ORACLE

# Protecting Critical OLTP Workloads in a Mixed Workload Environment
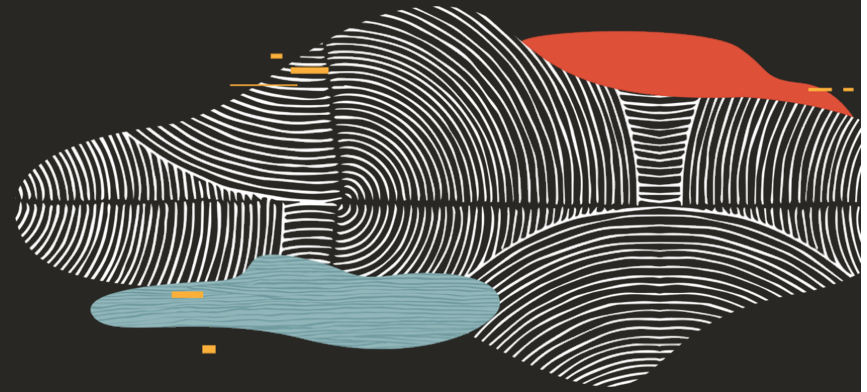
**Mihajlo Tekic**

**John Zimmerman**

Real-World Performance Team
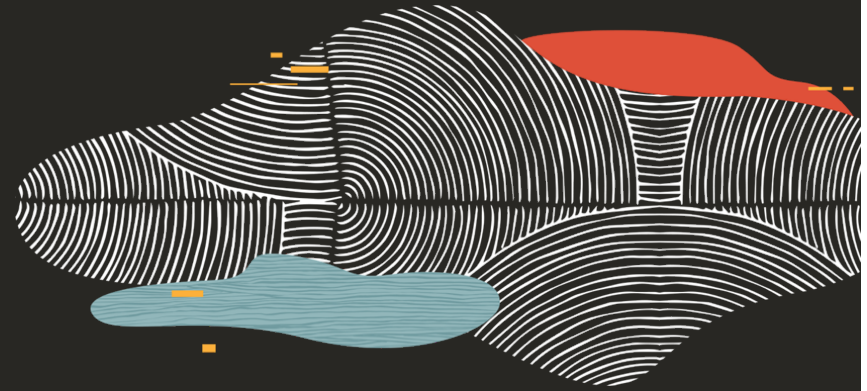Oracle Database Development

# Agenda

1     Types of Workloads

2     Challenges of Mixed Workloads

3     Ways to Manage Mixed Workloads

4     Demo

5     Lessons Learned

# Agenda

**Workload Characteristics**

## OLTP

- Many concurrent users
- SQL statements process a few rows at a time

## DW Queries

- Fewer concurrent users
- Data-intensive queries processing many rows

## Data Loading and Processing
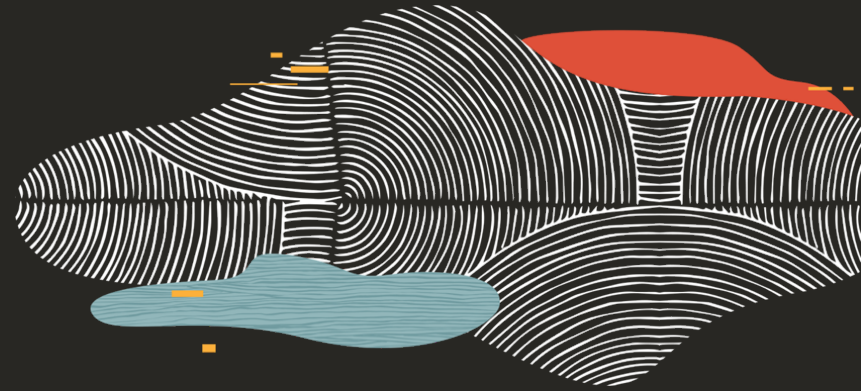
- Fewer concurrent processes
- DML processing many rows

ORACLE®
REAL-WORLD PERFORMANCE

# Agenda

# Competition for Resources

ORACLE
REAL-WORLD PERFORMANCE

## Memory

- Avoid competing for memory
- Competing for memory ends badly

**Network and Storage**

- Oracle uses CPU when performing network and storage I/O
- Limiting CPU naturally limits network and storage I/O
- Focus on CPU

