

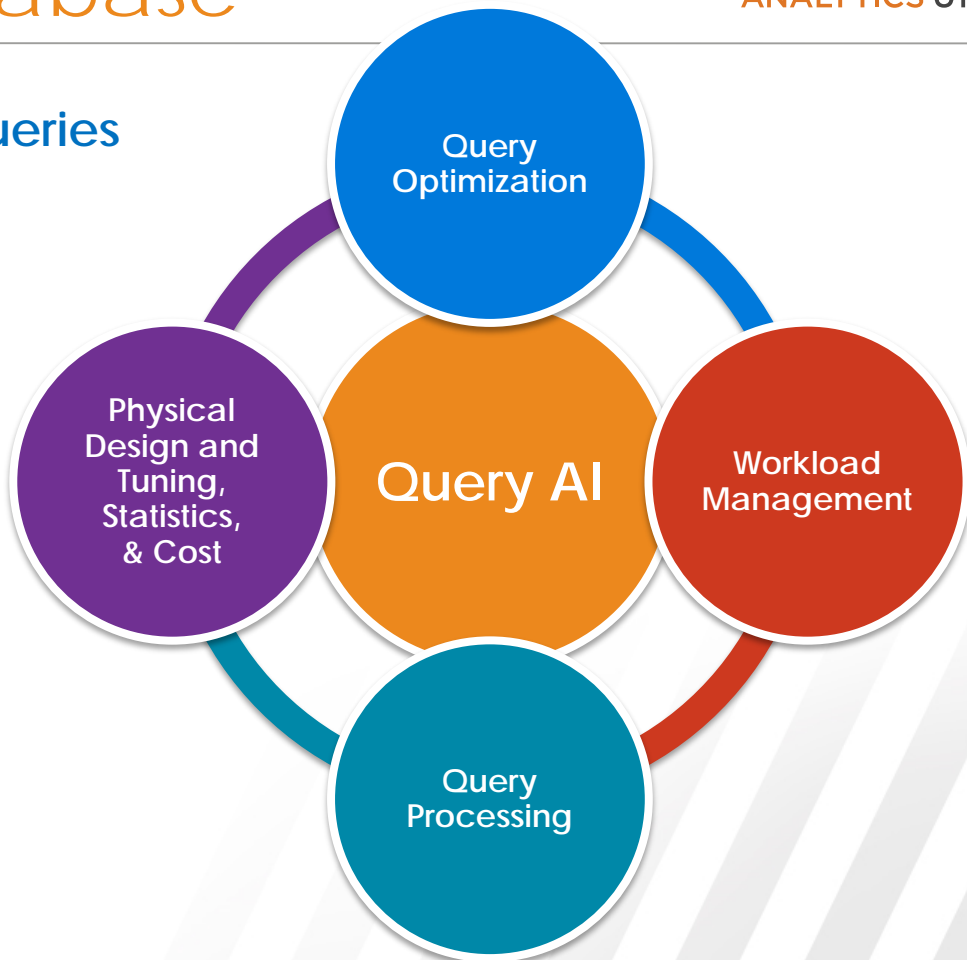


Controlling the Workload Mix using TASM in an Autonomous World

Douglas Brown, Technical Innovation Office (TIO)
Ruth Fenwick, Database Engineering

- **The future of Autonomics in TASM**
 - Artificial Intelligence
 - Machine Learning
- **Flex Throttle Review**
- **Arrival Rate Metering (ARM)**
- **Service Level Goals (SLG)**
- **Questions**

- Artificial Intelligence for Queries
 - Physical design and tuning
 - Statistics
 - Costing
 - Optimization
 - Workload Management
 - Query Processing
- Learning at the Edge
- Focused ML Models
- Automation
- Closed loop



AI Evolution



ARTIFICIAL INTELLIGENCE

Early theory and accomplishments stirs interest and excitement

1950s



MACHINE LEARNING

Predictive accuracy improves with more data

K-means clustering, Bayesian Networks, Support Vector Machines

mid 2000s



DEEP LEARNING

Allows more-complex problems to be tackled, and others to be solved with higher accuracy, with less cumbersome manual fine-tuning

Recurrent Neural Networks, Convolutional Deep Neural Networks, Generative Adversarial Network, Deep Feed Forward

mid 2010s

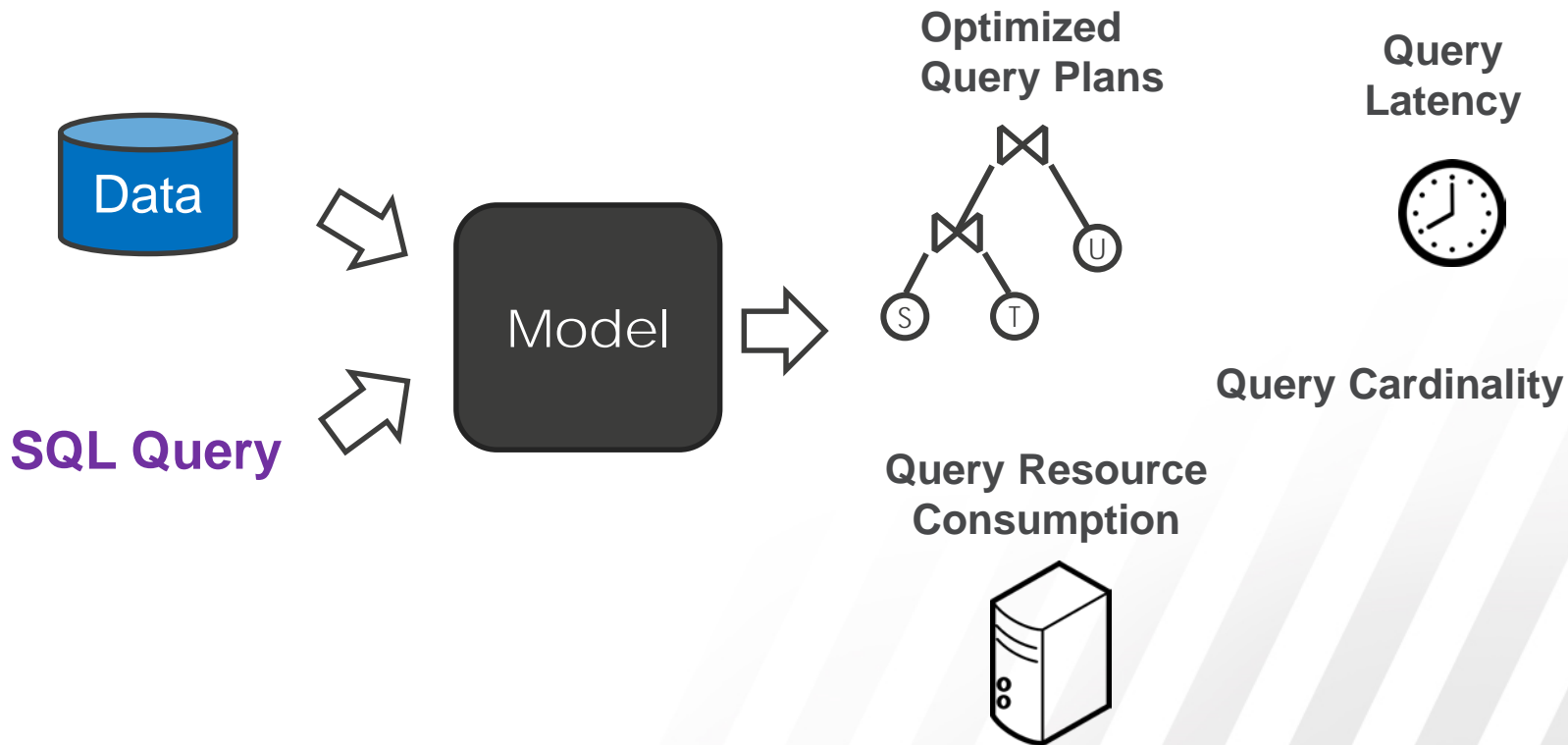
What makes Deep Learning so intriguing?

