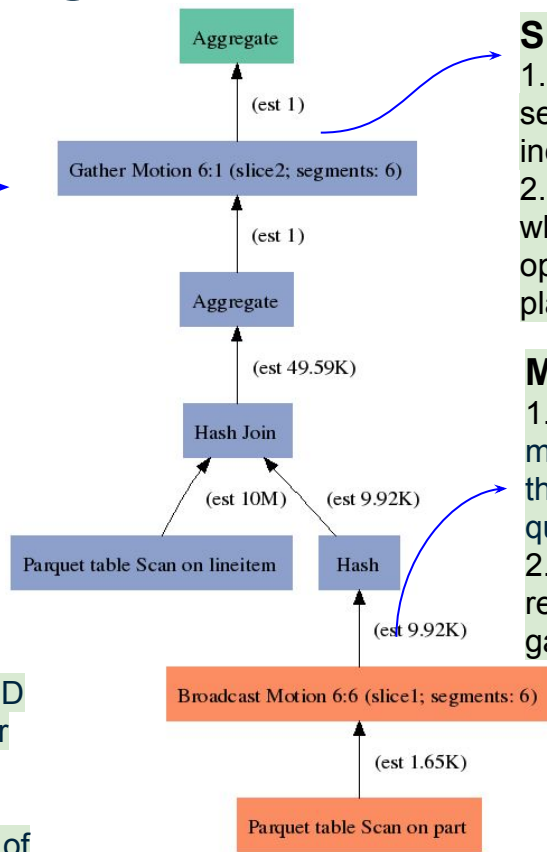


Apache HAWQ Query Processing



Virtual Segment:

1. a resource unit for QD and Resource Manager
2. an execution unit
3. VSEG number determines the degree of parallelism of a query dynamically.



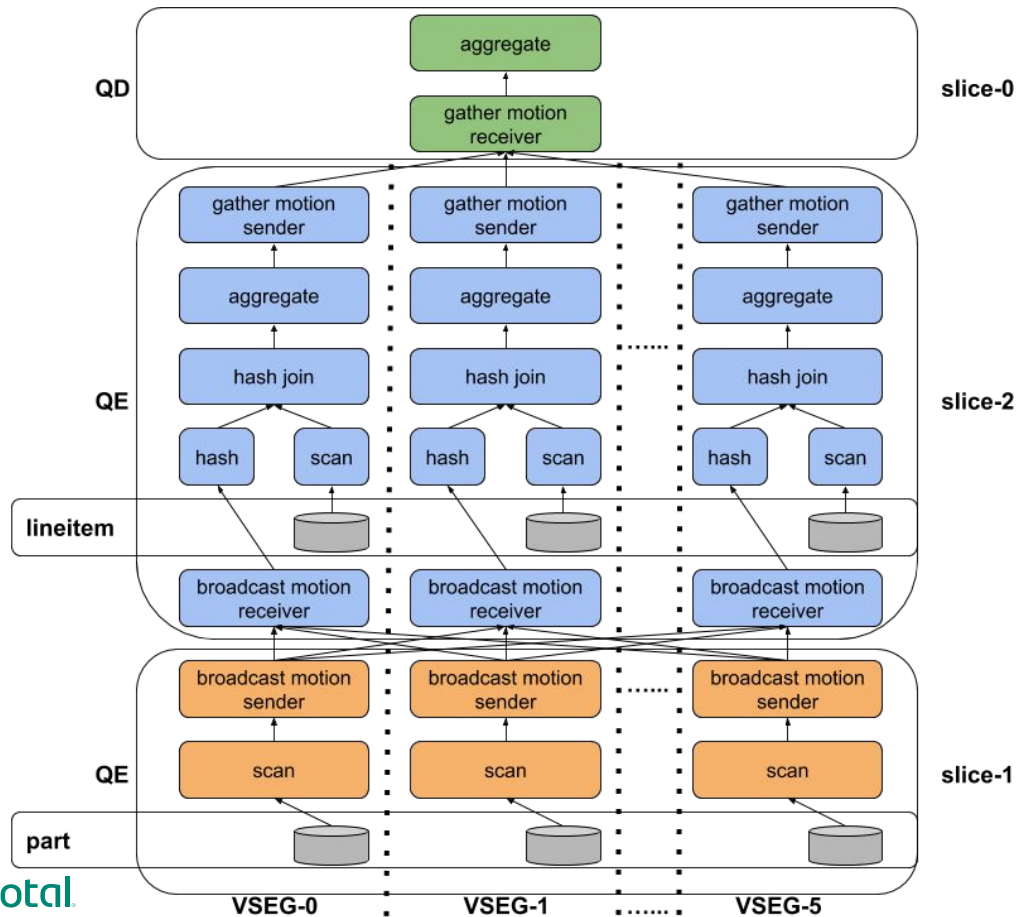
Slice:

1. a portion of the plan that segments can work on independently.
2. a query plan is sliced wherever a motion operation occurs in the plan.