## CPU

- The Oracle database uses a process based architecture: when you connect, a dedicated foreground process is started to serve your calls

- To perform efficiently, a process:

  1. needs to get on CPU as quickly as possible
  2. needs to stay on CPU as long as possible:
  3. should minimize voluntary sleeps
  4. should experience as few involuntary sleeps as possible

**CPU Resources and OLTP Workloads**

- As CPU Utilization increases, the chance of a process getting scheduled on CPU decreases
- This has a noticeable impact on OLTP performance at 60-70% CPU utilization

| CPU Utilization | Chance of getting scheduled |
| --- | --- |
| 50% | 1 in 2 |
| 66% | 1 in 3 |
| 80% | 1 in 5 |
| 90% | 1 in 10 |

**ORACLE®**
REAL-WORLD PERFORMANCE

**The Mixed Workload Dilemma—Opposing resource management goals**

| Workload | Goal | CPU Strategy |
|---|---|---|
| OLTP | Fast Response Time | Minimize |
| Analytical Queries | Throughput and Response Time | Maximize |
| Data Processing | Throughput | Maximize |

ORACLE®
REAL-WORLD PERFORMANCE

# Agenda

**Ways to Manage Mixed Workloads**

—

## Multiple Databases

- Virtual Machines

- Instance Caging

- Multitenant

## Single Database

- RAC Services

- Database Resource Manager

**ORACLE**
REAL-WORLD PERFORMANCE

## Multiple Databases: Virtual Machines

- Allocate virtual machines for each workload

- Workloads cannot use more than allocated CPUs

**Server**

| OLTP VM 24 CPUs | Query VM 18 CPUs | ETL VM 6 CPUs |

ORACLE®
REAL-WORLD PERFORMANCE

## Multiple Databases: Instance Caging

- Use instance caging to control the number of processes on CPU

- Use the CPU_COUNT parameter to control

- Enable a Database Resource Manager plan

**Server**

| OLTP | Query | ETL |
|------|-------|-----|
| CPU_COUNT=24 | CPU_COUNT=18 | CPU_COUNT=6 |

**ORACLE**
REAL-WORLD PERFORMANCE

## Multiple Databases: Multitenant

- Enable a CDB resource plan

- Use instance caging to limit CPU of individual PDBs

- Or use Shares or Limits

**Server**

**OLTP PDB**
CPU_COUNT=24

**Query PDB**
CPU_COUNT=18

**ETL PDB**
CPU_COUNT=6

**ORACLE**
REAL-WORLD PERFORMANCE

**Single Database: RAC Services**

- Use services for different workloads
- Map services for different workloads to different nodes in a cluster

| OLTP Node | Query Node | ETL Node |
|-----------|------------|----------|

**Single Database: DBRM Consumer Groups**

- Create a DBRM plan
- Map workloads to different consumer groups
- Use Shares or Limits

**Server**

OLTP

Query

ETL

ORACLE®
REAL-WORLD PERFORMANCE

**Database Resource Manager**

## Shares

Divide resources between workloads using ratios

## Limits

Set hard limits on CPU utilization for each workload

## Parallel Queuing

Control the number of PX processes used by each workload

These can be combined to develop a more sophisticated CPU utilization strategy

ORACLE
REAL-WORLD PERFORMANCE

**Shares**

| Workload | Shares |
|----------|--------|
| OLTP | 7 |
| Queries | 3 |
| Total | 10 |