- Ability to create value with little or no domain knowledge required
- Ability to incorporate data from across multiple, seemingly unrelated sources
- Ability to tolerate very noisy data

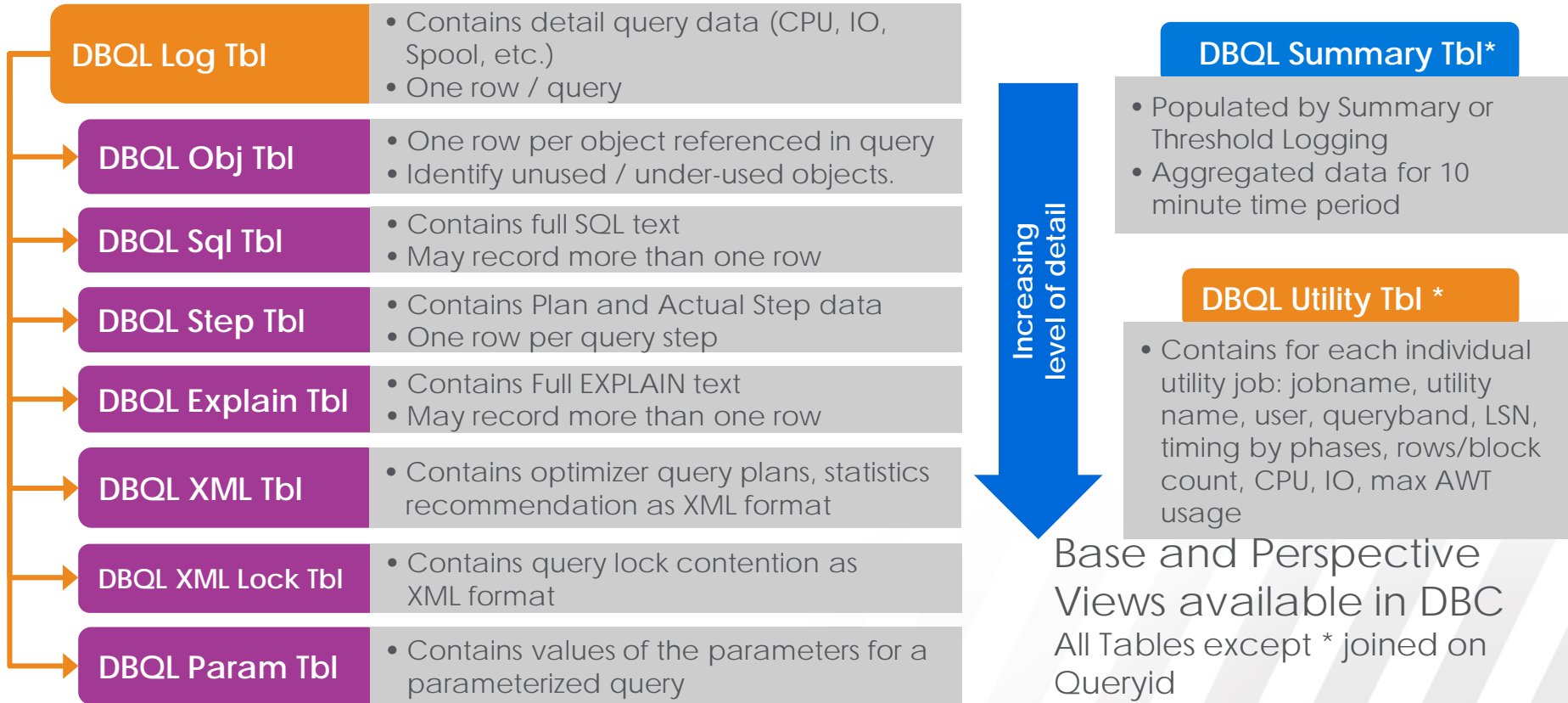


Our Vision

Automate **workload management** in the context of **machine learning** and **deep learning**

