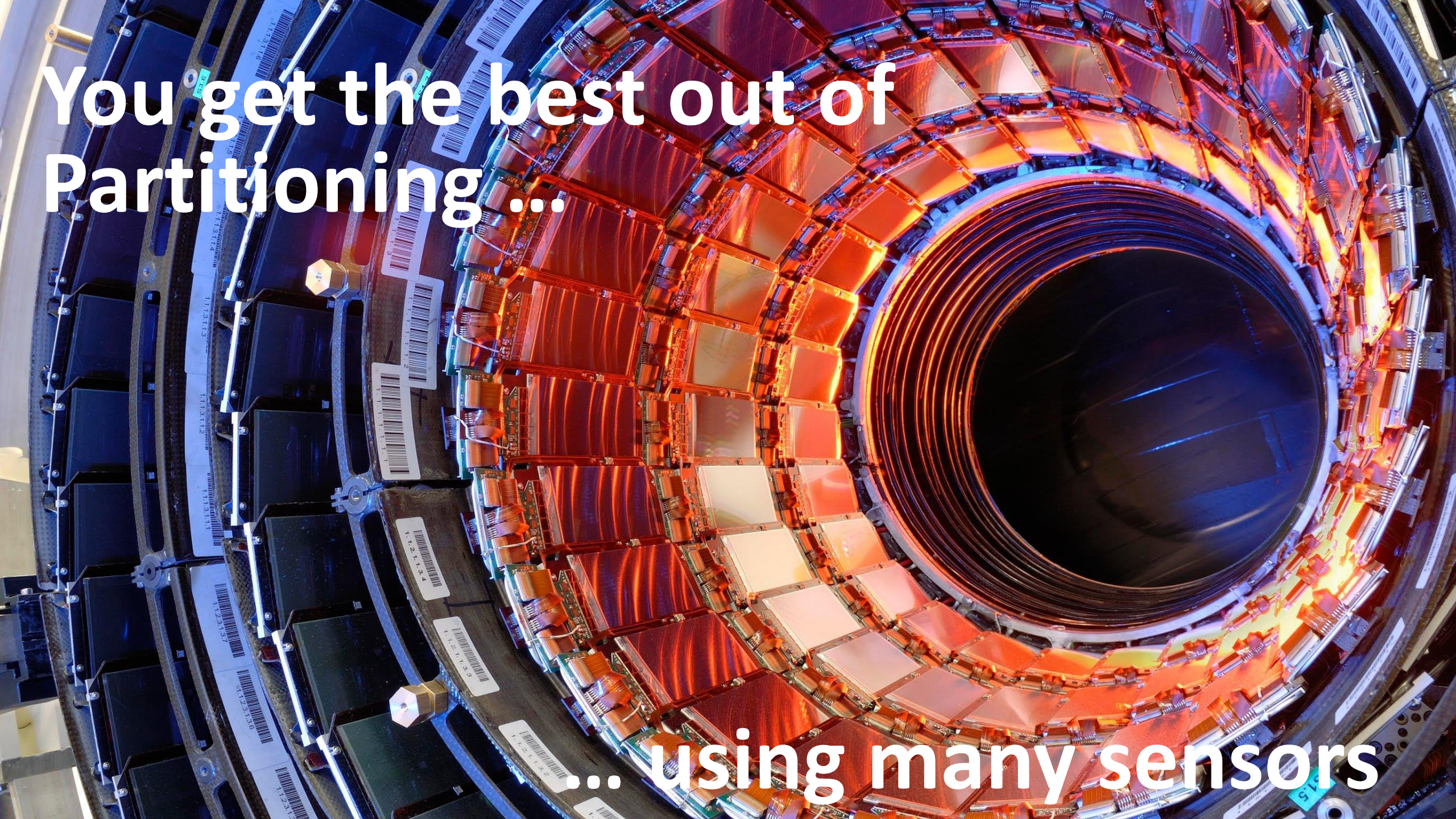


ORACLE®



You get the best out of
Partitioning ...

... using many sensors

Get the best out of Oracle Partitioning

A practical guide and reference

Thomas Teske

Thomas.teske@oracle.com

[@ThomasTeskeORCL](#)

Credits: Hermann Bär, PM for partitioning, created this extensive presentation. It has been extended to include Oracle 18c, the PET trick by Thomas Teske.

Please share your feedback (comments, suggestions, errors) with me – thank you!



Safe Harbor Statement

The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Oracle's products remains at the sole discretion of Oracle.

Before we start ..

- Oracle wants to hear from you!
 - There's still lots of ideas and things to do
 - Input steers the direction
- Let us know about
 - Interesting use cases and implementations
 - Enhancement requests
 - Complaints



Content

- [Flexibility has its price](#)
 - [What's New in 18?](#)
 - [What's New in 12.2?](#)

 - [Partitioning Summary](#)
 - [Partitioning Benefits](#)
 - [Partitioning Concepts](#)
- [Partitioning Methods](#)
 - [Range Partitioning](#)
 - [Hash Partitioning](#)
 - [List Partitioning](#)
 - [Interval Partitioning](#)
 - [Range versus Interval](#)
 - [Reference Partitioning](#)
 - [Interval Reference Partitioning](#)
 - [Virtual Column Based Partitioning](#)
- [Indexing of Partitioned Tables](#)
 - [Local Indexing](#)
 - [Global Non-Partitioned Indexing](#)
 - [Global Partitioned Indexing](#)
 - [Indexing for unique constraints and primary keys](#)
 - [Partial Indexing](#)
 - [Unusable versus partial Indexes](#)
 - [Partitioning for Performance](#)
 - [Composite Partitioning](#)
 - [Multi-Column Range Partitioning](#)
 - [Range Partitioned Hash Cluster](#)



More content

- Partition Maintenance
 - Maintenance on multiple Partitions
 - Cascading Truncate and Exchange for Reference Partitioning
 - Online Move Partition
 - Asynchronous Global Index Maintenance
- Stats Management for Partitioning
- Attribute Clustering/ Zone Maps
- Tips and tricks
 - Think about partitioning strategy
 - Physical and logical attributes

 - Eliminate hot spots
 - Smart partial exchange

 - Exchange with PK and UK
 - Partition Elimination Table (PET)



Flexibility has its price



Flexibility with Oracle Partitioning

- One million partitions –the more the better?
- Online operations – the holy grail?
- PMOPs over DML all the time?



One millions partitions – the more the better?

- More is not always better
 - Every partition represents metadata in the dictionary
 - Every partition increases the metadata footprint in the SGA
 - Large number of partitions can impact performance of catalog views
- Find your personal balance between the number of partitions and its average size
 - There is nothing wrong about single-digit GB sizes for a segment on “normal systems”
 - Consider more partitions \geq 5GB segment size



Online operations – the holy grail?

Plan for best time window

- Online operations are introduced to eliminate application transparency and therefore business impact
 - Not introduced to stop thinking about application workflow and design



Online operations – the holy grail?

Plan for best time window

- Online operations are introduced to eliminate application transparency and therefore business impact
 - Not introduced to stop thinking about application workflow and design
- **Concurrency is the main challenge**



Online operations – the holy grail?

Plan for best time window

- Online operations are introduced to sustain application transparency and therefore business impact
 - Not introduced to stop thinking about application workflow and design
- **Concurrency is the main challenge**
 - Minimize concurrent DML operations if possible
 - Require additional disk space and resources for journaling
 - Journal will be applied recursively after initial bulk move
 - The larger the journal, the longer the runtime
 - Concurrent DML has impact on compression efficiency
 - Best compression ratio with initial bulk move



PMOPs over DML all the time?

- Partition maintenance operations are a fast and efficient way to load or unload data



PMOPs over DML all the time?

- Partition maintenance operations are a fast and efficient way to load or unload data
- .. but it has its price:
 - Recursive DML to update partition metadata
 - Most commonly linear to number of involved partitions (tables and indexes), with exceptions
 - Cursor invalidation
 - Working hard on doing more fine-grained invalidation and incremental metadata invalidation/refresh



PMOPs over DML all the time?

- Partition maintenance operations are a fast and efficient way to load or unload data
- .. but it has its price:
 - Recursive DML to update partition metadata
 - Most commonly linear to number of involved partitions (tables and indexes), with exceptions
 - Cursor invalidation
 - Working hard on doing more fine-grained invalidation and incremental metadata invalidation/refresh
- DML is a viable alternative
 - Especially for smaller data volumes



What's New in Oracle Database 18c ?

**NEW IN
18^c**



The Dilemma continues

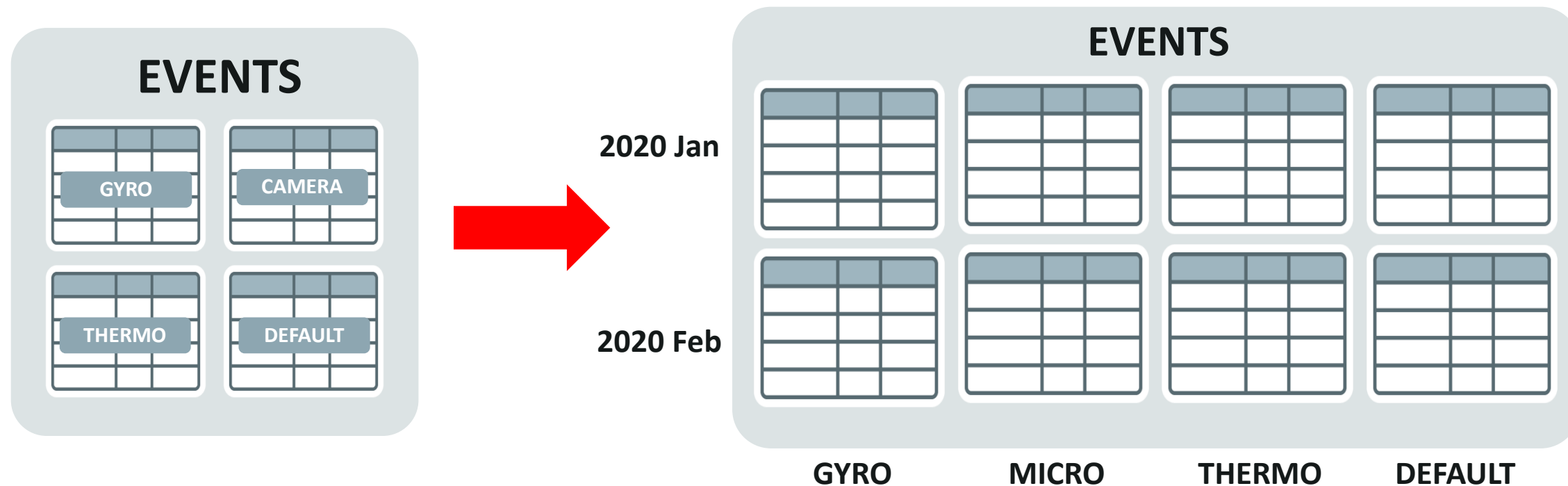
- Not everybody thinks big and starts small ...
 - ... so tables can start off small as non-partitioned ones
 - ... and they grow and grow
 - ... and they are used in a different way than expected
 - ... and their maintenance becomes a problem
 - ... and performance can get impacted

How to convert such tables without downtime?

- Now I have partitioning ...
 - ... but I chose the “wrong” type/granularity (for whatever reason)



Online Table Conversion of Partitioned Tables



- Completely non-blocking (online) DDL for table and indexes

Online Table Conversion of Partitioned Tables

- Indexes are converted and kept online throughout the conversion process
- Default indexing rules to provide minimal to no access change behavior
 - Almost identical than rules for conversion of non-partitioned table
 - Differences:
 - Local indexes stay local if any of the partition keys of the two dimensions is included
 - Global prefixed partitioned indexes will be converted to local partitioned indexes
- Full flexibility for indexes, following today's rules
 - Override whatever you want to see being changed



Online Table Conversion of Partitioned Tables

Syntax Example

```
CREATE TABLE EVENTS ( run_id      NUMBER,
                      sensor_type VARCHAR2 (50), ... )
PARTITION BY LIST ( ... )

ALTER TABLE EVENTS MODIFY
PARTITION BY RANGE ( run_id      )
SUBPARTITION BY LIST ( sensor_type )...
UPDATE INDEXES
  (i1_run_id GLOBAL,
   i2_sensor LOCAL,
   i3 GLOBAL PARTITION BY RANGE ( ... )
     (PARTITION p0100 VALUES LESS THAN (100000),
      PARTITION p1500 VALUES LESS THAN (1500000),
      PARTITION pmax  VALUES LESS THAN (MAXVALUE)))
ONLINE;
```



Asynchronous Global Index Maintenance

- Usable global indexes after DROP and TRUNCATE PARTITION without the need of index maintenance
 - Affected partitions are known internally and filtered out at data access time
- DROP and TRUNCATE become fast, metadata-only operations
 - Significant speedup and reduced initial resource consumption
- Delayed Global index maintenance
 - Deferred maintenance through ALTER INDEX REBUILD | COALESCE
 - Automatic cleanup using a scheduled job



Asynchronous Global Index Maintenance

Before

```
SQL> select count(*) from pt partition for (9999);
```

COUNT(*)
25341440

```
Elapsed: 00:00:01.00
SQL> select index_name, status, orphaned_entries from user_indexes;
```

INDEX_NAME	STATUS	ORPHANED_ENTRIES
I1_PT	VALID	NO

```
Elapsed: 00:00:01.04
SQL>
SQL> alter table pt drop partition for (9999) update indexes;
```

Table altered.

```
Elapsed: 00:02:04.52
```

```
SQL>
SQL> select index_name, status, orphaned_entries from user_indexes;
```

INDEX_NAME	STATUS	ORPHANED_ENTRIES
I1_PT	VALID	NO

```
Elapsed: 00:00:00.10
```

After

```
SQL> select count(*) from pt partition for (9999);
```

COUNT(*)
25341440

```
Elapsed: 00:00:00.98
SQL> select index_name, status, orphaned_entries from user_indexes;
```

INDEX_NAME	STATUS	ORPHANED_ENTRIES
I1_PT	VALID	NO

```
Elapsed: 00:00:00.33
SQL>
SQL> alter table pt drop partition for (9999) update indexes;
```

Table altered.

```
Elapsed: 00:00:00.04
```

```
SQL>
SQL> select index_name, status, orphaned_entries from user_indexes;
```

INDEX_NAME	STATUS	ORPHANED_ENTRIES
I1_PT	VALID	YES

```
Elapsed: 00:00:00.05
```

Asynchronous Global Index Maintenance

- Initial implementation of maintenance package
 - Always use INDEX COALESCE CLEANUP
 - Rely on parallelism of index
- Enhancements added to latest release
 - Choice of INDEX COALESCE CLEANUP or “classical” index cleanup
 - Choice of parallelism for maintenance operation
- **Classical cleanup recommended for more frequent index cleanup**
 - Seems to be the more common customer use case, thus the new default
- Functionality available for 12.1 through bug 24515918



What's New in Oracle Database 12c Release 2 ?

**NEW IN
12.2**



What's New in 12.2

- New Core Functionality Features

- Auto-list partitioning
- Multi-column list partitioning
- Partitioned external tables

- New Performance Features

- Online partition maintenance operations
- Online table conversion to partitioned table
- Reduced cursor invalidations for DDL's

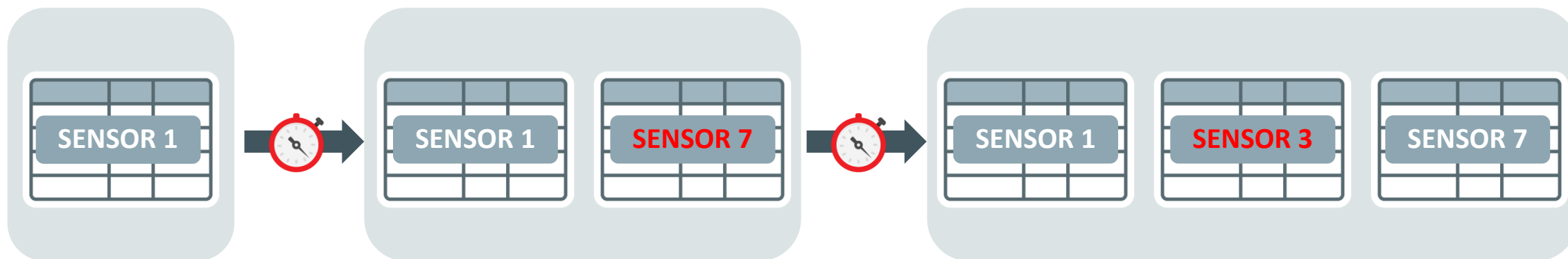
- New Manageability Features

- Filtered partition maintenance operations
- Read only partitions
- Create table for exchange

New Core Functionality Features

Auto- List Partitioning

Auto-List Partitioning



- Partitions are created automatically as data arrives
 - Extension to LIST partitioning
 - Every distinct partition key value will be stored in separate partition

Details of Auto-List Strategy

- Automatically creates new list partitions that contain one value per partition
 - Only available as top-level partitioning strategy in 12.2.0.1
- No notion of default partition
- System generated partition names for auto-created partitions
 - Use FOR VALUES clause for deterministic [sub]partition identification
- Can evolve list partitioning into auto-list partitioning
 - Only requirement is having no DEFAULT partition
 - Protection of your investment into a schema



Auto-List Partitioned Table

Syntax example

```
CREATE TABLE EVENTS ( sensor_type  VARCHAR2 (50) ,  
                       channel      VARCHAR2 (50) , ...)  
PARTITION BY LIST (sensor_type) AUTOMATIC  
( partition p1 values ('GYRO') );
```



Auto-List is not equivalent to List + DEFAULT

- Different use case scenarios
- List with DEFAULT partitioning
 - Targeted towards multiple large distinct list values plus “not classified”
- Auto-list partitioning
 - Expects ‘critical mass of records’ per partition key value
 - Could be used as pre-cursor state for using List + DEFAULT



Auto-List is not equivalent to List + DEFAULT

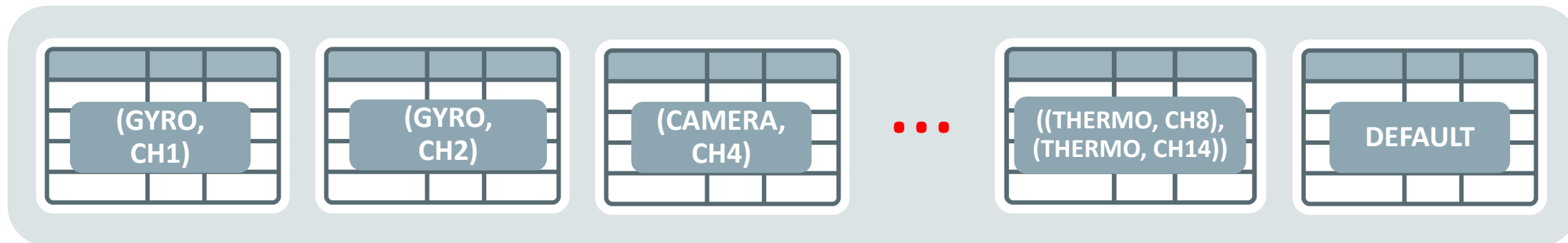
- Different use case scenarios
- List with DEFAULT partitioning
 - Targeted towards multiple large distinct list values plus “not classified”
- Auto-list partitioning
 - Expects ‘critical mass of records’ per value
 - Could be used as pre-cursor state for using List + DEFAULT
- **.. Plus they are functionally conflicting and cannot be used together**
 - Either you get a new partition for a new partition key value
 - .. Or “dump” it in the catch-it-all bucket



Multi-Column List Partitioning

Trick related to MCLP:
[Partition Elimination Table \(PET\)](#)

Multi-Column List Partitioning



- Data is organized in lists of multiple values (multiple columns)
 - Individual partitions can contain sets of multiple values
 - Functionality of DEFAULT partition (catch-it-all for unspecified values)
- Ideal for segmentation of distinct value tuples, e.g. (sensor_type, channel, ...)

Details of Multi-Column List Strategy

- Allow specification of more than one column as partitioning key
 - Up to 16 partition key columns
 - Each set of partitioning keys must be unique
- Notation of one DEFAULT partition
- Functional support
 - Supported as both partition and sub-partition strategy
 - Support for heap tables
 - Support for external tables
 - Supported with Reference Partitioning and Auto-List



Multi-Column List Partitioned Table

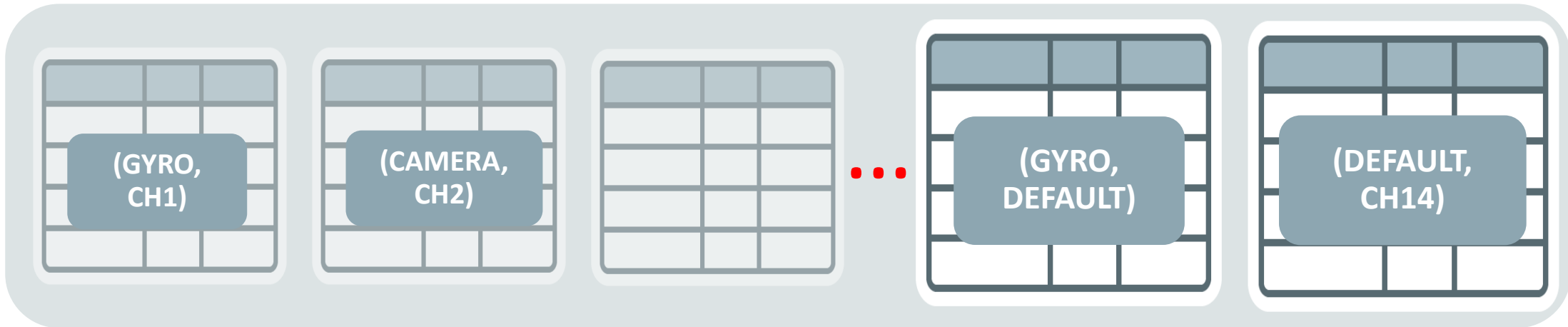
Syntax example

```
CREATE TABLE EVENTS( sensor_type  VARCHAR2(50),
                      channel      VARCHAR2(50), ...)
PARTITION BY LIST (sensor type, channel)
( partition p1 values ('GYRO','CH1'),
  partition p2 values ('GYRO','CH2'),
  partition p3 values ('CAMERA','CH4'),
  ...
  partition p44 values (('THERMO','CH8'),
                       ('THERMO','CH14')),
  partition p45 values (DEFAULT)
);
```



Multi-Column List Partitioning

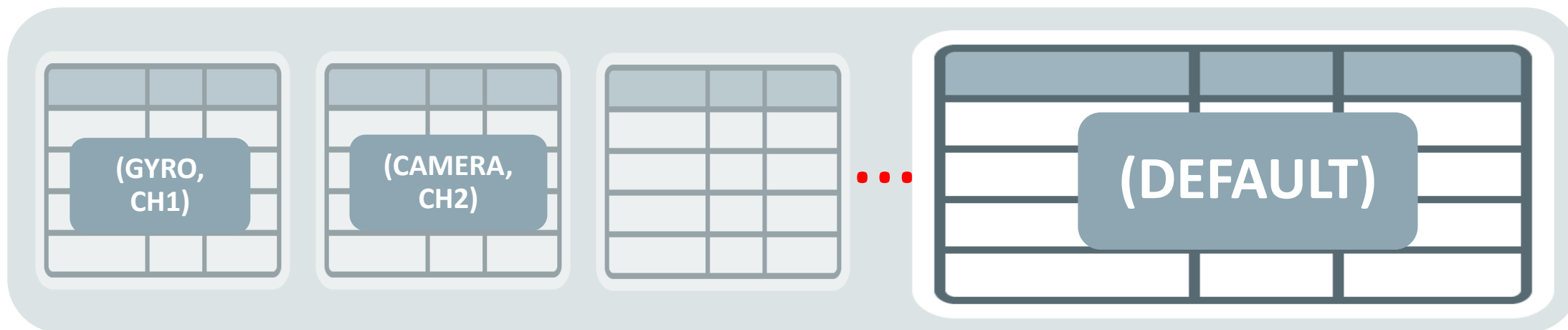
What if there was a DEFAULT per column?



- Where do we store (GYRO, CH12) ????

Multi-Column List Partitioning

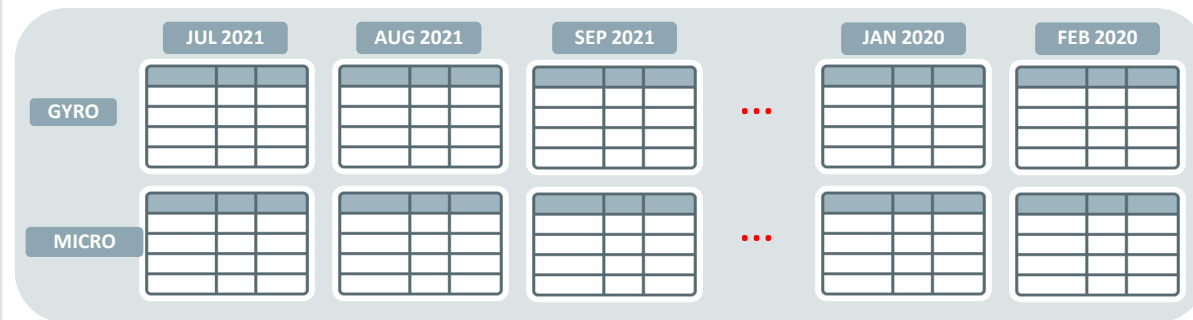
What if there was a DEFAULT per column?



- Where do we store (GYRO, CH12) ????
 - In the one-and-only DEFAULT partition

Multi-column List Partitioning prior to 12.2

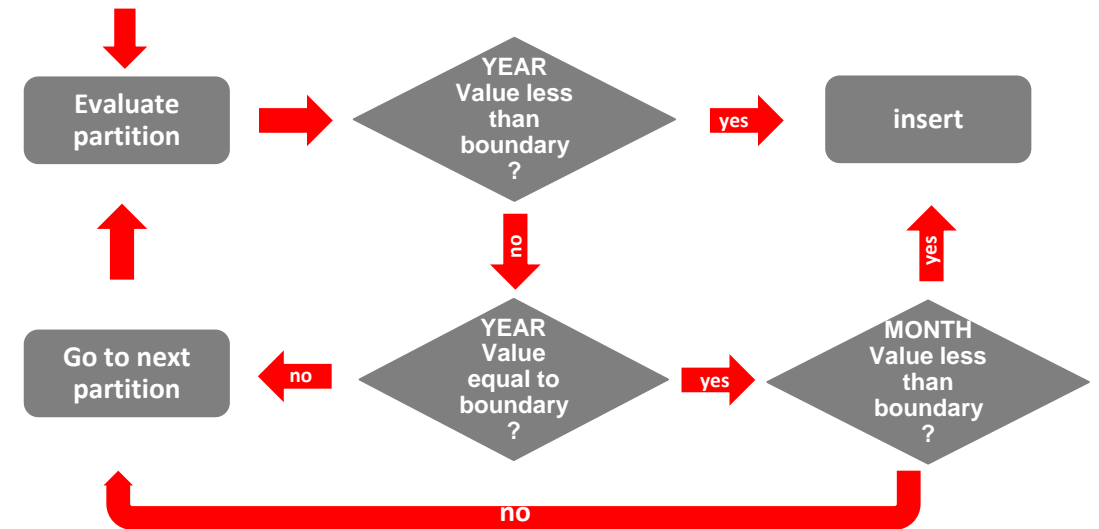
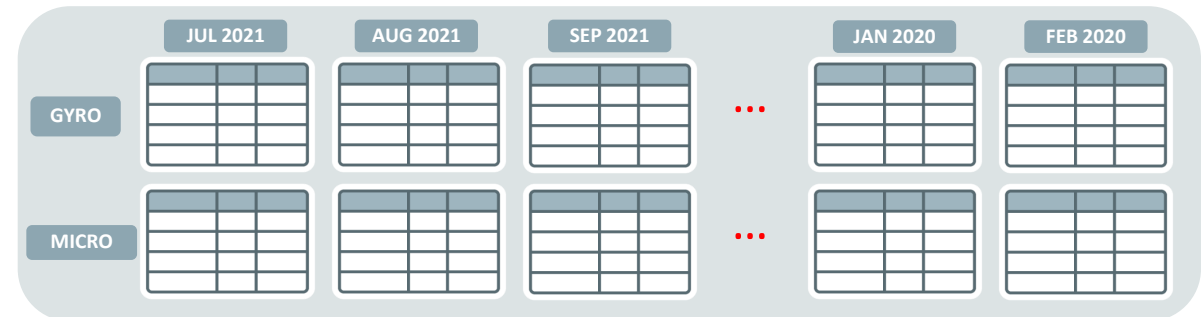
- List – List partitioning
 - Almost equivalent
 - Only two columns as keys (two levels)
 - Conceptual symmetrical



Multi-column List Partitioning prior to 12.2

- List – List partitioning
 - Almost equivalent
 - Only two columns as key (two levels)
 - Conceptual symmetrical

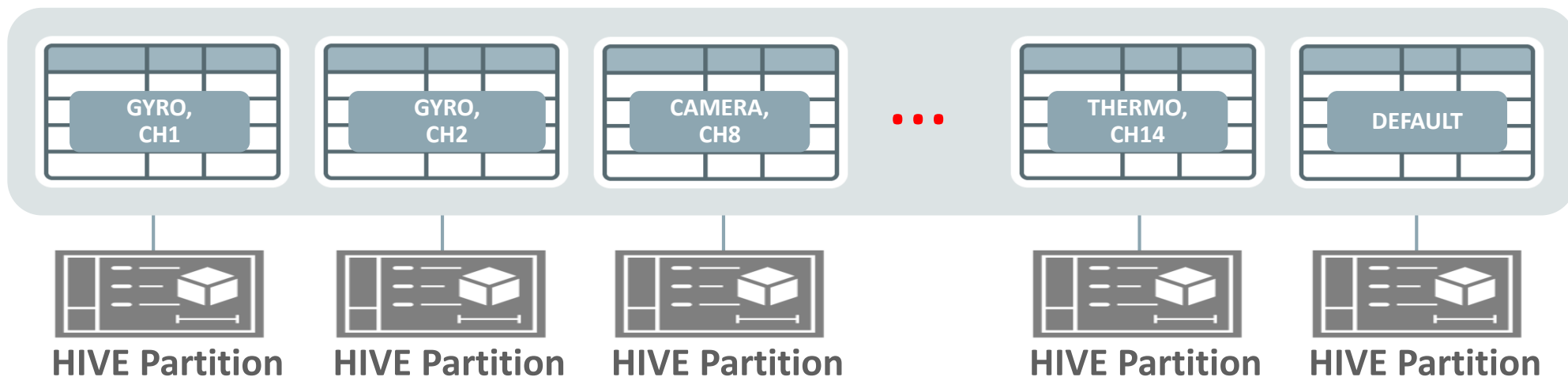
- Multi-column range partitioning
 - NOT equivalent
 - Hierarchical evaluation of predicates only in case of disambiguity



Partitioned External Tables

[External tables documentation](#)

Partitioned External Tables



- External tables can be partitioned, using all partitioning techniques
 - Multi-column partitioning optimally suited for partitioned HIVE tables
- Partition pruning and limited partition maintenance
 - Support of add partition, drop partition, exchange partition

New Performance Features

Online Partition Maintenance Operations

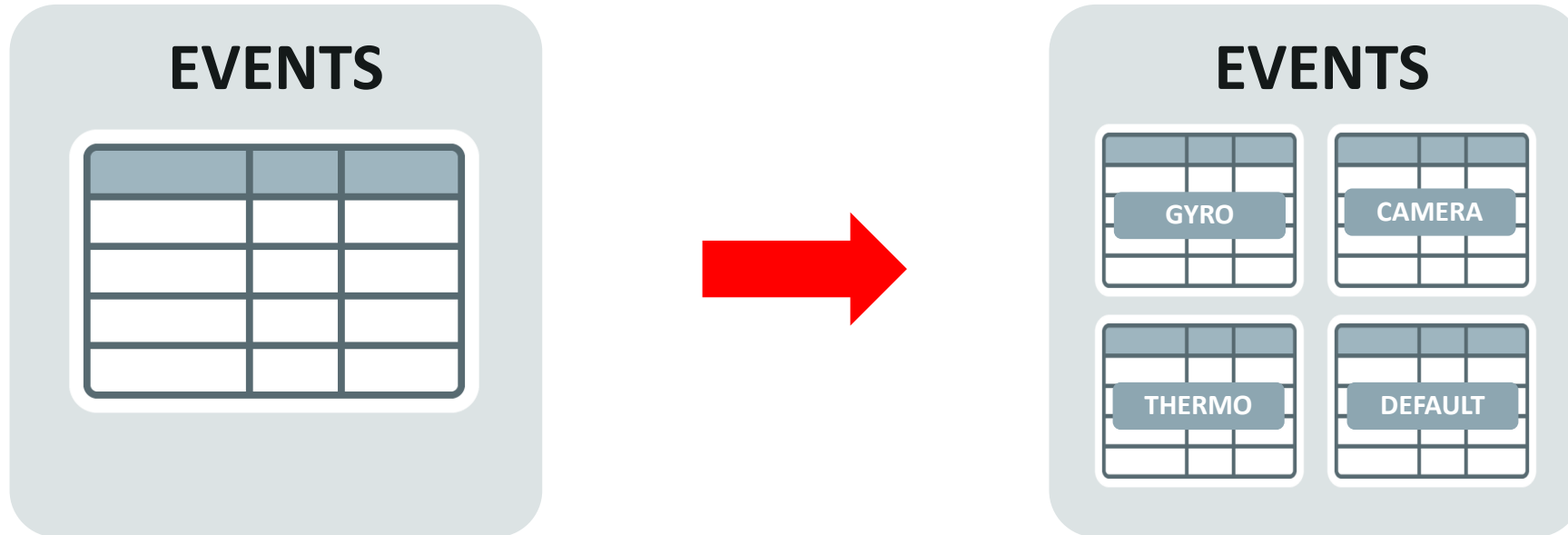
More Online DDL functionality

- Alter table modify non-partitioned table to partitioned table
- Alter table move online for heap tables
- Alter table split partition online



Online Table Conversion to Partitioned Table

Online Table Conversion



- Completely non-blocking (online) DDL

Online Table Conversion

Syntax Example

```
CREATE TABLE EVENTS ( sensor_grp  VARCHAR2 (50),  
                       channel    VARCHAR2 (50), ... );
```

```
ALTER TABLE EVENTS MODIFY  
PARTITION BY LIST ( sensor_grp )  
  (partition p1 values ('GYRO_GRP'),  
   partition p2 values ('CAMERA_GRP'),  
   partition p3 values ('THERMO_GRP'),  
   partition p4 values (DEFAULT))
```

```
UPDATE INDEXES ONLINE;
```



Online Table Conversion

Indexing

- Indexes are converted and kept online throughout the conversion process
- Full flexibility for indexes, following today's rules
- Default indexing rules to provide minimal to no access change behavior
 - Global partitioned indexes will retain the original partitioning shape.
 - Non-prefixed indexes will become global non-partitioned indexes.
 - Prefixed indexes will be converted to local partitioned indexes.
 - Bitmap indexes will become local partitioned indexes



Reduced Cursor Invalidations for DDL's

Reduced Cursor Invalidations for DDL's

- Reduces the number of hard parses caused by DDL's
 - If hard parses are unavoidable, workload is spread over time
- New optional clause “[DEFERRED | IMMEDIATE] INVALIDATION” for several DDL's
 - If DEFERRED, Oracle will avoid invalidating dependent cursors when possible
 - If IMMEDIATE, Oracle will immediately invalidate dependent cursors
 - If neither, CURSOR_INVALIDATION parameter controls default behavior
- Supported DDL's:
 - Create, drop, alter index
 - Alter table column operations
 - Alter table segment operations
 - Truncate table



Reduced Cursor Invalidations for DDL's

Syntax Example

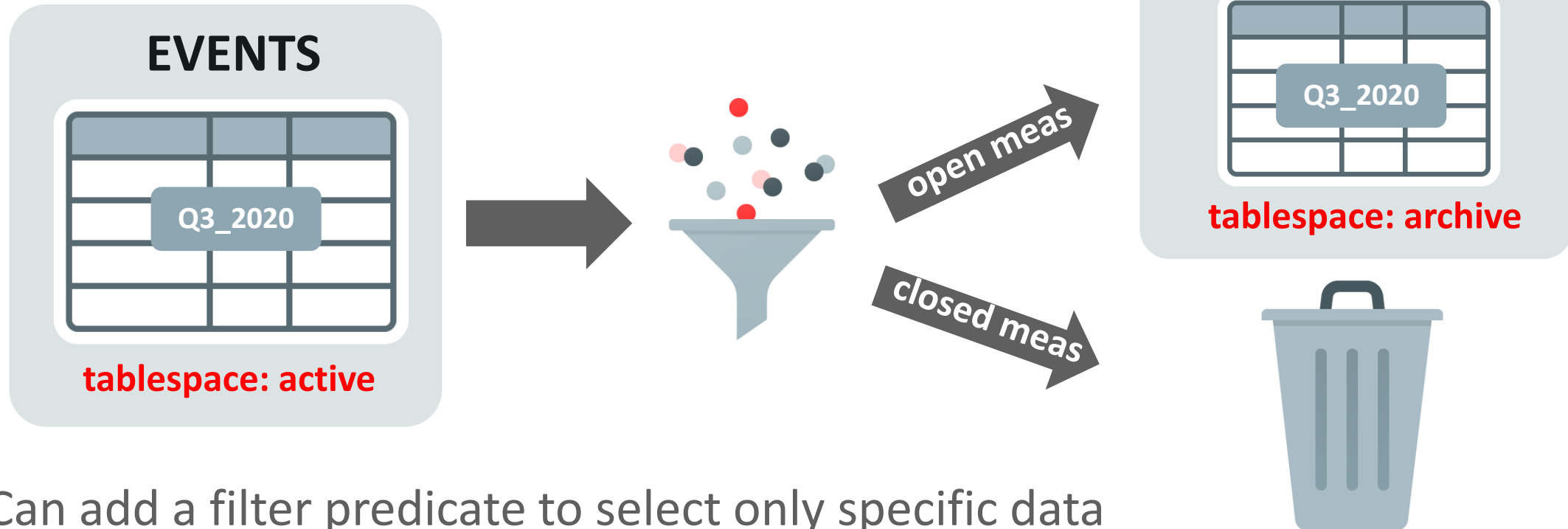
```
DROP INDEX meas_campaign DEFERRED INVALIDATION;
```

New Manageability Features

Filtered Partition Maintenance Operations

Filtered Partition Maintenance Operations

Move Partition Example



- Can add a filter predicate to select only specific data
- Combines data maintenance with partition maintenance

Details of Filtered Partition Maintenance Operations

- Can specify a single table filter predicate to MOVE, SPLIT and MERGE operations
 - Specification must be consistent across all partition maintenance
 - Specification needs to clearly specify the data of interest
- Specification will be added to the recursively generated CTAS command for the creation of the various new partition or sub-partitions segments
- Filter predicates work for both offline and new online PMOP's



Filtered Partition Maintenance Operations

Move Partition Syntax Example

```
ALTER TABLE orders MOVE PARTITION q3_2020  
TABLESPACE archive  
INCLUDING ROWS WHERE meas state = 'open';
```



Filtered Partition Maintenance Operations

Move Partition Syntax Example

```
ALTER TABLE orders MOVE PARTITION q3_2020  
TABLESPACE archive  
INCLUDING ROWS WHERE meas state = 'open';
```

.. and what about online?



Filtered Partition Maintenance Operation

DML Behavior for online operations

- Filter condition is NOT applied to ongoing concurrent DML

```
INCLUDING ROWS WHERE meas_state = 'open'
```



Filtered Partition Maintenance Operation

DML Behavior for online operations

- Filter condition is NOT applied to ongoing concurrent DML

```
INCLUDING ROWS WHERE meas_state = 'open'
```

- Inserts will always go through

```
INSERT VALUES (meas_state = 'closed')
```



Filtered Partition Maintenance Operation

DML Behavior for online operations

- Filter condition is NOT applied to ongoing concurrent DML

```
INCLUDING ROWS WHERE meas_state = 'open'
```

- Inserts will always go through

```
INSERT VALUES (meas_state = 'closed')
```

- Deletes on included data will always go through

```
DELETE WHERE meas_state = 'open'
```

- Deletes on deleted data are void

```
DELETE WHERE meas_state = 'closed'
```



Filtered Partition Maintenance Operation

DML Behavior for online operations

- Filter condition is NOT applied to ongoing concurrent DML

```
INCLUDING ROWS WHERE meas_state = 'open'
```

- Inserts will always go through

```
INSERT VALUES (meas_state = 'closed')
```

- Deletes on included data will always go through

```
DELETE WHERE meas_state = 'open'
```

- Deletes on deleted data are void

```
DELETE WHERE meas_state = 'closed'
```

- Updates on included data always goes through

```
UPDATE set order_status = 'closed'  
WHERE meas_state = 'open'
```

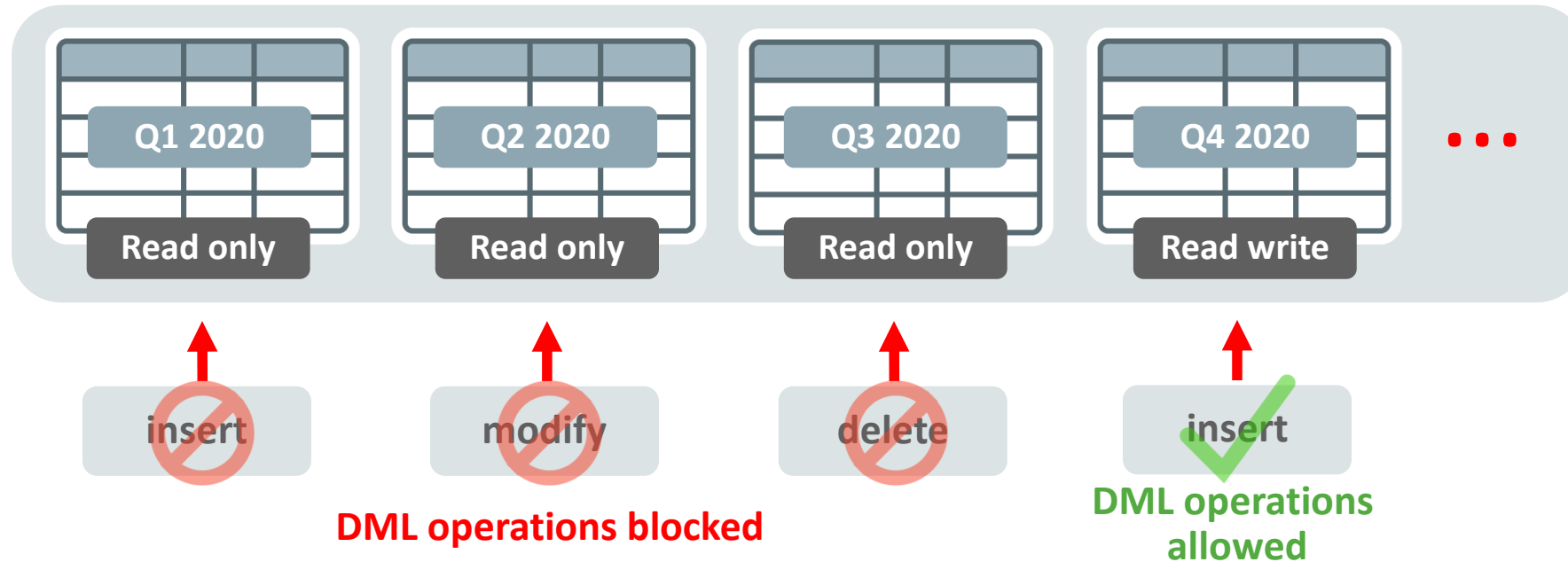
- Updates on excluded data are void

```
UPDATE set order_status = 'open'  
WHERE meas_state = 'closed'
```



Read Only Partitions

Read Only Partitions



- Partitions and sub-partitions can be set to read only or read write
- Any attempt to alter data in a read only partition will result in an error
- Ideal for protecting data from unintentional DML by any user or trigger

Details of Read Only Partitions

- Read only attribute guarantees data immutability
 - “SELECT <column_list> FROM <table>” will always return the same data set after a table or [sub]partition is set to read only
- If not specified, each partition and subpartition will inherit read only property from top level parent
 - Modifying lower level read only property will override higher level property
 - Alter tablespace has highest priority and cannot be overwritten
- Data immutability does not prevent all structural DDL for a table
 - ADD and MODIFY COLUMN are allowed and do not violate data immutability of existing data
 - Others like DROP/RENAME/SET UNUSED COLUMN are forbidden
 - DROP [read only] PARTITION forbidden, too - - violates data immutability of the table



Read Only Partitions

```
CREATE TABLE events ( event_id number,  
                       evnt_date DATE, ... ) read only  
  
PARTITION BY RANGE(event_date)  
( partition q1_2020 values less than ('2020-04-01'),  
  partition q2_2020 values less than ('2020-07-01'),  
  partition q3_2020 values less than ('2020-10-01'),  
  partition q4_2020 values less than ('2021-01-01') read write  
);
```

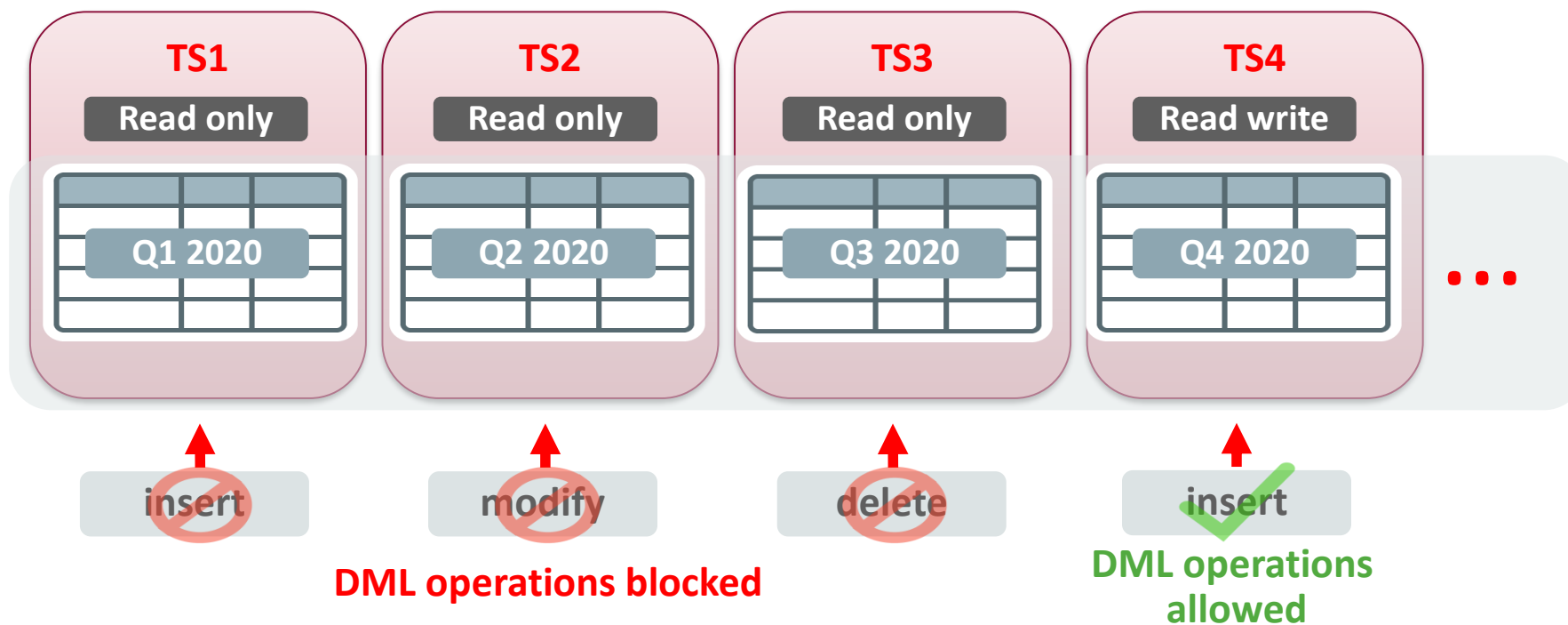


Read only tablespace vs. read only partitions

Read only partitions introduced in Oracle Database 12c Release 2

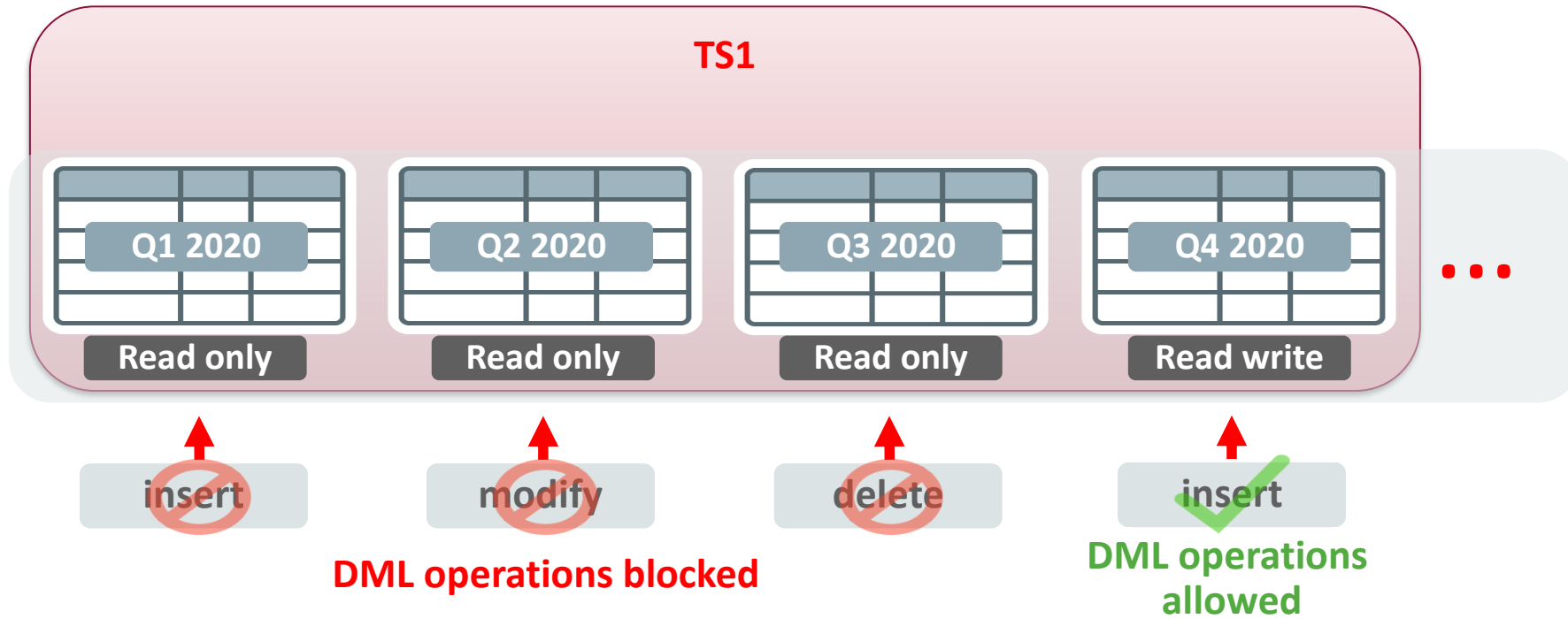


Read Only Tablespaces and Partitions



- Partitions and sub-partitions can be placed in read only tablespaces
- Any attempt to alter data in a read only tablespace will result in an error

Read Only Partitions



- **Partitions and sub-partitions can be set to read only or read write**
- Any attempt to alter data in a read only partition will result in an error

Read Only Object vs. Read Only Tablespace

- Read Only Tablespaces protect physical storage from updates
 - DDL operations that are not touching the storage are allowed
 - E.g. ALTER TABLE SET UNUSED, DROP TABLE
 - No guaranteed data immutability
- Read Only Objects protect data from updates
 - ‘Data immutability’
 - Does not prevent changes on storage
 - E.g. ALTER TABLE MOVE COMPRESS, ALTER TABLE MERGE PARTITIONS



Read Only Partitions

- Read only attribute guarantees data immutability
 - “SELECT <column_list> FROM <table>” will always return the same data set after a table or [sub]partition is set to read only
- Data immutability does not prevent all structural DDL for a table
 - ADD and MODIFY COLUMN are allowed and do not violate data immutability of existing data
 - Others like DROP/RENAME/SET UNUSED COLUMN are forbidden
 - DROP [read only] PARTITION forbidden, too - - violates data immutability of the table



Create Table for Exchange

Create Table for Exchange

- Simple DDL command
- Ensures both the semantic and internal table shape are identical so partition exchange command will always succeed
- Operates like a special CREATE TABLE AS SELECT operation
- Always creates an empty table



Create Table for Exchange

Syntax Example

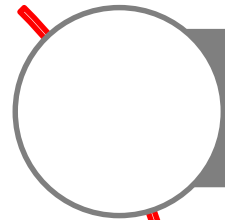
```
CREATE TABLE events_cp TABLESPACE ts_boson  
FOR EXCHANGE WITH events;
```



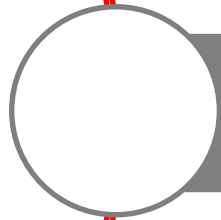
Partitioning Summary



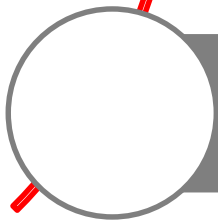
What is Oracle Partitioning?



