# Hot to Apply Modeling and Optimization to Select the Appropriate Cloud Platform

Boris Zibitsker

BEZNext

bzibitsker @beznext. com

# About Speaker



Boris Zibitsker

Dr. Boris Zibitsker is a CEO of BEZNext. His focus is on the development of performance assurance, performance engineering, dynamic performance management and long-term capacity planning software tools for big data, data warehouse and cloud applications.

He is a member of SPEC Big Data Research Group.
Boris consults with many Fortune 500 companies, and he manages Capstone projects for graduate students in MS in Analytics at University of Chicago.
Boris a Honorable Doctor of BGUIR and during last 5 years he was a co-chairman of Big Data Advanced Analytics Conference.

# Abstract

Organizations want to take advantage of the flexibility and scalability of Cloud platforms. By migrating to the Cloud, they hope to develop and implement new applications faster with lower cost.  Amazon AWS, Microsoft Azure, Google, IBM, Oracle and others Cloud providers support different DBMS like Snowflake, Redshift, Teradata Vantage, and others. These platforms have different architecture, mechanism of allocation and management of resources, and sophistication of DBMS optimizers which affect performance, scalability and cost.  As a result, the response time, CPU Service Time and the number of I/Os for the same query, accessing the similar table in the Cloud could be significantly different than On Prem.

In order to select the appropriate Cloud platform, we use a modeling and optimization. First, we perform a Workload Characterization for On Prem Data Warehouse. Each Data Warehouse workload represents a specific line of business and includes activity of many users generating concurrently simple and complex queries accessing data from different tables. Each workload has different demand for resources and different Response Time and Throughput Service Level Goals.
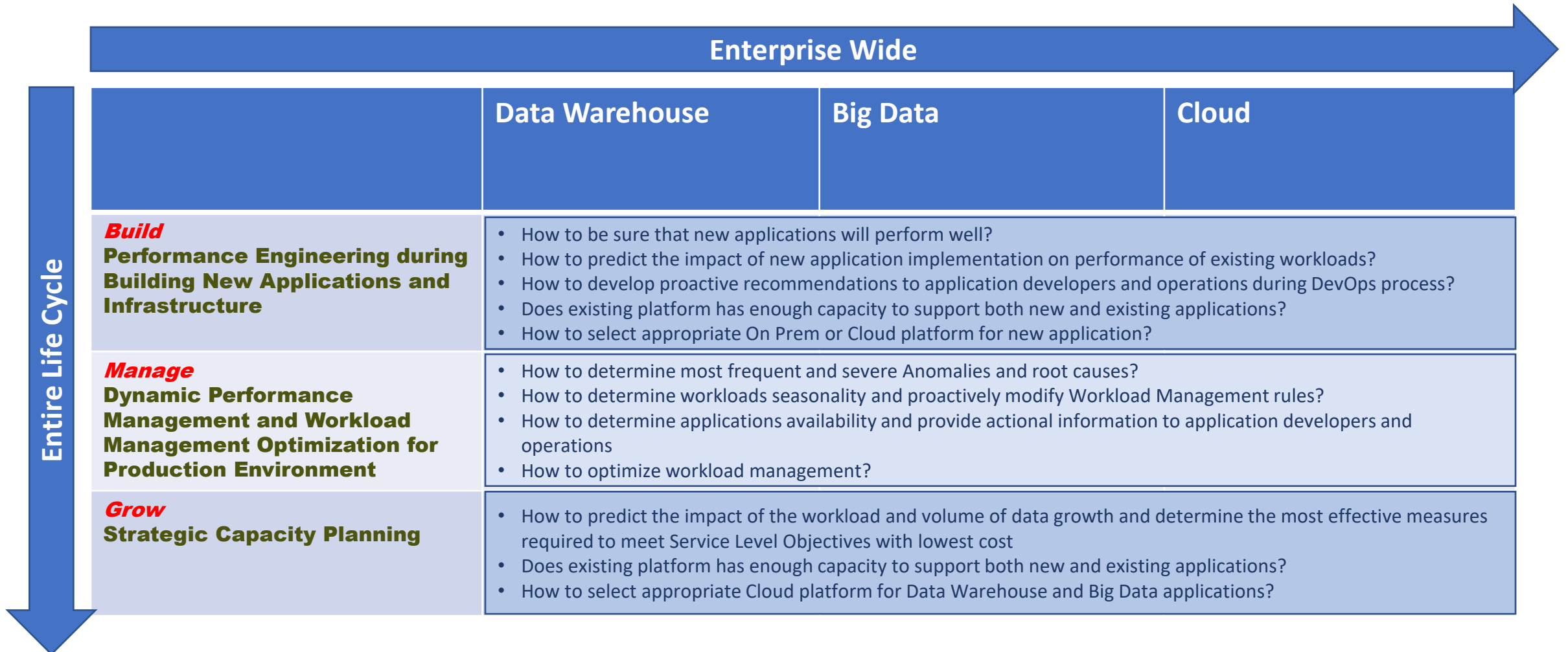
- In this paper we will review results of the workload characterization for On Prem Data Warehouse environment.

- Secondly, we must collect measurement data for standard TPC-DS benchmark tests performed in AWS Vantage, Redshift and Snowflake Cloud platform for different sizes of the data sets and different number of concurrent users.

- During third step we use the results of the workload characterization and measurement data collected during the benchmark to modify BEZNext On Prem Closed Queueing model to model individual Clouds.

- And finally, during the fourth step we use the Model to take into consideration differences in concurrency, priorities and resource allocation to different workloads. BEZNext Capacity Planning optimization algorithms incorporate Graduate search mechanism to find the AWS instance type and minimum number of instances which will be required to meet SLGs for each of the workloads. Publicly available information about the cost of the different AWS instances is used to predict the cost of supporting workloads in the Cloud month by month during next 12 months.