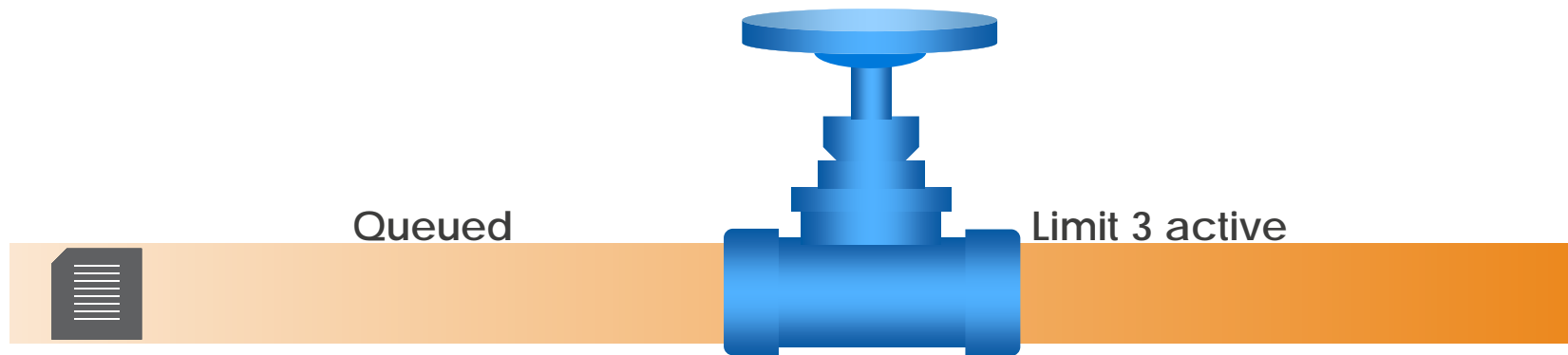


What is stored in Database for Learning



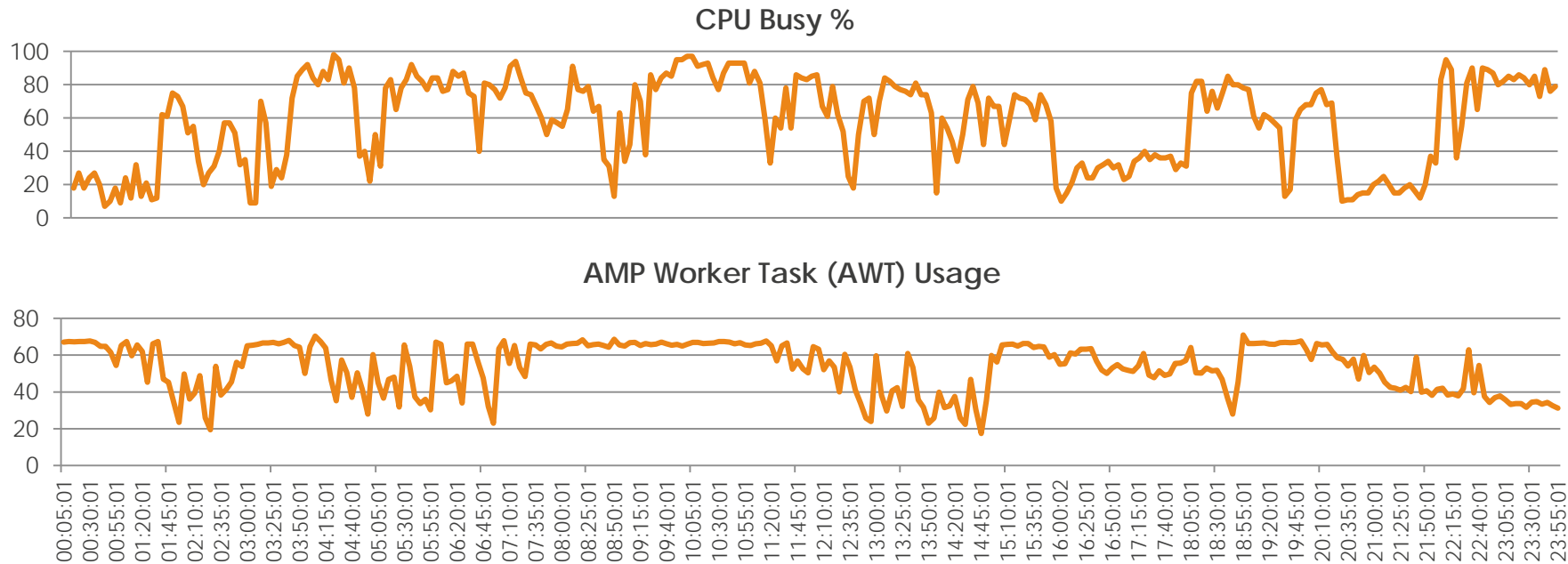
- **Why we need ML/AI models**
 - Throttle Analysis – How much should I throttle each request to guarantee optimal system performance?
 - Which workload(s) should I throttle?
 - Which type of throttle should I use: System, WD, Group, Flex, Arrival Rate Meters?
 - What is a recommended throttle setting?
 - Are my throttles under-performing?
 - SLG analysis – Will I meet my SLA?
 - How much budget does each query need?
 - Do I have enough capacity to meet my SLA?
 - How much does it cost to run my query?

What is a Throttle?



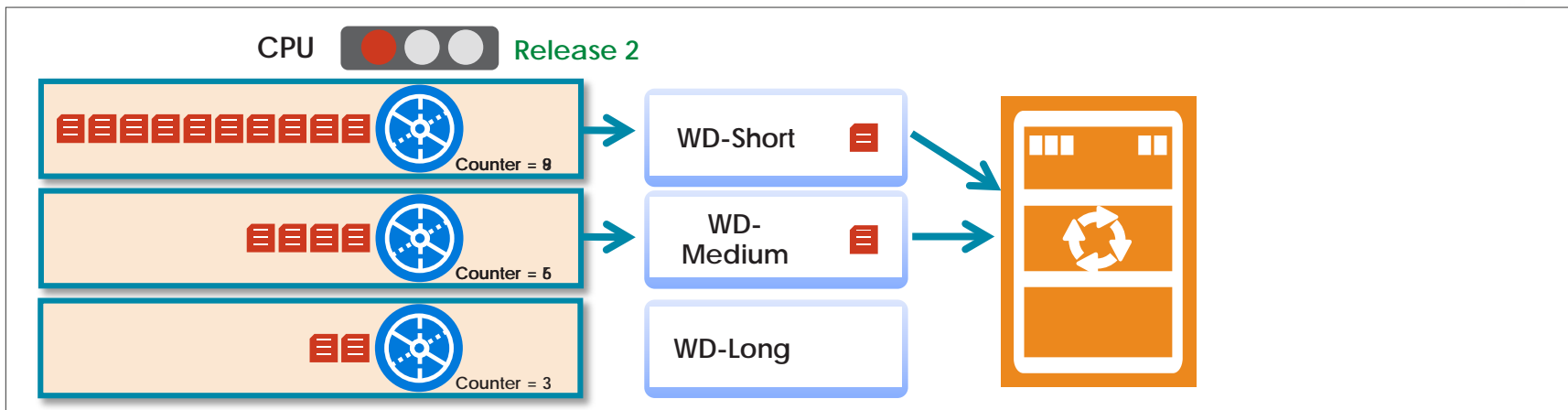
TASM Flex Throttles

What Problem Do They Solve?