Powerful functionality to logically divide objects into smaller pieces



Key requirement for large databases needing high performance and high availability



Driven by business requirements

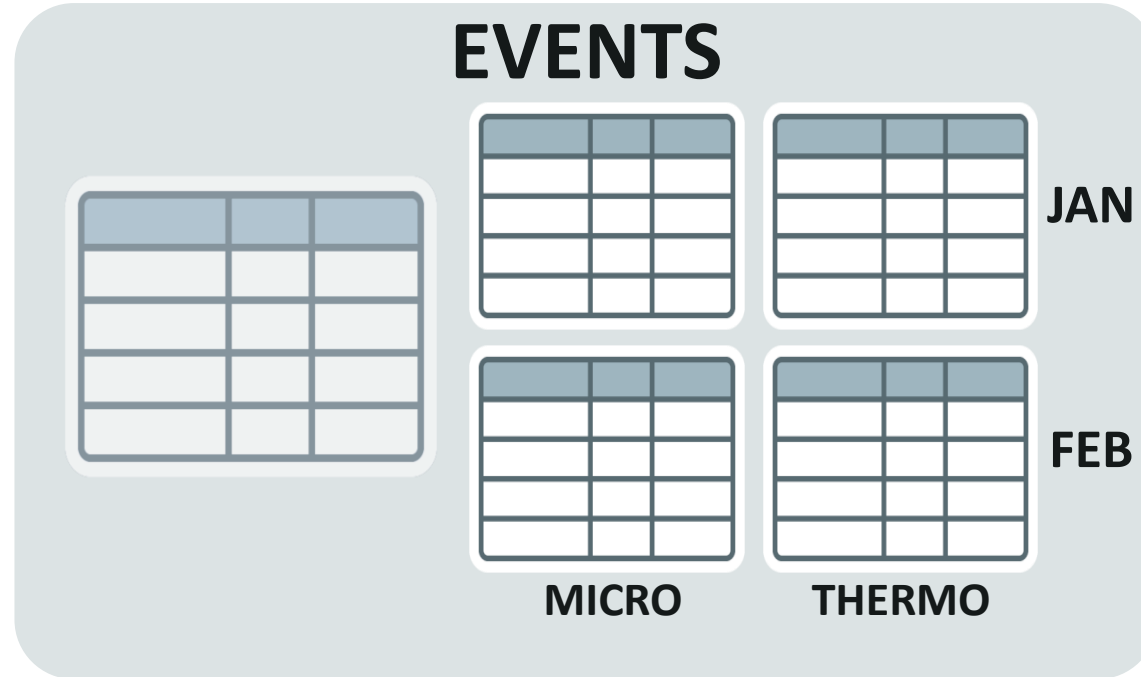
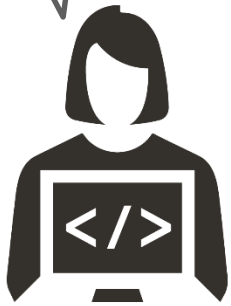
Why use Oracle Partitioning?

- ↑ Performance – lowers data access times
- ↑ Availability – improves access to critical information
- ↓ Costs – leverages multiple storage tiers
- ✓ Easy Implementation – requires no changes to applications and queries
- ✓ Mature Feature – supports a wide array of partitioning methods
- ✓ Well Proven – used by thousands of Oracle customers



The two Personalities of Partitioning

```
SELECT *  
FROM  
EVENTS;
```

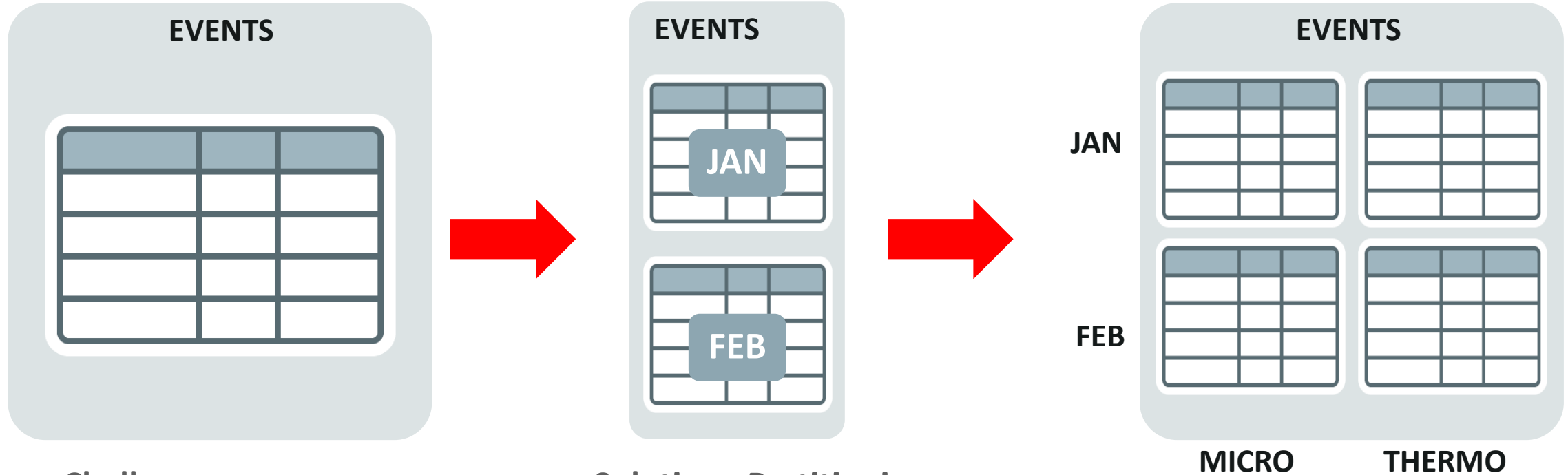


```
MOVE PARTITION  
COMPRESS  
READ ONLY;
```



How does Partitioning work?

Enables large databases and indexes to be split into smaller, more manageable pieces



Challenges:

Large tables are difficult to manage

Solution: Partitioning

- Divide and conquer
- Easier data management
- Improve performance

Partitioning Benefits



Increased Performance

Only work on the data that is relevant

Partitioning enables data management operations such as...

- Data loads, joins and pruning,
- Index creation and rebuilding,
- Optimizer statistics management,
- Backup and recovery

... at partition level instead of on the entire table

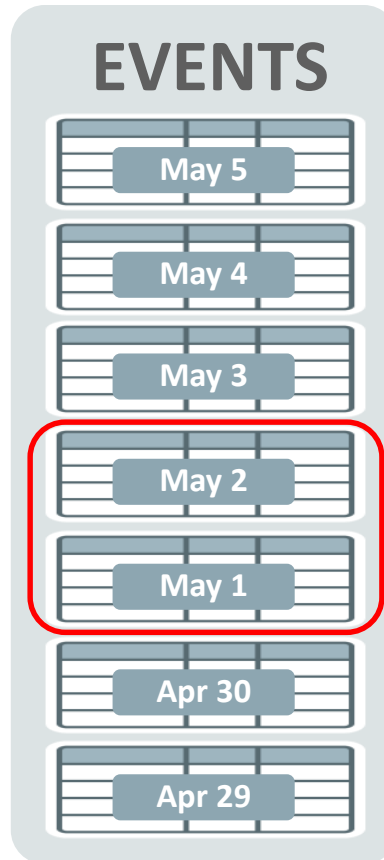
Result: Order of magnitude gains on performance



Increased Performance - Example

Partition Pruning

What are the total
EVENTS for May 1-2?

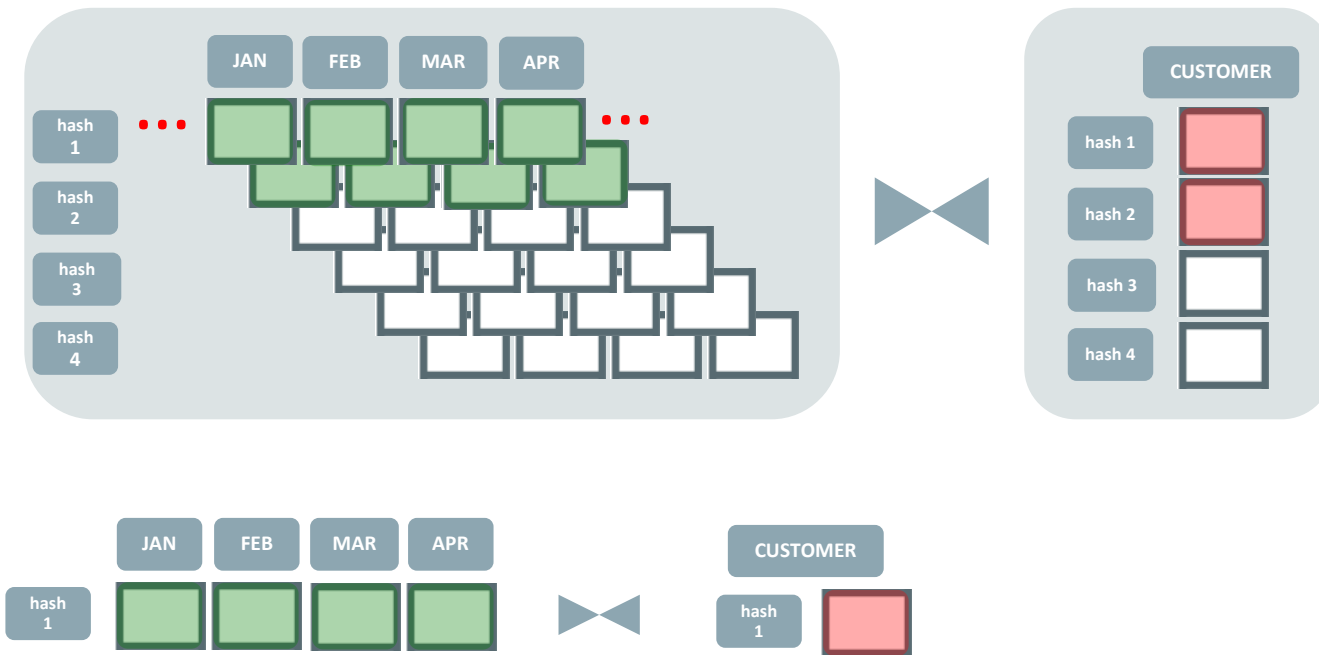


- Partition elimination
 - Dramatically reduces amount of data retrieved from storage
 - Performs operations only on relevant partitions
 - Transparently improves query performance and optimizes resource utilization



Increased Performance - Example

Partition-wise joins



A large join is divided into multiple smaller joins, executed in parallel

- # of partitions to join must be a multiple of DOP
- Both tables must be partitioned the same way on the join column

Decreased Costs

Store data in the most appropriate manner

Partitioning finds the balance between...

- data importance,
- storage performance,
- storage reliability,
- storage form

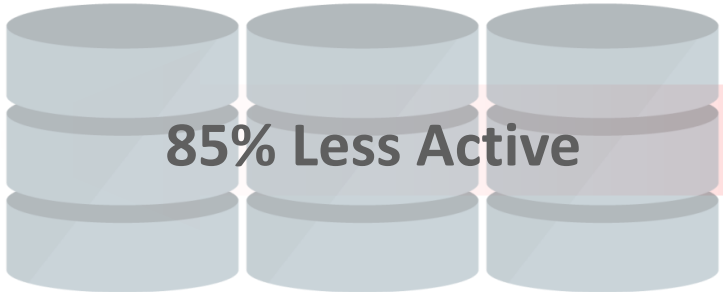
... allowing you to leverage multiple storage tiers

Result: Reduce storage costs by 2x or more



Decreased Costs - Example

Partition for Tiered Storage



Low End Storage Tier



Mid Storage Tier



High End Storage Tier

Increased Availability

Individual partition manageability

Partitioning reduces...

- Maintenance windows
- Impact of scheduled downtime and failures,
- Recovery times

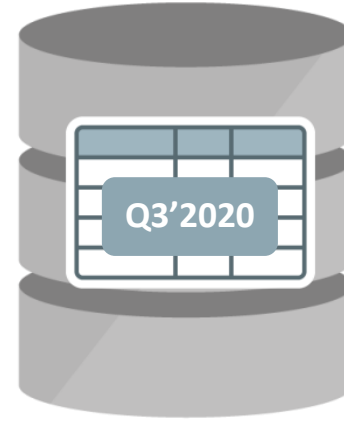
... if critical tables and indexes are partitioned

Result: Improves access to critical information



Increased Availability - Example

Partition for Manageability/Availability



Other partitions visible and usable

Easy Implementation

Transparent to applications

- Partitioning requires NO changes to applications and queries
 - Adjustments might be necessary to fully exploit the benefits of Partitioning



Mature, Well Proven Functionality

Over a decade of development

- Used by tens of thousands of Oracle customers
- Supports a wide array of partitioning methods



Oracle Partitioning

	Core functionality	Performance	Manageability
Oracle 8.0	Range partitioning Local and global Range indexing	Static partition pruning	Basic maintenance: ADD, DROP, EXCHANGE
Oracle 8i	Hash partitioning Range-Hash partitioning	Partition-wise joins Dynamic partition pruning	Expanded maintenance: MERGE
Oracle 9i	List partitioning		Global index maintenance
Oracle 9i R2	Range-List partitioning	Fast partition SPLIT	
Oracle 10g	Global Hash indexing		Local Index maintenance
Oracle 10g R2	1M partitions per table	Multi-dimensional pruning	Fast DROP TABLE
Oracle 11g	Virtual column based partitioning More composite choices Reference partitioning		Interval partitioning Partition Advisor Incremental stats mgmt
Oracle 11g R2	Hash-* partitioning Expanded Reference partitioning	“AND” pruning	Multi-branch execution (aka table or-expansion)
Oracle 12c R1	Interval-Reference partitioning	Partition Maintenance on multiple partitions Asynchronous global index maintenance Zone maps	Online partition MOVE Cascading TRUNCATE Partial indexing
Oracle 12c R2	Auto-list partitioning Multi-column list partitioning Partitioned external tables	Online partition maintenance operations Online table conversion to partitioned table Reduced cursor invalidations for DDL's	Filtered partition maintenance Read only partitions Create table for exchange



Partitioning Concepts



def Par•ti•tion

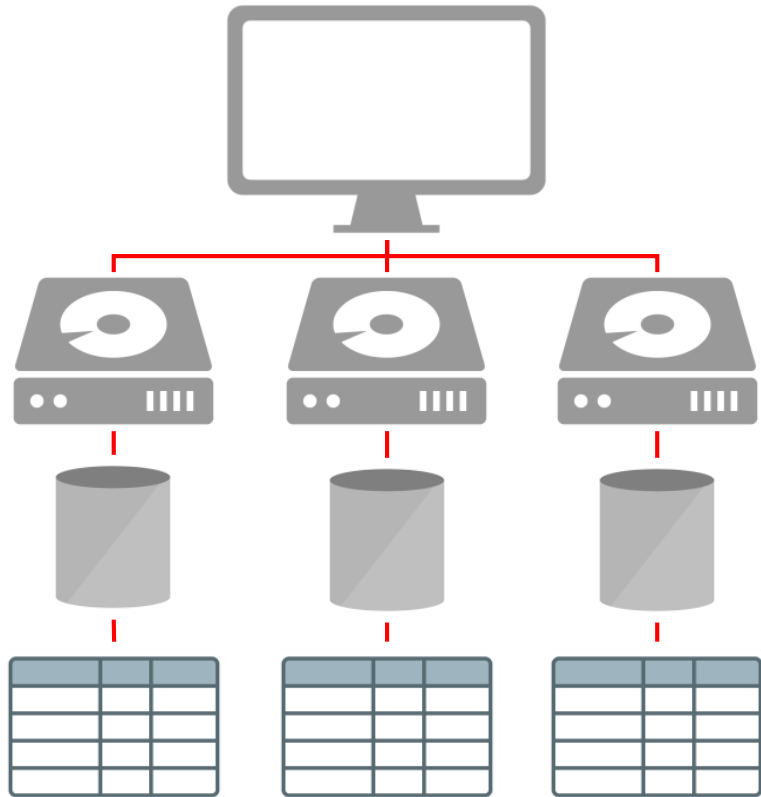
To divide (something) into parts

– “Merriam Webster Dictionary”



Physical Partitioning

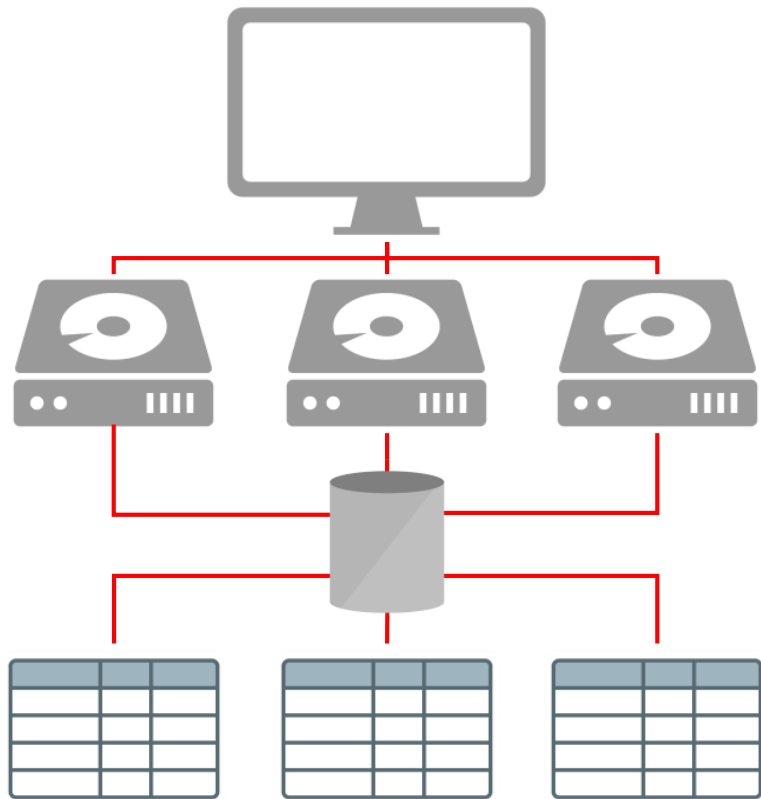
Shared Nothing Architecture



- Fundamental system setup requirement
 - Node owns piece of DB
 - Enables parallelism
- Number of partitions is equivalent to minimum required parallelism
 - Always needs HASH or random distribution
- Equally sized partitions per node required for proper load balancing

Logical Partitioning

Shared Everything Architecture - Oracle



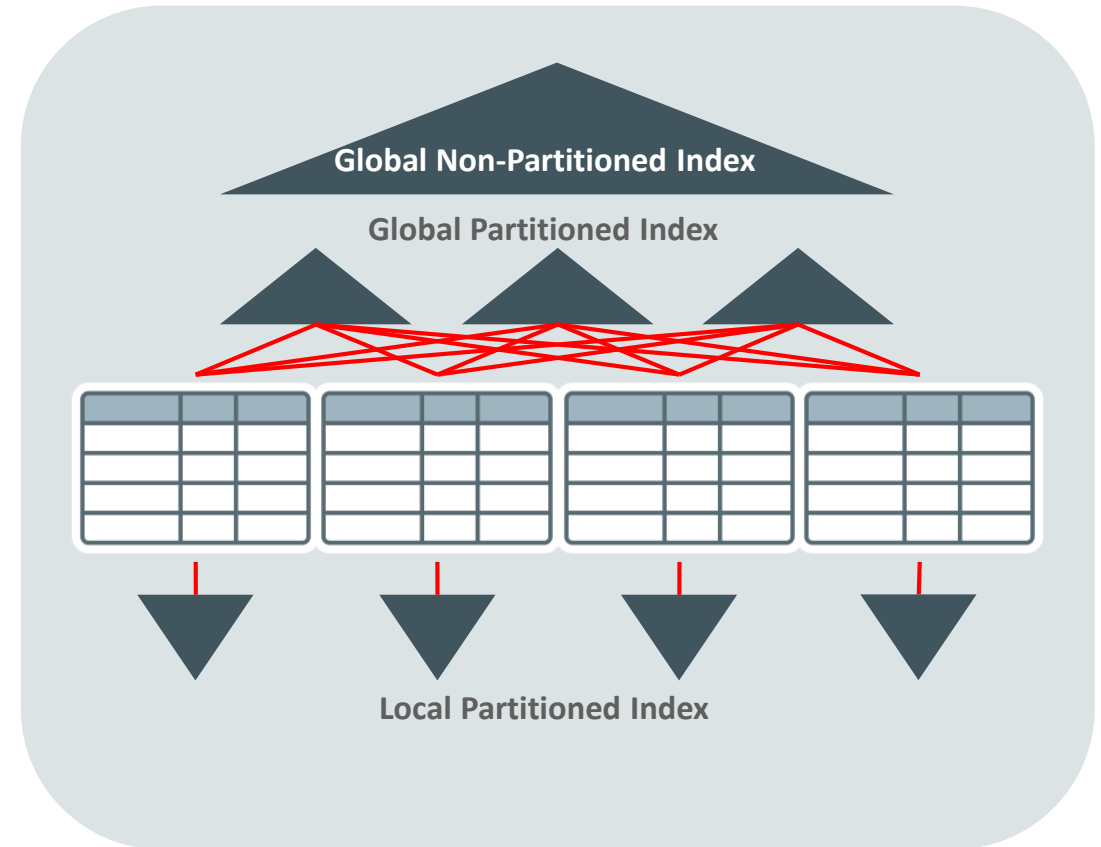
- Does not underlie any constraints
 - SMP, MPP, Cluster, Grid does not matter
- Purely based on the business requirement
 - Availability, Manageability, Performance
- Beneficial for every environment
 - Provides the most comprehensive functionality

Partitioning Methods



What can be partitioned?

- Tables
 - Heap tables
 - Index-organized tables
- Indexes
 - Global Indexes
 - Local Indexes
- Materialized Views
- Hash Clusters



Partitioning Methods

Single-level partitioning

- Range
- List
- Hash

Composite-level partitioning

- [Range | List | Hash | Interval] – [Range | List | Hash]

Partitioning extensions

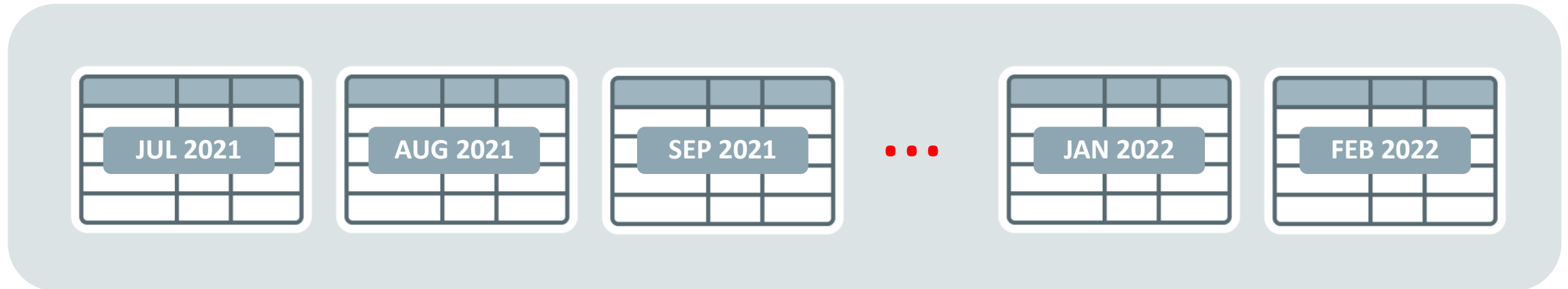
- Interval
- Reference
- Interval Reference
- Virtual Column Based
- Auto



Range Partitioning

Introduced in Oracle 8.0

Range Partitioning



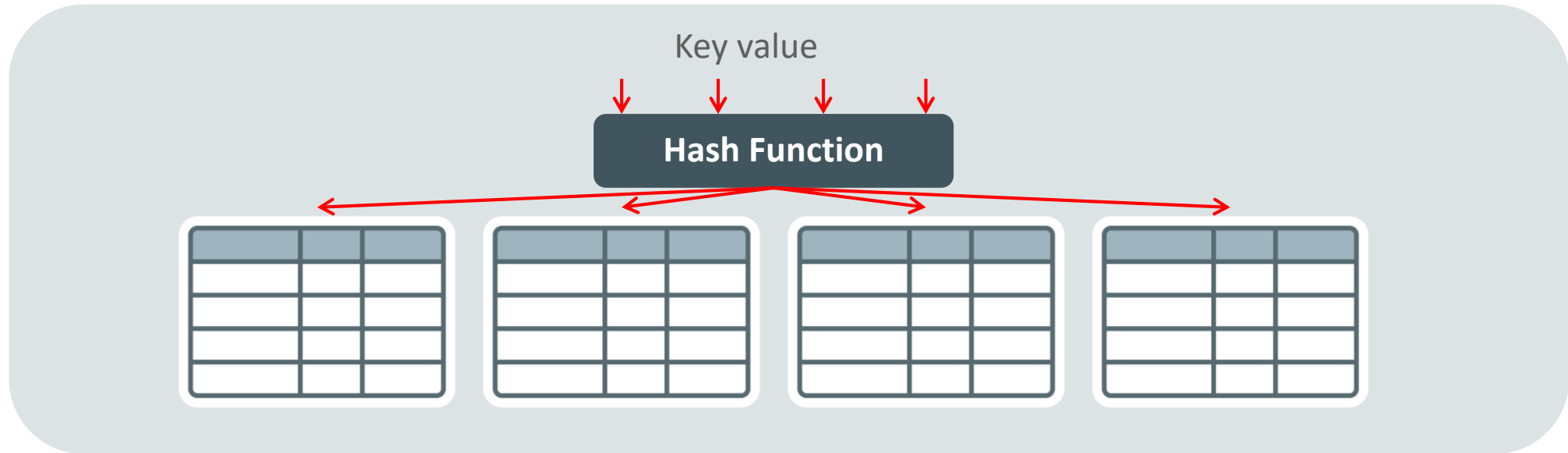
- Data is organized in ranges
 - Lower boundary derived by upper boundary of preceding partition
 - Split and merge as necessary
 - No gaps
- Ideal for chronological data

Hash Partitioning

Introduced in Oracle 8i (8.1)



Hash Partitioning



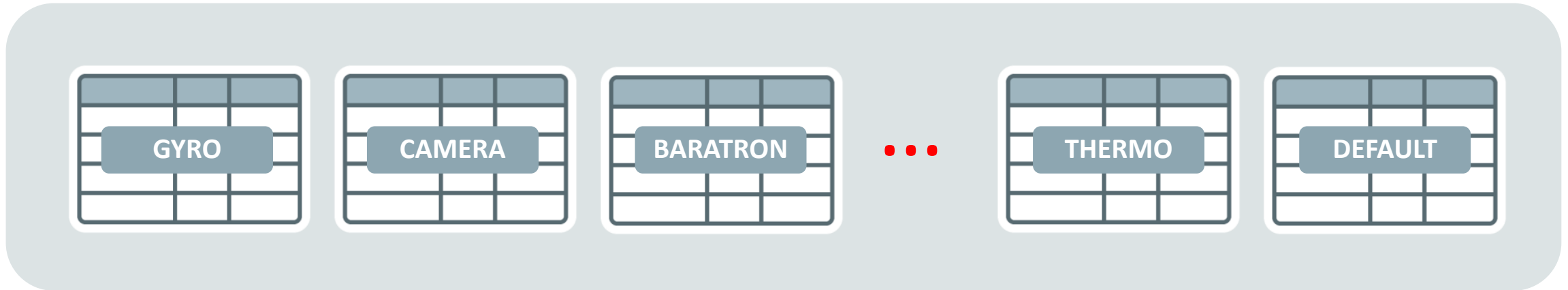
- Data is placed based on hash value of partition key
 - Number of hash buckets equals number of partitions
- Ideal for equal data distribution
 - Number of partitions should be a power of 2 for equal data distribution

List Partitioning

Introduced in Oracle 9i (9.0)



List Partitioning



- Data is organized in lists of values
 - One or more unordered distinct values per list
 - Functionality of DEFAULT partition (Catch-it-all for all unspecified values)
 - Check contents of DEFAULT partition – create new partitions as per need
- Ideal for segmentation of distinct values, e.g. region

Interval Partitioning

Introduced in Oracle 11g Release 1 (11.1)

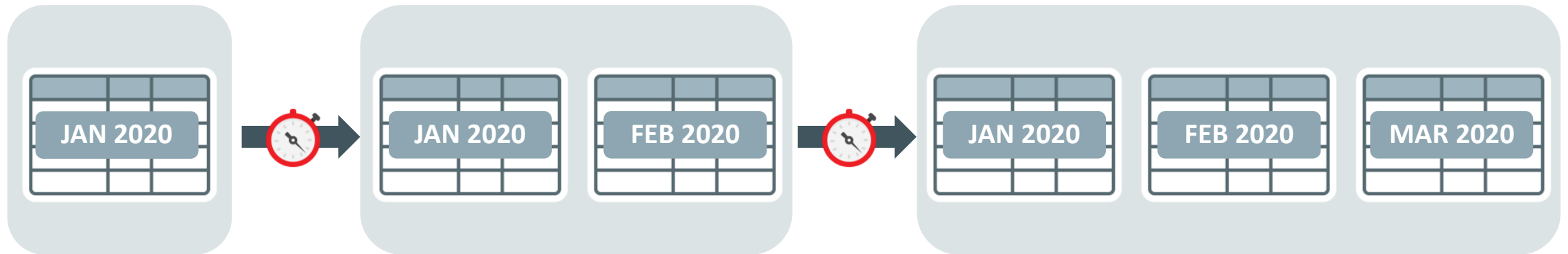
Interval Partitioning

- Extension to Range Partitioning
- Full automation for equi-sized range partitions
- Partitions are created as metadata information only
 - Start Partition is made persistent
- Segments are allocated as soon as new data arrives
 - No need to create new partitions
 - Local indexes are created and maintained as well

No need for any partition management



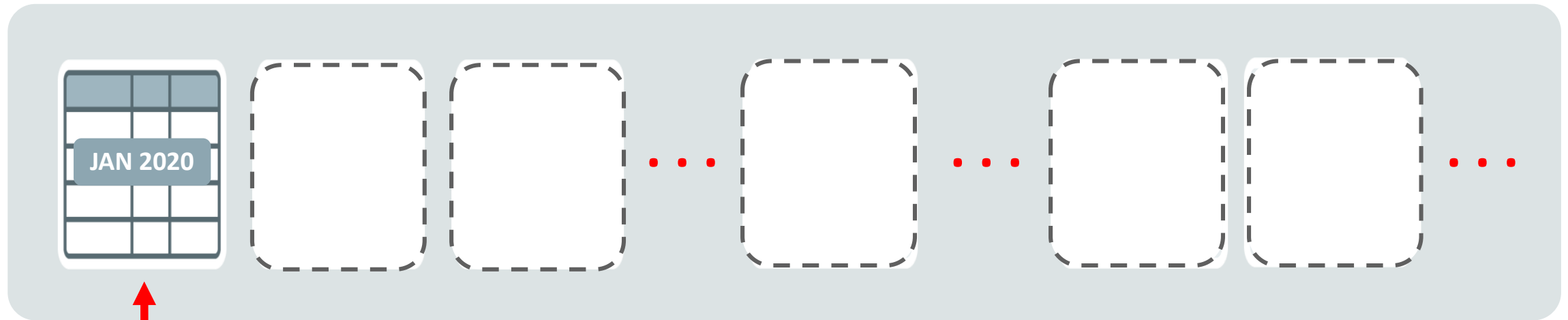
Interval Partitioning



- Partitions are created automatically as data arrives
 - Extension to RANGE partitioning

Interval Partitioning

As easy as One, Two, Three...

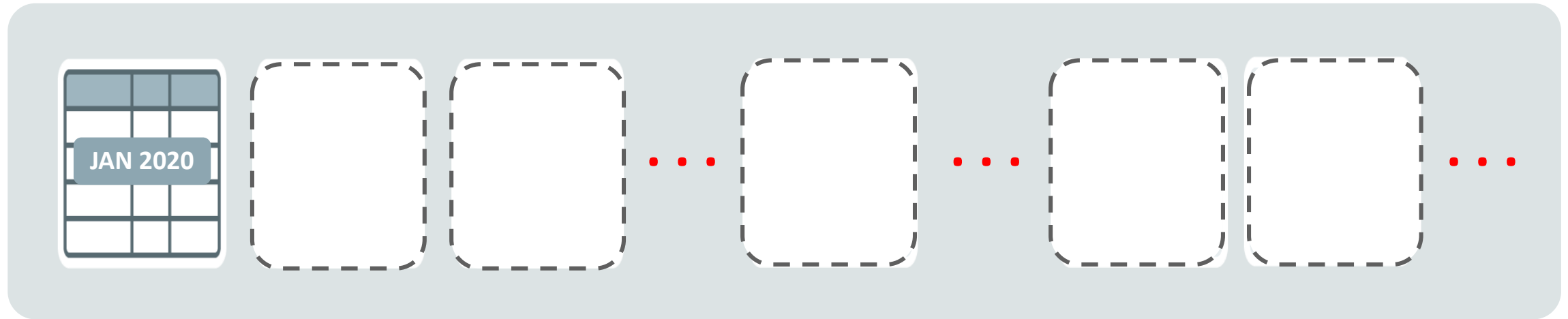


**First
partition
is created**

```
CREATE TABLE EVENTS (meas_date DATE, ...)  
PARTITION BY RANGE (meas_date)  
INTERVAL (NUMTOYMINTERVAL(1, 'month')  
(PARTITION p_first VALUES LESS THAN ('01-FEB-2020'));
```

Interval Partitioning

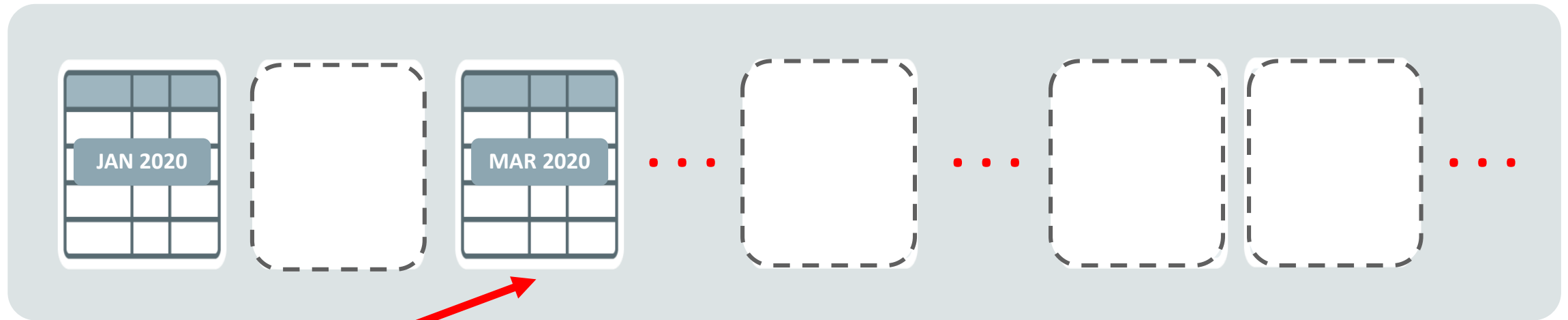
As easy as One, Two, Three...



Other partitions only exist in table metadata

Interval Partitioning

As easy as One, Two, Three...

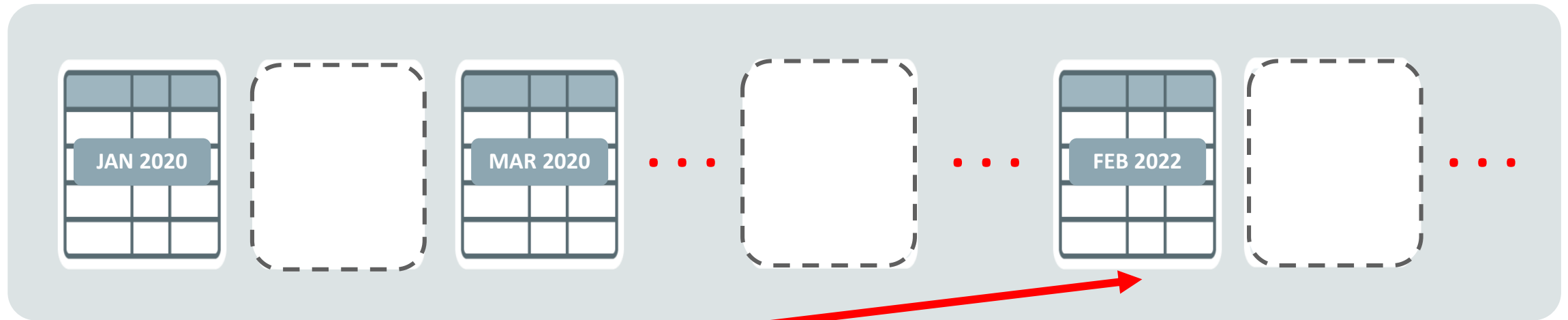


**New partition is
automatically instantiated**

```
INSERT INTO EVENTS (meas_date DATE, ...)
VALUES ('15-MAR-2020', ...);
```

Interval Partitioning

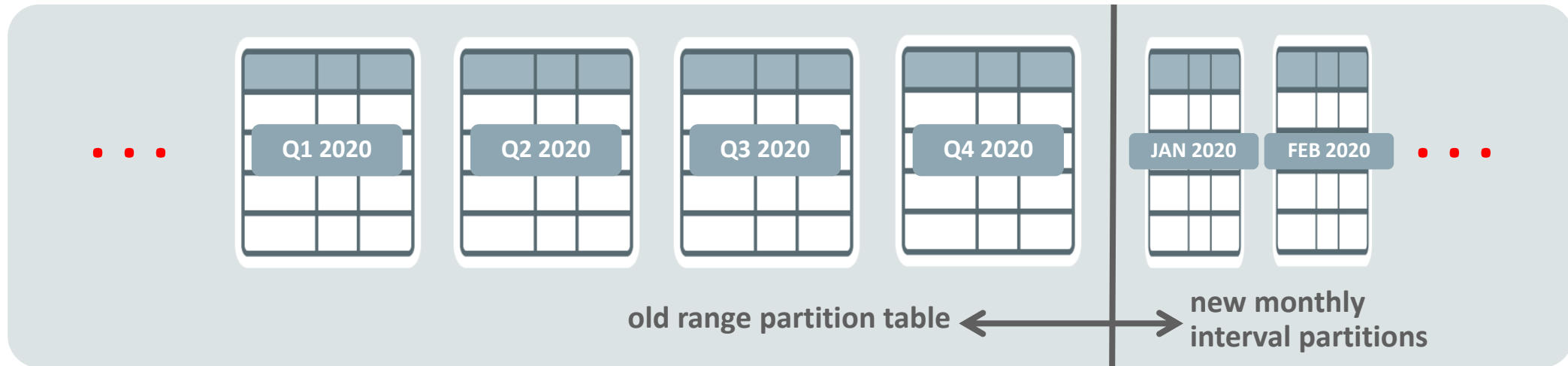
As easy as One, Two, Three...



**Whenever data for
a new partition arrives**

```
INSERT INTO EVENTS ( meas_date DATE, ... )  
VALUES ('04-FEB-2022', ...);
```

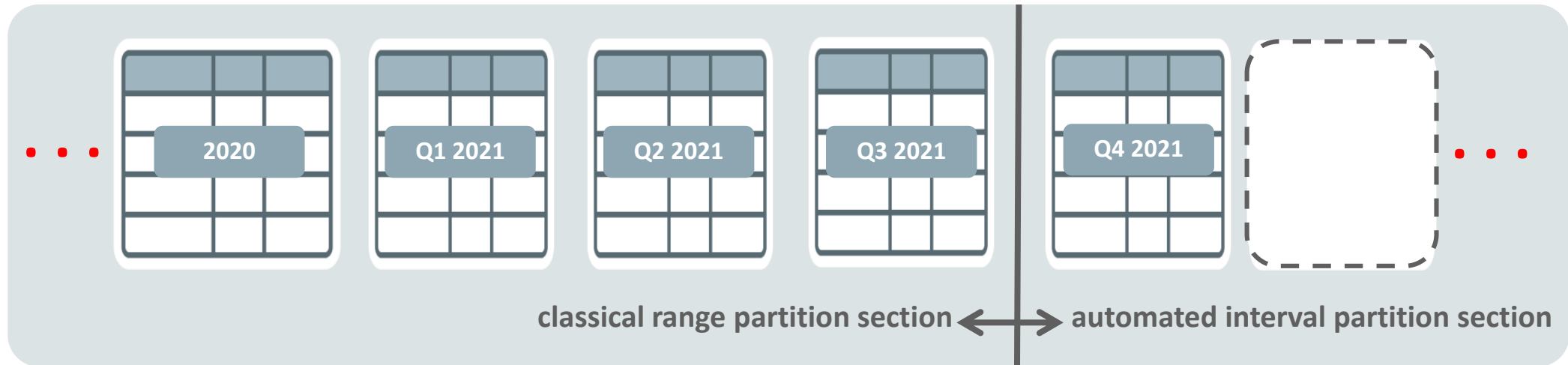
Interval Partitioning



- Range partitioned tables can be extended into interval partitioned tables
 - Simple metadata command
 - Investment protection

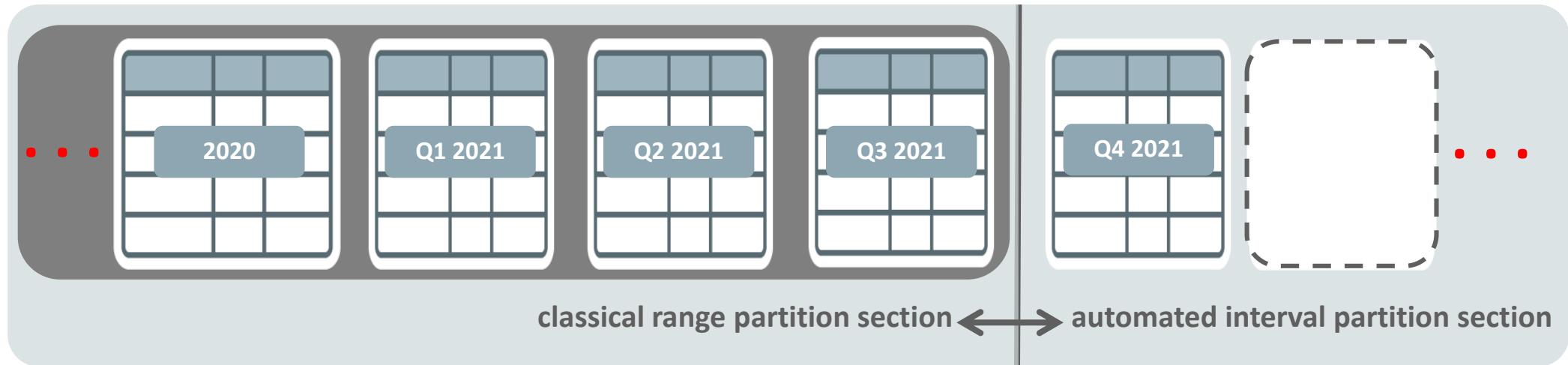
```
ALTER TABLE EVENTS  
SET INTERVAL (NUMTOYMINTERVAL (1, 'month'));
```


Interval Partitioning



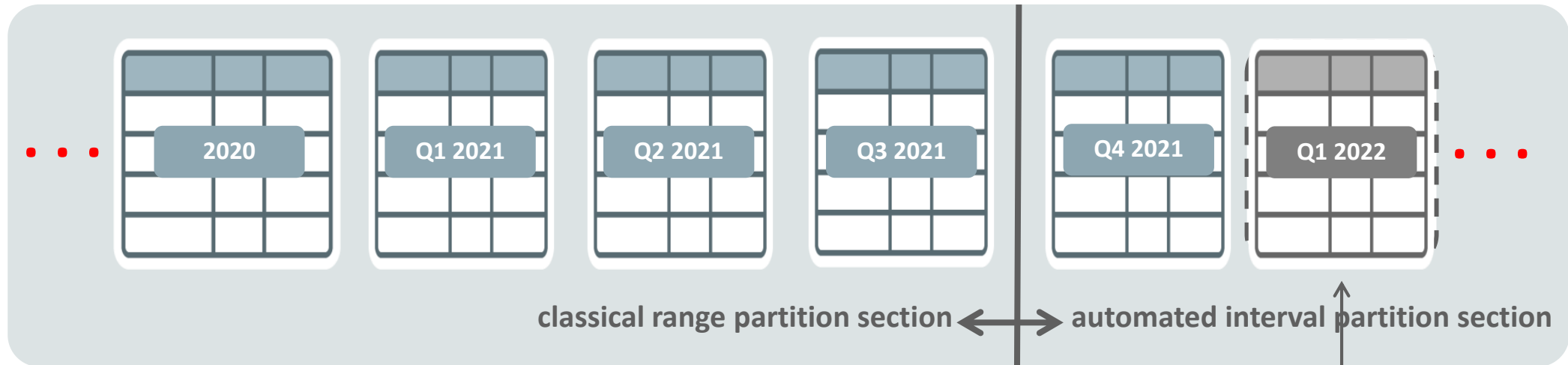
- Interval partitioned table has classical range and automated interval section
 - Automated new partition management plus full partition maintenance capabilities: ***“Best of both worlds”***

Interval Partitioning



1. Merge and move old partitions for ILM

Interval Partitioning



1. Merge and move old partitions for ILM
2. Insert new data
 - Automatic partition instantiation

Values ('13-JAN-2022')

Deferred Segment Creation vs Interval Partitioning

Interval Partitioning

- Maximum number of one million partitions are pre-defined
 - Explicitly defined plus interval-based partitions
- No segments are allocated for partitions without data
 - New record insertion triggers segment creation
- Ideal for “ever-growing” tables

“Standard” Partitioning with deferred segment creation

- Only explicitly defined partitions are existent
 - New partitions added via DDL
- No segments are allocated for partitions without data
 - New record insertion triggers segment creation when data matches pre-defined partitions
- Ideal for sparsely populated pre-defined tables



Difference Between Range and Interval

Interval Partitioning

- Full automation for equi-sized range partitions
- Partitions are created as metadata information only
 - Start Partition is made persistent
- Segments are allocated as soon as new data arrives
 - No need to create new partitions
 - Local indexes are created and maintained as well
- Interval Partitioning is almost a transparent extension to range partitioning
 - .. But interval implementation introduces some subtle differences



Interval versus Range Partitioning

Partition bounds

- Interval partitions have lower and upper bound
 - No infinite upper bound (MAXVALUES)
- Range partitions only have upper bounds
 - Lower bound derived by previous partition
 - Upper bound infinite (MAXVALUES)

Partition naming

- Interval partitions cannot be named in advance
 - Use the PARTITION FOR (<value>) clause
- Range partitions must be named



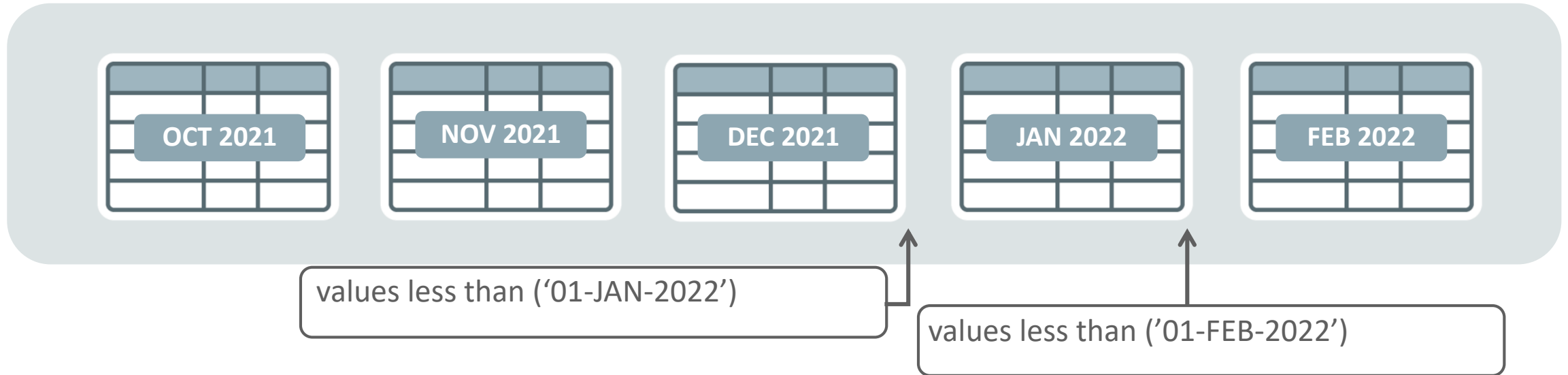
Interval versus Range Partitioning, cont.