# Outline

- Introduction to Cloud Selection and Performance Assurance

- Data Collection and Workload Characterization

- Predicting Minimum Cloud Configurations Required to meet SLGs

- Predicting Cost

- Performance Assurance of Data Warehouse Workloads in the Cloud

- Summary

# Introduction

# BEZNext Performance Assurance Software and Services

| Enterprise Wide | | |
|---|---|---|

| | Data Warehouse | Big Data | Cloud |
|---|---|---|---|
| **Build**<br>**Performance Engineering during Building New Applications and Infrastructure** | • How to be sure that new applications will perform well?<br>• How to predict the impact of new application implementation on performance of existing workloads?<br>• How to develop proactive recommendations to application developers and operations during DevOps process?<br>• Does existing platform has enough capacity to support both new and existing applications?<br>• How to select appropriate On Prem or Cloud platform for new application? | | |
| **Manage**<br>**Dynamic Performance Management and Workload Management Optimization for Production Environment** | • How to determine most frequent and severe Anomalies and root causes?<br>• How to determine workloads seasonality and proactively modify Workload Management rules?<br>• How to determine applications availability and provide actional information to application developers and operations<br>• How to optimize workload management? | | |
| **Grow**<br>**Strategic Capacity Planning** | • How to predict the impact of the workload and volume of data growth and determine the most effective measures required to meet Service Level Objectives with lowest cost<br>• Does existing platform has enough capacity to support both new and existing applications?<br>• How to select appropriate Cloud platform for Data Warehouse and Big Data applications? | | |

**Entire Life Cycle**

6

# Criteria of Cloud Platform Selection

## MULTIPLE CRITERIA

Performance

Scalability

Cost

Security

Elasticity

Deployment Flexibility

Ecosystem Integration

Database Management

Analytic and Database Functionality

## WE WILL FOCUS ON

Performance and

Cost

# Major Steps

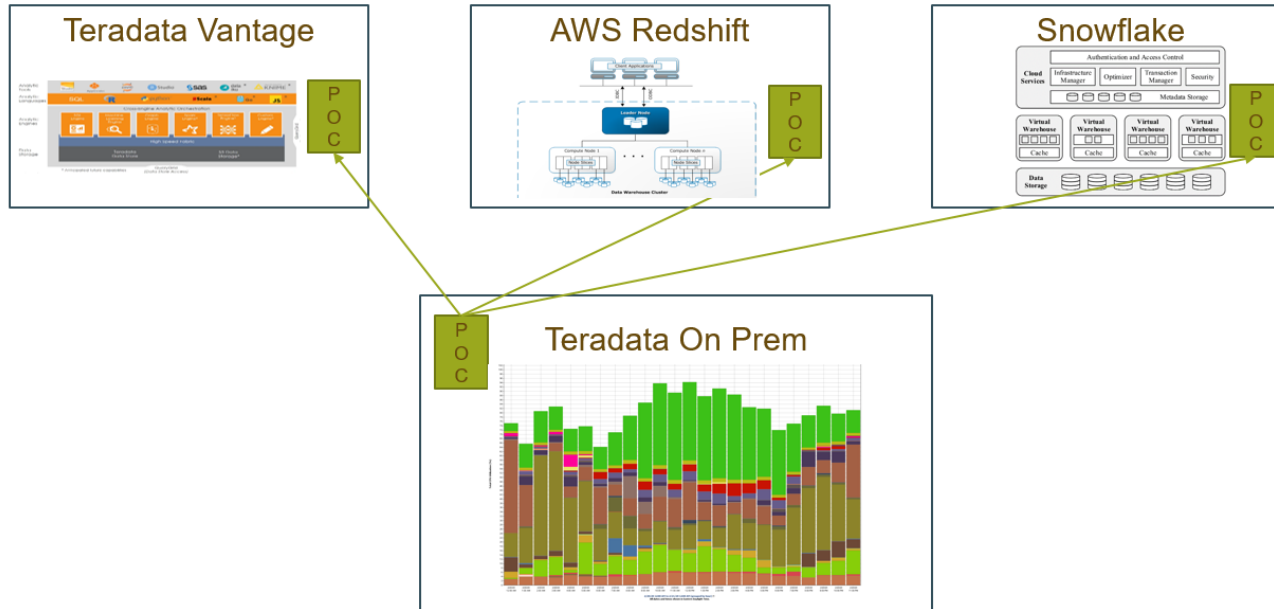1. Data Collection and Workload Characterization

2. Modeling & Optimization

Predicting the Minimum Configuration Required
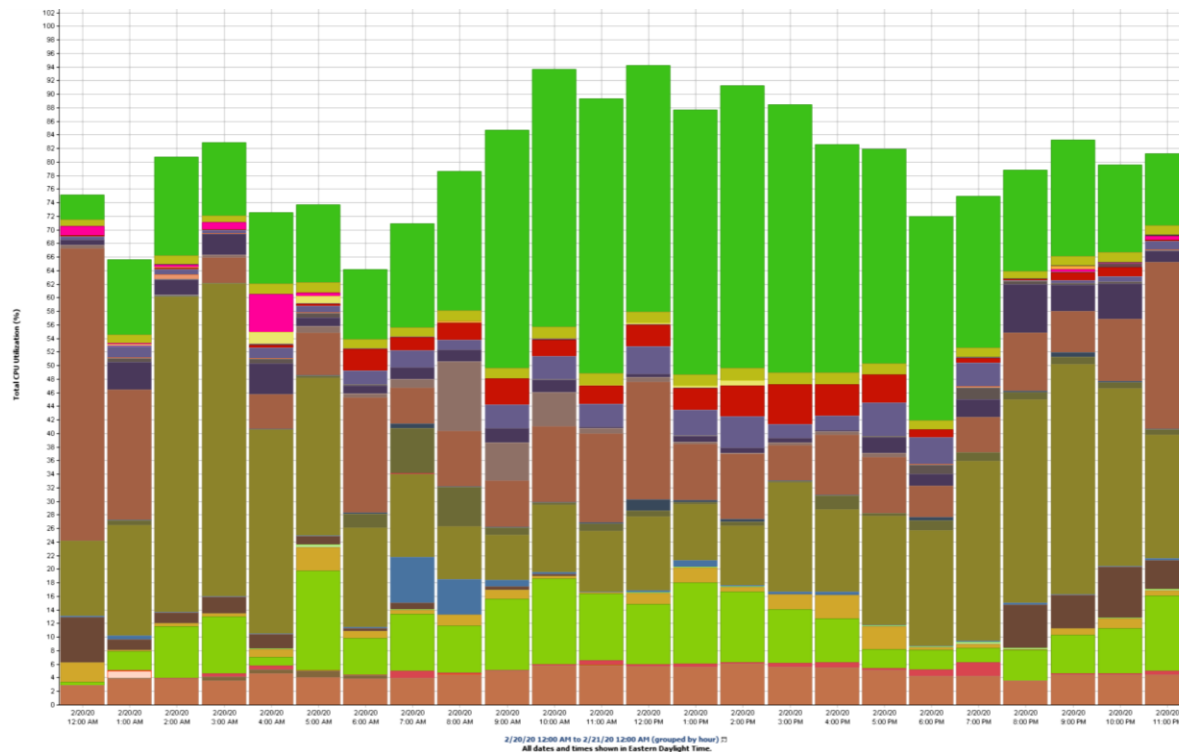
3. Predicting Cost

