

A large teal abstract graphic on the left side of the slide, featuring a curved shape that resembles a stylized globe or a cloud. The text 'teradata.' is overlaid on this graphic in white.

teradata.

Teradata Database

Vantage

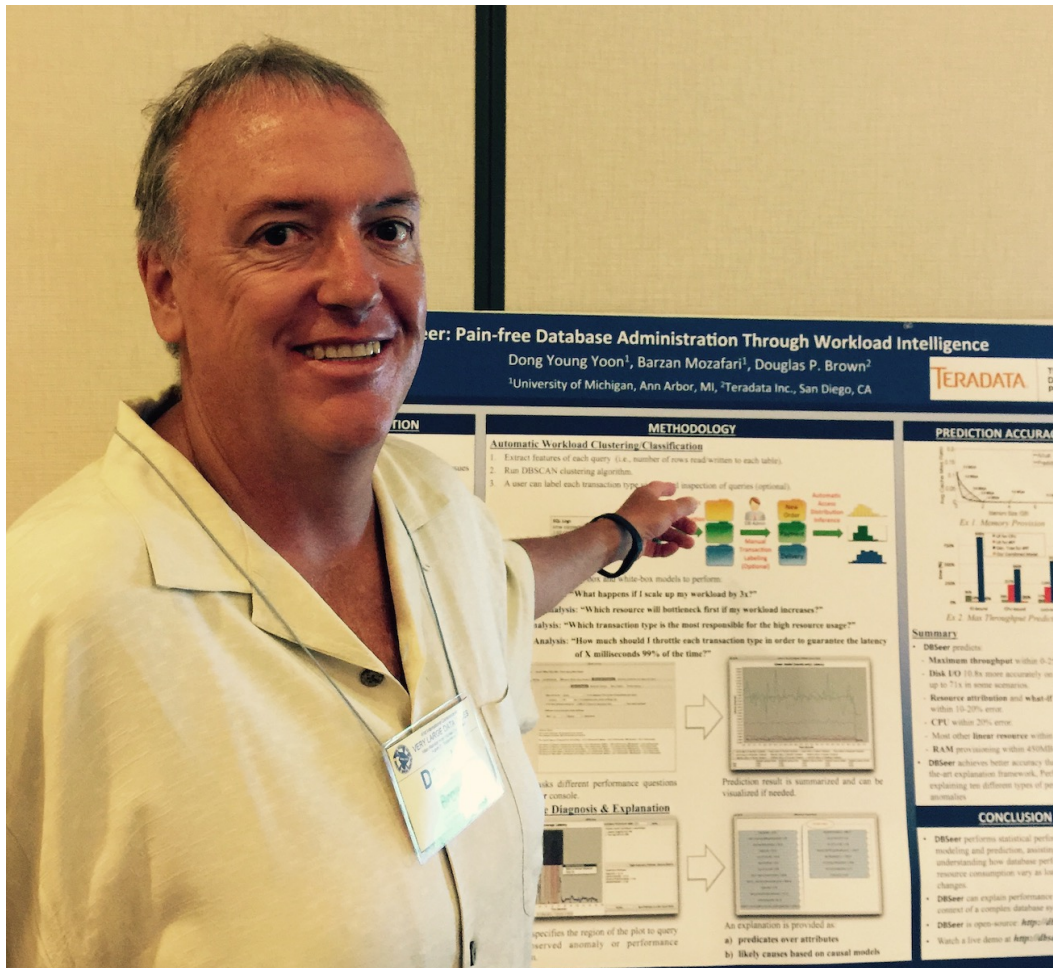
Douglas Brown

May 12, 2022

Agenda

- Who Am I - Background
- Teradata Database High Level Overview
 - Key Components in the Teradata Architecture
- Teradata MPP architecture
- Data Management
- Query Execution
- Optimizer
- Workload Management
- Technology Research projects

Doug Brown Engineering Fellow - Technology Innovation Office (TIO)



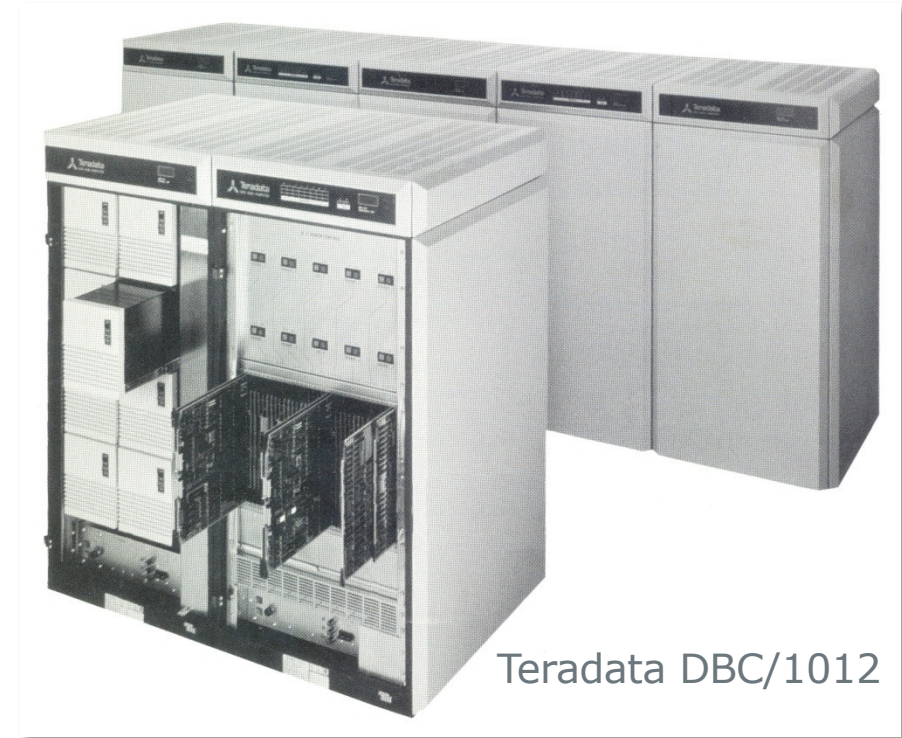
Doug is currently a member of the Technology Innovation Office (TIO) in San Diego. He is presently the lead researcher and Product Owner (PO) for the workload management components of the Teradata database aka Teradata Active Systems Management (TASM).

Doug has been working in the Teradata database engineering team since 1989 and has contributed over 75 US Patents.

1979

Birth of Teradata

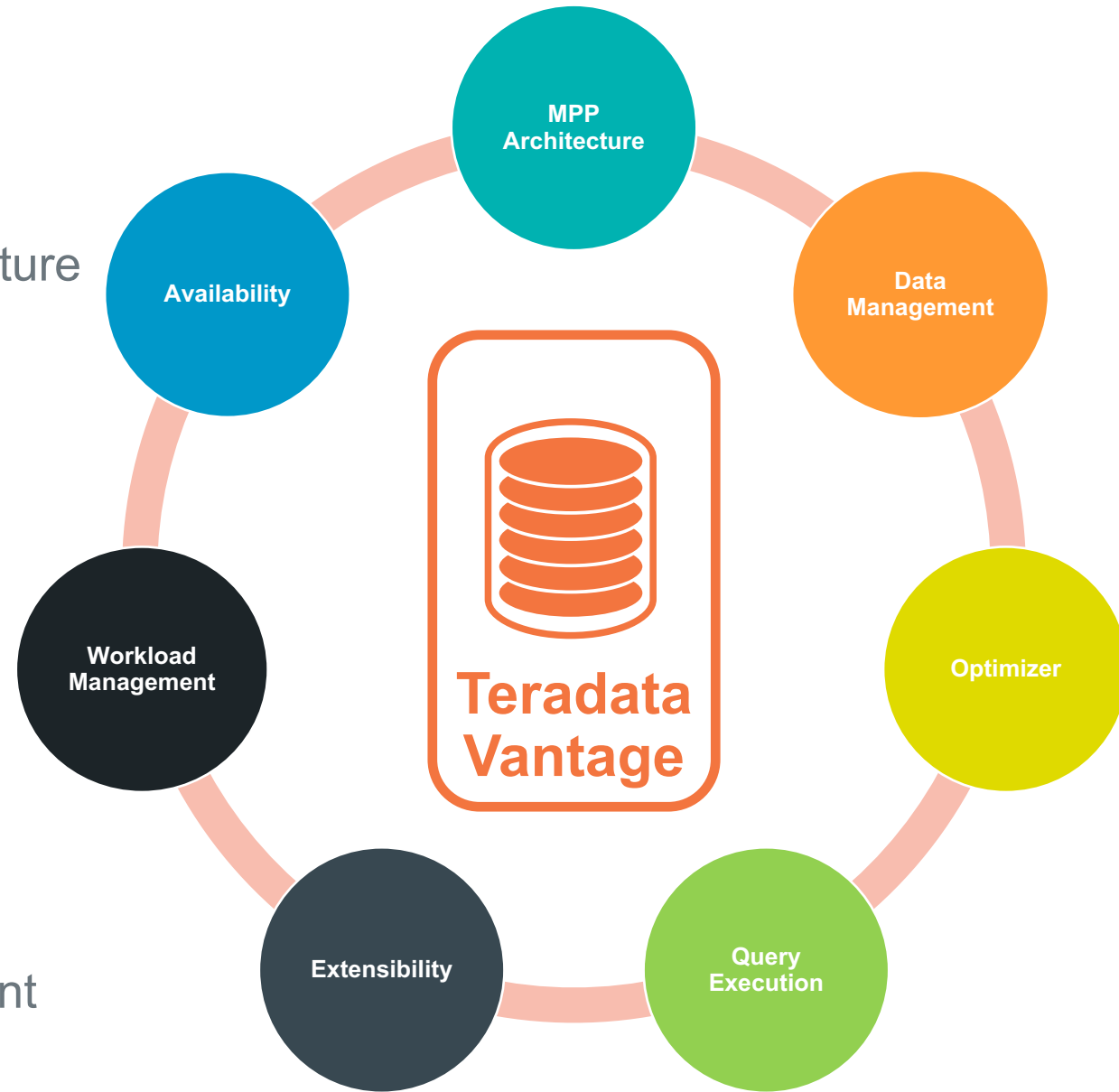
- Microprocessor - 8086
- Personal computer - Commodore 64, Apple I, TRS80
- Disk drive - 200MB weighs 30lbs and requires a washing machine sized cabinet
- Enterprise computing - Blue or Bunch mainframe
- Client/Server computing - didn't exist
- Internet and e-mail - Only if you are a university researcher or DoD



Teradata was founded to solve enterprise Terabyte problems utilizing microprocessors and commodity disks.

Key Components/Concepts

- MPP Shared-nothing Software-based Architecture
 - Hash-based with partitioning (row and column)
- Data Management
 - Relational, semi-structured, unstructured
- Cost-based Optimizer
- Query Execution
 - SQL Engine, Machine Learning Engines
- Extensibility
 - UDTs, UDFs, procedures, foreign servers
- High-concurrency mixed workload management
- Reliability and availability



Teradata Data Warehouse Technology Is Unique

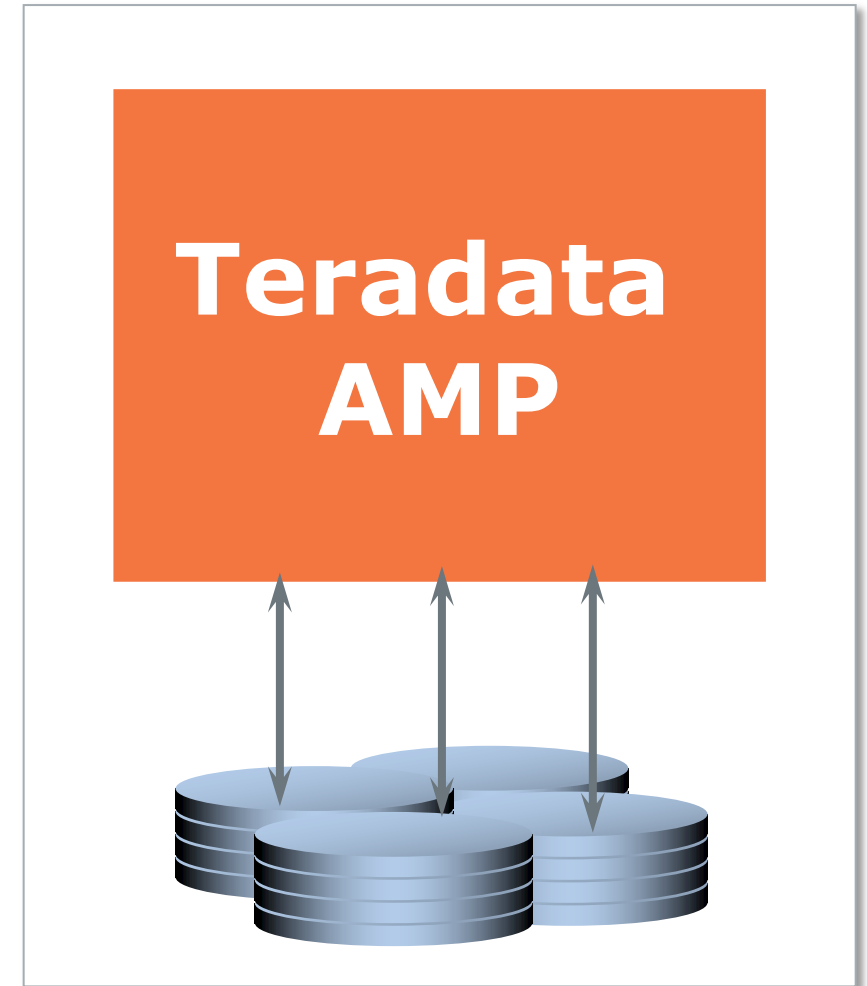


- Parallelism built-in from the ground up
- Dedicated to automatic management and operation
- Easy and fast data movement in and out
- Committed to the highest levels of reliability and availability
- High concurrency mixed workload management
- Parallel extensibility for high performance dynamic functionality
- ***Unequaled scalability in every dimension***



Teradata Access Module Processor (AMP)

- A single instance of Teradata has many database execution engine virtual processors (vprocs) called Access Module Processors (AMP)
- Each owns a portion of the data (1/# of AMPs)
- Performs all operations on data (select, insert, update, delete, index, join, aggregate,...)
- Each AMP operates independently on steps sent from Dispatcher



Teradata Parsing Engine (PE)

- **SQL Parser**

- Analyzes SQL statement for proper syntax
- Resolves object names
- Validates that the user has the appropriate access rights to the objects

- **Optimizer**

- Uses hardware configuration and statistics to generate a step by step execution plan
- Calculates cost for all reasonable choices
- Automatically rewrites the query if the cost can be improved

- **Query Step Dispatcher**

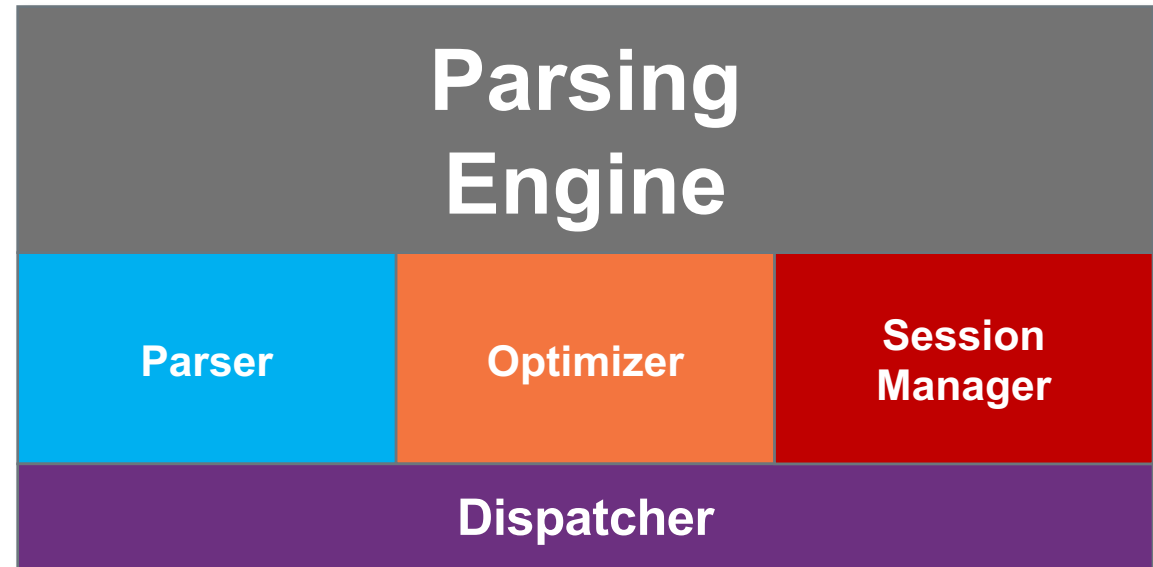
- Sends the steps in the execution plan to a single AMP, a group of AMPs, or all AMPs

- **Session Manager**

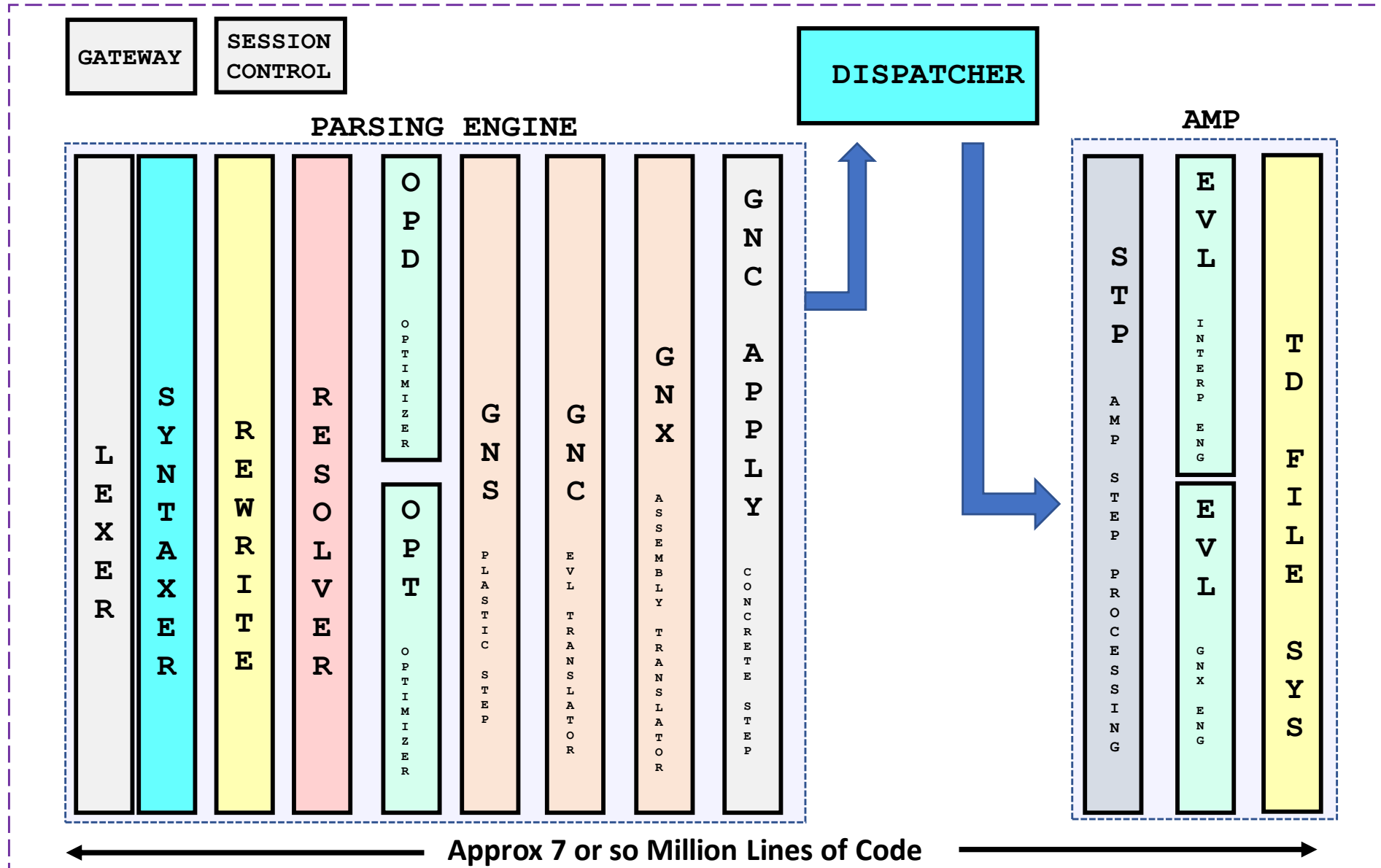
- Manages each session and communication to the client

- **Input Data Conversion**

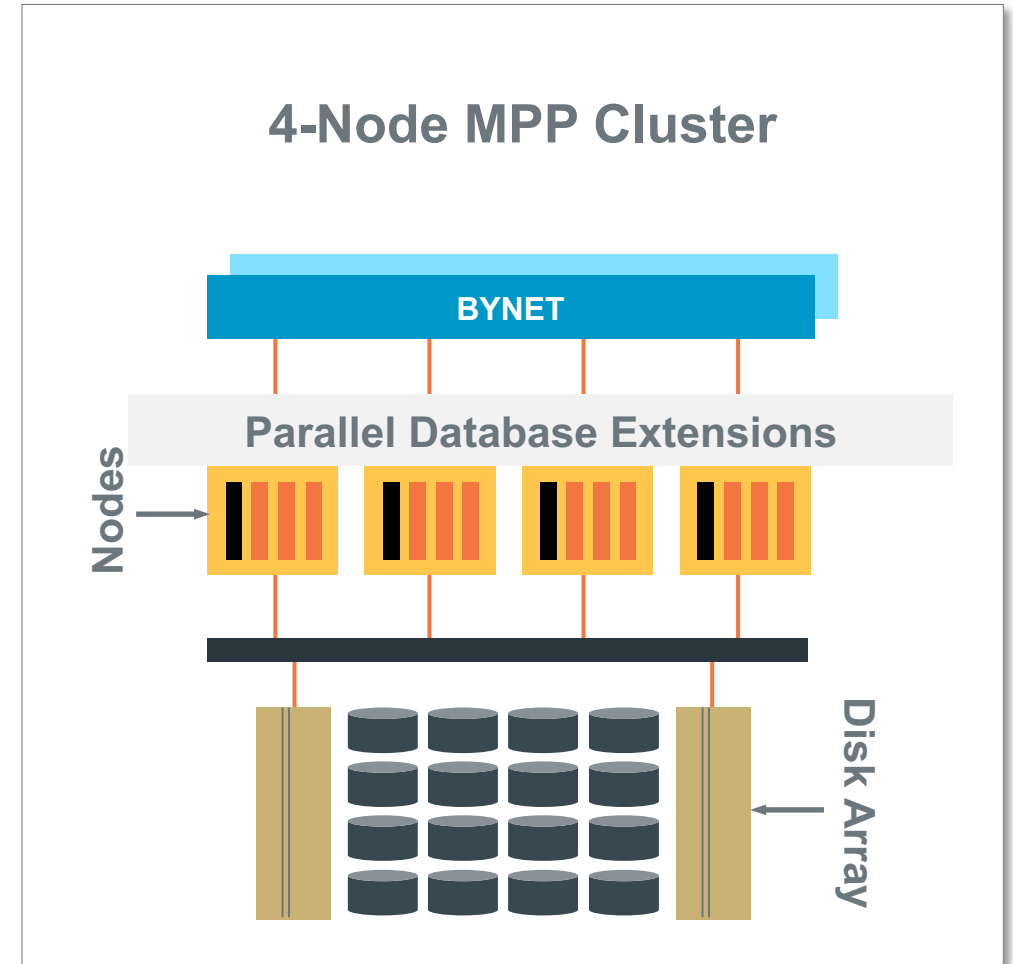
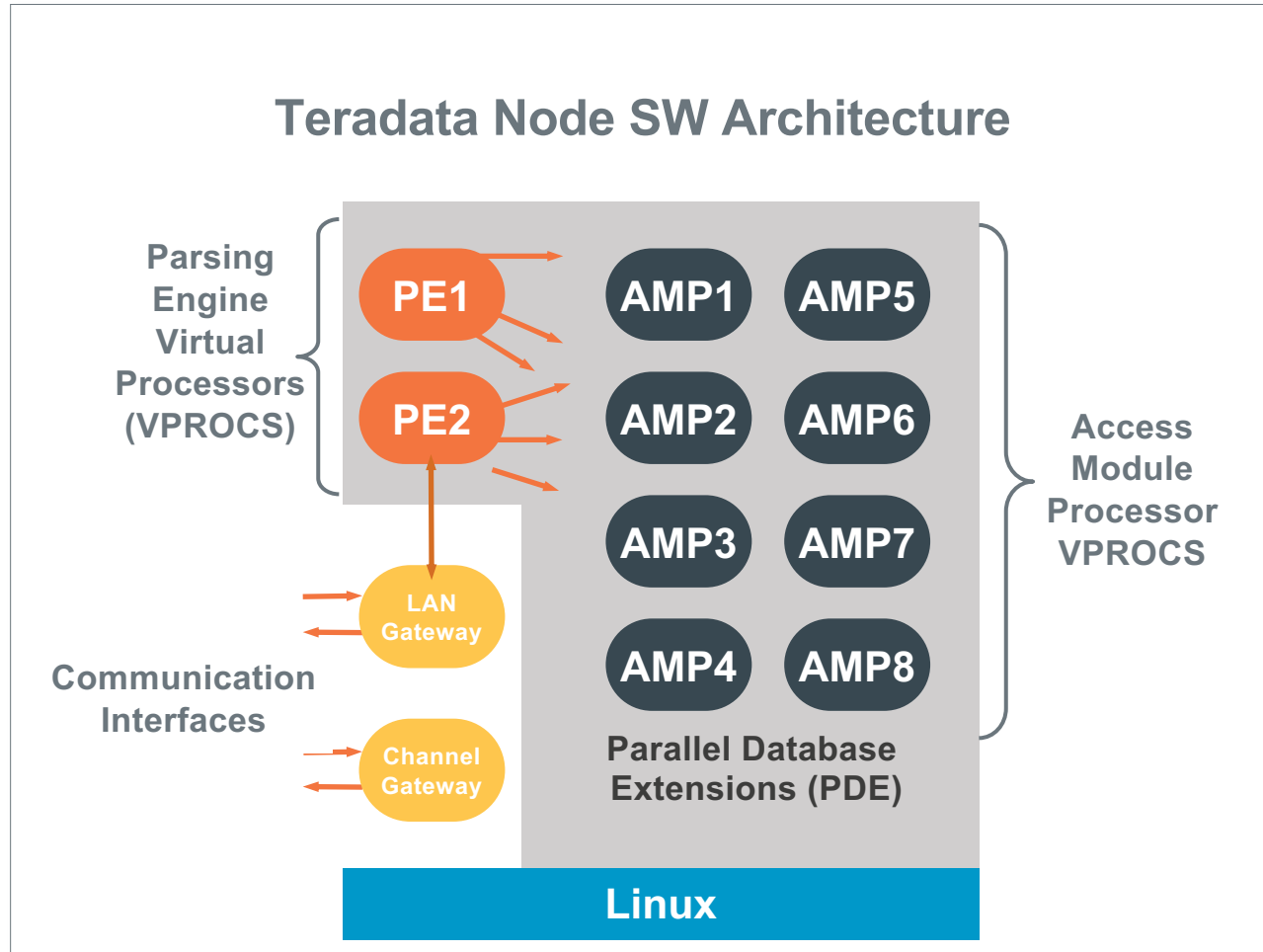
- EBCDIC to ASCII conversion



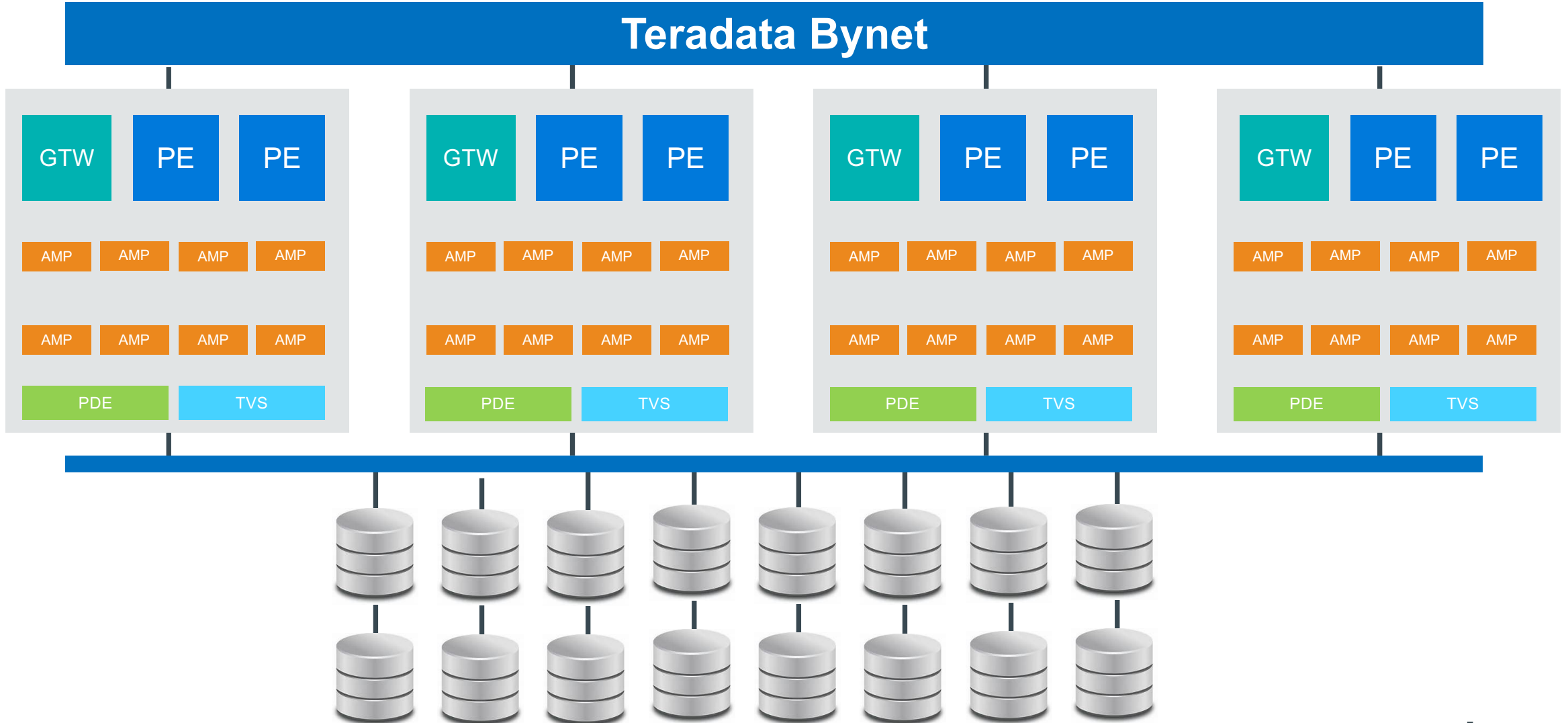
TERADATA "STACK"



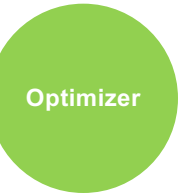
Parallel Database Extensions (PDE) Ties It All Together



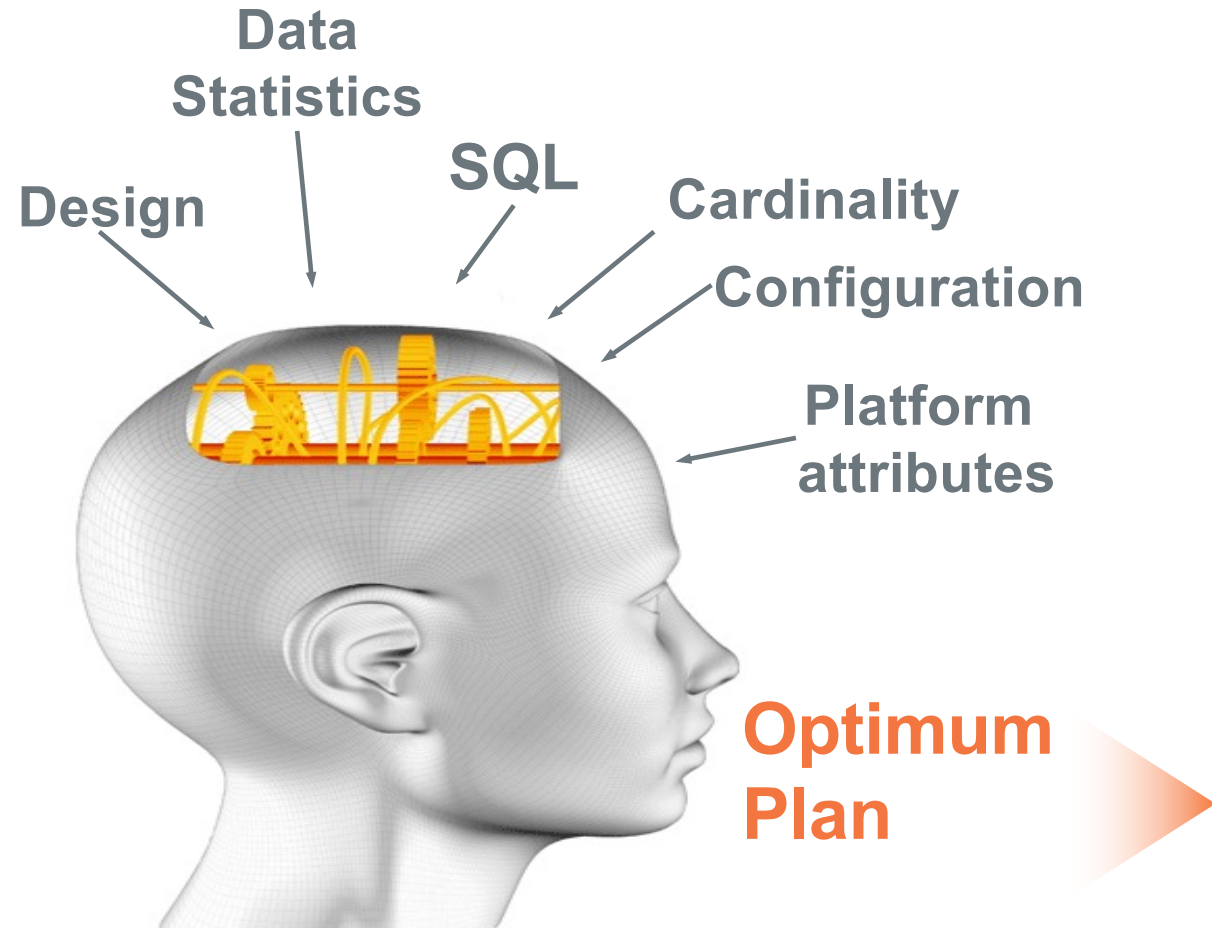
Unlimited Parallelism



Teradata Optimizer

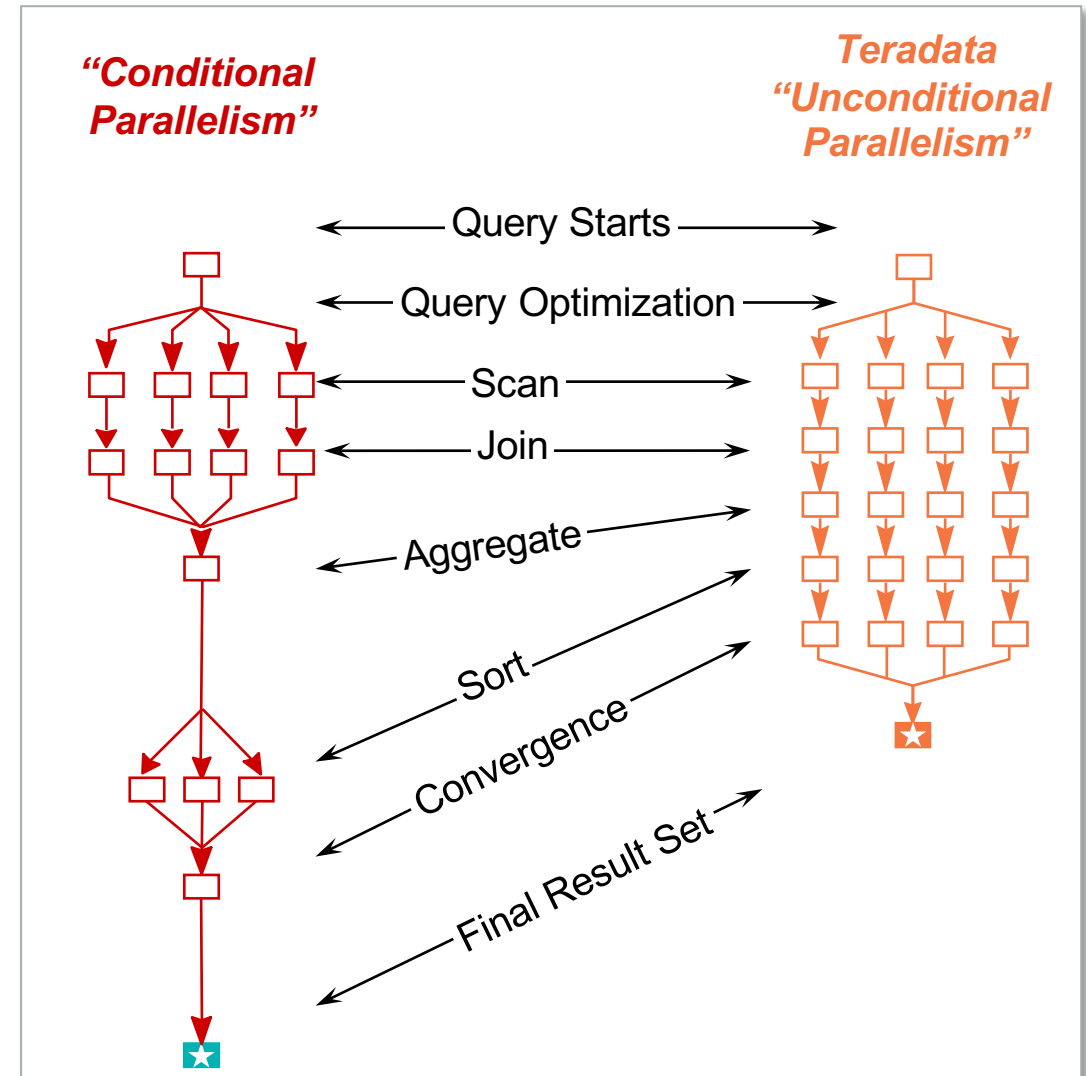


- Cost based query planning
 - Balance costs of
 - Cardinalities
 - Sorts and joins
 - Redistribution
 - Disk I/Os
 - Networks and nodes
 - Compression
 - Lowest resource costs
 - Fastest elapsed time
- Automatic rewrites
- Automatic optimization
 - No hints
 - No degrees of parallelism