Flex Throttle Example

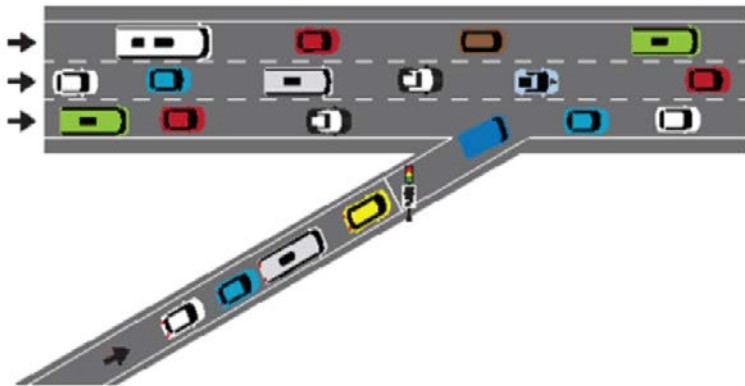
Workload Throttles	Flex ON/OFF	Workload Priority	Throttle Limit	Throttle Counter
WD-Short	ON	High	8	9
WD-Medium	ON	Medium	5	6
WD-Long	OFF	Low	3	3



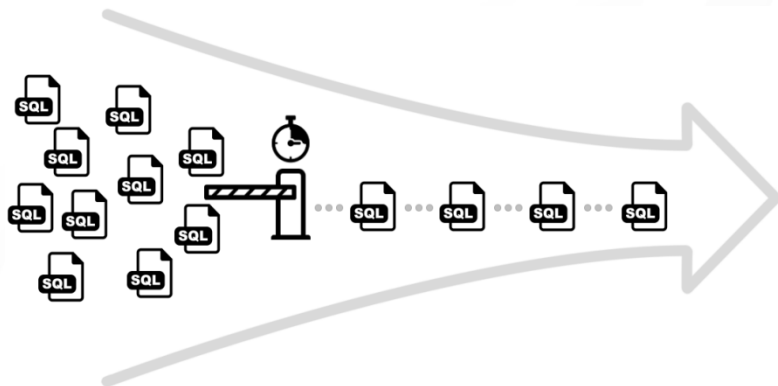
Arrival Rate Meter

- New TASM rule type
- Regulate the flow of SQL requests being admitted by TASM
- Specify maximum rate during a specific time unit.
- Example: 3 queries per hour for queries with estimated processing time of more than 1 hour.

Freeway with Ramp Metering



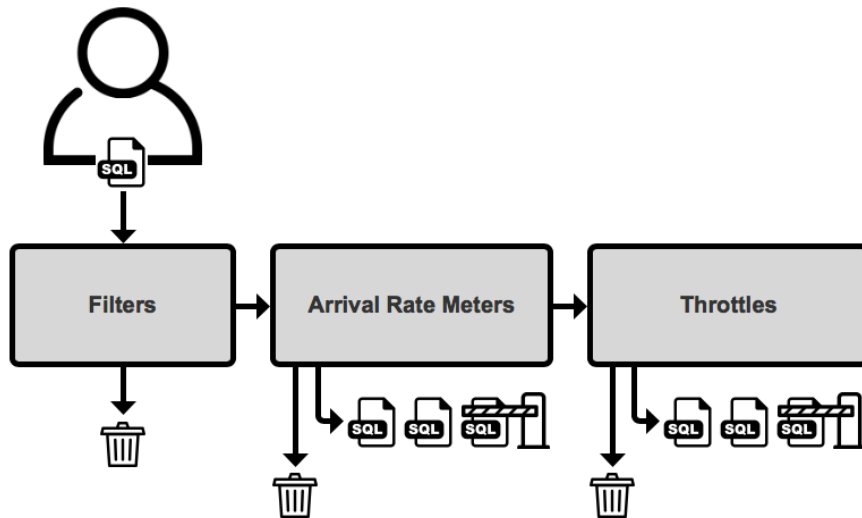
TASM Arrival Rate Meter



Arrival Rate Meter - Processing

Order of processing:

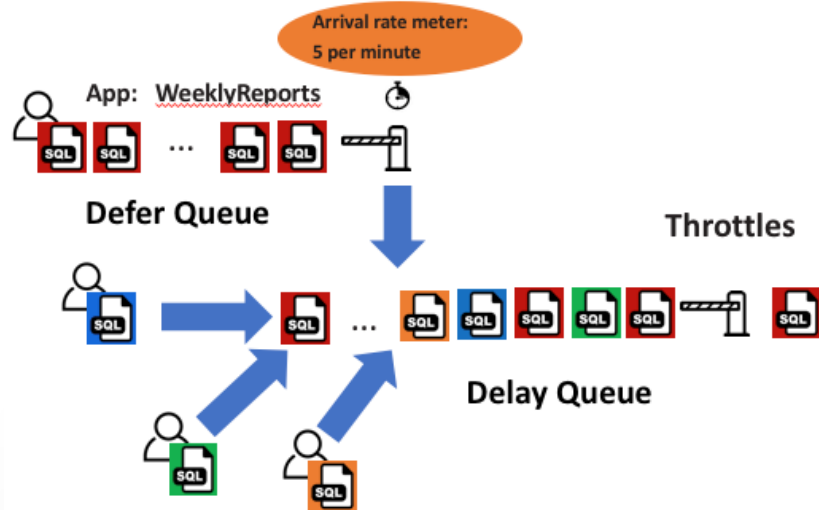
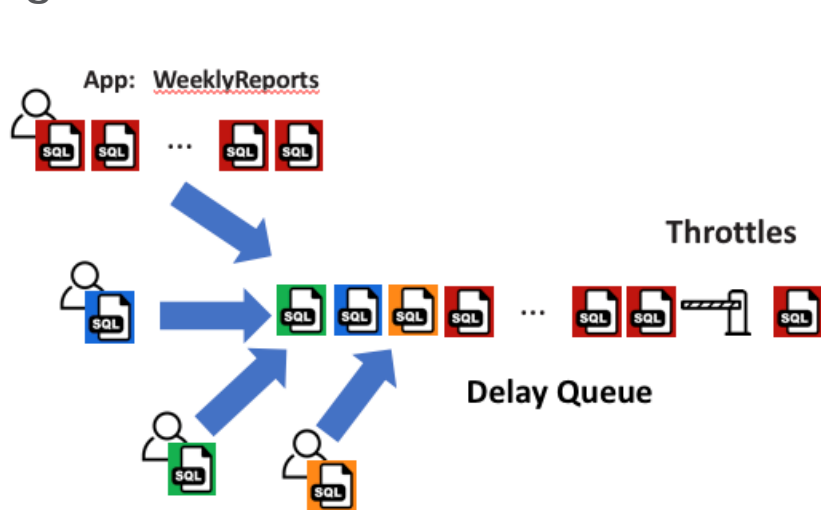
1. Filters
2. Arrival Rate Meters (Defer Queue)
3. Throttles (Delay Queue)



Arrival Rate Meter – Use case 1

Problem: Flood of requests from specific application monopolizing the front of the throttle delay queue.

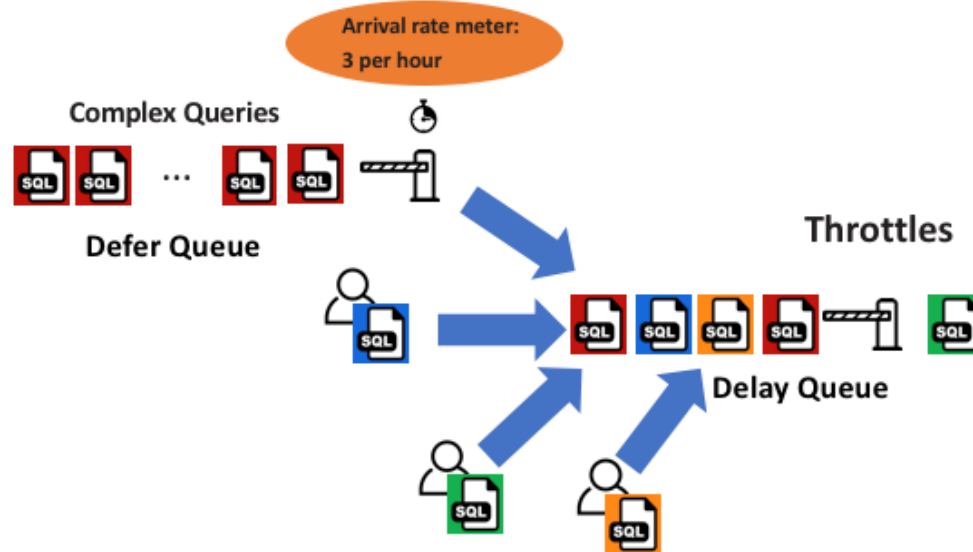
Solution: Create an arrival rate meter for this application so that other requests can get in between each time unit.



Arrival Rate Meter – Use case 2

Problem: Complex queries use too much resources.

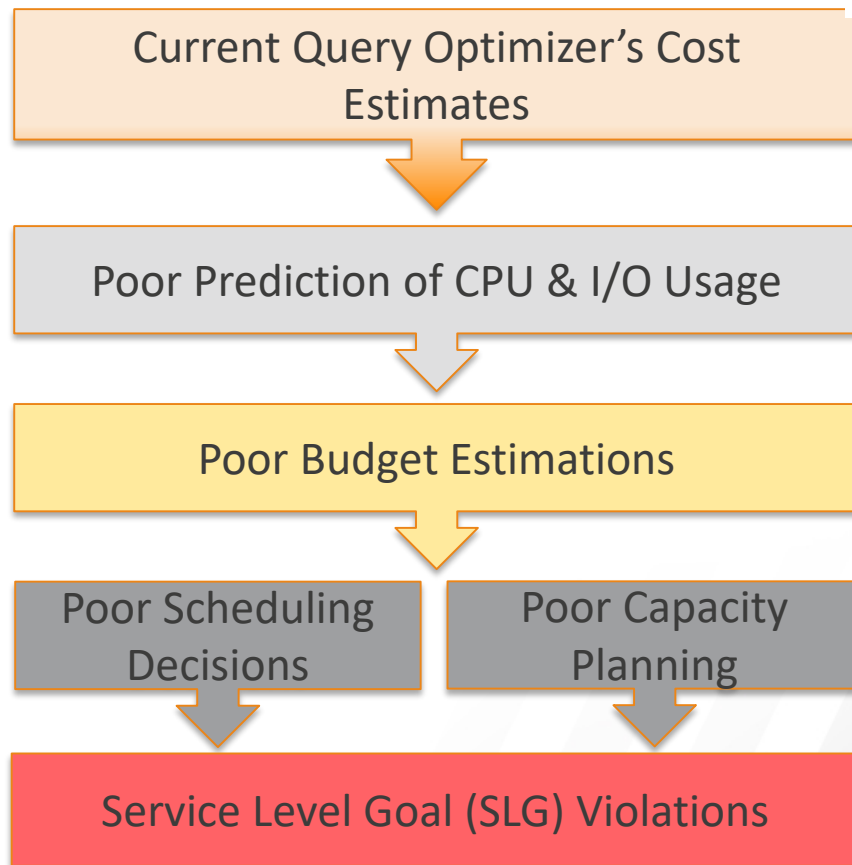
Solution: Create an arrival rate meter for complex queries to limit resource usage.