4. Cloud Platform Selection

# Data Collection and Workload Characterization On Prem and in the Cloud

# Data Collection On Prem and Cloud Platforms



Teradata Vantage

AWS Redshift

Snowflake

Teradata On Prem

- Production workload On Prem

- POC Benchmark with representative Queries on each Platform

- Standard TPC DS Benchmark on each Platform

# CPU Utilization by Production On Prem Workloads During 24 Hours



**Hourly Profiles for Each Workload are built during Workload Characterization:**

- Performance Profile
- Resource Consumption Profile
- Data Usage Profile

# CPU Utilization by Workloads Selected for Cloud

# Applying Modeling and Optimization to Determine the Minimum Configuration Required to meet SLG for each Workload

# Modeling Approach

On Prem Queueing Network Model for MPP Data Warehouse reflects Hardware,   Software configurations, Workload Management parameters and workload characterization results

Cloud Models are built by modifying parameters of On Prem models to reflect specific architecture of each Cloud DBMS platform and results of the benchmarks



Teradata Vantage

AWS Redshift

Snowflake

Teradata On Prem

# Determining the Minimum Cloud Configuration Capable Meeting SLGs

## GRADIENT OPTIMIZATION



## APPROACH

1. Apply Gradient Optimization to Workload Management

2. Use modeling to predict which workload will violate SLGs the most

3. Determine which resource will be the bottleneck

4. Apply Gradient method and iterative modeling to find the minimum hardware configuration required to meet SLGs for all workloads

# Predicting performance of Data Warehouse workloads in Teradata Vantage environment

- AWS Instance selection
- Limited scalability, but sophisticated optimizer and workload management

# Vantage Workload Management Optimization



OPTIMIZATION OF CAPACITY PLANNING, PERFORMANCE MANAGEMENT AND WORKLOAD MANAGEMENT

LIMIT CONCURRENCY REDUCE CONTENTION BUT INCREASE the # of REQUESTS WAITING for the THREAD

PREDICTED IMPACT OF INCREASING PRIORITY FOR SALES

The same approach can be applied to predict the Impact of Changing Classifications and answer questions like: What if Users SQL required less than 1 sec CPU Time will be Processed with Higher Priority...

# Predicting performance of Data Warehouse workloads in AWS Redshift environment

- Limited number of nodes/instances in Redshift cluster

- Future release of Redshift will use AQUA to accelerate Redshift queries by running data intensive tasks such as filtering and aggregation, compression and others closer to the storage layer.

- We did not model the impact of AQUA on Redshift performance

Current Redshift Architecture

Advanced Query Accelerator (AQUA) for Amazon Redshift

# Predicting performance of Data Warehouse workloads in Snowflake environment

- Snowflake automatically scales out and scale up

- First scenario - running each workload in dedicated Virtual Warehouse

- Second scenario - running all workloads in one Virtual Warehouse

# Examples of Modeling and Optimization Predicting Min Configurations Required to meet SLGs for Each Cloud, 2$^{nd}$ shift during next 10 Month

MIGRATE 4  WORKLOADS TO THE CLOUD

EXPECTED GROWTH IN NUMBER OF USERS - 12% ANNUALLY

EXPECTED GROWTH IN VOLUME OF DATA PROCESSED - 10% ANNUALLY

# Vantage 2nd Shift / Month

MIGRATE 4  WORKLOADS INTO ONE VIRTUAL WAREHOUSE:

# Vantage: Recommended minimum configuration for 2nd shift during next 12 months

# As a result of changing configuration according recommendation  and workload priorities, Vantage Response Time will meet SLGs During next 12 months

Service Level Goal

2

Adding Instances

3

4

1

Migration from
On Prem to Vantage

## Optimized Workload Management Priorities

| | | Sales | Marketing | Finance | BI |
|---|---|---|---|---|---|
| Current Priority | Current | 24.99 | 38.54 | 39.66 | 16.58 |
| Optimized Priority | Month 1 | 17.45 | 29.09 | 18.25 | 50.76 |
| Optimized Priority | Month 2 | 17.26 | 28.68 | 18.13 | 49.73 |
| Optimized Priority | Month 3 | 13.83 | 22.92 | 14.59 | 39.52 |
| Optimized Priority | Month 4 | 12.05 | 21.05 | 12.14 | 35.08 |
| Optimized Priority | Month 5 | 10.84 | 18.52 | 11.47 | 31.38 |
| Optimized Priority | Month 6 | 10.97 | 18.44 | 11.47 | 31.38 |
| Optimized Priority | Month 7 | 10.98 | 19.57 | 11.36 | 31.36 |
| Optimized Priority | Month 8 | 10.96 | 18.55 | 11.35 | 31.36 |
| Optimized Priority | Month 9 | 11.07 | 18.7 | 11.52 | 31.38 |
| Optimized Priority | Month 10 | 10.98 | 18.75 | 11.22 | 31.46 |
| Optimized Priority | Month 11 | 11.02 | 18.76 | 11.35 | 31.07 |
| Optimized Priority | Month 12 | 11.03 | 18.76 | 11.35 | 31.07 |