Motion:

1. an operation involves moving tuples between the segments during query processing.
2. three types: redistribution, broadcast, gather motion.

Virtual Segment



- Resource allocation unit
- Query execution unit
- Variable virtual segment number
- Place on any physical segment

Summary

- High Performance
- Storage computing separation
- Fine-grained resource management
- Elastic query execution engine
- Stateless segment



Cloud Database

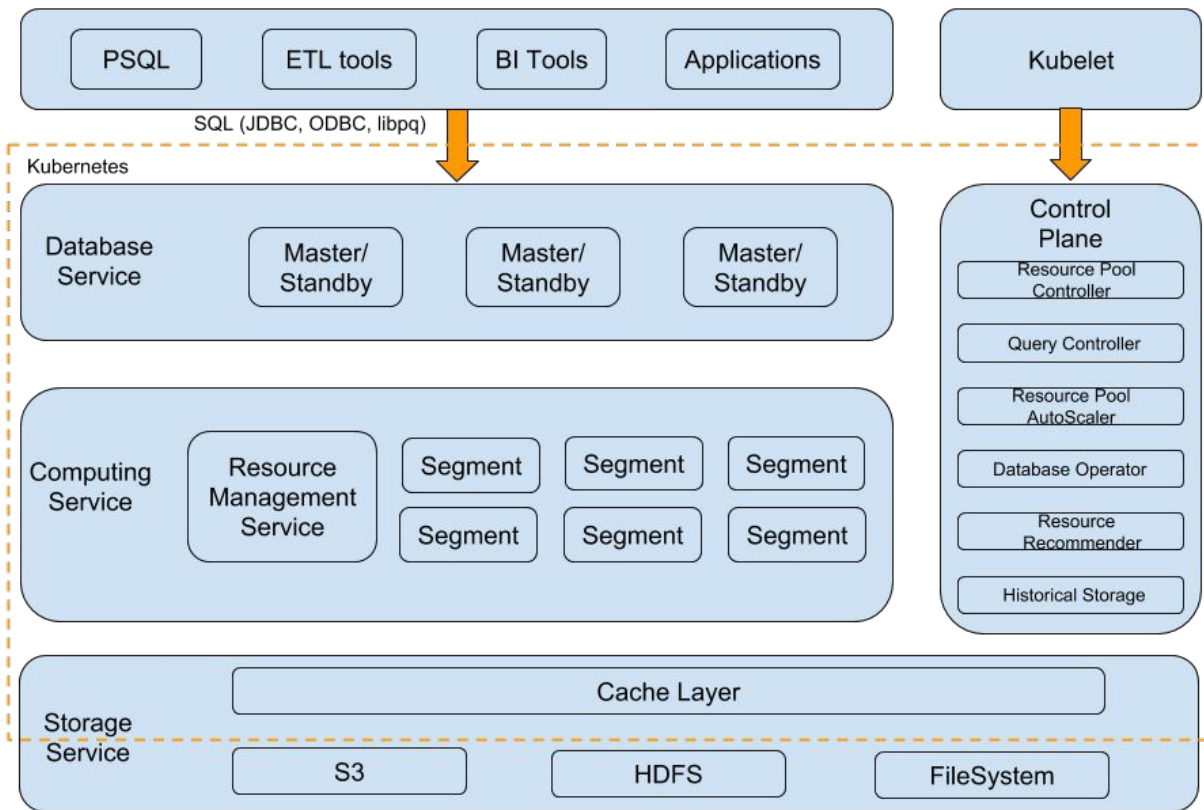
Requirement

- Database as Service
- Efficient Resource Management
- Infrastructure Agnostic
- DBA Free

Deployment & Operation

- Container VS. Virtual Machine
- Kubernetes
 - Service discovery
 - Load balancing
 - Horizontal and Vertical auto scaling
 - Rolling upgrade
 - Monitor and metrics collection
 -