- Partition merge
 - Multiple non-existent interval partitions are silently merged
 - Only two adjacent range partitions can be merged at any point in time
- Number of partitions
 - Interval partitioned tables have always one million partitions
 - Non-existent partitions “exist” through INTERVAL clause
 - No MAXVALUE clause for interval partitioning
 - Maximum value defined through number of partitions and INTERVAL clause
 - Range partitioning can have up to one million partitions
 - MAXVALUE clause defines most upper partition



Interval Versus Range Partitioning

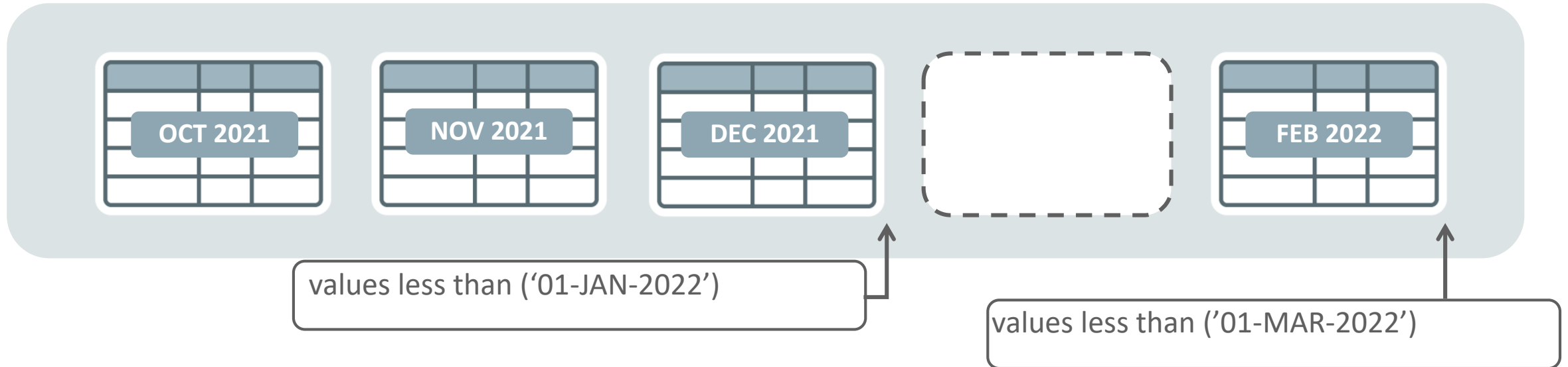
Partition Bounds for **Range Partitioning**



- Partitions only have upper bounds
 - Lower bound derived through upper bound of previous partition

Interval Versus Range Partitioning

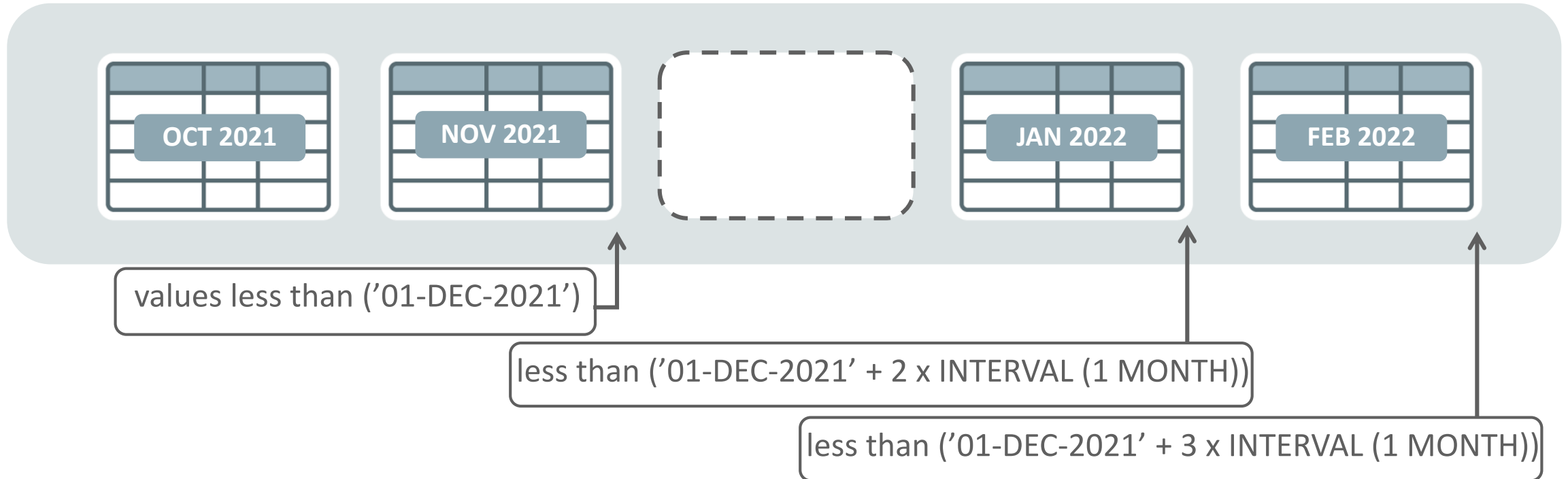
Partition Bounds for **Range Partitioning**



- Drop of previous partition moves lower boundary
 - “Feb 2022” now spawns 01-JAN-2022 to 28-FEB-2022

Interval Versus Range Partitioning

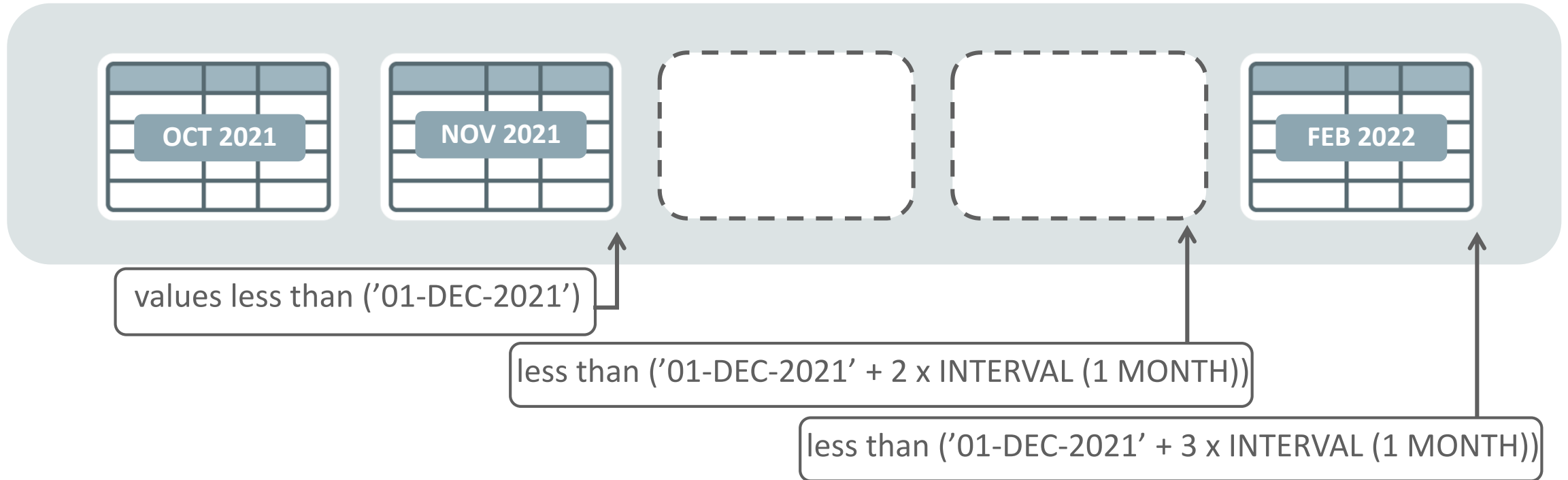
Partition Bounds for **Interval Partitioning**



- Partitions have upper and lower bounds
 - Derived by INTERVAL function and last range partition

Interval Versus Range Partitioning

Partition Bounds for **Interval Partitioning**



- Drop does not impact partition boundaries
 - “Feb 2022” still spawns 01-FEB-2022 to 28-FEB-2022

Interval versus Range Partitioning

Partition Naming

- Range partitions **can** be named
 - System generated name if not specified

```
SQL> alter table t add partition values less than(20);  
Table altered.  
SQL> alter table t add partition P30 values less than(30);  
Table altered.
```

- Interval partitions **cannot** be named
 - Always system generated name

```
SQL> alter table t add partition values less than(20);  
*  
ERROR at line 1: ORA-14760: ADD PARTITION is not permitted  
on Interval partitioned objects
```

- Use new deterministic PARTITION FOR () extension

```
SQL> alter table t1 rename partition for (9) to p_10;  
Table altered.
```



Interval Versus Range Partitioning

Partition Merge – **Range Partitioning**

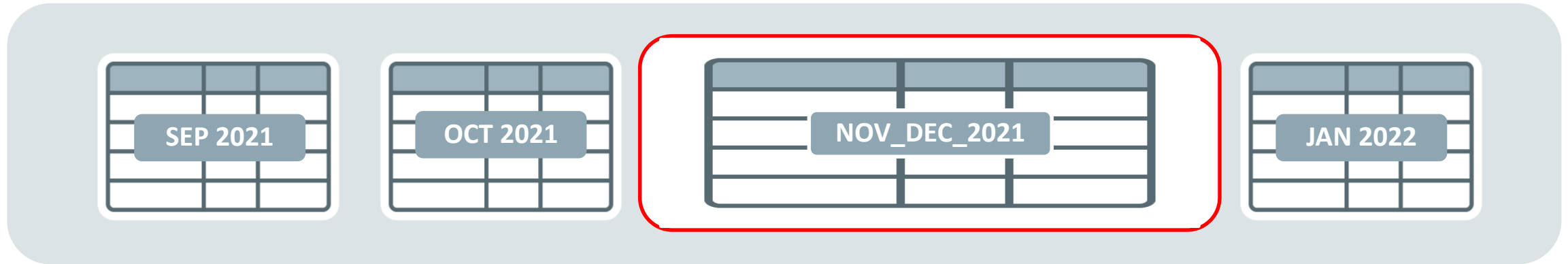


```
MERGE PARTITIONS NOV_2021, DEC_2021 INTO PARTITION NOV_DEC_2021
```

- Merge two adjacent partitions for range partitioning
 - Upper bound of higher partition is new upper bound
 - Lower bound derived through upper bound of previous partition

Interval Versus Range Partitioning

Partition Merge – **Range Partitioning**



```
MERGE PARTITIONS NOV_2021, DEC_2021 INTO PARTITION NOV_DEC_2021
```

- New segment for merged partition is created
 - Rest of the table is unaffected

Interval Versus Range Partitioning

Partition Merge – **Interval Partitioning**

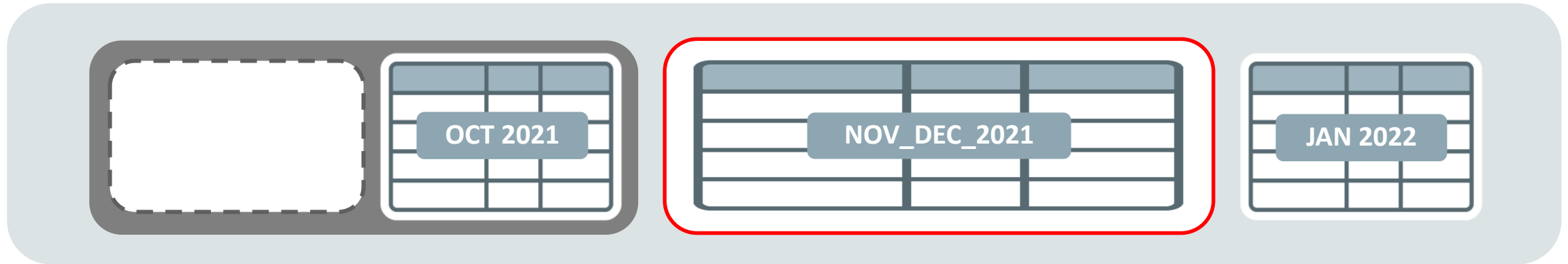


```
MERGE PARTITIONS NOV_2021, DEC_2021 INTO PARTITION NOV_DEC_2021
```

- Merge two adjacent partitions for interval partitioning
 - Upper bound of higher partition is new upper bound
 - Lower bound derived through lower bound of first partition

Interval Versus Range Partitioning

Partition Merge – **Interval Partitioning**



```
MERGE PARTITIONS NOV_2021, DEC_2021 INTO PARTITION NOV_DEC_2021
```

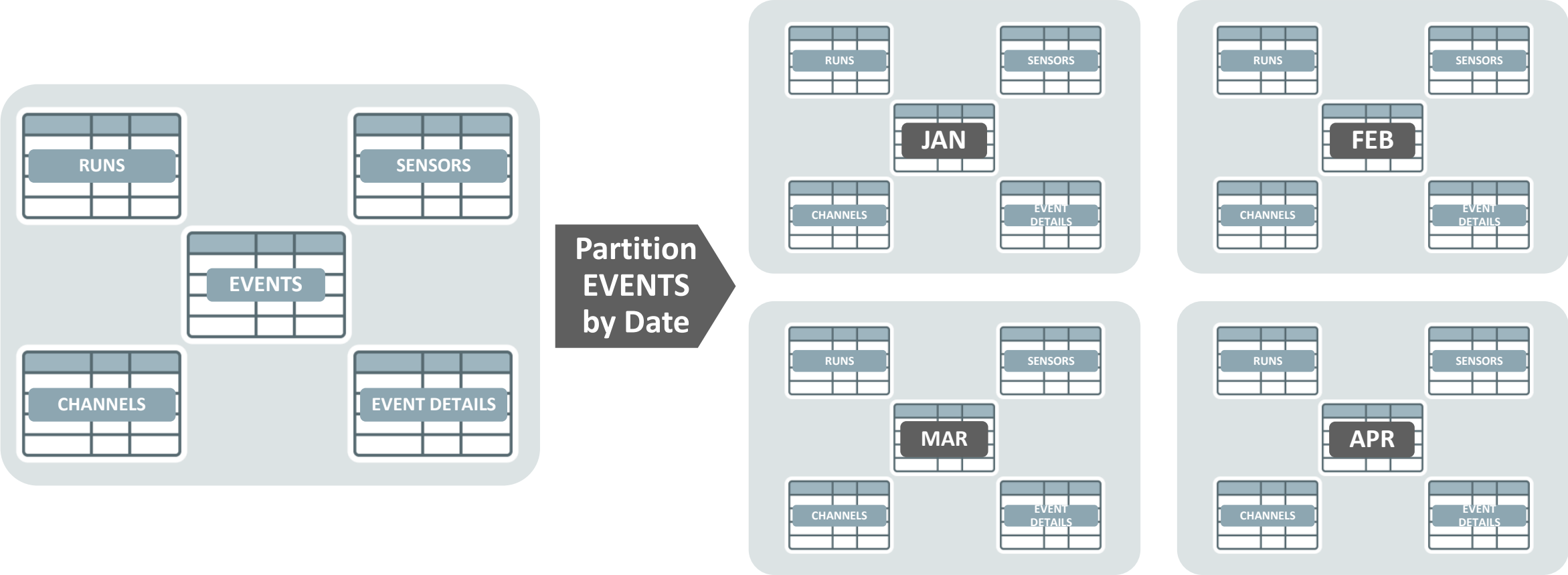
- New segment for merged partition is created
 - Holes before highest non-interval partition will be silently “merged” as well
 - Interval only valid beyond the highest non-interval partition

Reference Partitioning

Introduced in Oracle 11g Release 1 (11.1)

Reference Partitioning

Inherit partitioning strategy



Reference Partitioning

Business Problem

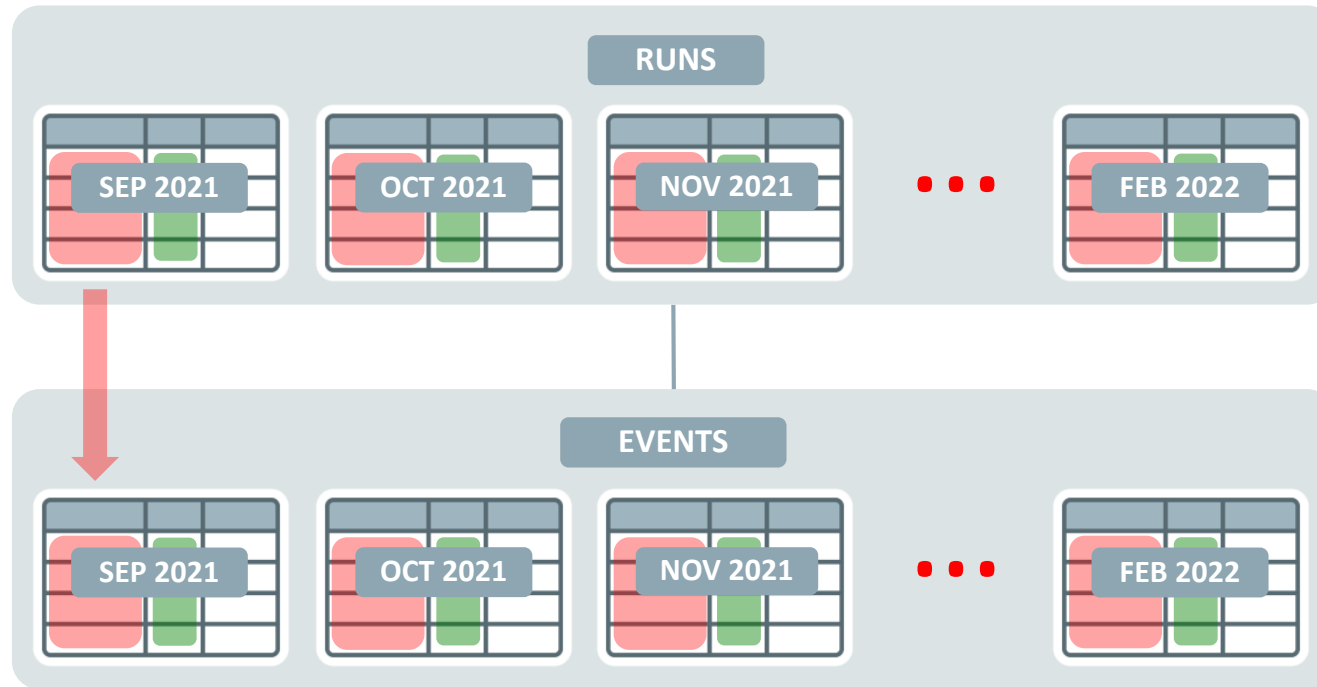
- Related tables benefit from same partitioning strategy
 - Sample 3NF order entry data model
- Redundant storage of same information solves problem
 - Data and maintenance overhead

Solution

- Oracle Database 11g introduces Reference Partitioning
 - Child table inherits the partitioning strategy of parent table through PK-FK
 - Intuitive modelling
- Enhanced Performance and Manageability



Without Reference Partitioning

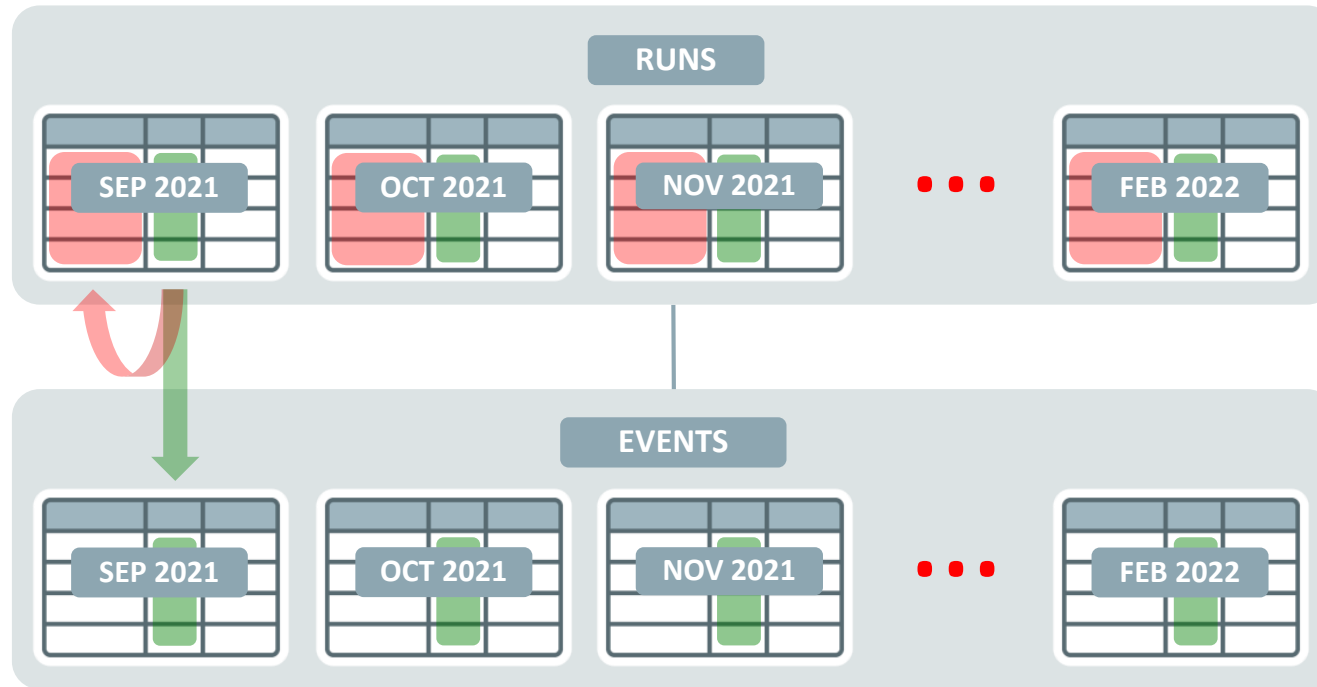


```
RANGE (run_date)  
Primary key run_id
```

- Redundant storage
- Redundant maintenance

```
RANGE (run_date)  
Foreign key run_id
```

With Reference Partitioning



```
RANGE (run_date)  
Primary key run_id
```

- Partitioning key inherited through PK-FK relationship

```
RANGE (run_date)  
Foreign key run_id
```

Reference Partitioning

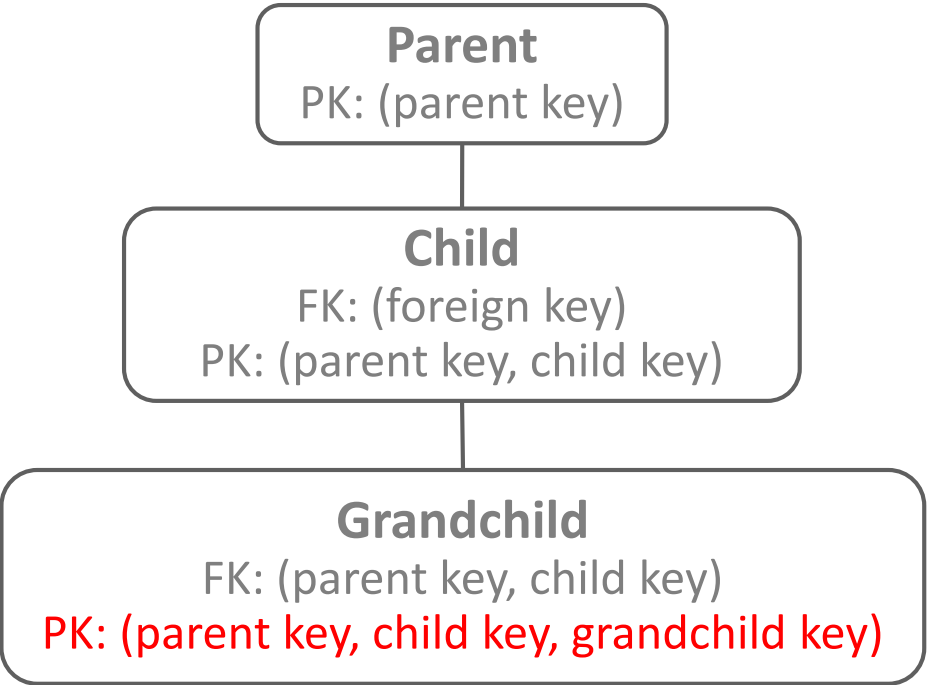
Use Cases

- Traditional relational model
 - Primary key inherits down to all levels of children and becomes part of an (elongated) primary key definition
- Object oriented-like model
 - Several levels of primary-foreign key relationship
 - Primary key on each level is primary key + “object ID”
- Reference Partitioning optimally suited to address both modeling techniques

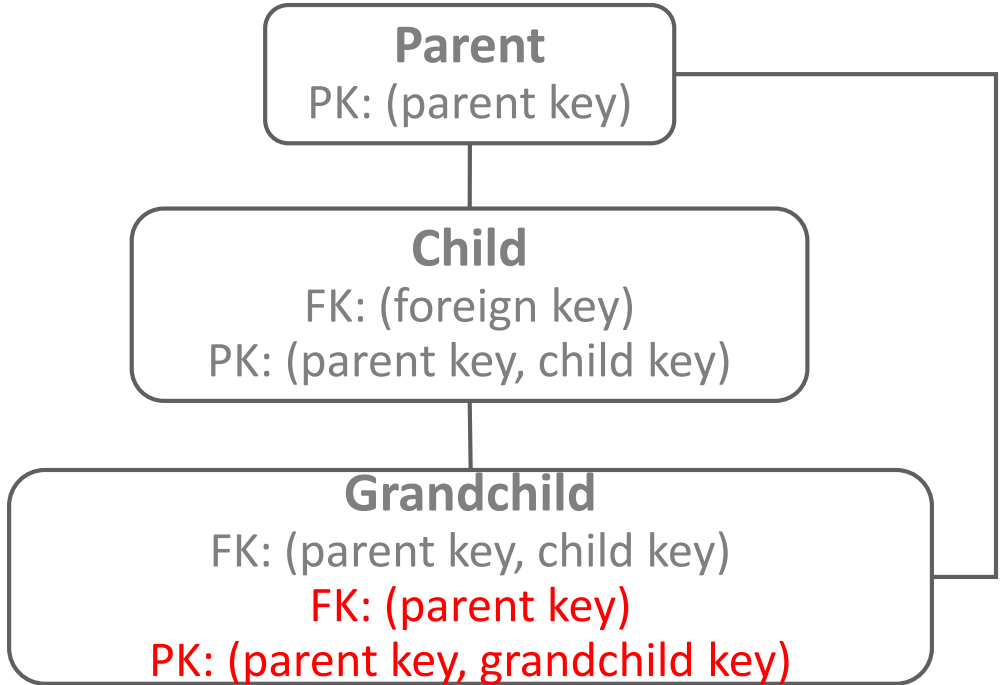


Reference Partitioning

Relational Model



“Object-like” model



Reference Partitioning

Example

```
create table project (project_id number not null,  
                    project_number varchar2(30),  
                    project_name varchar2(30), ...  
                    constraint proj_pk primary key (project_id))  
partition by list (project_id)  
(partition p1 values (1),  
 partition p2 values (2),  
 partition pd values (DEFAULT));
```

```
create table project_customer (project_cust_id number not null,  
                              project_id number not null,  
                              cust_name varchar2(30),  
                              constraint pk_proj_cust primary key  
                                (project_id, project_cust_id),  
                              constraint proj_cust_proj_fk foreign key  
                                (project_id) references project(project_id))  
partition by reference (proj_cust_proj_fk);
```



Reference Partitioning

Example, cont.

```
create table proj_cust_address (project_cust_addr_id number not null,  
                               project_cust_id number not null,  
                               project_id number not null,  
                               cust_address varchar2(30),  
                               constraint pk_proj_cust_addr primary key  
                                 (project_id, project_cust_addr_id),  
                               constraint proj_c_addr_proj_cust_fk foreign key  
                                 (project_id, project_cust_id)  
                                 references project_customer  
                                 (project id, project cust id))  
partition by reference (proj_c_addr_proj_cust_fk);
```



Reference Partitioning

Some metadata

Table information

```
SQL> SELECT table_name, partitioning_type, ref_ptn_constraint_name
       FROM   user_part_tables
       WHERE  table_name IN ('PROJECT','PROJECT_CUSTOMER','PROJ_CUST_ADDRESS');
```

TABLE_NAME	PARTITION	REF_PTN_CONSTRAINT_NAME
PROJECT	LIST	
PROJECT_CUSTOMER	REFERENCE	PROJ_CUST_PROJ_FK
PROJ_CUST_ADDRESS	REFERENCE	PROJ_C_ADDR_PROJ_FK

Partition information

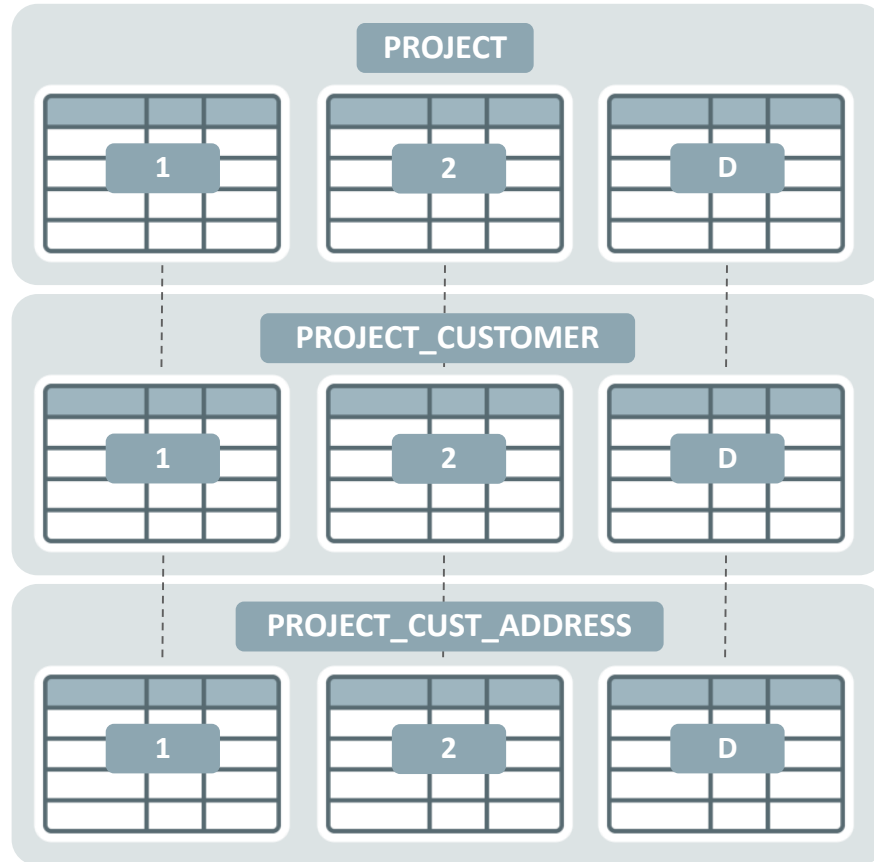
```
SQL> SELECT table_name, partition_name, high_value
       FROM   user_tab_partitions
       WHERE  table_name in ('PROJECT','PROJECT_CUSTOMER')
       ORDER BY table_name, partition_position;
```

TABLE_NAME	PARTITION_NAME	HIGH_VALUE
PROJECT	P1	1
PROJECT	P2	2
PROJECT	PD	DEFAULT
PROJECT_CUSTOMER	P1	
PROJECT_CUSTOMER	P2	
PROJECT_CUSTOMER	PD	



Reference Partitioning

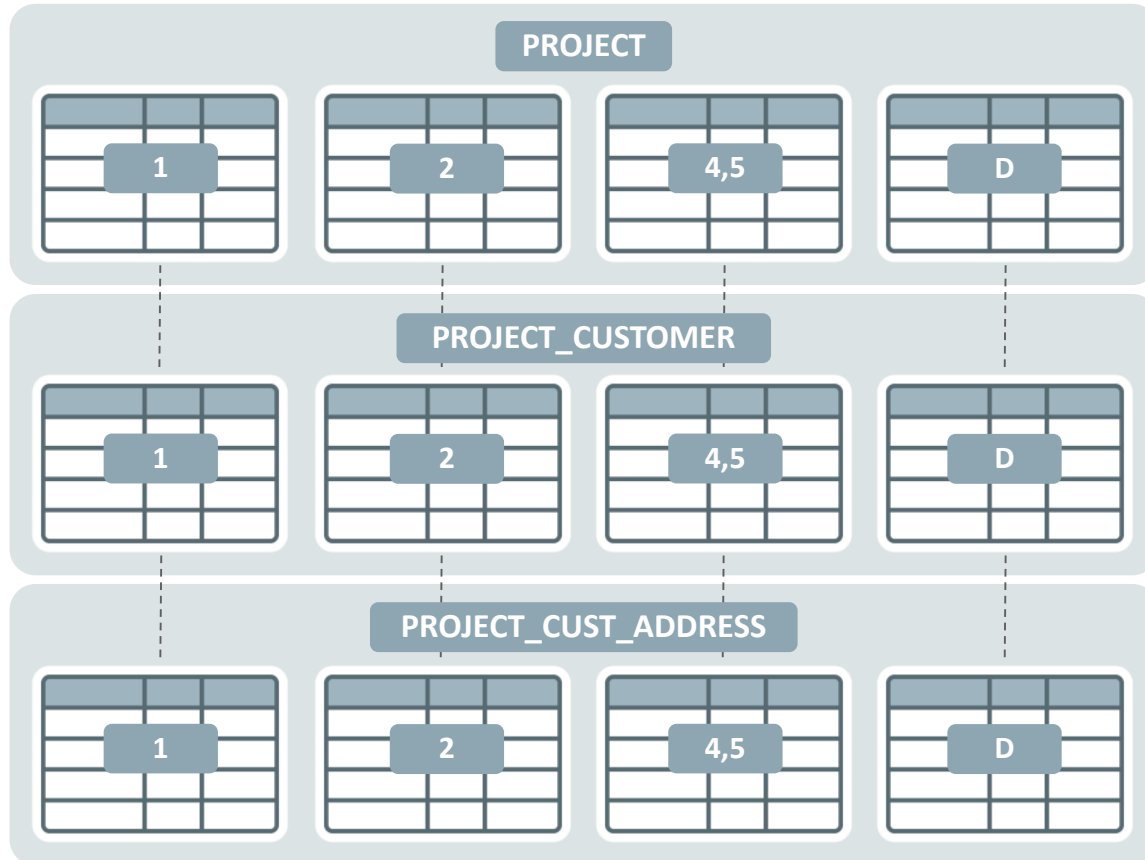
Partition Maintenance



```
ALTER TABLE project  
SPLIT PARTITION pd VALUES (4,5) INTO  
(PARTITION pd, PARTITION p45);
```

Reference Partitioning

Partition Maintenance

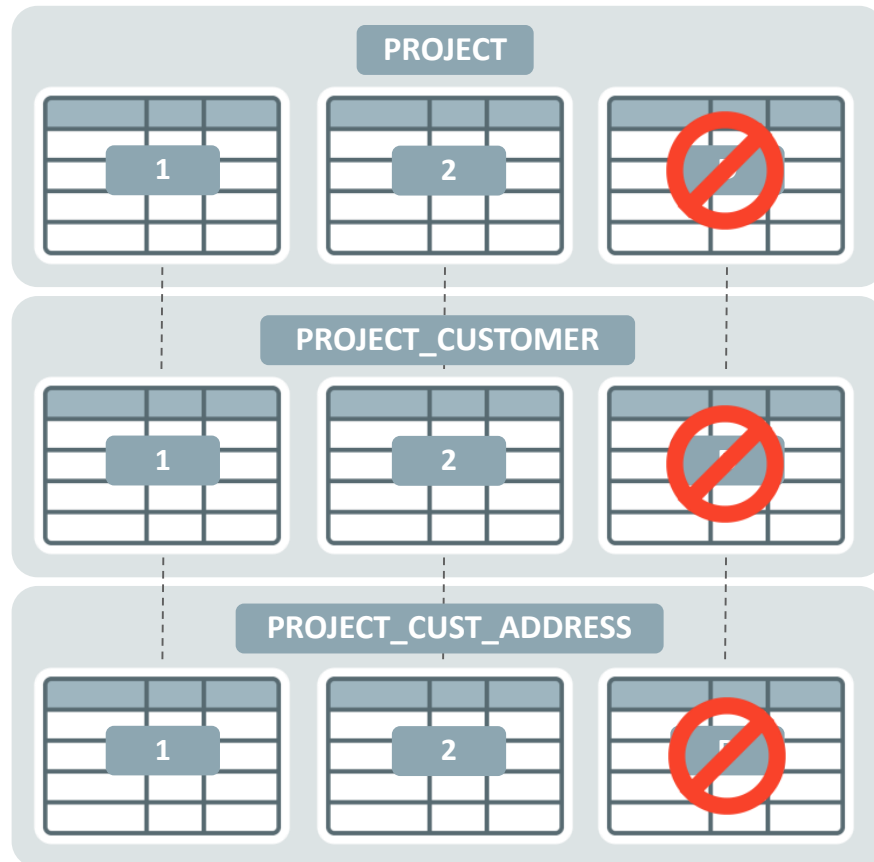


```
ALTER TABLE project  
SPLIT PARTITION pd VALUES (4,5) INTO  
(PARTITION pd, PARTITION p45);
```

- PROJECT partition PD will be split
 - “Default” and (4,5)
- PROJECT_CUSTOMER will split its dependent partition
 - Co-location with equivalent parent record of PROJECT
 - Parent record in (4,5) means child record in (4.5)
- PROJECT_CUST_ADDRESS will split its dependent partition
 - Co-location with equivalent parent record of PROJECT_CUSTOMER
- One-level lookup required for both placements

Reference Partitioning

Partition Maintenance



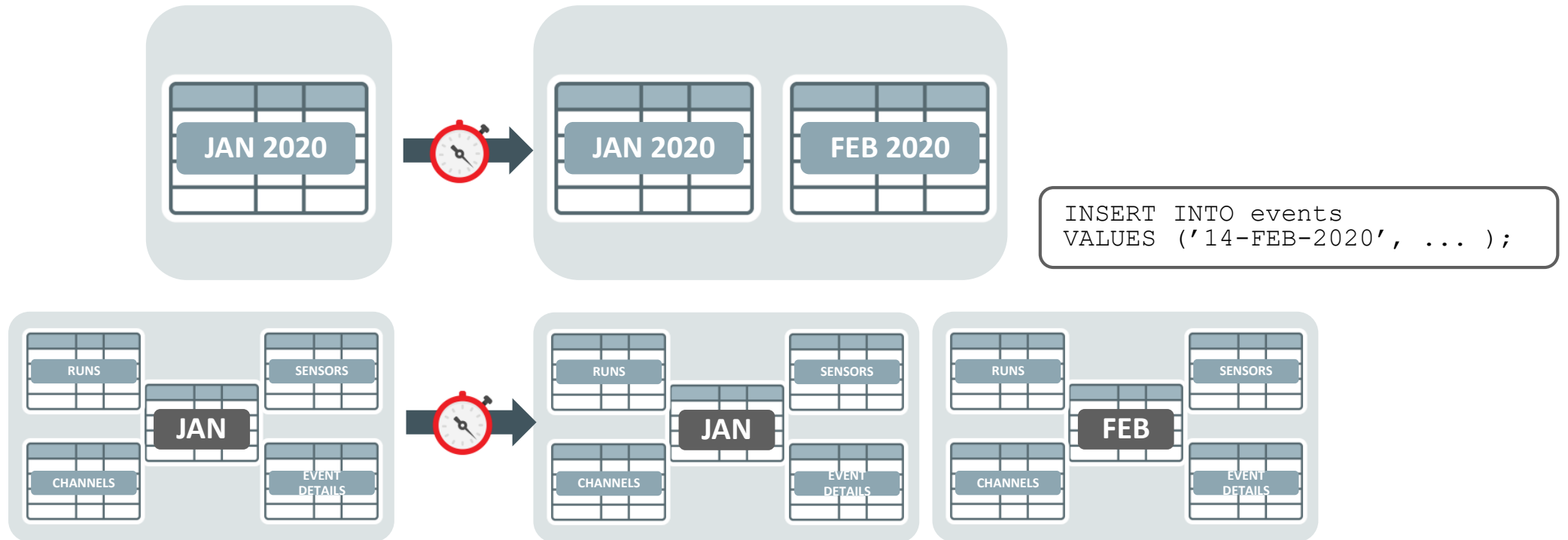
```
ALTER TABLE project_cust_address  
DROP PARTITION pd;
```

- PROJECT partition PD will be dropped
 - PK-FK is guaranteed not to be violated
- PROJECT_CUSTOMER will drop its dependent partition
- PROJECT_CUST_ADDRESS will drop its dependent partition
- Unlike “normal” partitioned tables, PK-FK relationship stays enabled
 - You cannot arbitrarily drop or truncate a partition with the PK of a PK-FK relationship
- Same is true for TRUNCATE
 - Bottom-up operation

Interval Reference Partitioning

Introduced in Oracle 12c Release 1 (12.1)

Interval-Reference Partitioning



- New partitions are automatically created when new data arrives
- All child tables will be automatically maintained
- Combination of two successful partitioning strategies for better business modeling

Interval-Reference Partitioning

```
SQL> REM create some interval-referenced tables ..
SQL> create table intRef_p (pkcol number not null, col2 varchar2(200),
 2                          constraint pk_intref primary key (pkcol))
 3 partition by range (pkcol) interval (10)
 4 (partition p1 values less than (10));
```

Table created.

```
SQL>
SQL> create table intRef_c1 (pkcol number not null, col2 varchar2(200), fkcol number not null,
 2                          constraint pk_c1 primary key (pkcol),
 3                          constraint fk_c1 foreign key (fkcol) references intRef_p(pkcol) ON DELETE CASCADE)
 4 partition by reference (fk_c1);
```

Table created.

```
SQL>
SQL> create table intRef_c2 (pkcol number primary key not null, col2 varchar2(200), fkcol number not null,
 2                          constraint fk_c2 foreign key (fkcol) references intRef_p(pkcol) ON DELETE CASCADE)
 3 partition by reference (fk_c2);
```

Table created.



Interval-Reference Partitioning

- New partitions only created when data arrives
 - No automatic partition instantiation for complete reference tree
 - Optimized for sparsely populated reference partitioned tables
- Partition names inherited from already existent partitions
 - Name inheritance from direct relative
 - Parent partition p100 will result in child partition p100
 - Parent partition p100 and child partition c100 will result in grandchild partition c100



Virtual Column Based Partitioning

Introduced in Oracle 11g Release 1 (11.1)

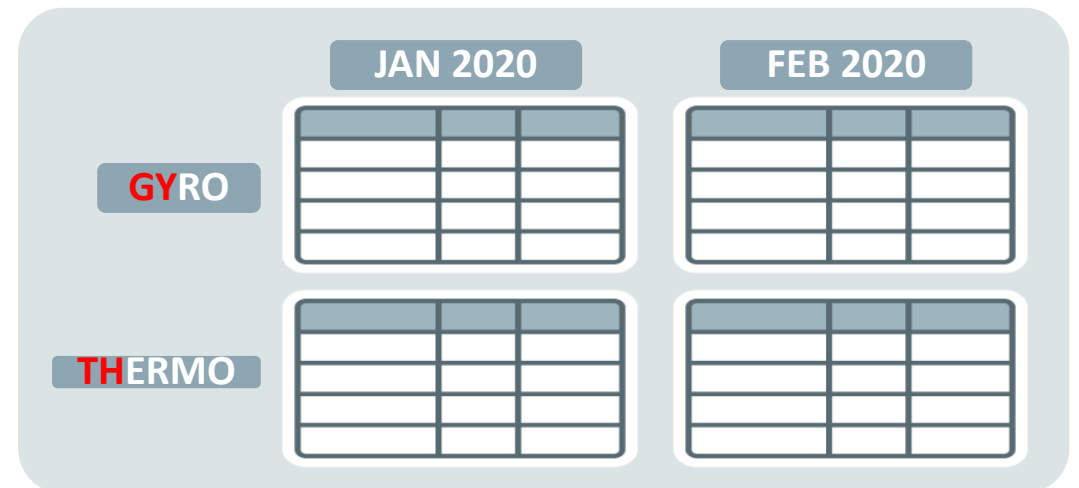


Virtual Column Based Partitioning

EVENTS

EVENT_ID	EVENT_DATE	SENSOR_ID	...	REGION AS	(SUBSTR(EVENT_ID, 6, 2))
9834	GY-14	12-JAN-2020		65920	GY
8300	TH-97	14-FEB-2020		39654	TH
3886	TH-02	16-JAN-2020		4529	GY
2566	GY-94	19-JAN-2020		15327	TH
3699	GY-63	02-FEB-2020		18733	TH

- REGION requires no storage
- Partition by ORDER_DATE, REGION



Virtual Columns

Example

- Base table with all attributes ...

```
CREATE TABLE accounts  
(acc_no      number(10)    not null,  
 acc_name    varchar2(50) not null, ...
```

12500	Adams	
12507	Blake	
12666	King	
12875	Smith	

Virtual Columns

Example

- Base table with all attributes ...
... is extended with the virtual (derived) column

```
CREATE TABLE accounts
(acc_no      number(10)    not null,
 acc_name    varchar2(50) not null, ...
 acc_branch  number(2)     generated always as
              (to_number(substr(to_char(acc_no),1,2)))
```

12500	Adams	12
12507	Blake	12
12666	King	12
12875	Smith	12

Virtual Columns

Example

- Base table with all attributes ...

... is extended with the virtual (derived) column
... and the virtual column is used as partitioning key

```
CREATE TABLE accounts
(acc_no      number(10)    not null,
 acc_name    varchar2(50) not null, ...
 acc_branch  number(2)     generated always as
              (to_number(substr(to_char(acc_no),1,2)))
 partition by list (acc_branch) ...
```

12500	Adams	12
12507	Blake	12
12666	King	12
12875	Smith	12

...