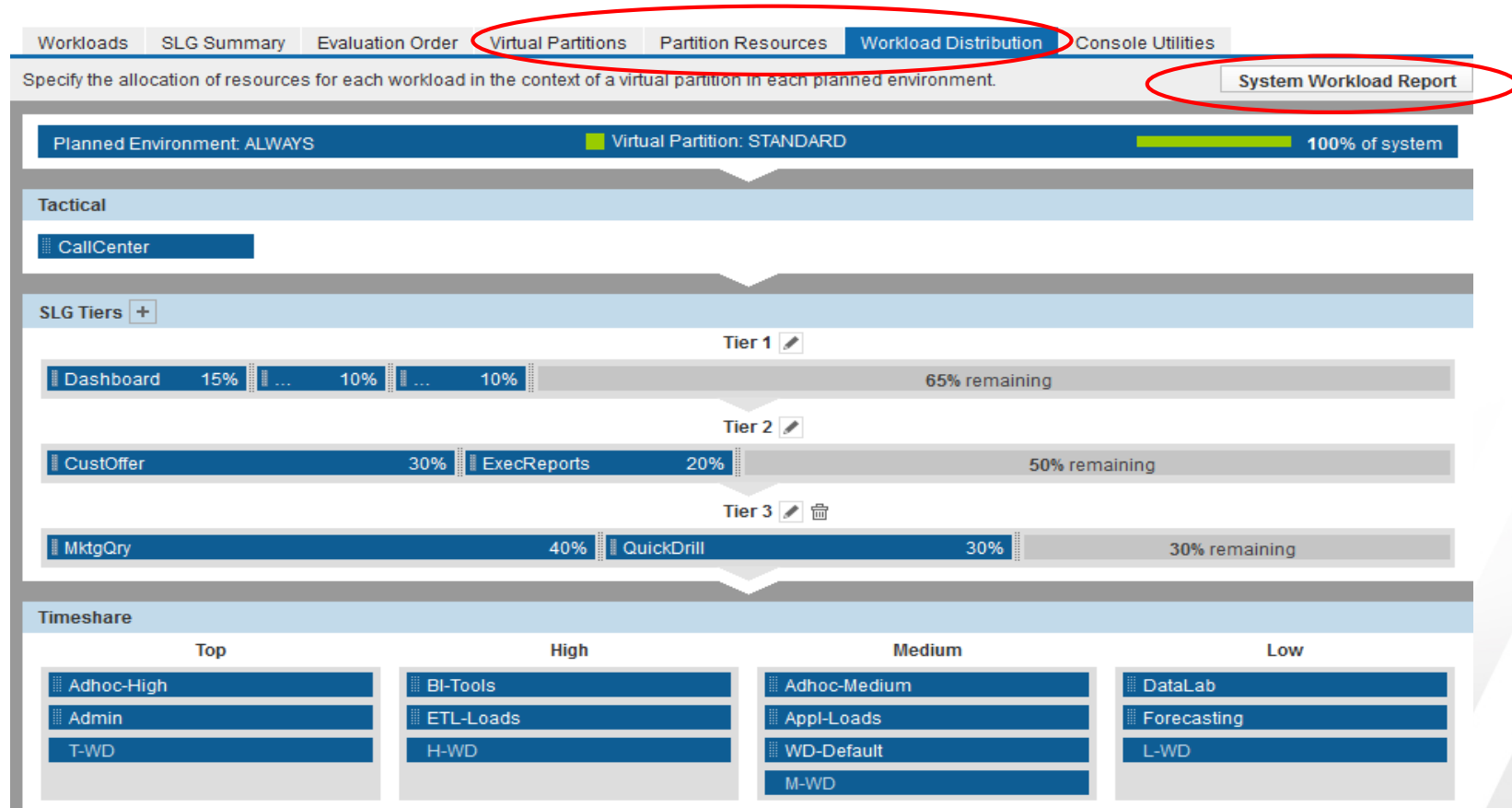
Service Level Goals (SLGs)

- Purpose:
 - Increase consistency and predictability
 - Make it easier to keep important jobs higher priority
 - Easier configuration
- Replace the current methodology of having to create Tier WDs and assign specific shares.
- User defines response time goal and DBS executes query to meet that goal.
- TASM determines the rate needed to execute to meet SLG.

Inaccurate Predictions Bigger Problems!



SLG Tiers



SLG Tier Reporting

System Workload Report





































View workload resource allocations across all virtual partitions. Translation of workload allocations to system resource percentages does not consider workload hard limits. To sort on a second column, hold the Control/Command key.

Planned Environment: **Always**

Virtual Partitions

Standard 90% Ne... 10%

Workloads

WORKLOAD		VIRTUAL PARTI. ▲	TIER ▲	% OF TIER	% OF SYSTEM	% OF SYSTEM
slg1		NewVP	Tier 1	13.6	1.4	
slg2		NewVP	Tier 2	26	2.2	
ts-top		NewVP	TS1 Timeshare Top		2.2	
ts-top2		NewVP	TS1 Timeshare Top		2.2	
ts-high		NewVP	TS2 Timeshare High		1.1	
ts-med		NewVP	TS3 Timeshare Medium		0.6	
ts-low		NewVP	TS4 Timeshare Low		0.3	
webApp1		Standard	Tier 1	10	9	
webApp2		Standard	Tier 1	10	9	
dashboard		Standard	Tier 1	15	13.5	
CustOffer		Standard	Tier 2	30	17.6	
ExecRpts		Standard	Tier 2	20	11.7	
MktQry		Standard	Tier 3	40	11.7	
QuickDrill		Standard	Tier 3	30	8.8	
SLGapp1		Standard	Tier 4	5	0.4	
SLG5app1		Standard	Tier 5	5	0.4	
T-WD		Standard	TS1 Timeshare Top		2.9	
wd1		Standard	TS2 Timeshare High		1.4	

23 rows total

Close

Profiling Service Level Goals

Response Time Distribution

- Needs DBQL Data
- Shows the distribution of queries by time

ResUsage

- Needs ResUsage Data
- Shows CPU, IO, AWT

SLG

- Need TASM Data with SLGs defined
- Shows # of queries met/not met SLG

Deployed as Jupyter Python Notebook @ [hub](#)

Response Time Distribution

“Response Time Distribution” View

- For a customer with no SLA, this helps in setting up SLG and Service Level Percentage (SLP)
- Each TASM planned environment has a different graph.
- Lets you know how RT behaved for same WD in different environments.