## OLTP Busy

OLTP

Query

## Query Busy

Query

OLTP

## Both Busy

OLTP

Query

ORACLE®
REAL-WORLD PERFORMANCE

| Workload | Limits |
|----------|--------|
| OLTP | |
| Queries | 30% |

**Limits**

## OLTP Busy

OLTP

Query

## Query Busy

OLTP

Query

## Both Busy

OLTP

Query

ORACLE®
REAL-WORLD PERFORMANCE

# Agenda

1    Types of Workloads

2    Challenges of Mixed Workloads

3    Ways to Manage Mixed Workloads

4    Demo

5    Lessons Learned

**Demo Workloads**

## OLTP

- Lots of users playing online games
- Short transaction times
- Sensitive to high server CPU utilization
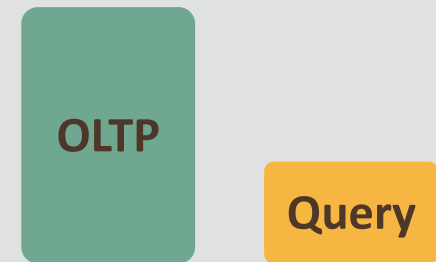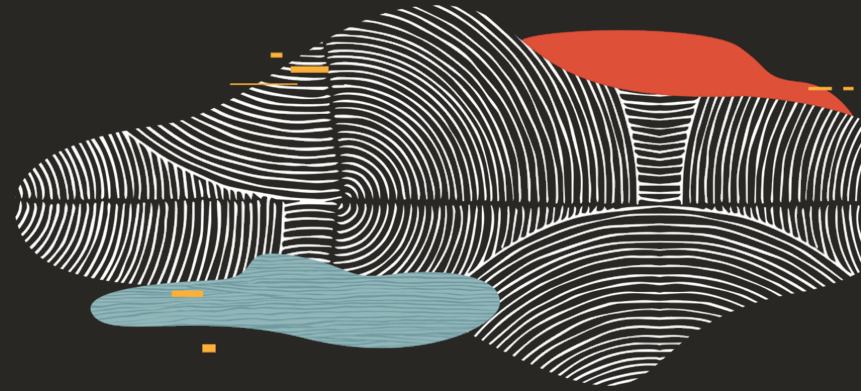- This is the workload we want to protect

## Queries

- 16 users running analytical queries which take a few seconds each
- Parallel execution with parallel degree 8
- Some queries perform table scans from disk and some scans are from memory

## ETL

- Single user performing an ELT strategy
  1. Load Data
  2. Remove duplicates
  3. Transformations
  4. Aggregation
- Parallel execution with parallel degree 16

ORACLE®
REAL-WORLD PERFORMANCE

**OLTP Workload**

- We can control the workload by changing Think Time

    Decreasing Think Time increases demand
    Increasing Think Time decreases demand