32320	Jones	32
32407	Clark	32
32758	Hurd	32
32980	Kelly	32

Virtual Columns

Partition Pruning

- Conceptual model considers virtual columns as **visible** and **used** attributes
- Partition pruning currently only works with predicates on the virtual column (partition key) itself
 - No transitive predicates
 - Enhancement planned for the future



Partitioning for performance

Composite Partitioning

Range-Hash introduced in Oracle 8i

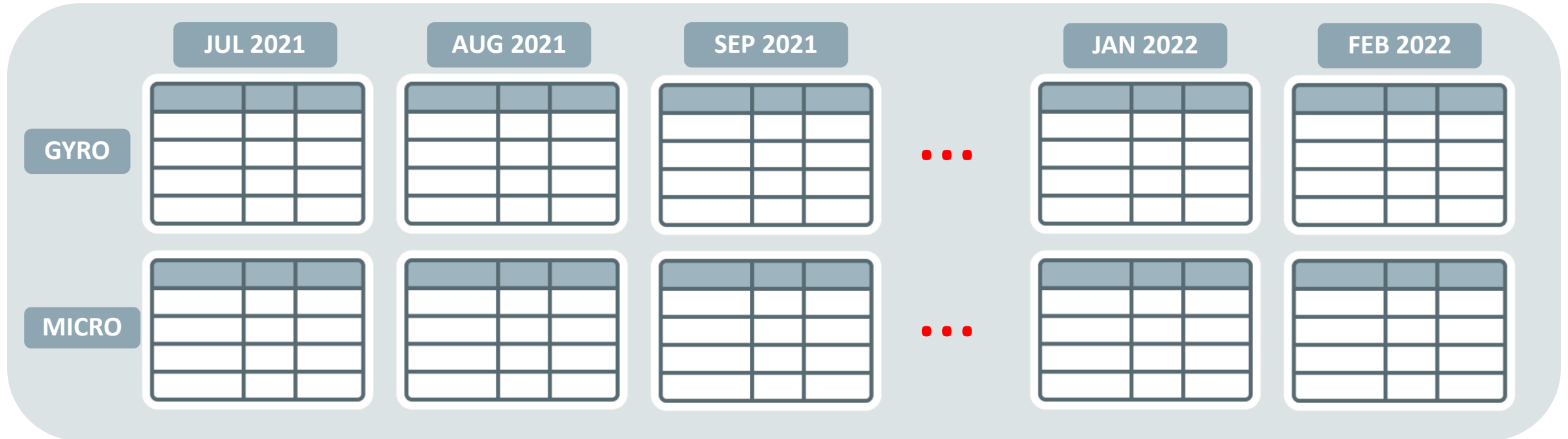
Range-List introduced in Oracle 9i Release 2

[Interval | Range | List | Hash]-[Range | List | Hash] introduced in Oracle 11g
Release 1|2

*Hash-Hash in 11.2



Composite Partitioning



- Data is organized along two dimensions
 - Record placement is deterministically identified by dimensions
 - Example RANGE-LIST

Composite Partitioning

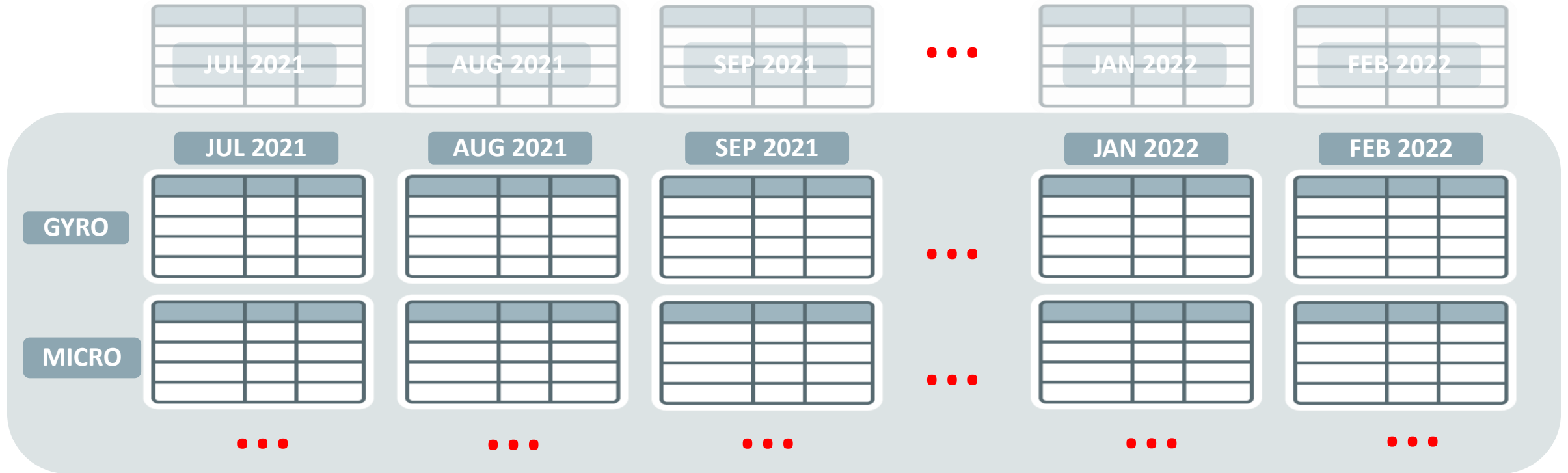
Concept



```
CREATE TABLE EVENTS ..PARTITION BY RANGE (time_id)
```

Composite Partitioning

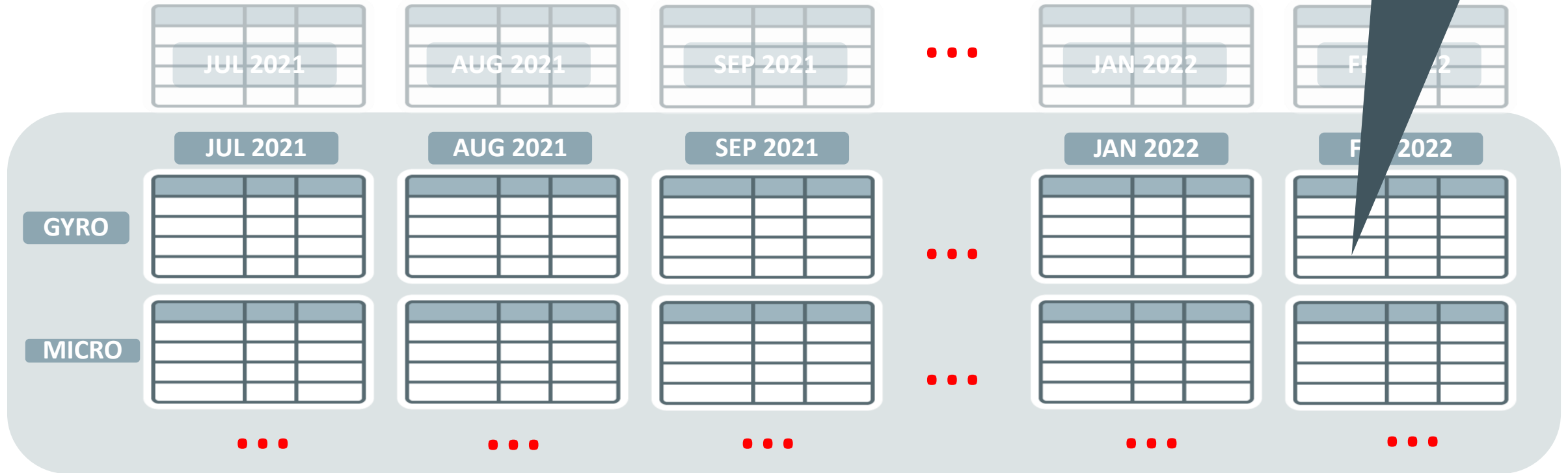
Concept



```
CREATE TABLE EVENTS ..PARTITION BY RANGE (time_id)
                    SUPARTITION BY LIST (sensor_type)
```

Composite Partitioning

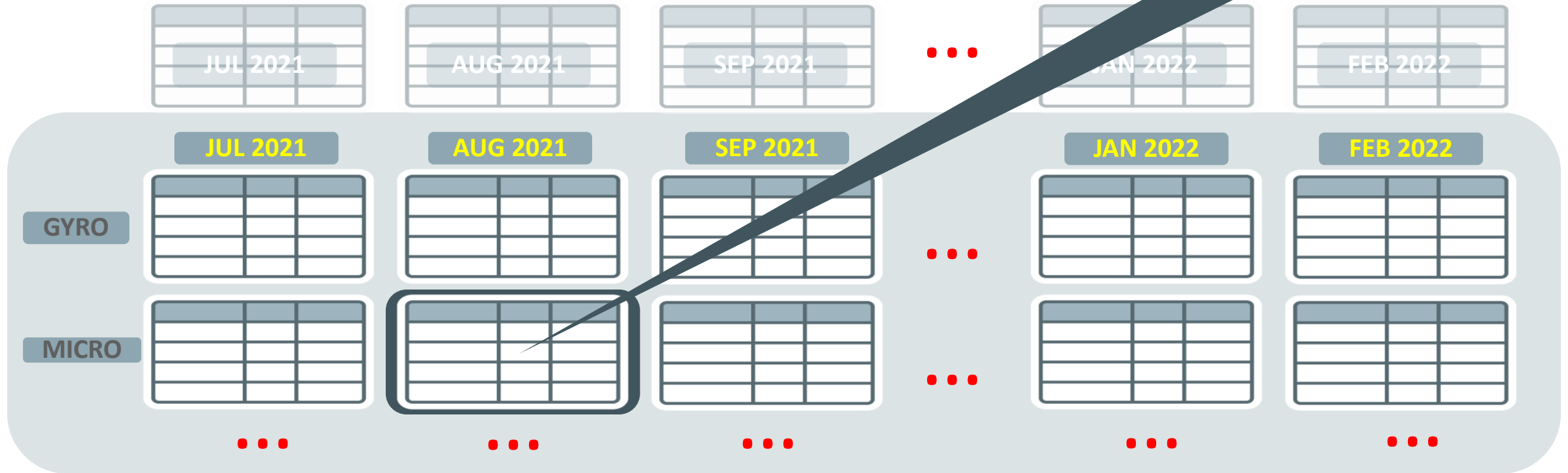
Concept



```
CREATE TABLE EVENTS ..PARTITION BY RANGE (time_id)
                    SUPARTITION BY LIST (sensor_type)
```

Composite Partitioning

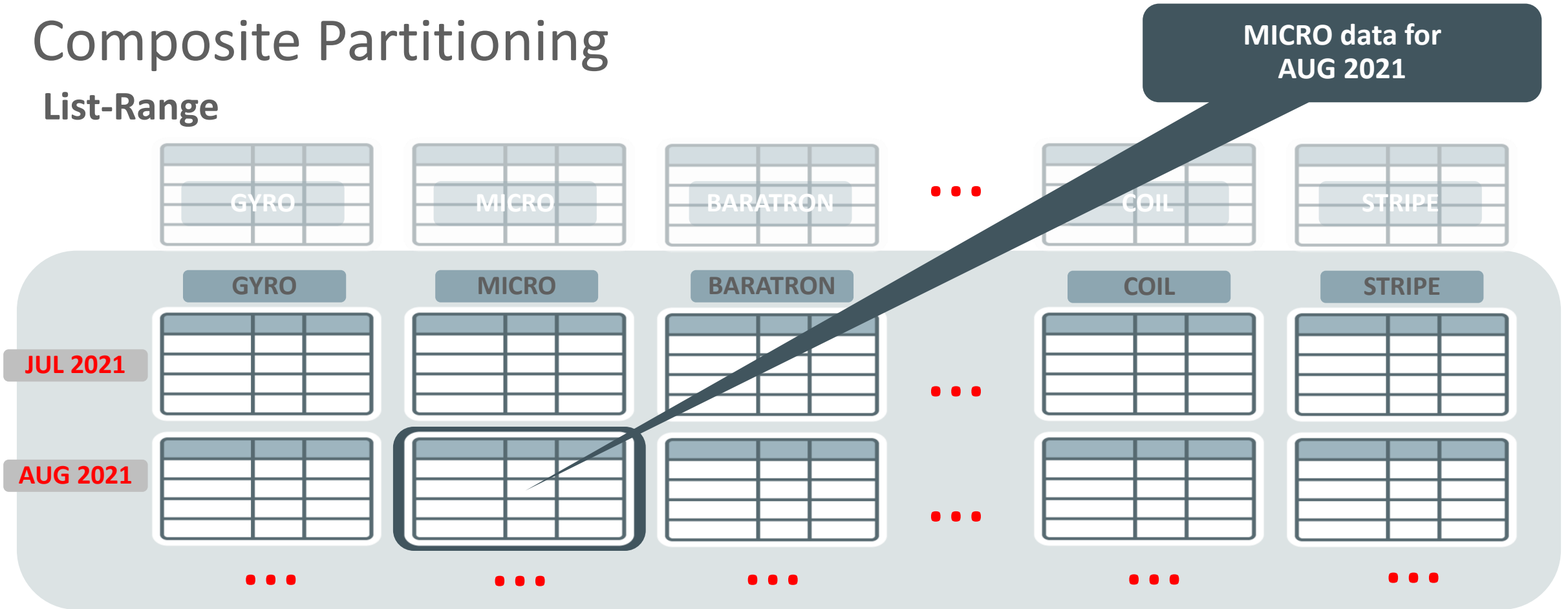
Range-List



```
CREATE TABLE EVENTS ..PARTITION BY RANGE (time_id)  
SUPARTITION BY LIST (sensor_type)
```

Composite Partitioning

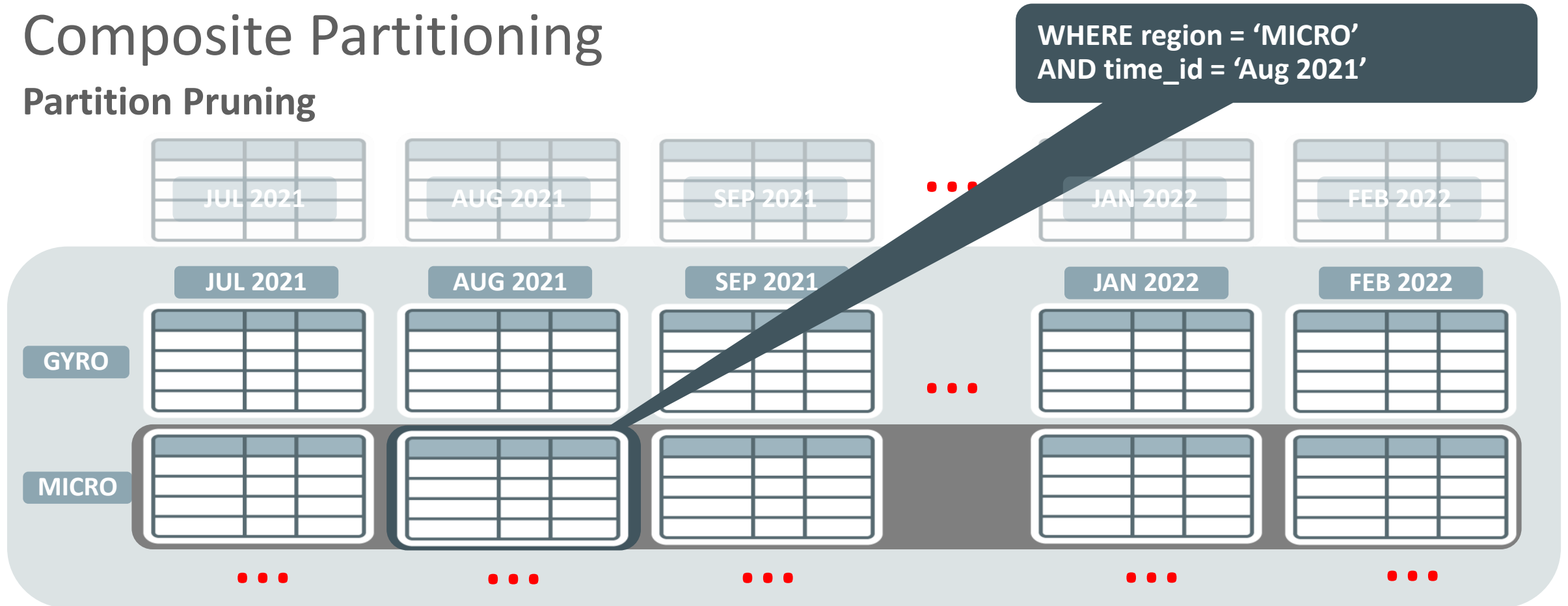
List-Range



```
CREATE TABLE EVENTS ..PARTITION BY LIST (sensor_type)
SUPARTITION BY RANGE (time_id)
```


Composite Partitioning

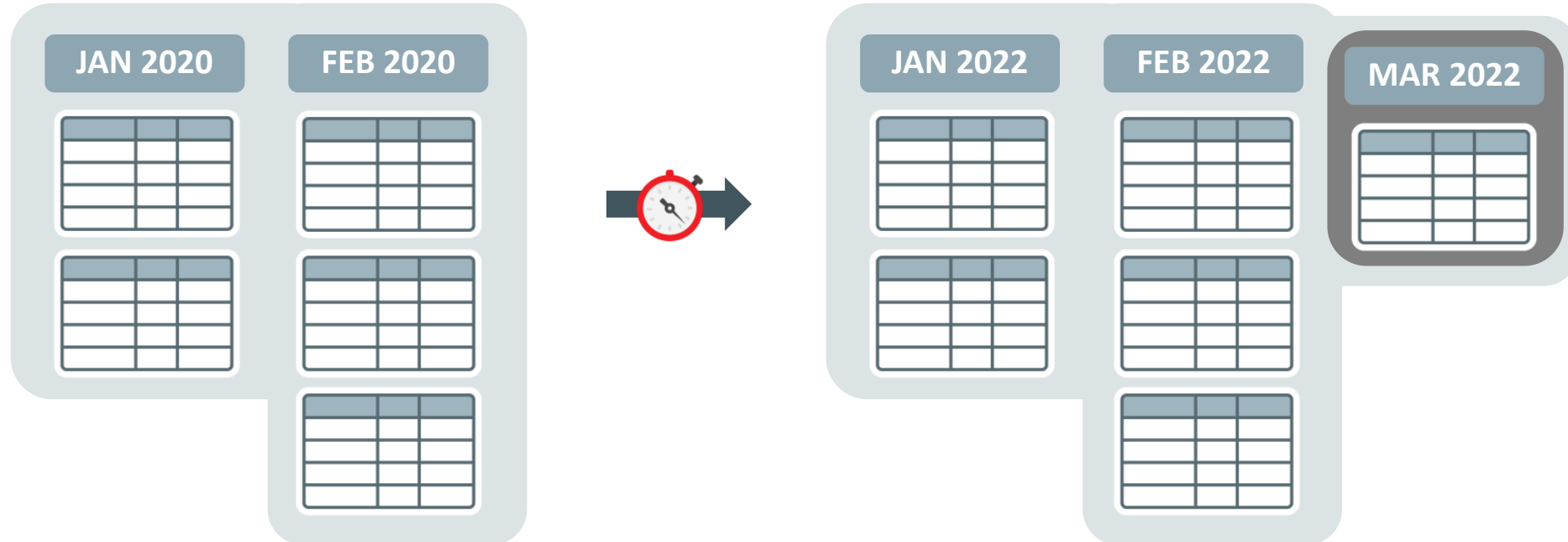
Partition Pruning



- Partition pruning is independent of composite order
 - Pruning along one or both dimensions
 - Same pruning for RANGE-LIST and LIST_RANGE

Composite Interval Partitioning

Add Partition



- Without subpartition template, only **one** subpartition will be created
 - Range: MAXVALUE
 - List: DEFAULT
 - Hash: one hash bucket

Composite Interval Partitioning

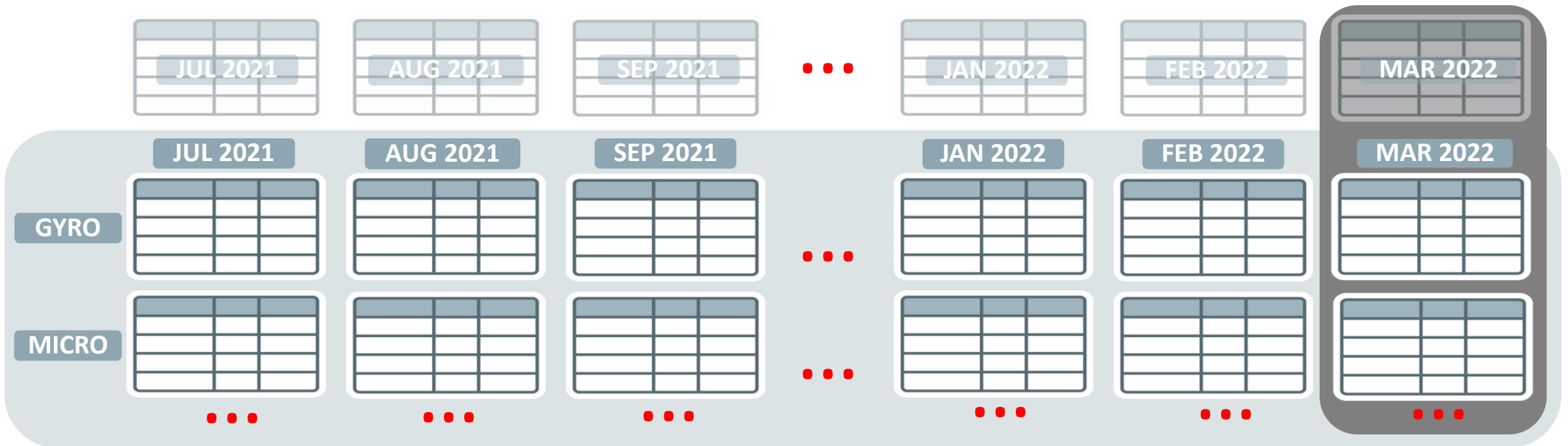
Subpartition template

- Subpartition template defines shape of **future** subpartitions
 - Can be added and/or modified at any point in time
 - No impact on existing [sub]partitions
- Controls physical attributes for subpartitions as well
 - Just like the default settings for a partitioned table does for partitions
- Difference Interval and Range Partitioning
 - Naming template only for Range
 - System-generated names for Interval



Composite Partitioning

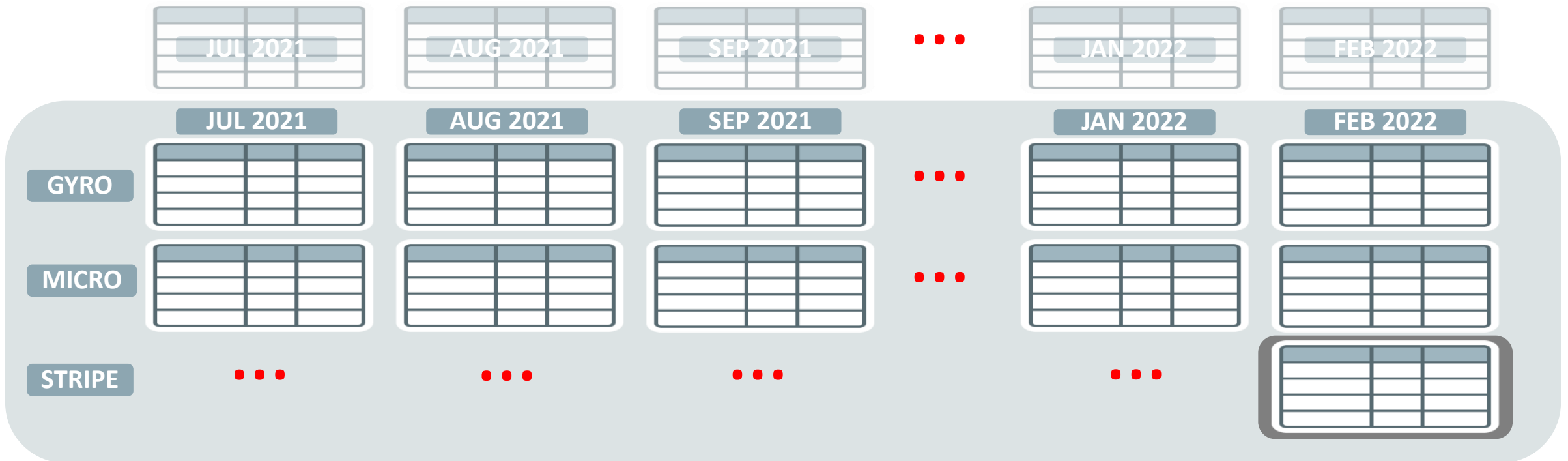
Add Partition



- ADD PARTITION always on top-level dimension
 - Identical for all newly added subpartitions
 - RANGE-LIST: new time_id range
 - LIST-RANGE: new list of region values

Composite Partitioning

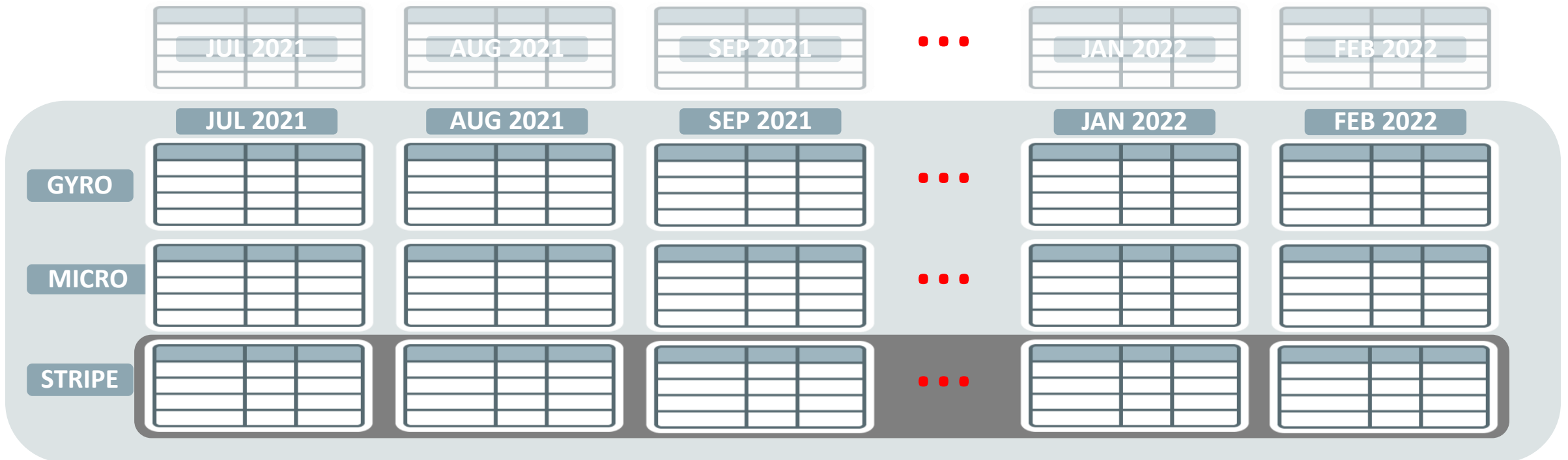
Add Subpartition



- **ADD SUBPARTITION** only for one partition
 - Asymmetric, only possible on subpartition level
 - Impact on partition-wise joins

Composite Partitioning

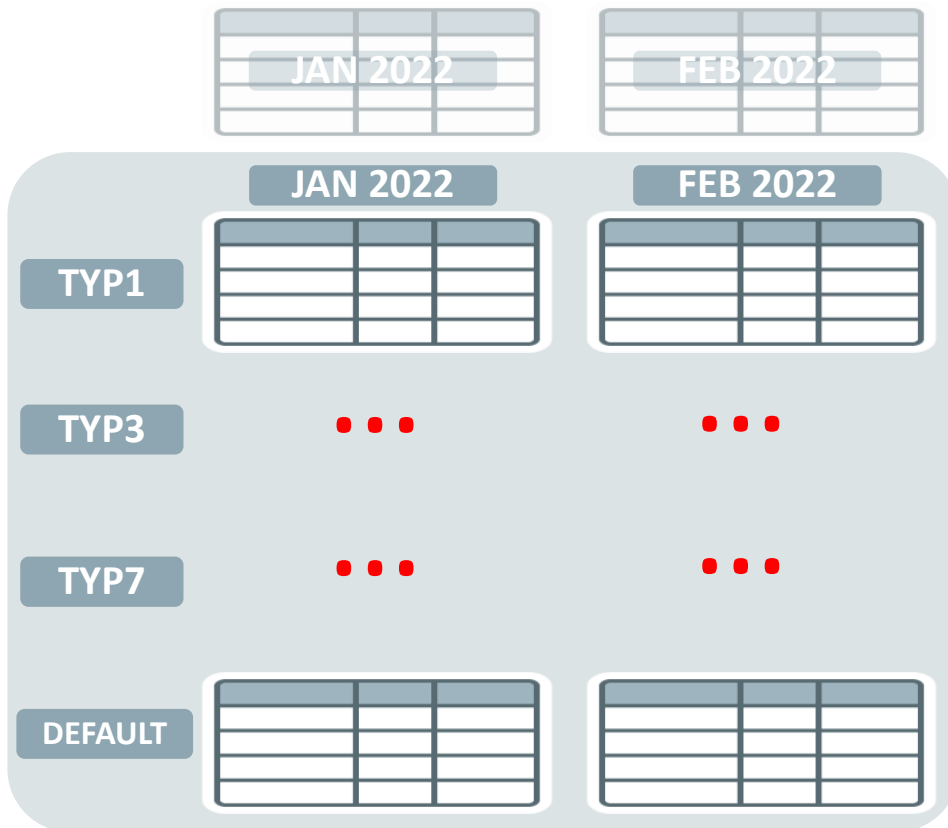
Add Subpartition



- ADD SUBPARTITION for all partitions
 - N operations necessary (for each existing partition)
 - Adjust subpartition template for future partitions

Composite Partitioning

Asymmetric subpartitions

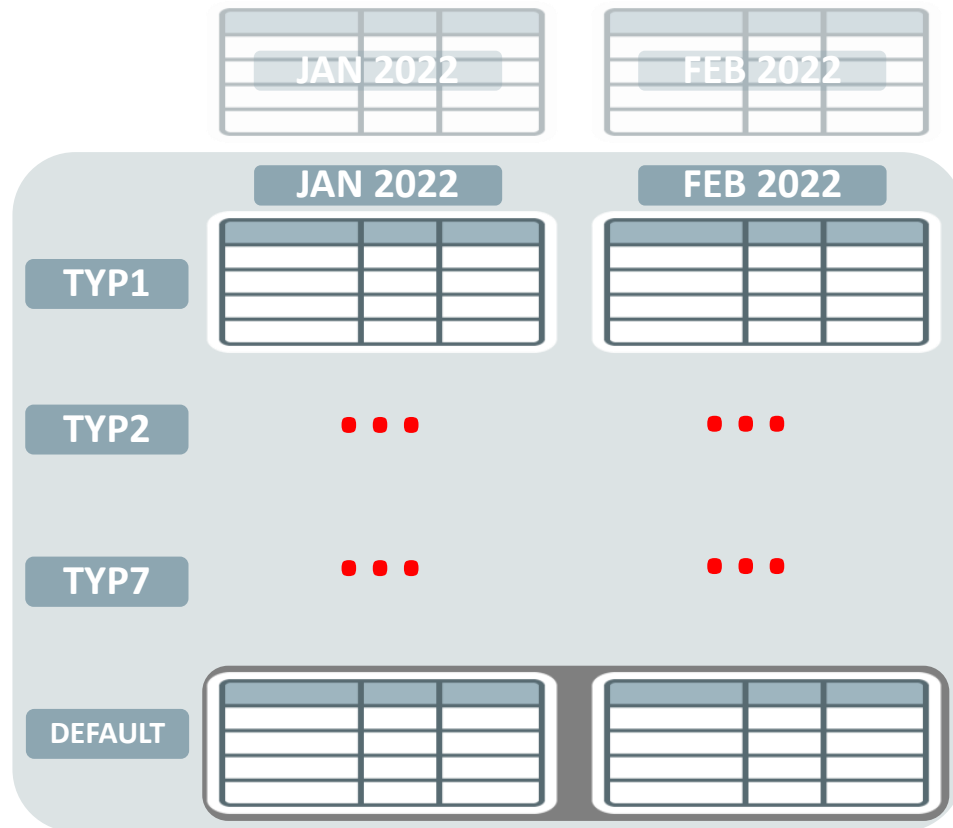


- Number of subpartitions varies for individual partitions
 - Most common for LIST subpartition strategies

```
CREATE TABLE EVENTS..  
PARTITION BY RANGE (time_id)  
SUPARTITION BY LIST (sensor_type)
```

Composite Partitioning

Asymmetric subpartitions

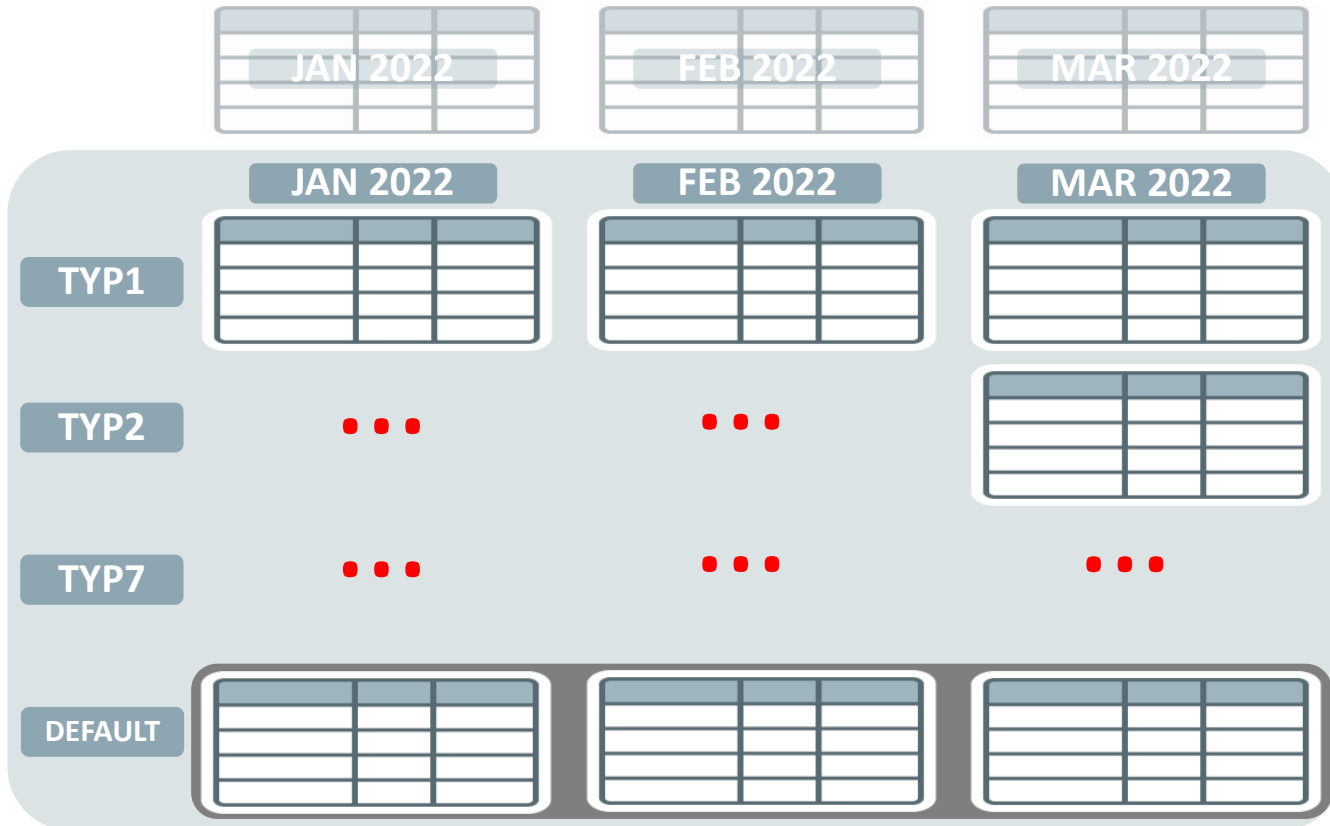


- Number of subpartitions varies for individual partitions
 - Most common for LIST subpartition strategies
- Zero impact on partition pruning capabilities

```
SELECT .. FROM events  
WHERE model = 'TYP7';
```


Composite Partitioning

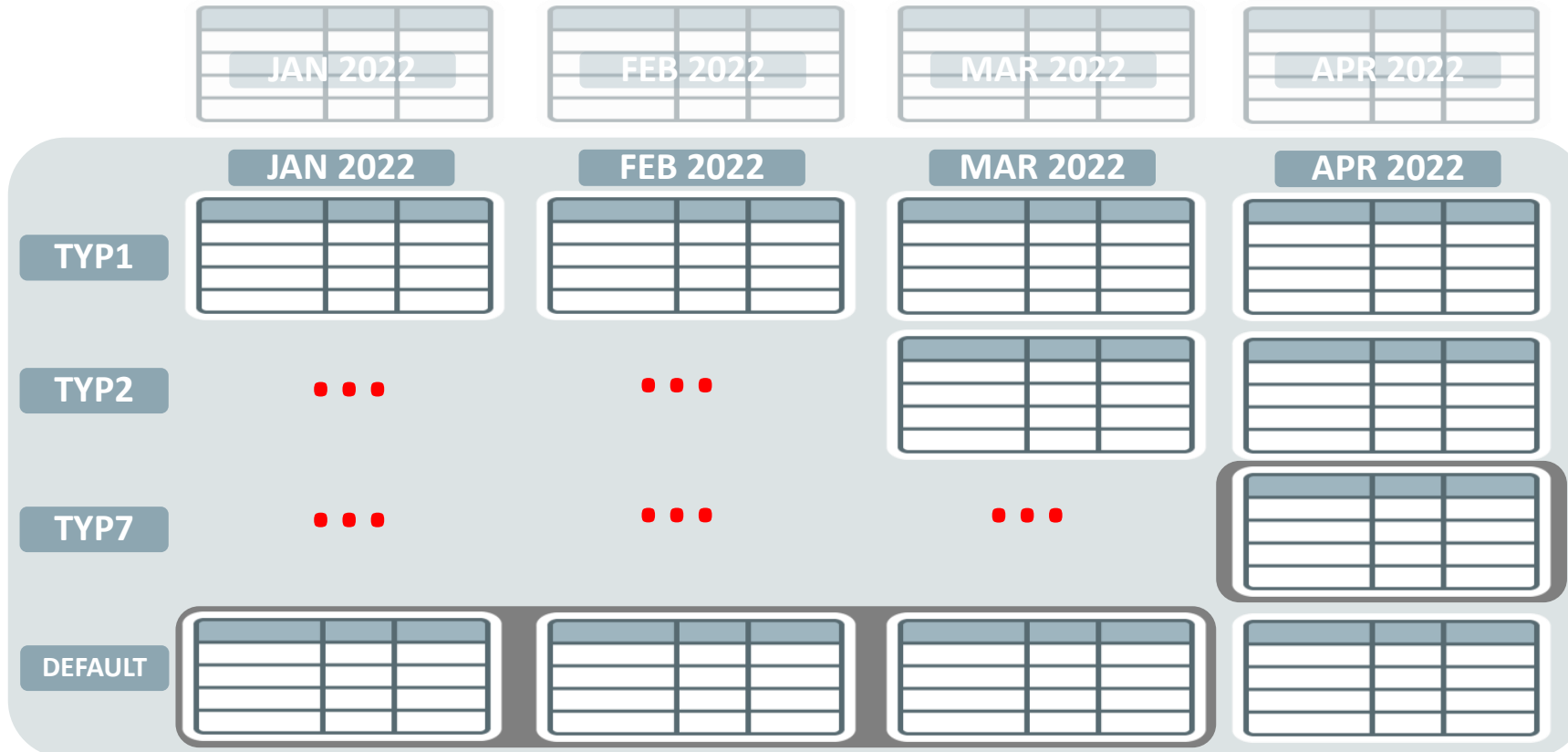
Asymmetric subpartitions



```
SELECT .. FROM events  
WHERE model = 'TYP7';
```

Composite Partitioning

Asymmetric subpartitions



```
SELECT .. FROM events  
WHERE model = 'TYP7';
```

Composite Partitioning

- Always use appropriate composite strategy
 - Top-level dimension mainly chosen for Manageability
 - E.g. add and drop time ranges
 - Sub-level dimension chosen for performance or manageability
 - E.g. load_id, customer_id
 - Asymmetry has advantages but should be thought through
 - E.g. different time granularity for different regions
 - Remember the impact of asymmetric composite partitioning



Multi-Column Range Partitioning

Introduced in Oracle 8i (8.1)



Multi-column Range Partitioning

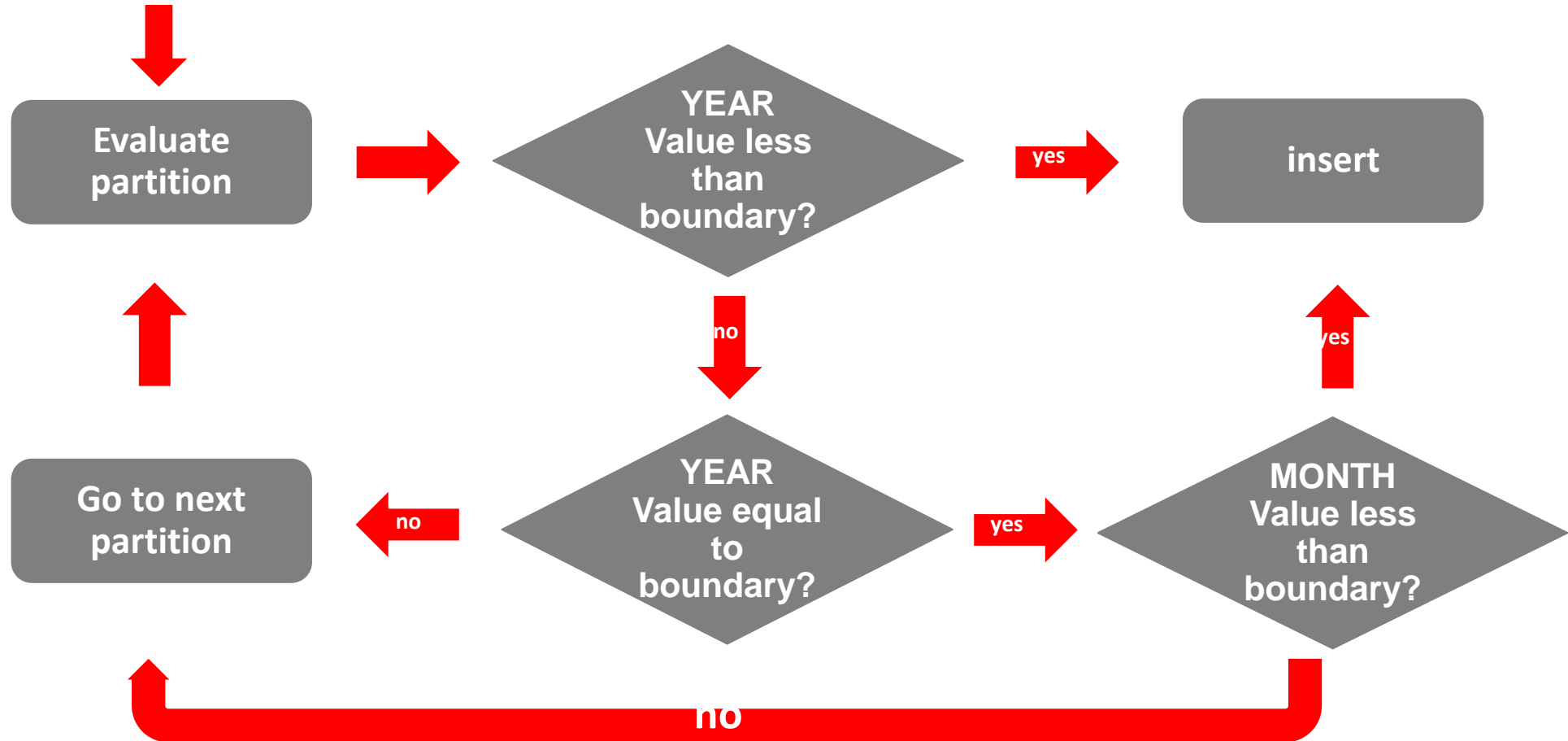
Concept

- Partitioning key is composed of several columns and subsequent columns define a higher granularity than the preceding one
 - E.g. (YEAR, MONTH, DAY)
 - It is NOT an n-dimensional partitioning
- Major watch-out is difference of how partition boundaries are evaluated
 - For simple RANGE, the boundaries are **less than** (exclusive)
 - Multi-column RANGE boundaries are **less than or equal**
 - The n^{th} column is investigated only when all previous (n-1) values of the multicolumn key exactly match the (n-1) bounds of a partition



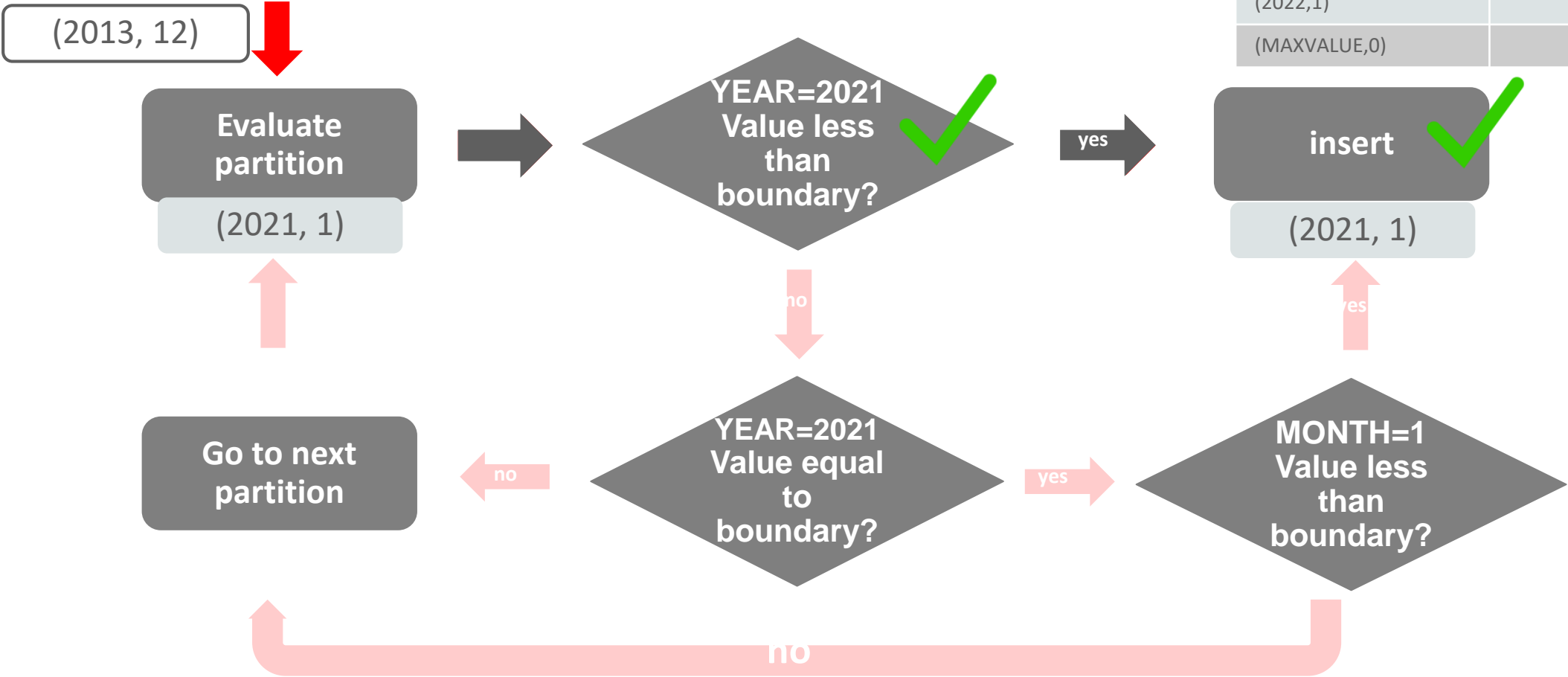
Multi-Column Range Partition

Sample Decision Tree (YEAR, MONTH)



Multi-Column Range Partition

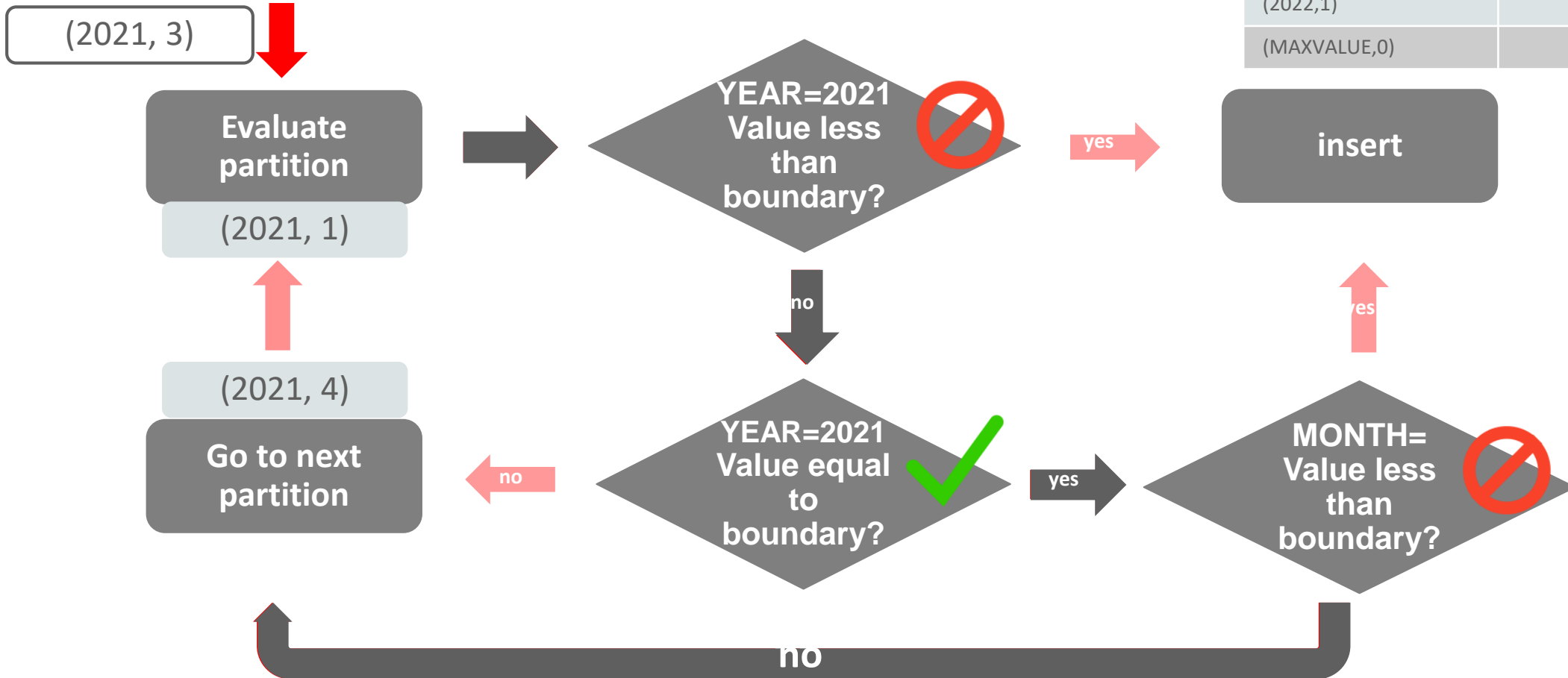
Example



(YEAR,MONTH) Boundaries	Values
(2021,1)	(2013, 12)
(2021,4)	
(2021,7)	
(2021,10)	
(2022,1)	
(MAXVALUE,0)	