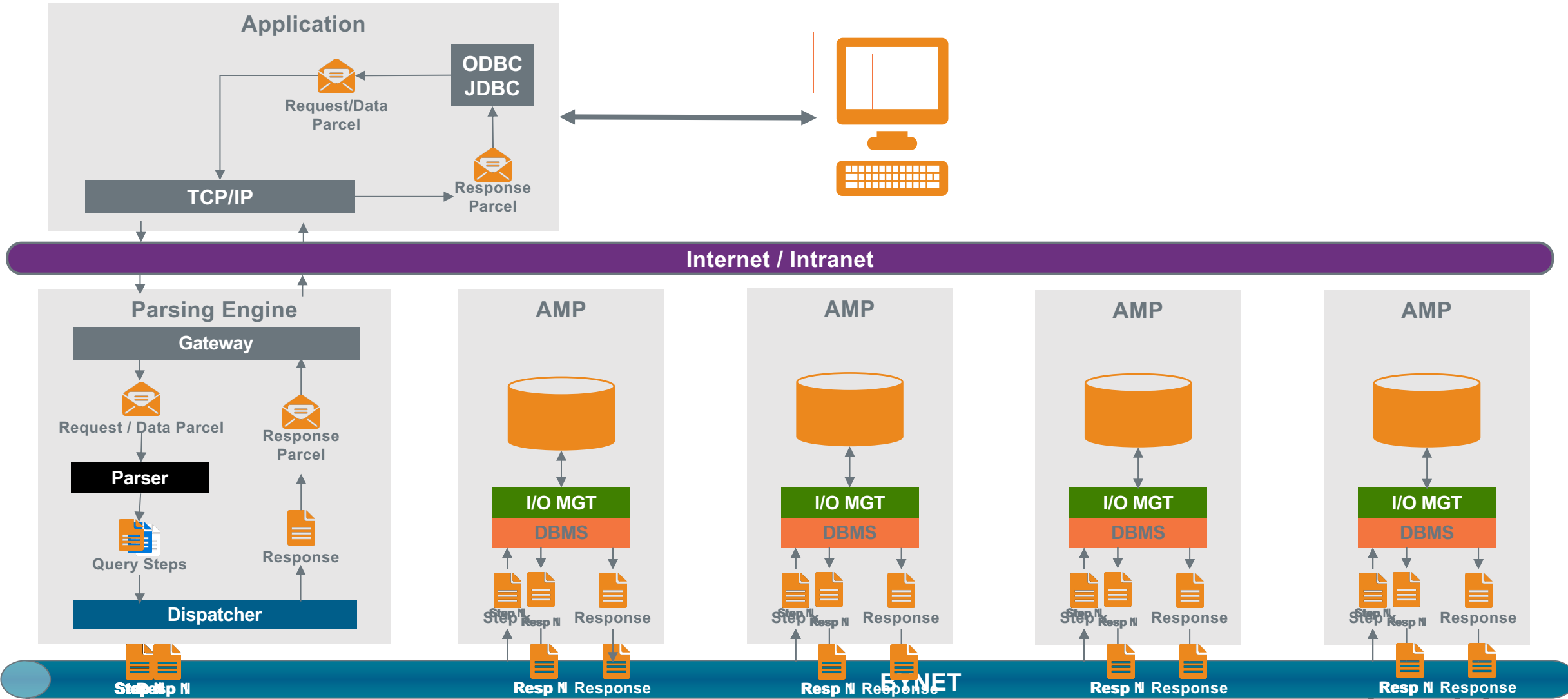
Query Parallelization

- Query parsing, management is fully distributed across the nodes
 - No head node/coordinator node
- All operations fully parallel
 - No single threaded operations
 - Scans, Joins, Index access, Aggregation, Sort, Insert, Update, Delete
 - Ordered Analytics
 - Extensibility functions
 - Result return



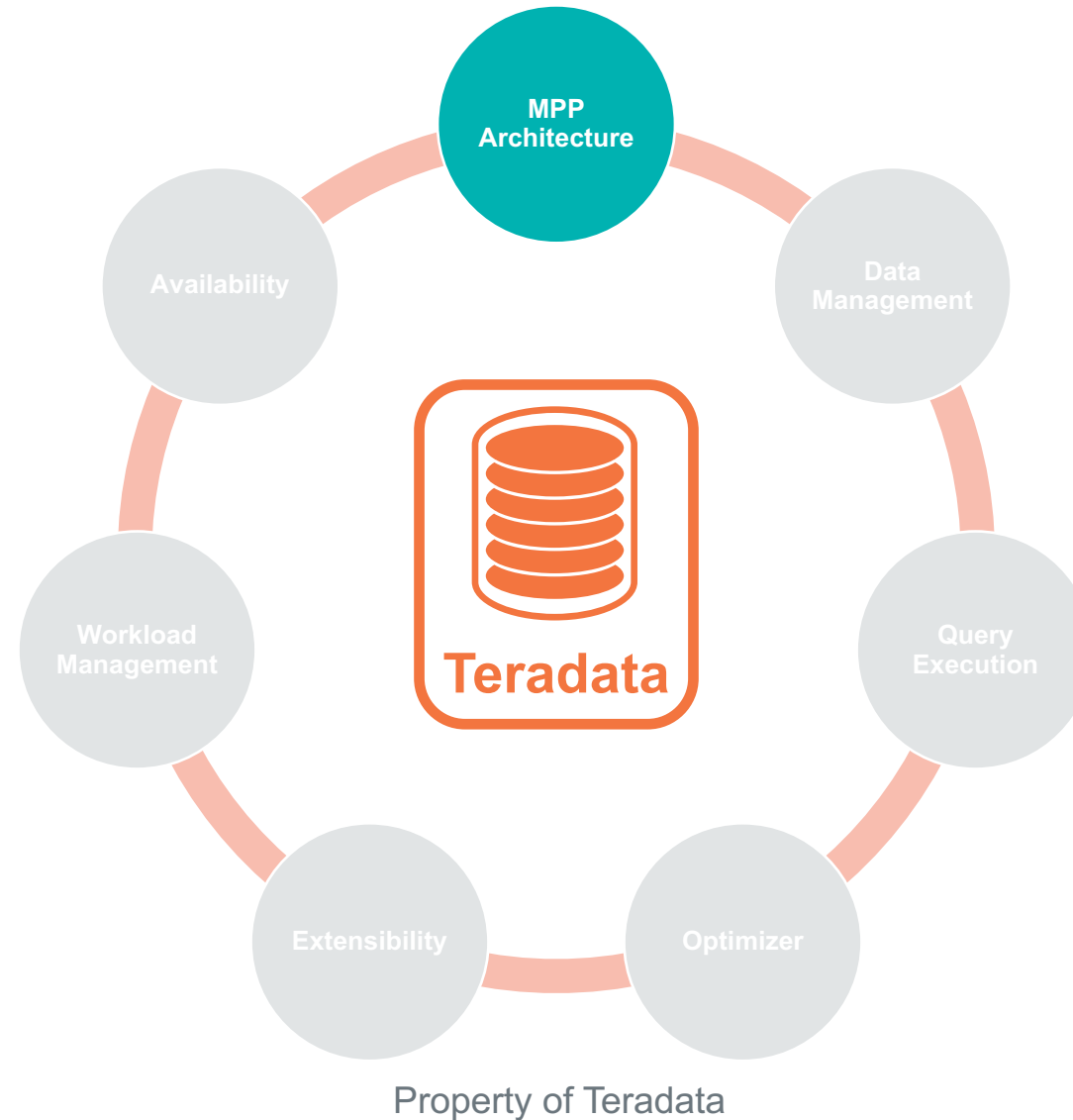
Query Execution

The Life of a Teradata Query



What Makes Teradata Unique

Key Components/Concepts



Teradata MPP Architecture

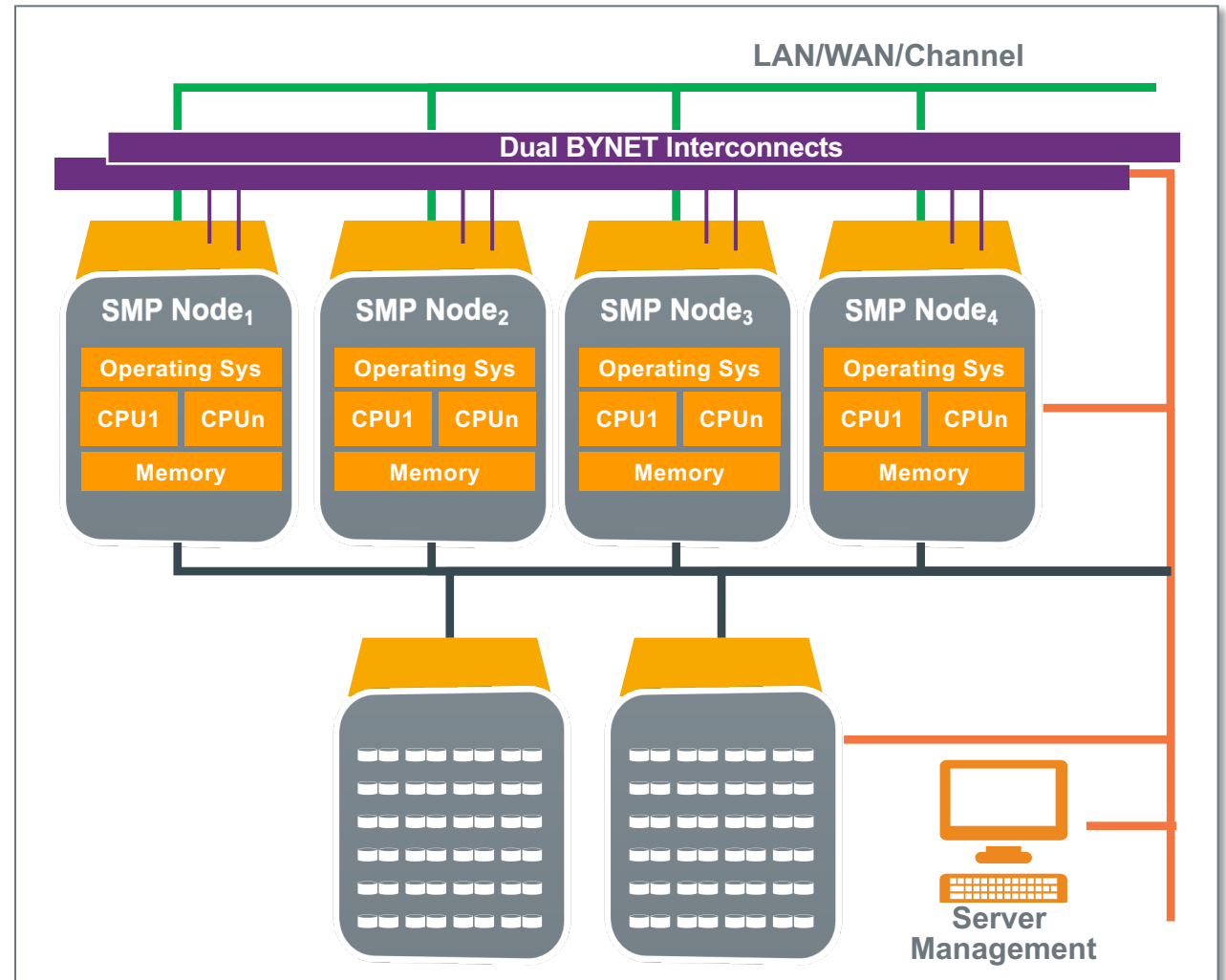
Key Concepts

- As the system scales (e.g. add more nodes), it should be no additional administrative overhead.
- As it scales, there should be linear improvement of performance.
- The rest of the world is catching on to our MPP concept. We've done it from the day one.
- Teradata no longer assumes that we always provide the hardware/platform that our database runs on.



Teradata MPP Server Architecture

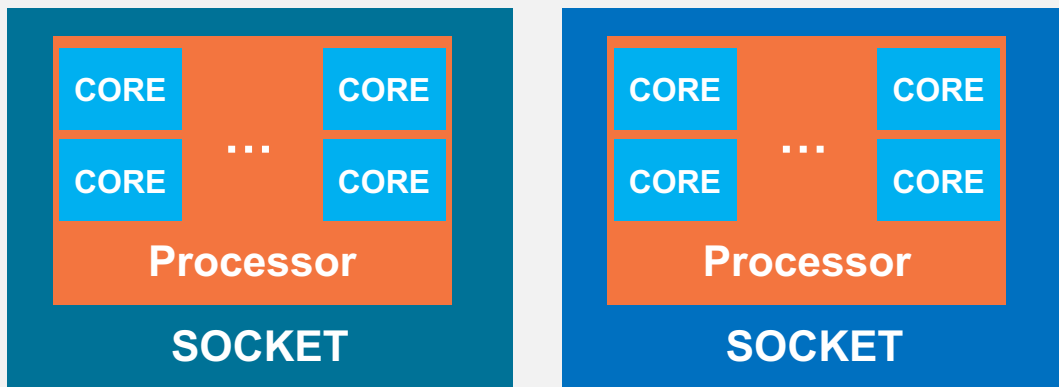
- Nodes
 - Incrementally scalable to 2048 nodes
- Operating System
 - Linux (SUSE)
- Storage
 - Independent I/O
 - Scales per node
- BYNET Interconnect
 - Fully scalable bandwidth
- Connectivity
 - Fully scalable
 - Traffic spread across all nodes
 - Channel – ESCON/FICON
 - LAN, WAN
- Server Management
 - One console to view the entire system



Industry-Leading Intel® Xeon Technology

Teradata leverages leading Intel technology to deliver industry-leading performance:

- Leverage Moore's Law and use massive numbers of transistors to build multiple core processors
- High-performance Intel Xeon processors
 - We stay very current with the Intel roadmap



For example:

- Twelve-Core Processor
- Two Processors = 24 Cores/node
- Two Hyper-threads per Core = 48 virtual cores

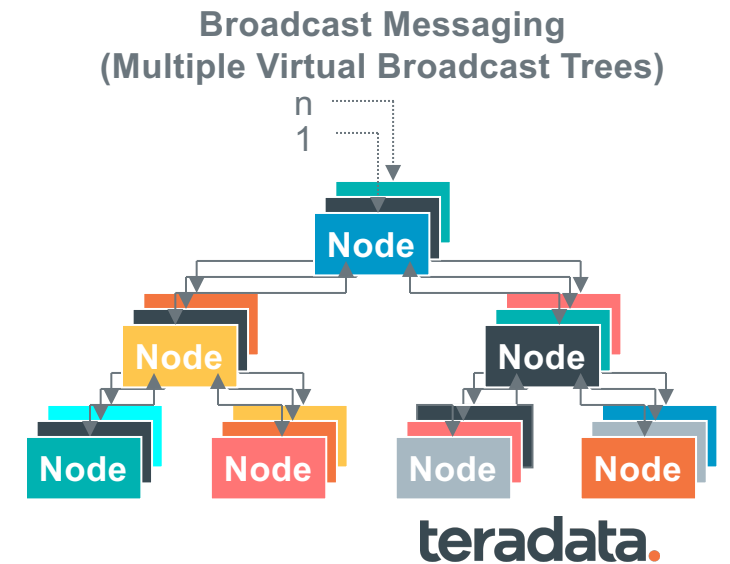
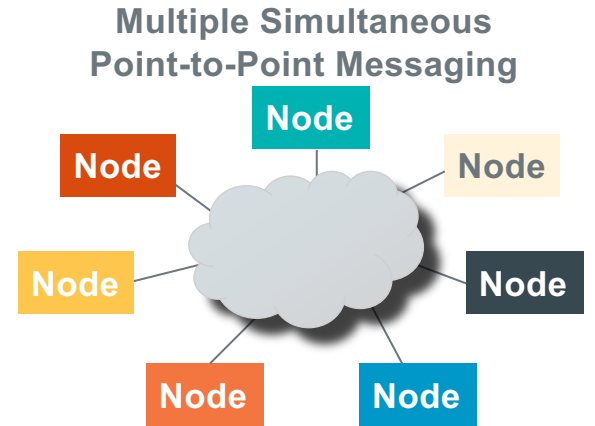
- **Optimized for Teradata Active Data Warehouse Performance**

- Linearly scalable bandwidth
- Bynet Low Latency Interface (BLLI) – Streamlined communication protocol
- Bynet software provides unique Teradata features: broadcast message support, row merge support, multi-fabric message traffic shaping, software guaranteed message delivery (point to point and broadcast)
- Broadcast reliability is implemented in software in the form of multiple virtual broadcast trees.
- Supports several physical fabrics; InfiniBand , Ethernet (1Gb & 10Gb), Native Bynet.

- **Proven High-Availability**

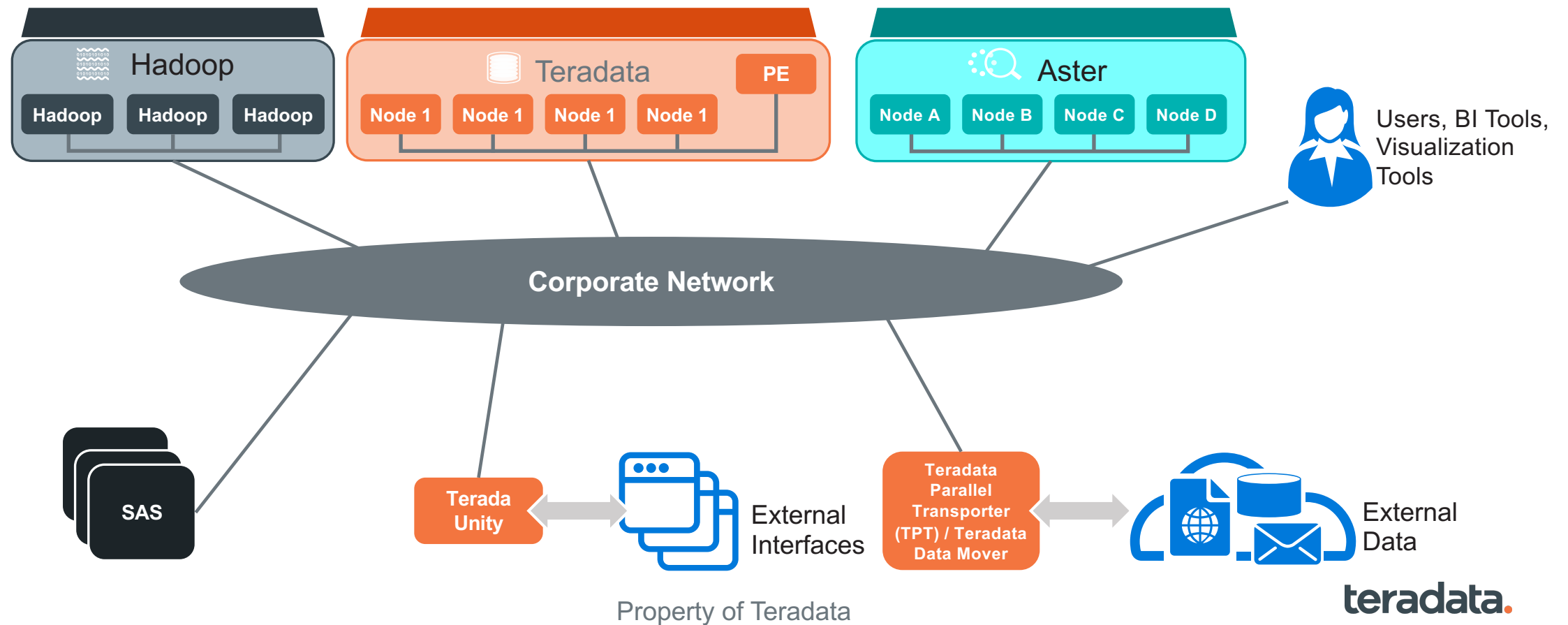
- Each fabric is fault tolerant (multiple paths, redundant power & cooling)
- 2 active and independent fabrics (no single point of failure)

The **Teradata Optimizer** chooses between Point-to-Point and Broadcast Messaging to select the most effective communication.



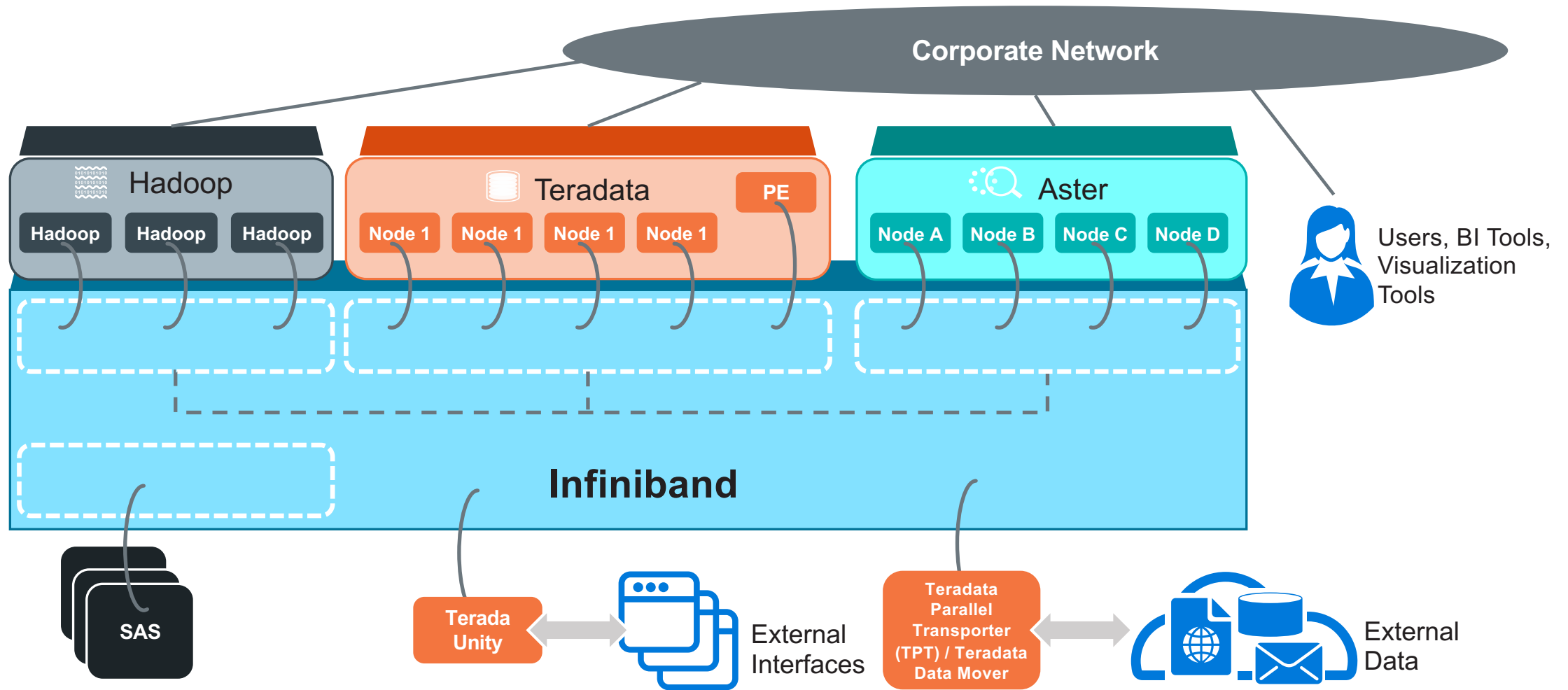
Network Connected Teradata Unified Data Architecture

MPP
Architecture

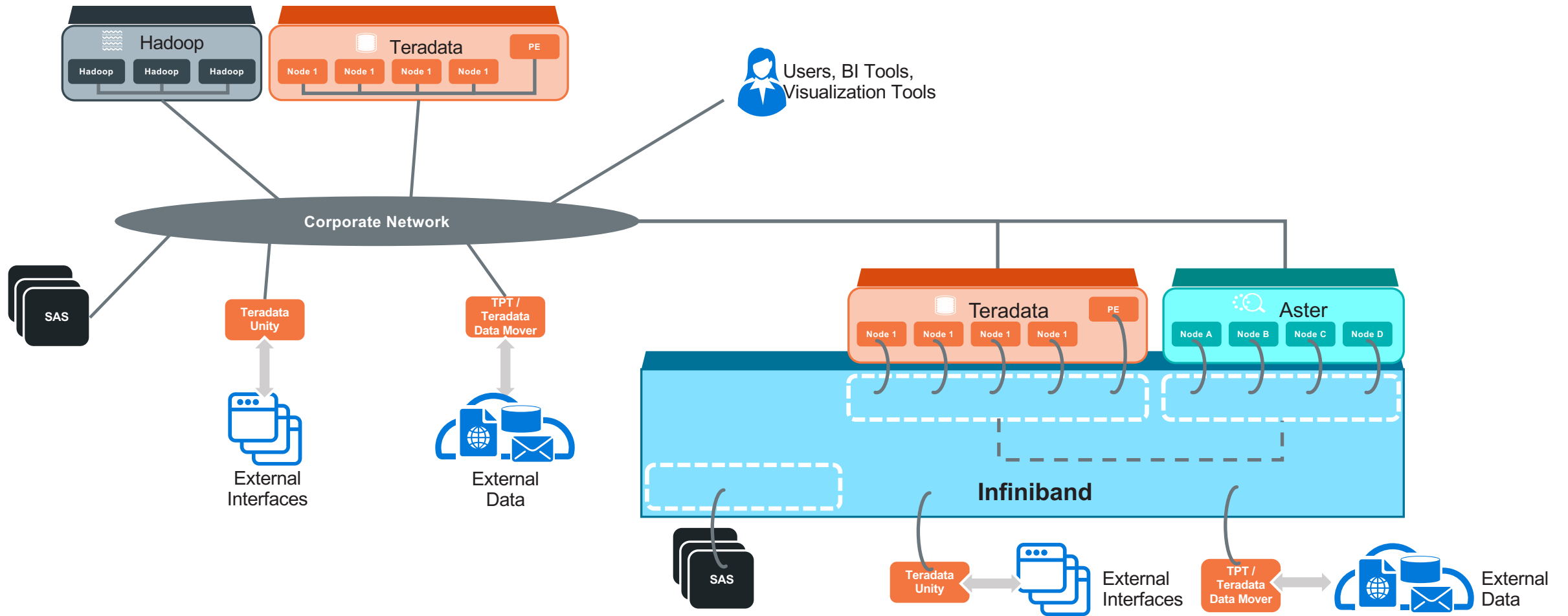


Infiniband

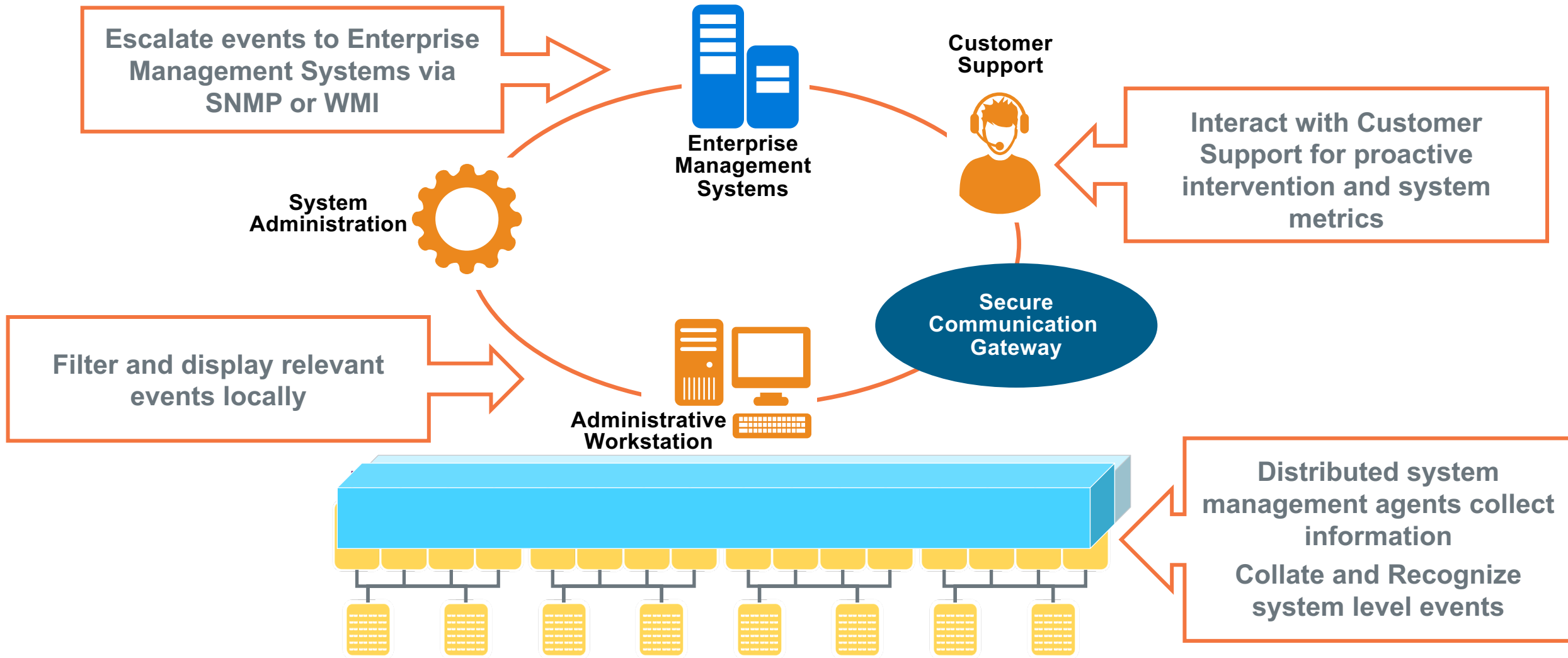
Integrated Communication for Teradata Unified Data Architecture



Infiniband Hybrid Architecture



Teradata Systems Management



Summary

Teradata is Teradata

Teradata Workload-Specific Platform Family



On-Premises

Teradata Cloud



Managed Cloud

Teradata Database on VMware®



Virtual Environment

Teradata Database on

...

Microsoft Azure

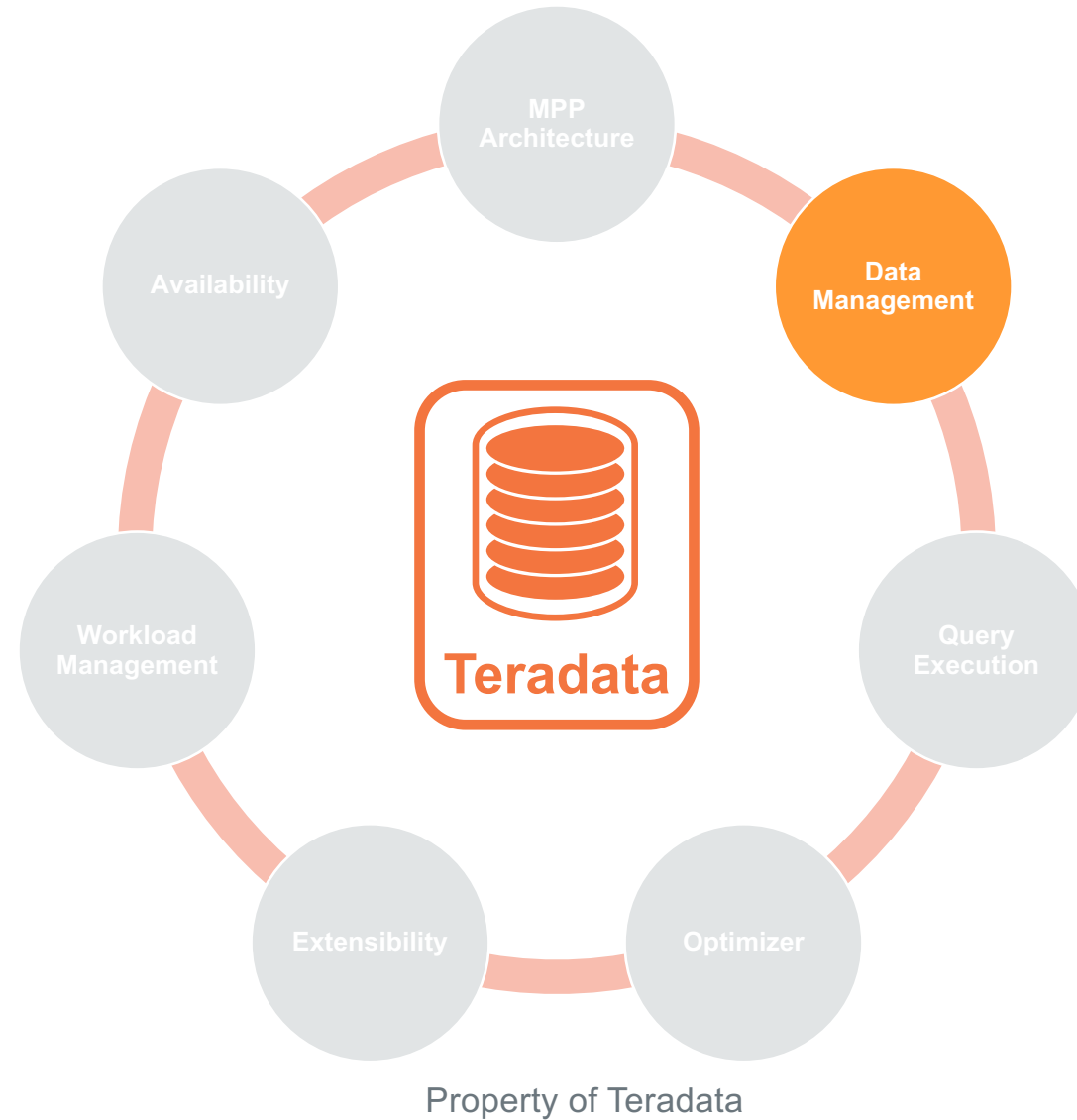


Google Cloud Platform

Public Cloud

What Makes Teradata Unique

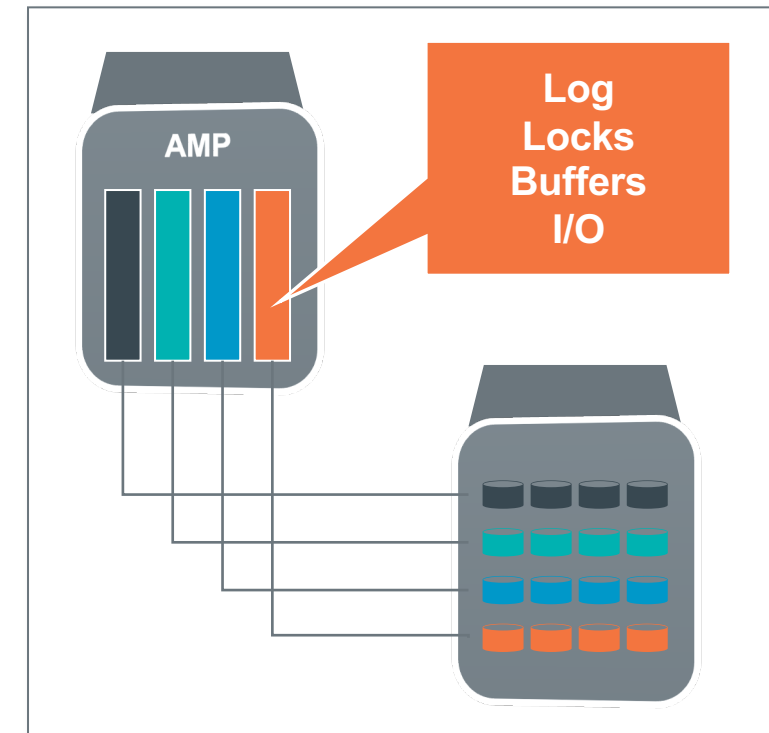
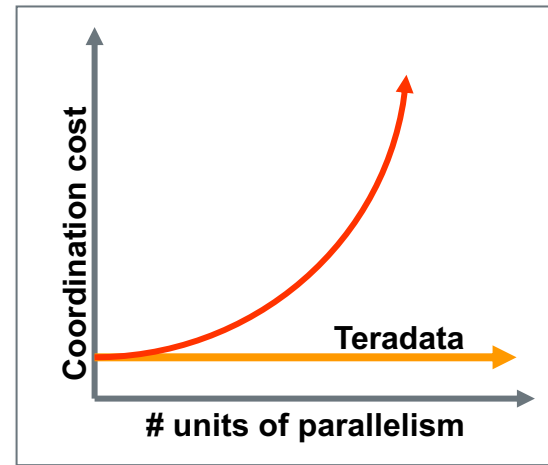
Key Components/Concepts



Shared Nothing – Managing the Data

Key Concepts

- Basis of Teradata scalability
 - Each AMP owns an equal slice of the disk
 - Only that AMP reads that slice
- No single point of control for any operation
 - I/O, Buffers, Locking, Logging, Dictionary
 - Nothing centralized
 - Exponential communication costs avoided



teradata.