# Redshift 2nd Shift

MIGRATE 4 WORKLOADS TO THE REDSHIFT:

CLIENT REPORTING, USER REPORTING, PROGRAM INTEGRITY AND PROVIDER ECONOMICS

# As a result of changing configuration according recommendation, Redshift Response Time will meet SLGs during next 12 months
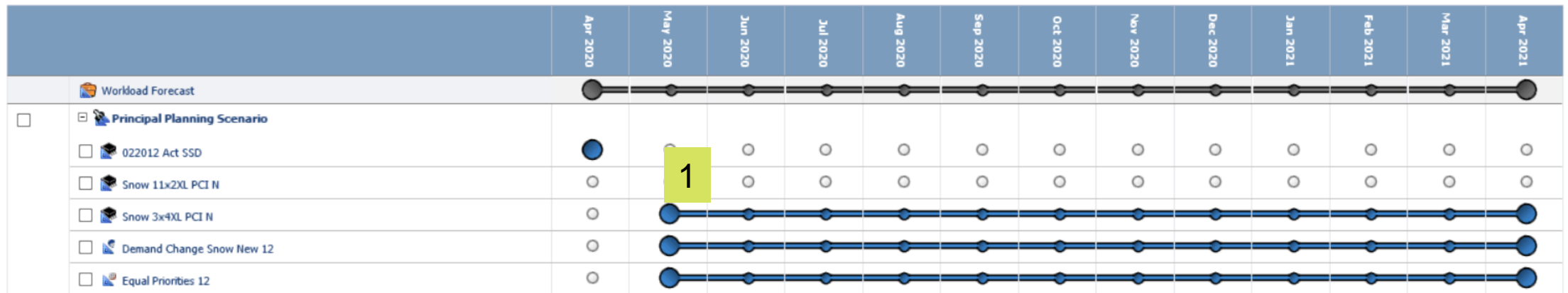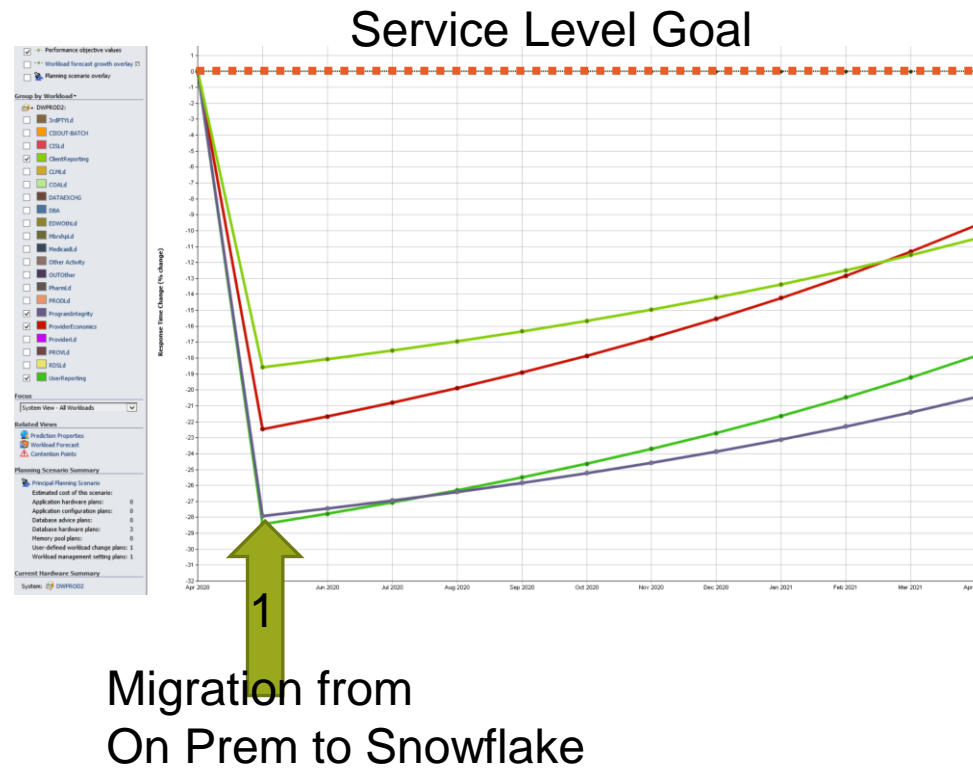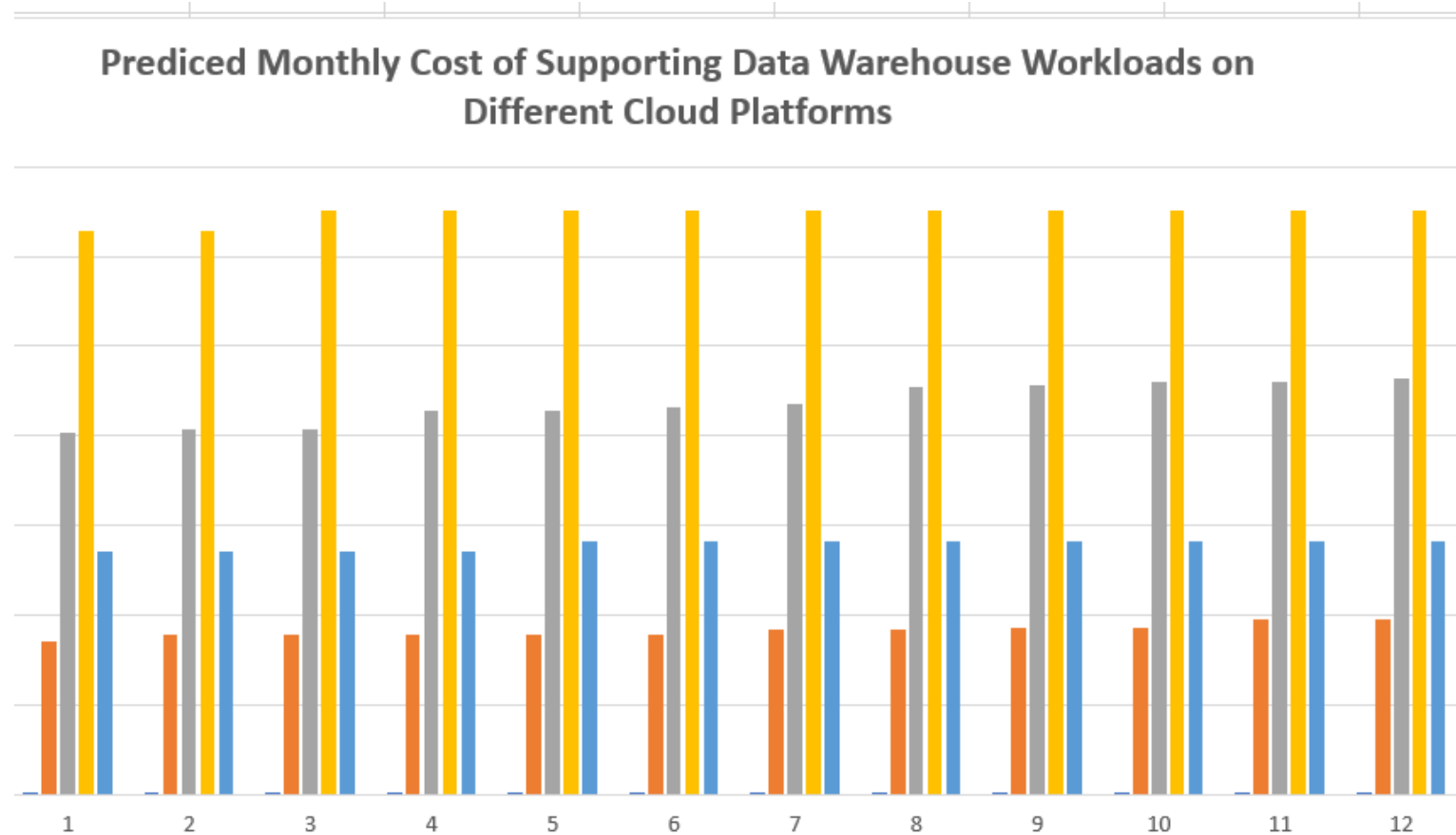
# Predicted Redshift Relative Response Time Change