Multi-Column Range Partition

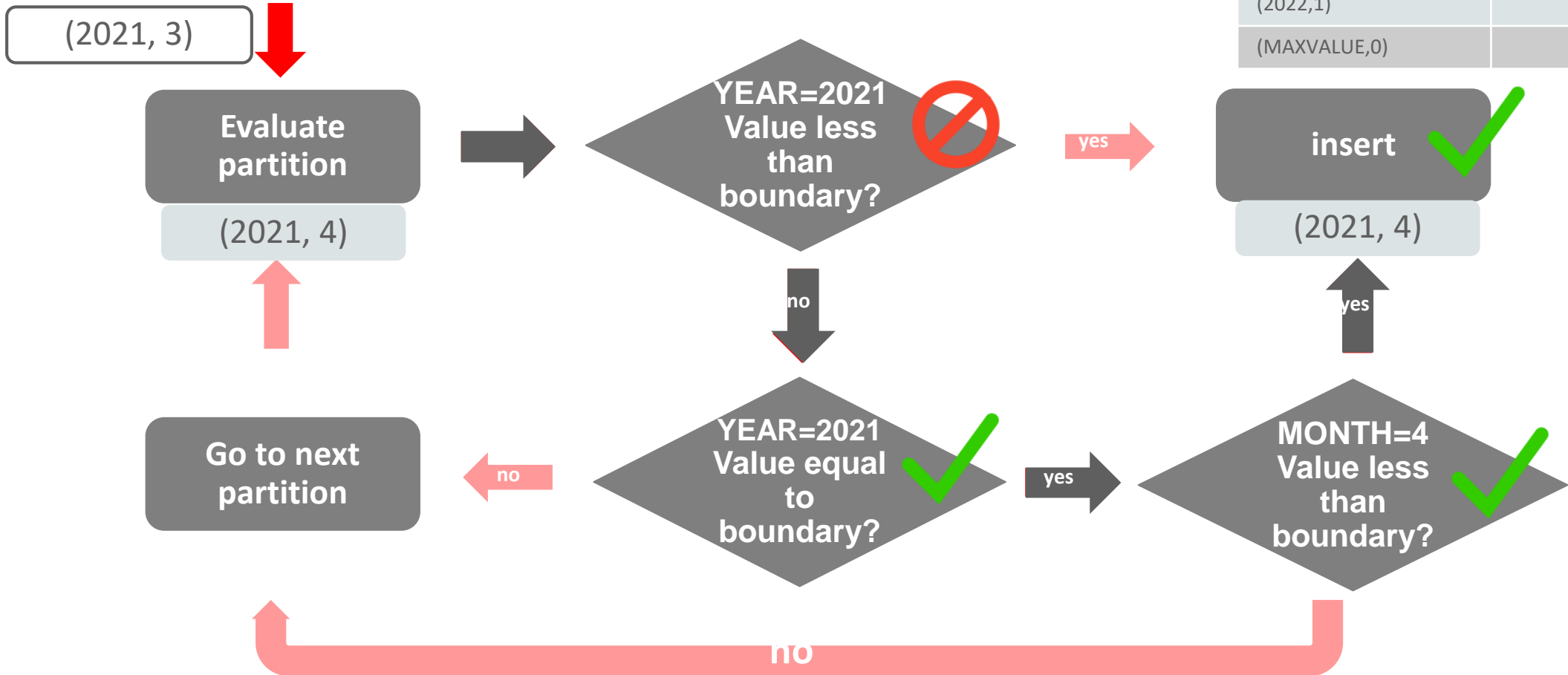
Example Cont'd



(YEAR,MONTH) Boundaries	Values
(2021,1)	(2013, 12)
(2021,4)	
(2021,7)	
(2021,10)	
(2022,1)	
(MAXVALUE,0)	

Multi-Column Range Partition

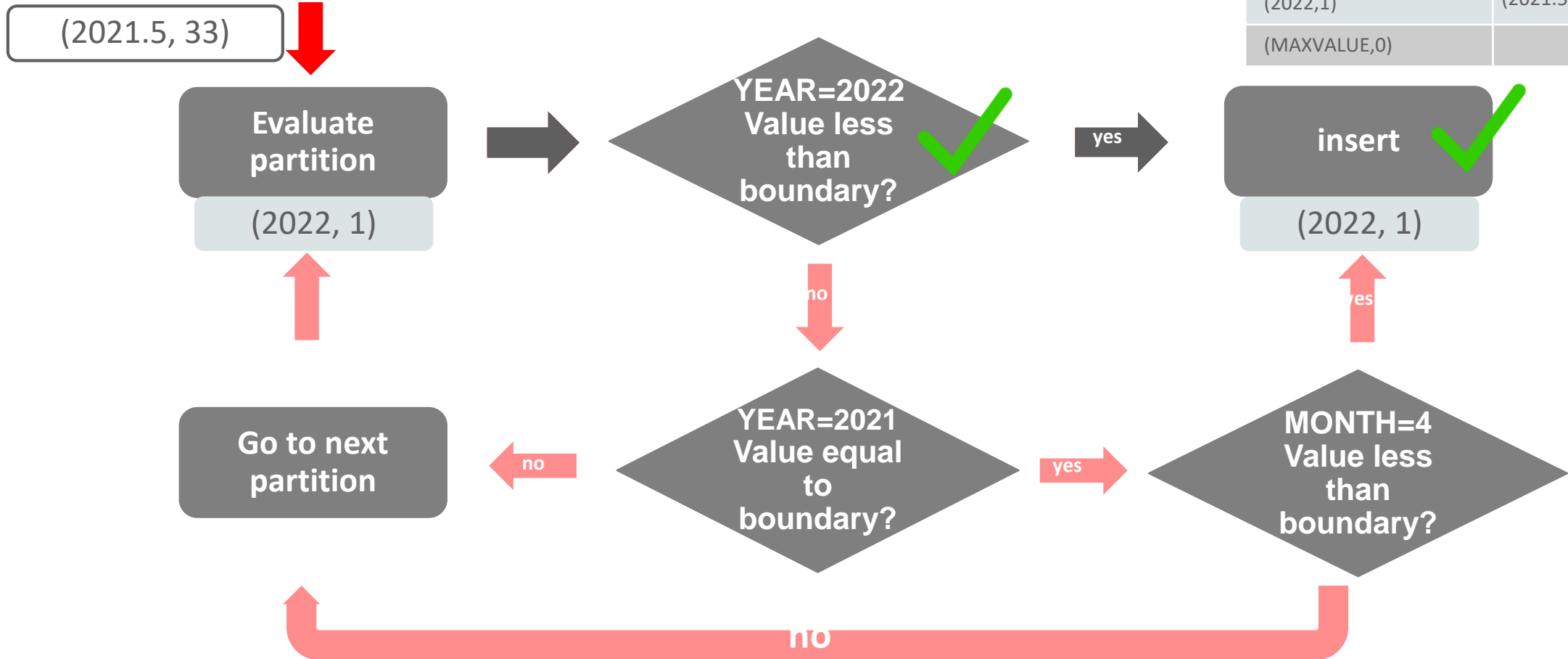
Example Cont'd



(YEAR,MONTH) Boundaries	Values
(2021,1)	(2013, 12)
(2021,4)	(2021, 3)
(2021,7)	
(2021,10)	
(2022,1)	
(MAXVALUE,0)	

Multi-Column Range Partition

Example Cont'd



(YEAR,MONTH) Boundaries	Values
(2021,1)	(2013, 12)
(2021,4)	(2021, 3)
(2021,7)	
(2021,10)	
(2022,1)	(2021.5, 33)
(MAXVALUE,0)	

Multi-Column Range Partitioning

Some things to bear in mind

- ✓ • Powerful partitioning mechanism to add a third (or more) dimensions
 - Smaller data partitions
- Pruning works also for trailing column predicates without filtering the leading column(s)
- △ • Boundaries are not enforced by the partition definition
 - Ranges are consecutive
- Logical ADD partition can mean SPLIT partition in the middle of the table

Multi-Column Range Partition

A slightly different real-world scenario

- Multi-column range used to introduce a third (non-numerical) dimension

```
CREATE TABLE events (event_id number, site_id CHAR(2), start_date date)
PARTITION BY RANGE (site_id, start_date)
SUBPARTITION BY HASH (event_id) SUBPARTITIONS 16
(PARTITION 11_2020 VALUES LESS THAN ('L1', to_date('01-JAN-2021', 'dd-mon-yyyy')),
PARTITION 11_2021 VALUES LESS THAN ('L1', to_date('01-JAN-2022', 'dd-mon-yyyy')),
PARTITION 12_2020 VALUES LESS THAN ('L2', to_date('01-JAN-2021', 'dd-mon-yyyy')),
PARTITION 13_2020 VALUES LESS THAN ('L3', to_date('01-JAN-2021', 'dd-mon-yyyy')),
PARTITION x3_2021 VALUES LESS THAN ('X1', to_date('01-JAN-2022', 'dd-mon-yyyy')),
PARTITION x4_2020 VALUES LESS THAN ('X4', to_date('01-JAN-2021', 'dd-mon-yyyy'))
);
```

Character SITE_ID has to be defined in an ordered fashion



Multi-Column Range Partition

A slightly different real-world scenario

- Multi-column range used to introduce a third (non-numerical) dimension

```
CREATE TABLE events (event_id number, site_id CHAR(2), start_date date)
PARTITION BY RANGE (site_id, start_date)
SUBPARTITION BY HASH (event_id) SUBPARTITIONS 16
(PARTITION 11_2020 VALUES LESS THAN ('L1', to_date('01-JAN-2021', 'dd-mon-yyyy')),
PARTITION 11_2021 VALUES LESS THAN ('L1', to_date('01-JAN-2022', 'dd-mon-yyyy')),
PARTITION 12_2020 VALUES LESS THAN ('L2', to_date('01-JAN-2021', 'dd-mon-yyyy')),
PARTITION 12_2021 VALUES LESS THAN ('L2', to_date('01-JAN-2022', 'dd-mon-yyyy')),
PARTITION x1_2020 VALUES LESS THAN ('X1', to_date('01-JAN-2021', 'dd-mon-yyyy')),
PARTITION x1_2021 VALUES LESS THAN ('X1', to_date('01-JAN-2022', 'dd-mon-yyyy'))
);
```

Non-defined SITE_ID will follow the LESS THAN probe and always end in the lowest partition of a defined SITE_ID



Multi-Column Range Partition

A slightly different real-world scenario

- Multi-column range used to introduce a third (non-numerical) dimension

```
CREATE TABLE events(prod_id number, site_id CHAR(2), start_date date)
PARTITION BY RANGE (site_id, start_date)
SUBPARTITION BY HASH (prod id) SUBPARTITIONS 16
(PARTITION 11_2020 VALUES LESS THAN ('L1',to_date('01-JAN-2014','dd-mon-yyyy')),
PARTITION 11_2021 VALUES LESS THAN ('L1',to_date('01-JAN-2020','dd-mon-yyyy')),
PARTITION 12_2020 VALUES LESS THAN ('L2',to_date('01-JAN-2014','dd-mon-yyyy')),
PARTITION x1_2021 VALUES LESS THAN ('X1',to_date('01-JAN-2020','dd-mon-yyyy')),
PARTITION x4_2020 VALUES LESS THAN ('X4',to_date('01-JAN-2014','dd-mon-yyyy')),
PARTITION x4_2021 VALUES LESS THAN ('X4',to_date('01-JAN-2020','dd-mon-yyyy'))
);
```

Future dates will always go in the lowest partition of the next higher SITE_ID or being rejected



Multi-Column Range Partition

A slightly different real-world scenario

- Multi-column range used to introduce a third (non-numerical) dimension

```
create table events(prod_id number, site_id CHAR(2),start_date date)
partition by range (site_id, start_date)
subpartition by hash (prod id) subpartitions 16
(partition below_l1 values less than ('L1',to_date('01-JAN-1492','dd-mon-yyyy')),
partition l1_2013 values less than ('L1',to_date('01-JAN-2014','dd-mon-yyyy')),
partition l1_2021 values less than ('L1',to_date('01-JAN-2020','dd-mon-yyyy')),
partition l1_max values less than ('L1',MAXVALUE),
partition below_x1 values less than ('X1',to_date('01-JAN-1492','dd-mon-yyyy')),
...
partition x4_max values less than ('X4',MAXVALUE),
partition pmax values less than (MAXVALUE,MAXVALUE));
```

**Introduce a dummy 'BELOW_...' partition
to catch "lower" nondefined SITE_ID**



Multi-Column Range Partition

A slightly different real-world scenario

- Multi-column range used to introduce a third (non-numerical) dimension

```
create table events(prod_id number, site_id CHAR(2),start_date date)
partition by range (site_id, start_date)
subpartition by hash (prod_id) subpartitions 16
(partition below_l1 values less than ('L1',to_date('01-JAN-1492','dd-mon-yyyy')),
 partition l1_2020 values less than ('L1',to_date('01-JAN-2021','dd-mon-yyyy')),
 partition l1_2021 values less than ('L1',to_date('01-JAN-2022','dd-mon-yyyy')),
 partition l1_max values less than ('L1',MAXVALUE),
 partition below_x1 values less than ('X1',to_date('01-JAN-1492','dd-mon-yyyy')),
 ...
 partition x4_max values less than ('X4',MAXVALUE),
 partition pmax values less than (MAXVALUE,MAXVALUE));
```

**Introduce a MAXVALUE 'X_FUTURE' partition
to catch future dates**



Multi-Column Range Partition

A slightly different real-world scenario

- Multi-column range used to introduce a third (non-numerical) dimension

```
create table events(prod_id number, site_id CHAR(2),start_date date)
partition by range (site_id, start_date)
subpartition by hash (prod_id) subpartitions 16
(partition below_l1 values less than ('L1',to_date('01-JAN-1492','dd-mon-yyyy')),
 partition l1_2020 values less than ('L1',to_date('01-JAN-2021','dd-mon-yyyy')),
 partition l1_2021 values less than ('L1',to_date('01-JAN-2022','dd-mon-yyyy')),
 partition l1_max values less than ('L1',MAXVALUE),
 partition below_x1 values less than ('X1',to_date('01-JAN-1492','dd-mon-yyyy')),
 ...
 partition x4_max values less than ('X4',MAXVALUE),
 partition pmax values less than (MAXVALUE,MAXVALUE));
```

If necessary, catch the open-ended SITE_ID (leading key column)

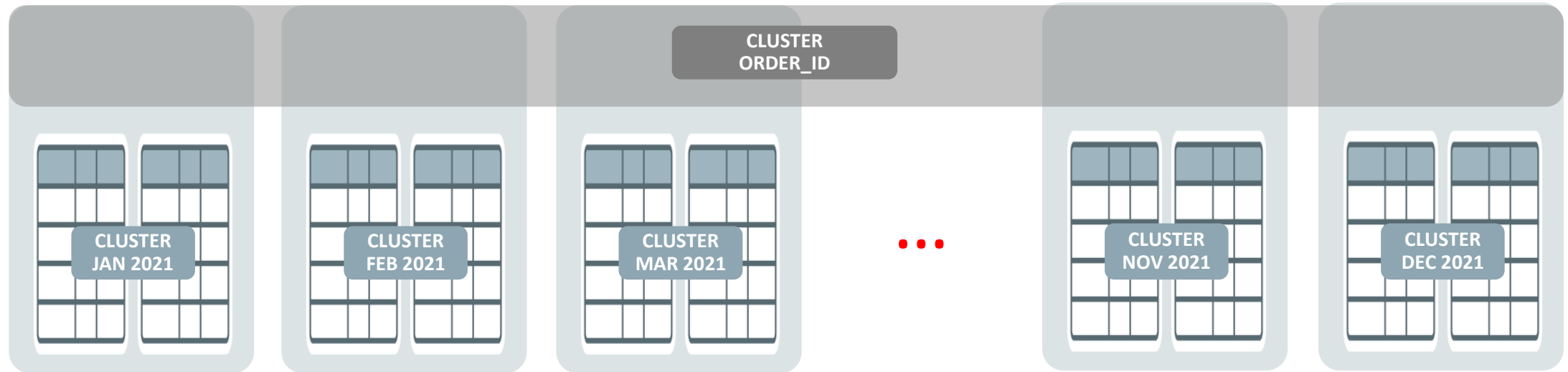


Range-Partitioned Hash Cluster

Introduced in Oracle 12c Release 1 (12.1.0.2)



Range-Partitioned Hash Cluster



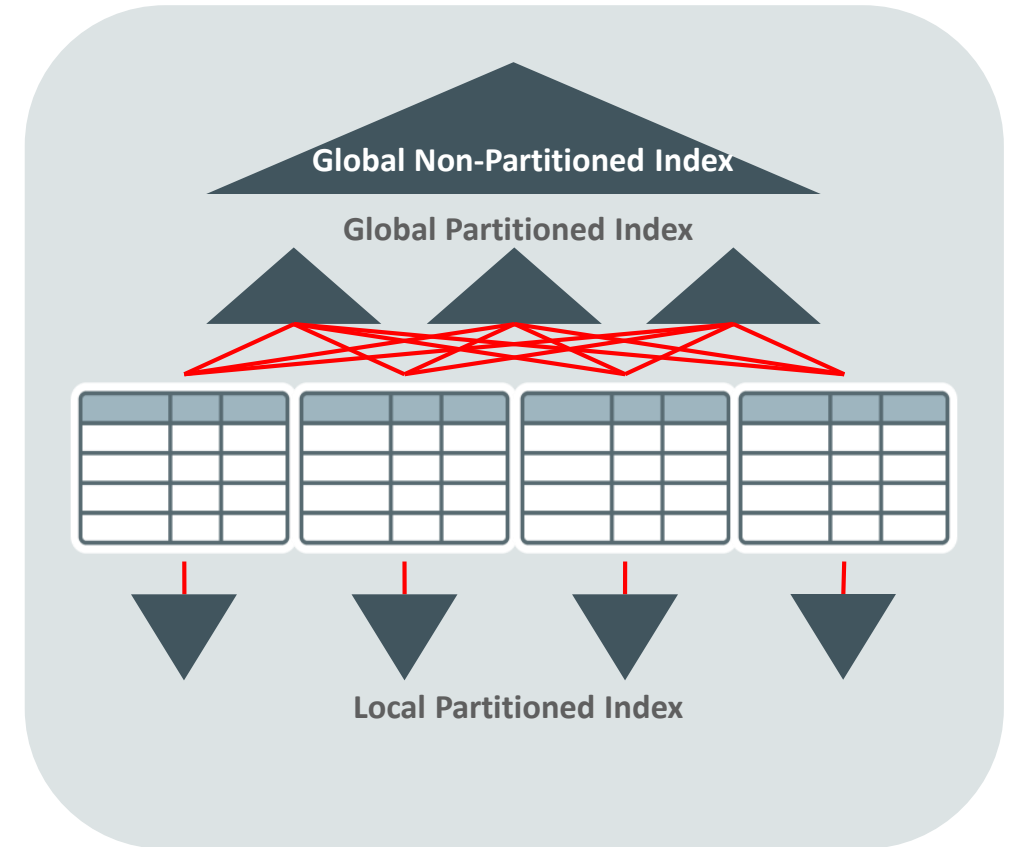
- Single-level range partitioning
 - No composite partitioning
 - No index clusters

Indexing of Partitioned Tables



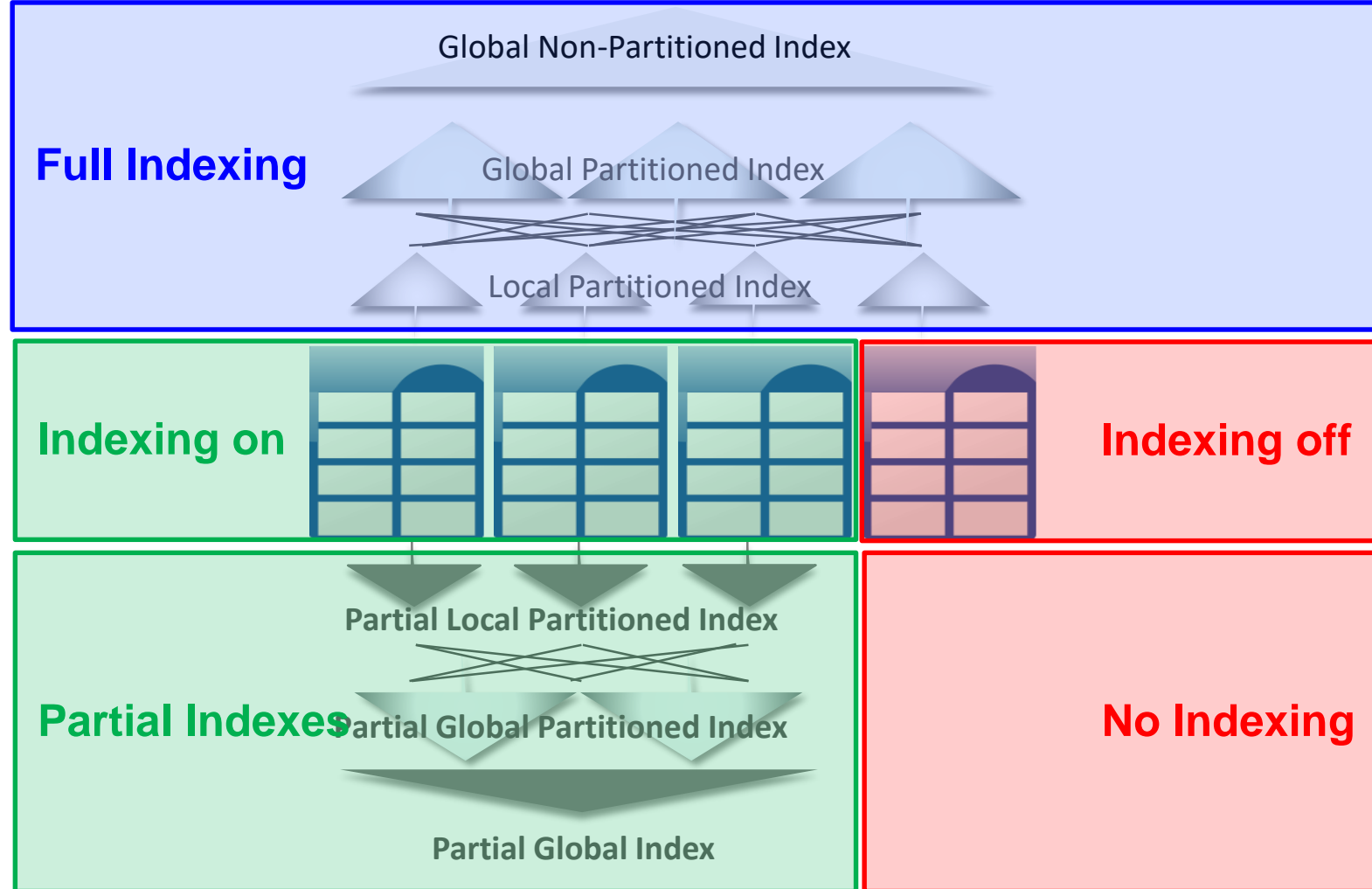
Indexing of Partitioned Tables

- GLOBAL index points to rows in any partition
 - Index can be partitioned or not
- LOCAL index is partitioned same as table
 - Index partitioning key can be different from index key

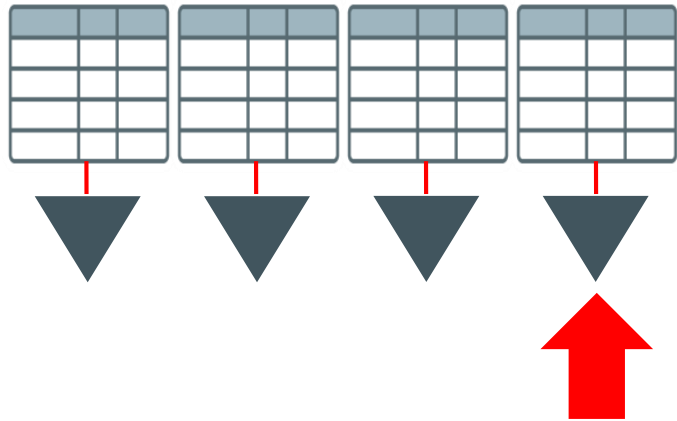


Indexing of Partitioned Tables

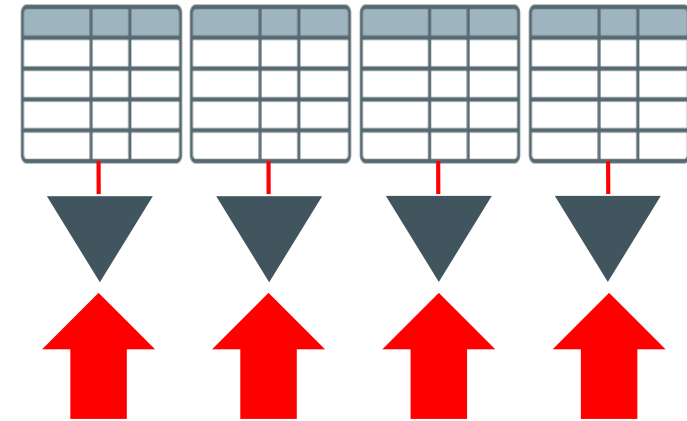
- Partial indexes span only some partitions
- Applicable to local and global indexes
- Complementary to full indexing
- Full support of online index maintenance



Data Access – Local Index and Global Partitioned Index



- Partitioned index access with single partition pruning



- Partitioned index access without any partition pruning

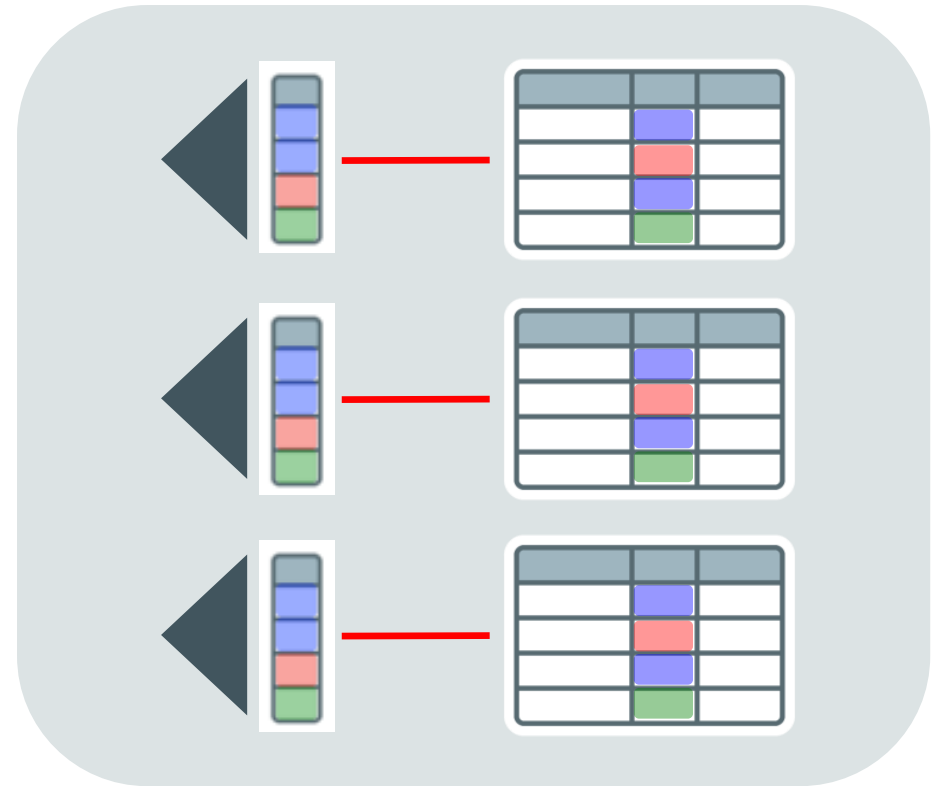
Data Access – Local Index and Global Partitioned Index

- Number of index probes identical to number of accessed partitions
 - No partition pruning leads to a probe into all index partitions
- Not optimally suited for OLTP environments
 - No guarantee to always have partition pruning
 - **Exception:** global hash partitioned indexes for DML contention alleviation
 - Most commonly small number of partitions
- Pruning on global partitioned indexes based on the index prefix
 - Index prefix identical to leading keys of index



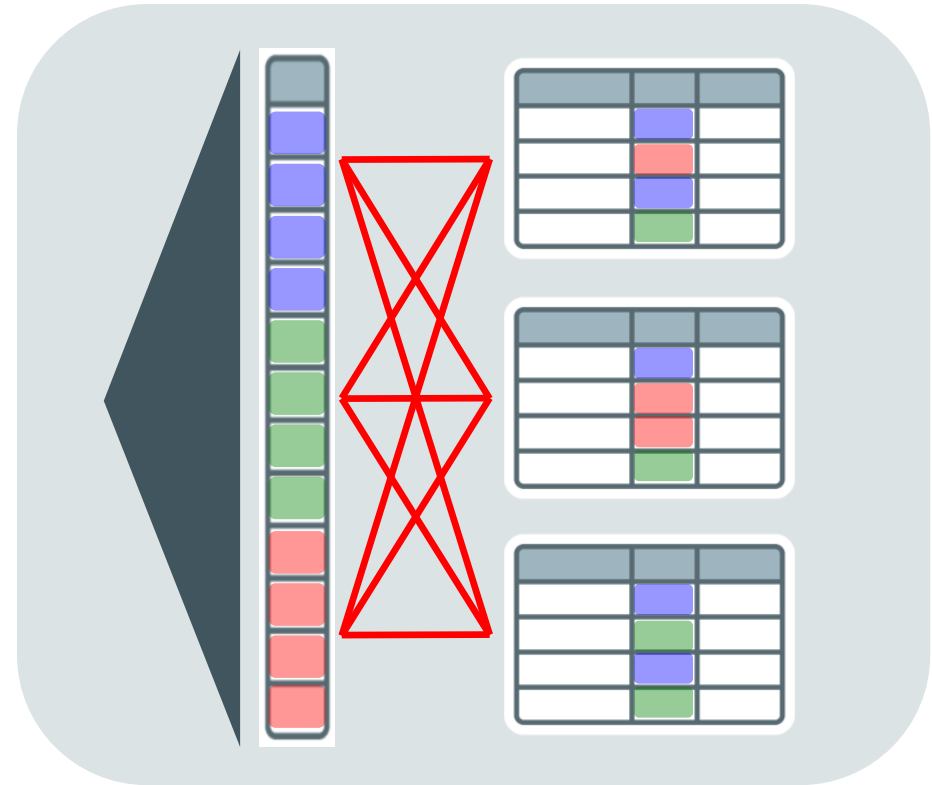
Local Index

- Index is partitioned along same boundaries as data
 - B-tree or bitmap
- Pros
 - Easy to manage
 - Parallel index scans
- Cons
 - Less efficient for retrieving small amounts of data



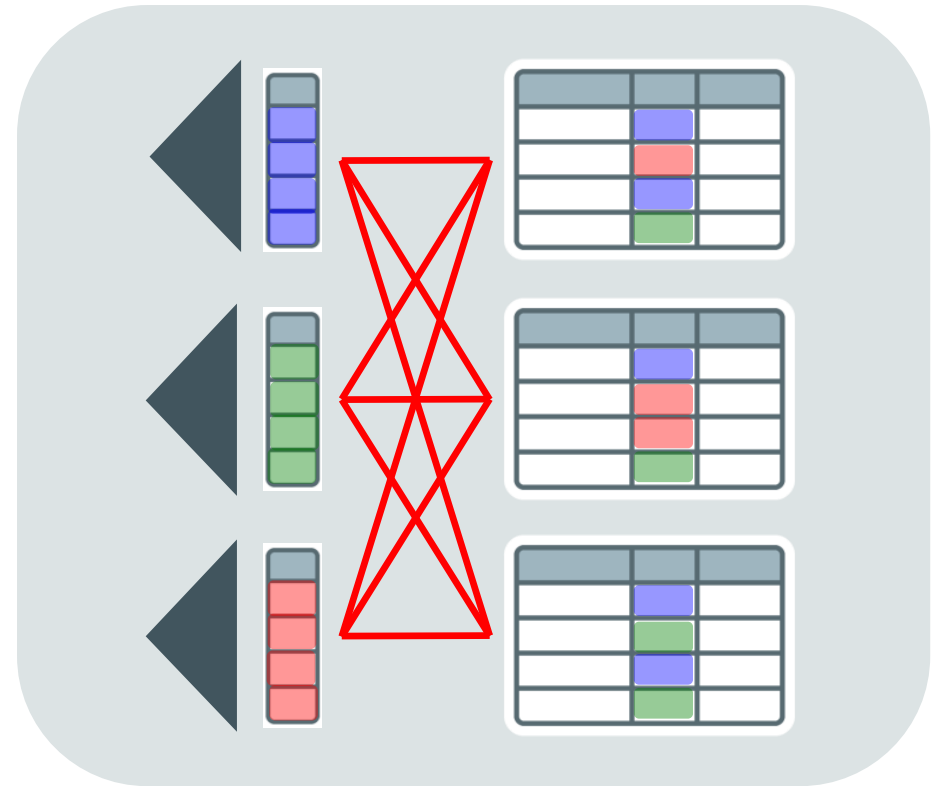
Global Non-Partitioned Index

- One index b-tree structure that spans all partitions
- Pros
 - Efficient access to any individual record
- Cons
 - Partition maintenance always involves index maintenance



Global Partitioned Index

- Index is partitioned independently of data
 - Each index structure may reference any and all partitions.
- Pros
 - Availability and manageability
- Cons
 - Partition maintenance always involves index maintenance



Index Maintenance and Partition Maintenance

- Online index maintenance available for both global and local indexes
 - Global index maintenance since Oracle 9i, local index maintenance since Oracle 10g
- Fast index maintenance for both local and global indexes for DROP and TRUNCATE
 - Asynchronous global index maintenance added in Oracle 12c Release 1
- Index maintenance necessary for both local and global indexes for all other partition maintenance operations



Index Maintenance and Partition Maintenance

- Online index maintenance available for both global and local indexes
 - Global index maintenance since Oracle 9i, local index maintenance since Oracle 10g
- Fast index maintenance for both local and global indexes for DROP and TRUNCATE
 - Asynchronous global index maintenance added in Oracle 12c Release 1
- Index maintenance necessary for both local and global indexes for all other partition maintenance operations
- **Decision for partition maintenance with index maintenance should be always performance versus availability**
 - Rebuild of index always faster when more than 5%-10% of data are touched
- Consider partial indexing for both old and new data
 - Not all data has to be indexed to begin with



Indexing for unique constraints and primary keys

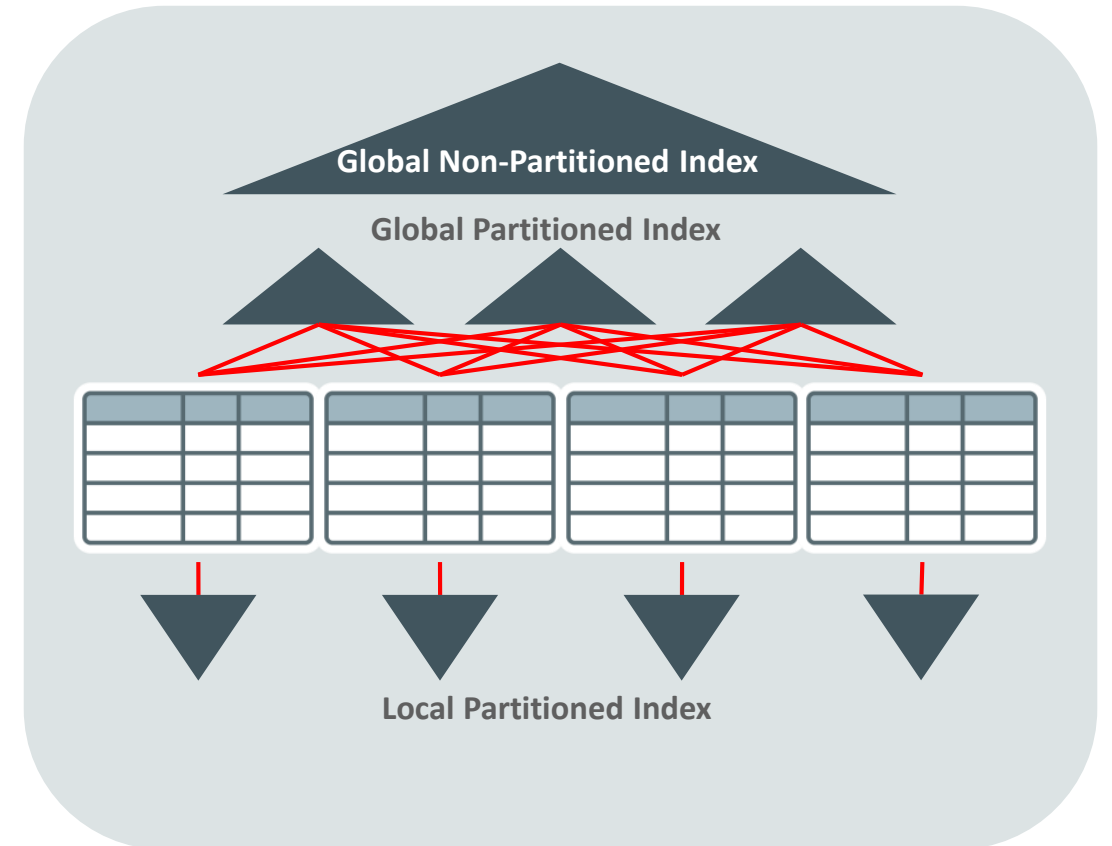
Unique Constraints/Primary Keys

- Unique constraints are enforced with unique indexes
 - Primary key constraint adds NOT NULL to column
 - Table can have only one primary key (“unique identifier”)
- Partitioned tables offer two types of indexes
 - Local indexes
 - Global index, both partitioned and non-partitioned
- Which one to pick?
 - Do I even have a choice?



Index Partitioning

- GLOBAL index points to rows in any partition
 - Index can be partitioned or not
 - Partition maintenance affects entire index
- LOCAL index is partitioned same as table
 - Index partitioning key can be different from index key
 - Index partitions can be maintained separately



Unique Constraints/Primary Keys

Applicability of Local Indexes

- Local indexes are equi-partitioned with the table
 - Follow autonomy concept of a table partition
 - “I only care about myself”
- Requirement for local indexes to enforce uniqueness
 - Partition key column(s) to be a subset of the unique key



Unique Constraints/Primary Keys, cont.

Applicability of Local Indexes

- Local indexes are equi-partitioned with the table
 - Follow autonomy concept of a table partition
 - “I only care about myself”
- Requirement for local indexes to enforce uniqueness
 - Partition key column(s) to be a subset of the unique key



```
PARTITION BY (col1), PK(col1)
```



```
PARTITION BY (col1), PK(col2)
```

Unique Constraints/Primary Keys, cont.

Applicability of Global Indexes

- Global indexes do not have any relation to the partitions of a table
 - By definition, a global index contains data from all partitions
 - True for both partitioned and non-partitioned global indexes
- Global index can always be used to enforce uniqueness



PARTITION BY (col1), PK(col1)



PARTITION BY (col1), PK(col2)

Partial Indexing

Introduced in Oracle 12c Release 1 (12.1)



Enhanced Indexing with Oracle Partitioning

Indexing prior to Oracle Database 12c

- Local indexes
- Non-partitioned or partitioned global indexes
- Usable or unusable index segments
 - Non-persistent status of index, no relation to table



Enhanced Indexing with Oracle Partitioning

Indexing with Oracle Database 12c

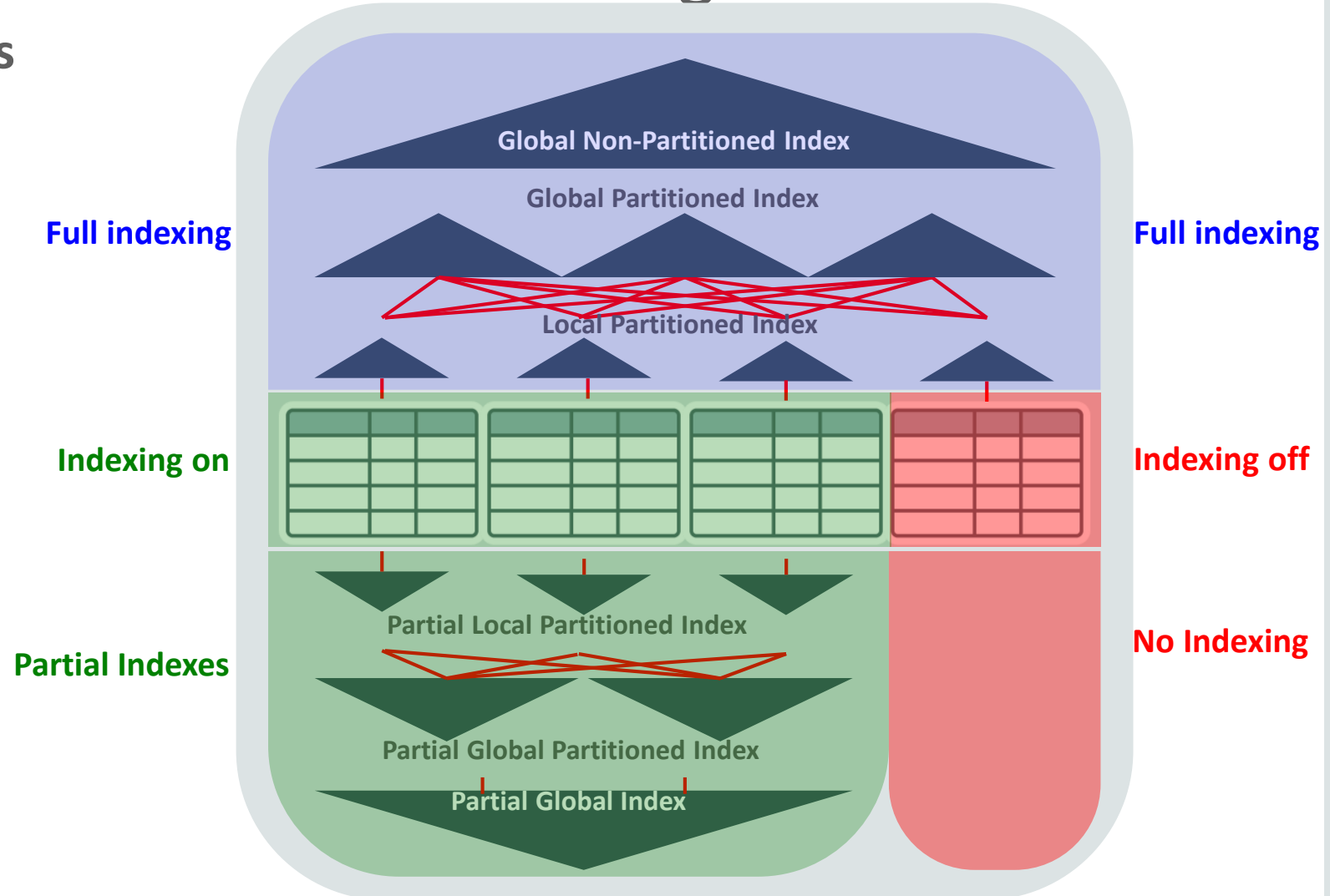
- Local indexes
- Non-partitioned or partitioned global indexes
- Usable or unusable index segments
 - Non-persistent status of index, no relation to table
- Partial local and global indexes
 - Partial indexing introduces table and [sub]partition level metadata
 - Leverages usable/unusable state for local partitioned indexes
 - Policy for partial indexing can be overwritten



Enhanced Indexing with Oracle Partitioning

Partial Local and Global Indexes

- Partial indexes span only some partitions
- Applicable to local and global indexes
- Complementary to full indexing
- Enhanced business modeling



Enhanced Indexing with Oracle Partitioning

Partial Local and Global Indexes

Before

```
SQL> create table pt (col1, col2, col3, col4)
  2  indexing off
  3  partition by range (col1)
  4  interval (1000)
  5  (partition p100 values less than (101) indexing on,
  6  partition p200 values less than (201) indexing on,
  7  partition p300 values less than (301) indexing on);
```

Table created.

```
SQL> REM partitions and its indexing status
SQL> select partition_name, high_value, indexing
  2  from user_tab_partitions where table_name='PT';
```

PARTITION_NAME	HIGH_VALUE	INDEXING
P100	101	ON
P200	201	ON
P300	301	ON
SYS_P1256	1301	OFF

After

```
SQL> REM local indexes
SQL> create index i_l_partpt on pt(col1) local indexing partial;
SQL> create index i_l_pt on pt(col4) local;

SQL> REM global indexes
SQL> create index i_g_partpt on pt(col2) indexing partial;
SQL> create index i_g_pt on pt(col3);
```

```
SQL> REM index status
SQL> select index_name, partition_name, status, null
  2  from user_ind_partitions where index_name in ('I_L_PARTPT','I_L_PT')
  3  union all
  4  select index_name, indexing, status, orphaned_entries
  5  from user_indexes where index_name in ('I_G_PARTPT','I_G_PT');
```

INDEX_NAME	PARTITION_NAME	STATUS	ORPHAN
I_L_PARTPT	P100	USABLE	
I_L_PARTPT	P200	USABLE	
I_L_PARTPT	P300	USABLE	
I_L_PARTPT	SYS_P1257	UNUSABLE	
I_L_PT	P200	USABLE	
I_L_PT	P300	USABLE	
I_L_PT	SYS_P1258	USABLE	
I_L_PT	P100	USABLE	
I_G_PT	FULL	VALID	NO
I_G_PARTPT	PARTIAL	VALID	NO

10 rows selected.



Enhanced Indexing with Oracle Partitioning

Partial Local and Global Indexes

- Partial global index excluding partition 4

```
SQL> explain plan for select count(*) from pt where col2 = 3;
```

Explained.

```
SQL> select * from table(dbms_xplan.display);
```

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT		1	22	54 (12)	00:00:01		
1	SORT AGGREGATE		1	22				
2	VIEW	VW_TE_2	2		54 (12)	00:00:01		
3	UNION-ALL							
* 4	TABLE ACCESS BY GLOBAL INDEX ROWID BATCHED	PT	1	26	2 (0)	00:00:01	ROWID	ROWID
* 5	INDEX RANGE SCAN	I_G_PARTPT	1		1 (0)	00:00:01		
6	PARTITION RANGE SINGLE		1	26	52 (12)	00:00:01	4	4
* 7	TABLE ACCESS FULL	PT	1	26	52 (12)	00:00:01	4	4

Predicate Information (identified by operation id):

```
4 - filter("PT"."COL1"<301)
5 - access("COL2"=3)
7 - filter("COL2"=3)
```



Unusable versus Partial Indexes



Unusable Indexes

- Unusable index partitions are commonly used in environments with fast load requirements
 - “Save” the time for index maintenance at data insertion
 - Unusable index segments do not consume any space (11.2)
- Unusable indexes are ignored by the optimizer

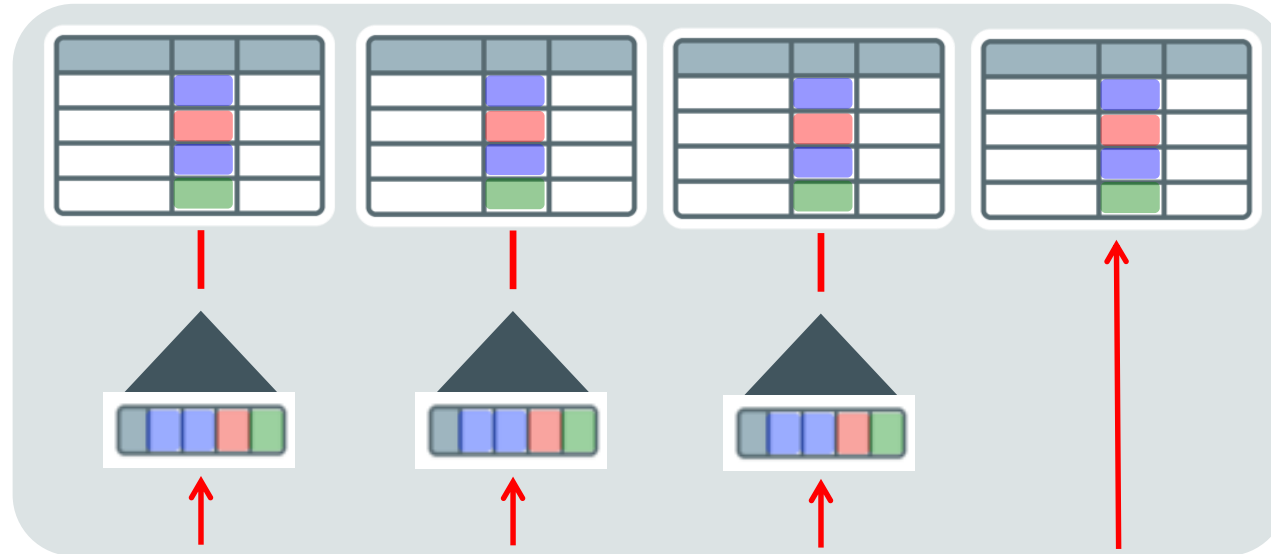
```
SKIP_UNUSABLE_INDEXES = [TRUE | FALSE ]
```

- Partitioned indexes can be used by the optimizer even if some partitions are unusable
 - Prior to 11.2, static pruning and only access of usable index partitions mandatory
 - With 11.2, intelligent rewrite of queries using UNION ALL



Table-OR-Expansion

Multiple SQL branches are generated and executed



- Intelligent UNION ALL expansion in the presence of partially unusable indexes
 - Transparent internal rewrite
 - Usable index partitions will be used
 - Full partition access for unusable index partitions

Table-OR-Expansion

Sample Plan - Multiple SQL branches are generated and executed

```
select count(*) from toto where name = 'FOO' and rn between 1300 and 1400
```

Plan hash value: 2830852558

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				27M(100)			
1	SORT AGGREGATE		1	21				
2	VIEW	VW_TE_2	2		27M (3)	92:15:22		
3	UNION-ALL							
4	PARTITION RANGE SINGLE		1	20	2 (0)	00:00:01	14	14
5	TABLE ACCESS BY LOCAL INDEX ROWID	TOTO	1	20	2 (0)	00:00:01	14	14
* 6	INDEX RANGE SCAN	I_TOTO	1		1 (0)	00:00:01	14	14
7	PARTITION RANGE SINGLE		1	22	27M (3)	92:15:22	15	15
* 8	TABLE ACCESS FULL	TOTO	1	22	27M (3)	92:15:22	15	15

Predicate Information (identified by operation id):

```
6 - access("NAME"='FOO')
8 - filter(("NAME"='FOO' AND "TOTO"."RN">=1400))
```

27 rows selected.