Teradata Hash Maps - Basics

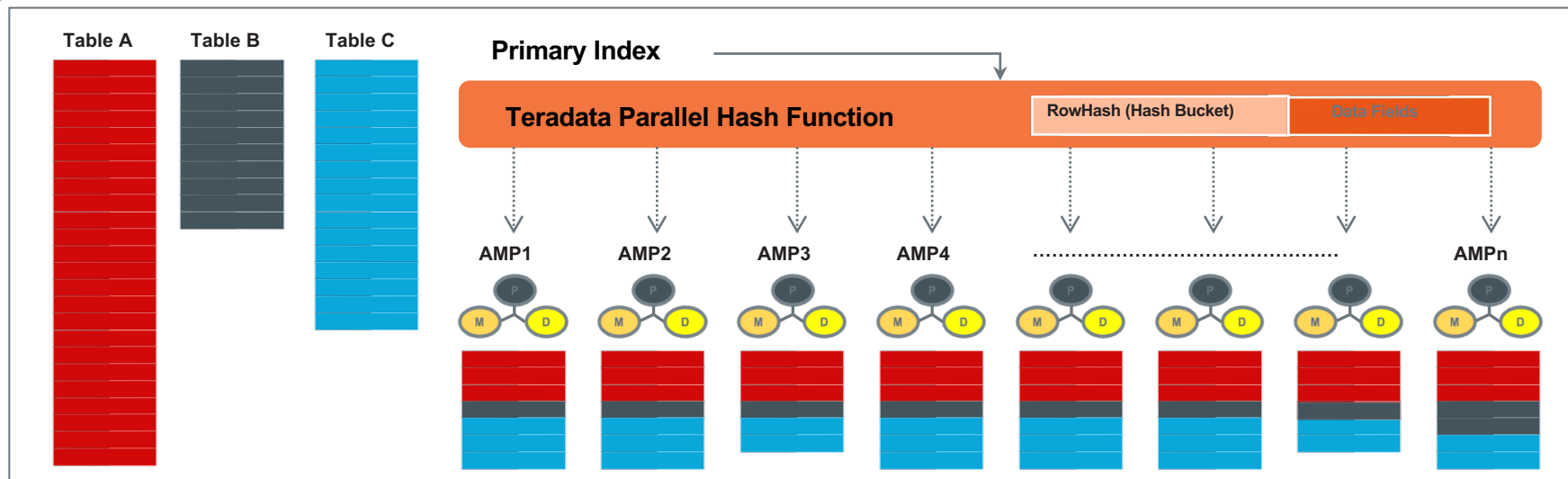
- A Teradata Hash Map
 - A hash map tells Teradata where data belongs in a parallel system.
- Primary Index data values are hashed
 - Hashing algorithm
 - Divide by 1M hash buckets
- AMPs own Hash Buckets
 - The "Hash Map" maps each of the 1M buckets to whatever number of AMPs exist within the system.
 - Rows are assigned to the AMPs via the hash bucket
- Spreads data evenly
 - For performance
 - Reduce skew
 - Enables co-located table rows



Teradata Data Management

Rows automatically distributed evenly by **hash partitioning**

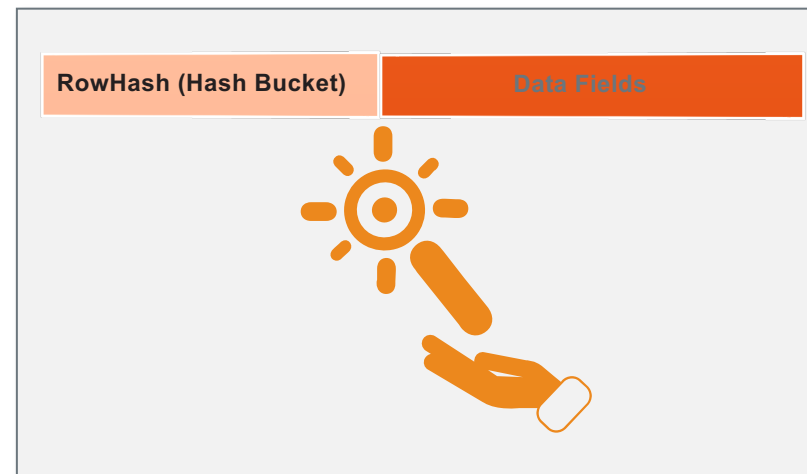
- Even distribution results in scalable performance
- Done in real-time as data are loaded, appended, or changed.
- Hash map defined and maintained by the system
 - 2^{32} hash codes, 1,048,576 buckets distributed to AMPs
- Primary Index (PI) column(s) are hashed
- Hash is always the same - for the same values
- No reorgs, repartitioning, space management



Property of Teradata

The Magic of Row Hash

- Store rows in a table in row hash order (logically)
- Direct access via Primary Index
 - For continuous update or single lookup
 - Hash the value
 - Go direct to AMP (routed by BYNET)
 - MI (in memory) -> CI -> data block
 - Very sparse N* tree
 - MI and CI very small structures
 - Guarantees no more than 2 IOs to get to any record via PI
 - CI cached if frequent access
 - “First Index is free”
- Sort merge join with no sort
 - Joins which never require scan of the large table
 - E.g. Star Join doesn't need to hash join or spool the fact table and sort it to join



Defining a Table in Teradata

Create the table

- Standard SQL syntax

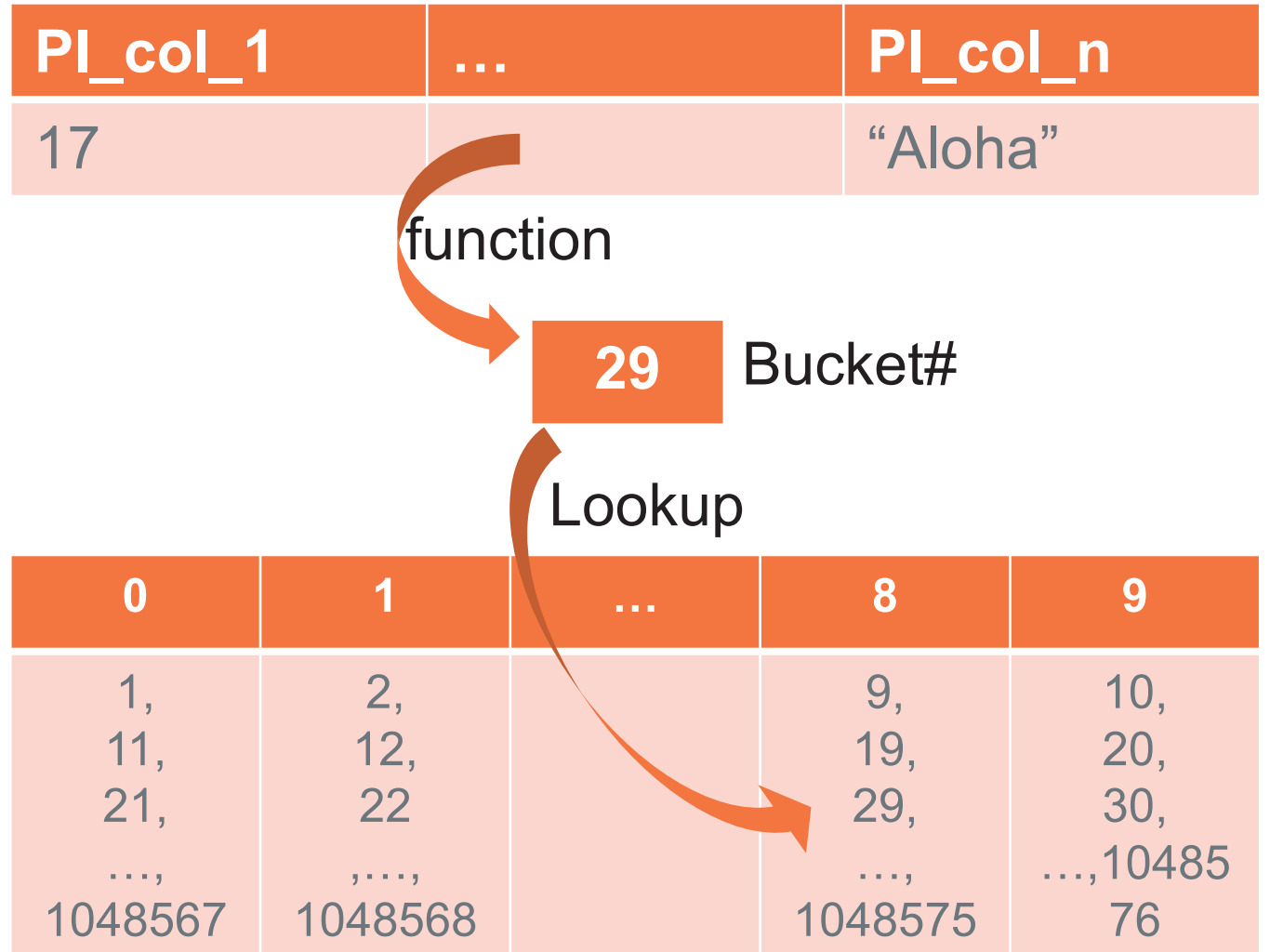
Define the primary index

- Extra line at end of table definition

```
CREATE TABLE LineItem (  
    OrderKey INTEGER NOT NULL,  
    PartKey INTEGER NOT NULL,  
    SupplierKey INTEGER,  
    LineNumber INTEGER,  
    Quantity INTEGER NOT NULL,  
    ExtendedPrice DECIMAL(13,2),  
    Discount DECIMAL(13,2),  
    Tax DECIMAL(13,2),  
    Comment VARCHAR(50)  
)  
PRIMARY INDEX (OrderKey);
```

Teradata Hash Maps - Basics

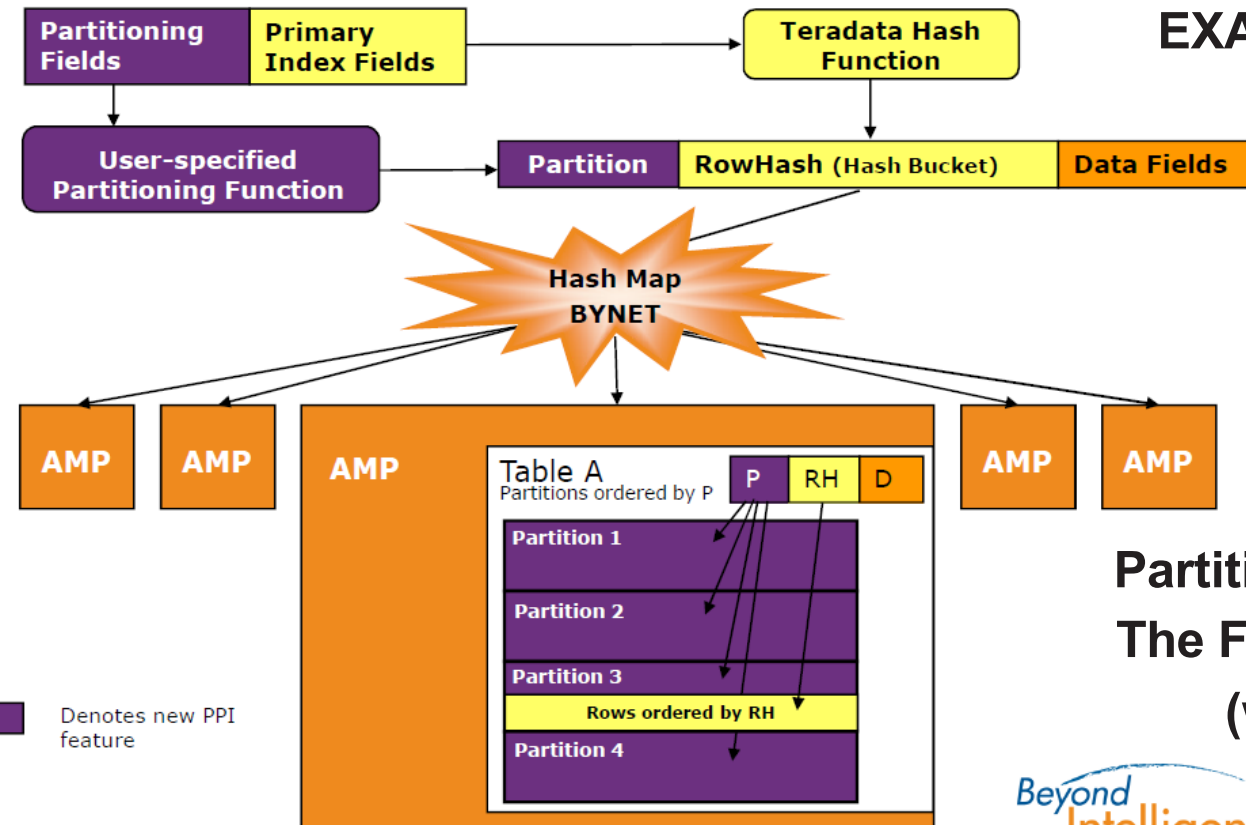
- Teradata Hash Maps
 - Key to the distribution of a table in an MPP System
 - The “PI” Fields are concatenated and run through a mathematical function that returns a “bucket”, a number between 1 and 1,048,576 (20 bits)
 - The “Hash Map” maps each of the buckets to whatever number of AMP’s exist within the system



Partitioned Primary Index (PPI)



Partitioned Primary Index One Table for all Queries – Recent and Historical



File System

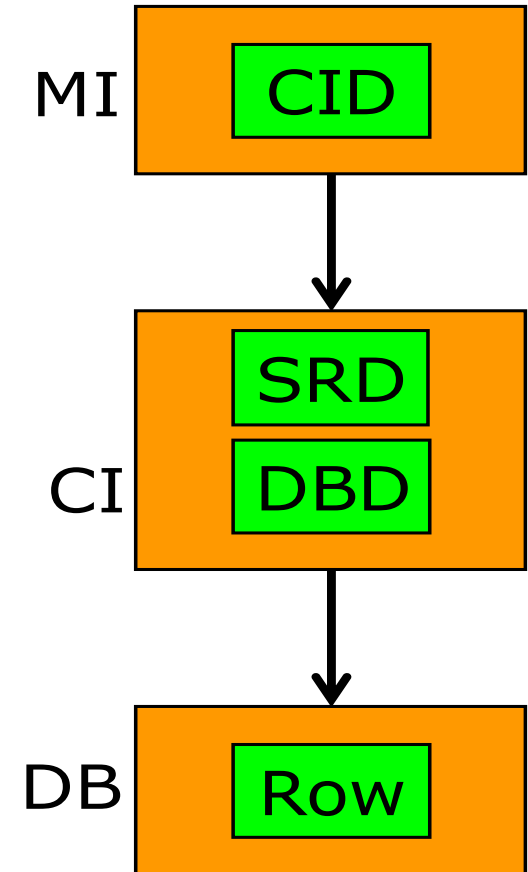
- Teradata wrote a new rule book
 - Old one written by IBM 40 years ago, used by most mainstream DBMSs today - except Teradata
- File system built of raw slices
- Rows stored in blocks
 - Variable length
 - Grow and shrink on demand
 - Rows located dynamically
 - May be moved to reclaim space, defrag
 - Maximum block size is configurable
 - System default or per table
 - 8K to 1MB
 - Change dynamically
- Blocks stored in 12MB allocation units - “cylinders”
- Indexes are just rows in tables



File System Structure

File System Index Structure

- B*Tree
- MI Master Index
 - > CID Cylinder Index Descriptor
- CI Cylinder Index
 - > SRD Subtable Refarray Descriptor
 - > DBD Data Block Descriptor
- DB Data Block
 - > ROW The actual physical row
- 3 level physical tree
 - > But MI always in memory



LOGICAL ROW ADDRESSES

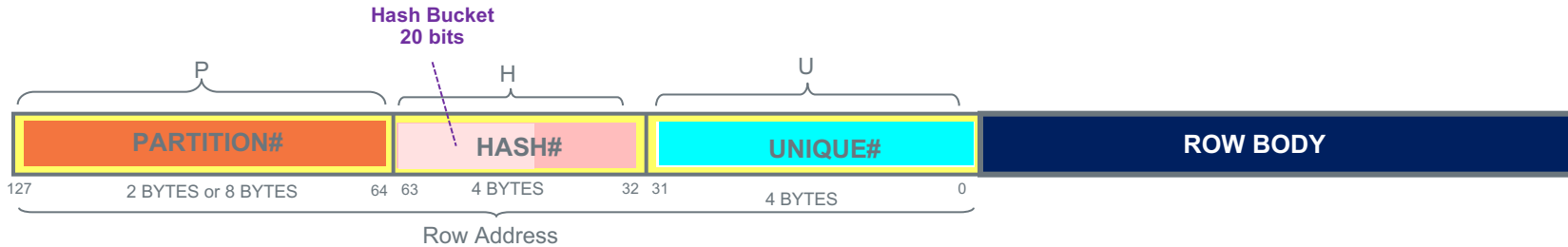


Table Row Address:

- Usage of optional Partition# field is dependent upon row-usage
 - Traditional Primary Index Table Rows: doesn't use partition# field
 - PPI-2 table rows : uses two bytes
 - PPI-8 table rows : uses 4 bytes
- Usage of the Hash# field:
 - For table rows, this typically contains computed hash associated with the primary index
- Usage of the Unique# field:
 - For table rows, this is utilized to assign unique row address for the case in which there are matching hash values.
- **IMPORTANT!:** For partitioned tables, a “partition boundary” is a “row address boundary”, not a physical storage boundary. Partitions may cross data block and/or cylinder boundaries

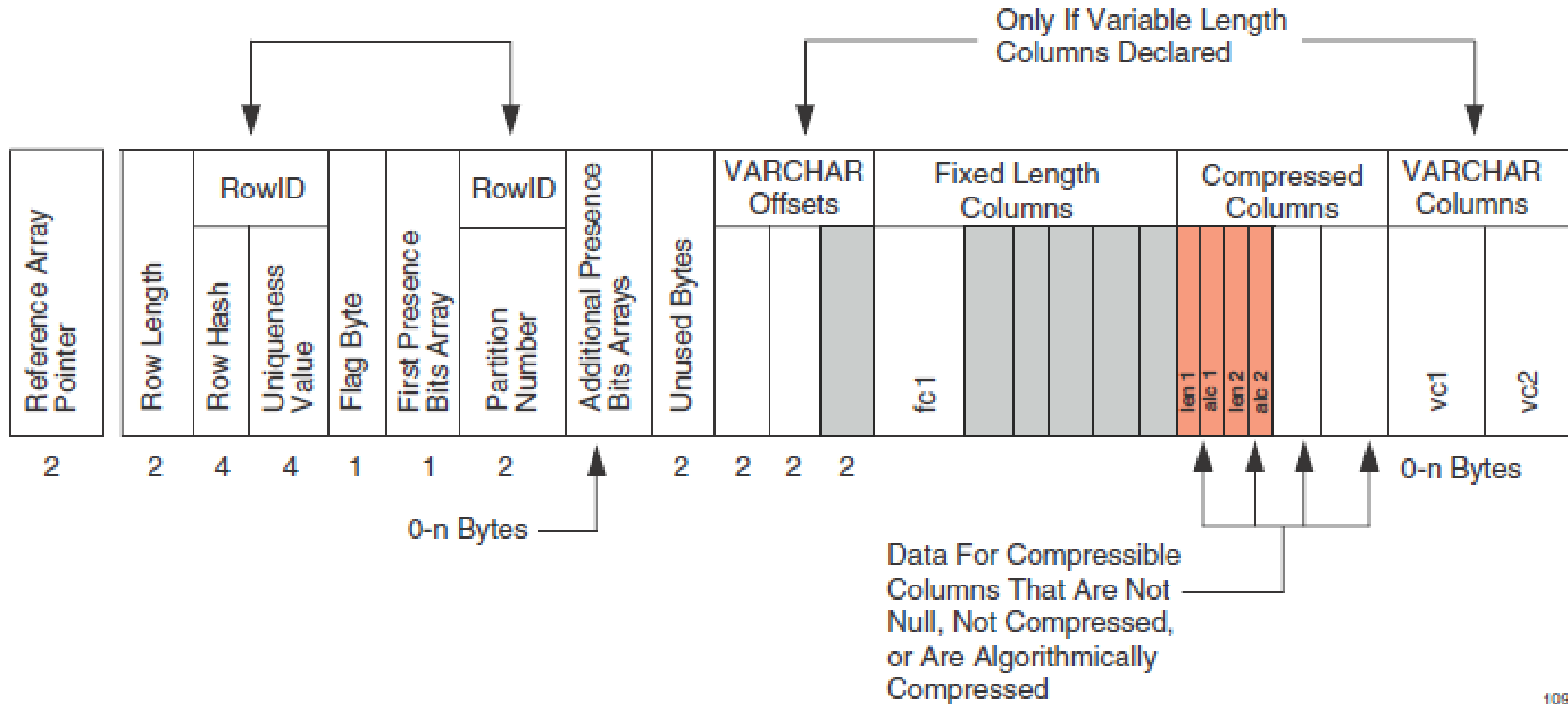
Spool File Row Address:

- May have table-like row address assignments or artificially assigned row addresses

Special Case RTREE-index sub-table rows

- Partition#(4 BYTES) and Hash#(4 BYTES) fields used as a single 8-BYTE field containing the “Virtual Block Identifier”
- Unique# field used to identify rows within a given Virtual Block

Base Table Row Formats



1094A097

Space Allocation

- Space allocation is entirely dynamic
 - No tablespaces or journal spaces or any pre-allocation
 - Spool (temp) and tables share space pool, no fixed reserved allocations
- If no cylinder free, combine partial cylinders
 - Dynamic and automatic
 - Background compaction based on tunable threshold
- Quotas control disk space utilization
 - Increase quota (trivial online command) to allow user to use more space



Partitioned Primary Index (PPI) Comparison

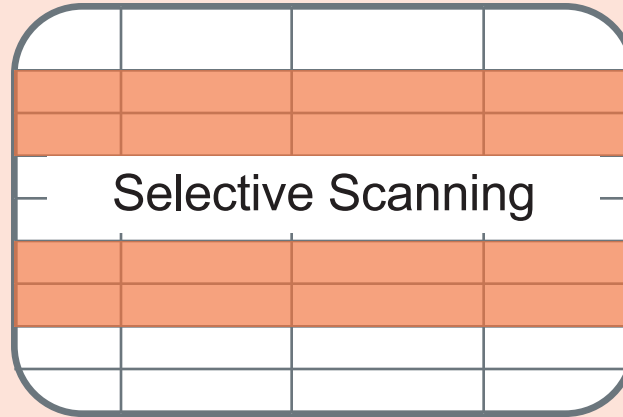
No Partitions

Looking for a customer's transactions across all time and states



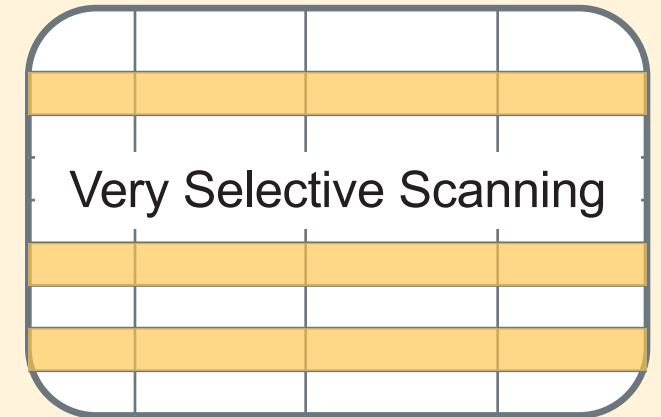
Single Level Partition

Looking for two particular weeks



Multi-level Partitions

Looking for two states in two weeks



Teradata Columnar

Row Store – Employee Table

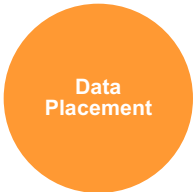
Emp ID	Lastname	Firstname	Salary
1	Smith	Joe	40000
2	Jones	Mary	50000
3	Johnson	Cathy	44000

Column Store – Employee Table

Emp ID	Lastname	Firstname	Salary
1	Smith	Joe	40000
2	Jones	Mary	50000
3	Johnson	Cathy	44000

- Key advantages:
 - Better query **performance** – read only the columns you need
 - Smaller **data size** – better compression on homogeneous column values
- Query execution optimized for column processing
- Full hybrid columnar
 - Row tables, column tables, mixed tables, and mixed index and table
 - Physical design optimization option

Multi-Level Partitioned Primary Index (PPI)



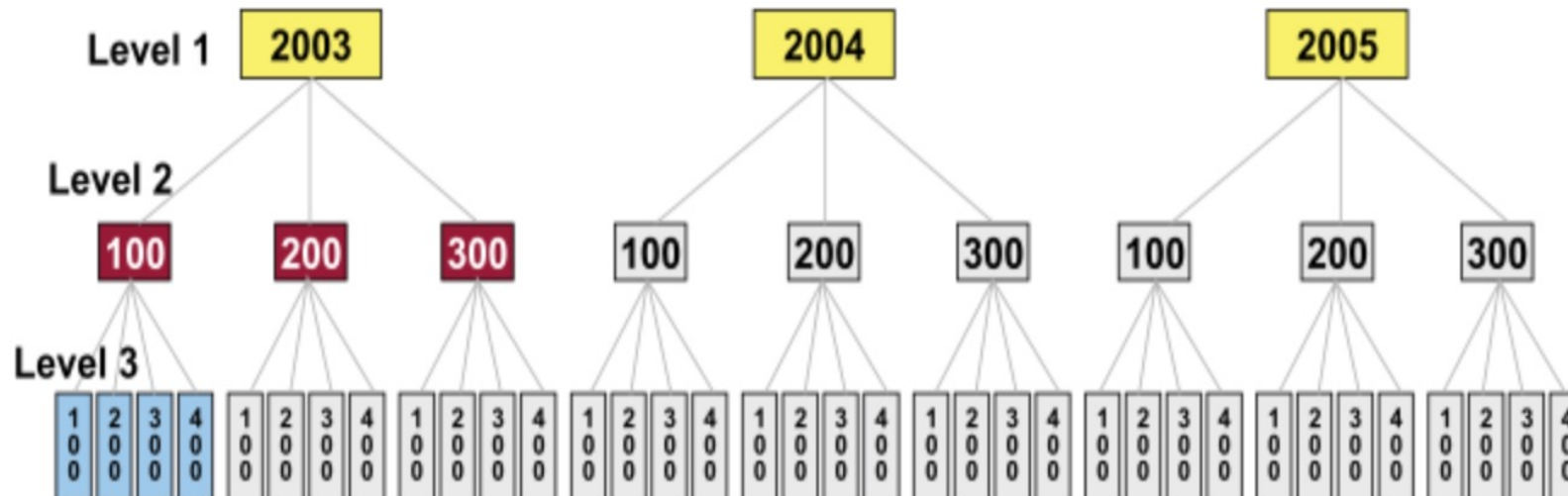
```

CREATE TABLE Sales
(storeid INTEGER NOT NULL,
productid INTEGER NOT NULL,
salesdate DATE FORMAT 'yyyy-mm-dd' NOT NULL,
totalrevenue DECIMAL(13,2),
totalsold INTEGER,
note VARCHAR(256))
UNIQUE PRIMARY INDEX (storeid, productid, salesdate)
PARTITION BY (
  RANGE_N(salesdate BETWEEN DATE '2003-01-01' AND DATE '2005-12-31'
    EACH INTERVAL '1' YEAR),
  RANGE_N(storeid BETWEEN 1 AND 300 EACH 100),
  RANGE_N(productid BETWEEN 1 AND 400 EACH 100));
    
```

Input to partition function			Sales						
L1	L2	L3	salesdate	storeid	productid	total revenue	totalsold	note	
1	1	1	2003-04-15	96	10	4158	42	Good day	
1	1	2	2003-07-06	71	184	1972	68	Marginal	
2	2	3	2004-11-09	175	241	3055	47	Slow day	
1	1	4	2003-12-24	82	363	1261	13	Promotion	

Multi-Level Partitioned Primary Index (PPI)

Values at Level 1 x Level 2 x Level 3
= 3 x 3 x 4
= 36 (must be less than 65,535)

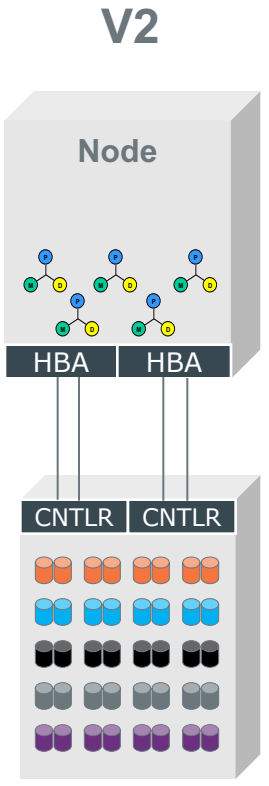
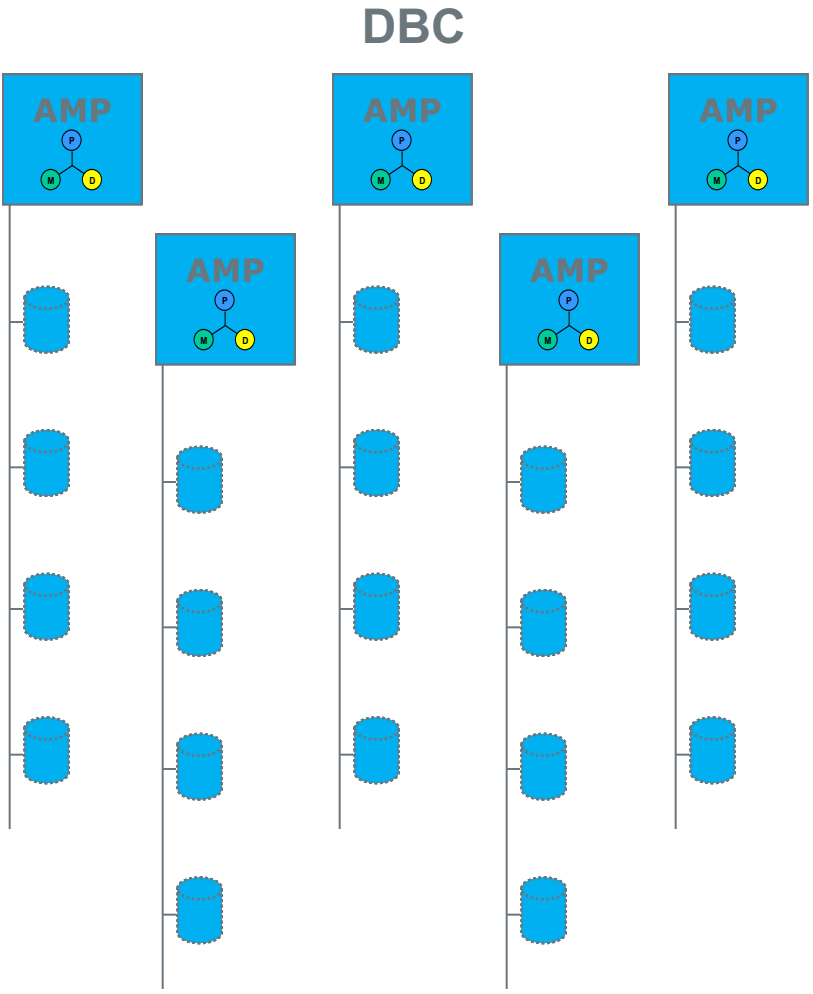


How Rows for the Sales Table Are Grouped on an AMP

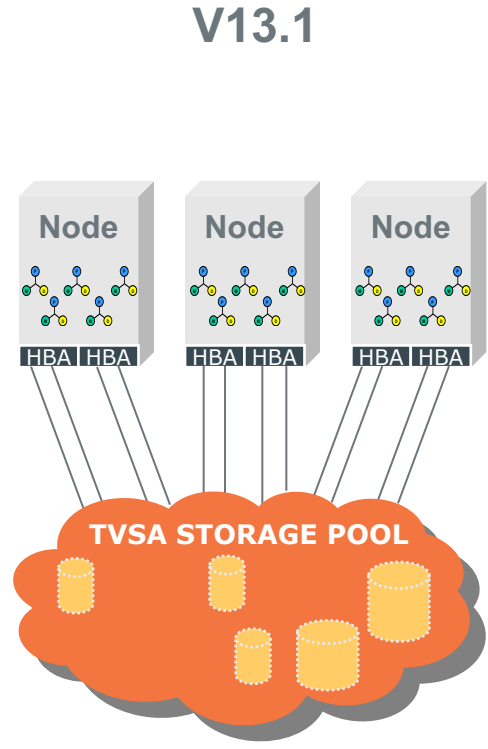
Compression

	Multi Value	Algorithmic	Data Block
What it is	<ul style="list-style-type: none">• Dictionary-based compression that replaces values specified by the user	<ul style="list-style-type: none">• Make your own - It applies a compression / decompression algorithm to a column of data	<ul style="list-style-type: none">• It compresses all types of data in a data block before it's stored on a disk
Suited for	<ul style="list-style-type: none">• Data with repeated values	<ul style="list-style-type: none">• Data with well-known attributes	<ul style="list-style-type: none">• Data with a low frequency of access (cold data)
Type of Data	<ul style="list-style-type: none">• NULLs + Common values	<ul style="list-style-type: none">• UTF-16 to UTF-8, long strings	<ul style="list-style-type: none">• Works on any/all data
Characteristic	<ul style="list-style-type: none">• Zero CPU overhead	<ul style="list-style-type: none">• Can write custom procedures	<ul style="list-style-type: none">• Built-in library
Compression Ratio	<ul style="list-style-type: none">• 30-40%	<ul style="list-style-type: none">• 50% (2x)	<ul style="list-style-type: none">• 3X or more