- 4000ms represents low demand

- 500ms represents expected peak demand

**ORACLE**
REAL-WORLD PERFORMANCE

1. Start the OLTP workload, setting Think Time to 4000 ms, representing low demand

2. About 9000 transactions per second with

3. OLTP CPU is about 3%

4. Sub-millisecond response time in the database

1. Adjust Think Time to 2000 ms, representing medium demand

2. TPS doubles

3. OLTP CPU doubles

4. Response time in the database is almost unchanged

1. Adjust Think Time to 1000 ms, representing high demand

2. TPS doubles

3. OLTP CPU doubles

4. Response time in the database remains sub-millisecond

1. Adjust Think Time to 500 ms, representing peak demand

2. TPS doubles

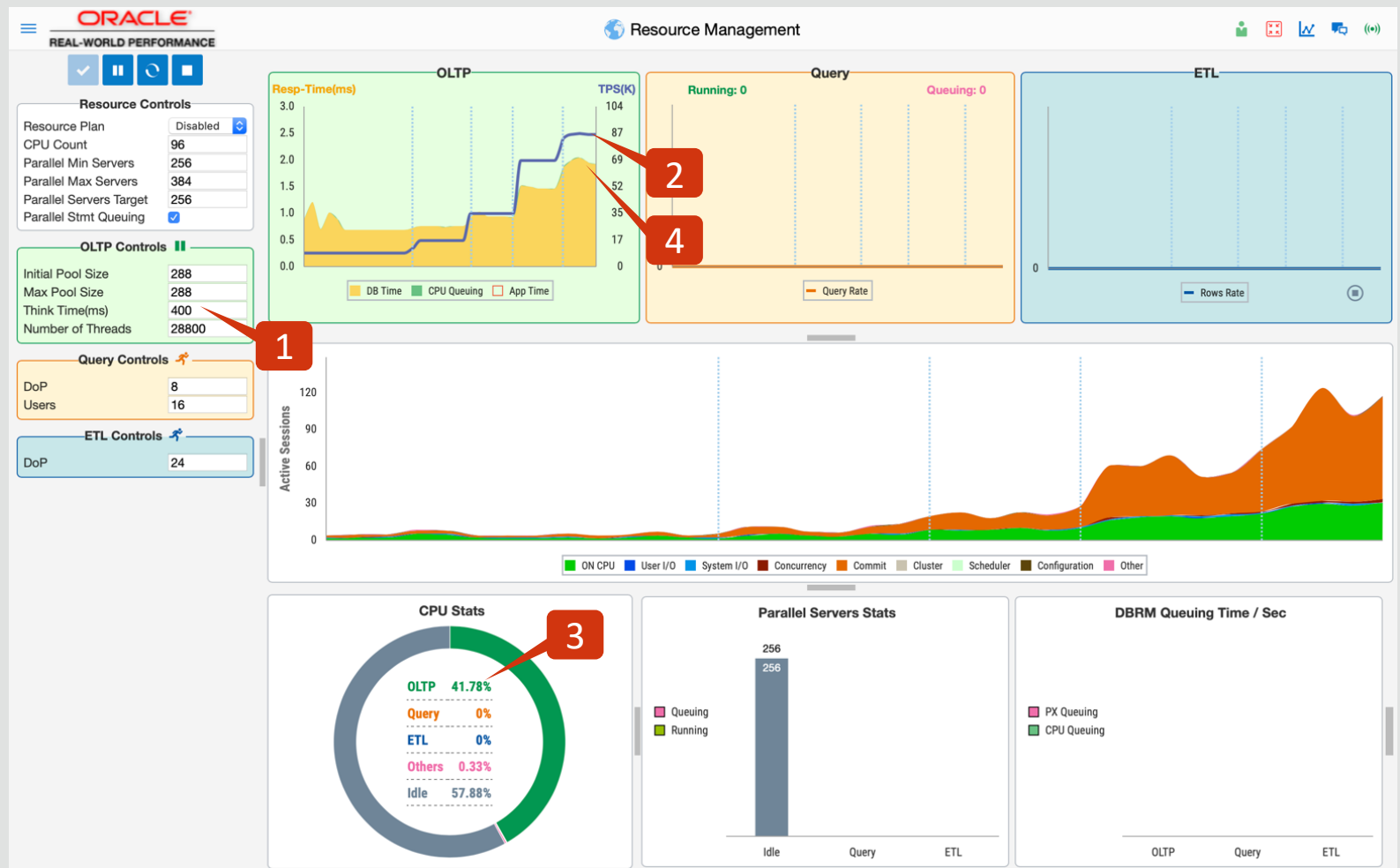3. OLTP CPU doubles

4. Response time in the database also increases

1. Adjust Think Time to 400 ms to check for headroom

2. TPS increases

3. OLTP CPU increases

4. Response time in the database also increases

1. Application Time indicates response time in the database plus time waiting for a connection. There is a short spike in waiting for a connection when demand increases.

1. Adjust Think Time to 333 ms to check for headroom

2. TPS increases

3. OLTP CPU increases

4. Response time in the database also increases

1. Some waiting for connections is now present at all times, degrading overall response times

**Query Workload**

- We can control the workload by changing Users

    Increasing Users increases demand
    Decreasing Users decreases demand

- 1 User represents low demand

- 8 Users represents expected peak demand

ORACLE®
REAL-WORLD PERFORMANCE

1. Reset Think Time to 4000 ms to represent low demand

2. About 9000 transactions per second with

3. OLTP CPU is about 3%

4. Sub-millisecond response time in the database

1. Start the Query workload, setting Users to 1, representing low demand

2. OLTP throughput is unchanged

3. Query CPU is around 10%

4. Sub-millisecond response time in the database

1. Adjust Users to 8, representing peak demand

2. OLTP throughput is unchanged

3. Query CPU is around 40%

4. Response time in the database increases

1. Adjust Parallel Servers Target to 16

2. OLTP throughput is unchanged

3. Query CPU is around 10%

4. Response time in the database improves

**ETL Workload**

- We can control the workload by changing Degree of Parallelism (DoP)