Partitioning for Performance



Partitioning for Performance

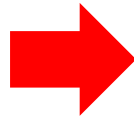
- Partitioning is transparently leveraged to improve performance
- Partition pruning
 - Using partitioning metadata to access only partitions of interest
- Partition-wise joins
 - Join equi-partitioned tables with minimal resource consumption
 - Process co-location capabilities for RAC environments
- Partition-Exchange loading
 - “Load” new data through metadata operation



Partitioning for Performance

Partition Pruning

What are the total
EVENTS for May 1-2?



EVENTS		
		May 5
		May 4
		May 3
		May 2
		May 1
		Apr 30
		Apr 29
		Apr 28
		Apr 27

- Partition elimination
 - Dramatically reduces amount of data retrieved from storage
 - Performs operations only on relevant partitions
 - Transparently improves query performance and optimizes resource utilization

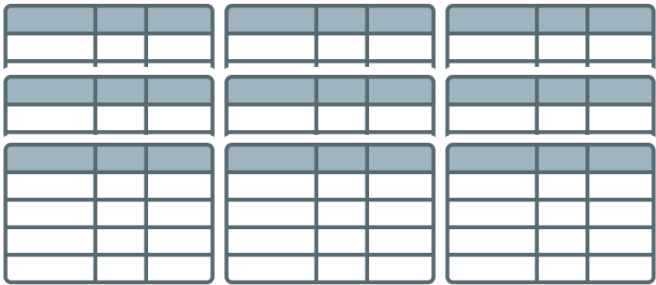


Partition Pruning

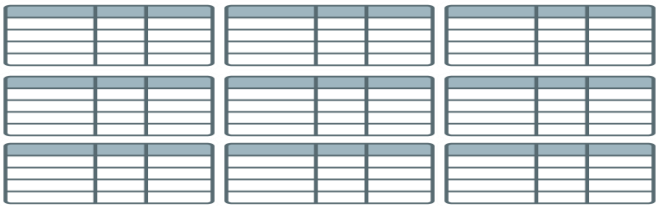
- Works for simple and complex SQL statements
- Transparent to any application
- Two flavors of pruning
 - Static pruning at compile time
 - Dynamic pruning at runtime
- Complementary to Exadata Storage Server
 - Partitioning prunes logically through partition elimination
 - Exadata prunes physically through storage indexes
 - Further data reduction through filtering and projection



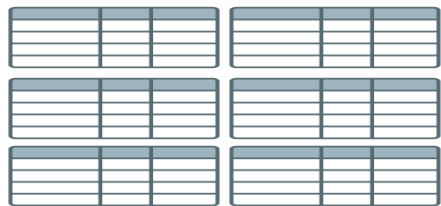
Performance Features Multiply the Benefits /example



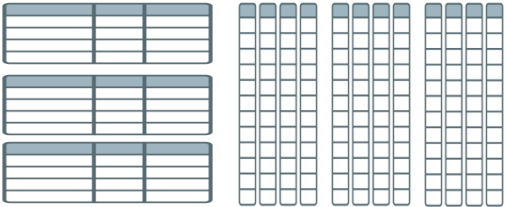
100 TB of User Data



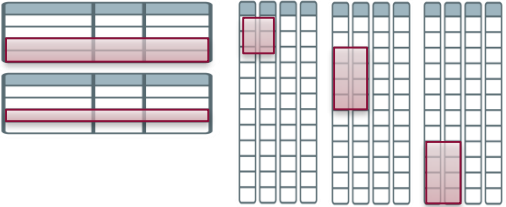
10 TB of User Data
With 10x Compression



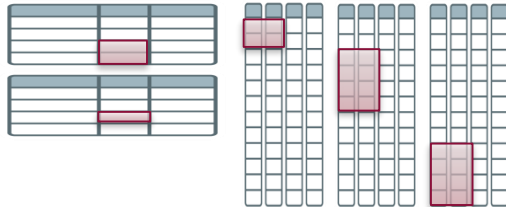
2TB of User Data
With Partition Pruning



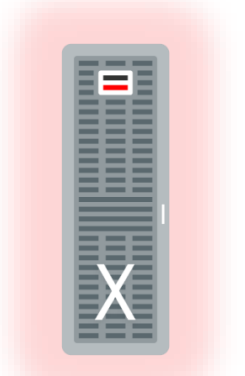
2 TB of User Data
1TB on disk, 1TB in-memory



100 GB of User Data
With Storage Indexes
and Zone Maps



30 GB of User Data
With Smart Scan



Sub second Scan
No Indexes



Static Partition Pruning

```
SELECT avg( luminosity ) FROM EVENTS  
WHERE times id  
BETWEEN '01-MAR-2021' and '31-MAY-2021';
```

2021-JAN

2021-FEB

2021-MAR

2021-APR

2021-MAY

2021-JUN

- Relevant Partitions are known at compile time
 - Look for actual values in PSTART/PSTOP columns in the plan
- Optimizer has most accurate information for the SQL statement



Static Pruning

Sample Plan

```
SELECT avg( luminosity )
FROM atlas.EVENTS s, atlas.times t
WHERE s.time_id = t.time_id
AND s.time_id between TO_DATE('01-JAN-2021', 'DD-MON-YYYY')
                    and TO_DATE('01-JAN-2020', 'DD-MON-YYYY')
```

Plan hash value: 2025449199

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				3 (100)			
1	SORT AGGREGATE		1	12				
2	PARTITION RANGE ITERATOR		313	3756	3 (0)	00:00:01	9	13
* 3	TABLE ACCESS FULL	EVENTS	313	3756	3 (0)	00:00:01	9	13

Predicate Information (identified by operation id):

```
3 - filter("S"."TIME_ID"<=TO_DATE(' 2020-01-01 00:00:00', 'syyy-mm-dd hh24:mi:ss'))
```

22 rows selected.



Static Pruning

Sample Plan

```
SELECT avg( luminosity )
FROM atlas.EVENTS s, atlas.times t
WHERE s.time_id = t.time_id
AND s.time_id between TO_DATE('01-JAN-2021', 'DD-MON-YYYY')
    and TO_DATE('01-JAN-2020', 'DD-MON-YYYY')
```

Plan hash value: 2025449199

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				3 (100)			
1	SORT AGGREGATE		1	12				
2	PARTITION RANGE ITERATOR		313	3756	3 (0)	00:00:01	9	13
* 3	TABLE ACCESS FULL	EVENTS	313	3756	3 (0)	00:00:01	9	13

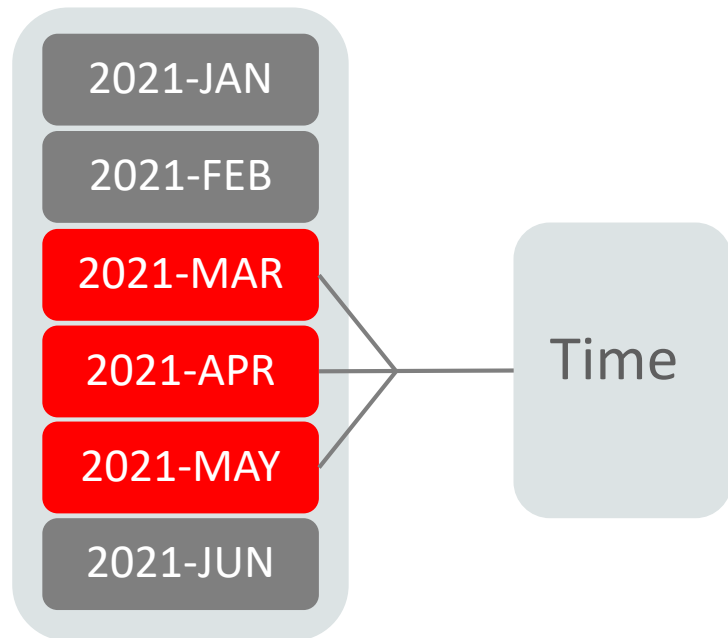
Predicate Information (identified by operation id):

```
3 - filter("S"."TIME_ID"<=TO_DATE(' 2020-01-01 00:00:00', 'syyyymm-dd hh24:mi:ss'))
```

22 rows selected.



Dynamic Partition Pruning



```
SELECT avg( luminosity )  
FROM EVENTS s, times t  
WHERE t.time_id = s.time_id  
AND t.calendar_month_desc IN  
      ('MAR-2021', 'APR-2021', 'MAY-2021');
```

- Advanced Pruning mechanism for complex queries
- Relevant partitions determined at runtime
 - Look for the word 'KEY' in PSTART/PSTOP columns in the Plan

Dynamic Partition Pruning

Sample Plan – Nested Loop

```
SELECT avg( luminosity )
FROM atlas.EVENTS s, atlas.times t
WHERE s.time_id = t.time_id
AND t.calendar_month_desc in ('MAR-2021', 'APR-2021', 'MAY-2021')
```

Plan hash value: 1350851517

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				13 (100)			
1	SORT AGGREGATE		1	28				
2	NESTED LOOP		2	56	13 (0)	00:00:01		
* 3	TABLE ACCESS FULL	TIMES	2	32	13 (8)	00:00:01		
4	PARTITION RANGE ITERATOR		2	24	0 (0)		KEY	KEY
* 5	TABLE ACCESS FULL	EVENTS	2	24	0 (0)		KEY	KEY

Predicate Information (identified by operation id):

- 3 - filter(("T"."CALENDAR_MONTH_DESC"='MAR-2021' OR "T"."CALENDAR_MONTH_DESC"='APR-2021' OR "T"."CALENDAR_MONTH_DESC"='MAY-2021'))
- 5 - filter("T"."TIME_ID"="S"."TIME_ID")

26 rows selected.



Dynamic Partition Pruning

Sample Plan – Nested Loop

```
SELECT avg( luminosity )
FROM atlas.EVENTS s, atlas.times t
WHERE s.time_id = t.time_id
AND t.calendar_month_desc in ('MAR-2021', 'APR-2021', 'MAY-2021')
```

Plan hash value: 1350851517

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				13 (100)			
1	SORT AGGREGATE		1	28				
2	NESTED LOOP		2	56	13 (0)	00:00:01		
* 3	TABLE ACCESS FULL	TIMES	2	32	13 (8)	00:00:01		
4	PARTITION RANGE ITERATOR		2	24	0 (0)		KEY	KEY
* 5	TABLE ACCESS FULL	EVENTS	2	24	0 (0)		KEY	KEY

Predicate Information (identified by operation id):

```
3 - filter(("T"."CALENDAR_MONTH_DESC"='MAR-2021' OR "T"."CALENDAR_MONTH_DESC"='APR-2021'
          OR "T"."CALENDAR_MONTH_DESC"='MAY-2021'))
5 - filter("T"."TIME_ID"="S"."TIME_ID")
```

26 rows selected.



Dynamic Partition Pruning

Sample Plan - Subquery pruning

```
SELECT avg( luminosity )
FROM atlas.EVENTS s, atlas.times t
WHERE s.time_id = t.time_id
AND t.calendar_month_desc in ('MAR-2021', 'APR-2021', 'MAY-2021')
```

Plan hash value: 2475767165

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				2000K(100)			
1	SORT AGGREGATE		1	28				
* 2	HASH JOIN		24M	646M	2000K(100)	06:40:01		
* 3	TABLE ACCESS FULL	TIMES	2	32	43 (8)	00:00:01		
4	PARTITION RANGE SUBQUERY		10G	111G	1166K(100)	03:53:21	KEY(SQ)	KEY(SQ)
5	TABLE ACCESS FULL	EVENTS	10G	111G	1166K(100)	03:53:21	KEY(SQ)	KEY(SQ)

Predicate Information (identified by operation id):

- 2 - access("S"."TIME_ID"="T"."TIME_ID")
- 3 - filter(("T"."CALENDAR_MONTH_DESC"='MAR-2021' OR "T"."CALENDAR_MONTH_DESC"='APR-2021' OR "T"."CALENDAR_MONTH_DESC"='MAY-2021'))

26 rows selected.



Dynamic Partition Pruning

Sample Plan - Bloom filter pruning

```
SELECT avg( luminosity )
FROM atlas.EVENTS s, atlas.times t
WHERE s.time_id = t.time_id
AND t.calendar_month_desc in ('MAR-2021', 'APR-2021', 'MAY-2021')
```

Plan hash value: 365741303

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				19 (100)			
1	SORT AGGREGATE		1	28				
* 2	HASH JOIN		2	56	19 (100)	00:00:01		
3	PART JOIN FILTER CREATE	:BF0000	2	32	13 (8)	00:00:01		
* 4	TABLE ACCESS FULL	TIMES	2	32	13 (8)	00:00:01		
5	PARTITION RANGE JOIN-FILTER		960	11520	5 (0)	00:00:01	:BF0000	:BF0000
6	TABLE ACCESS FULL	EVENTS	960	11520	5 (0)	00:00:01	:BF0000	:BF0000

Predicate Information (identified by operation id):

```
2 - access("S"."TIME ID"="T"."TIME ID")
4 - filter(("T"."CALENDAR_MONTH_DESC"='MAR-2021' OR "T"."CALENDAR_MONTH_DESC"='APR-2021'
OR "T"."CALENDAR_MONTH_DESC"='MAY-2021'))
```

27 rows selected.



“AND” Pruning

```
FROM events s, times t ...  
WHERE s.time id = t.time id ..  
AND t.fiscal year in (2021,2020)  
AND s.time id  
  between TO DATE('01-JAN-2021','DD-MON-YYYY')  
  and TO DATE('01-JAN-2022','DD-MON-YYYY')
```

Dynamic pruning

Static pruning

- All predicates on partition key will be used for pruning
 - Dynamic and static predicates will now be used combined
- Example:
 - Star transformation with pruning predicate on both the FACT table and a dimension



“AND” Pruning

Sample Plan

Plan hash value: 552669211

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT		1	24	17 (12)	00:00:01		
1	SORT AGGREGATE		1	24				
* 2	HASH JOIN		204	4896	17 (12)	00:00:01		
3	PART JOIN FILTER CREATE	:BF0000	185	2220	13 (8)	00:00:01		
* 4	TABLE ACCESS FULL	TIMES	185	2220	13 (8)	00:00:01		
5	PARTITION RANGE AND		313	3756	3 (0)	00:00:01	KEY (AP)	KEY (AP)
* 6	TABLE ACCESS FULL	EVENTS	313	3756	3 (0)	00:00:01	KEY (AP)	KEY (AP)

Predicate Information (identified by operation id):

```
2 - access("S"."TIME_ID"="T"."TIME_ID")
4 - filter("T"."TIME_ID"<=TO_DATE(' 2020-01-01 00:00:00', 'syyy-mm-dd hh24:mi:ss') AND
          ("T"."FISCAL_YEAR"=2021 OR "T"."FISCAL_YEAR"=2020) AND "T"."TIME_ID">=TO_DATE(' 2021-01-01
          00:00:00', 'syyy-mm-dd hh24:mi:ss'))
6 - filter("S"."TIME_ID"<=TO_DATE(' 2020-01-01 00:00:00', 'syyy-mm-dd hh24:mi:ss'))
```

22 rows selected.



Ensuring Partition Pruning

Don't use functions on partition key filter predicates

```
SELECT avg( luminosity )  
FROM atlas.EVENTS s, atlas.times t  
WHERE s.time_id = t.time_id  
AND TO_CHAR(s.time_id, 'YYYYMMDD') between '20210101' and '20220101'
```

Plan hash value: 672559287

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				6 (100)			
1	SORT AGGREGATE		1	12				
2	PARTITION RANGE ALL		2	24	6 (17)	00:00:01	1	16
* 3	TABLE ACCESS FULL	EVENTS	2	24	6 (17)	00:00:01	1	16

Predicate Information (identified by operation id):

```
3 - filter((TO_CHAR(INTERNAL_FUNCTION("S"."TIME_ID"), 'YYYYMMDD') >= '20210101' AND  
TO_CHAR(INTERNAL_FUNCTION("S"."TIME_ID"), 'YYYYMMDD') <= '20220101'))
```

23 rows selected.



Ensuring Partition Pruning

Don't use functions on partition key filter predicates

```
SELECT avg( luminosity )
FROM atlas.EVENTS s, atlas.times t
WHERE s.time_id = t.time_id
AND TO_CHAR(s.time_id, 'YYYYMMDD') between '20210101' and '20220101'
```

```
SELECT avg( luminosity )
FROM atlas.EVENTS s, atlas.times t
WHERE s.time_id = t.time_id
AND s.time_id between TO_DATE('20210101', 'YYYYMMDD') and TO_DATE('20220101', 'YYYYMMDD')
```

Plan hash value: 2025449199

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				3 (100)			
1	SORT AGGREGATE		1	12				
2	PARTITION RANGE ITERATOR		313	3756	3 (0)	00:00:01	9	13
* 3	TABLE ACCESS FULL	EVENTS	313	3756	3 (0)	00:00:01	9	13

Predicate Information (identified by operation id):

```
3 - filter("S"."TIME_ID"<=TO_DATE(' 2020-01-01 00:00:00', 'syyy-mm-dd hh24:mi:ss'))
```

22 rows selected.

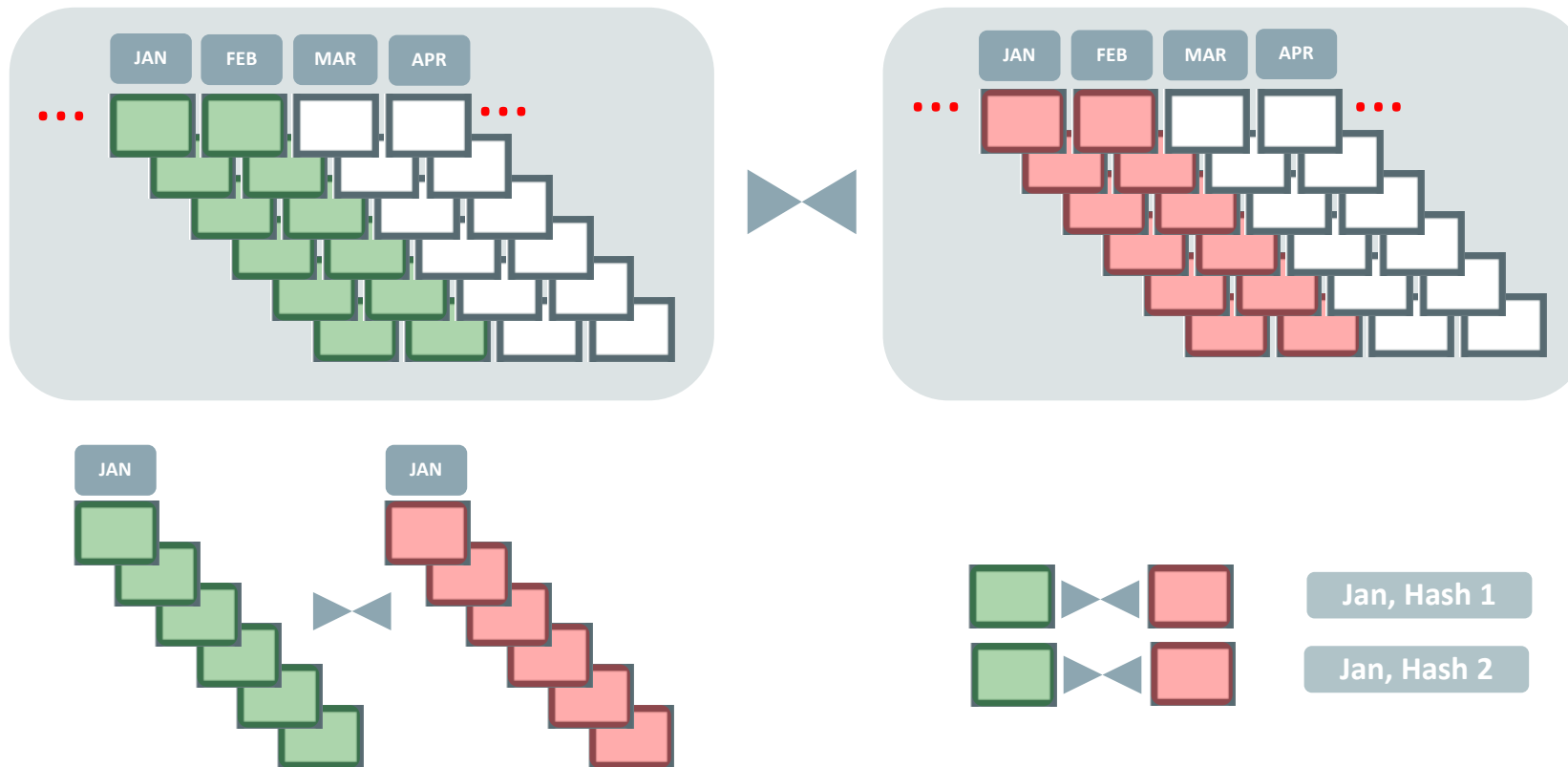
Pstart	Pstop
1	16
1	16

101' AND
01'))



Partition-wise Joins

Partition pruning and PWJ's "at work"

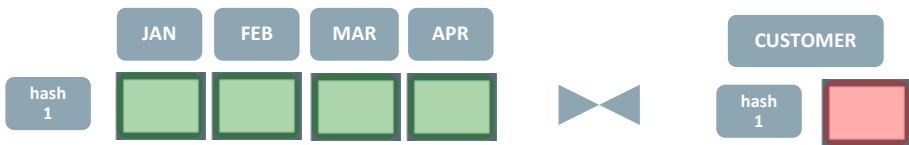
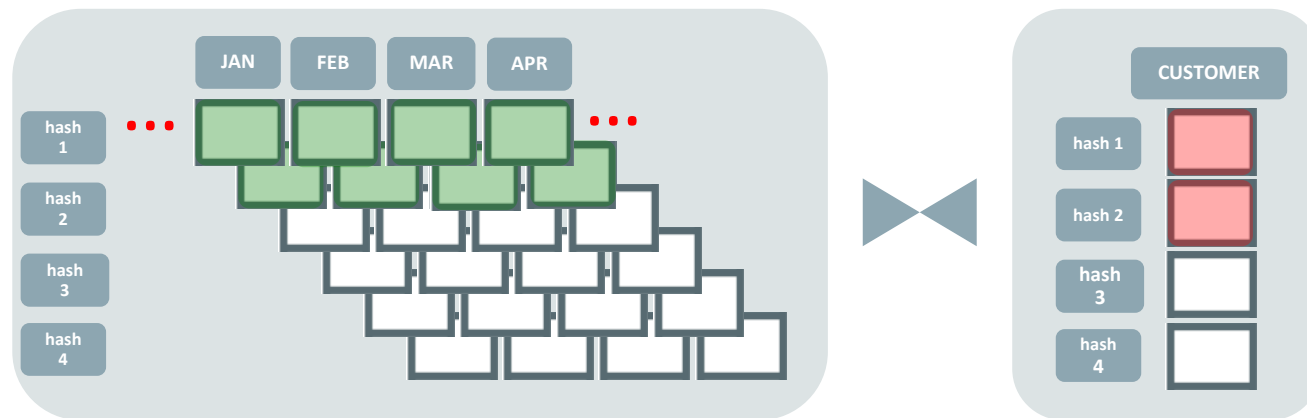


Large join is divided into multiple smaller joins, executed in parallel

- # of partitions to join must be a multiple of DOP
- Both tables must be partitioned the same way on the join column

Partition-wise Joins

Partition pruning and PWJ's "at work"



Large join is divided into multiple smaller joins, executed in parallel

- # of partitions to join must be a multiple of DOP
- Both tables must be partitioned the same way on the join column

Partition Purging and Loading

- Remove and add data as metadata only operations
 - Exchange the metadata of partitions
- Exchange standalone table w/ arbitrary single partition
 - Data load: standalone table contains new data to being loaded
 - Data purge: partition containing data is exchanged with empty table
- Drop partition alternative for purge
 - Data is gone forever



Partitioning Maintenance



Partition Maintenance

Fundamental Concepts for Success

- While performance seems to be the most visible one, don't forget about the rest, e.g.
 - Partitioning must address all business-relevant areas of Performance, Manageability, and Availability
- Partition autonomy is crucial
 - Fundamental requirement for any partition maintenance operations
 - Acknowledge partitions as metadata in the data dictionary



Partition Maintenance

Fundamental Concepts for Success

- Provide full partition autonomy
 - Use local indexes whenever possible
 - Enable partition all table-level operations for partitions, e.g. TRUNCATE, MOVE, COMPRESS
- Make partitions visible and usable for database administration
 - Partition naming for ease of use
- Maintenance operations must be partition-aware
 - Also true for indexes
- Maintenance operations must not interfere with online usage of a partitioned table



Aspects of Data Management

Addressable with Partition Maintenance Operations

- Fast population of data
 - EXCHANGE
 - Per-partition direct path load
- Fast removal of data
 - DROP, TRUNCATE, EXCHANGE
- Fast reorganization of data
 - MOVE, SPLIT, MERGE



Partition Maintenance

Table Partition Maintenance Operations

```
ALTER TABLE ADD PARTITION(S)
ALTER TABLE DROP PARTITION(S)
ALTER TABLE EXCHANGE PARTITION
ALTER TABLE MODIFY PARTITION
    [PARALLEL] [ONLINE]
ALTER TABLE MOVE PARTITION [PARALLEL] [ONLINE]
ALTER TABLE RENAME PARTITION
ALTER TABLE SPLIT PARTITION [PARALLEL] [ONLINE]
ALTER TABLE MERGE PARTITION(S) [PARALLEL]
ALTER TABLE COALESCE PARTITION [PARALLEL]
ALTER TABLE ANALYZE PARTITION
ALTER TABLE TRUNCATE PARTITION(S)
Export/Import [by partition]
Transportable tablespace [by partition]
```

Index Maintenance Operations

```
ALTER INDEX MODIFY PARTITION
ALTER INDEX DROP PARTITION(S)
ALTER INDEX REBUILD PARTITION
ALTER INDEX RENAME PARTITION
ALTER INDEX RENAME
ALTER INDEX SPLIT PARTITION
ALTER INDEX ANALYZE PARTITION
```

All partitions remain available all the time

- DML Lock on impacted partitions
- Move partition online no lock at all



Partition Maintenance on Multiple Partitions

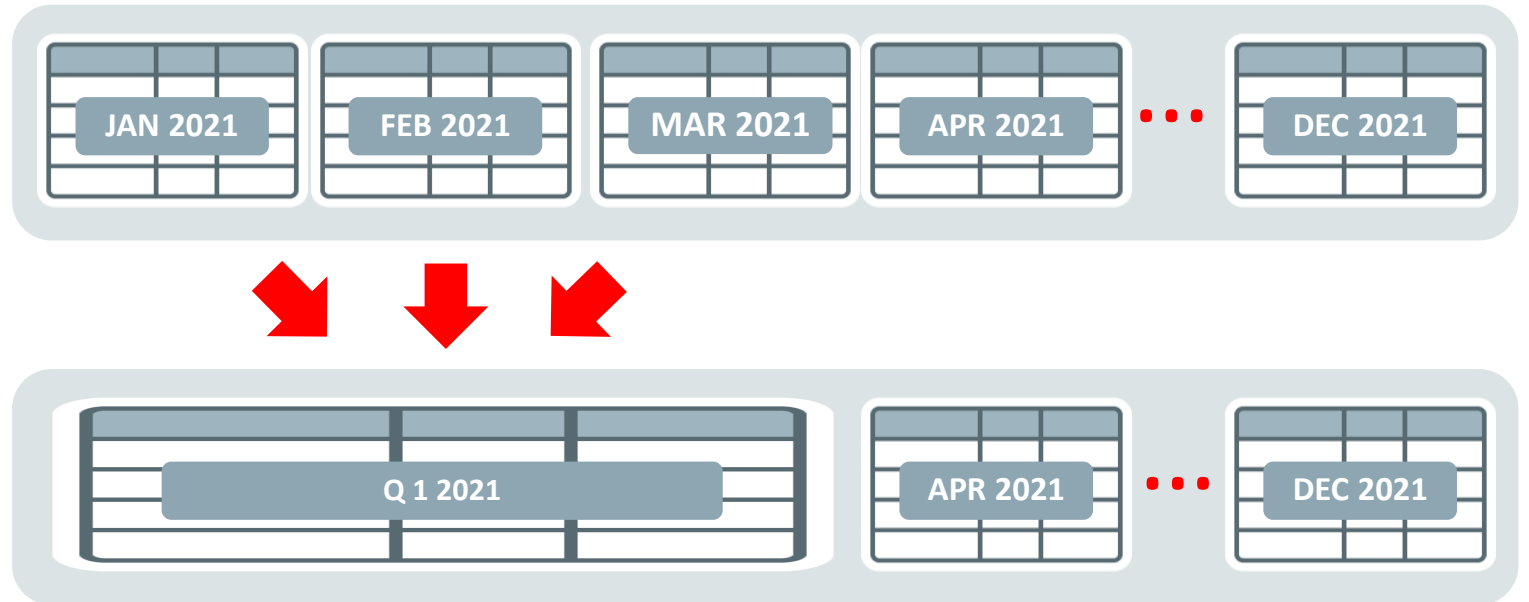
Introduced in Oracle 12c Release 1 (12.1)



Enhanced Partition Maintenance Operations

Operate on multiple partitions

- Partition Maintenance on multiple partitions in a single operation
- Full parallelism
- Transparent maintenance of local and global indexes



```
ALTER TABLE events
MERGE PARTITIONS Jan2021, Feb2021, Mar2021
INTO PARTITION Q1_2021 COMPRESS FOR ARCHIVE HIGH;
```


Enhanced Partition Maintenance Operations

Operate on multiple partitions

- Specify multiple partitions in order

```
SQL > alter table pt merge partitions part05, part15, part25  
      into partition p30;
```

Table altered.

- Specify a range of partitions

```
SQL > alter table pt merge partitions part10 to part30  
      into partition part30;
```

Table altered.

```
SQL > alter table pt split partition p30 into  
2  (partition p10 values less than (10),  
3  partition p20 values less than (20),  
4  partition p30);
```

Table altered.

- Works for all PMOPS
 - Supports optimizations like fast split



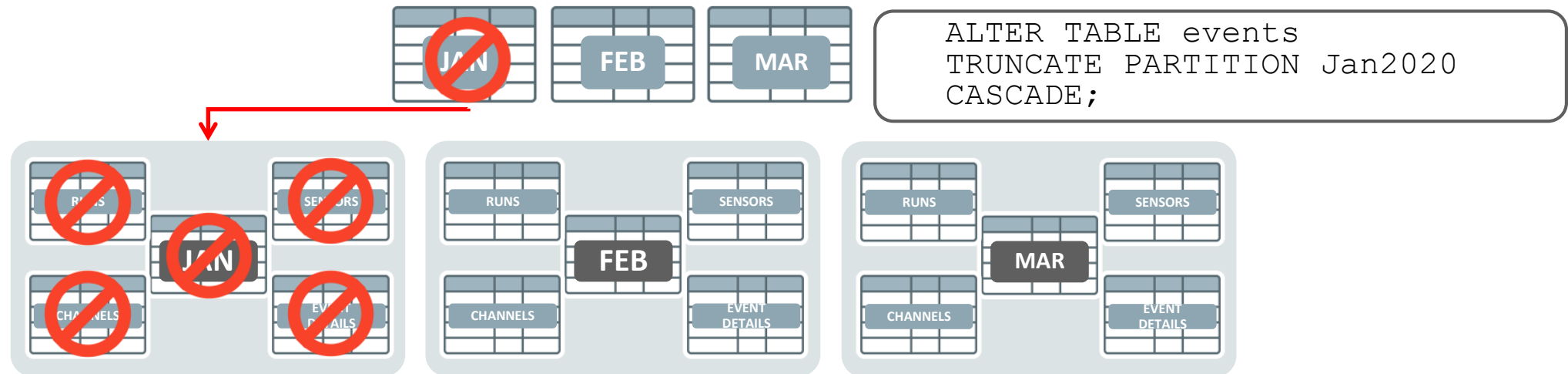
Cascading Truncate and Exchange for Reference Partitioning

Introduced in Oracle 12c Release 1 (12.1)



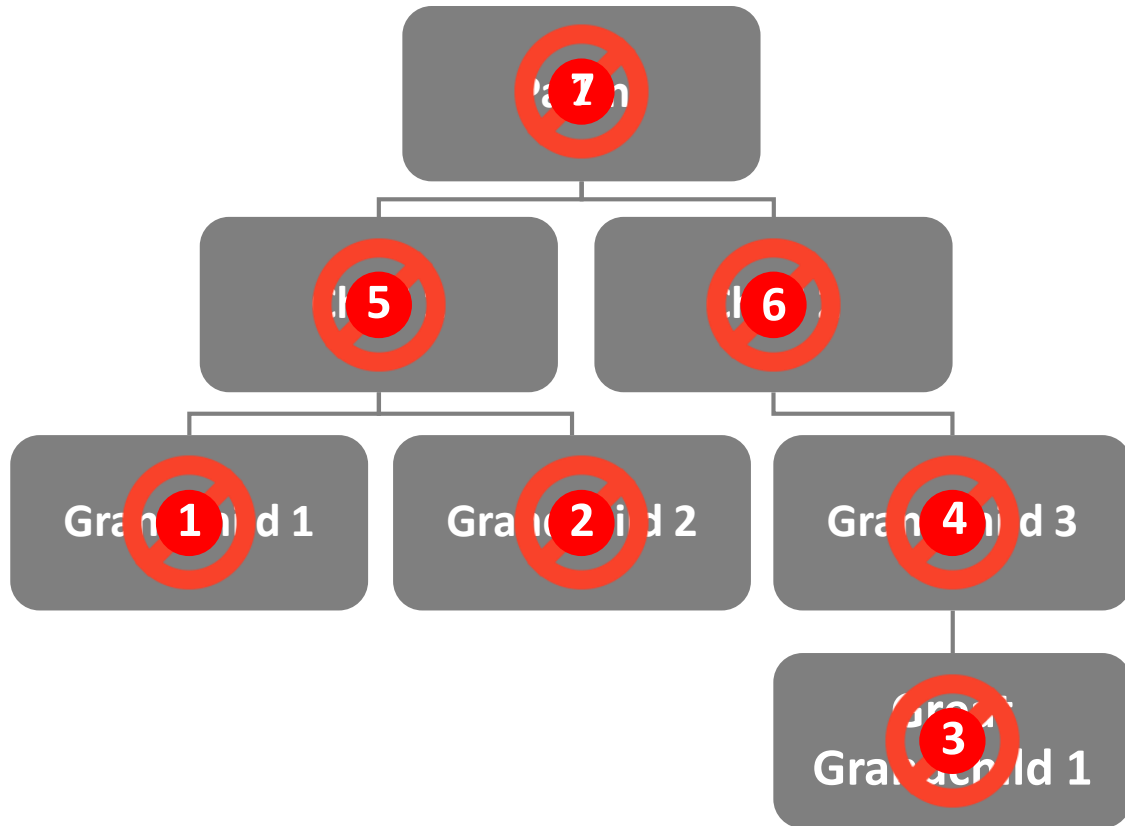
Advanced Partitioning Maintenance

Cascading TRUNCATE and EXCHANGE PARTITION

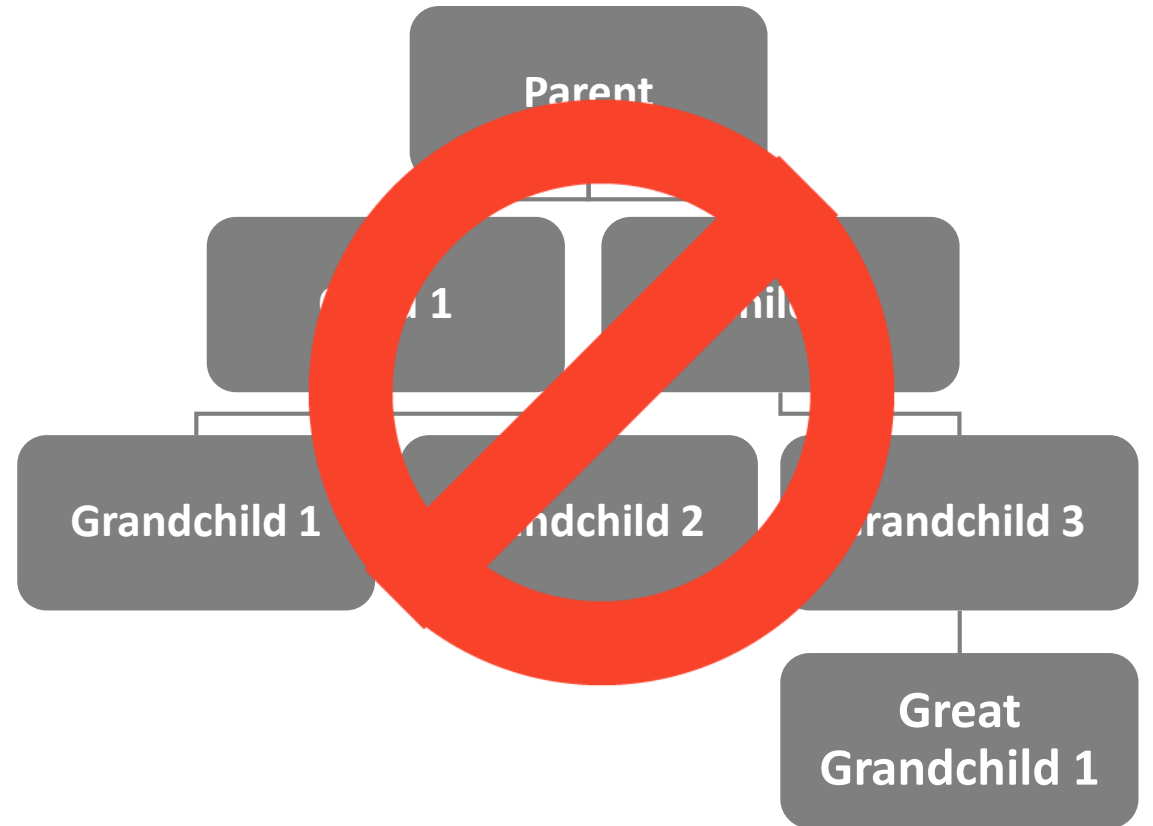


- Cascading TRUNCATE and EXCHANGE for improved business continuity
- Single atomic transaction preserves data integrity
- Simplified and less error prone code development

Cascading TRUNCATE PARTITION



- Proper bottom-up processing required
- Seven individual truncate operations



- One truncate operation

Cascading TRUNCATE PARTITION

```
SQL> create table intRef_p (pkcol number not null, col2 varchar2(200),
 2                          constraint pk_intref primary key (pkcol))
 3 partition by range (pkcol) interval (10)
 4 (partition p1 values less than (10));
```

Table created.

```
SQL>
SQL> create table intRef_c1 (pkcol number not null, col2 varchar2(200), fkcol number not null,
 2                          constraint pk_c1 primary key (pkcol),
 3                          constraint fk_c1 foreign key (fkcol) references intRef_p(pkcol) ON DELETE CASCADE)
 4 partition by reference (fk_c1);
```

Table created.



Cascading TRUNCATE PARTITION

```
SQL> create table intRef_p (pkcol number(2) constraint pkc1 primary key,
3 partition by range (pkcol) into (partition p1 values less than (1000),
4 (partition p2 values less than (2000)));
```

Table created.

```
SQL>
SQL> create table intRef_c1 (pkcol number(2) constraint pkc2 primary key,
3 constraint fkc1 foreign key (pkcol) references intRef_p (pkcol),
4 partition by reference (fk_c1) reference intRef_p partition p1);
```

Table created.

```
SQL> select * from intRef_p;
```

PKCOL	COL2
333	data for truncate - p
999	data for truncate - p

```
SQL> select * from intRef_c1;
```

PKCOL	COL2	FKCOL
1333	data for truncate - c1	333
1999	data for truncate - c1	999

```
SQL> alter table intRef_p truncate partition for (999) cascade update indexes;
```

Table truncated.

```
SQL> select * from intRef_p;
```

PKCOL	COL2
333	data for truncate - p

```
SQL> select * from intRef_c1;
```

PKCOL	COL2	FKCOL
1333	data for truncate - c1	333

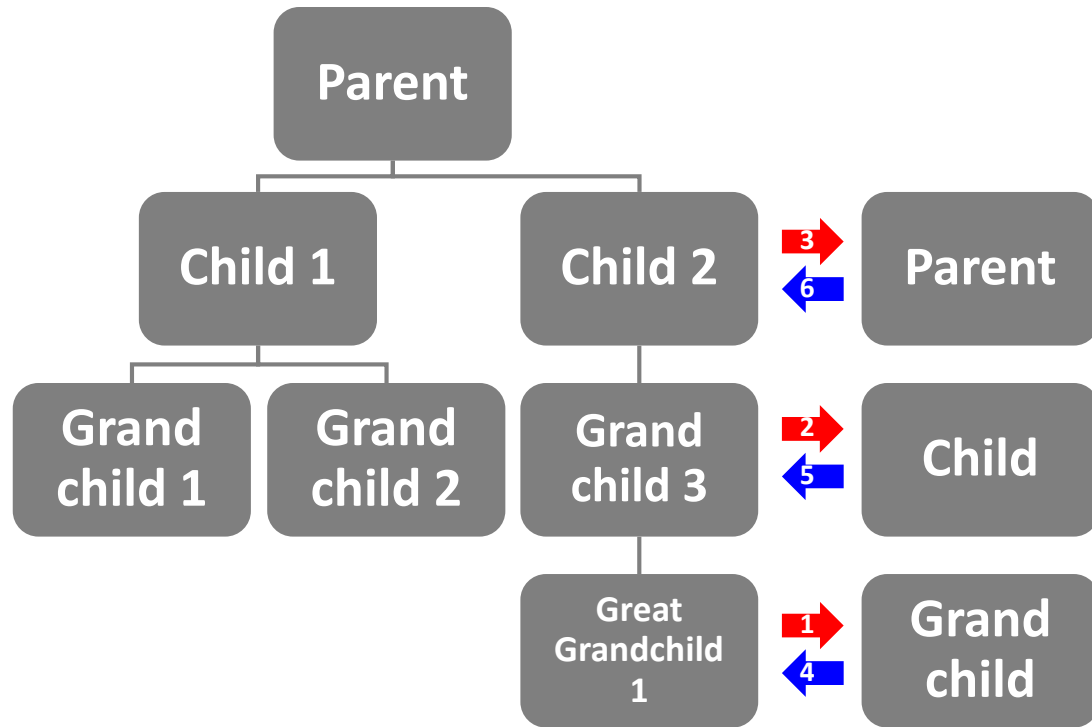


Cascading TRUNCATE PARTITION

- CASCADE applies for whole reference tree
 - Single atomic transaction, all or nothing
 - Bushy, deep, does not matter
 - Can be specified on any level of a reference-partitioned table
- ON DELETE CASCADE for all foreign keys required
- Cascading TRUNCATE available for non-partitioned tables as well
 - Dependency tree for non-partitioned tables can be interrupted with disabled foreign key constraints
- Reference-partitioned hierarchy must match for target and table to-be-exchanged
- For bushy trees with multiple children on the same level, each child on a given level must reference to a different key in the parent table
 - Required to unambiguously pair tables in the hierarchy tree

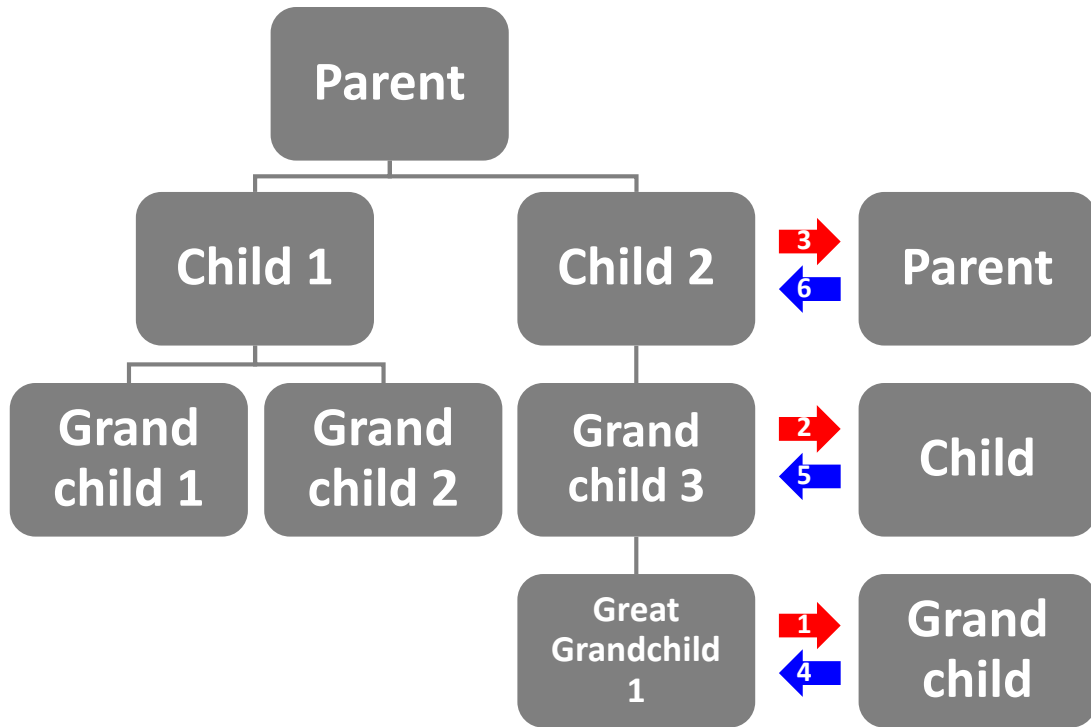


Cascading EXCHANGE PARTITION

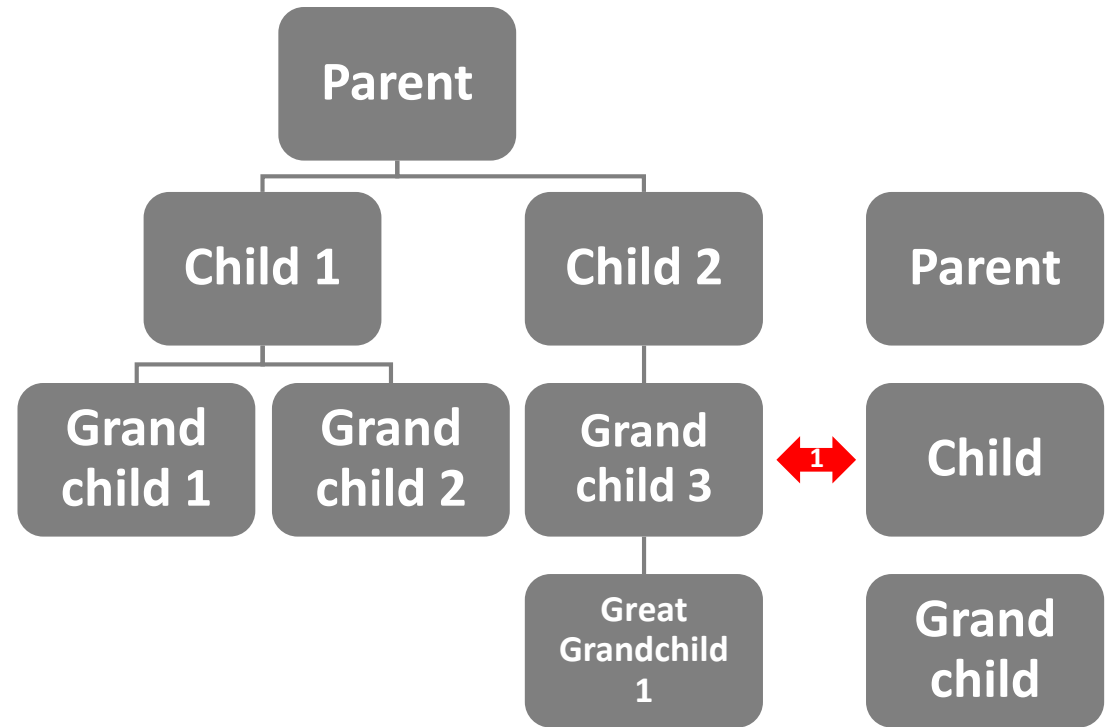


- Exchange (clear) out of target bottom-up
- Exchange (populate) into target top-down

Cascading EXCHANGE PARTITION



- Exchange (clear) out of target bottom-up
- Exchange (populate) into target top-down



- Exchange complete hierarchy tree
- One exchange operation

Cascading EXCHANGE PARTITION

```
SQL> create table intRef_p (pkcol number not null, col2 varchar2(200),
 2          constraint pk_intref primary key (pkcol))
 3 partition by range (pkcol) interval (10)
 4 (partition p1 values less than (10));

SQL> create table intRef_c1 (pkcol number not null, col2 varchar2(200), fkcol number not null,
 2          constraint pk_c1 primary key (pkcol),
 3          constraint fk_c1 foreign key (fkcol) references intRef_p(pkcol) ON DELETE CASCADE)
 4 partition by reference (fk_c1);

SQL> create table intRef_gc1 (col1 number not null, col2 varchar2(200), fkcol number not null,
 2          constraint fk_gc1 foreign key (fkcol) references intRef_c1(pkcol) ON DELETE CASCADE)
 3 partition by reference (fk_gc1);
```