TERADATA FILE SYSTEM: The Evolution of TVSA

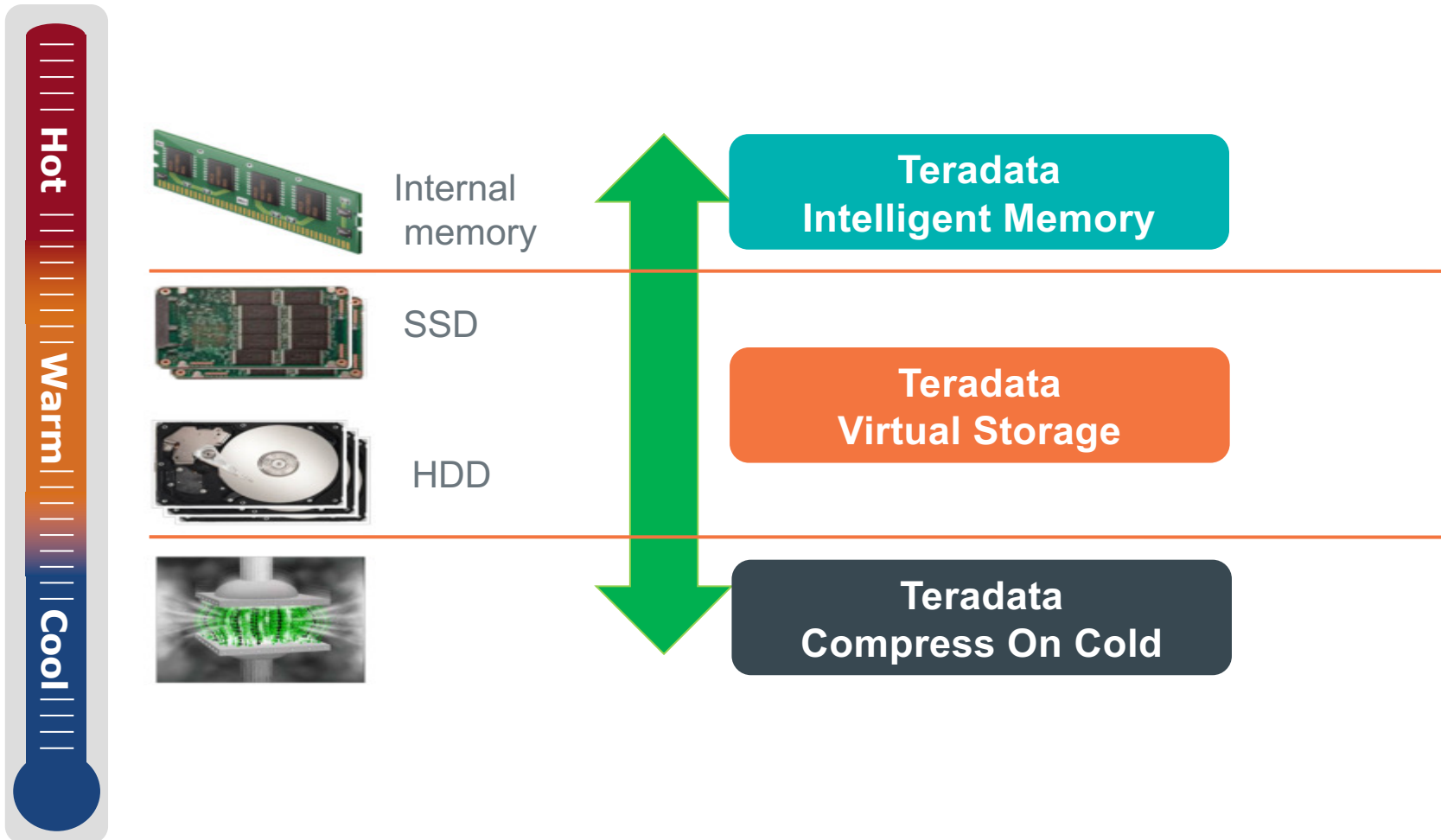


- Nodes virtualized from AMPs
- Disks consolidated in Arrays
- Disks are still tied to AMPs



Disks are virtualized from Amps

Integrated Management of Hybrid Storage



- Automatic
- Transparent
- No DBA effort
- No SQL changes
- Maintain user access to ALL data for analysis
- Avoid separate systems and copies of data for each use case

Summary

Automate and Eliminate

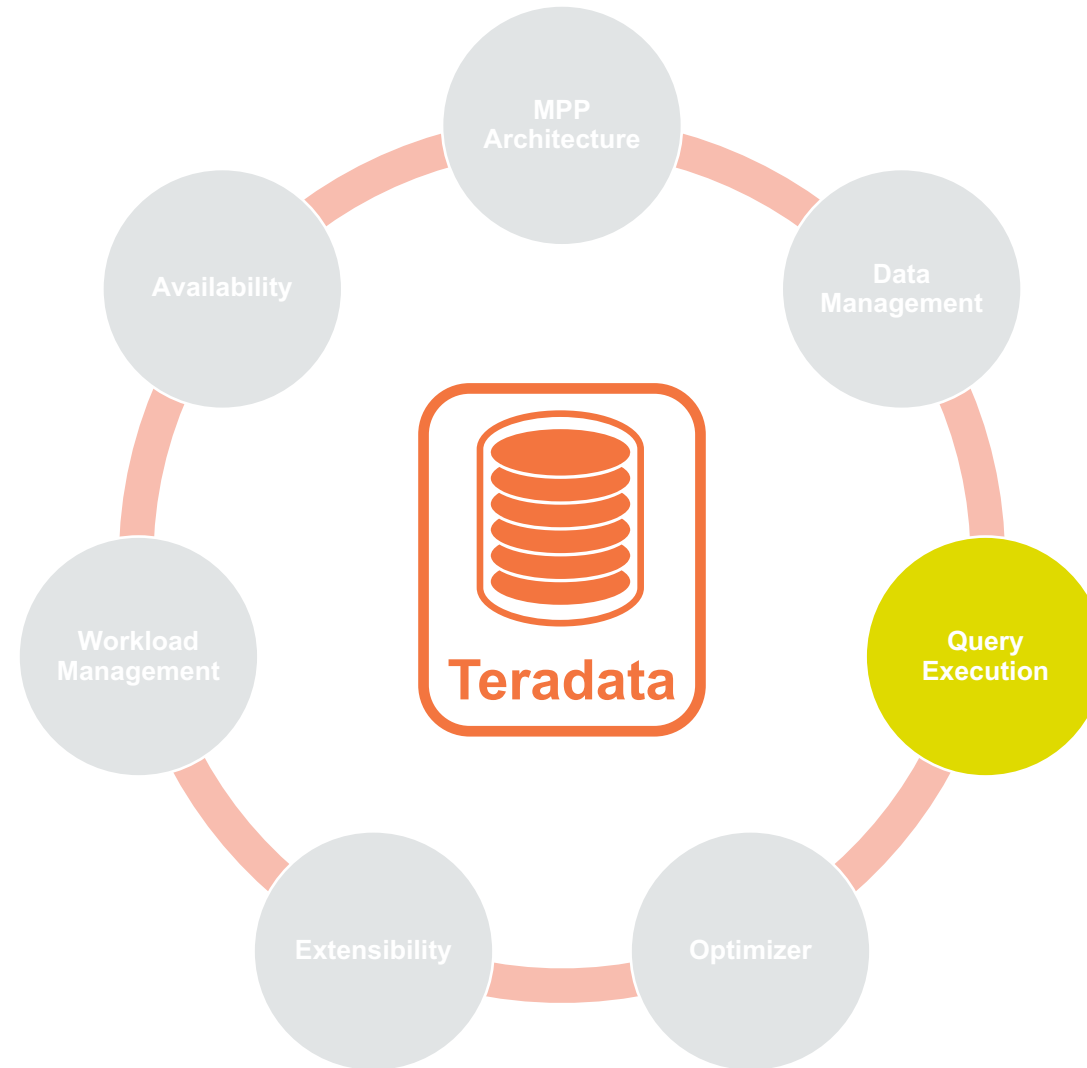
- Database should manage the data, not DBA or the user
- Data management should be as invisible as possible
 - Define and forget
- Enables user self-service

- No reorgs
 - Don't even have a reorg utility
- No index rebuilds
- No re-partitioning
- No detailed space management
- Easy database and table definition
- Minimum ongoing maintenance
 - All performed automatically



What Makes Teradata Unique

Key Components/Concepts



Property of Teradata

Overview: Query Execution

- Teradata query execution is designed to make every analytic query scalable
- Enabling in-database execution is a primary goal
- Parallelization is built-in throughout all operations and is fully automatic



What's on a Node

- **Gateway**

- Connect sessions to outside world
- Balance external traffic workload across nodes

- **Parsing Engine (PE)**

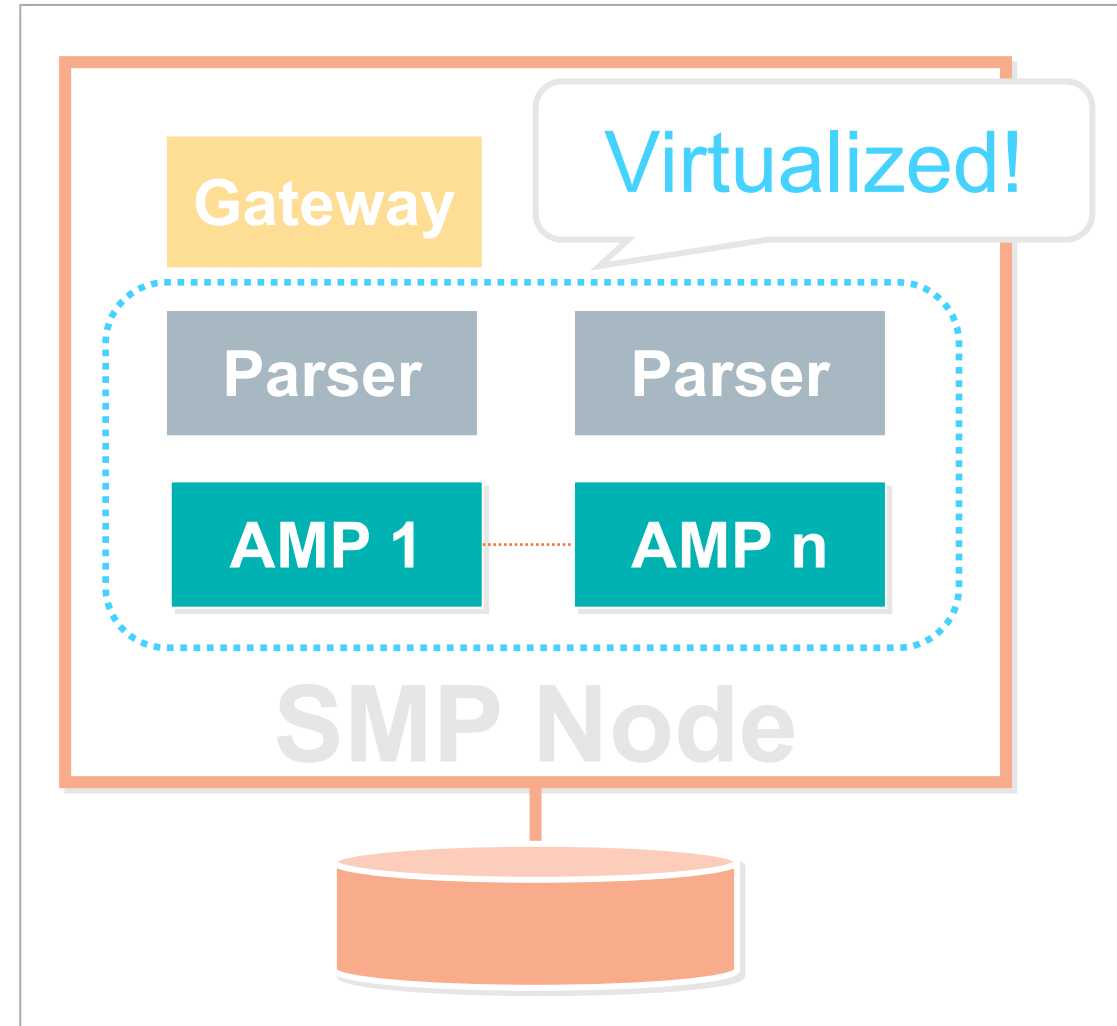
- Parse & Optimize
- Dispatcher to AMPs

- **AMP (Access Module Processor)**

- Execution engine
- Logs & locks
- Data dictionary
- I/O management

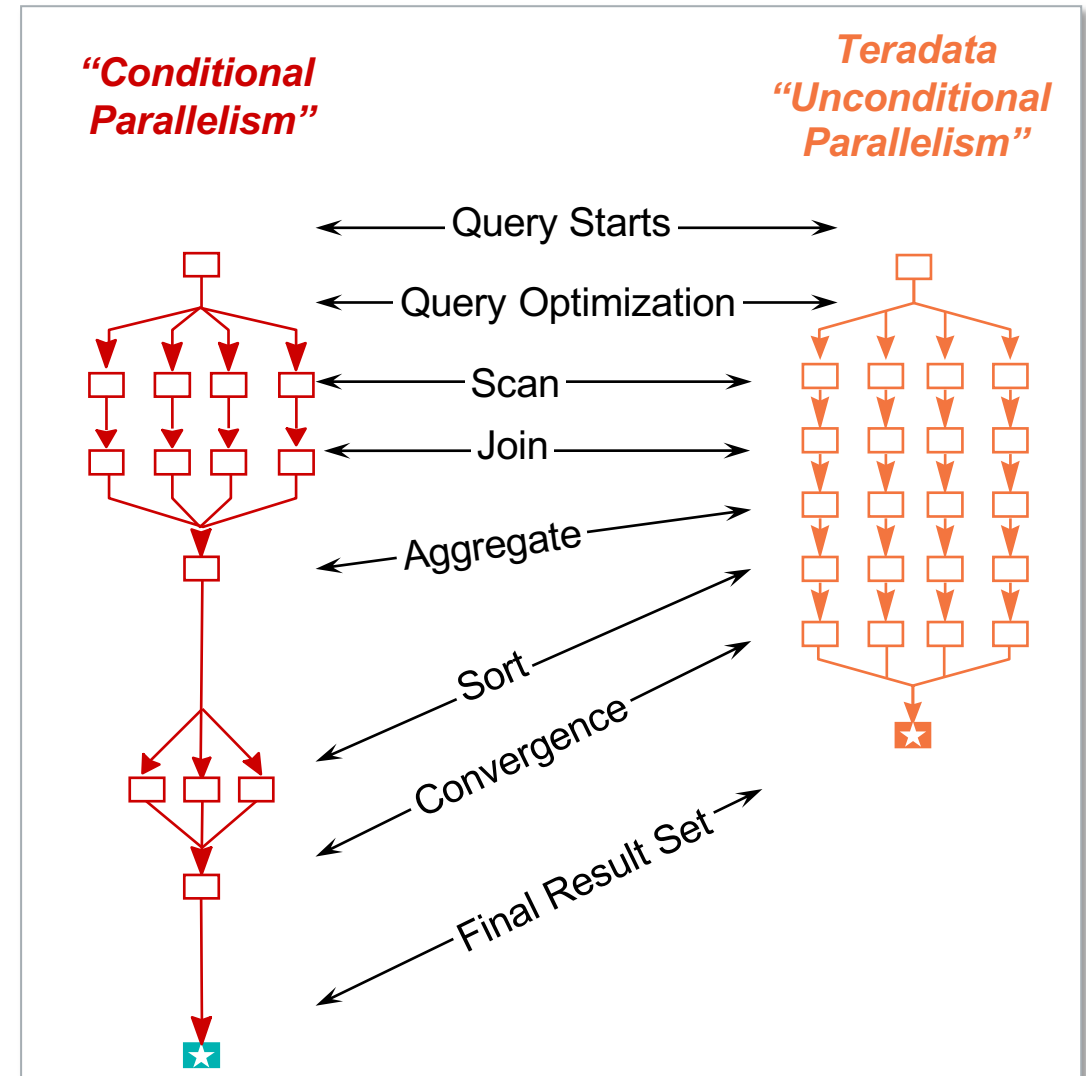
- **“Vprocs”**

- Virtual “processors” sharing one physical node



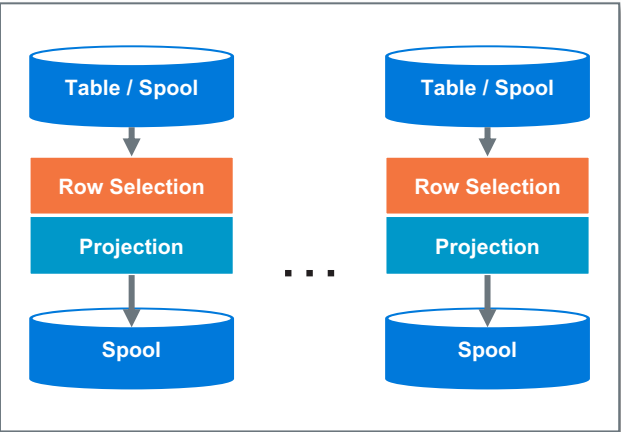
Query Parallelization

- Query parsing, management is fully distributed across the nodes
 - No head node/coordinator node
- All operations fully parallel
 - No single threaded operations
 - Scans, Joins, Index access, Aggregation, Sort, Insert, Update, Delete
 - Ordered Analytics
 - Extensibility functions
 - Result return



Data Access Methods

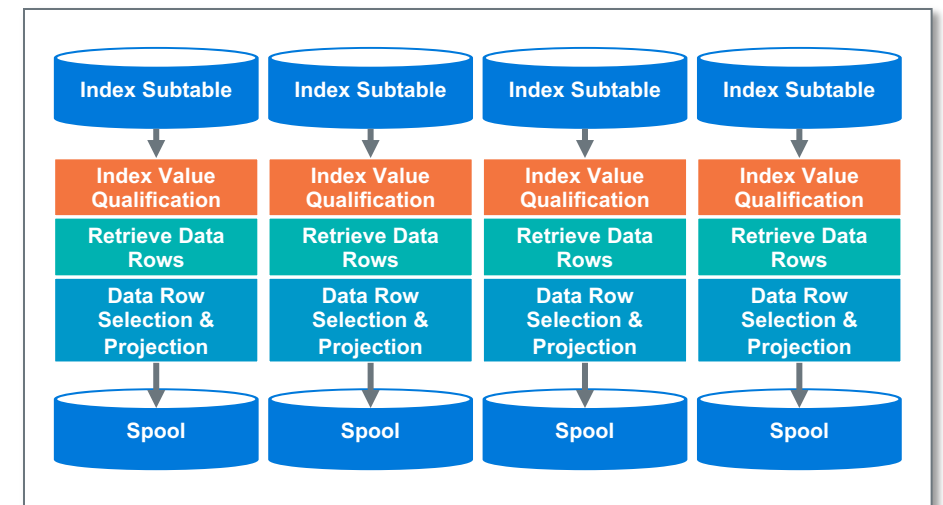
- Scan
 - Read all rows of a table
 - Executed on all AMPs simultaneously, each AMP reads the rows hashed to that AMP
 - A table or spool will be read, rows will be qualified (Row Selection)
 - Unneeded columns will be removed (Projection)
 - Result will be written to a spool file
 - Only interaction between the AMPs is for each of them to report completion
- Primary Index
 - Unique (UPI) and Non-Unique (NUPI)
 - Allows efficient direct access by value or via join
- Additional Indexes
 - Allow optimization of access
 - Used to support workloads that are high frequency, high SLA
 - Selected for use automatically by the query optimizer
 - PI, PPI and high performance scan eliminate need for many indexes



Secondary Index: Non Unique Secondary Index (NUSI)

Non Unique Secondary Index (NUSI)

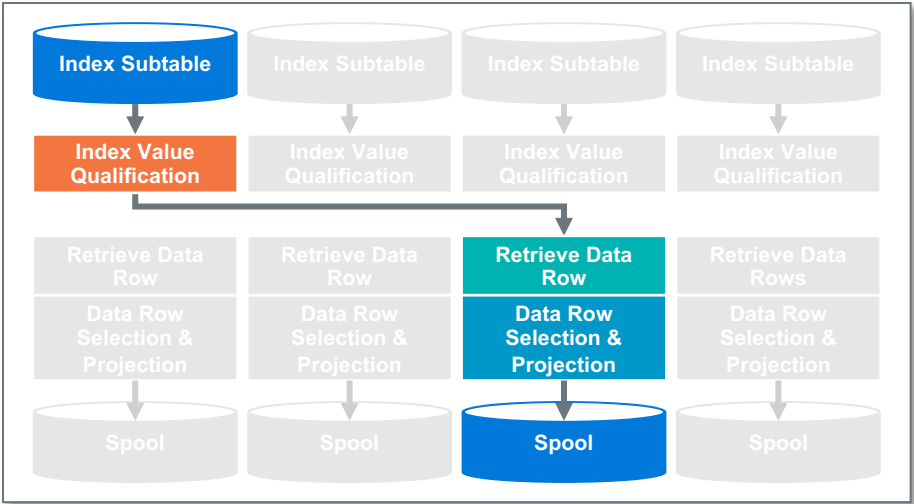
- Designed for accessing one or a small number of values where many rows have the same value
- Index entries are stored on the same AMP with the data rows
- An index entry is a value followed by a list of rowids of rows containing that value
- Ordered by hash of value being indexed
- Execution via a NUSI reads the index first, then uses the rowids to retrieve the data rows
 - Executed on All-AMPs simultaneously



Secondary Index: Unique Secondary Index (USI)

Unique Secondary Index (USI)

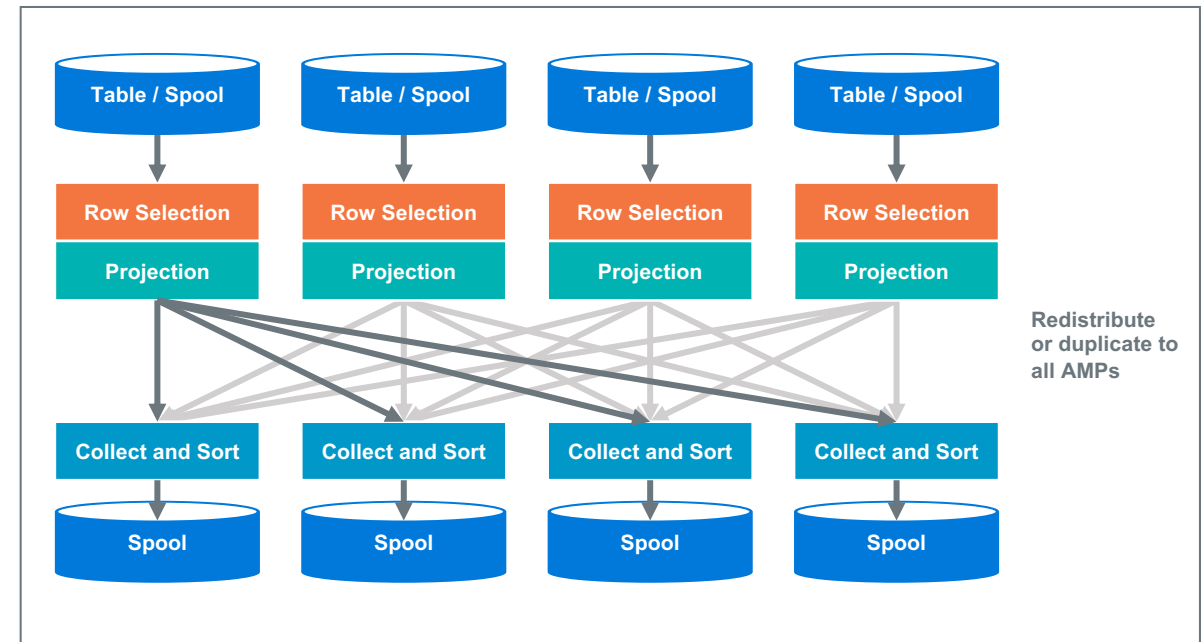
- Designed for accessing one or a small number of rows by value
 - Enforces Uniqueness
- Index entry is stored on hash AMP of value
- Index entry contains the value and the single rowid
- Execution via a USI accesses the index entry on the hash AMP(s) of the index value(s), then sends a message to the AMP of the rowid to retrieve the row
 - Executed on the minimum number of AMPs possible



Join Preparation

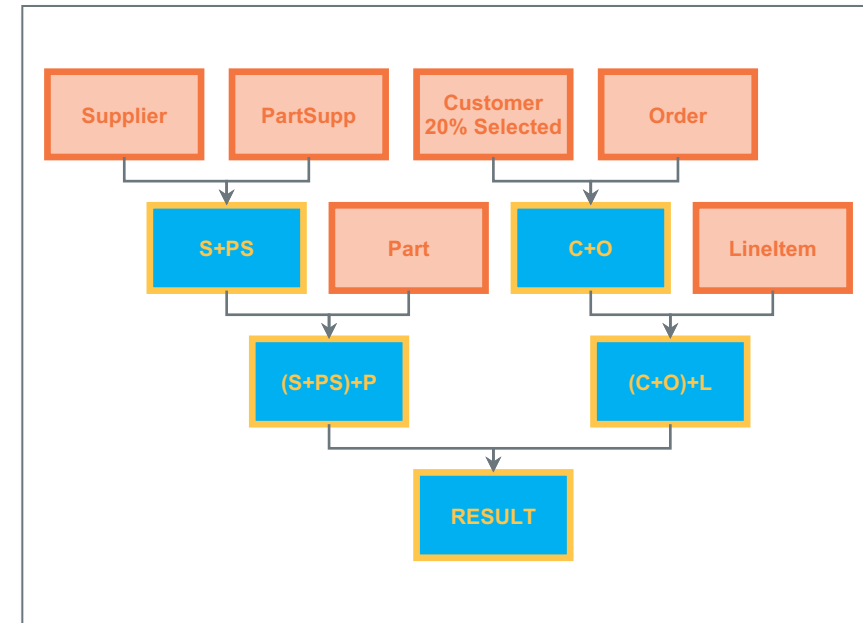
- Join preparation works in parallel across all AMPs
 - Selection to eliminate rows
 - Projection to eliminate columns
 - Eliminate bytes as early in the plan as possible
- Leverages BYNET to get data co-located for join
 - Redistribute : Hash by some set of attributes, often the join columns
 - Duplicate: Make a full copy of the prepared set on every AMP
- Collect results in receivers and write to spool or feed directly to first stage of sort if needed

- Uses Spool for intermediate results
 - Intermediate spools can be released after use (marked “last use” in explain)



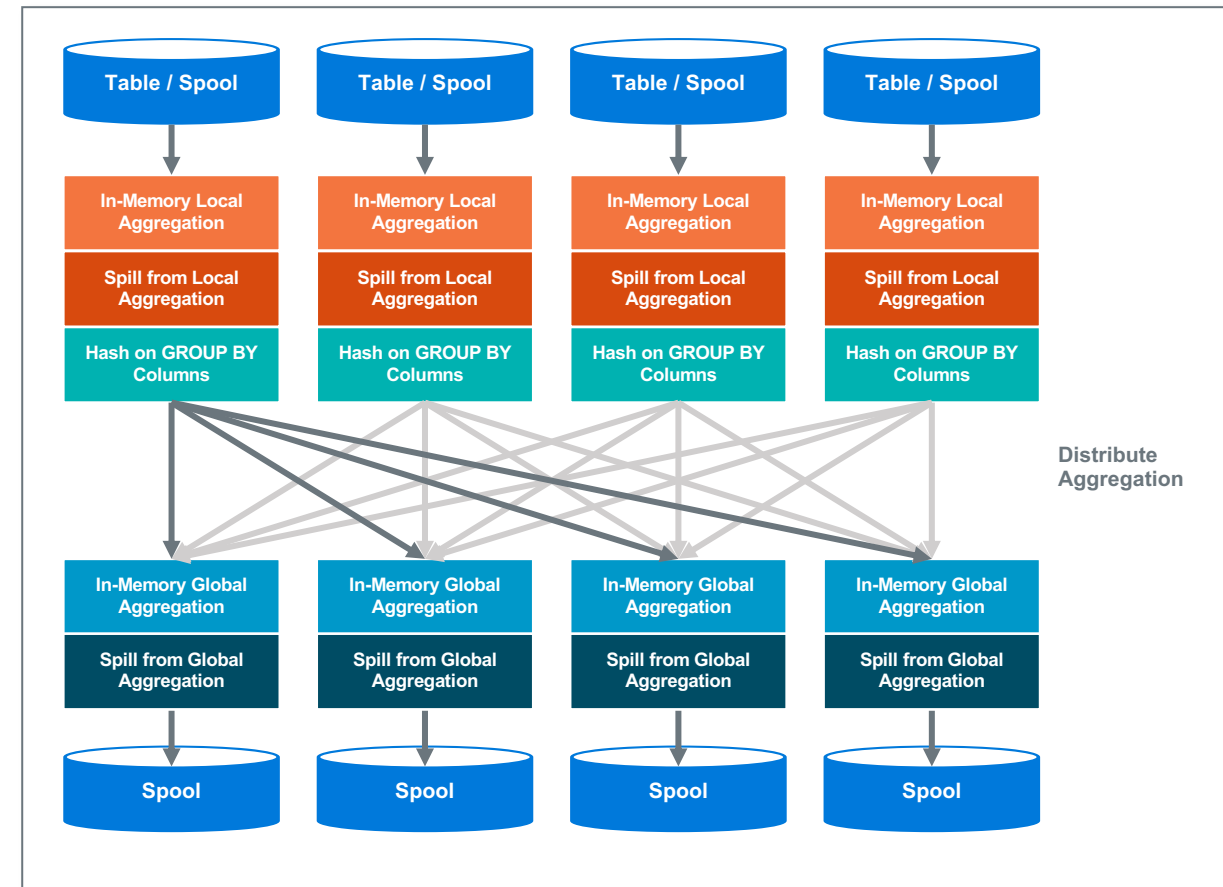
Joins

- Joins in Teradata all work in parallel and scale linearly
 - Must have data co-located in an AMP to join
- Join operates on two relations producing a join result which can be spooled locally or piped into redistribution and or sort as required
 - Queries with more than two tables are made into a series of two relation joins
- Many types of join operations available
 - Merge Join
 - Hash Join
 - Product Join
 - Nested Join
 - Outer Join versions of each
 - Exclusion and Inclusion Join for sub-query processing
- Each AMP performs it's part of the join independently
 - Communicate only when completely done



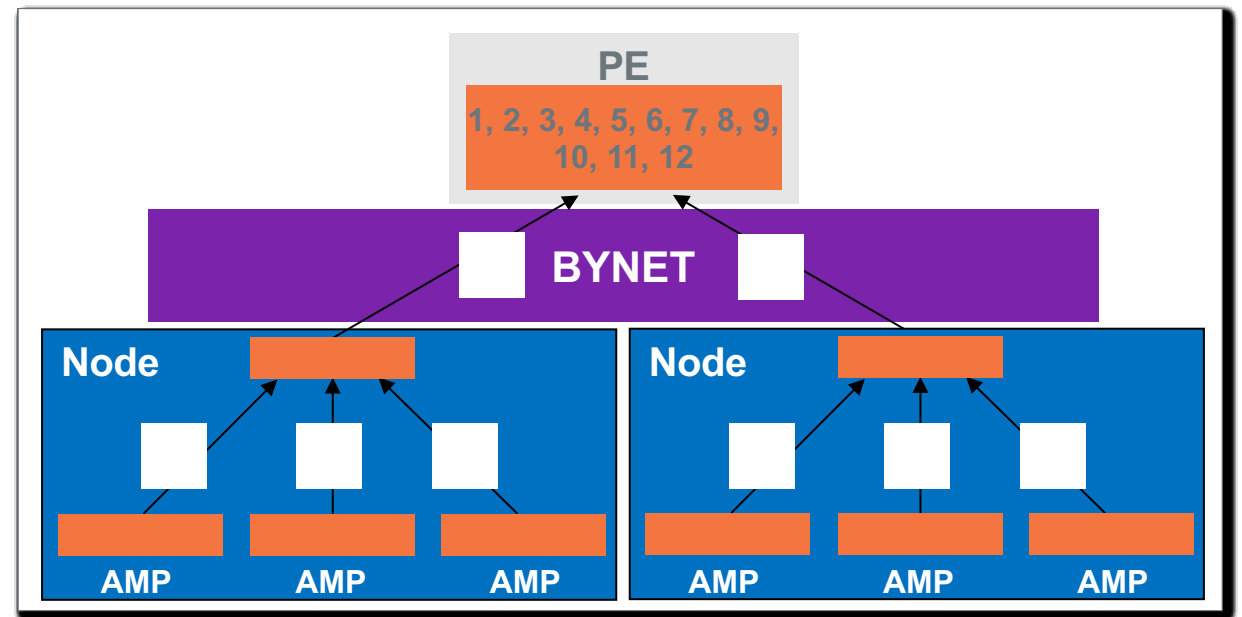
Parallel Aggregation Assisted by the Optimizer

- Aggregation makes full parallel use of all AMPs
 - Final result and final aggregation distributed across the platform
 - Linearly scalable across any data size
- Perform local aggregation resulting in (at most) one record per group per AMP
 - Skip if we know that little or no reduction will occur
- Redistribute by hash on GROUP BY columns to get all rows for a group together and distribute groups evenly across all AMPs
 - Skip if we know already distributed by GROUP BY columns
- Receive redistributed rows and perform Global Aggregation
- Caches at both local and global level do aggregation in-memory unless number of groups is very large
 - Cache spill that adapts to data pattern (e.g. most frequent keys don't spill)



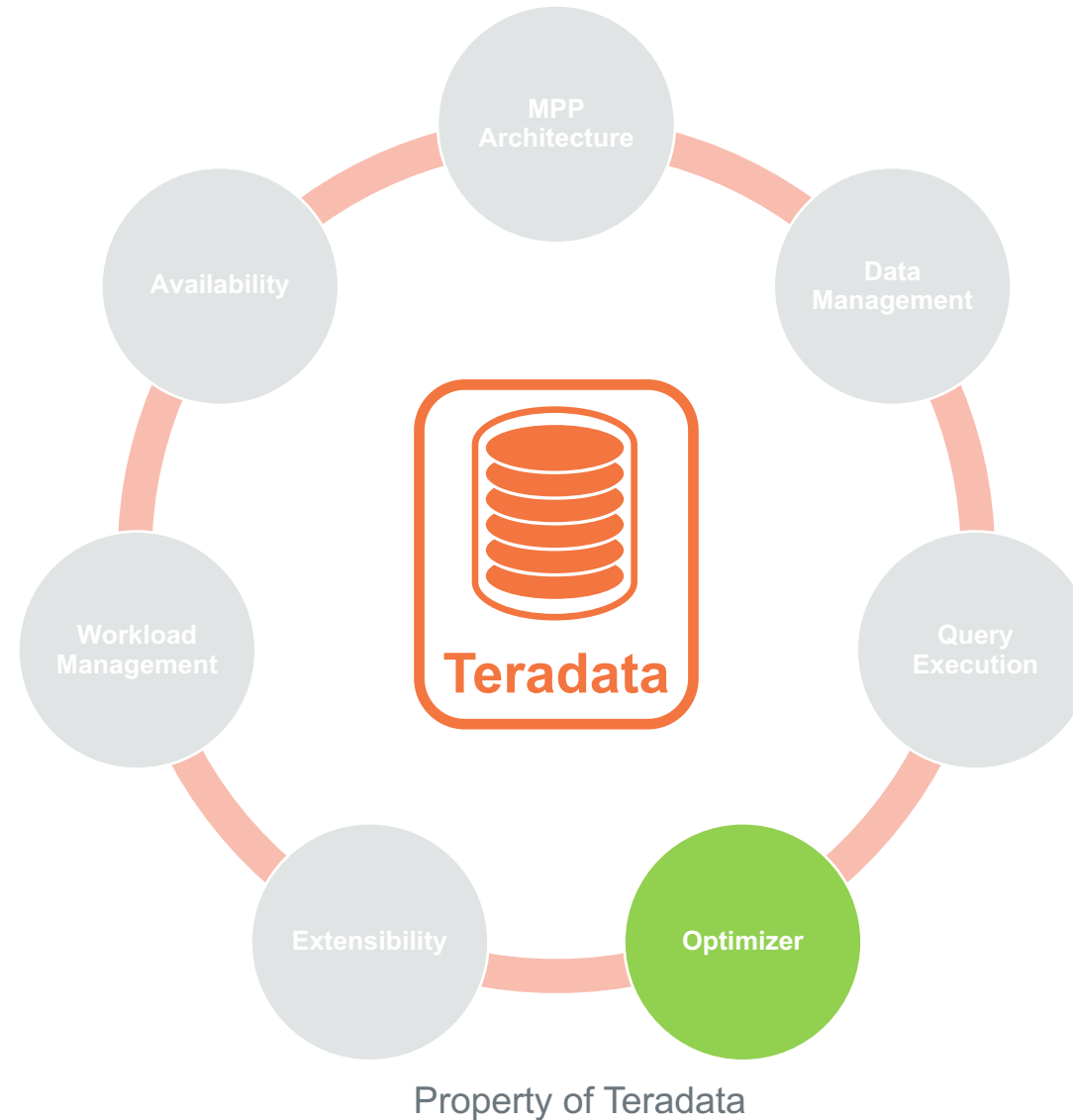
Merge of an Answer Set

- Each AMP sorts and spools its part of answer set
 - Query complete reported to application
 - Application performs fetch, interface (JDBC, ODBC,...) requests buffer of rows
 - Each node merges its contributing AMPs answer set into a buffer for the BYNET
 - BYNET merges the sorted buffers and fills a buffer to return to the interface
-
- All levels merge in parallel, and data is transferred using point-to-point messages, maintaining linear scalability as the size of the answer set or number of nodes increases
 - The "big sort" penalty has never existed ($N(\log(N))$ algorithm, N divided by # of AMPS)
 - Final answer set never has to be brought together into a single node for a large final sort – only merge a buffer at a time on demand