  Increasing DoP increases demand
  Decreasing DoP decreases demand

- Degree of Parallelism of 24 represents expected demand

ORACLE®
REAL-WORLD PERFORMANCE

1. Reset Parallel Servers Target to 256

1. Start the ETL workload, setting DoP to 24, representing expected demand

2. OLTP throughput is unchanged

3. ETL CPU is around 25%

4. Response time in the database increases

**Increasing The Demand For OLTP**

- We can control the OLTP workload by changing Think Time

  Decreasing Think Time increases demand
  Increasing Think Time decreases demand

- 4000ms represents low demand

- 500ms represents expected peak demand

ORACLE®
REAL-WORLD PERFORMANCE

1. Adjust Think Time to 2000 ms, representing medium demand
2. TPS doubles
3. OLTP CPU doubles
4. Response time in the database increases

1. Adjust Think Time to 1000 ms, representing high demand

2. TPS doubles

3. OLTP CPU doubles

4. Response time in the database increases

1. Adjust Think Time to 1000 ms, representing high demand

2. TPS doubles

3. OLTP CPU doubles

4. Response time in the database increases

1. Adjust Think Time to 500 ms, representing expected peak demand

2. TPS increases a little but becomes erratic

3. OLTP CPU increases

4. Response time in the database increases

1. Significant waiting for connections is now present at all times, degrading overall response times

**Enabling Database Resource Manager**

- We can reduce competition for CPU by enabling Database Resource Manager

  Limit the CPU utilization for Query
  Limit the CPU utilization for ETL
  Increase the CPU available for OLTP

ORACLE®
REAL-WORLD PERFORMANCE

1. Enable Database Resource Manager, limiting Query to 15% CPU and limiting ETL to 10% CPU

2. TPS increases

3. OLTP CPU increases

4. Response time in the database improves

1. Waiting for connections is significantly reduced, dramatically improving overall response times

**Increasing Competition For CPU**

- We can increase competition for CPU by increasing the limits on CPU utilization

  Increasing the limit on CPU utilization for Query increases competition
  Increasing the limit on CPU utilization for ETL increases competition
  Increasing CPU utilization may degrade OLTP performance

ORACLE®
REAL-WORLD PERFORMANCE

1. Increase Limit on Query to 20% CPU

2. TPS is unchanged

3. OLTP CPU is about the same

4. Response time in the database increases a little

1. Waiting for connections increases a little

1. Increase Limit on Query to 25% CPU

2. TPS becomes erratic

3. OLTP CPU is about the same

4. Response time in the database increases

1. Waiting for connections is now present at all times, degrading overall response times

**Decreasing Competition For CPU**

- We can decrease competition for CPU by decreasing the limits on CPU utilization

    Decreasing the limit on CPU utilization for Query decreases competition
    Decreasing the limit on CPU utilization for ETL decreases competition
    Decreasing CPU utilization may improve OLTP performance

- We can check for headroom by changing Think Time

    Decreasing Think Time increases demand
    DBRM can limit the competition for CPU due to Query and ETL