Cascading EXCHANGE PARTITION

```
SQL> REM create some PK-FK equivalent table construct for exchange
SQL> create table XintRef_p (pkcol number not null, col2 varchar2(200),
 2                          constraint xpk_intref primary key (pkcol));

SQL> create table XintRef_c1 (pkcol number not null, col2 varchar2(200), fkcol number not null,
 2                          constraint xpk_c1 primary key (pkcol),
 3                          constraint xfk_c1 foreign key (fkcol) references XintRef_p(pkcol) ON DELETE CASCADE);

SQL> create table XintRef_gc1 (col1 number not null, col2 varchar2(200), fkcol number not null,
 2                          constraint xfk_gc1 foreign key (fkcol) references XintRef_c1(pkcol) ON DELETE CASCADE);
```



Cascading EXCHANGE PARTITION

```
SQL> select * from intRef_p;
```

PKCOL	COL2
333	p333 - data BEFORE exchange - p
999	p999 - data BEFORE exchange - p

```
SQL> select * from intRef_c1;
```

PKCOL	COL2	FKCOL
1333	p333 - data BEFORE exchange - c1	333
1999	p999 - data BEFORE exchange - c1	999

```
SQL> select * from intRef_gc1;
```

COL1	COL2	FKCOL
1333	p333 - data BEFORE exchange - gc1	1333
1999	p999 - data BEFORE exchange - gc1	1999

```
SQL> select * from XintRef_p;
```

PKCOL	COL2
333	p333 - data AFTER exchange - p

```
SQL> select * from XintRef_c1;
```

PKCOL	COL2	FKCOL
1333	p333 - data AFTER exchange - c1	333

```
SQL> select * from XintRef_gc1;
```

COL1	COL2	FKCOL
1333	p333 - data AFTER exchange - gc1	1333



Cascading EXCHANGE PARTITION

```
SQL> alter table intRef_p exchange partition for (333) with table XintRef_p cascade update indexes;  
Table altered.
```



Cascading EXCHANGE PARTITION

```
SQL> select * from intRef_p;
```

```
PKCOL COL2
```

```
-----  
333 p333 - data AFTER exchange - p  
999 p999 - data BEFORE exchange - p
```

```
SQL> select * from intRef_c1;
```

```
PKCOL COL2
```

```
FKCOL
```

```
-----  
1333 p333 - data AFTER exchange - c1 333  
1999 p999 - data BEFORE exchange - c1 999
```

```
SQL> select * from intRef_gc1;
```

```
COL1 COL2
```

```
FKCOL
```

```
-----  
1333 p333 - data AFTER exchange - gc1 1333  
1999 p999 - data BEFORE exchange - gc1 1999
```

```
SQL> select * from XintRef_p;
```

```
PKCOL COL2
```

```
-----  
333 p333 - data BEFORE exchange - p
```

```
SQL> select * from XintRef_c1;
```

```
PKCOL COL2
```

```
FKCOL
```

```
-----  
1333 p333 - data BEFORE exchange - c1 333
```

```
SQL> select * from XintRef_gc1;
```

```
COL1 COL2
```

```
FKCOL
```

```
-----  
1333 p333 - data BEFORE exchange - gc1 1333
```



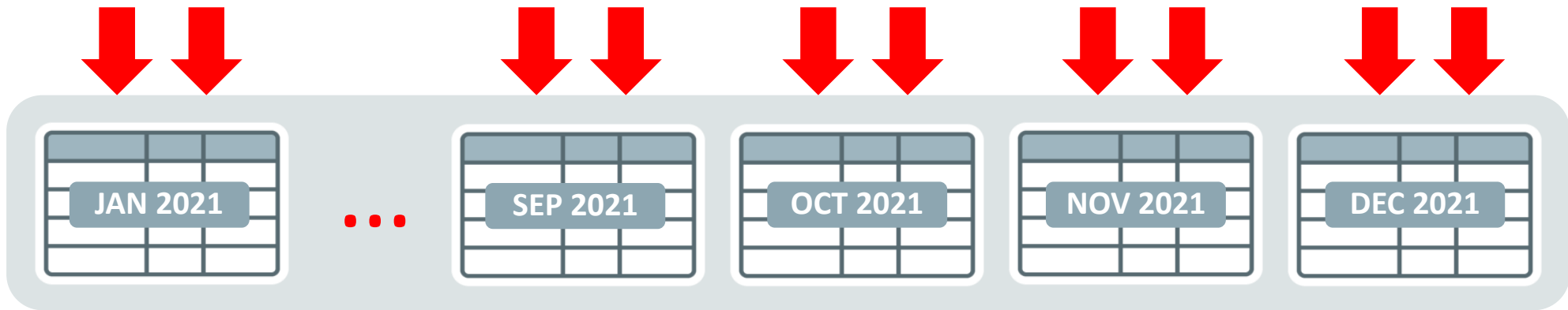
Online Move Partition

Introduced in Oracle 12c Release 1 (12.1)



Enhanced Partition Maintenance Operations

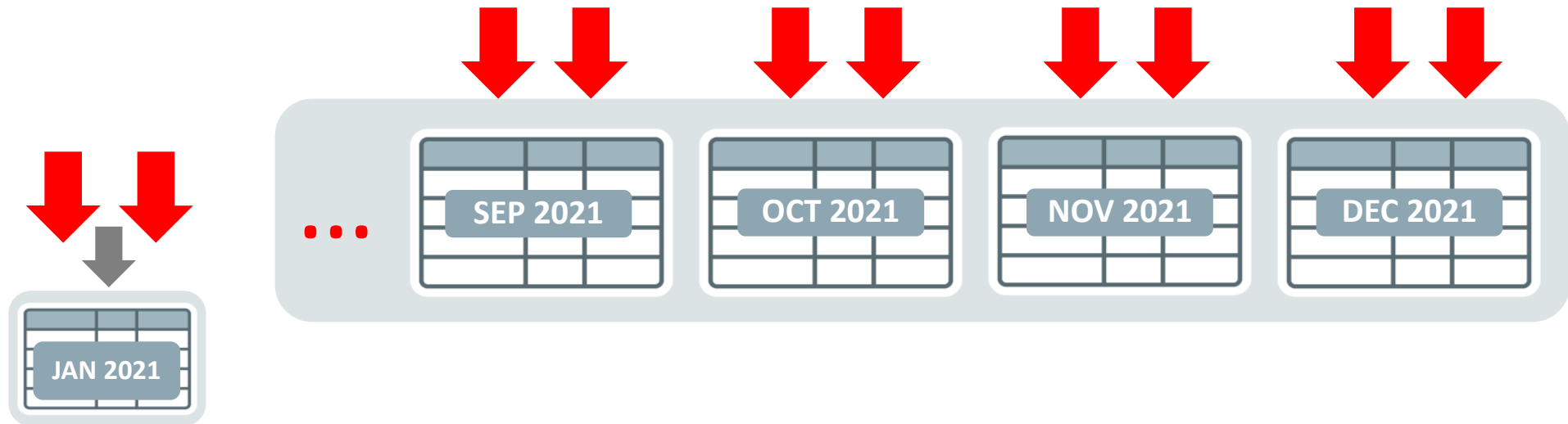
Online Partition Move



- Transparent MOVE PARTITION ONLINE operation
- Concurrent DML and Query
- Index maintenance for local and global indexes

Enhanced Partition Maintenance Operations

Online Partition Move



- Transparent MOVE PARTITION ONLINE operation
- Concurrent DML and Query
- Index maintenance for local and global indexes

Enhanced Partition Maintenance Operations

Online Partition Move – Best Practices

- Minimize concurrent DML operations if possible
 - Require additional disk space and resources for journaling
 - Journal will be applied recursively after initial bulk move
 - The larger the journal, the longer the runtime
- Concurrent DML has impact on compression efficiency
 - Best compression ratio with initial bulk move



Asynchronous Global Index Maintenance

Introduced in Oracle 12c Release 1 (12.1)



Enhanced Partition Maintenance Operations

Asynchronous Global Index Maintenance

- Usable global indexes after DROP and TRUNCATE PARTITION without index maintenance
 - Affected partitions are known internally and filtered out at data access time
- DROP and TRUNCATE become fast, metadata-only operations
 - Significant speedup and reduced initial resource consumption
- Delayed Global index maintenance
 - Deferred maintenance through ALTER INDEX REBUILD | COALESCE
 - Automatic cleanup using a scheduled job



Enhanced Partition Maintenance Operations

Asynchronous Global Index Maintenance

Before

```
SQL> select count(*) from pt partition for (9999);

COUNT(*)
-----
25341440

Elapsed: 00:00:01.00
SQL> select index_name, status, orphaned_entries from user_indexes;

INDEX_NAME                STATUS  ORPHANED_ENTRIES
-----
I1_PT                      VALID   NO

Elapsed: 00:00:01.04
SQL>
SQL> alter table pt drop partition for (9999) update indexes;

Table altered.
Elapsed: 00:02:04.52
SQL>
SQL> select index_name, status, orphaned_entries from user_indexes;

INDEX_NAME                STATUS  ORPHANED_ENTRIES
-----
I1_PT                      VALID   NO

Elapsed: 00:00:00.10
```

After

```
SQL> select count(*) from pt partition for (9999);

COUNT(*)
-----
25341440

Elapsed: 00:00:00.98
SQL> select index_name, status, orphaned_entries from user_indexes;

INDEX_NAME                STATUS  ORPHANED_ENTRIES
-----
I1_PT                      VALID   NO

Elapsed: 00:00:00.33
SQL>
SQL> alter table pt drop partition for (9999) update indexes;

Table altered.
Elapsed: 00:00:00.04
SQL>
SQL> select index_name, status, orphaned_entries from user_indexes;

INDEX_NAME                STATUS  ORPHANED_ENTRIES
-----
I1_PT                      VALID   YES

Elapsed: 00:00:00.05
```



Statistics Management for Partitioning



Statistics Gathering

- You must gather Optimizer statistics
 - Using dynamic sampling is not an adequate solution
 - Statistics on global and partition level recommended
 - Subpartition level optional
- Run all queries against empty tables to populate column usage
 - This helps identify which columns automatically get histograms created on them
- Optimizer statistics should be gathered after the data has been loaded but before any indexes are created
 - Oracle will automatically gather statistics for indexes as they are being created



Statistics Gathering

- By default DBMS_STATS gathers the following stats for each table
 - global (table level), partition level, sub-partition level
- Optimizer uses global stats if query touches two or more partitions
- Optimizer uses partition stats if queries do partition elimination and only one partition is necessary to answer the query
 - If queries touch two or more partitions the optimizer will use a combination of global and partition level statistics
- Optimizer uses sub-partition level statistics only if your queries do partition elimination and one sub-partition is necessary to answer query

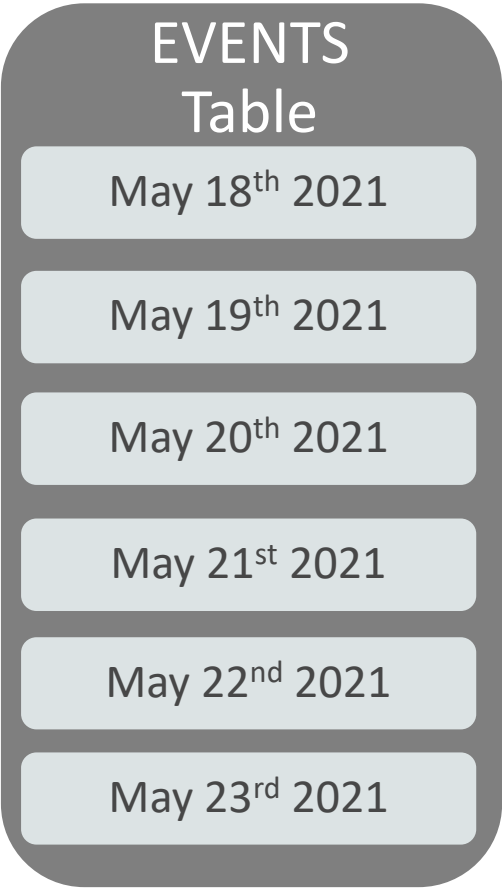


Efficient Statistics Management

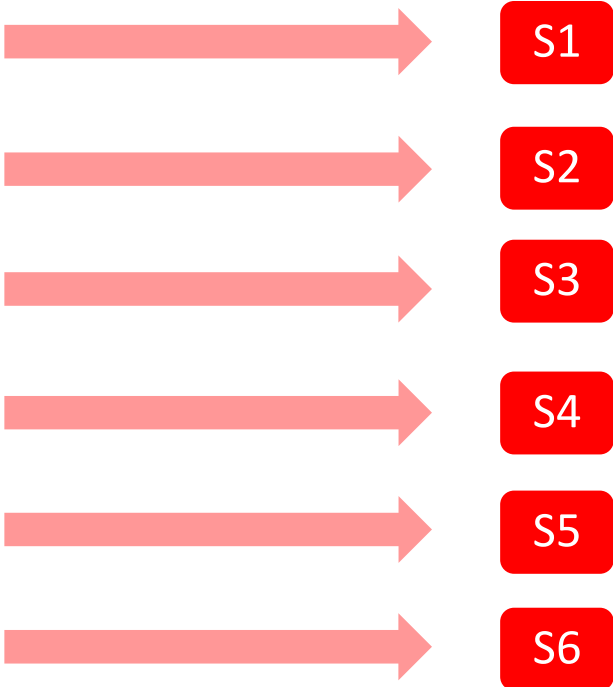
- Use `AUTO_SAMPLE_SIZE`
 - The only setting that enables new efficient statistics collection
 - Hash based algorithm, scanning the whole table
 - Speed of sampling, accuracy of compute
- Enable incremental global statistics collection
 - Avoids scan of all partitions after changing single partitions
 - Prior to 11.1, scan of all partitions necessary for global stats
 - Managed on per table level
 - Static setting
 - Create synopsis for non-partitioned table to being exchanged (Oracle Database 12c)



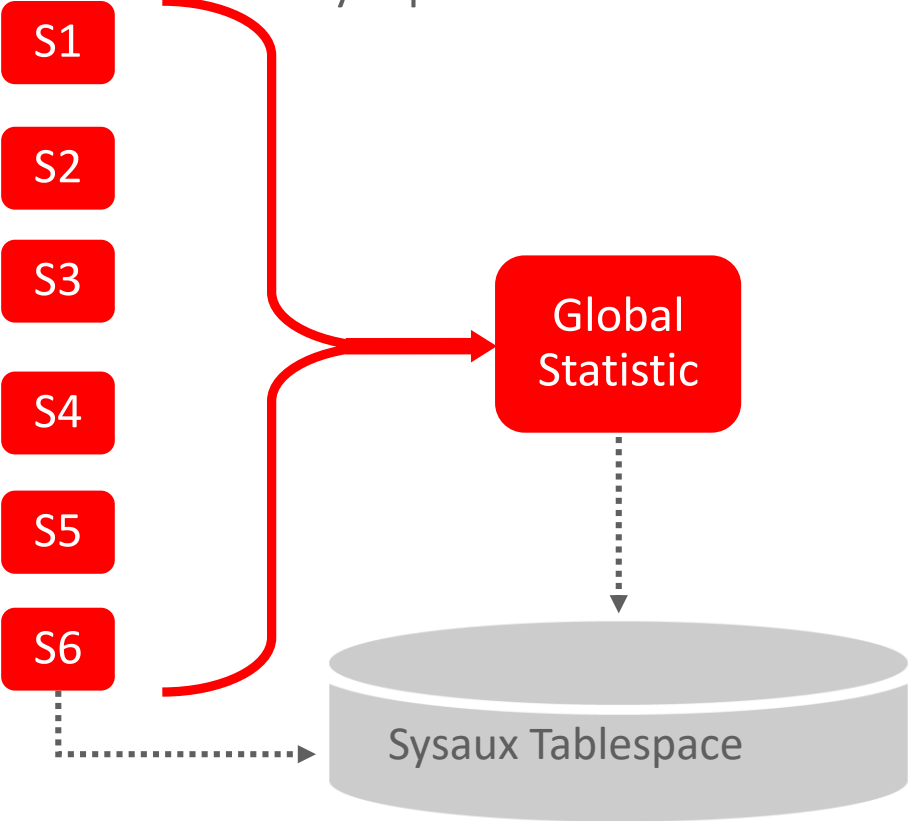
Incremental Global Statistics



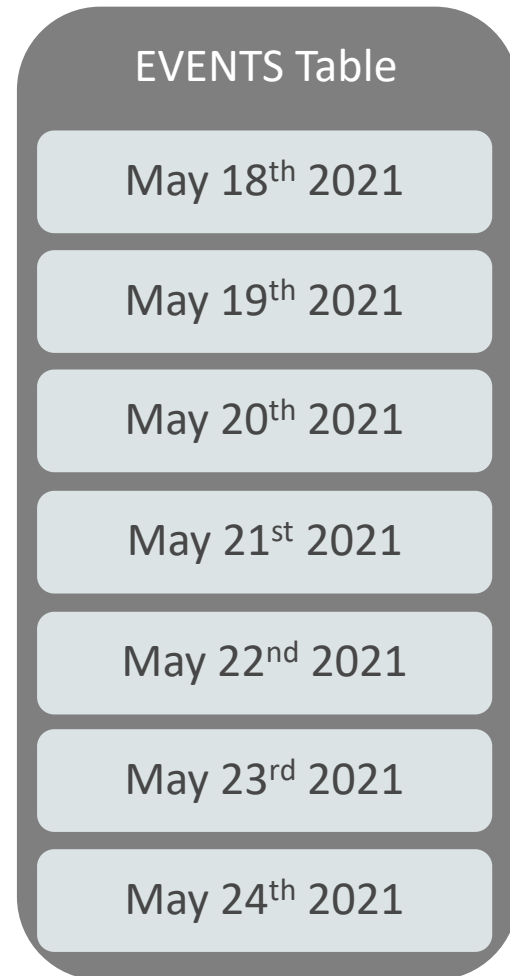
1. Partition level stats are gathered & synopsis created



2. Global stats generated by aggregating partition synopsis



Incremental Global Statistics Cont'd



3. A new partition is added to the table and data is loaded

4. Gather partition statistics for new partition



Incremental Global Statistics Cont'd

EVENTS Table

May 18 th 2021
May 19 th 2021
May 20 th 2021
May 21 st 2021
May 22 nd 2021
May 23 rd 2021
May 24 th 2021

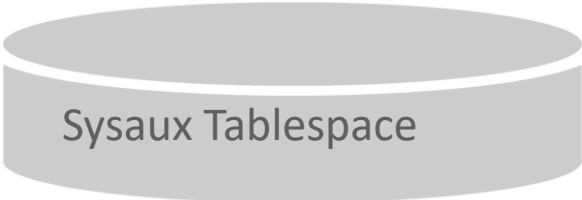
5. Retrieve synopsis for each of the other partitions from Sysaux



- S1
- S2
- S3
- S4
- S5
- S6
- S7

6. Global stats generated by aggregating the original partition synopsis with the new one

Global
Statistic



Step necessary to gather accurate statistics

- Turn on incremental feature for the table

```
EXEC DBMS_STATS.SET_TABLE_PREFS('ATLAS','EVENTS','INCREMENTAL','TRUE');
```

- After load gather table statistics using GATHER_TABLE_STATS

- No need to specify parameters

```
EXEC DBMS_STATS.GATHER_TABLE_STATS('ATLAS','EVENTS');
```

- The command will collect statistics for partitions and update the global statistics based on the partition level statistics and synopsis
- Possible to set incremental to true for all tables
 - Only works for already existing tables

```
EXEC DBMS_STATS.SET_GLOBAL_PREFS('INCREMENTAL','TRUE');
```



Attribute Clustering and Zone Maps

Introduced in Oracle 12c Release 1 (12.1.0.2)

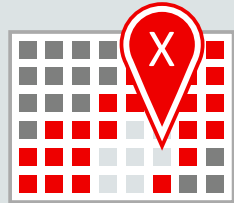


Zone Maps with Attribute Clustering



Attribute Clustering

Orders data so that columns values are stored together on disk



Zone maps

Stores min/max of specified columns per zone

Used to filter un-needed data during query execution

- Combined Benefits
- Improved query performance and concurrency
 - Reduced physical data access
 - Significant IO reduction for highly selective operations
- Optimized space utilization
 - Less need for indexes
 - Improved compression ratios through data clustering
- Full application transparency
 - Any application will benefit



Attribute Clustering

Concepts and Benefits

- Orders data so that it is in close proximity based on selected columns values: “attributes”
- Attributes can be from a single table or multiple tables
 - e.g. from fact and dimension tables
- Significant IO pruning when used with zone maps
- Reduced block IO for table lookups in index range scans
- Queries that sort and aggregate can benefit from pre-ordered data
- Enable improved compression ratios
 - Ordered data is likely to compress more than unordered data



Attribute Clustering for Zone Maps

Ordered rows

```
ALTER TABLE EVENTS  
ADD CLUSTERING BY  
LINER ORDER (category);  
  
ALTER TABLE EVENTS  
MOVE;
```

Category	Country
BOYS	AR
BOYS	JP
BOYS	SA
BOYS	US
GIRLS	AR
GIRLS	JP
GIRLS	SA
GIRLS	US
MEN	AR
MEN	JP
MEN	SA
MEN	US
WOMEN	AR
WOMEN	JP
WOMEN	SA
WOMEN	US

- Ordered rows containing category values BOYS, GIRLS and MEN.
- *Zone maps* catalogue regions of rows, or *zones*, that contain particular column value ranges.
- By default, each zone is up to 1024 blocks.
- For example, we only need to scan this zone if we are searching for category “GIRLS”. We can skip all other zones.

Attribute Clustering

Basics

- Two types of attribute clustering
 - LINEAR ORDER BY
 - Classical ordering
 - INTERLEAVED ORDER BY
 - Multi-dimensional ordering
- Simple attribute clustering on a single table
- Join attribute clustering
 - Cluster on attributes derived through join of multiple tables
 - Up to four tables
 - Non-duplicating join (PK or UK on joined table is required)



Attribute Clustering

Example

- CLUSTERING BY LINEAR ORDER (category, country)
- CLUSTERING BY INTERLEAVED ORDER (category, country)

Category	Country
BOYS	AR
BOYS	JP
BOYS	SA
BOYS	US
GIRLS	AR
GIRLS	JP
GIRLS	SA
GIRLS	US
MEN	AR
MEN	JP
MEN	SA
MEN	US
WOMEN	AR
WOMEN	JP
WOMEN	SA
WOMEN	US

- LINEAR ORDER

		Country						
Category	10	AR WOMEN	11	JP WOMEN	14	SA WOMEN	15	US WOMEN
	8	AR MEN	9	JP MEN	12	SA MEN	13	US MEN
	2	AR GIRLS	3	JP GIRLS	6	SA GIRLS	7	US GIRLS
	0	AR BOYS	1	JP BOYS	4	SA BOYS	5	US BOYS

- INTERLEAVED ORDER



Attribute Clustering

Basics

- Clustering directive specified at table level
 - ALTER TABLE ... ADD CLUSTERING ...
- Directive applies to new data and data movement
- Direct path operations
 - INSERT APPEND, MOVE, SPLIT, MERGE
 - Does not apply to conventional DML
- Can be enabled and disabled on demand
 - Hints and/or specific syntax



Zone Maps

Concepts and Basics

- Stores minimum and maximum of specified columns
 - Information stored per zone
 - [Sub]Partition-level rollup information for partitioned tables for multi-dimensional partition pruning
- Analogous to a coarse index structure
 - Much more compact than an index
 - Zone maps filter out what you don't need, indexes find what you do need
- Significant performance benefits with complete application transparency
 - IO reduction for table scans with predicates on the table itself or even a joined table using join zone maps (a.k.a. "hierarchical zone map")
- Benefits are most significant with ordered data
 - Used in combination with attribute clustering or data that is naturally ordered



Zone Maps

Basics

- Independent access structure built for a table
 - Implemented using a type of materialized view
 - For partitioned and non-partitioned tables
- One zone map per table
 - Zone map on partitioned table includes aggregate entry per [sub]partition
- Used transparently
 - No need to change or hint queries
- Implicit or explicit creation and column selection
 - Through Attribute Clustering: CREATE TABLE ... CLUSTERING
 - CREATE MATERIALIZED ZONEMAP ... AS SELECT ...



Attribute Clustering With Zone Maps

- CLUSTERING BY LINEAR ORDER (category, country)
- Zone map benefits are most significant with ordered data

Category	Country
BOYS	AR
BOYS	JP
BOYS	SA
BOYS	US
GIRLS	AR
GIRLS	JP
GIRLS	SA
GIRLS	US
MEN	AR
MEN	JP
MEN	SA
MEN	US
WOMEN	AR
WOMEN	JP
WOMEN	SA
WOMEN	US

Pruning with:

```
SELECT ..  
FROM table  
WHERE category =  
  'BOYS';
```

```
SELECT ..  
FROM table  
WHERE category =  
  'BOYS';  
AND country = 'US';
```

- LINEAR ORDER



Attribute Clustering With Zone Maps

- CLUSTERING BY INTERLEAVED ORDER (category, country)
- Zone map benefits are most significant with ordered data

		Country			
		10	11	14	15
Category	10	AR WOMEN	JP WOMEN	SA WOMEN	US WOMEN
	8	AR MEN	JP MEN	12 SA MEN	13 US MEN
	2	AR GIRLS	3 JP GIRLS	6 SA GIRLS	7 US GIRLS
	0	AR BOYS	1 JP BOYS	4 SA BOYS	5 US BOYS

- INTERLEAVED ORDER

Pruning with:

```
SELECT ..  
FROM table  
WHERE category =  
  'BOYS';
```

```
SELECT ..  
FROM table  
AND country = 'US';
```

```
SELECT ..  
FROM table  
WHERE category =  
  'BOYS'  
AND country = 'US';
```



Zone Maps

Staleness

- DML and partition operations can cause zone maps to become fully or partially stale
 - Direct path insert does not make zone maps stale
- Single table 'local' zone maps
 - Update and insert marks impacted zones as stale (and any aggregated partition entry)
 - No impact on zone maps for delete
- Joined zone map
 - DML on fact table equivalent behavior to single table zone map
 - DML on dimension table makes dependent zone maps fully stale



Zone Maps

Refresh

- Incremental and full refresh, as required by DML
 - Zone map refresh does require a materialized view log
 - Only stale zones are scanned to refresh the MV
 - For joined zone map
 - DML on fact table: incremental refresh
 - DML on dimension table: full refresh
- Zone map maintenance through
 - DBMS_MVIEW.REFRESH()
 - ALTER MATERIALIZED ZONEMAP <xx> REBUILD;



Example – Dimension Hierarchies

events

id	sensor_id	device_id	reading
1	3	23	2.00
2	88	55	43.75
3	31	99	33.55
4	33	62	23.12
5	21	11	38.00
6	33	21	5.00
7	44	71	10.99

Note: a zone typically contains many more rows than show here.
This is for illustrative purposes only.

devices

device_id	Exp_zone	Exp_type
23	Atlas	Cryo
102	CMS	Therm
55	Atlas	Charge
1	Area51	X-field
62	Atlas	z-search

```
CREATE TABLE orders ( ... )  
CLUSTERING orders  
JOIN locations ON (orders.device_id = locations.device_id)  
BY INTERLEAVED ORDER (locations.state, locations.county)  
WITH MATERIALIZED ZONEMAP ...
```



Example – Dimension Hierarchies

events

id	sensor_id	device_id	amount
1	3	23	2.00
2	88	55	43.75
3	31	99	33.55
4	33	62	23.12
5	21	11	38.00
6	33	21	5.00
7	44	71	10.99

Scan
Zone



devices

device_id	Exp_zone	Exp_type
23	Atlas	Cryo
102	CMS	Therm
55	Atlas	Charge
1	Area51	X-field
62	Atlas	z-search

Note: a zone typically contains many more rows than show here.
This is for illustrative purposes only.

```
SELECT SUM(amount)
FROM orders
JOIN locations ON (orders.location.id = locations.location.id)
WHERE exp_ome = 'Atlas';
```



Example – Dimension Hierarchies

events

id	sensor_id	device_id	amount
1	3	23	2.00
2	88	55	43.75
3	31	99	33.55
4	33	62	23.12
5	21	11	38.00
6	33	21	5.00
7	44	71	10.99

Scan
Zone



devices

device_id	Exp_zone	Exp_type
23	Atlas	Cryo
102	CMS	Therm
55	Atlas	Charge
1	Area51	X-field
62	Atlas	z-search

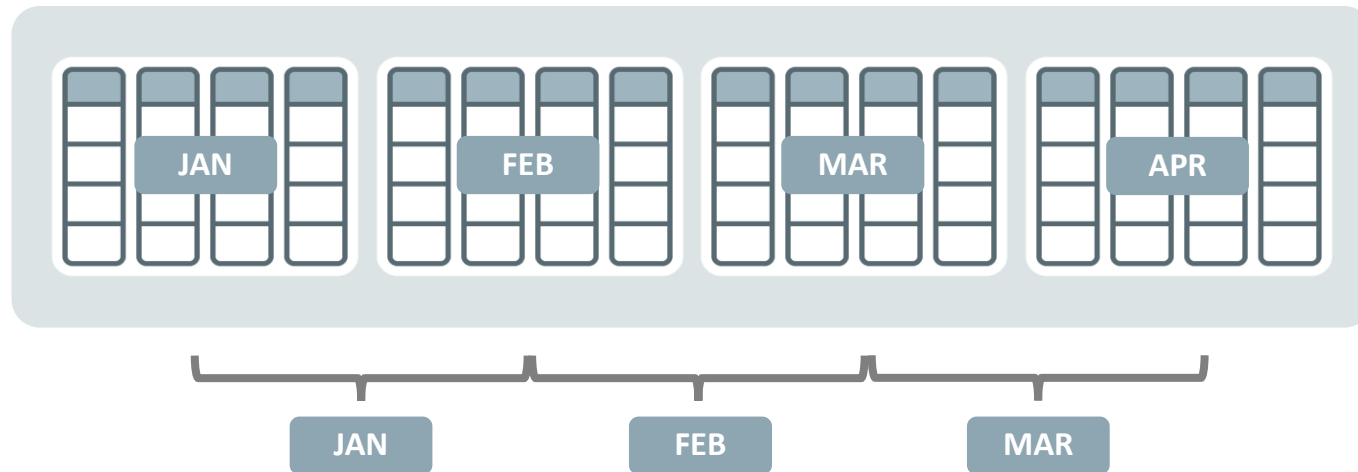
Note: a zone typically contains many more rows than show here.
This is for illustrative purposes only.

```
SELECT SUM(amount)
FROM orders
JOIN locations ON (orders.location.id = locations.location.id)
WHERE state = 'California'
AND county = 'Charge'
```



Zone Maps and Partitioning

Partition Key:
EVENT_DATE



Zone map column
EVENT_DATE
correlates with
partition key
RUN_DATE

Zone map:
RUN_DATE

- Zone maps can prune partitions for columns that are not included in the partition (or sub-partition) key

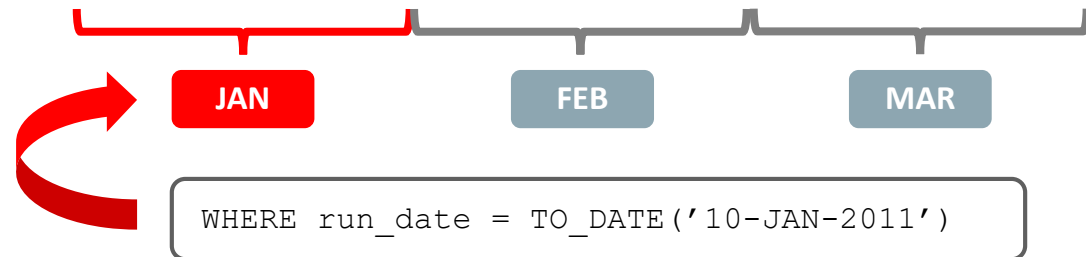
Zone Maps and Partitioning

Partition Key:
EVENT_DATE



MAR and APR partitions
are pruned

Zone map:
RUN_DATE



- Zone maps can prune partitions for columns that are not included in the partition (or sub-partition) key

Zone Maps and Storage Indexes

- Attribute clustering and zone maps work transparently with Exadata storage indexes
 - The benefits of Exadata storage indexes continue to be fully exploited
- In addition, zone maps (when used with attribute clustering)
 - Enable additional and significant IO optimization
 - Provide an alternative to indexes, especially on large tables
 - Join and fact-dimension queries, including dimension hierarchy searches
 - Particularly relevant in star and snowflake schemas
 - Are able to prune entire partitions and sub-partitions
 - Are effective for both direct and conventional path reads
 - Include optimizations for joins and index range scans
 - Part of the physical database design: explicitly created and controlled by the DBA



Partitioning tips and tricks



Partitioning Tips and Tricks

- Think about partitioning strategy
- Physical and logical attributes
- Eliminate hot spots
- Smart partial exchange
- Exchange with PK and UK
- Partition Elimination Table (PET)



Think about your partitioning strategy



Choosing your Partitioning Strategy

- Think about
 - your data
 - your usage
- What do you expect from Partitioning?
 - Query performance benefits
 - Load (or purge) performance benefits
 - Data management benefits



Choosing your Partitioning Strategy

Logical shape of the data

- How is data inserted into your system?
- How is data maintained in your system?
- How is data accessed in your system?



Choosing your Partitioning Strategy

Logical shape of the data

- How is data inserted into your system?
 - Time, location, tenant, business user, ...
 - Ranges, unrelated list of values, “just lots of them”, ...
- How is data maintained in your system?
- How is data accessed in your system?



Choosing your Partitioning Strategy

Logical shape of the data

- How is data inserted into your system?
 - Time, location, tenant, business user, ...
 - Ranges, unrelated list of values, “just lots of them”, ...
- How is data maintained in your system?
 - Moving window of active data, legal requirements, data “forever”, ...
 - Don't know yet
- How is data accessed in your system?



Choosing your Partitioning Strategy

Logical shape of the data

- How is data inserted into your system?
 - Time, location, tenant, business user, ...
 - Ranges, unrelated list of values, “just lots of them”, ...
- How is data maintained in your system?
 - Moving window of active data, legal requirements, data “forever”, ...
 - Don't know yet
- How is data accessed in your system?
 - Always full, with common FILTER predicates, always index access, ...
 - Don't know yet



Choosing your Partitioning Strategy

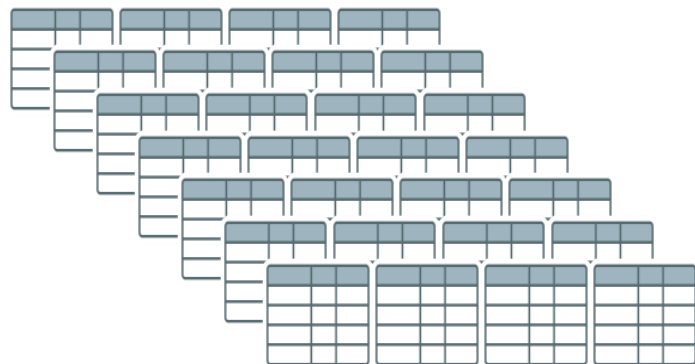
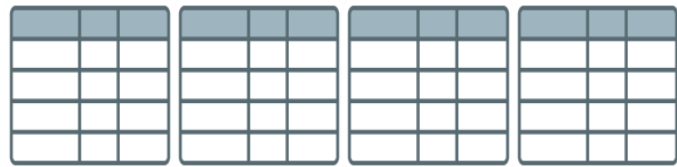
Performance improvements

- Query speedup
 - Partition elimination
 - Partition-wise joins
- DML speedup
 - Alleviation of contention points
- Data maintenance
 - DDL instead of DML



Choosing your Partitioning Strategy

Data Access – Full Table Access

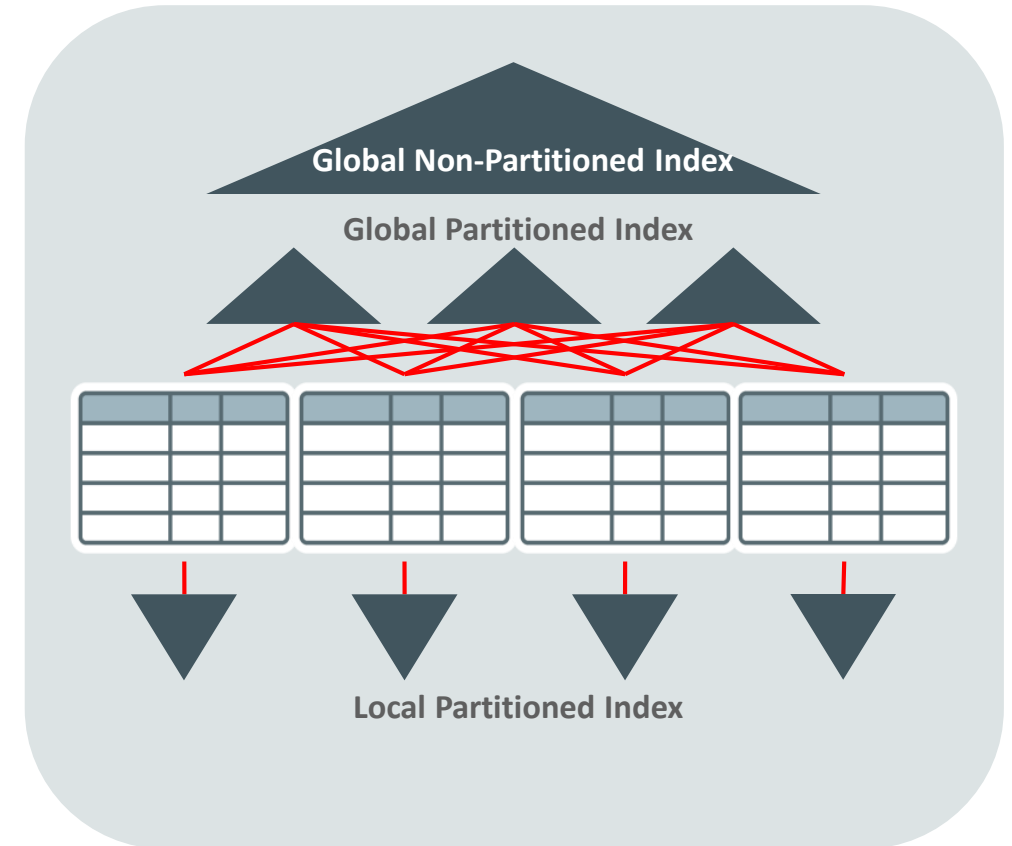


- I/O savings are linear to the number of pruned partitions
 - One of 10: ten times less IO
 - One of 100: hundred times less IO
- Runtime improvements depend on
 - Relative contribution of IO versus CPU work
 - Potential impact on subsequent operations

Choosing your Partitioning Strategy

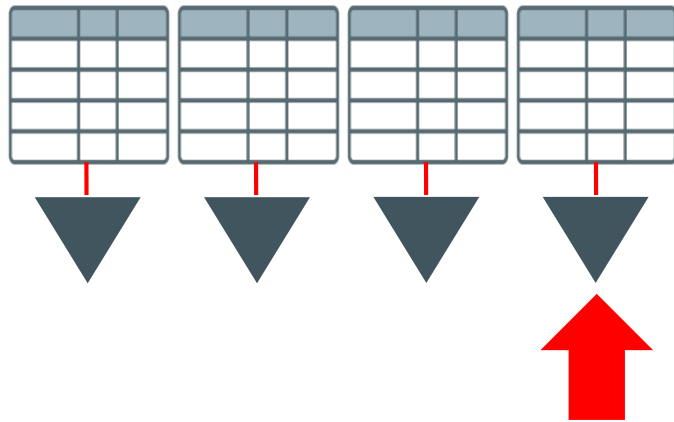
Indexing of partitioned tables

- GLOBAL index points to rows in any partition
 - Index can be partitioned or not
- LOCAL index is partitioned same as table
 - Index partitioning key can be different from index key

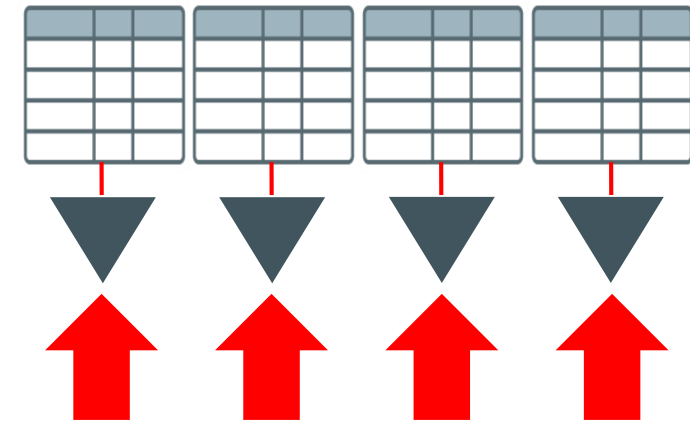


Choosing your Partitioning Strategy

Data Access – local index and global partitioned index



- Partitioned index access with single partition pruning



- Partitioned index access without any partition pruning

Local and Global Partitioned Indexes

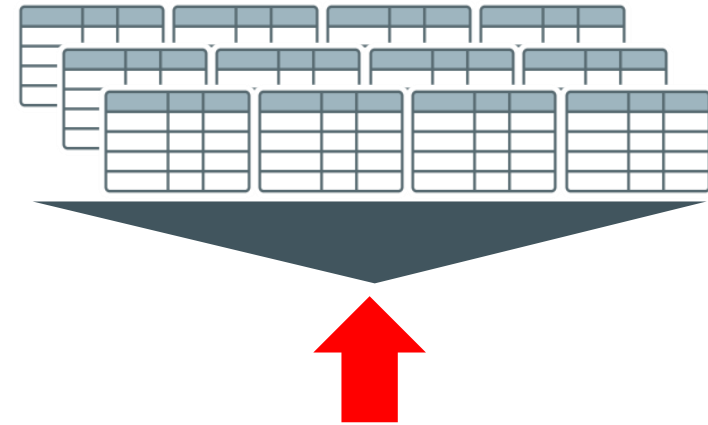
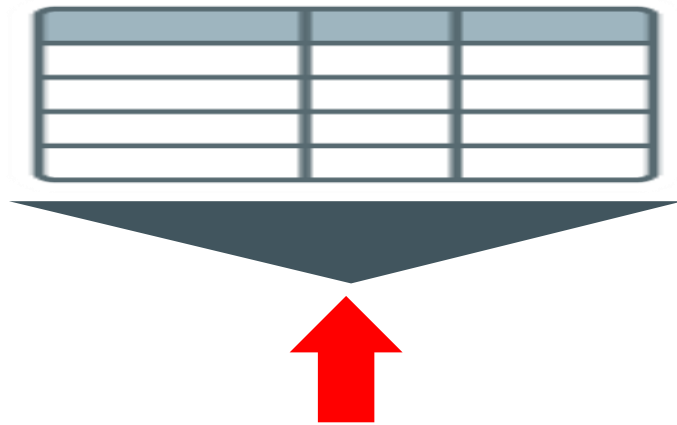
Data Access

- Number of index probes identical to number of accessed partitions
 - No partition pruning leads to a probe into all index partitions
- Not optimally suited for OLTP environments
 - No guarantee to always have partition pruning
 - Exception: global hash partitioned indexes for DML contention alleviation
 - Most commonly small number of partitions
- Pruning on global partitioned indexes based on the index prefix
 - Index prefix identical to leading keys of index



Choosing your Partitioning Strategy

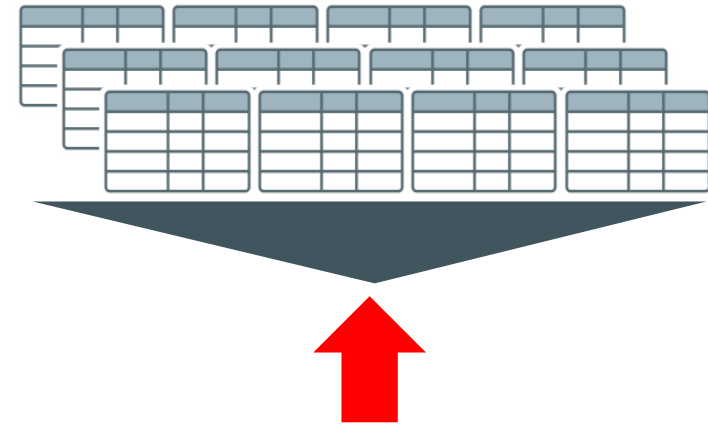
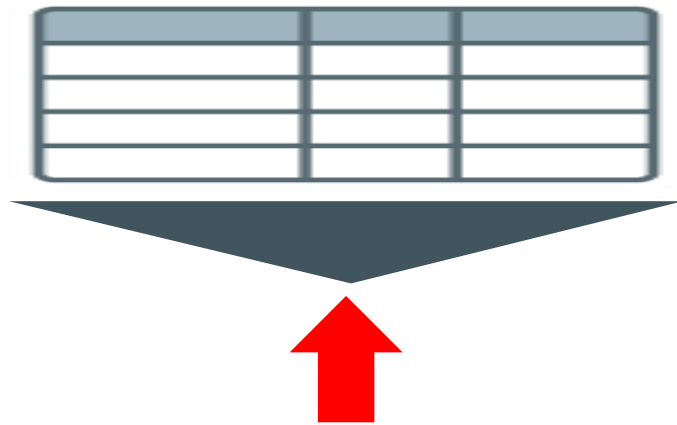
Global nonpartitioned index



- Can you see the difference?

Choosing your Partitioning Strategy

Global nonpartitioned index



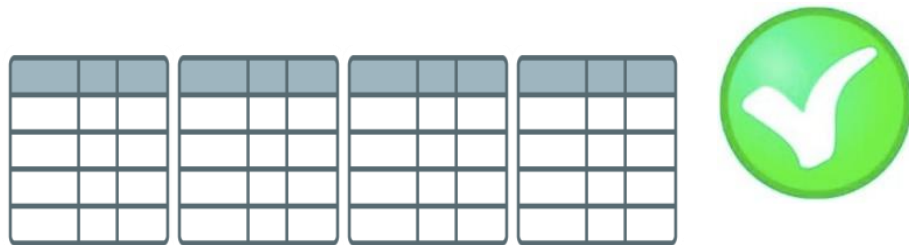
- Can you see the difference?
- There is more or less none*

* Some differences for index size, due to large rowid

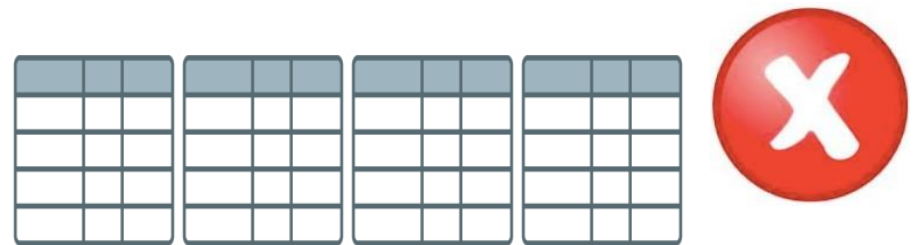
Global Indexes

Data Access

- No pruning for non-partitioned indexes
 - You always probe into a single index segment
- Global partitioned index prefix identical to leading keys of index
 - Pruning on index prefix, not partition key column(s)
- Most common in OLTP environments



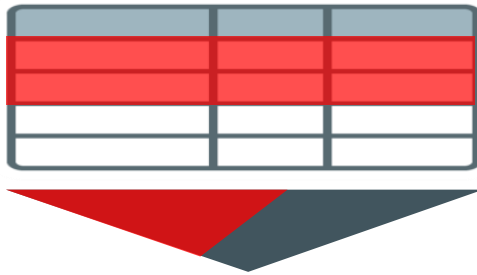
```
PARTITION BY (col1), idx(col1)
```



```
PARTITION BY (col1), idx(col2)
```

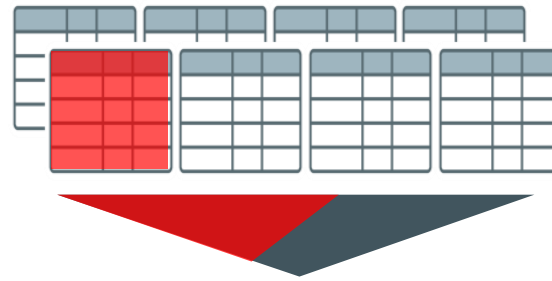

Choosing your Partitioning Strategy

Data Maintenance



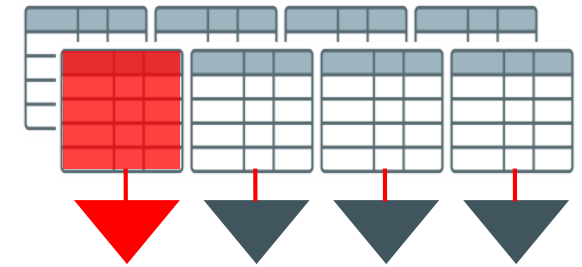
```
DELETE FROM ...  
WHERE ...
```

- Records get deleted
 - Index maintenance
 - Undo and redo



```
ALTER TABLE ... DROP PARTITION ...
```

- Partition gets dropped
 - Fast global index maintenance (12c)
 - Minimal undo



- Partition gets dropped
 - Local index gets dropped
 - Minimal undo

Local Indexes

Data Maintenance

- Incremental index creation possible
 - Initial unusable creation, rebuild of individual partitions
- Fast index maintenance for all partition maintenance operations that only touch one partition
 - Exchange, drop, truncate
- Partition maintenance that touches more than one partition require index maintenance
 - Merge, split creates new data segments
 - New index segments are created as well



Global Indexes

Data Maintenance

- Incremental index creation is hard, if not impossible
- “Fast” index maintenance for drop and truncate beginning with Oracle Database 12c
 - Fast actually means delayed index maintenance
- Partition maintenance except drop and truncate requires index maintenance
 - Conventional index maintenance equivalent to the DML operations that would represent the PMOP



How many partitions?

How many partitions?