What Makes Teradata Unique

Key Components/Concepts



Teradata Optimizer Goals

Primary Goals of the Optimizer:

- In-database
- Minimum tuning efforts
 - No hints
 - Only indexes and statistics
- Use the platform effectively and efficiently
- Turn a SQL statement into a series of steps for execution

17095 Via Del Campo, San Diego, CA

San Diego Airport, North Harbor Drive, San Diego, CA

GET DIRECTIONS

Suggested routes

I-15 S and CA-163 S	25.8 mi, 30 mins
● In current traffic: 31 mins	
I-15 S	28.4 mi, 32 mins
● In current traffic: 33 mins	
I-15 S and I-5 S	31.1 mi, 35 mins
● In current traffic: 36 mins	

Driving directions to San Diego International Airport

17095 Via Del Campo
San Diego, CA 92127

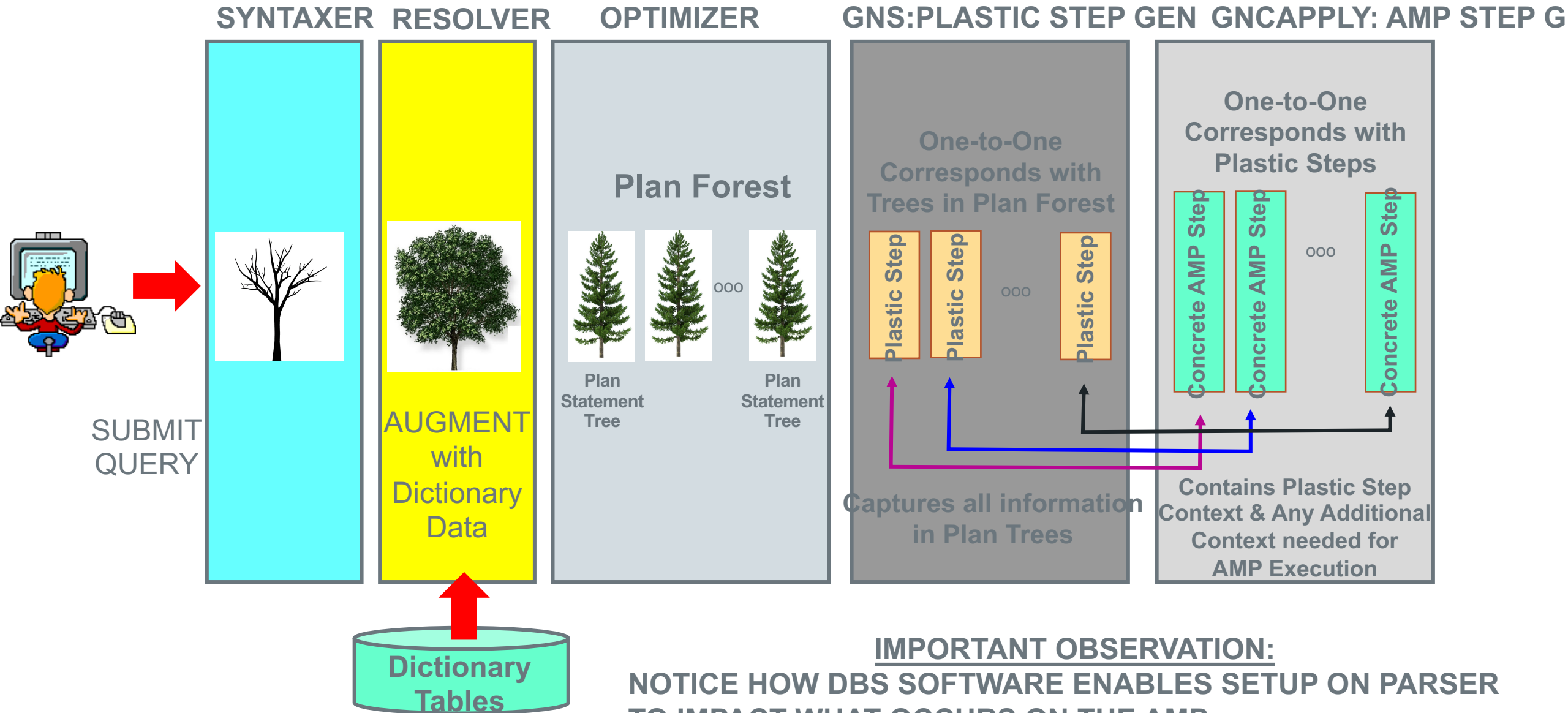
1. Head northeast 30 ft
2. Turn right toward Via Del Campo 236 ft

Running Queries on Teradata Database

- ANSI compliant SQL
 - Small number of Teradata specific analytic extensions
- High performance algorithms
 - Efficient use of resources not just counting on parallelism
 - Join, Aggregation, Sort, ...
 - Compiled expressions
- Complex query features
 - Derived tables, Case expressions, all forms of Sub-queries, Sample, ...
 - Big limits: 128 table joins, 128 nesting levels, ...
 - 1MB SQL/Views/Macros



CONCRETE STEP (AKA, AMP STEP) GENERATION

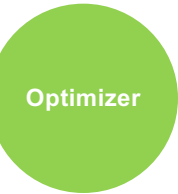


SUBMIT
QUERY

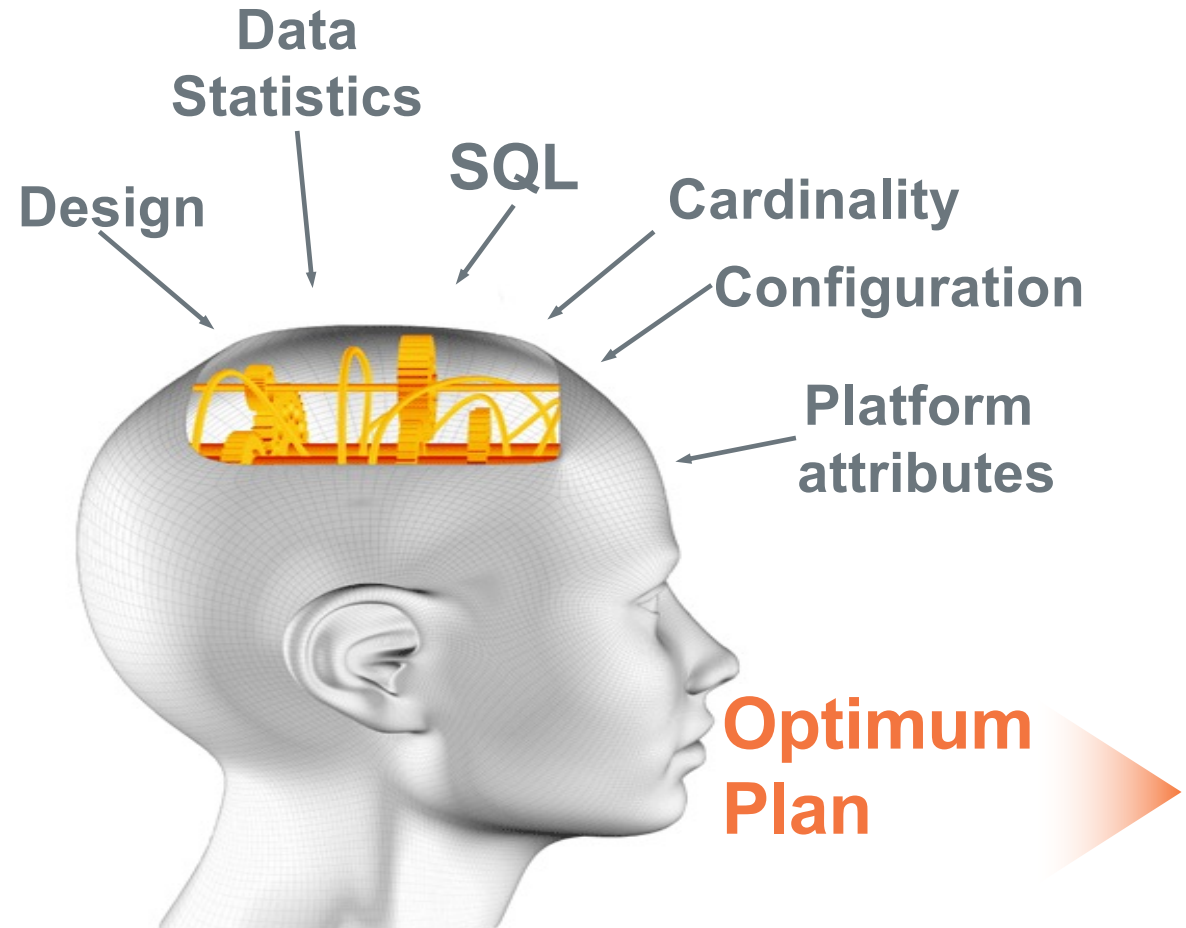
Dictionary
Tables

IMPORTANT OBSERVATION:
NOTICE HOW DBS SOFTWARE ENABLES SETUP ON PARSER
TO IMPACT WHAT OCCURS ON THE AMP

Teradata Optimizer

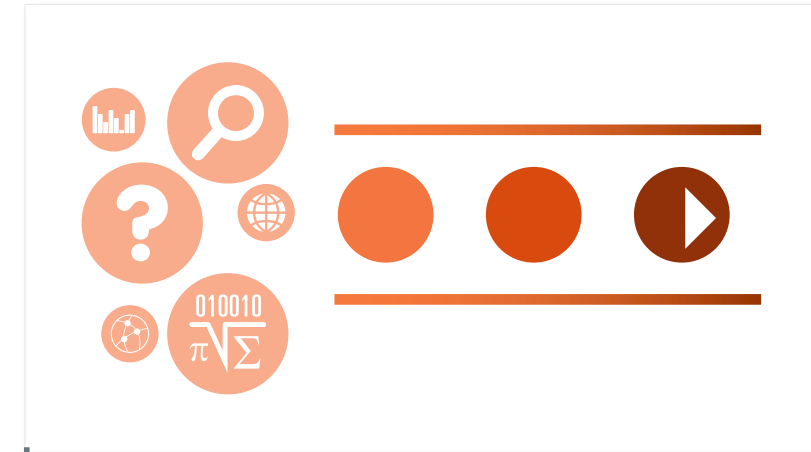


- Cost based query planning
 - Balance costs of
 - Cardinalities
 - Sorts and joins
 - Redistribution
 - Disk I/Os
 - Networks and nodes
 - Compression
 - Lowest resource costs
 - Fastest elapsed time
- Automatic rewrites
- Automatic optimization
 - No hints
 - No degrees of parallelism



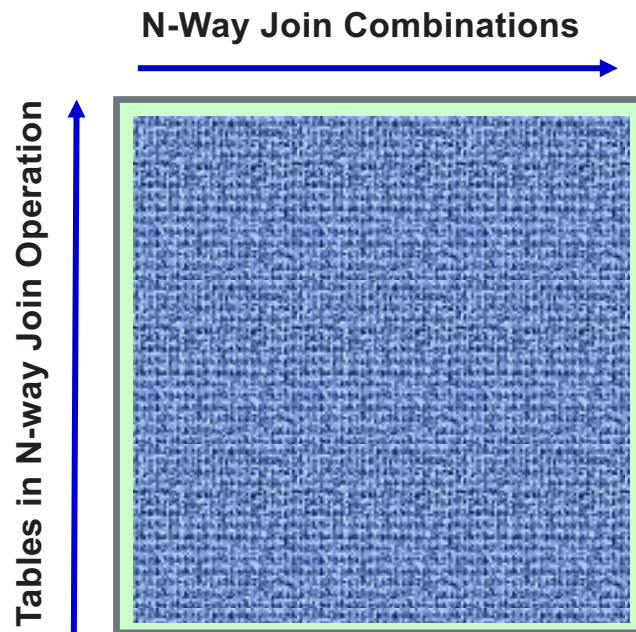
Planning

- Process of turning a query into a plan
- Optimizer deconstructs SQL into individual components
 - Relations, Join Conditions, Access Conditions
- Access Method: how to access each relation
 - Table Scan, Index Use, Bitmap Use
- Join Method: how pairs of table are joined
 - Merge Join, Product Join, Hash Join
 - Outer Joins, Inclusion and Exclusion Joins
- Join Order: sequence of table joins
- Join Geography: how rows are relocated prior to the join
 - Redistribute rows, duplicate rows, local spool
- Sort Order: order for merge join, ordered result, distinct

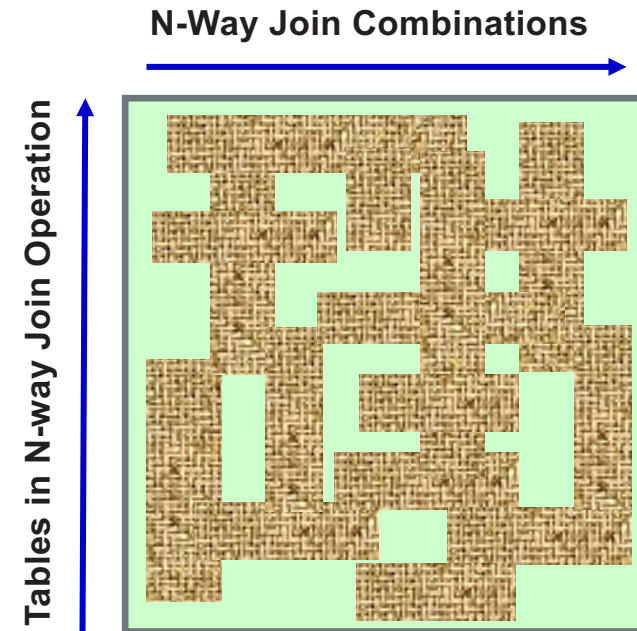


Join Planning

Classic System “R”
Exhaustive



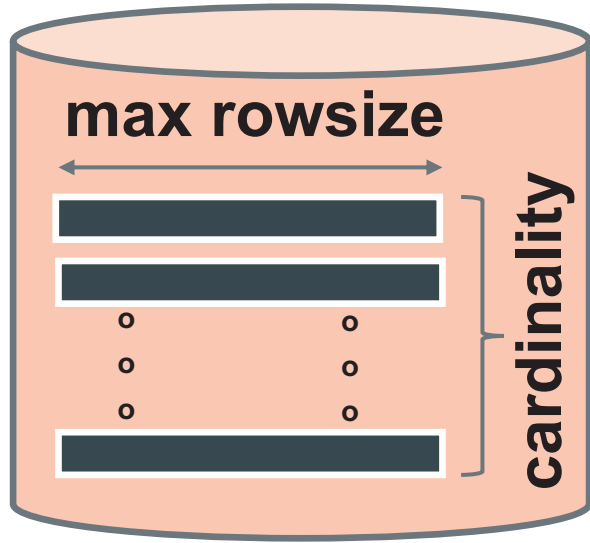
Teradata
Intelligent Plan Navigation



Optimizer Plan Selection: Comes Down to “Cost”

Cost \equiv *Time to complete an operation in ms*

Costing Building Blocks



Allows us to
Calculate



OPT PLAN COST

- ΔT : Disk I/O cost
- ΔT : Network (Redistrib) Cost
- ΔT : CPU Processing Cost

+



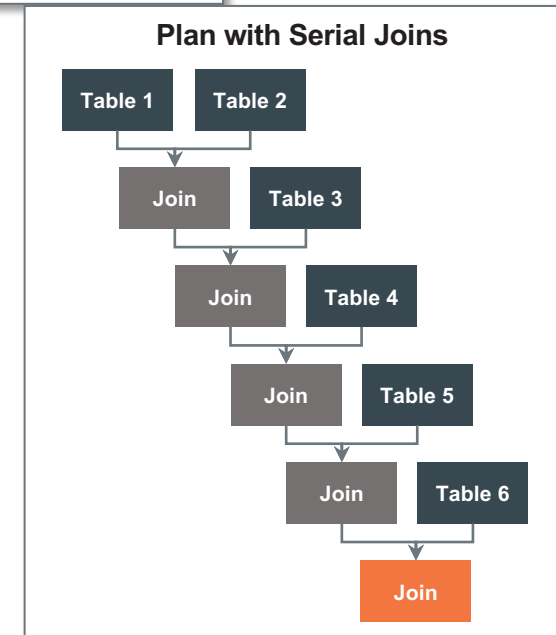
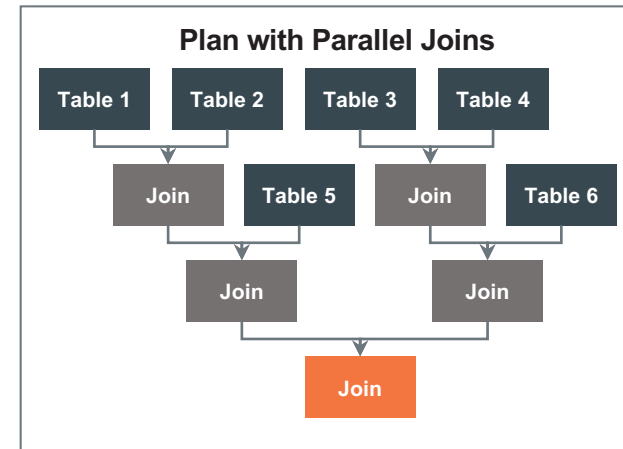
CPU
cycles-per-row
operation

**“Cardinality is the “hardest one”
This is where “Statistics”
comes into play**

Join Order

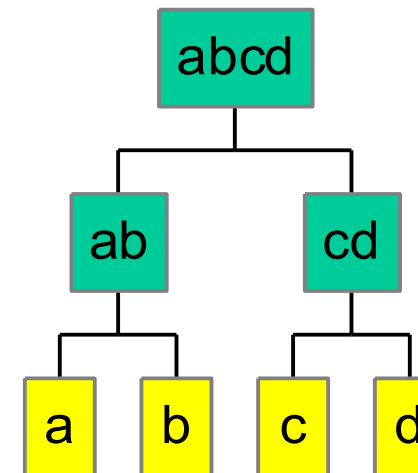
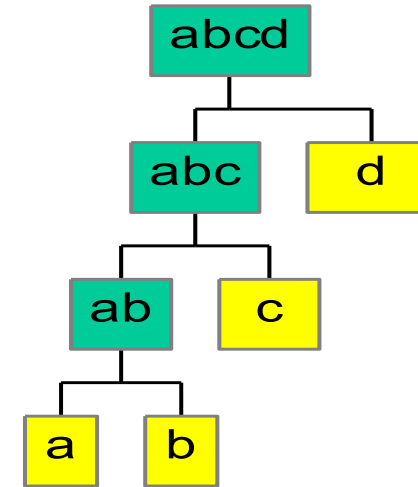


- The order of joins is determined by identifying the minimum cost ordering
- Look at size (with conditions applied), distribution, available access methods, available join methods, ...
- Create each order and determine cost, pick best
- Thousands or millions of combinations
 - Game theory used to eliminate unlikely combinations
 - Look-ahead limits to control reviewed combinations
 - Heuristics for special cases, starting points
- Dependent on good Statistics!
- And the Optimizer makes the decision – does the hard work of figuring it out

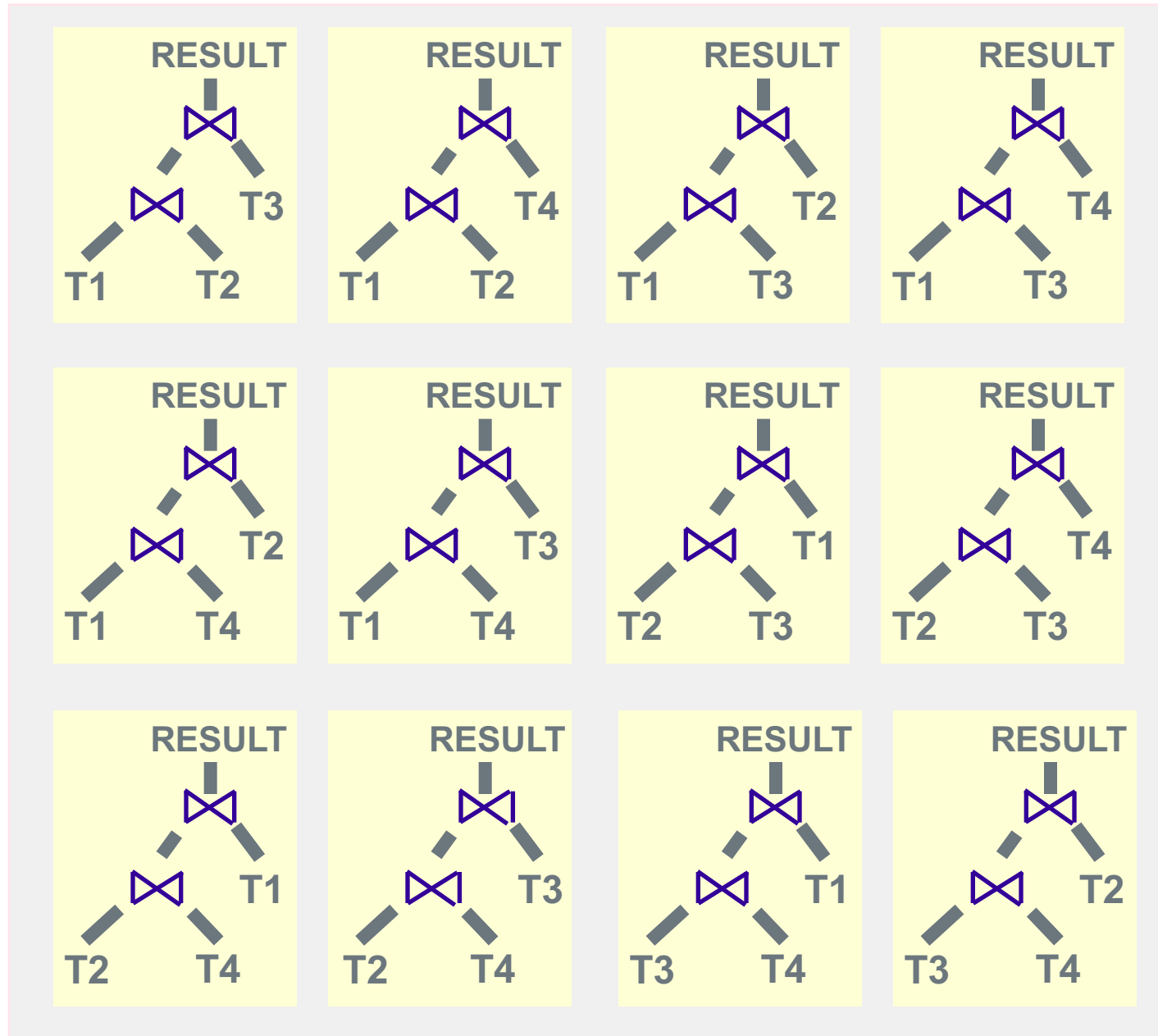


N-Way Join - Combinations

- Sequence of binary joins
- Not all possible combinations considered
 - 10 way join - 17.6 billion possibilities
 - 64 way join - $1.2 * 10^{124}$ possibilities
- Driven by join conditions
 - Look for “connected” relations
 - “Greedy” algorithm

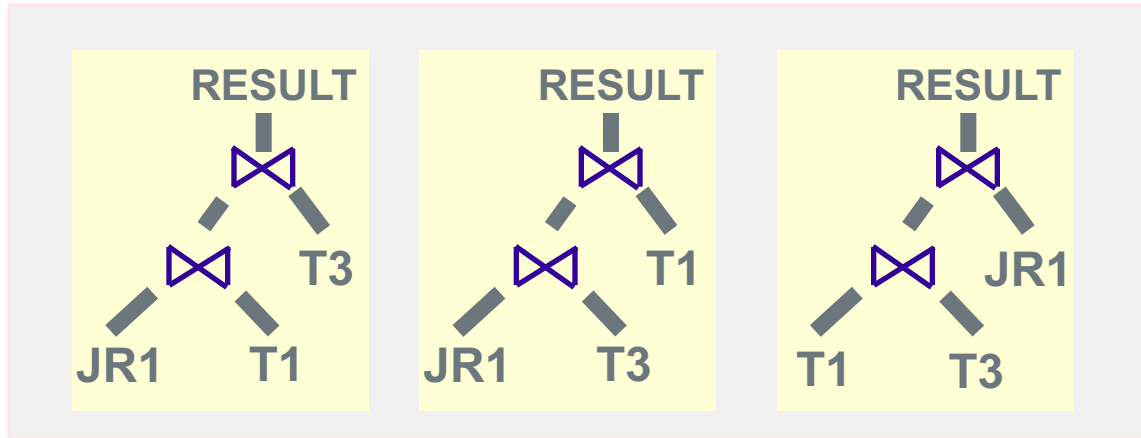


4-Way Join : LookaheadOne Join Planning

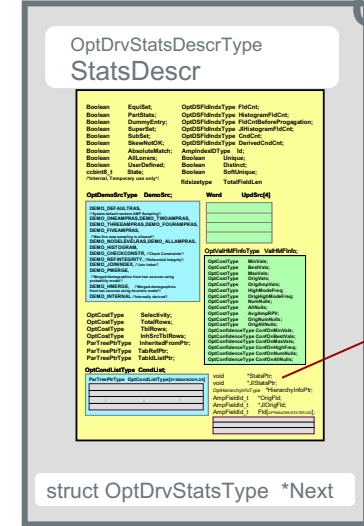


4-Way Join : LookaheadOne Join Planning

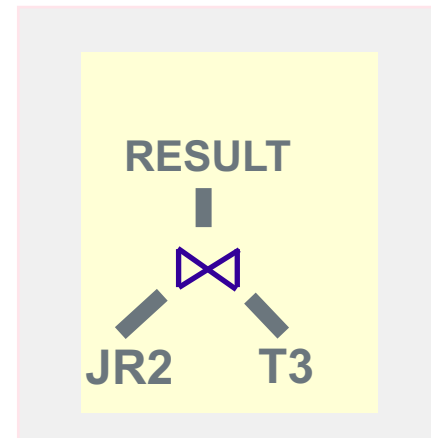
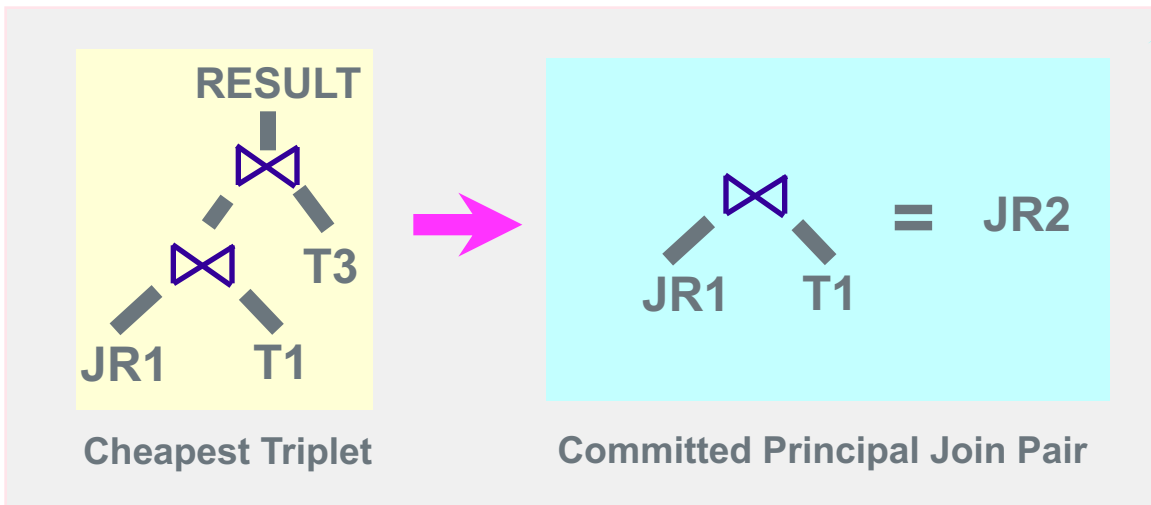
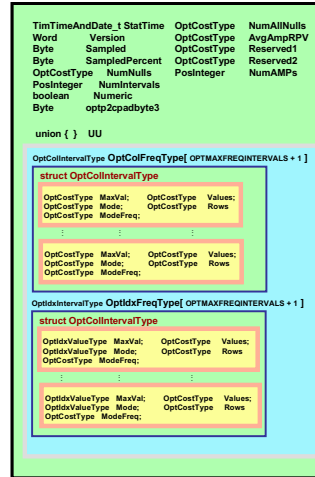
Need to Estimate Cardinality of Committed Pair



struct OptDrvStatsType



struct OptStatisticsType



Binary Join Generation

- Try all possible join types based on join conditions
 - Equality: nested join, merge join, hash join, product join
 - Other: product join
- For each join type, try all possible combinations of geography of tables
- Examples for Merge Join

direct left	direct right
hash left	hash right
hash left	local right
local left	hash right
local left	duplicate right
duplicate left	local right

Views and Physical Database Design

High performance views:

- No penalty for using views
- Enables wide use for ease of use, security, privacy
- Recommend using views rather than users accessing the base tables
- Compiled early so optimizer only sees base tables and complete condition sets

Physical database design affects planning:

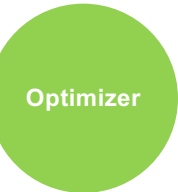
- Indexes, Join Indexes, partitions – provide access methods
- PI, Referential Integrity, data types – guide join planning
- Columnar, Compression, row definition – guide size and I/O estimates

Examples of Optimizer's Advanced Capabilities



- Aggregate Join Index (AJI)
- SATisfiability-Transitive Closure (SAT/TC)
- Query Re-Writes
 - Reorganize a query's structure for better optimization/performance
- Partial Redistribution Partial Duplication (PRPD)
 - Skew optimization
- Incremental Planning and Execution (IPE)
 - Plan in stages based on partial results
- In-Memory and Chipset Optimizations
- Temporal sequencing/normalization
- Geospatial indexing/bounding

Statistics



Three types of histograms are generated:

- Equal-height interval histograms:
 - Each interval has the same number of values
- High-biased interval histograms:
 - Each interval has no more than two values
- Compression histograms:
 - Up to 200 equal-height intervals, plus one or more high-biased intervals

**COLLECT STATISTICS ON Employee
COLUMN Department_Number**

26 Rows		Min value = 100			
8 Unique values,		Mode Value = 401			
0 Nulls		Mode Frequency = 7			
4 Intervals (in High-biased interval histogram)					
	Max	Mode	Mode Frequency	Non-Mode Value	Rows
(1)	201	100	1	-2	2
(2)	302	301	3	-2	1
(3)	402	401	7	-2	2
(4)	501	403	6	-2	4

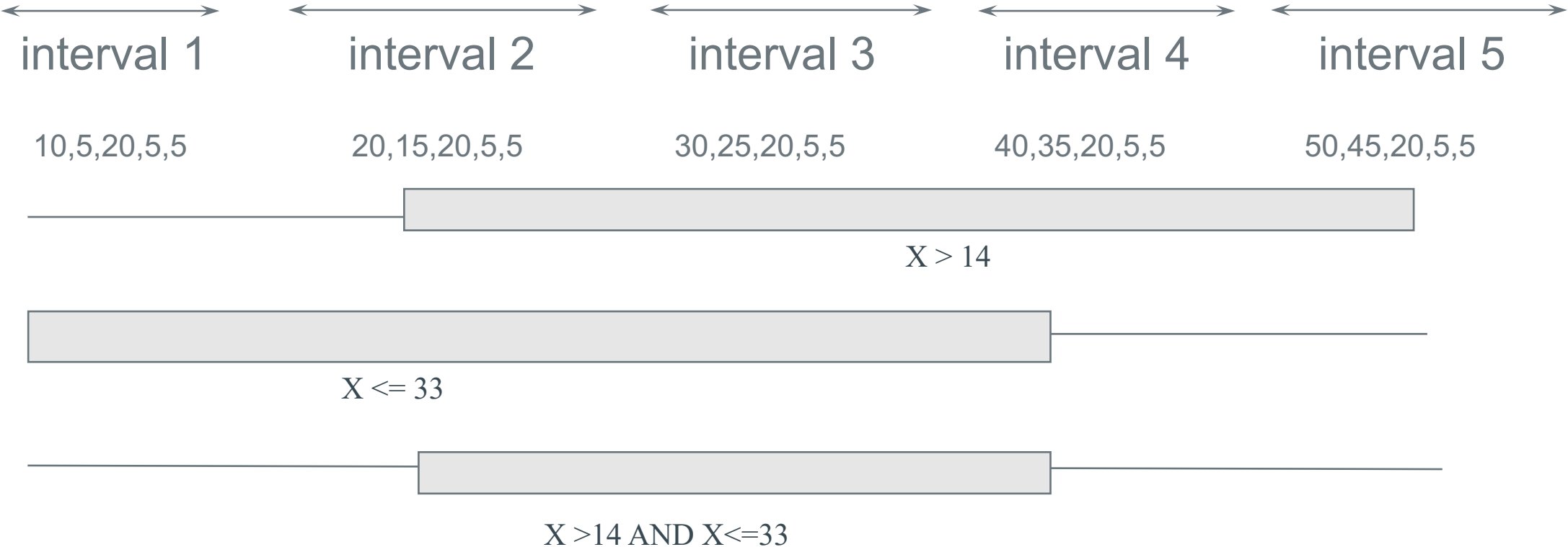
Statistics Collection

- High bias values separated to recognize skew
 - Optimizer will generate different plans for different values or ranges
- Statistics calculated across all rows on all AMPs with a fully scalable algorithm
- Statistics sampled when not available, dynamic choice of single AMP or all AMPs
- Statistics extrapolated when stale relative to data change and table size relative to when statistics were collected

Using Statistics to Estimate Number of Rows (cont.)