# Snowflake 2$^{nd}$ Shift / Month

MIGRATE 4  WORKLOADS INTO ONE SNOWFLAKE VIRTUAL WAREHOUSE:

CLIENT REPORTING, USER REPORTING, PROGRAM INTEGRITY AND PROVIDER ECONOMICS

# As a result of changing configuration according recommendation, Snowflake Response Time will meet SLGs during next 12 months

| | Apr 2020 | May 2020 | Jun 2020 | Jul 2020 | Aug 2020 | Sep 2020 | Oct 2020 | Nov 2020 | Dec 2020 | Jan 2021 | Feb 2021 | Mar 2021 | Apr 2021 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Workload Forecast | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● |
| Principal Planning Scenario | | | | | | | | | | | | | |
| 022012 Act SSD | ● | | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ |
| Snow 11x2XL PCI N | ○ | 1 | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ |
| Snow 3x4XL PCI N | ○ | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● |
| Demand Change Snow New 12 | ○ | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● |
| Equal Priorities 12 | ○ | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● | ● |

# Predicted Relative Response Time Change



Service Level Goal

Migration from
On Prem to Snowflake

# Predicted # of Instances and Instance Type required to meet SLGs for each Cloud / Shift / Month

| Prredicted Minimum Number of Instances per Cloud/ Shift / Month | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Month | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| **Vantage** | | | | | | | | | | | | |
| 1st Shift | | | | | | | | | | | | |
| 1st Shift Min # Instances | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 11 | 11 |
| 2nd Shift | | | | | | | | | | | | |
| 2nd Shift Min # Instances | 32 | 34 | 34 | 34 | 34 | 34 | 36 | 36 | 36 | 36 | 38 | 38 |
| 3rd Shift | | | | | | | | | | | | |
| 3rd Shift Min # Instances | 13 | 13 | 13 | 13 | 13 | 13 | 13 | 13 | 14 | 14 | 14 | 14 |
| **Redshift** | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| 1st Shift | | | | | | | | | | | | |
| Min # ra3 Instances | 52 | 52 | 52 | 54 | 54 | 54 | 56 | 56 | 58 | 58 | 58 | 60 |
| 2nd Shift | | | | | | | | | | | | |
| Min # ra3 Instances | 130 | 130 | 130 | 140 | 140 | 140 | 140 | 150 | 150 | 150 | 150 | 150 |
| 3rd Shift | | | | | | | | | | | | |
| Min # ra3 Instances | 72 | 74 | 74 | 76 | 76 | 78 | 78 | 80 | 80 | 82 | 82 | 82 |
| **Snowflake 4 workloads** | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| 1st Shift | | | | | | | | | | | | |
| Instance Type | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL |
| Min # Instances | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| 2nd Shift | | | | | | | | | | | | |
| Instance Type | 4XL | 4XL | 4XL | 4XL | 4XL | 4XL | 4XL | 4XL | 4XL | 4XL | 4XL | 4XL |
| Min # Instances | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 |
| 3rd Shift | | | | | | | | | | | | |
| Instance Type | 3XL | 3XL | 3XL | 3XL | 3XL | 3XL | 3XL | 3XL | 3XL | 3XL | 3XL | 3XL |
| Min # Instances | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| **Snowflake 3 workloads** | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| 1st Shift | | | | | | | | | | | | |
| Instance Type | XL | XL | XL | XL | XL | XL | XL | XL | XL | XL | XL | XL |
| Min # Instances | 6 | 6 | 6 | 6 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |
| 2nd Shift | | | | | | | | | | | | |
| Instance Type | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL | 2XL |
| Min # Instances | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| 3rd Shift | | | | | | | | | | | | |
| Instance Type | XL | XL | XL | XL | XL | XL | XL | XL | XL | XL | XL | XL |
| Min # Instances | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |

# Predicted Monthly Cost to Maintain SLGs on Different Cloud Platforms



Prediced Monthly Cost of Supporting Data Warehouse Workloads on Different Cloud Platforms

# Performance Engineering focusing on Devops
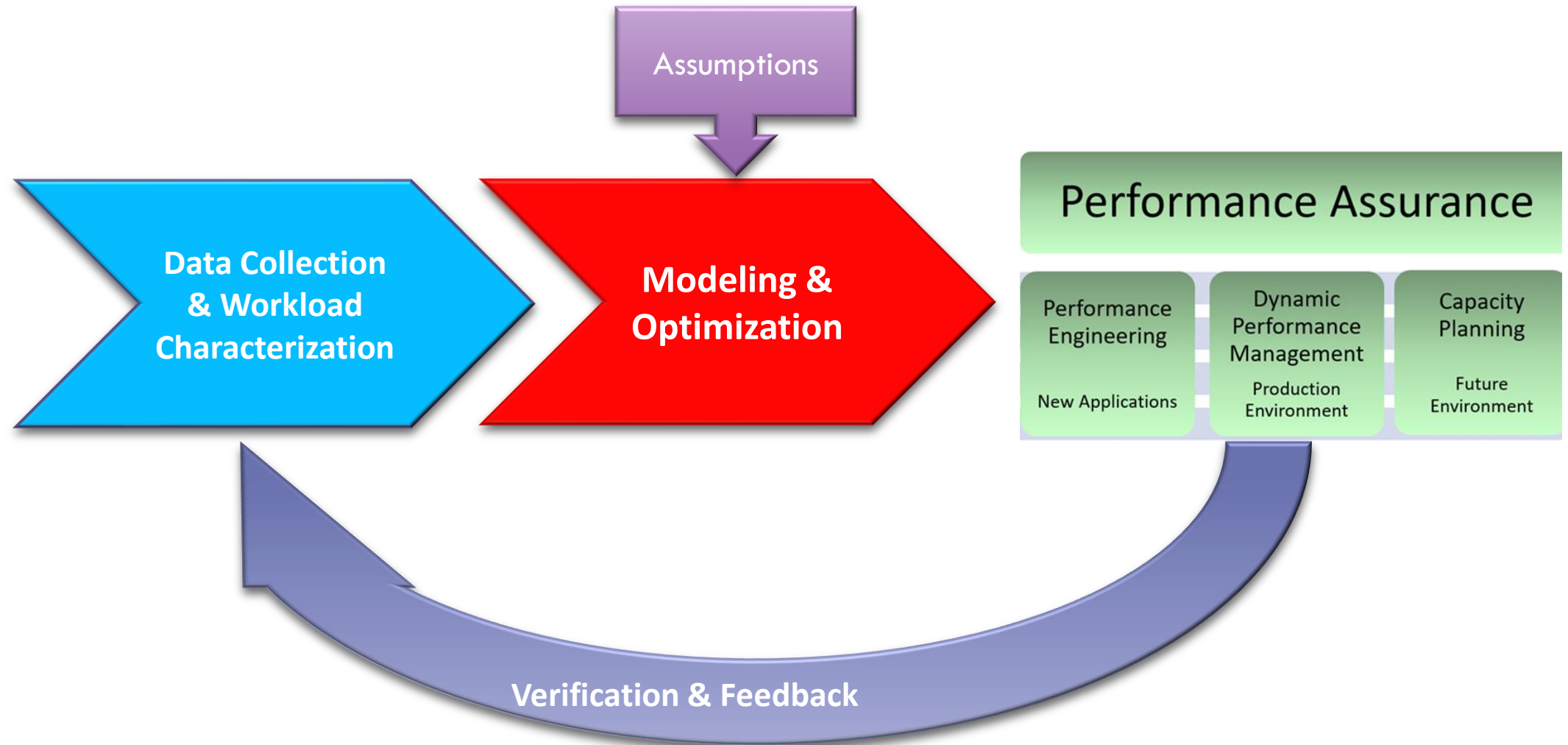
# ROLE OF MODELING DURING APPLICATION DEVELOPMENT

**App Development**

Modeling

| Plan | Code | Create | Test |

**Operations**

| Release | Deploy | Operate |

DevOps

Continuous Deployment

Continuous Delivery

Continuous Integration

Agile Development

- **Predict new applications implementation impact**
  - Predict how new application will perform in production environment
  - Identify Anomalies and their Root Causes during testing of new applications
  - Develop recommendations to Application Developers

- **Predict how new application will affect existing production applications**
  - Predict how implementation of new applications will affect Response Time and Throughput of existing applications
  - Develop capacity planning recommendations
  - Set up realistic expectations

# ROLE OF MODELING FOR OPERATIONS

**App Development**

Modeling

Plan | Code | Create | Test

**Operations**

Release | Deploy | Operate

← DevOps →

← Continuous Deployment →

← Continuous Delivery →

← Continuous Integration →

← Agile Development →

- **Develop Proactive Performance Management and Workload Management Recommendations**

  - Compare performance measurement results after implementation of the new application with expected

  - Develop proactive performance tuning recommendations

  - Develop proactive workload management recommendations

  - Reevaluate Capacity Planning recommendations

# MODELING IS A BASE FOR PERFORMANCE ASSURANCE FOR DEVOPS

# TEN STEPS OF APPLYING MODELING TO OPTIMIZE APPLICATION DEVELOPMENT AND OPERATIONAL DEVOPS DECISIONS



When RT will exceed SLG? What should be done proactively?