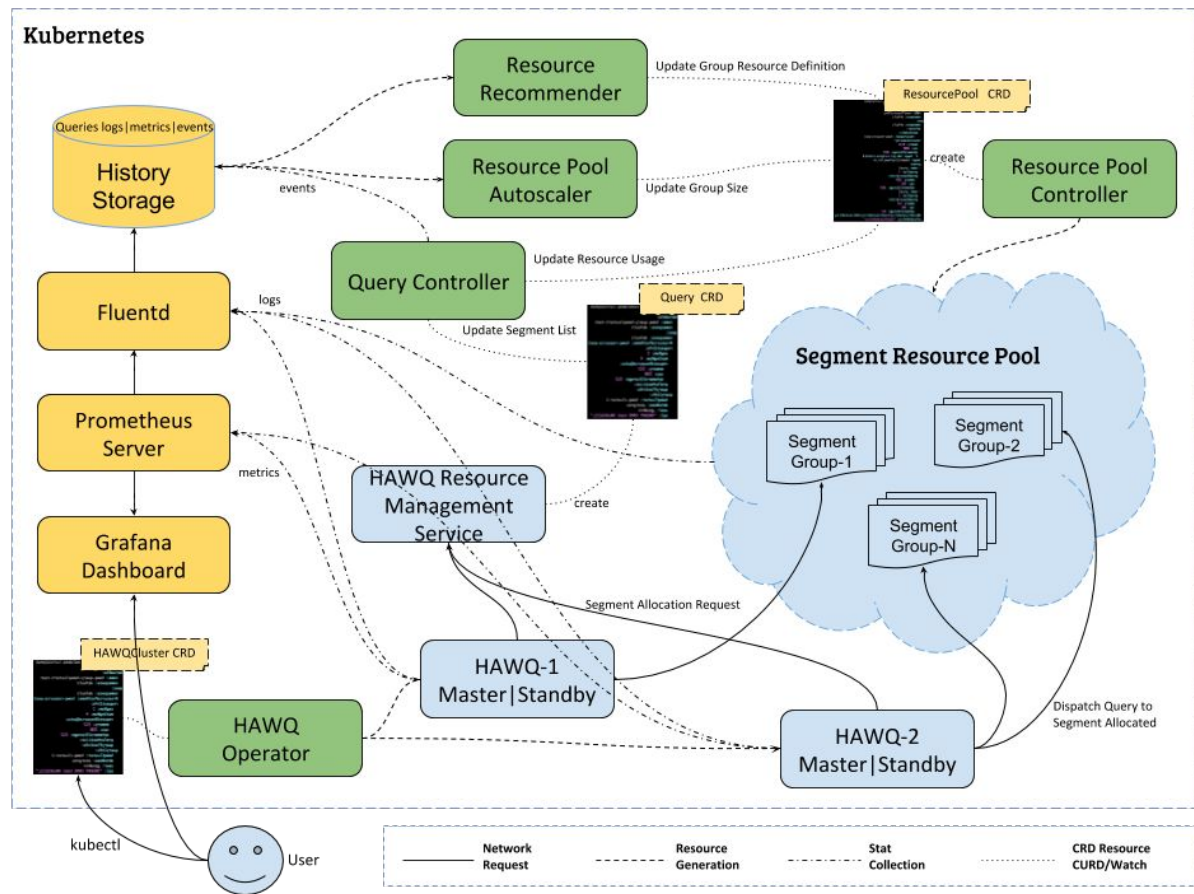
Architecture



Architecture

- Storage Service
 - Cloud Storage, Amazon S3, Hadoop.....
 - Unified Cache Lever by Alluxio
- Computing Service
 - Shared Segment Pool
 - Global Resource Management Service
- Database Service
 - Master/Standby as Database
 - Get Segments for Query on Demand
- Control Plane
 - Operator/Controllers as DBA

Apache HAWQ on Kubernetes



Custom Resource

```
1 apiVersion: pivotaldata.io/v1alpha1
2 kind: HAWQCluster
3 metadata:
4   name: hawq-cluster-demo
5 spec:
6   clusterPort: 30001
7   userGUCs:
8     hawq_master_address_port: 5432
9   masterResourceLimit:
10    cpu: 250
11    memory: 1024M
12   standbyResourceLimit:
13    cpu: 250
14    memory: 1024M
15 status:
16   clusterStatus: Ready
```