- Interval Information

Max, Mode, Frequency, Number of values \neq Mode, Number of rows \neq Mode

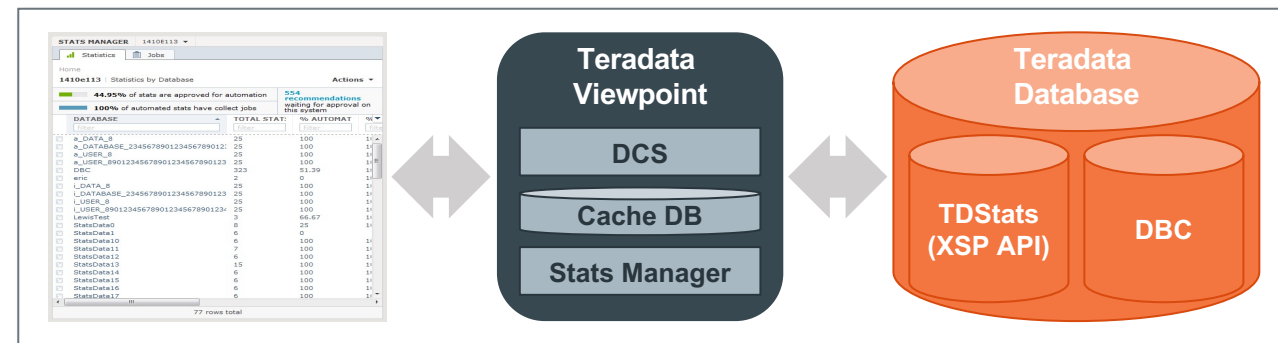


Auto-Stats Collection



Automates and provides intelligence to DBA tasks related to Optimizer Statistics Collections where such tasks include:

- Identifying and collecting missing statistics that are needed for query optimization
- Detecting stale statistics and promptly refreshing them
- Identifying and removing unused or unimportant statistics from routine maintenance
- Prioritizing the list of pending collections such that important and stale statistics are given precedence
- Executing needed collections in the background during scheduled time periods
- Dynamically issuing collections in response to key events, most notably the completion of bulk load operations
- Configuring the system to ensure that resource usage incurred by statistics collections is properly regulated and throttled through TASM or other mechanisms



Explain

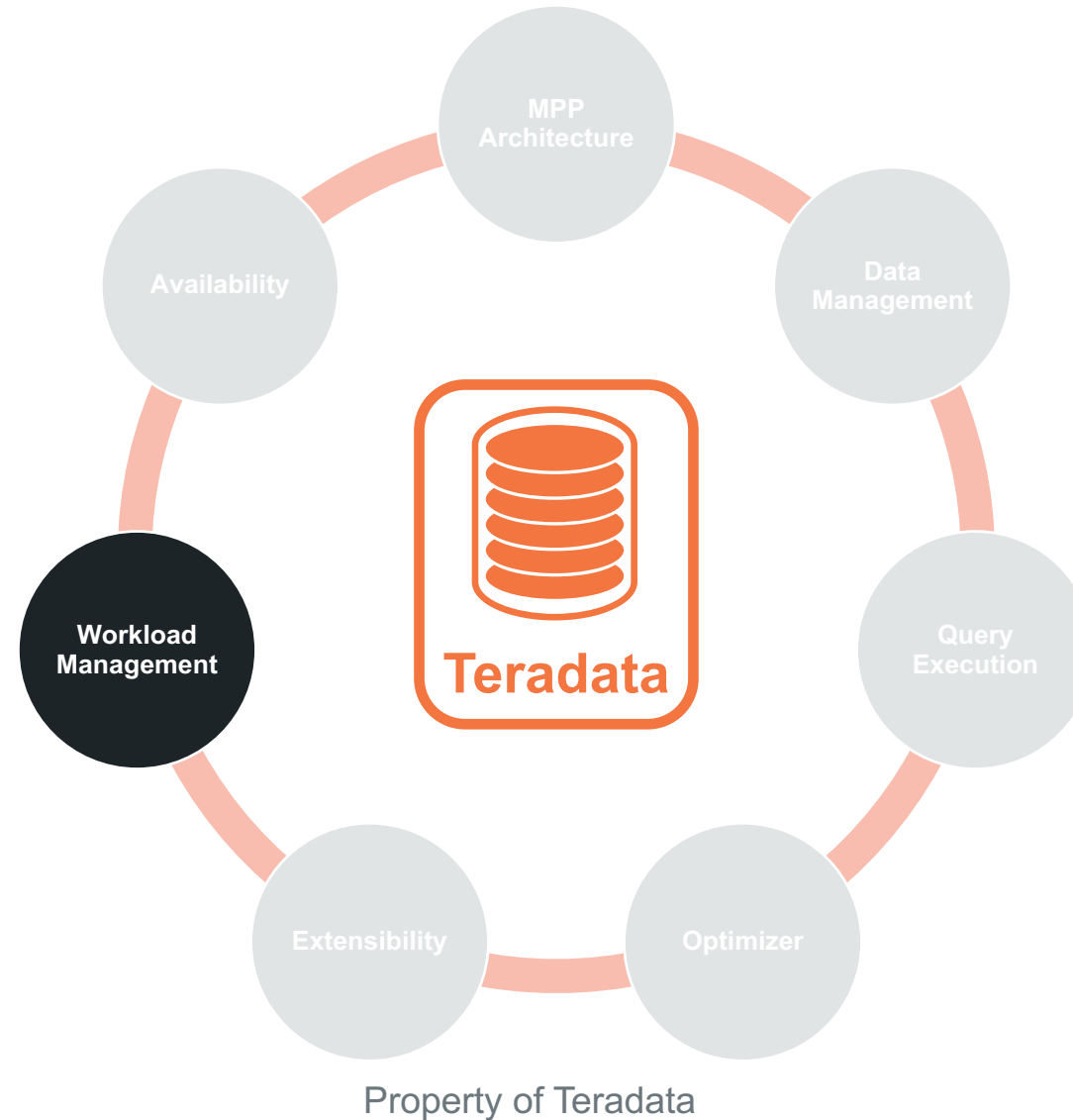
- Two table join with NUSI access in one table and FK to FK join. and aggregation and sort

BLOG

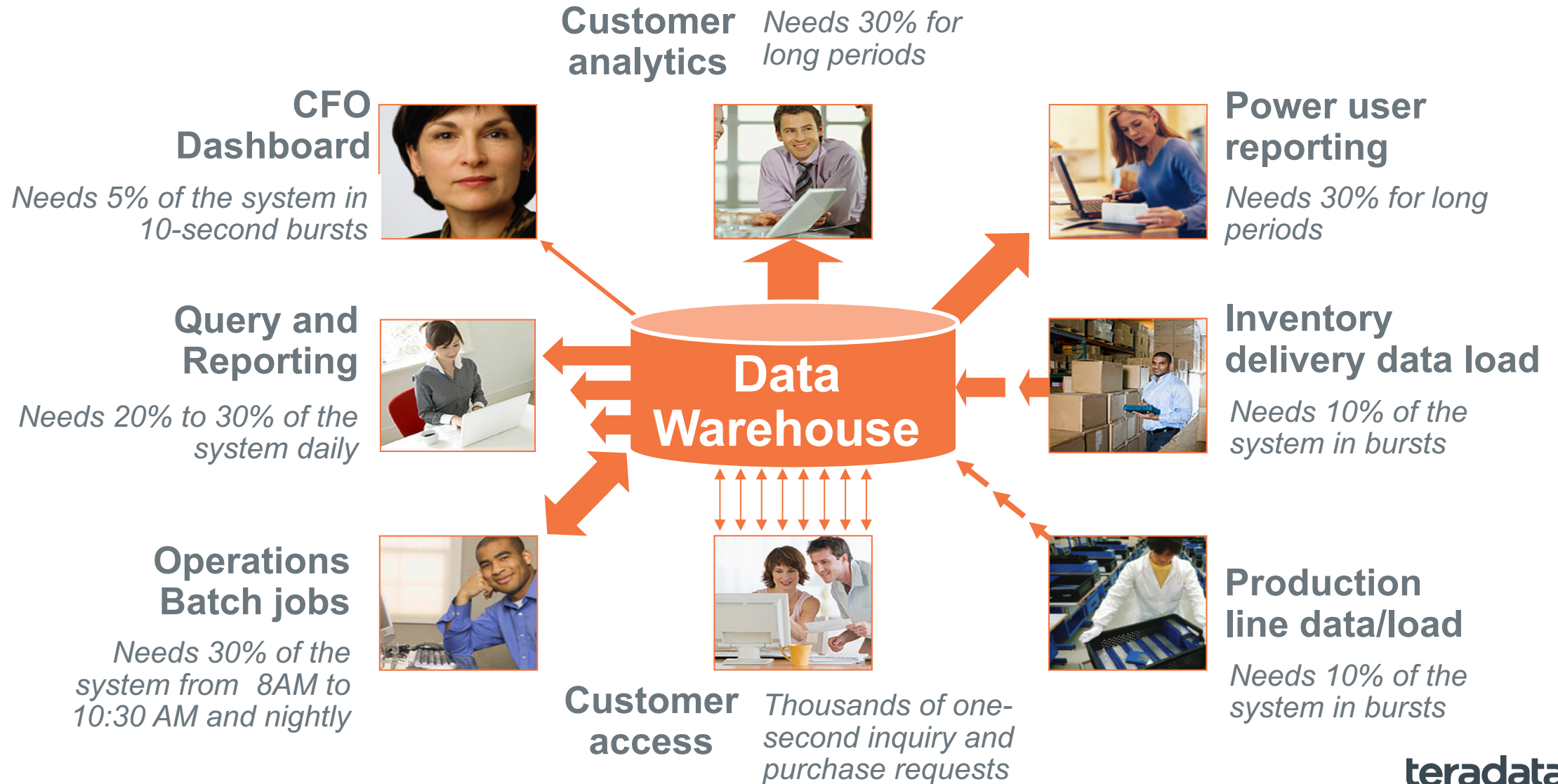
- <https://www.teradata.com/Blogs/How-to-Repurpose-Successful-Database-Techniques-inside-Teradata-Vantage>

What Makes Teradata Unique

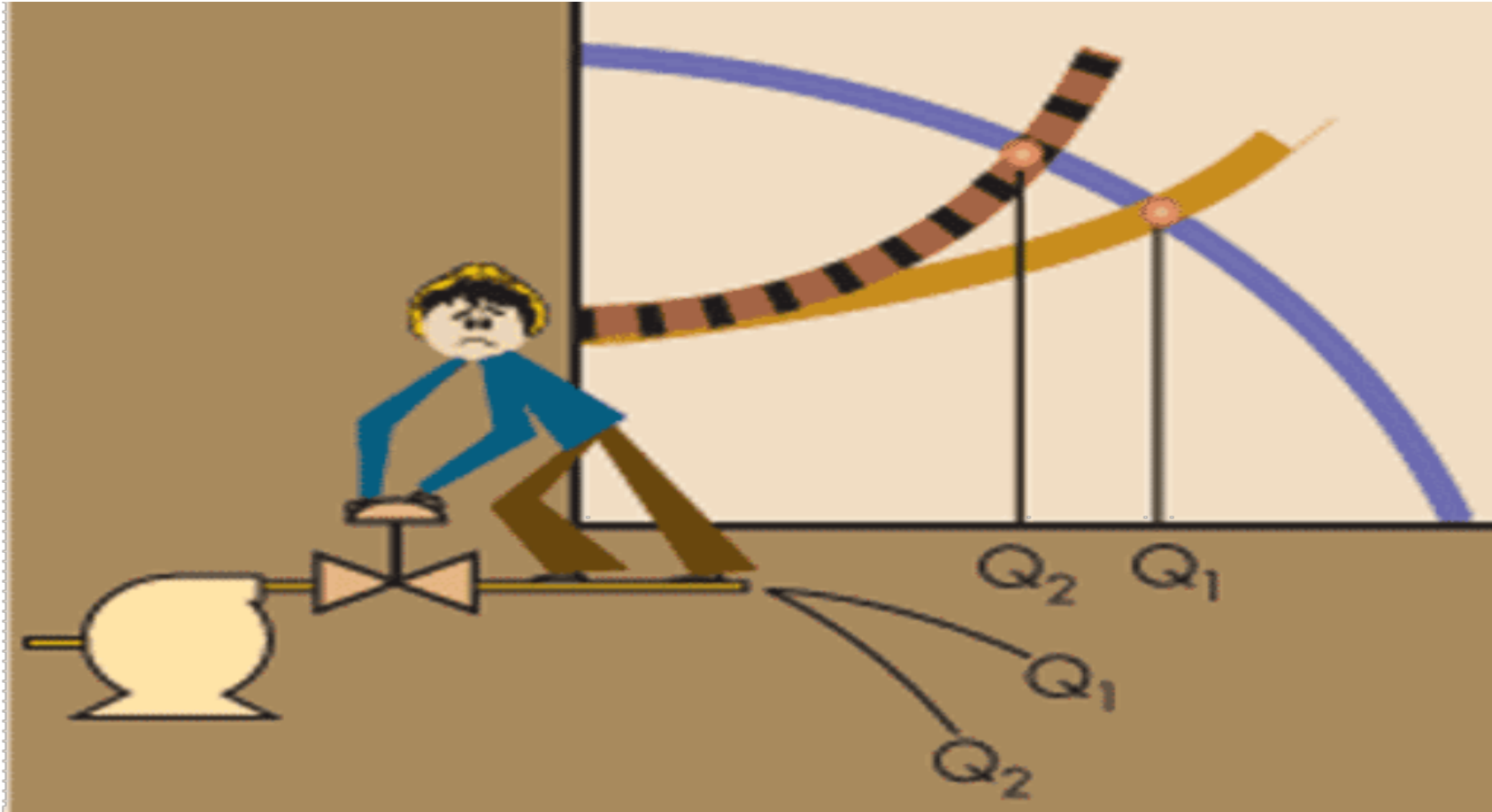
Key Components/Concepts



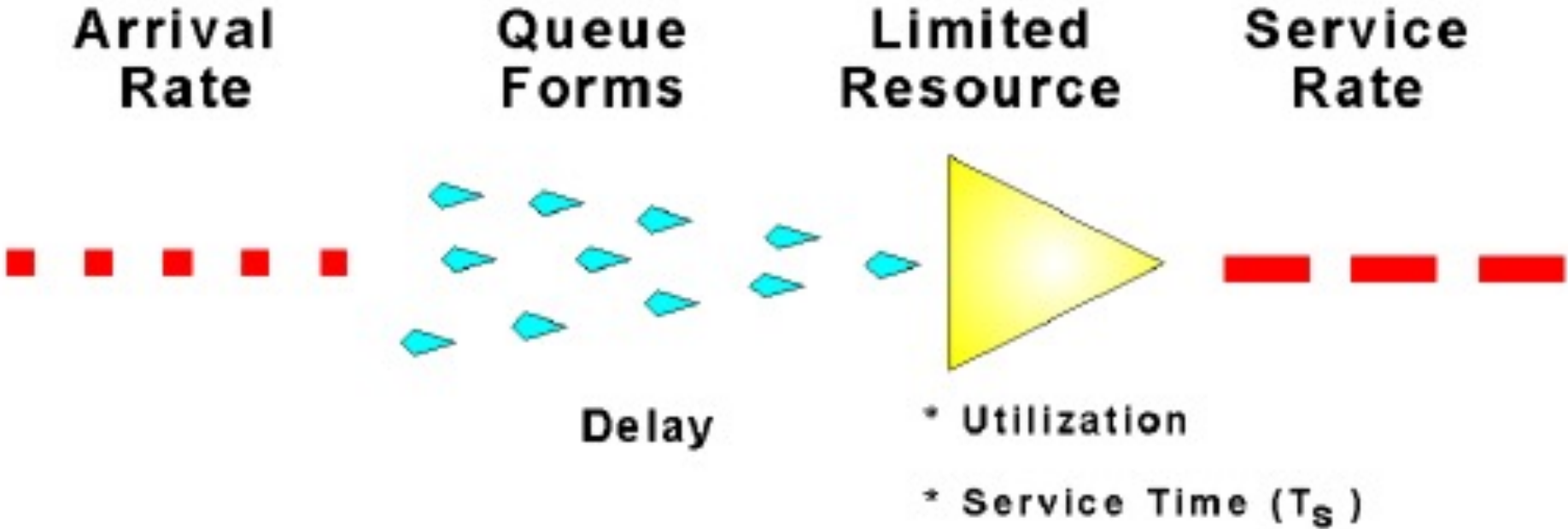
Many Demanding Users



Basic Workload Management Controls

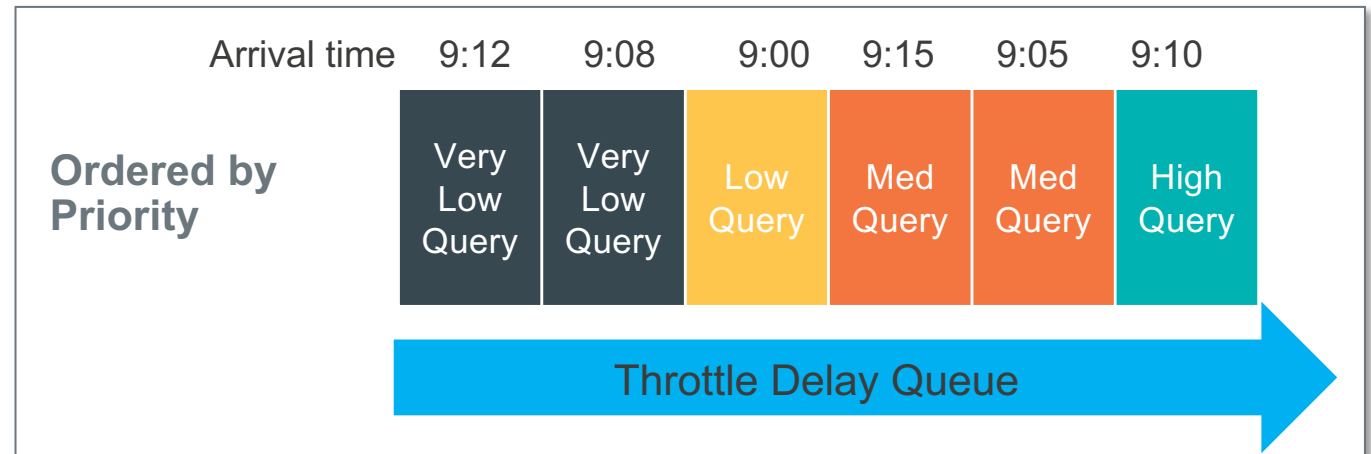
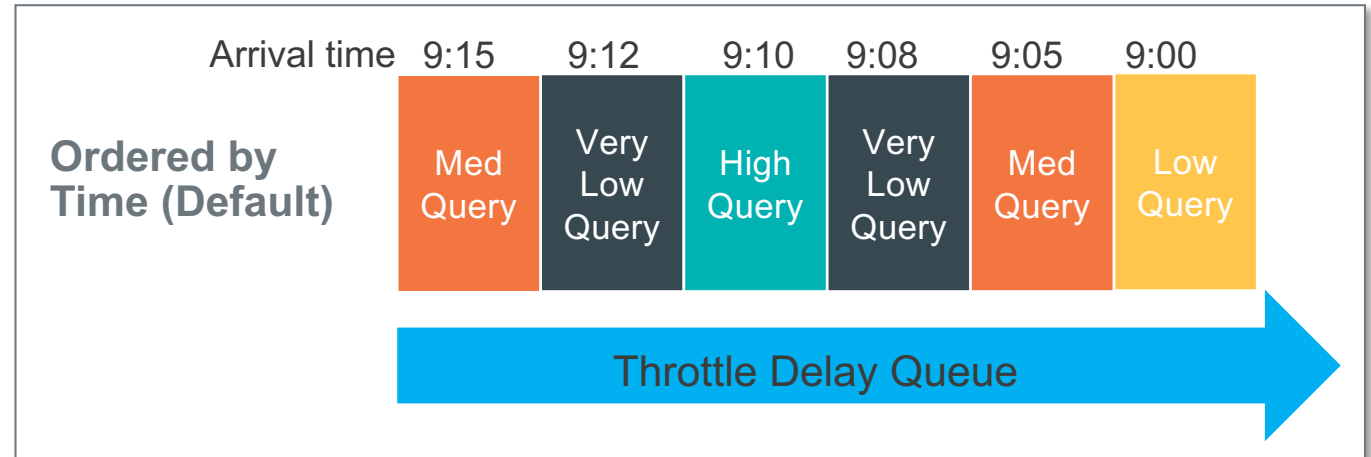


Workload Management Queuing



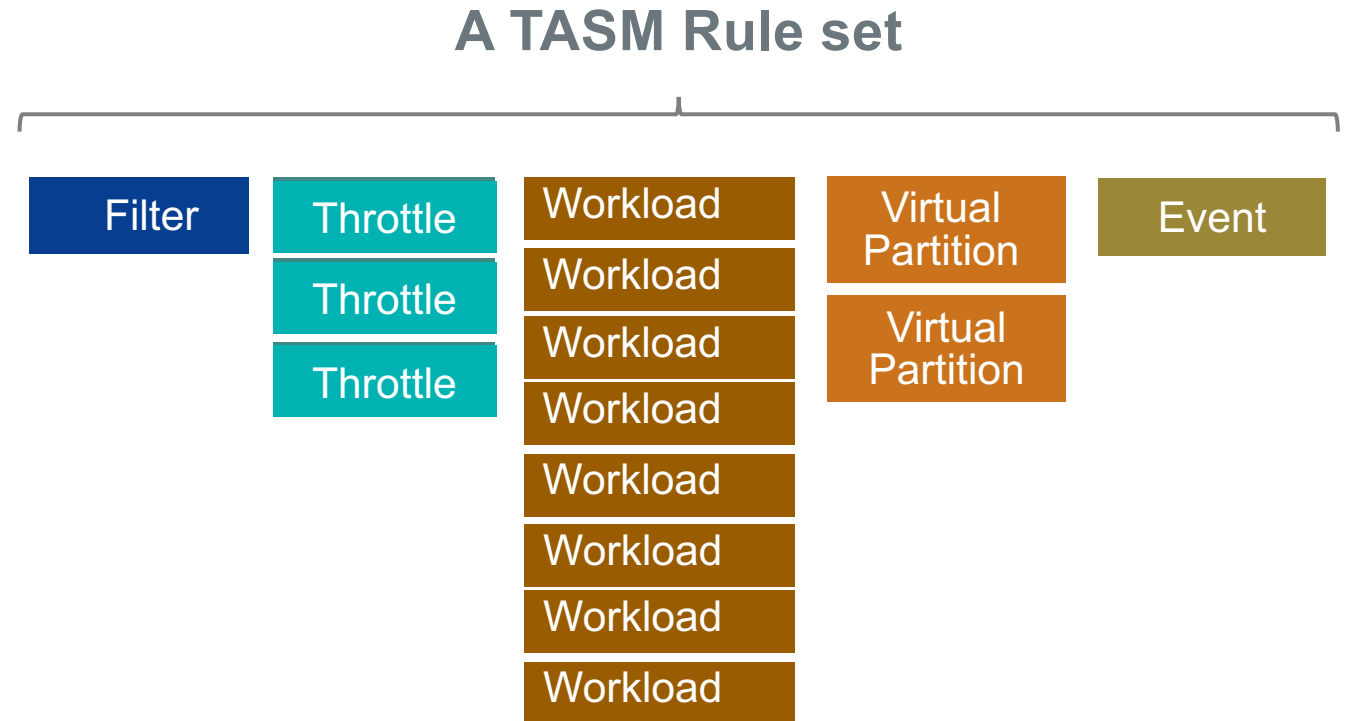
Prioritized Throttle Delay Queue Option

- Default delay queue ordering is first-in-first-out (FIFO)
- New option to reorder the delay queue by priority
- Can be used to make sure the highest priority requests are released first

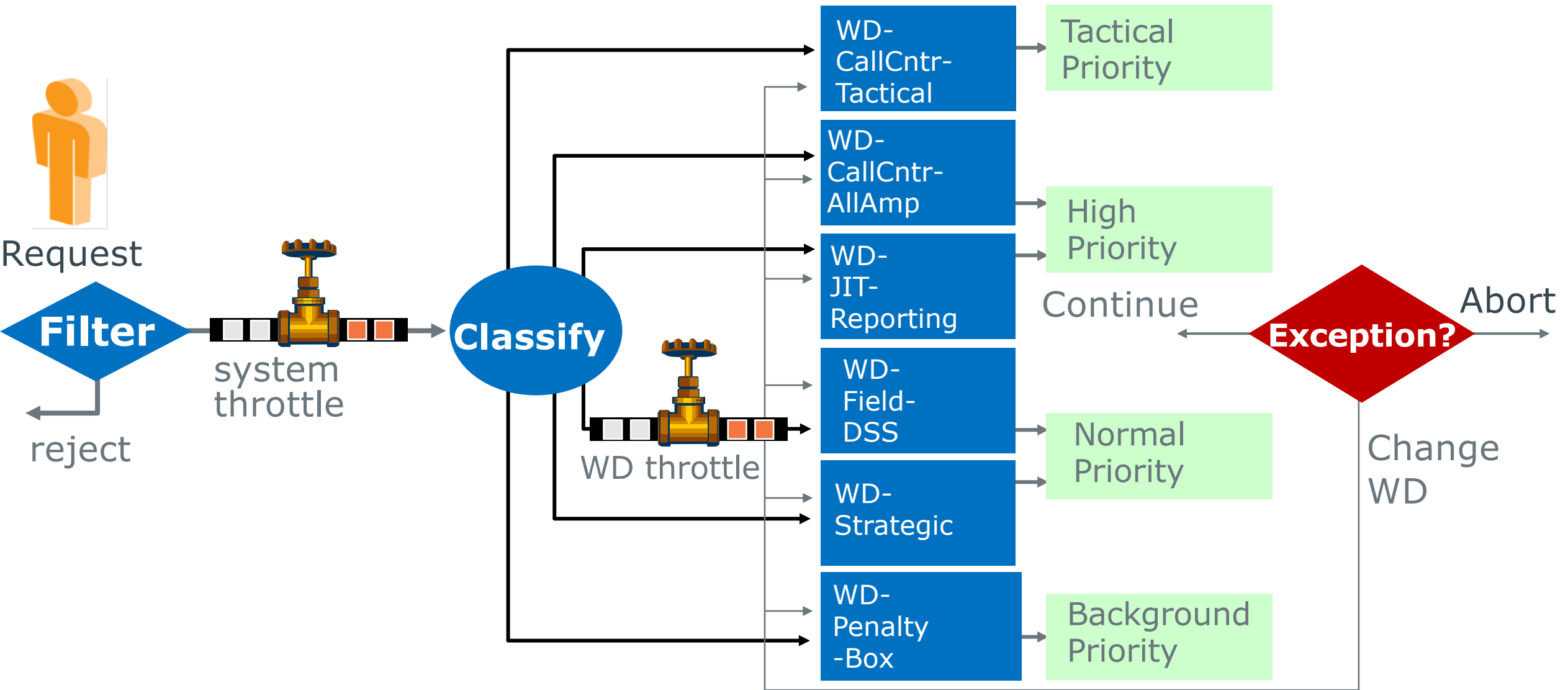


How Workload Management Goals are Defined

- A rule set is a bundle of settings and definitions
- Created within Viewpoint Workload Designer
- Only one rule set may be active at a time



TASM: Up to 6 Methods of Management



A Wide Variety of Workload Classification Criteria

Requests are mapped to workloads, system throttles, or system filters if the query's characteristics match the object's classification criteria

WHO Criteria

- User
- Account
- ClientID
- ClientAddr
- Profile
- Application
- QueryBand

“Source” criteria

WHERE Criteria

- Tables
- Databases
- Views
- Stored procedure
- Functions/methods
- QueryGrid server

“Target” criteria

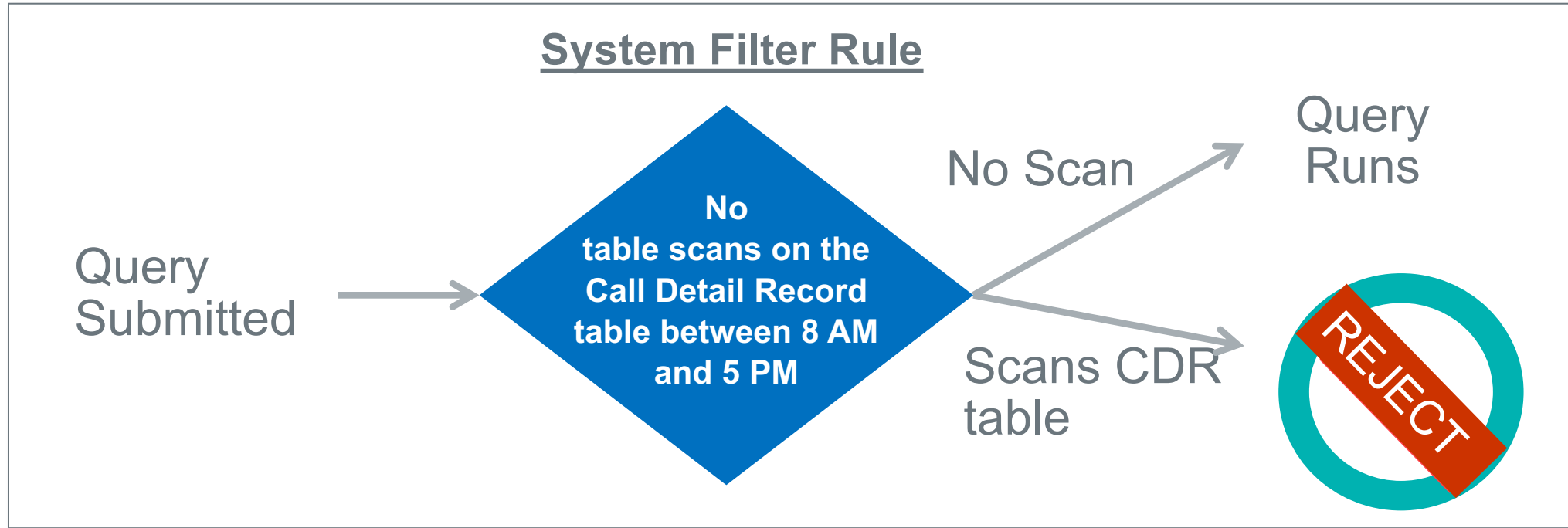
WHAT Criteria

- Estimated time
- Type of join
- Estimated rows
- AMP limits
- Load utility type
- Statement type
- Final row count
- Percent of table accessed

“Query characteristics”

System Filter Rules

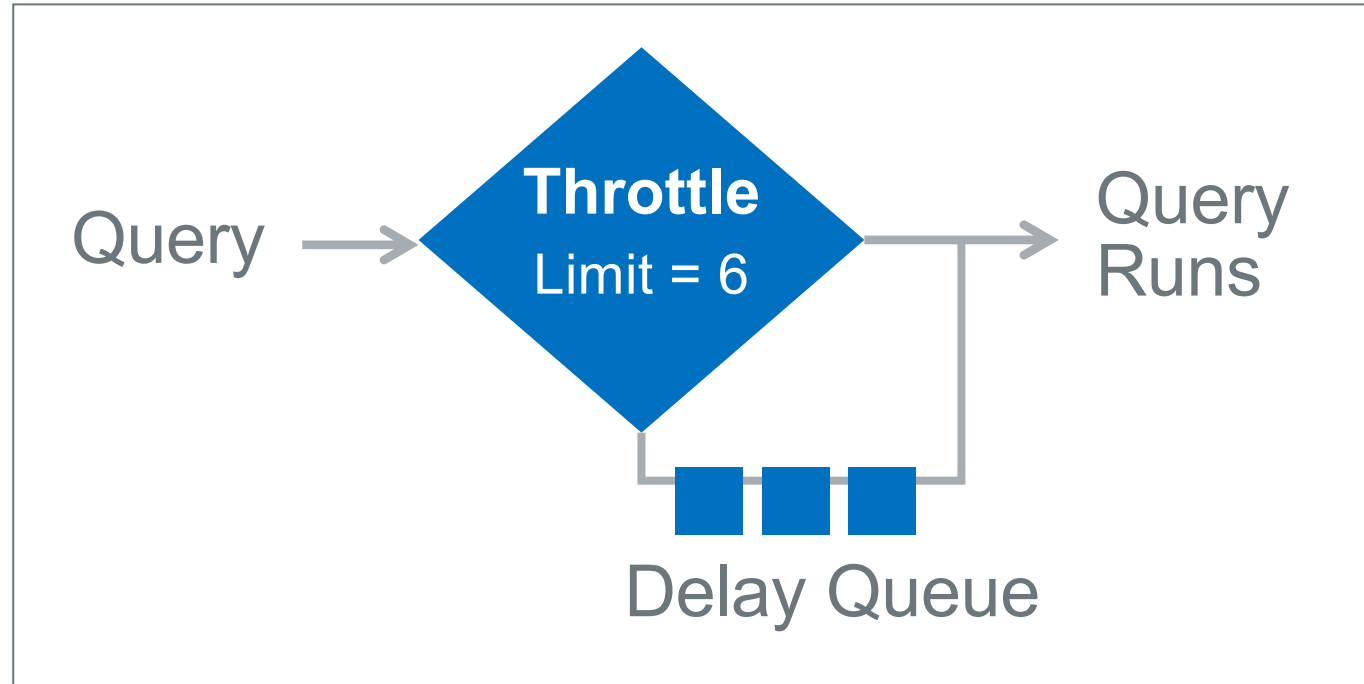
Reject Queries with Specific characteristics, Before They Begin Running



- Prevents inappropriate queries from executing
- Can contribute to user education
- “Warning mode” will report, but not reject the query

Throttle Rules

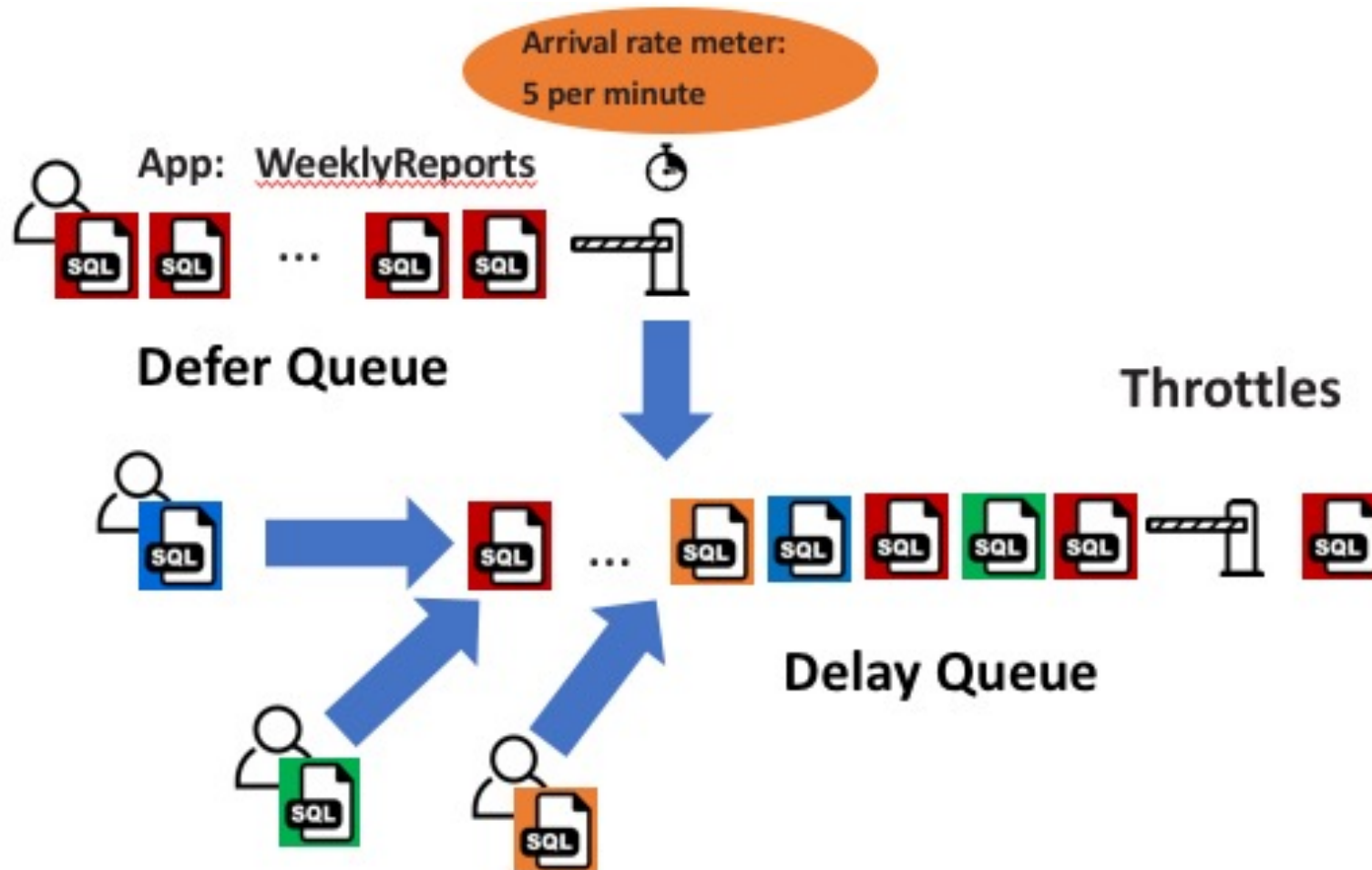
Soften Processing Peaks by Delaying Some Queries at Busy Times



- You select a query limit for each throttle rule
- A counter is kept to ensure the query limit is not exceeded
- Queries that would exceed the throttle limit are placed in a delay queue

Arrival Rate Meter – New in TD17.10

Create an Arrival Rate Meter for this application so that other requests can get in between each time unit



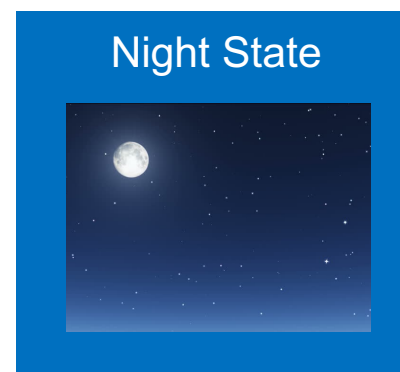
System Events Trigger State Changes

- A “State” encapsulates a set of “working values”
- Different TASM “states” can be defined based on:
 - Time periods or day-of-the-week
 - Level of system busy-ness (CPU utilization levels, for example)
 - When a particular type of processing completes (such as nightly loads)
- **A move from one state to another triggers setup definition changes**



*Online apps get
high priority*

*Batch work gets
low priority*



*Online apps get
low priority*

*Batch work gets
high priority*

What is a System Event?