What does this view have?

Input: Time Window, WD

Output: % of queries in each bin of response time

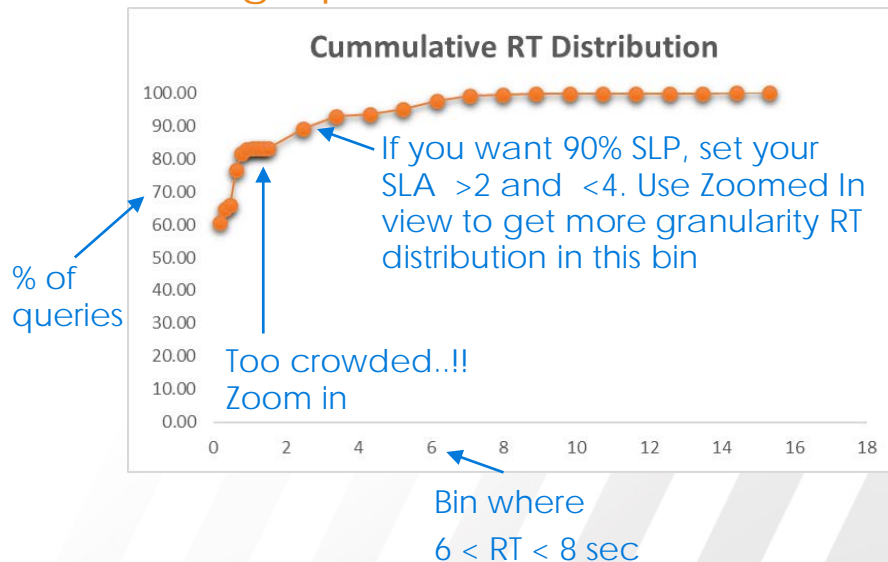
X-Axis: Response time

Y-Axis: % queries in that bin.

Additional Information in table:

- How many queries in that bin met SLA.
- how many % met with exception and without exception.
- how many WD each tier had and how many had SLA defined.

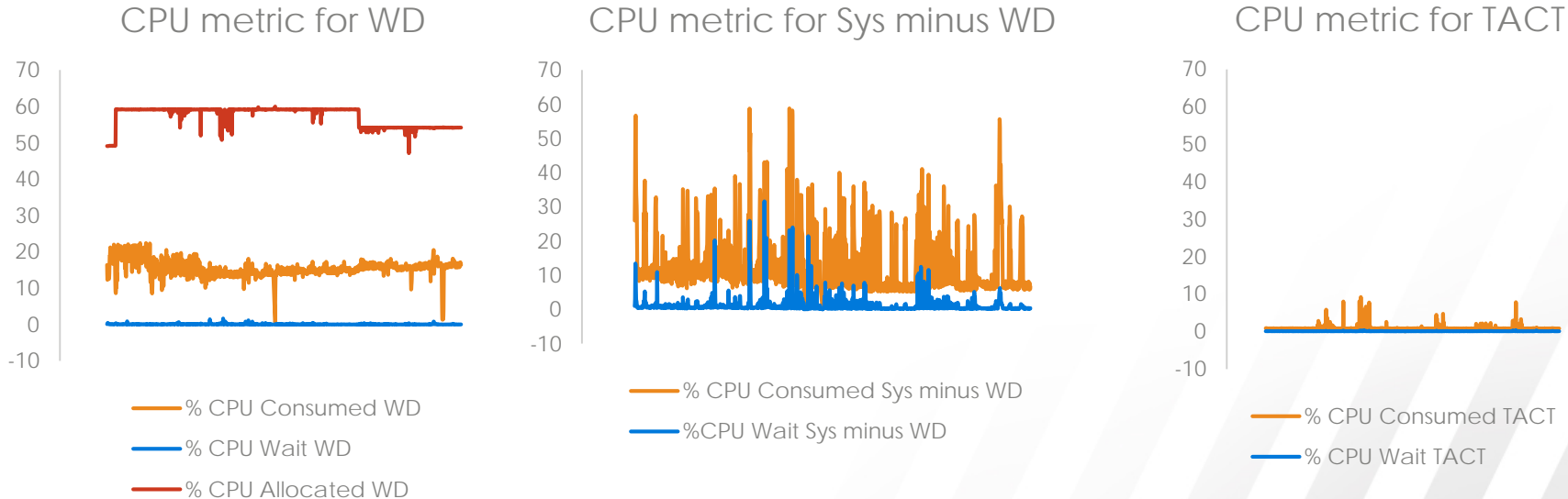
How graph looks like?



ResUsage Metric Graph

Two Reasons:

- SLA behavior
- Resources Consumed
- What resources correspond to the SLG behavior reported
- Helps in analysis of queries if missed SLG - what was the load on system resources.



ResUsage metric graph

X-Axis: Time

Y-Axis: Metric value(each graph has different Y axis)

SLG by Tiers

Why "Tier Wise" view?

- Top level view
- Simplified aggregated information for all tiers
- User can identify which Tier is bad
- Not to overwhelm user with all existing WDs
- Drill down more with other views to identify bad WD

What does this view have?