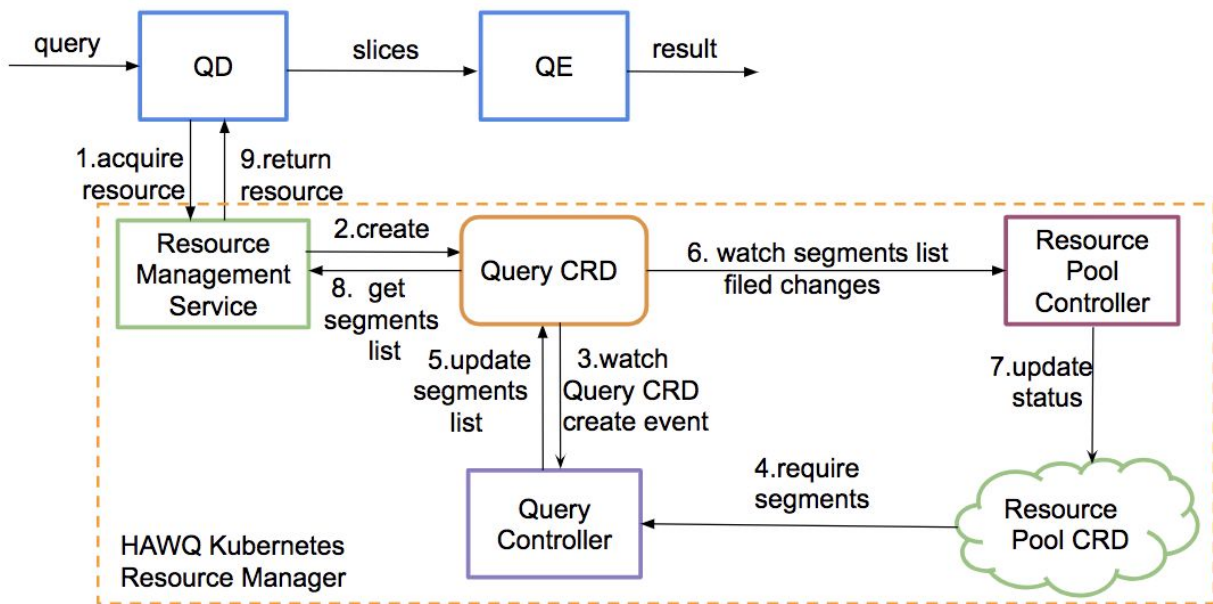
```
1 apiVersion: pivotaldata.pivotaldata.io/v1alpha1
2 kind: HAWQQuery
3 metadata:
4   name: q-10.4.1.95-074367
5 spec:
6   queryInfo:
7     dataBase: "demo"
8     hawqCluster: 10.4.1.95
9     sql: "SELECT COUNT(*) FROM lineitem;"
10    user: "hawq"
11   requestInfo:
12     maxSegNum: 1
13     preferHostList: ""
14     queryPlanInfo: ""
15     requestResourceQuota:
16       cpu: 200
17       ephemeralStorage: 512
18       memory: 512
19     segNum: 1
20   resourcePoolName: hawq-resource-pool
21 status:
22   actualUsedResource: null
23   events: null
24   queryExecuteStatus: "SUCCESS"
25   segmentList:
26   - ip: 10.4.1.93
27     name: group1-0
28   startTimeStamp: 2018-11-09T10:00:17.284390568Z
29   stopTimeStamp: 2018-11-09T10:00:25.193212063Z
```

```
1 apiVersion: pivotaldata.pivotaldata.io/v1alpha1
2 kind: HAWQResourcePool
3 metadata:
4   name: hawq-resource-pool
5 spec:
6   allocatePolicy: FastAllocatePolicy
7   autoScale:
8     maxSegmentNum: 4
9     scaleDownInterval: 30
10    scaleDownRatio: 0.25
11    scaleUpRatio: 0.75
12   groups:
13   - groupResourceLimit:
14     cpu: 250
15     ephemeralStorage: 1024
16     memory: 1024
17     groupSize: 1
18     name: group1
19   - groupResourceLimit:
20     cpu: 250
21     ephemeralStorage: 512
22     memory: 512
23     groupSize: 1
24     name: group2
25   image: hawqbeijing/hawq:demo
26   resourceCapacity:
27     cpu: 4000
28     ephemeralStorage: 3072
29     memory: 6144
30 status:
31   availableResourceStatus:
32     cpu: 500
33     ephemeralStorage: 1536
34     memory: 1536
35   resourceGroupStatus:
36   - allocatedPodNum: 0
37     availablePodNum: 1
38     name: group1
39   - allocatedPodNum: 0
40     availablePodNum: 1
41     name: group2
42   resourcePoolStatus: ready
```


Controller

- HAWQ Operator
- Resource Pool Controller
- Resource Pool AutoScaler
- Resource Recommender
- Query Controller

Apache HAWQ on Kubernetes



Lesson Learned

Architecture

- Service Oriented Architecture
 - Monolithic → Micro Service
- Resource Centric
 - Abstract Component as Resource
 - Service for Resource Usage
 - Controller for Resource Management

Containerization

- Container != Image
- Container != VM
- Container = Fine-grained Resource

Resource Management

- Traditional Database
 - Fixed resources
 - Balance resource usage among queries
- Cloud Database
 - Dynamic resources
 - Maximize resource sharing
 - Maximize resource utilization for each query

Resource Monitoring & Tuning

- Database
 - Variant Query Workload
 - Data Size
 -
- Query Similarity and Classification
- Query Resource Monitoring
 - Pod Runtime Metrics
 - Application Logs
 - Kubernetes Events
 -
- Intelligent Resource Tuning
 - Resource Pool Definition
 - Horizontal & Vertical

Kubernetes Ecosystem

- Log Collection
 - Fluentd
- Monitoring and Metrics Collection
 - Prometheus
- Visualization
 - Grafana
-

Others

- Management Utilities
- Imperative VS. Declarative
- Pod Priority
-



Thank you! Questions?