- Triggers a change in state
- A wide variety of events are available
- Most highly-used system events:
 - AMP worker task availability
 - System CPU usage levels
 - New: I/O usage event

Settings that can be changed:

Workload priorities
Throttle query limits
Filter rules on/off
Workload Management Capacity
on Demand

Sample of system-wide events:

Excessive use of CPU
Start/End of load window
Hardware component disabled
Workload arrival rates
Service level goals being missed
Delay queue depth or time in the delay queue

Resource Usage Data

ResUsageSpma

- Contains data for system-wide information by node

ResUsageSvpr

- Contains data for system-wide virtual processor information

ResUsageScpu

- Contains data specific to the CPUs within the Nodes

ResUsageShst

- Contains data specific to the host channels and TCP/IP networks communicating with the Teradata database

ResUsageSps

- Contains data regarding workload behavior statistics for utilities and SQL operations

ResUsageSawt

- Contains data regarding Amp Worker Task (AWT) behavior statistics for utilities and SQL operations

ResUsageIpma

- Contains data for system-wide information by node
- Generally used by Teradata engineers

ResUsageIvpr

- Contains data for system-wide virtual processor information
- Generally used by Teradata engineers

ResUsageSldv

- Contains data for system-wide logical device data collected from storage devices

ResUsageSpdsk

- Contains data specific to physical disk storage statistics

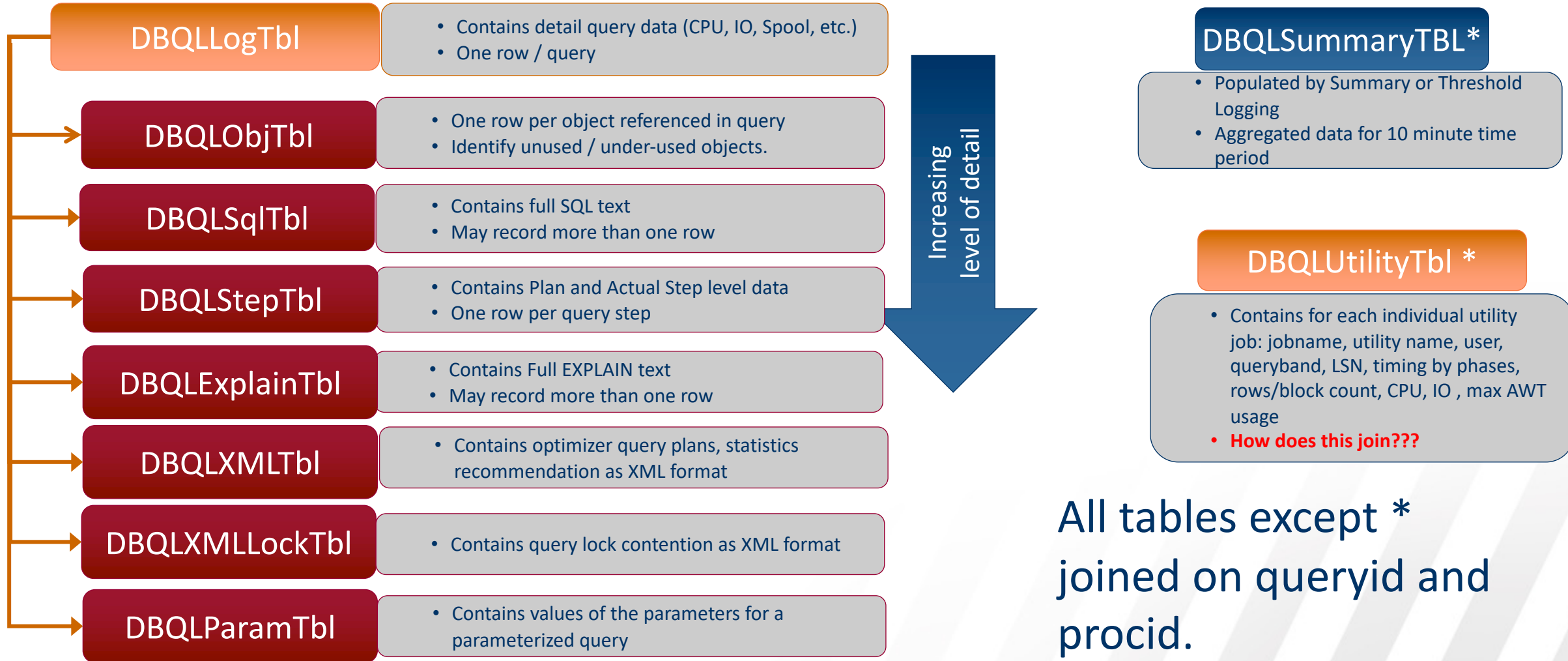
ResUsageSvdsk

- Contains data specific to amp-level disk storage statistics

Collection and Logging Standard

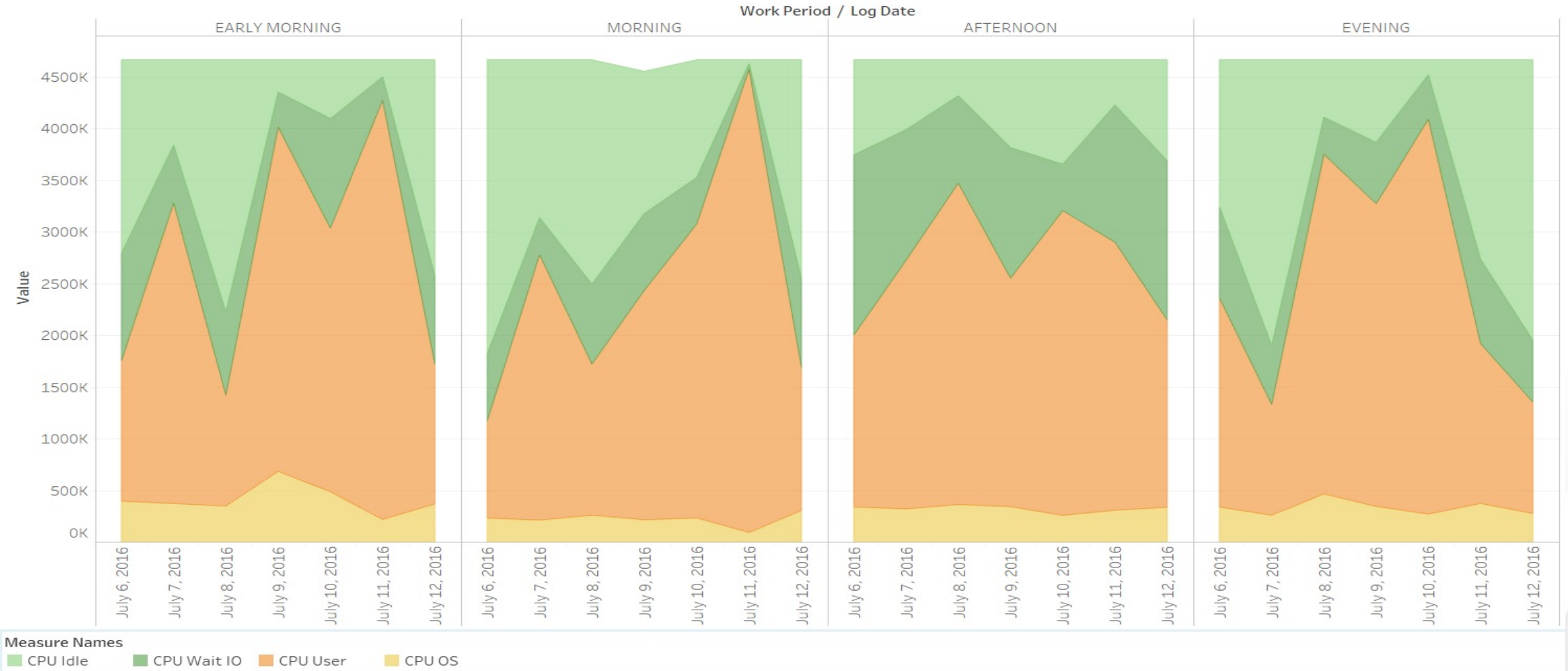
- ✓ 60 - 600 second collection interval
- ✓ 60 - 600 second logging interval
- ✓ Use Active Filter Enable for AWT and SPS
- ✓ Use Summary mode for vdsk if enabled

DBQL Data

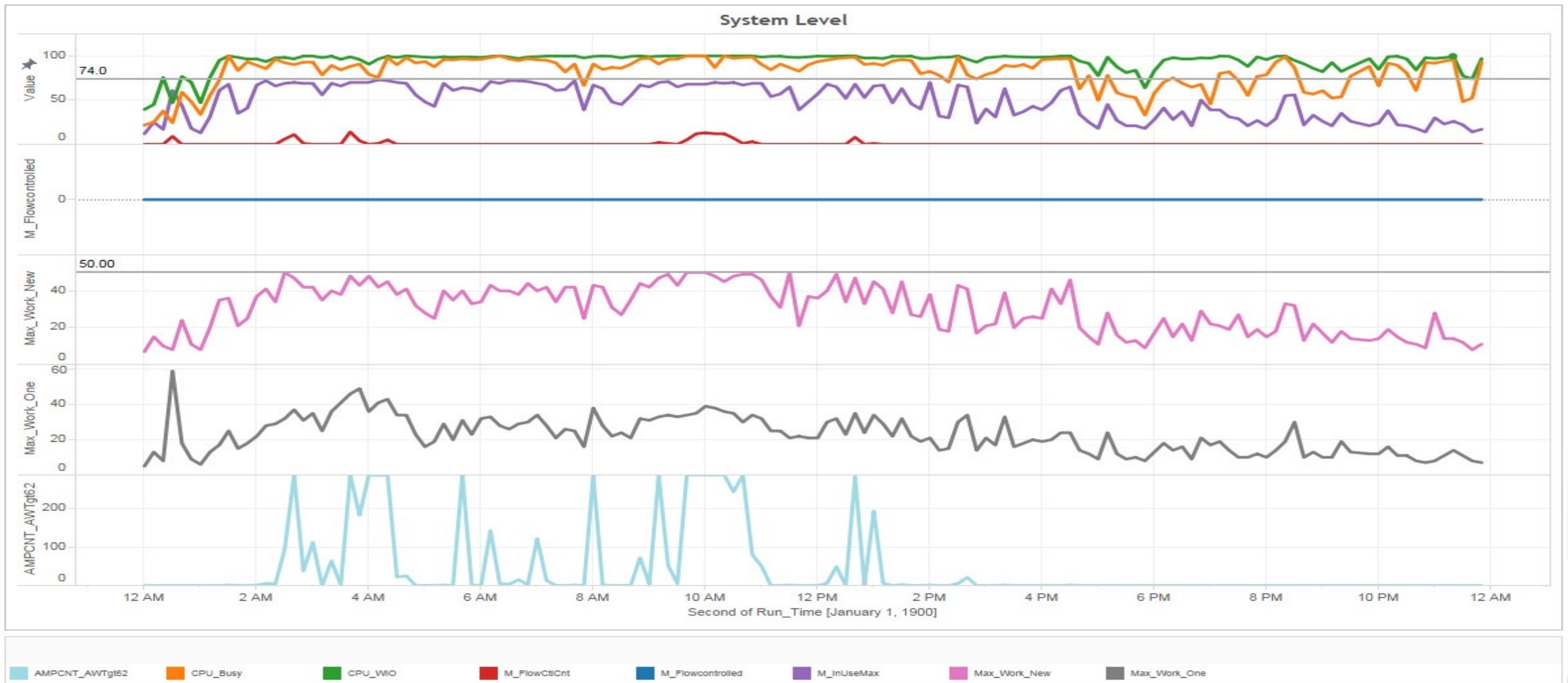


System Performance Sample

CPUChart

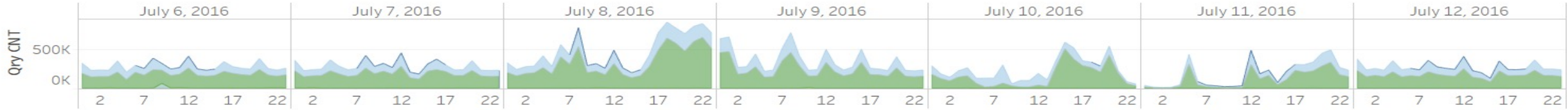


System Performance Sample



Application Performance Sample

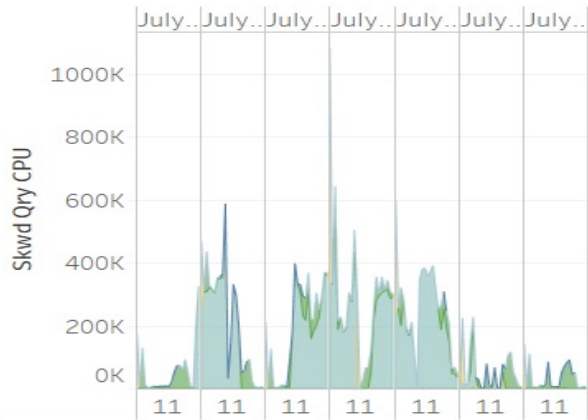
QueryCount



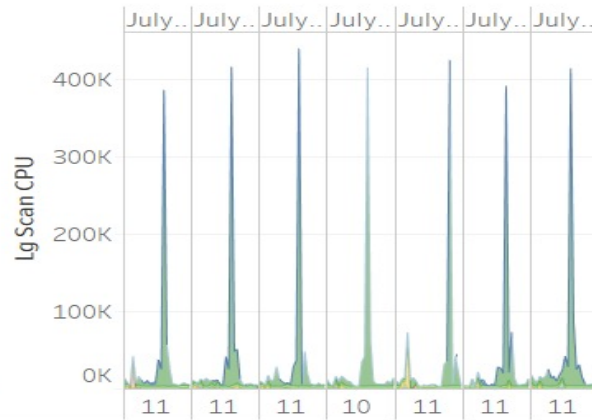
CPUSum



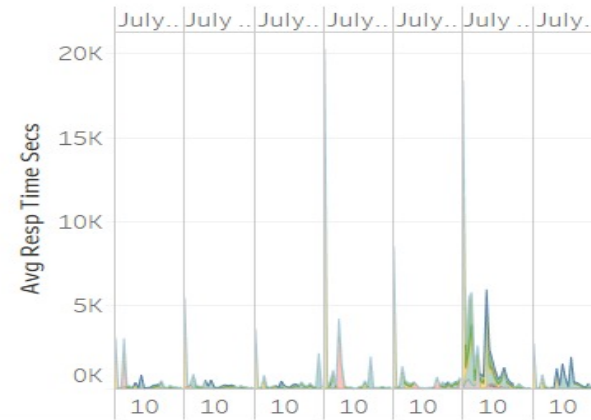
SKewCPU



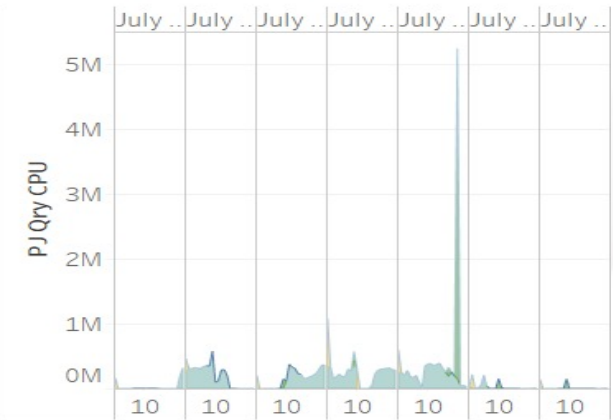
LgScanCPU



AvgRespTimeSecs



ProductJoin



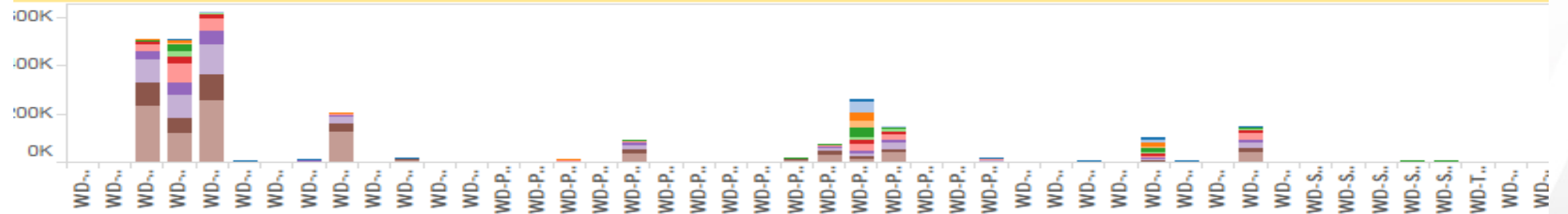
**EXEC TrendAnalysisDtRng_Rpt (date-7, date-1);
(Macro in Bonus Section)**

Workload Management Sample

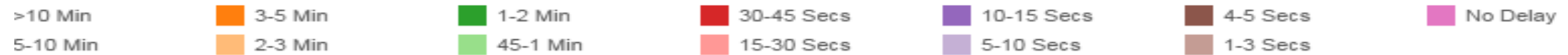
Percent of Overall delays by Day / Hour

nth, Da..	00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23
y 16, 2017	1%																							
y 15, 2017	2%	4%	2%	1%	0%	2%	5%	10%	11%	6%	1%	3%	3%	3%	2%	1%	0%	0%	0%	2%	0%	1%	8%	3%
y 14, 2017	0%	1%	0%	0%	3%	0%	0%	2%	5%	3%	4%	3%	8%	2%	4%	3%	2%	4%	4%	10%	2%	3%	6%	0%
y 13, 2017	0%	1%	3%	0%	1%	1%	2%	2%	0%	0%	0%	0%	1%	0%	0%	0%	0%	0%	14%	11%	0%	0%	0%	0%
y 12, 2017	0%	2%	3%	0%	8%	12%	3%	4%	3%	4%	5%	7%	2%	0%	0%	0%	0%	2%	2%	5%	0%	0%	0%	0%
y 11, 2017	0%	2%	3%	1%	3%	14%	22%	15%	10%	3%	1%	2%	0%	1%	0%	0%	0%	0%	1%	1%	0%	0%	0%	0%
y 10, 2017	0%	3%	3%	2%	2%	12%	11%	17%	16%	13%	15%	12%	4%	3%	2%	4%	2%	1%	1%	3%	0%	0%	2%	1%
y 9, 2017	0%	2%	2%	0%	3%	4%	7%	6%	4%	7%	2%	2%	6%	16%	14%	7%	2%	2%	0%	1%	0%	0%	2%	0%
y 8, 2017	1%	1%	3%	0%	3%	2%	11%	24%	6%	11%	8%	6%	7%	11%	3%	3%	2%	1%	0%	3%	1%	0%	2%	0%

Delays by Workload Type



Delay Time



Priority Scheduler



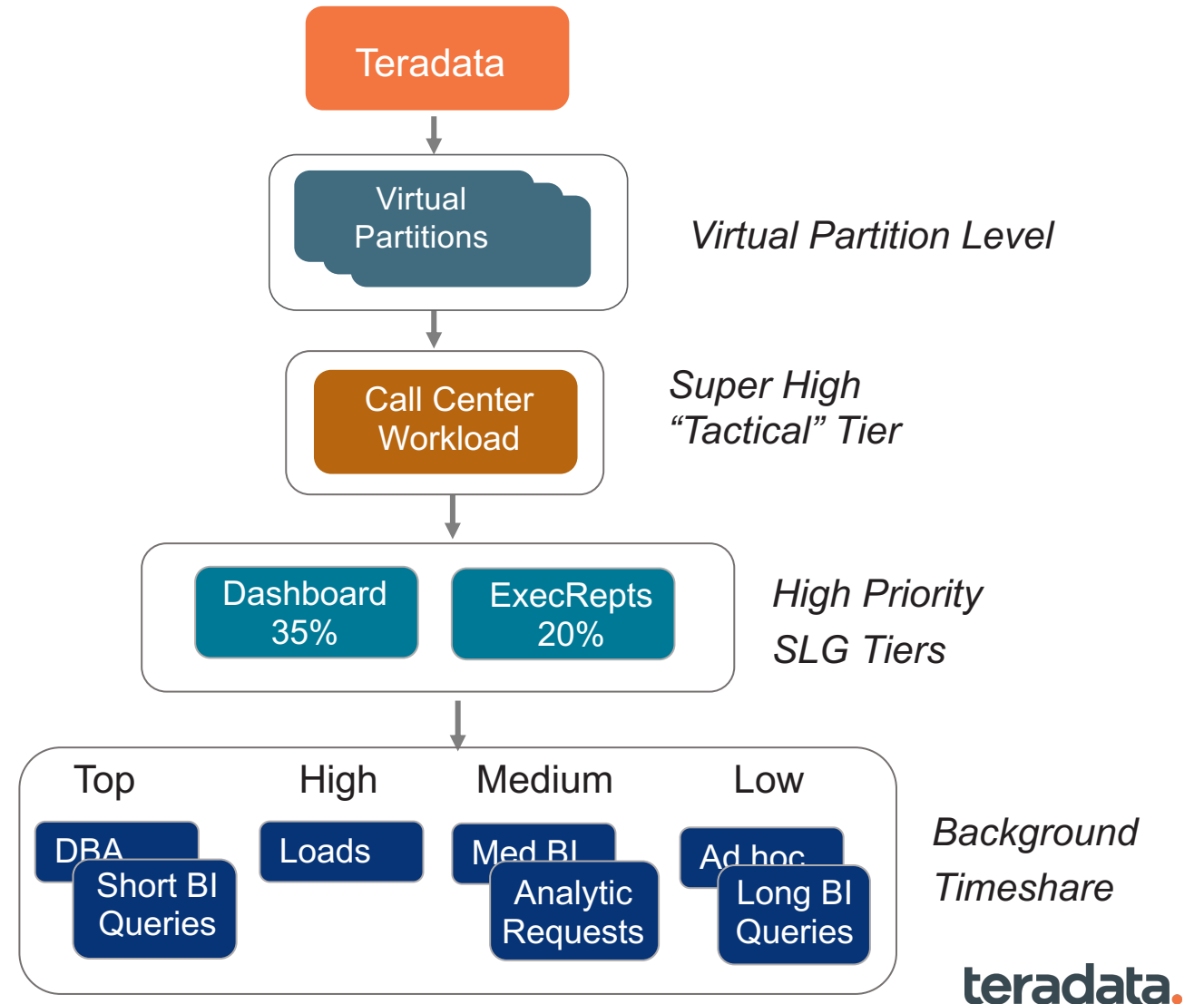
SLES11 – Virtual Partitions

Finance 60%	Marketing 30%	Sales 10%
Finance Workload Management Model	Marketing Workload Management Model	Sales Workload Mgt Model

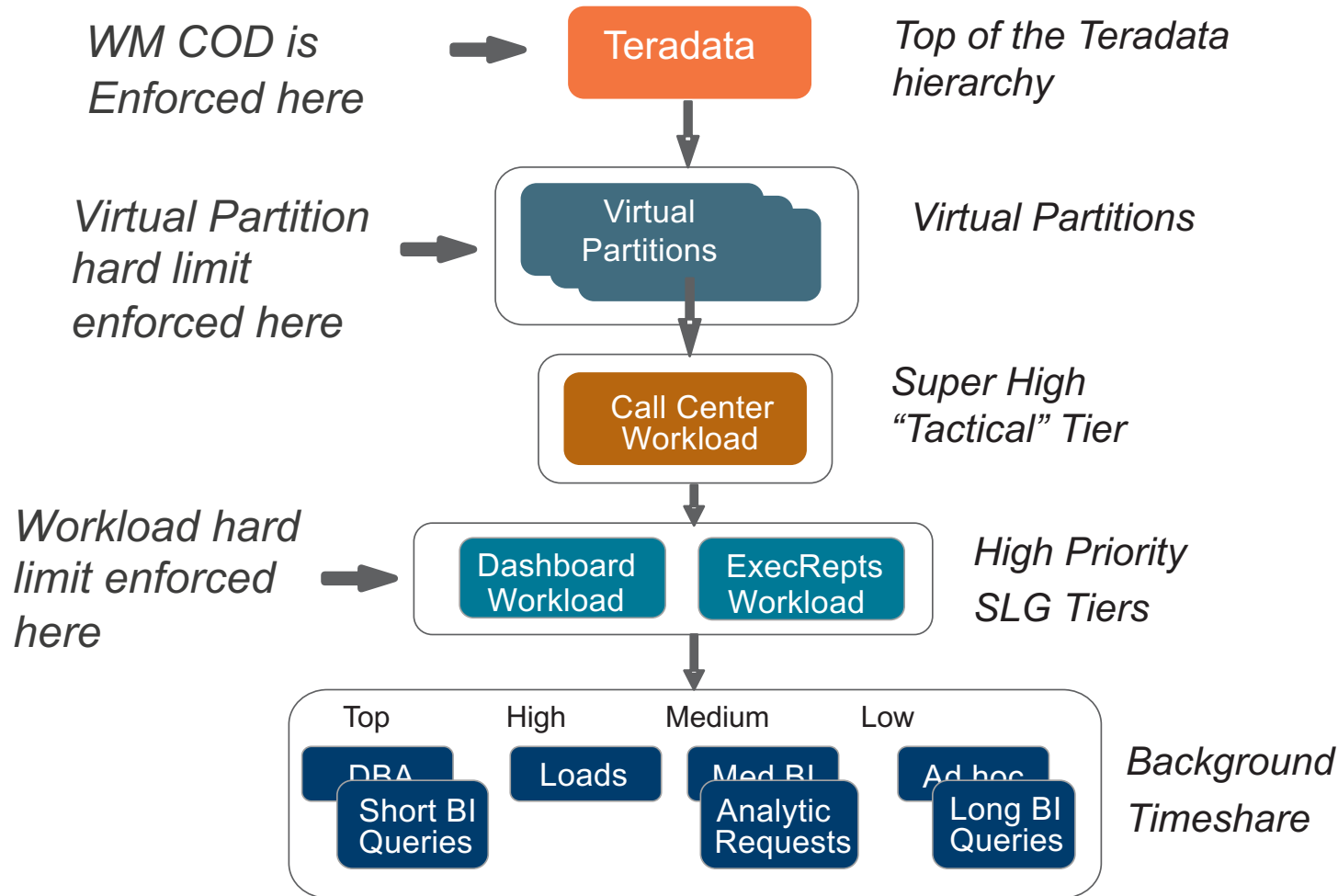
- Values indicate the target % of CPU or I/O for that partition
- Up to 10, default of 1
- TASM only feature

TASM Workloads are Placed in the Priority Hierarchy

- The “workload” is the priority component
- Hierarchical-based priority definitions
- Resources flow from the top tiers downward
- A special very high priority “Tactical Tier”
- Simple-to-use GUI tool
- Includes both CPU and I/O prioritization



Hard Limits at Multiple Levels



- WM COD puts a hard limit at the top of the Teradata priority hierarchy
- Defining a “fixed” VP enforces a hard limit at the VP’s defined allocation %
- SLG Tier workloads can have hard limits of any percent defined
 - Timeshare and tactical workload do not support hard limits

Viewpoint Workload Designer with SLES11

Workloads | SLG Summary | Evaluation Order | **Virtual Partitions** | Partition Resources | **Workload Distribution** | Console Utilities

Specify the allocation of resources for each workload in the context of a virtual partition in each planned environment. **System Workload Report**

Planned Environment: ALWAYS | Virtual Partition: STANDARD | 100% of system

Tactical

CallCenter

SLG Tiers +

Tier 1 ✎

Dashboard 15% | ... 10% | ... 10% | 65% remaining

Tier 2 ✎

CustOffer 30% | ExecReports 20% | 50% remaining

Tier 3 ✎ 🗑

MktgQry 40% | QuickDrill 30% | 30% remaining

Timeshare

Top	High	Medium	Low
Adhoc-High	BI-Tools	Adhoc-Medium	DataLab
Admin	ETL-Loads	Appl-Loads	Forecasting
T-WD	H-WD	WD-Default	L-WD
		M-WD	