Input: Time Window

Output: Per Tier %queries met/not met SLA

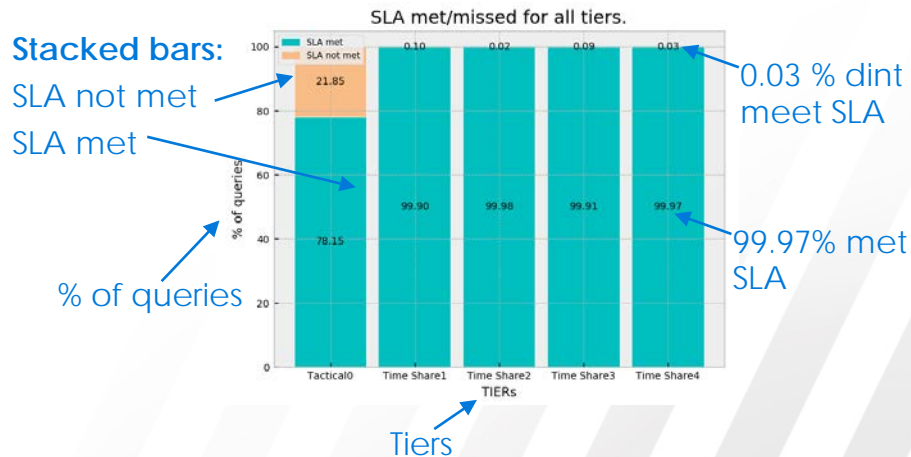
X-Axis: Tiers (not time)

Y-Axis: % queries met/not met SLA.

Additional Information in table:

- how many % met with exception and without exception.
- how many WD each tier had and how many had SLA defined.

Graph by Tiers



SLG by WD

“Specific Workload” view?

- Lowest level view
- To identify bad environment/settings for the WD
- Give two sides of information:
 - SLG met/not met by Planned Environment
 - ResUsage Information as discussed above

What does this view have?

Input: Time Window, WD

Output: Per WD setting, %queries met/not met SLA

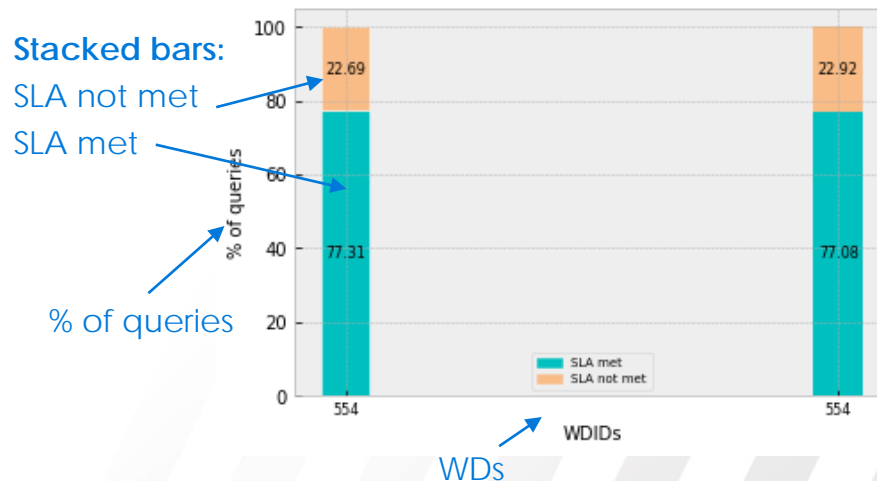
X-Axis: set of WD settings (not time)

Y-Axis: % queries met/not met SLA.

Additional Information in table:

- Setting information like planned env name, SLA, CPU rel share etc.

How graph looks like?



Thank You!

Rate This Session #

with the Teradata Analytics
Universe Mobile App

*1120 Controlling the Workload
Mix using TASM in an
Autonomous World*

Follow Me

Twitter @ RFenwick7

Questions/Comments

Email: **Doug.Brown@Teradata.com**
Ruth.Fenwick@Teradata.com

The background is a vibrant cosmic scene with a purple and orange nebula. A large, solid black arrow points upwards from the bottom center towards the top. A semi-transparent dark grey horizontal band is positioned across the middle of the image, containing the text 'Bonus Material' in white.

Bonus Material

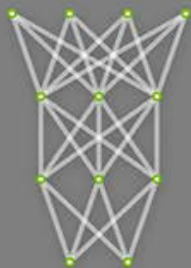
T E R A D A T A A N A L Y T I C S U N I V E R S E 2 0 1 8

DEEP LEARNING

TRAINING

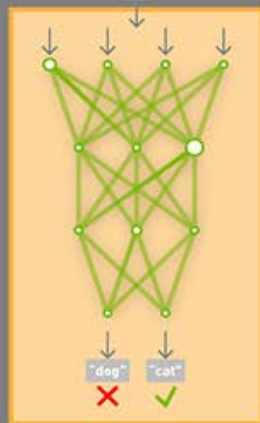
Learning a new capability
from existing data

Untrained
Neural Network
Model



Deep Learning
Framework

TRAINING
DATASET



Trained Model
New Capability



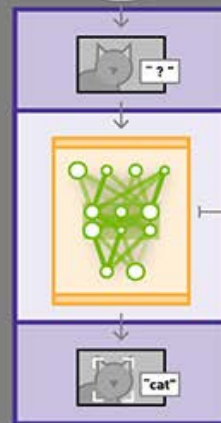
INFERENCE

Applying this capability
to new data

NEW
DATA



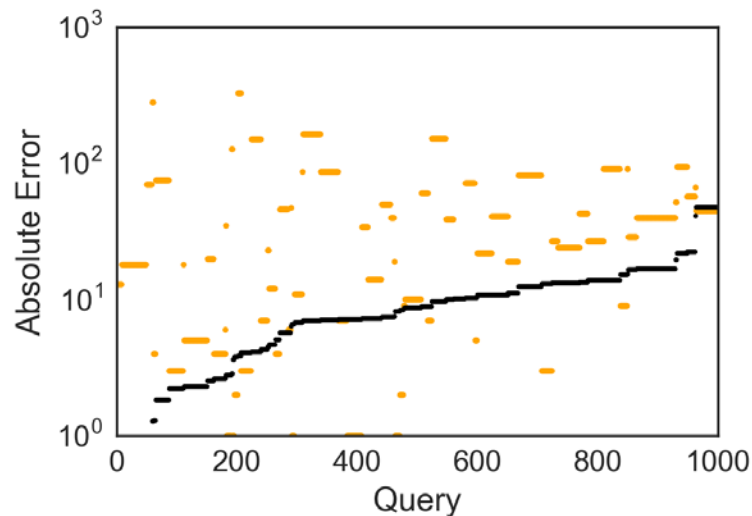
App or Service
Featuring Capability



Trained Model
Optimized for
Performance

Experimental Results

Deep Learning Model vs Teradata



Deep Learning Model vs Teradata