Response Time

SLG

Workload Growth

Delay Time (D)

Queueing Time (Q)

Service Time (S)

# FIRST STEP
## DATA COLLECTION DURING PERFORMANCE TESTING AND FOR PRODUCTION WORKLOADS

Data Collection during Performance Testing of New Application on Test System and for all workloads in Production Environments

**Measurement Data Types**

- Hardware and Software Configuration

- Response Time

- Throughput

- CPU Utilization and CPU Service Time per request

- Disk Utilization, I/O rate , #I/O operations per request and KB/Request, Channel Utilization

- Memory utilization

- Network utilization

- Level of concurrency

Test

New Application

Production

Sales

HR

Marketing

# SECOND STEP
## WORKLOAD CHARACTERIZATION

Test and Production Environments

# THIRD STEP
## ANOMALY AND ROOT CAUSE DETECTION

# FOURTH STEP
## WORKLOAD FORECASTING FOR NEW AND PRODUCTION WORKLOADS

Test and Production Environment

Expected Workload and Volume of Data Growth

# FIFTH STEP
## PREDICTING IMPACT OF EXPECTED WORKLOAD AND VOLUME OF DATA GROWTH IN PRODUCTION ENVIRONMENT

# SIX STEP
## PREDICTING IMPACT OF NEW APPLICATION IMPLEMENTATION

# SEVENTH STEP
## PREDICTING IMPACT OF THE WORKLOAD MANAGEMENT OPTIMIZATION
## WORKLOAD MANAGEMENT OPTIMIZATION WILL NOT BE SUFFICIENT TO MEET SLG

# EIGHTH STEP
## PREDICTING MINIMUM ON PREM UPGRADE REQUIRED TO MEET SLG AFTER NEW APPLICATION IMPLEMENTATION
### ADDITIONAL 14 NODES WILL BE REQUIRED TO MEET SLG

# NINTH STEP
## DETERMINING APPROPRIATE CLOUD PLATFORM FOR NEW APPLICATION

BEZNext Approach to Selection of the Appropriate Cloud

- ❑ Predict the minimum configuration required to meet SLGs

  - ▪ Instance type and # of instances which will be required Hour by Hour, Shift by Shift, Month by Month to meet SLGs for each of On Prem Production workload on each of the optional Cloud Platform

- ❑ Predict cost of running On Prem Data Warehouse Workloads on each of the optional Cloud Platforms

- ❑ Select Cloud platform capable to meet SLGs for all of the growing workloads with the lowest cost

Test

New Application

Production

Sales

New Application

HR

Marketing

# TENTH STEP
## AUTOMATIC RESULT VERIFICATION AND CREATION OF CONTINUOUS PERFORMANCE ASSURANCE PROCESS

# HOW TO OPTIMIZING DATA WAREHOUSE AND BIG DATA APPLICATIONS PERFORMANCE ON PREM AND IN THE CLOUD

# DYNAMIC PERFORMANCE MANAGEMENT FOR DATA WAREHOUSES, AND BIG DATA APPLICATIONS ON PREM AND IN THE CLOUD

- ❑ Set realistic Service Level Goals
  - ❑ Formal SLG are based on business requirements
  - ❑ Informal SLGs are based on analysis of historical data
  - ❑ Without SLG impossible to manage and plan effectively

- ❑ Data collection and Workload Aggregation
  - ❑ Automatic data collection across all platform and transforming to universal format reduce time required to coordinate and interpret data
  - ❑ WAG by line of business allows to present results of analysis and recommendation clear to business people and IT management

- ❑ Workload Characterization
  - ❑ Automatic generation of Performance, Resource utilization and Data usage by Line of Business/Workloads enables automation of identification of problems and their root causes and use of modeling and optimization to generate proactive recommendations, including determining:
    - ▪ The most frequent performance anomalies/problems and their root causes
    - ▪ Pattern and balance of performance and resource utilization
    - ▪ Application availability
    - ▪ Seasonality for each workload

- ❑ Evaluate solutions for fixing the problems

- ❑ Verify results

# DETERMINE MOST FREQUENT ANOMALIES AND ROOT CAUSES TO NARROW DOWN THE SCOPE OF PERFORMANCE TUNING
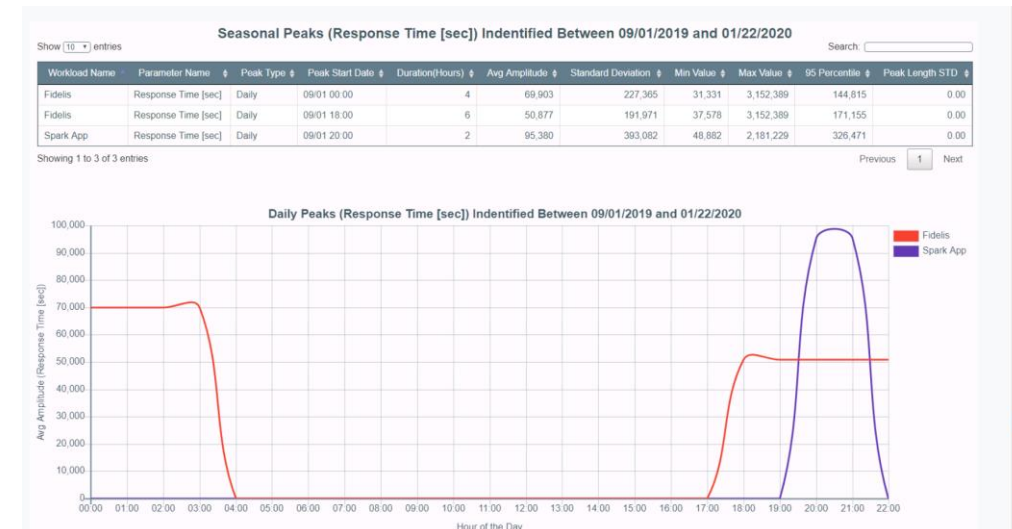
## Data Warehouse