Teradata 14.0 SLES 11 Priority Scheduler

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	■

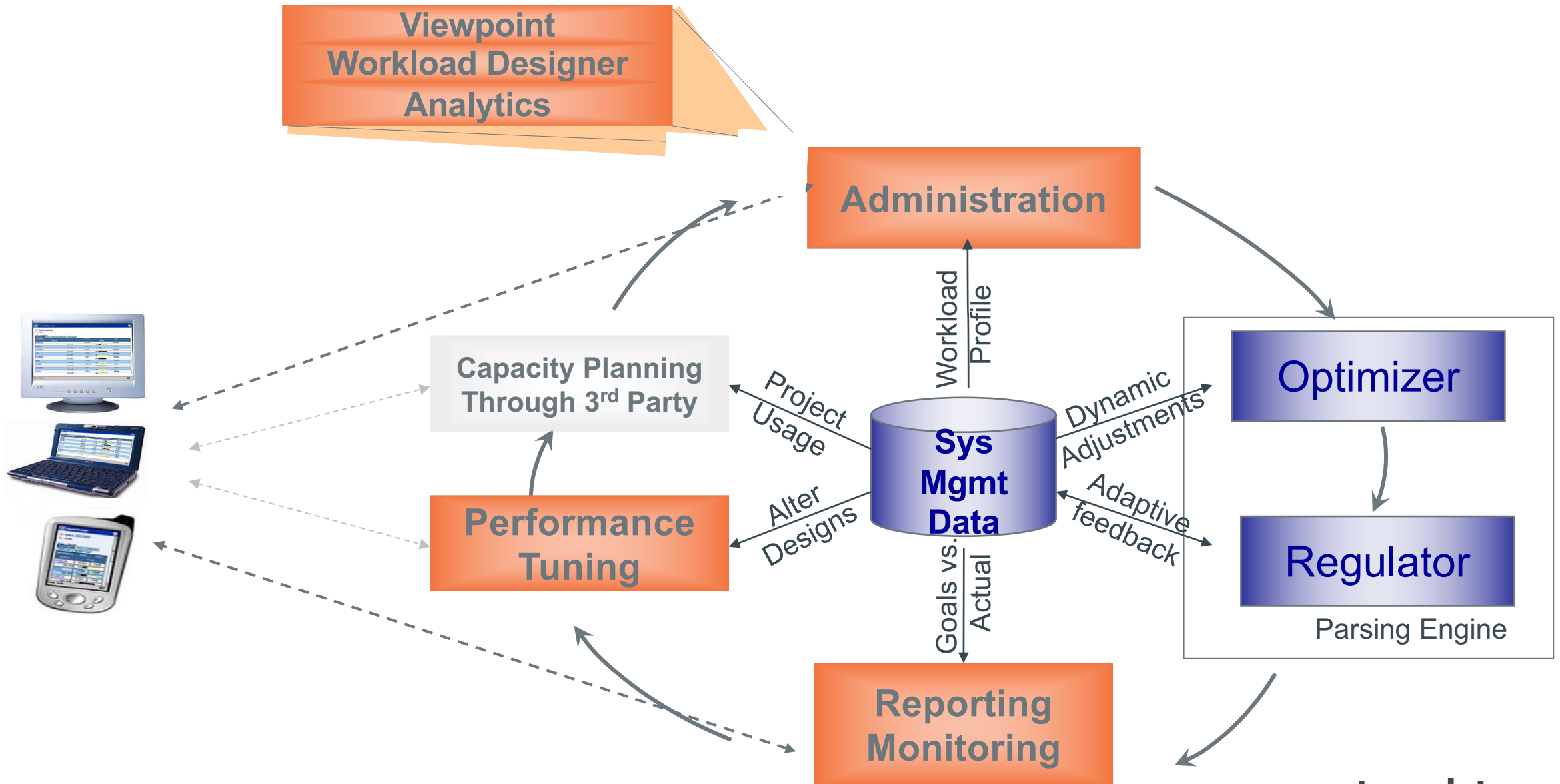
23 rows total

Close

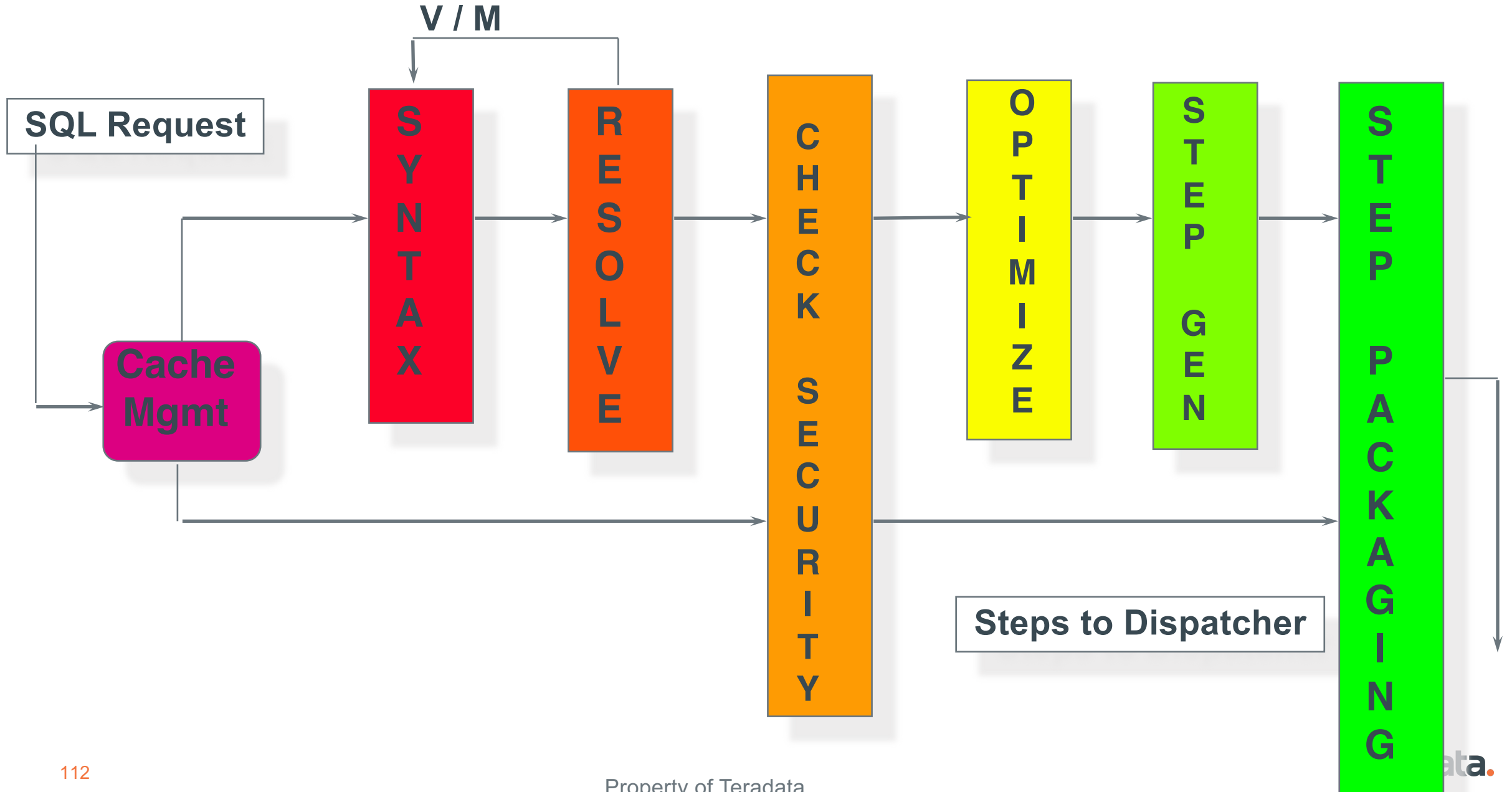
TASM Architecture



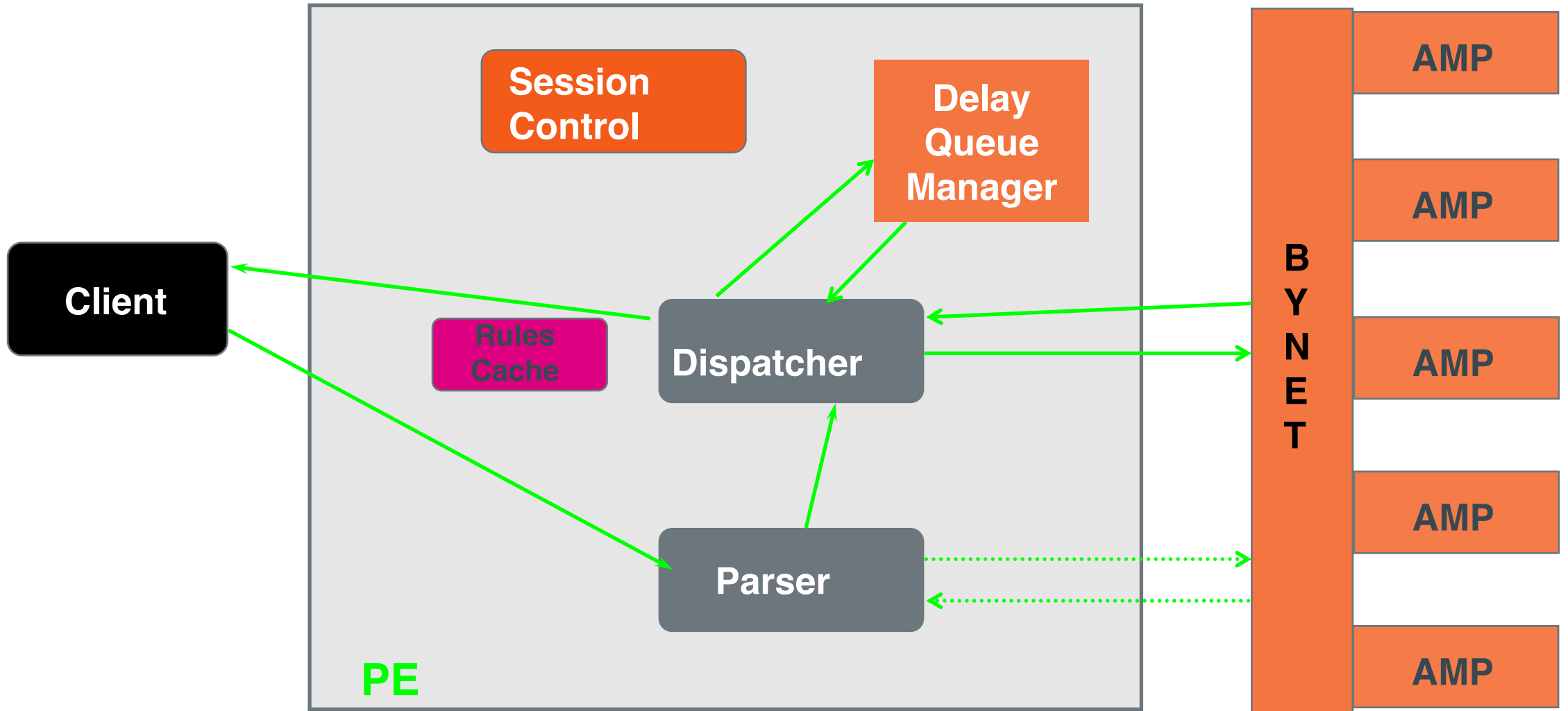
TASM – High level Architecture flow



Parser Flow

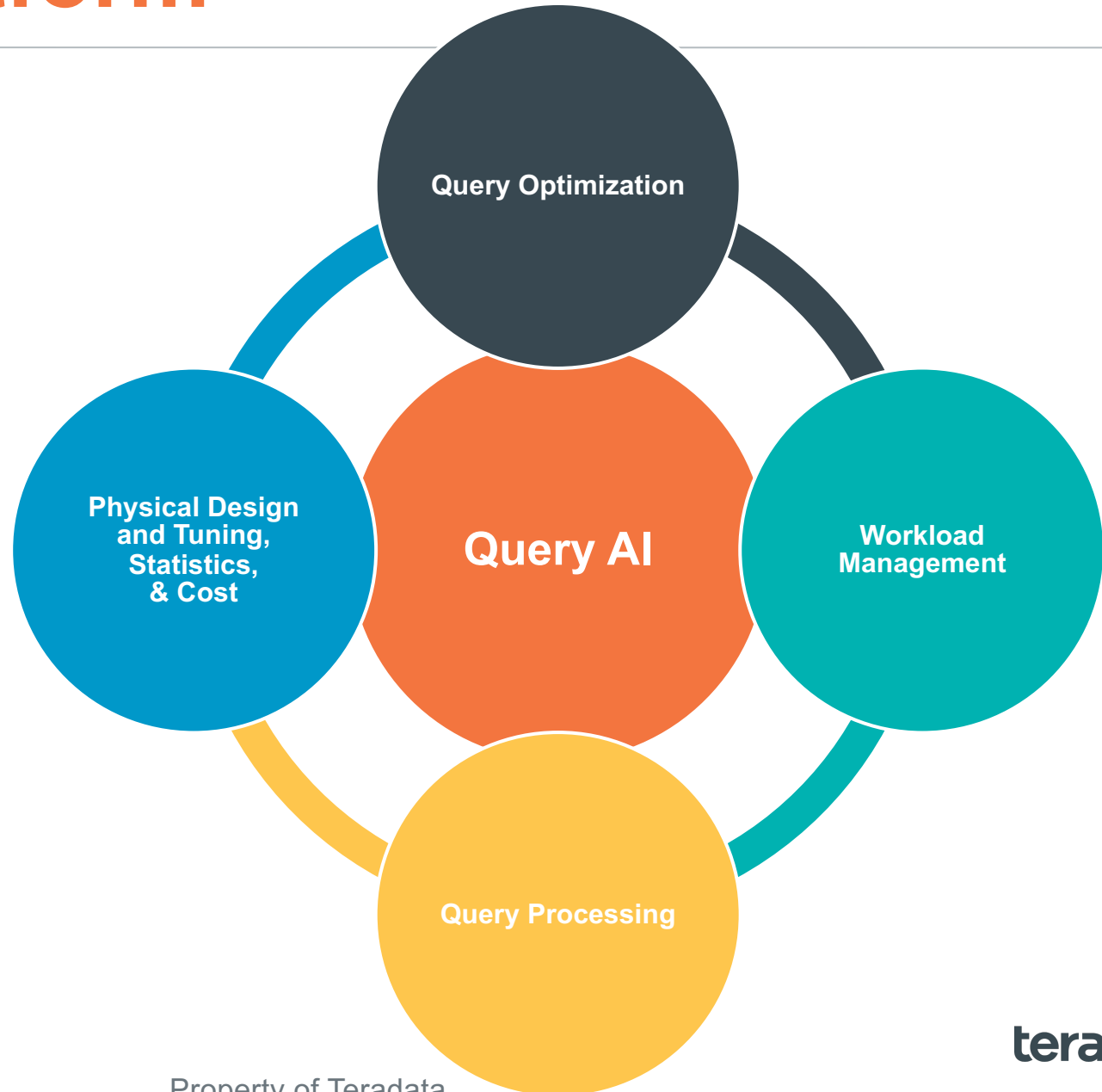


Parser Request Flow



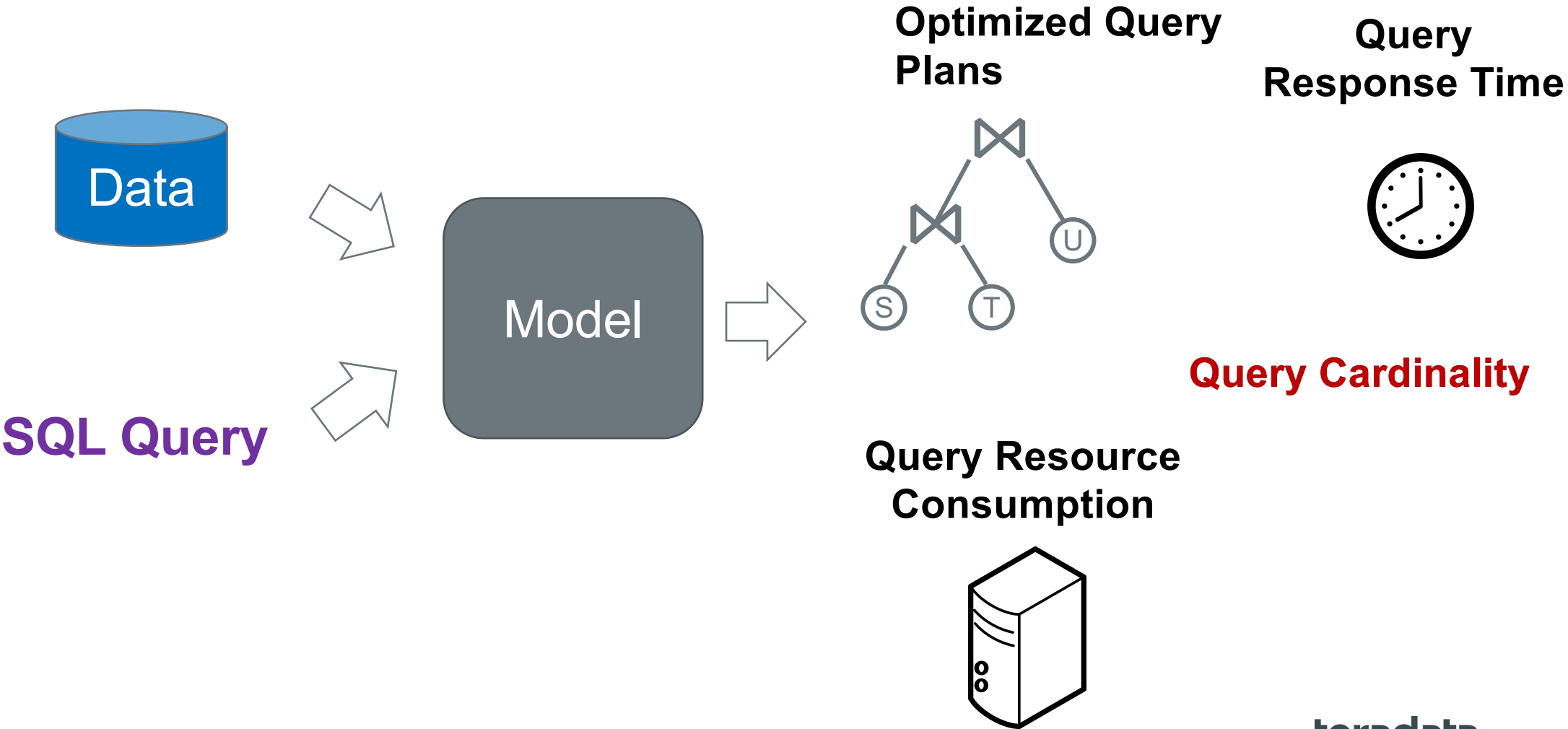
Autonomous Platform

- **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

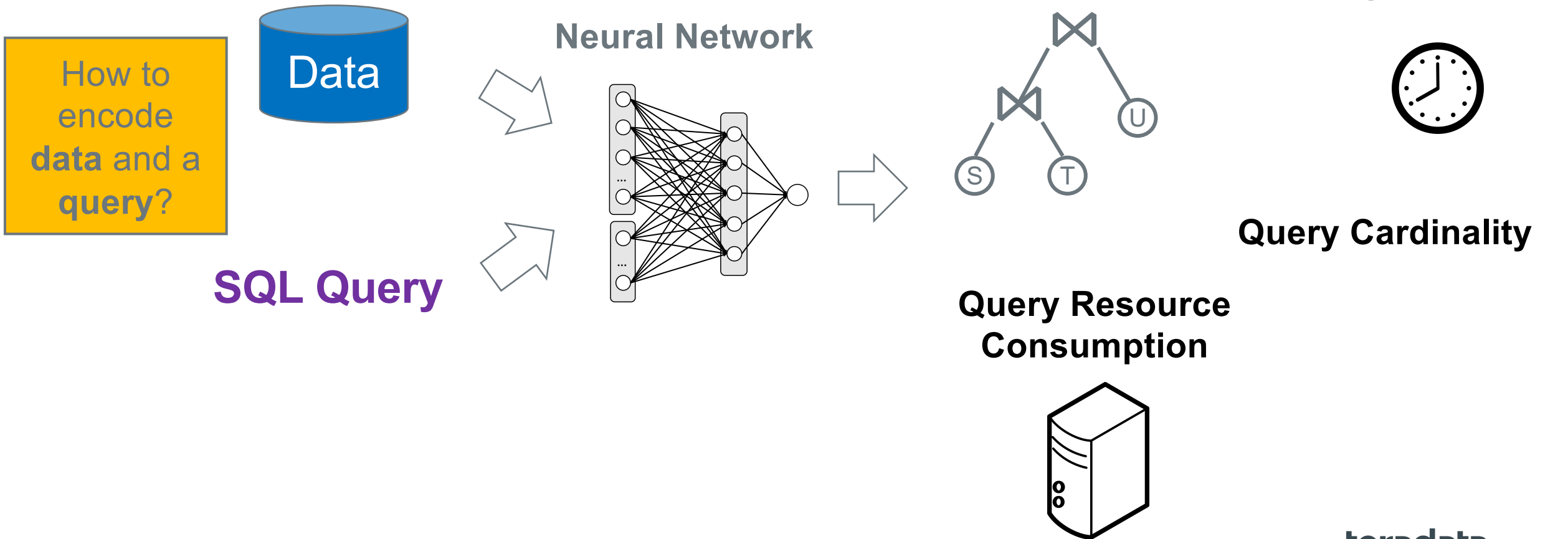


Our Vision

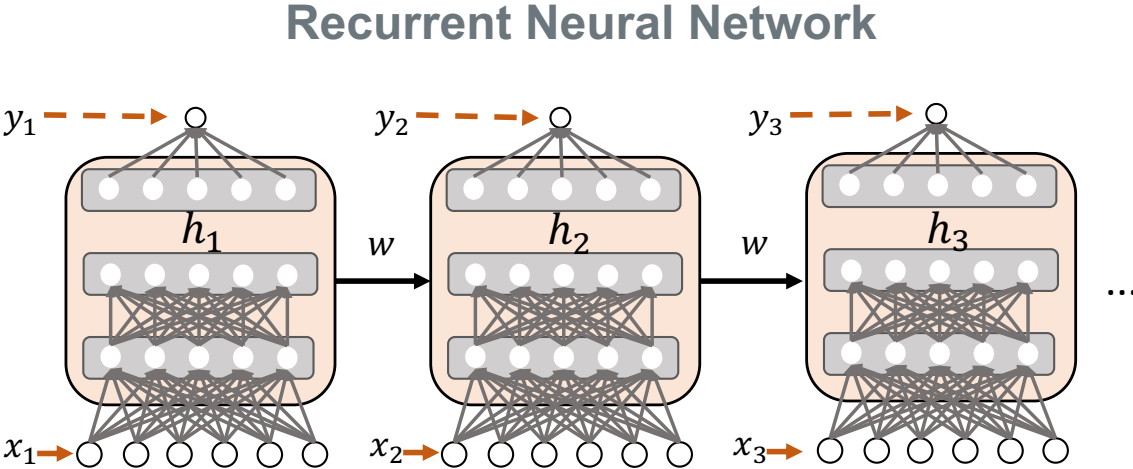
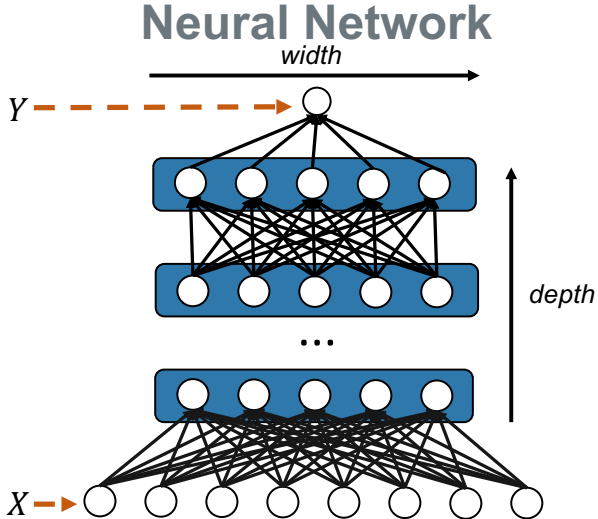
Can we rethink **Query Optimization** and **WLM** in the context of **Deep Learning**?



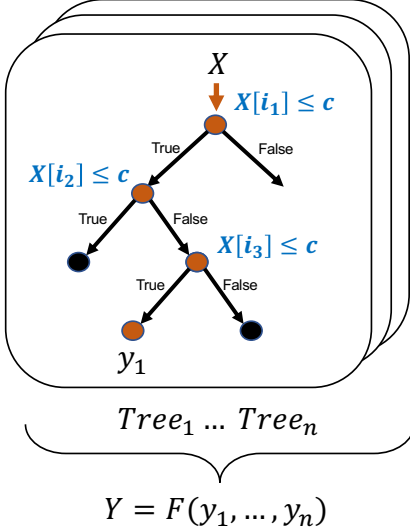
Deep Learning with Query Optimization



Models



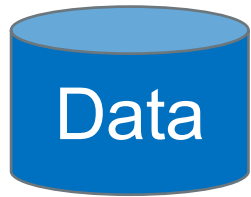
Random Forest



Encoding the Input

Example:

Orders ⋈ **Customers**

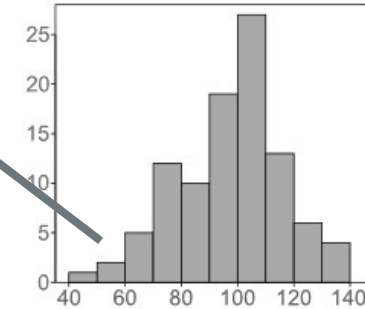


Basic Statistics:
1D-Histograms

Input Array



histogram of
Orders column



SQL Query

Encode Join Pred
Encode Selection

Need ability to
encode all
possible queries!

on customers and orders

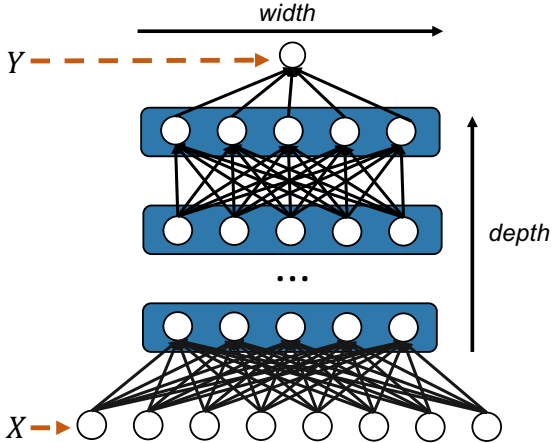
on customers and region

Encoding a Query

Example dataset:

- **Relations** : A, B, C
- Each has two columns:
 - A: a1 and a2
 - B: b1 and b2
 - C: c1 and c2

Neural Network



q: SELECT * FROM **A**,**B** WHERE **a₁** <= **10** and **a₂**=**b₁**

1	1	0	.1	1	1	1	0	0	1	0
---	---	---	----	---	---	---	---	---	---	---

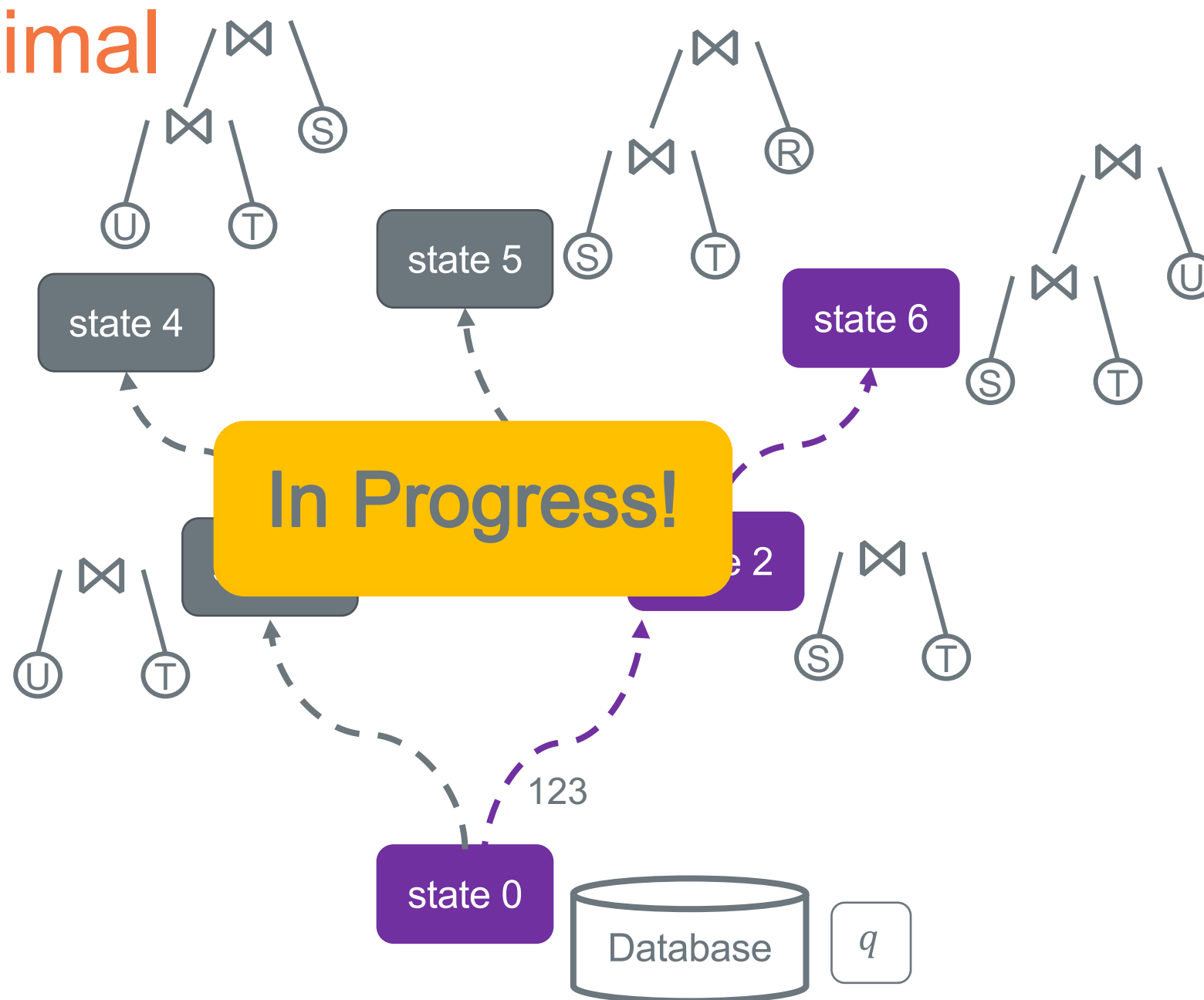
A	B	C	a₁	a₂	b₁	b₂	c	c₂	a₂=b₁	b₂=c₁	
<hr/>			<hr/>				1	<hr/>			

relation

selpred

joinpred teradata.

Learning Optimal Paths



Example: $S \bowtie T \bowtie U$

Property of Teradata

Input

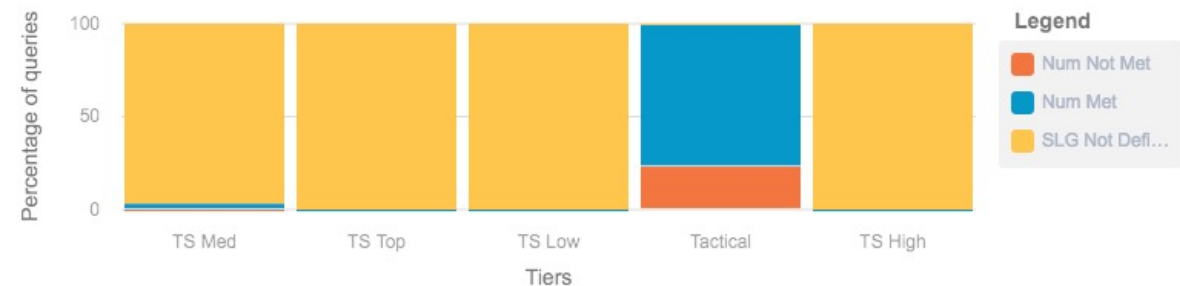
Provide the input to analyze the SLG

Rest Server IP Address 10.25.188.63:1080	Host italy	TD user nt186028	Password ...
DBQL_DB DBC	RESUSAGE_DB DBC	TDWM_DB TDWM	PDCR_DB PDCRINFO
Window Start 10/08/2018 00:00:00	Window End 10/14/2019 23:59:59		

SUBMIT

All Tiers Details

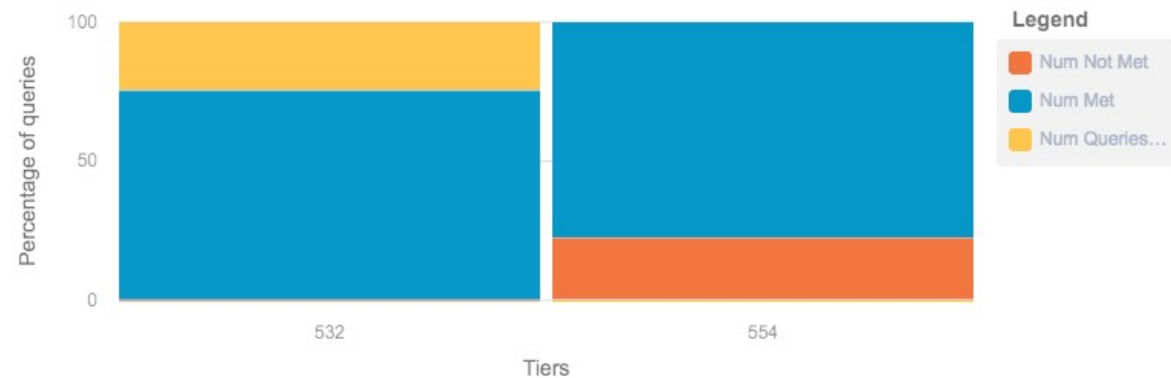
Per tier SLG analysis



VIEW DATA IN TABULAR FORM

Tactical Tier Analysis

Per WDID SLG analysis



VIEW DATA IN TABULAR FORM

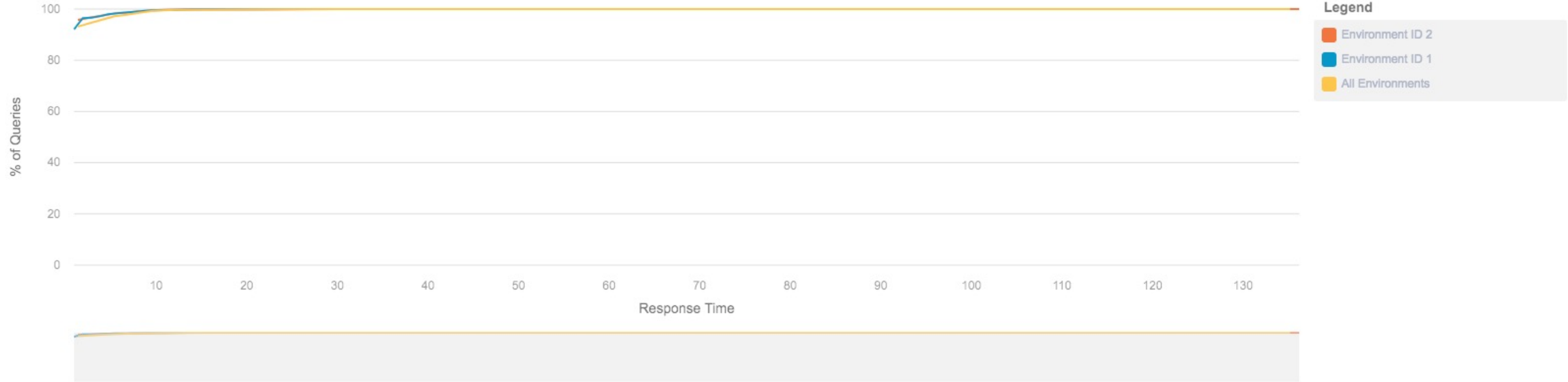
554 WDID details

SLG level details for WDID

Attribute	Value
WDID	554
Tier Name	Tactical
Workload Name	TLES_RTOC_Query_Throttled
Total Queries	786901
Number of Queries Without SLG	0

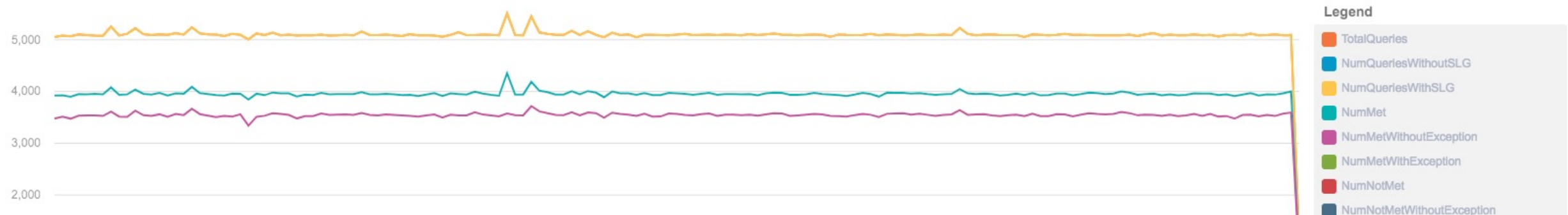
Response time distribution for WDID 554

Response time distribution



Query Statistics for WDID 554

SLG related statistics



Resource utilization for WDID 554

CPU, IO, AWT and other resource details

CPU Utilization

AWT

I/O

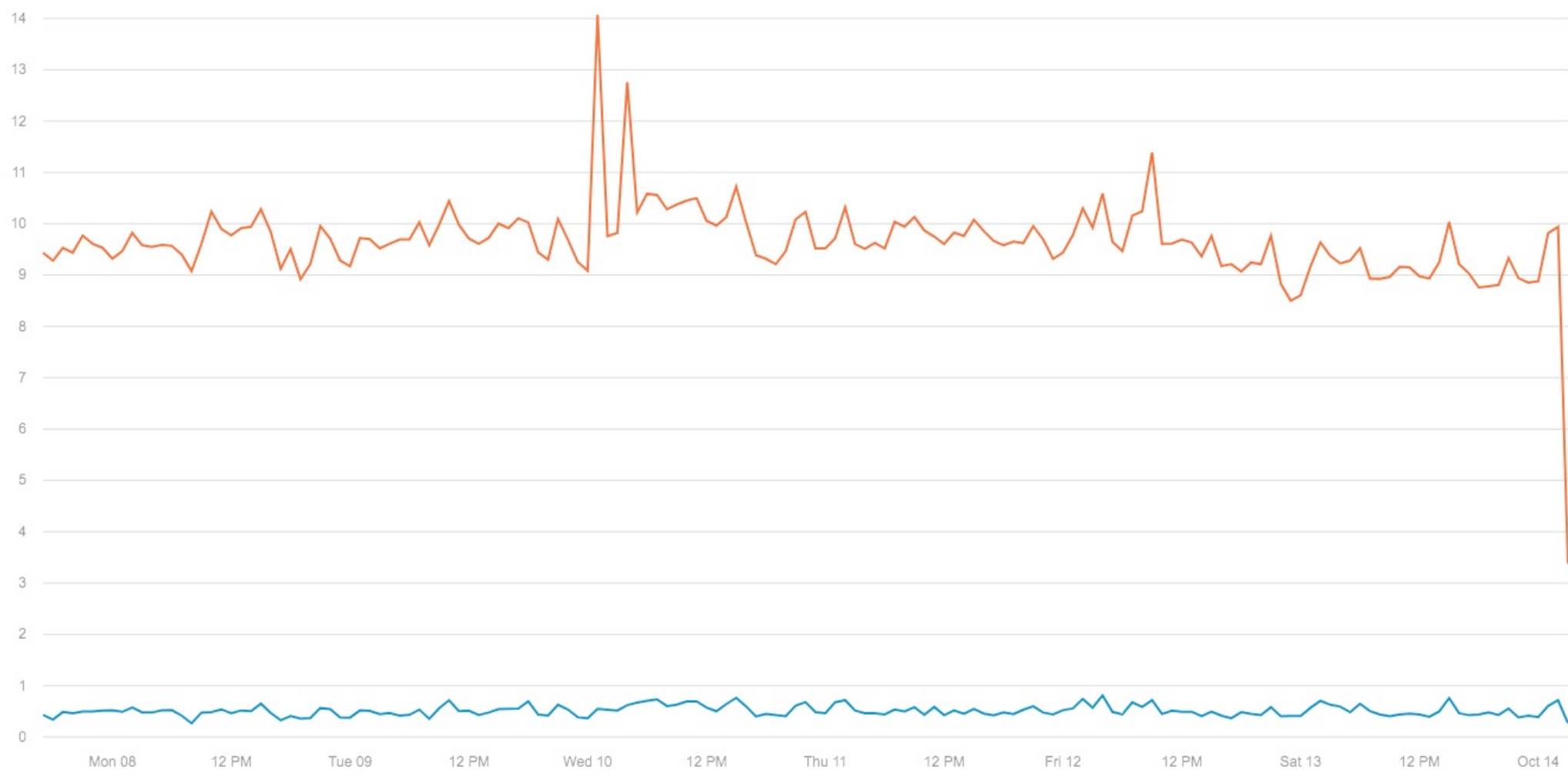
Request Count

Flow Control

Lock Wait

Other Wait

WD: 554



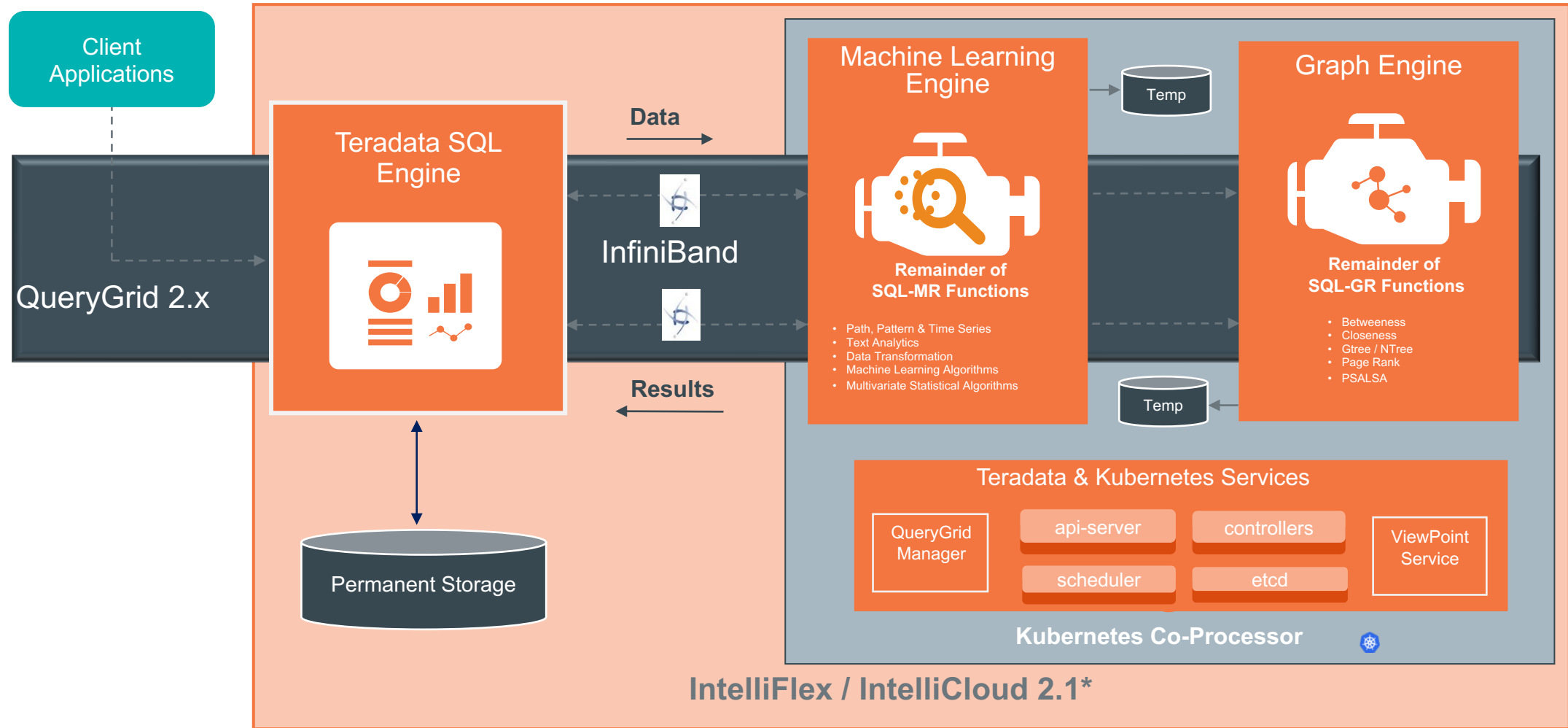
Legend

- CPU_USE_PCT_WD
- CPU_WAIT_PCT_WD

Teradata Vantage – Information Resources

- YouTube Teradata Database 101
- www.teradata.com
- <https://www.teradatauniversitynetwork.com>
- <https://www.teradata.com/Resources/TEN>

Teradata Vantage



Thank you.

teradata.

©2018 Teradata