**ORACLE®**
REAL-WORLD PERFORMANCE
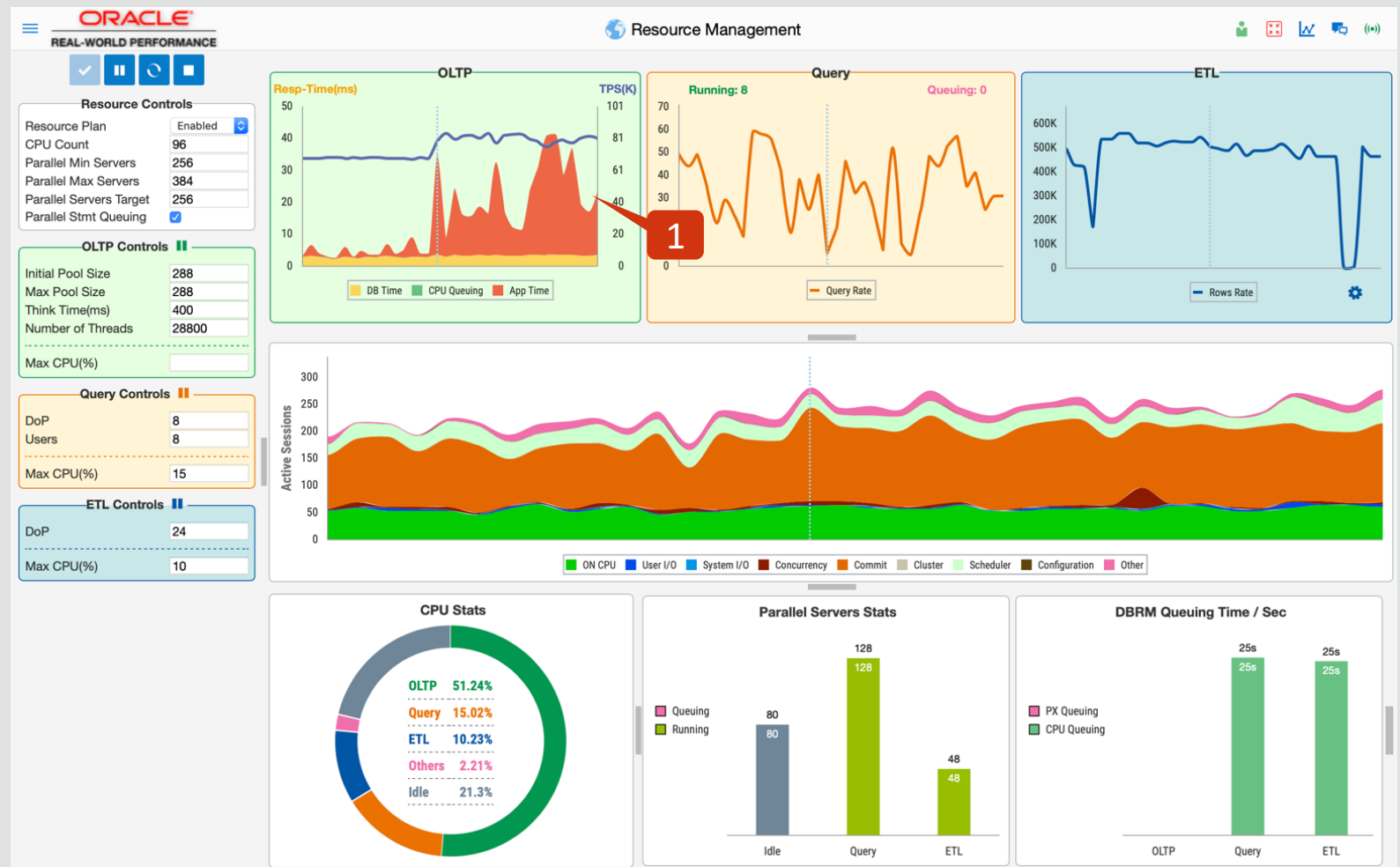
1. Reset Limit on Query to 15% CPU

1. Waiting for connections is significantly reduced, improving overall response times

1. Adjust Think Time to 400 ms to check for headroom

2. TPS increases but becomes erratic

3. OLTP CPU increases
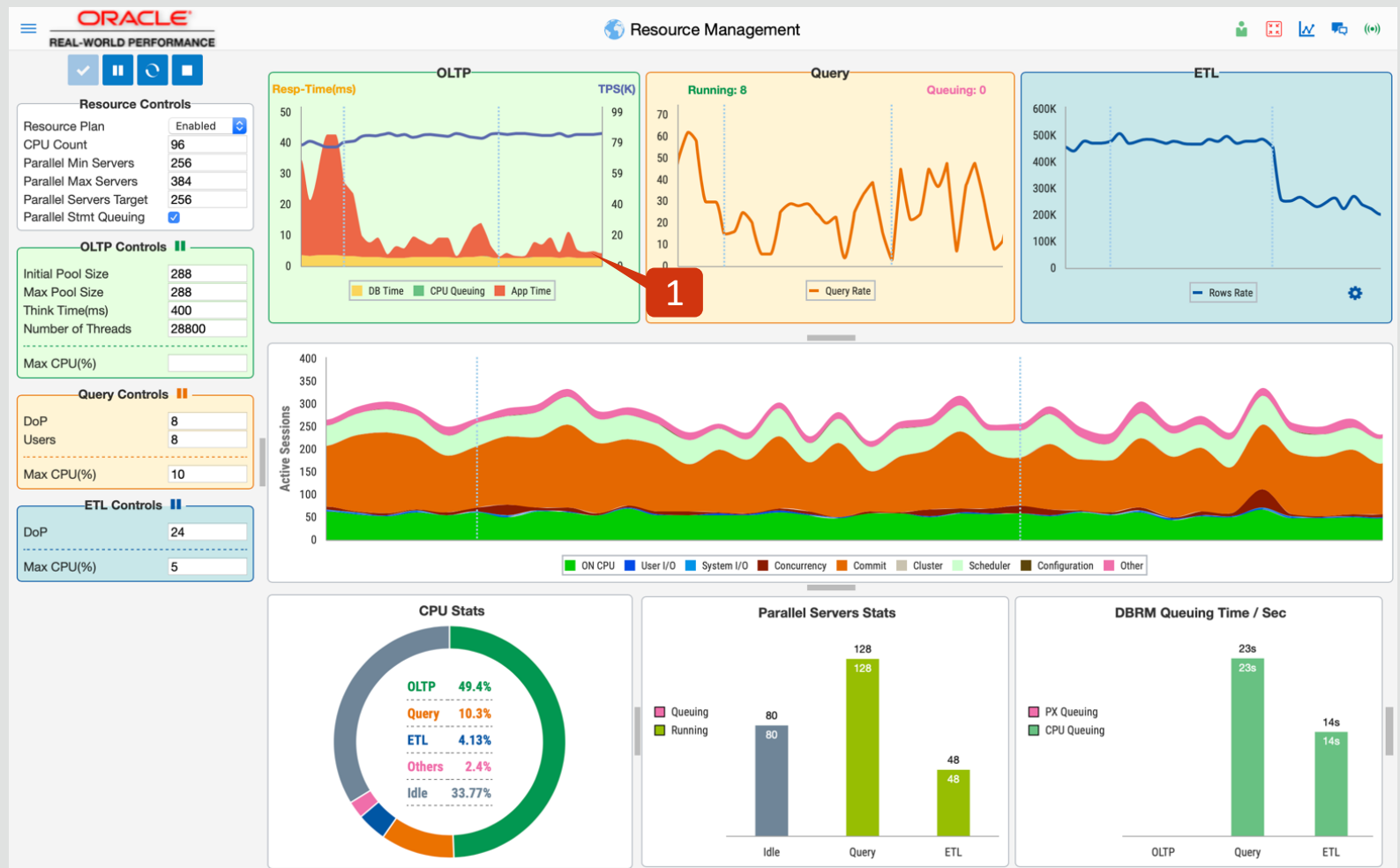
4. Response time in the database also increases a little

1. Significant waiting for connections is now present at all times, degrading overall response times. There is little headroom for increased demand from OLTP.

1. Reduce Limit on Query to 10% CPU and Limit on ETL to 5% CPU

2. TPS is consistent

3. OLTP CPU increases a little

4. Response time in the database improves a little

1. Waiting for connections is significantly reduced, dramatically improving overall response times

1. Reduce Limit on Query to 5% CPU

2. TPS is consistent

3. OLTP CPU is about the same

4. Response time in the database improves a little

1. Waiting for connections is reduced, improving overall response times

# Agenda

1  Types of Workloads

2  Challenges of Mixed Workloads

3  Ways to Manage Mixed Workloads

4  Demo

5  Lessons Learned

**How Do You Determine the Workload CPU Limits?**

1. Do not limit the critical OLTP workload

2. Determine the peak OLTP CPU Utilization % by itself

3. Use a simple formula to estimate the CPU available for other workloads

   60~70 – ((OLTP CPU %) * 1.3) = CPU for others

ORACLE®
REAL-WORLD PERFORMANCE

**Example**

- OLTP peak CPU Utilization is 30% by itself

- Lower limit
  60 – ((30* 1.3) = Other workloads can use up to ~20% CPU

- Upper limit
  70 – ((30* 1.3) = Other workloads can use up to ~30% CPU

ORACLE®
REAL-WORLD PERFORMANCE

**Lessons Learned**

1. Need to limit system CPU utilization to protect OLTP workloads
2. Avoid CPU contention from other workloads
3. Limits give you the highest degree of control

ORACLE®
REAL-WORLD PERFORMANCE

# Thank You

——

**Mihajlo Tekic**
**John Zimmerman**

Real-World Performance Team
Oracle Database Development