## Big Data

# DETERMINE SEASONAL RESOURCE UTILIZATION PEAKS
## TO OPTIMIZE WORKLOAD MANAGEMENT AND RESOURCE ALLOCATION RULES ON PREM AND IN THE CLOUD (TASM, YARN, ALLOCATION AND DEALLOCATION RESOURCES IN THE CLOUD)
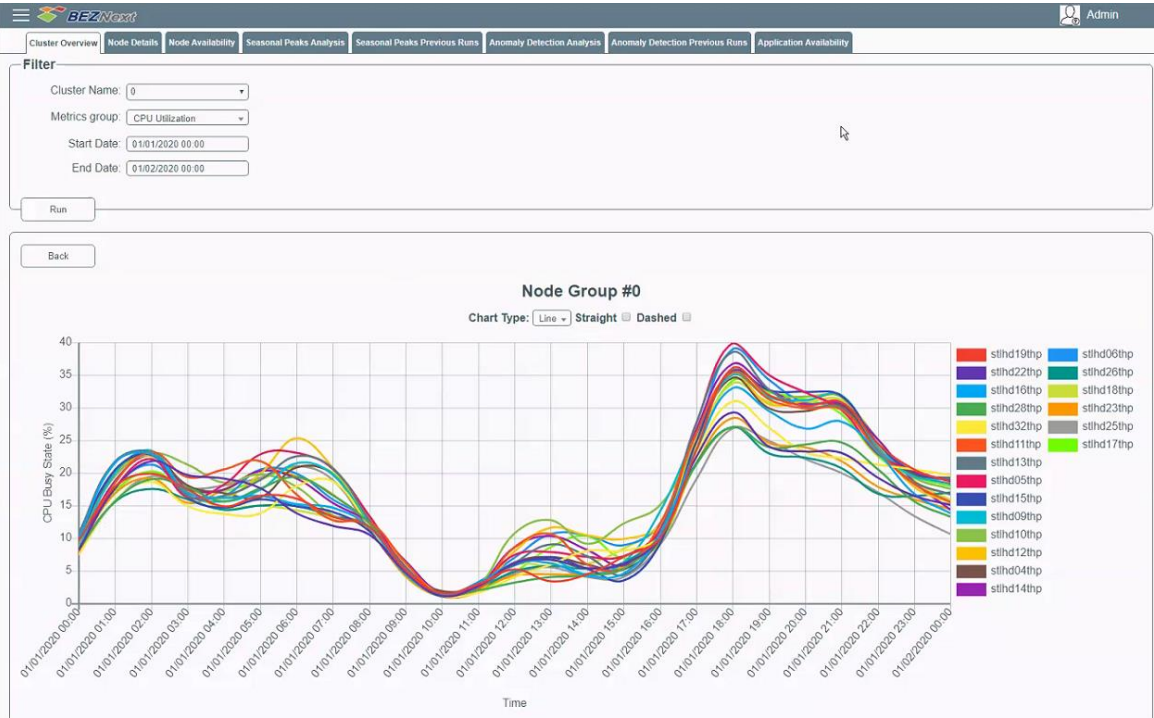
Data Warehouse

Big Data

# ANALYZE BIG DATA NODE UTILIZATION VARIABILITY
## RECOMMEND CHANGING OF DATA AND APPLICATIONS PLACEMENT TO IMPROVE RESOURCE UTILIZATION BALANCE

Difference in Big Data high and low nodes utilization    Big Data Top 20 Nodes Utilization in Time

# MEASURE APPLICATIONS AVAILABILITY
## IDENTIFY APPLICATIONS WITH THE HIGHEST FREQUENCY OF FAILURES, WASTING RESOURCES - CANDIDATES FOR TUNING

**Filter**

Cluster: BD01

Start Date: 01/01/2020 00:00

End Date: 01/07/2020 00:00

Run

**Results**

**Application Availability Summary**

Show 10 entries      Search:

| Application Name | Number of Failures | CPU time of Failed Runs | Number of Success Runs | CPU time of Success Runs |
|---|---|---|---|---|
| QueryResult.jar | 5 | 2144 | 982 | 733589 |
| HIVE-027d3c1c-7b97-4ee0-98cf-fafd7b2edc8a | 1 | 0 | 0 | 0 |
| HIVE-403c6d3c-fbac-448e-b044-f2b9b073468a | 1 | 0 | 0 | 0 |
| AUTH PROVIDER AFFILIATION FOR ADMITTING_AFFILIATION | 0 | 0 | 7 | 13375 |
| AUTH PROVIDER AFFILIATION FOR ATTENDING_AFFILIATION | 0 | 0 | 7 | 12943 |
| AUTH PROVIDER AFFILIATION FOR FACILITY_AFFILIATION | 0 | 0 | 6 | 12087 |
| AUTH PROVIDER AFFILIATION FOR PCP_AFFILIATION | 0 | 0 | 6 | 12932 |
| AUTH PROVIDER AFFILIATION FOR REFERRING_AFFILIATION | 0 | 0 | 6 | 11922 |
| CORE_DBO.PHONE.jar | 0 | 0 | 6 | 158 |
| distcp | 0 | 0 | 12 | 1562 |

Showing 1 to 10 of 115,878 entries      Previous 1 2 3 4 5 … 11588 Next

# SUMMARY

- We reviewed how modeling and optimization technology in predicting the minimum configurations required for each Cloud to meet SLGs for growing workloads during next 12 months and how to predict the corresponding cost.

- This approach can be used for other Cloud platform

- We also reviewed role of Performance Engineering for new applications and Dynamic Performance management of production workloads in the Cloud

# Thank you! Questions?

bzibitsker@beznext.com