It depends

Data Volume and Number of Partitions

- Imagine a 100TB table ...
 - With one million partitions, each partition is 100MB in size
- Imagine a 10TB table ...
 - With one million partitions, each partition is 10MB in size
- Imagine a 1TB table ...
 - With one million partitions, each partition is 1MB in size



Data Volume and Number of Partitions

- Imagine a 100TB table ...
 - With one million partitions, each partition is 100MB in size
- Imagine a 10TB table ...
 - With one million partitions, each partition is 10MB in size
- Imagine a 1TB table ...
 - With one million partitions, each partition is 1MB in size
- **How long does it take your system to read 1MB??**
 - Exadata full table scan rate is 25GB/sec ... (disk only, full rack X5-2)



Data Volume and Number of Partitions

- More is not always better
 - Every partition represents metadata in the dictionary
 - Every partition increases the metadata footprint in the SGA
- Find your personal balance between the number of partitions and its average size
 - There is nothing wrong about single-digit GB sizes for a segment on “normal systems”
 - Consider more partitions \geq 5GB segment size



Choosing your Partitioning Strategy

Customer Usage Patterns

- Range (Interval) still the most prevalent partitioning strategy
 - Almost always some time dependency
- List more and more common
 - Interestingly often based on time as well
 - Often as subpartitioning strategy
- Hash not only used for performance (PWJ, DML contention)
 - No control over data placement, but some understanding of it
 - Do not forget the power of two rule



Choosing your Partitioning Strategy

Extended Partitioning Strategies

- Interval Partitioning fastest growing new partitioning strategy
 - Manageability extension to Range Partitioning
- Reference Partitioning
 - Leverage PK/FK constraints for your data model
- Interval-Reference Partitioning (new in Oracle Database 12c)
- Virtual column based Partitioning
 - Derived attributes without little to no application change
- Any variant of the above



Physical and logical attributes

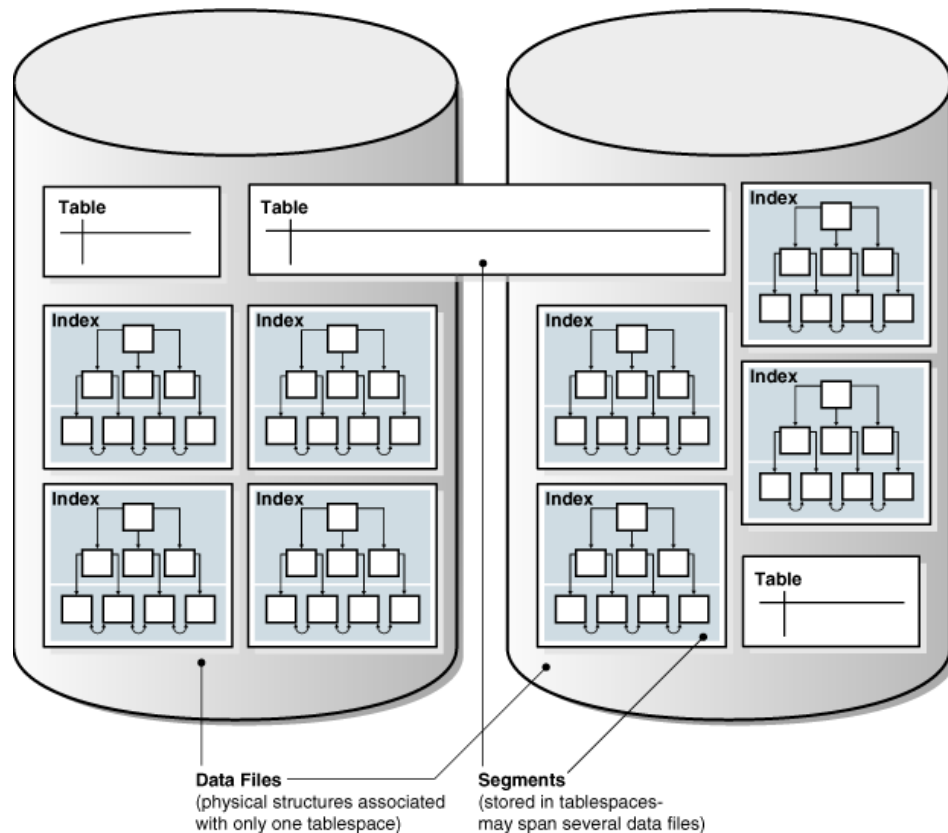
Physical and Logical Attributes

- Logical attributes
 - Partitioning setup
 - Indexing and index maintenance
 - Read only (in conjunction with tablespace separation)
- Physical attributes
 - Data placement
 - Segment properties in general



Nonpartitioned Tables

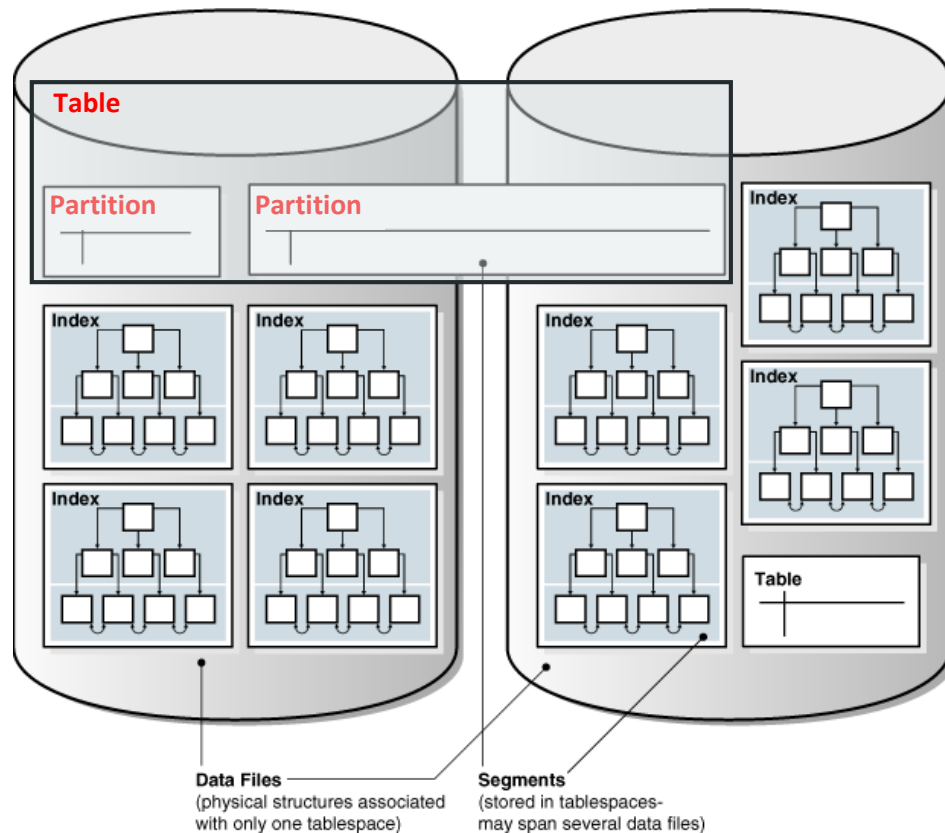
Physical and Logical Attributes



- Logical table properties
 - Columns and data types
 - Constraints
 - Indexes, ...
- Physical table properties
 - Table equivalent to segment
 - Tablespace
 - Compression, [Logging | nologging], ...
 - In-memory
 - Properties managed and changed on segment level

Partitioned Tables

Physical and Logical Attributes



- Logical **table** properties
 - Columns and data types
 - Constraints
 - **Partial** Indexes, ...
 - **Physical property directives**
- Physical **[sub]partition** properties
 - **[Sub]partition** equivalent to segment
 - Tablespace
 - Compression, [Logging | nologging], ...
 - In-memory
 - Properties managed and changed on segment level

Partitioned Tables

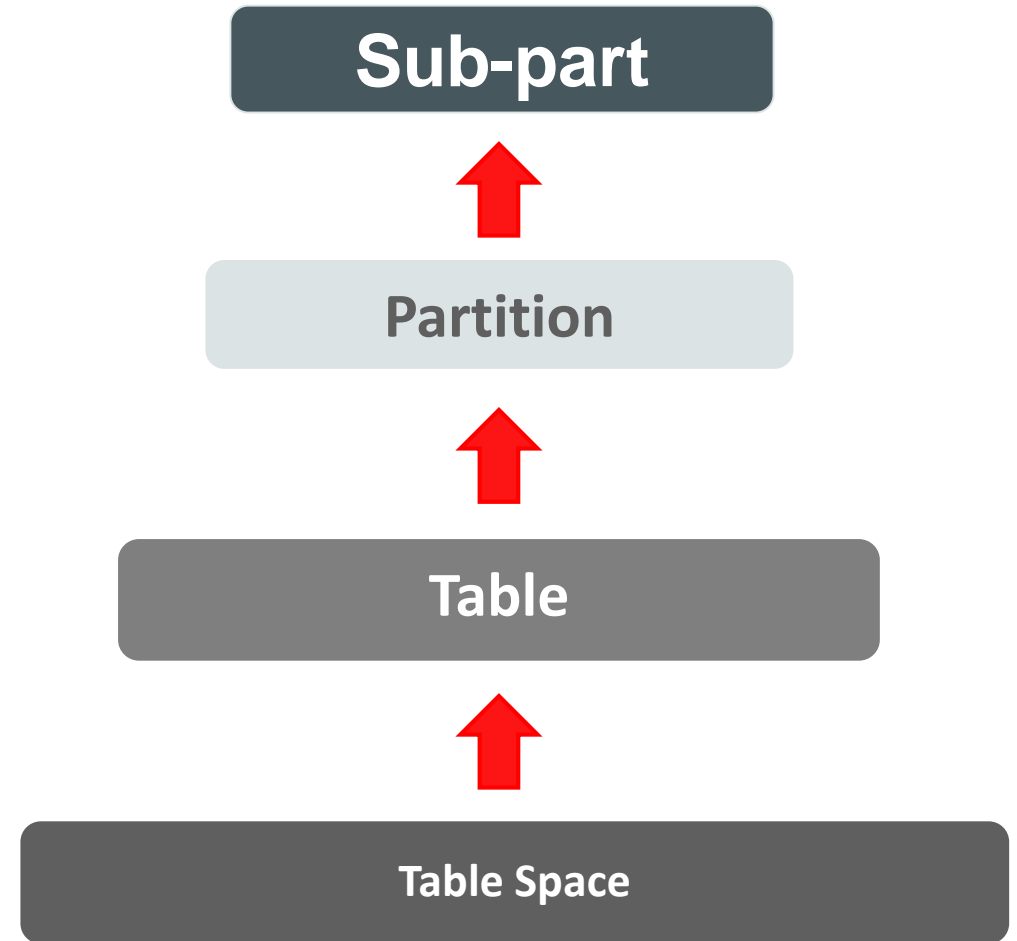
Physical and Logical Attributes

- Table is metadata-only and directive for future partitions
 - No physical segments on table level
 - Physical attributes become directive for new partitions, if specified
- Single-level partitioned table
 - Partitions are equivalent to segments
 - Physical attributes are managed and changed on partition level
- Composite-level partitioned tables
 - Partitions are metadata only and directive for future subpartitions
 - Subpartitions are equivalent to segments



Data Placement with Partitioned Tables

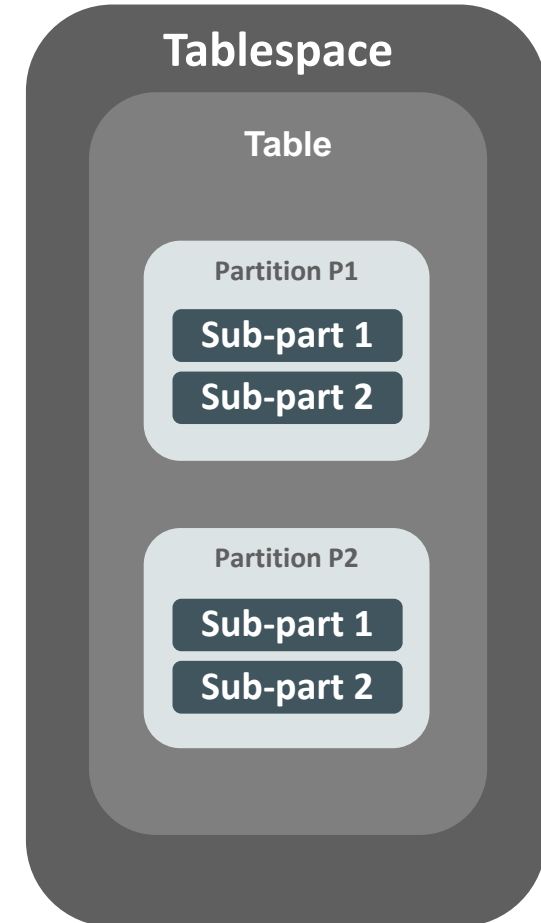
- Each partition or sub-partition is a separate object
- Specify storage attributes at each individual level
 - As placement policy for lower levels
 - For each individual [sub]partition
- If storage attributes are not specified standard hierarchical inheritance kicks in



Data Placement with Partitioned Tables

Special Case Interval Partitioning

- Interval Partitioning” pre-creates” all partitions
 - All 1 million [sub]partitions exist logically
- Physical storage is (almost) determined as well
- Partition placement
 - Inherited from table level
 - STORE IN () clause for round-robin partition placement
- Subpartition placement
 - Usage of subpartition template
 - Needs bug fix #8304261 (included in 11.2.0.3)
 - STORE IN clause currently is currently a no-op



Data Placement with Partitioned Tables

Subpartition template

- Introduced in Oracle Database 9 Release 2
 - Allows predefinition of subpartitions for **future** partitions
 - Stored as metadata in the data dictionary

```
CREATE TABLE stripe_regional_EVENTS
  (deptno number, item_no varchar2(20),
  txn_date date, txn_amount number, state varchar2(2))
PARTITION BY RANGE (txn_date)
SUBPARTITION BY LIST (state)
SUBPARTITION TEMPLATE
  (SUBPARTITION northwest VALUES ('OR', 'WA') TABLESPACE tbs_1
  SUBPARTITION southwest VALUES ('AZ', 'UT', 'NM') TABLESPACE tbs_2
  SUBPARTITION northeast VALUES ('NY', 'VM', 'NJ') TABLESPACE tbs_3
  SUBPARTITION southeast VALUES ('FL', 'GA') TABLESPACE tbs_4
  SUBPARTITION midwest VALUES ('SD', 'WI') TABLESPACE tbs_5
  SUBPARTITION south VALUES ('AL', 'AK') TABLESPACE tbs_6
  SUBPARTITION south VALUES (DEFAULT) TABLESPACE tbs_7
  )
(PARTITION q1_2021 VALUES LESS THAN ( TO_DATE('01-APR-2021', 'DD-MON-YYYY')),
(PARTITION q2_2021 VALUES LESS THAN ( TO_DATE('01-JUL-2021', 'DD-MON-YYYY')),
(PARTITION q3_2021 VALUES LESS THAN ( TO_DATE('01-OCT-2021', 'DD-MON-YYYY')),
(PARTITION q4_2021 VALUES LESS THAN ( TO_DATE('01-JAN-2020', 'DD-MON-YYYY')),
);
```

Subpartition
definition for all
future partitions

Subpartition applied
to every partition

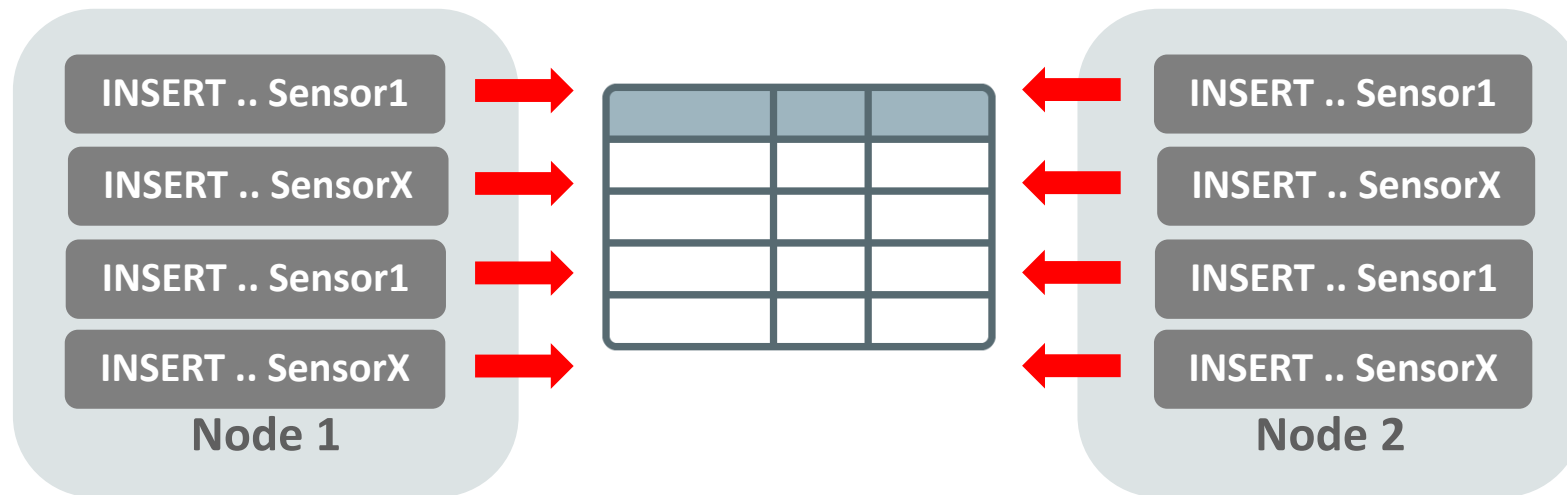


Using partitioning to eliminate hot spots

Using Partitioning to eliminate Hot Spots

Nonpartitioned table

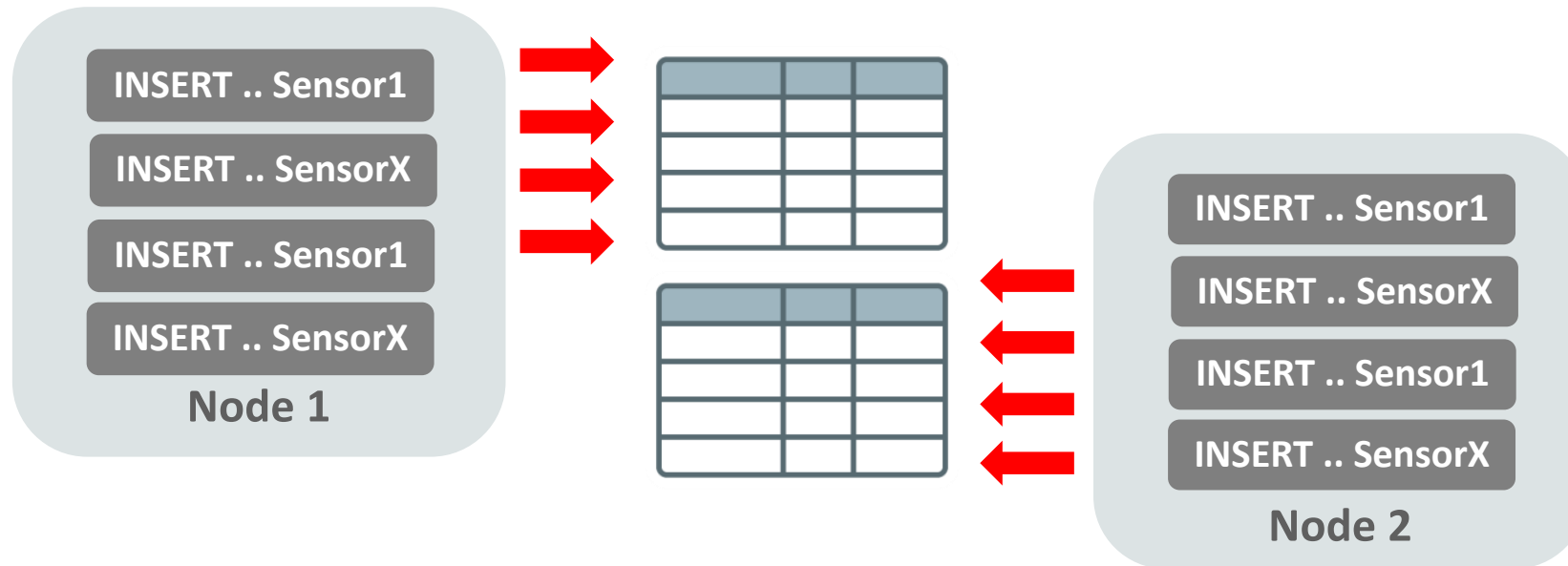
- On RAC, high DML workload causes high cache fusion traffic
 - Oracle calls this block pinging



Using Partitioning to eliminate Hot Spots

HASH partitioned table

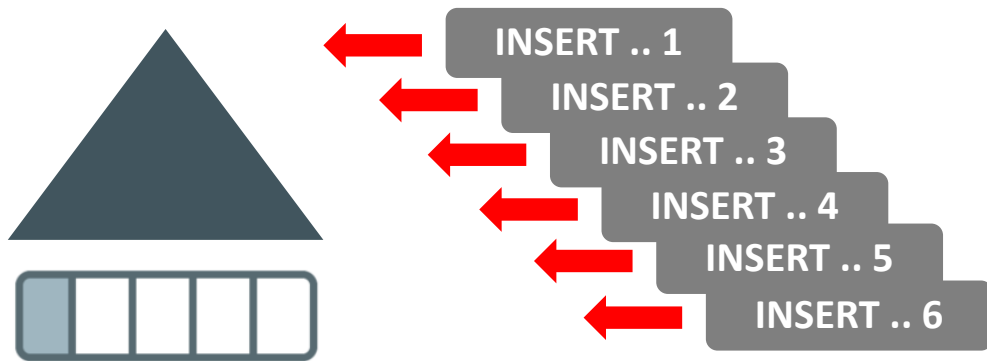
- On RAC, high DML workload causes high cache fusion traffic
 - Oracle calls this block pinging
- HASH (or LIST) partitioned table can alleviate this situation
 - Caveat: Normally needs some kind of “application partitioning” or “application RAC awareness”



Using Partitioning to eliminate Hot Spots

HASH partitioned index

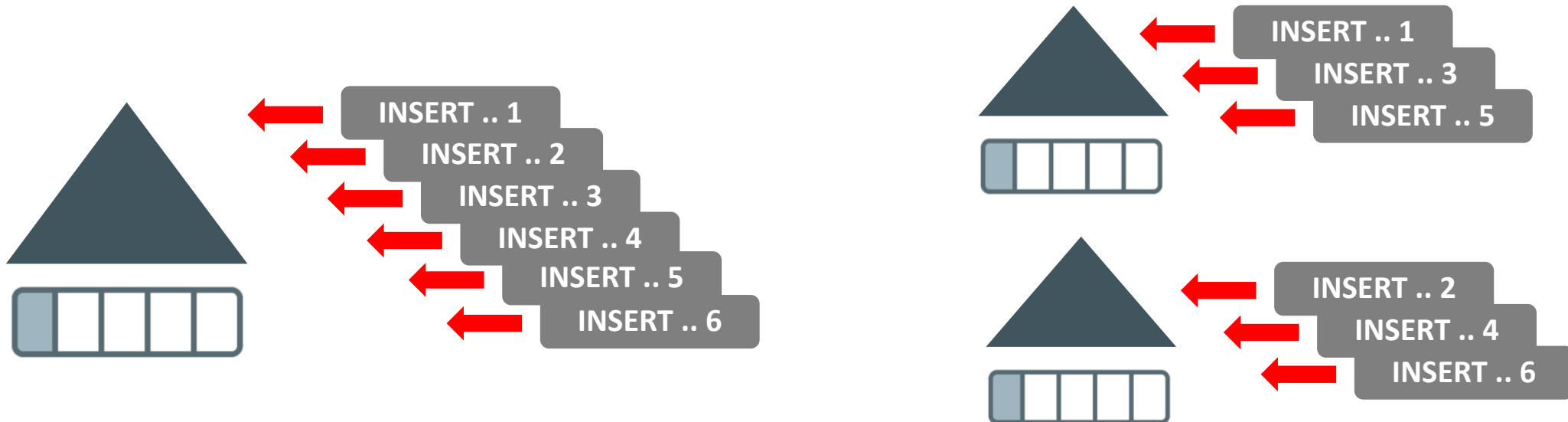
- High DML workload can create **hot spots** (contention) on index blocks
 - E.g. artificial (right hand growing) primary key index



Using Partitioning to eliminate Hot Spots

HASH partitioned index

- High DML workload can create **hot spots** (contention) on index blocks
 - E.g. artificial (right hand growing) primary key index
- With HASH partitioned index you get **warm spots**



Hot Spot Elimination – Use Case

Challenge

- Retail application using object-relational mapping
- Only “common” database functionality is used
- Every single row needs to be updated in a single transaction
- No bulk imports possible at all!
- Thousands of small SQL-Statements issued
- Sudden heavy peaks in user access
 - e.g. Cyber Monday, Christmas trade, special offers, ..
- **Experienced sporadic contention**



Hot Spot Elimination – Use Case

Performance without any application code change

Results from PoC (SKU data load)

- Reference system: 120 SKU's per second
- Exadata Machine (single node load)
 - 2,500 SKU's per second (20x faster)
- Exadata Machine X3-2 (two node load & without partitioning)
 - “only” 1,900 SKU's per second (slower than single node load !!!)
- Exadata Machine X3-2 (two node load & with proper partitioning)
 - 4,800 SKU's per second (40x faster)
- **Proper partitioning enables linear scaling**



Hot Spot Elimination – Use Case

How to (Alternative A, Hash Partitioning on store ID)

- HASH Partitioning creates <n> entry points into the table

```
CREATE TABLE <table_name> (  
  ID          NUMBER(10) NOT NULL,  
  Cn          ... )  
PARTITION BY HASH(ID) PARTITIONS <n>  
TABLESPACE <tablespace_name> STORAGE ( ... );  
  
CREATE UNIQUE INDEX <index_name> ON <table_name>  
(ID) LOCAL TABLESPACE <tablespace_name> STORAGE ( ... );  
  
INSERT INTO <table_name> (ID, ...)  
SELECT SEQ_ID.nextval, ... ;
```



Hot Spot Elimination – Use Case

How to (Alternative B, List Partitioning on instance #)

- Sequence SEQ_ID forces ID to be unique in each partition!
- List Partitioning completely separates the entry points per instance

```
CREATE TABLE <table_name> (  
  ID          NUMBER(10) NOT NULL,  
  Cn          ... ..  
  INSTANCE NUMBER NUMBER(1) DEFAULT sys_context('USERENV','INSTANCE') NOT NULL)  
PARTITION BY LIST (INSTANCE NUMBER)  
( PARTITION P1 VALUES (1),  
  PARTITION P2 VALUES (2),  
  ...  
  PARTITION Pn VALUES (n))  
TABLESPACE <tablespace_name> STORAGE ( ... );  
  
CREATE UNIQUE INDEX <index name> ON <table_name>  
(ID, INSTANCE NUMBER) LOCAL TABLESPACE <tablespace_name> STORAGE ( ... );  
  
INSERT INTO <table_name> (ID, ...) SELECT SEQ_ID.nextval, ... ;
```



Hot Spot Elimination – Use Case

How to (Enhanced alternative B, Hash Partitioning on instance #)

- Sequence SEQ_ID forces ID to be unique in each partition!

```
CREATE TABLE <table_name> (  
  ID          NUMBER(10) NOT NULL,  
  Cn          ...  
  INSTANCE NUMBER NUMBER(1) DEFAULT sys_context('USERENV','INSTANCE') NOT NULL)  
PARTITION BY LIST (INSTANCE_NUMBER)  
SUBPARTITION BY HASH (ID) SUBPARTITIONS <m>  
( PARTITION P1 VALUES(1),  
  PARTITION P2 VALUES(2),  
  ...  
  PARTITION Pn VALUES(n))  
TABLESPACE <tablespace_name> STORAGE ( ... );  
CREATE UNIQUE INDEX <index name> ON <table_name>  
(ID, INSTANCE NUMBER) LOCAL TABLESPACE <tablespace_name> STORAGE ( ... );  
INSERT INTO <table_name> (ID, ...) SELECT SEQ_ID.nextval, ... ;
```



Smart partial partition exchange

“Filtered partition maintenance”



Partition Exchange for Loading and Purging

- Remove and add data as metadata only operations
 - Exchange the metadata of partition and table
- Data load: standalone table contains new data to being loaded
- Data purge: partition containing data is exchanged with empty table
- Drop partition alternative for purge
 - Data is gone forever

<TABLE>



EVENTS Table
May 18 th 2021
May 19 th 2021
May 20 th 2021
May 21 st 2021
May 22 nd 2021
May 23 rd 2021



Smart Partial Partition Exchange

- Sounds easy but ...
- What to do if partition boundaries are not 100% aligned?
 - “Partial Purging”



Smart Partial Partition Exchange

Partial Purging

- Lock partition to being purged

```
LOCK TABLE ... PARTITION ...
```

EVENTS Table

May 18th 2021

May 19th 2021

May 20th 2021

May 21st 2021

May 22nd 2021

May 23rd 2021



Smart Partial Partition Exchange

Partial Purging

- Lock partition to being purged

```
LOCK TABLE ... PARTITION ...
```

- Create table containing remaining data set

```
CREATE TABLE ... AS SELECT WHERE ...
```

“REST”

EVENTS
Table

May 18th 2021

May 19th 2021

May 20th 2021

May 21st 2021

May 22nd 2021

May 23rd 2021



Smart Partial Partition Exchange

Partial Purging

- Lock partition to being purged

```
LOCK TABLE ... PARTITION ...
```

- Create table containing remaining data set

```
CREATE TABLE ... AS SELECT WHERE ...
```

- Create necessary indexes, if any

“REST”

EVENTS
Table

May 18th 2021

May 19th 2021

May 20th 2021

May 21st 2021

May 22nd 2021

May 23rd 2021



Smart Partial Partition Exchange

Partial Purging

- Lock partition to being purged

```
LOCK TABLE ... PARTITION ...
```

- Create table containing remaining data set

```
CREATE TABLE ... AS SELECT WHERE ...
```

- Create necessary indexes, if any

- Exchange partition

```
ALTER TABLE ... EXCHANGE PARTITION ...
```

May 18th 2021



EVENTS
Table

“REST”

May 19th 2021

May 20th 2021

May 21st 2021

May 22nd 2021

May 23rd 2021



Smart Partial Partition Exchange

Partial Purging

- Lock partition to being purged

```
LOCK TABLE ... PARTITION ...
```

- Create table containing remaining data set

```
CREATE TABLE ... AS SELECT WHERE ...
```

- Create necessary indexes, if any

- Exchange partition

```
ALTER TABLE ... EXCHANGE PARTITION ...
```

May 18th 2021



EVENTS
Table

“REST”

May 19th 2021

May 20th 2021

May 21st 2021

May 22nd 2021

May 23rd 2021



Exchange in the presence of unique and primary key constraints

Unique Constraints/Primary Keys

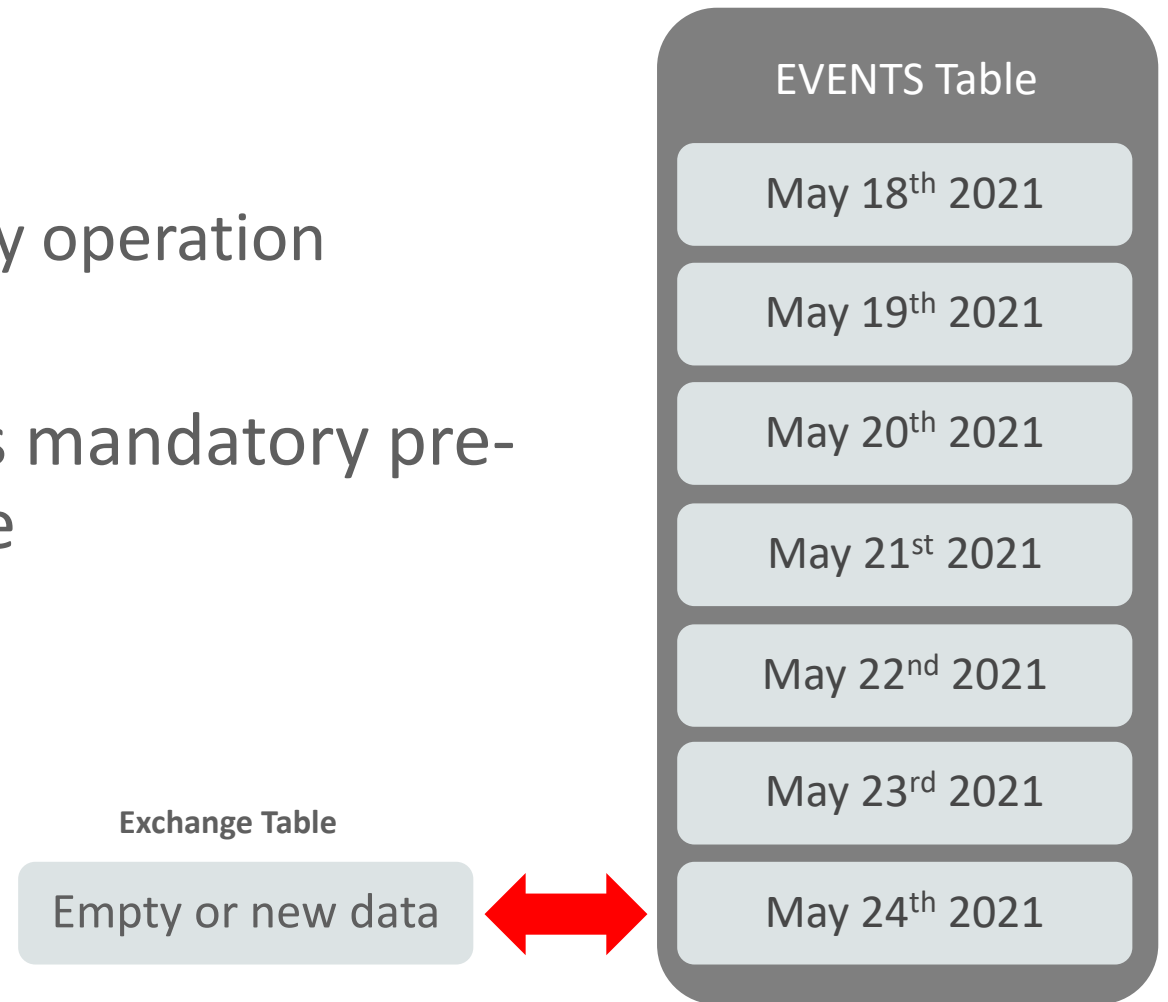
- Unique constraints are enforced with unique indexes
 - Primary key constraint adds NOT NULL to column
 - Table can have only one primary key (“unique identifier”)
- Partitioned tables offer two types of indexes
 - Local indexes
 - Global index, both partitioned and non-partitioned



Partition Exchange

A.k.a Partition Loading and Purging

- Remove and add data as metadata-only operation
 - Exchange the metadata of partitions
- Same logical shape for both tables is mandatory pre-requirement for successful exchange
 - Same number and data type of columns
 - Note that column name does not matter
 - Same constraints
 - Same number and type of indexes



Partition Exchange

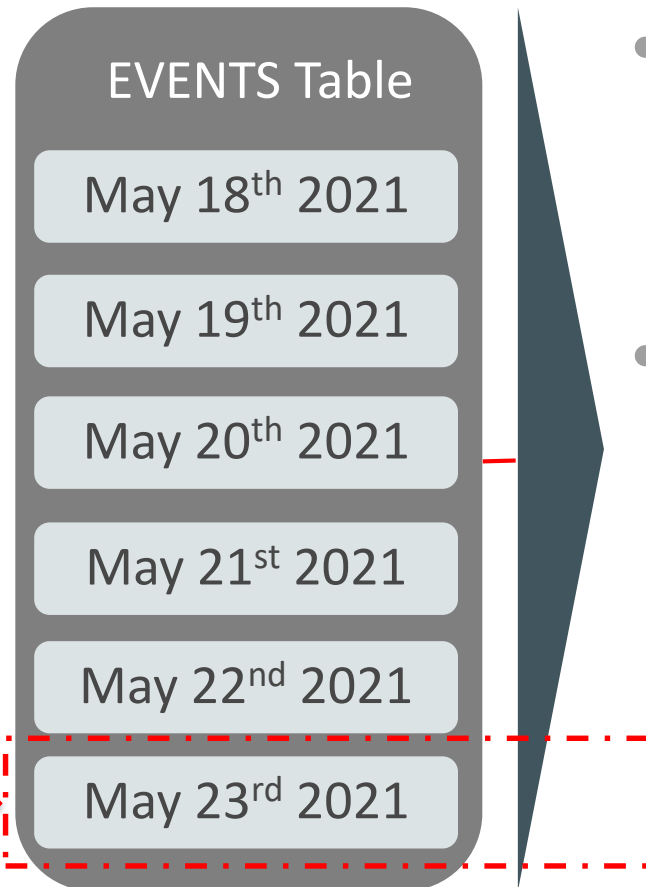
Local Indexes



- Any index on the exchange table is equivalent to a local partitioned index

Partition Exchange

Local Indexes



- Any index on the exchange table is equivalent to a local partitioned index
- What do I do when the PK index on the partitioned table needs global index enforcement?
 - Remember the requirement of logical equivalence ...

Partition Exchange and PK/Unique Constraint

The Dilemma

- Global indexes only exist for a partitioned table
 - But I need the index for the exchange table for uniqueness ...



Partition Exchange and PK/Unique Constraint

Not Really a Dilemma

- Global indexes only exist for a partitioned table
 - But I need the index for the exchange table for uniqueness ...
- Not generically true
 - Unique index only needed for **enabled** constraints
 - Enforcement for new or modified data through index probe



Partition Exchange and PK/Unique Constraint

Not Really a Dilemma

- Global indexes only exist for a partitioned table
 - But I need the index for the exchange table for uniqueness ...
- Not generically true
 - Unique index only needed for **enabled** constraints
 - Enforcement for new or modified data through index probe
 - **Disabled** constraint prevents data insertion

```
SQL> alter table tt add(constraint x unique (col1) disable validate);  
  
Table altered.  
  
SQL> insert into tt values(1,2);  
insert into tt values(1,2);  
*  
ERROR at line 1:  
ORA-25128: No insert/update/delete on table with constraint (SCOTT.X)  
disabled and validated
```



Partition Exchange and PK/Unique Constraint

The solution

- The partitioned target table
 - PK or unique constraint that is enforced by global index (partitioned or non-partitioned)
- The standalone table to be exchanged (“exchange table”)
 - Equivalent disabled validated constraint
 - No index for enforcement, no exchange problem



Partition Exchange and PK/Unique Constraint

A simple example

```
SQL > CREATE TABLE tx_simple
 2  (
 3      TRANSACTION_KEY          NUMBER,
 4      INQUIRY_TIMESTAMP        TIMESTAMP(6),
 5      RUN_DATE                 DATE
 6  )
 7  PARTITION BY RANGE (RUN_DATE)
 8  (
 9      PARTITION TRANSACTION_202105 VALUES LESS THAN (TO_DATE('20210601', 'yyyymmdd')),
10     PARTITION TRANSACTION_202106 VALUES LESS THAN (TO_DATE('20210701', 'yyyymmdd')),
11     PARTITION TRANSACTION_202107 VALUES LESS THAN (TO_DATE('20210801', 'yyyymmdd')),
12     PARTITION TRANSACTION_202108 VALUES LESS THAN (TO_DATE('20210901', 'yyyymmdd')),
13     PARTITION TRANSACTION_202109 VALUES LESS THAN (TO_DATE('20211001', 'yyyymmdd')),
14     PARTITION TRANSACTION_202110 VALUES LESS THAN (TO_DATE('20211101', 'yyyymmdd')),
15     PARTITION TRANSACTION_MAX VALUES LESS THAN (MAXVALUE)
16  )
17  /
```

Table created.



Partition Exchange and PK/Unique Constraint

A simple example

```
SQL > CREATE TABLE tx_simple
 2  (
 3      TRANSACTION_KEY          NUMBER,
 4      INQUIRY_TIMESTAMP        TIMESTAMP(6),
 5      RUN_DATE                  DATE
 6  )
 7  PARTITION BY RANGE (RUN_DATE)
 8  (
 9      PARTITION TRANSACTION_202105 VALUES LESS THAN (TO_DATE('20210601', 'yyyymmdd')),
10     PARTITION TRANSACTION_202106 VALUES LESS THAN (TO_DATE('20210701', 'yyyymmdd')),
11     PARTITION TRANSACTION_202107 VALUES LESS THAN (TO_DATE('20210801', 'yyyymmdd')),
12     PARTITION TRANSACTION_202108 VALUES LESS THAN (TO_DATE('20210901', 'yyyymmdd')),
13     PARTITION TRANSACTION_202109 VALUES LESS THAN (TO_DATE('20211001', 'yyyymmdd')),
14     PARTITION TRANSACTION_202110 VALUES LESS THAN (TO_DATE('20211101', 'yyyymmdd')),
15     PARTITION TRANSACTION_MAX VALUES LESS THAN (MAXVALUE)
16  )
17  /
```

Table created.

```
SQL > INSERT into tx_simple (
 2      select object_id, LAST_DDL_TIME,
 3             add months(TO_DATE('20210501', 'yyyymmdd'), mod(OBJECT_ID,
 4             12))
 5      from DBA_OBJECTS
 6      where object_id is not null)
 6  /
```

73657 rows created.



Partition Exchange and PK/Unique Constraint

A simple example

```
SQL > CREATE TABLE tx_simple
2   (
3     TRANSACTION_KEY      NUMBER,
4     INQUIRY_TIMESTAMP    TIMESTAMP(6),
5     RUN_DATE              DATE
6   )
7   PARTITION BY RANGE (RUN_DATE)
8   (
9     PARTITION TRANSACTION_202105 VALUES LESS THAN (TO_DATE('20210601', 'yyyymmdd')),
10    PARTITION TRANSACTION_202106 VALUES LESS THAN (TO_DATE('20210701', 'yyyymmdd')),
11    PARTITION TRANSACTION_202107 VALUES LESS THAN (TO_DATE('20210801', 'yyyymmdd')),
12    PARTITION TRANSACTION_202108 VALUES LESS THAN (TO_DATE('20210901', 'yyyymmdd')),
13    PARTITION TRANSACTION_202109 VALUES LESS THAN (TO_DATE('20211001', 'yyyymmdd')),
14    PARTITION TRANSACTION_202110 VALUES LESS THAN (TO_DATE('20211101', 'yyyymmdd')),
15    PARTITION TRANSACTION_MAX VALUES LESS THAN (MAXVALUE)
16  )
```

```
SQL > INSERT into tx_simple (
2     select object_id, LAST_DDL_TIME,
3     add months(TO_DATE('20210501', 'yyyymmdd'), mod(OBJECT_ID,
12))
```

```
SQL > CREATE UNIQUE INDEX tx_simple_PK ON tx_simple (TRANSACTION_KEY) nologging
2   GLOBAL PARTITION BY RANGE (TRANSACTION_KEY) (
3     PARTITION P_Max VALUES LESS THAN (MAXVALUE)
4   )
5 /
```

Index created.

```
SQL > ALTER TABLE tx_simple ADD ( CONSTRAINT tx_simple_PK PRIMARY KEY (TRANSACTION_KEY)
2   USING INDEX nologging);
```

Table altered.



Partition Exchange and PK/Unique Constraint

A simple example, cont.

```
SQL > create table DAILY_ETL_table
2     as
3     select * from tx_simple partition (TRANSACTION_202107);
```

Table created.

```
SQL > alter table daily_etl_table add ( constraint pk_etl primary key (transaction_key) disable validate);
```

Table altered.

```
SQL > alter table tx_simple
2     exchange partition TRANSACTION_202107
3     with table daily_ETL_table
4     including indexes
5     --excluding indexes
6     WITHOUT VALIDATION
7     UPDATE GLOBAL INDEXES
8     /
```

Table altered.



Partition Elimination Table (PET)

Partition Pruning

- Works for simple and complex SQL statements
- Transparent to any application
- Works for row store and in-memory store
- Two flavors of pruning
 - Static pruning at compile time
 - Dynamic pruning at runtime
- Complementary to Exadata Storage Server
 - Partitioning prunes logically through partition elimination
 - Exadata prunes physically through storage indexes
 - Further data reduction through filtering and projection



Static Partition Pruning

```
SELECT avg( luminosity ) FROM events
WHERE event_date BETWEEN '01-MAR-2020' and '31-MAY-2020'
AND sensor_type = '...'
AND event_type IN ( ... )
```

2020-JAN

2020-FEB

2020-MAR

2020-APR

2020-MAY

2020-JUN

- Relevant Partitions are known at compile time
 - Look for actual values in PSTART/PSTOP columns in the plan
- Optimizer has most accurate information for the SQL statement



Static Pruning

Sample Plan

```
SELECT avg( luminosity )
FROM atlas.events s, atlas.runs t
WHERE s.time_id = t.time_id
AND s.time_id between TO_DATE('07-JAN-2020', 'DD-MON-YYYY')
                    and TO_DATE('11-JAN-2020', 'DD-MON-YYYY');
```

Plan hash value: 2025449199

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				3 (100)			
1	SORT AGGREGATE		1	12				
2	PARTITION RANGE ITERATOR		313	3756	3 (0)	00:00:01	9	13
* 3	TABLE ACCESS FULL	EVENTS	313	3756	3 (0)	00:00:01	9	13

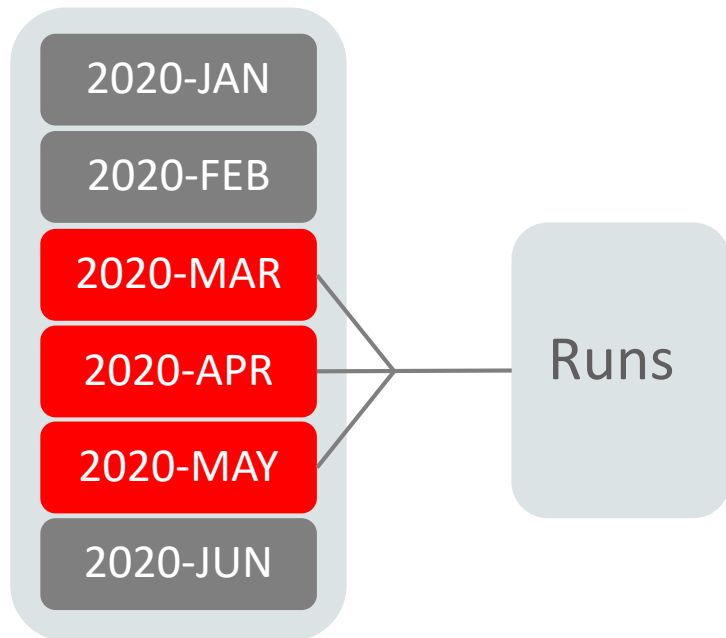
Predicate Information (identified by operation id):

```
3 - filter("S"."TIME_ID"<=TO_DATE(' 2020-01-11 00:00:00', 'syyy-mm-dd hh24:mi:ss'))
```

22 rows selected.



Dynamic Partition Pruning



```
SELECT avg( luminosity )  
FROM events s, runs t  
WHERE t.time_id = s.time_id  
AND t.calendar_month_desc IN  
      ('MAR-2020', 'APR-2020', 'MAY-2020');
```

- Advanced Pruning mechanism for complex queries
- Relevant partitions determined at runtime
 - Look for the word 'KEY' in PSTART/PSTOP columns in the Plan

Dynamic Partition Pruning

Sample Plan – Nested Loop

```
SELECT avg(luminosity)
FROM events s, runs t
WHERE s.time_id = t.time_id
AND t.calendar_month_desc in ('MAR-2020', 'APR-2020', 'MAY-2020')
```

Plan hash value: 1350851517

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				13 (100)			
1	SORT AGGREGATE		1	28				
2	NESTED LOOP		2	56	13 (0)	00:00:01		
* 3	TABLE ACCESS FULL	RUNS	2	32	13 (8)	00:00:01		
4	PARTITION RANGE ITERATOR		2	24	0 (0)		KEY	KEY
* 5	TABLE ACCESS FULL	EVENTS	2	24	0 (0)		KEY	KEY

Predicate Information (identified by operation id):

- 3 - filter(("T"."CALENDAR_MONTH_DESC"='MAR-2020' OR "T"."CALENDAR_MONTH_DESC"='APR-2020' OR "T"."CALENDAR_MONTH_DESC"='MAY-2020'))
- 5 - filter("T"."TIME_ID"="S"."TIME_ID")

26 rows selected.



Dynamic Partition Pruning

Sample Plan – Nested Loop

```
SELECT avg( luminosity )
FROM events s, runs t
WHERE s.time_id = t.time_id
AND t.calendar_month_desc in ('MAR-2020', 'APR-2020', 'MAY-2020')
```

Plan hash value: 1350851517

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				13 (100)			
1	SORT AGGREGATE		1	28				
2	NESTED LOOP		2	56	13 (0)	00:00:01		
* 3	TABLE ACCESS FULL	RUNS	2	32	13 (8)	00:00:01		
4	PARTITION RANGE ITERATOR		2	24	0 (0)		KEY	KEY
* 5	TABLE ACCESS FULL	EVENTS	2	24	0 (0)		KEY	KEY

Predicate Information (identified by operation id):

```
3 - filter(("T"."CALENDAR_MONTH_DESC"='MAR-2020' OR "T"."CALENDAR_MONTH_DESC"='APR-2020'
          OR "T"."CALENDAR_MONTH_DESC"='MAY-2020'))
5 - filter("T"."TIME_ID"="S"."TIME_ID")
```

26 rows selected.



Dynamic Partition Pruning

Sample Plan – Sub-query pruning

Please avoid optimizer hints!

```
SELECT /*+ FULL(s) USE HASH(s, t) CARDINALITY(s, 10000000000) */ avg_ ( luminosity )
FROM events s, runs t
WHERE s.time_id = t.time_id
AND t.calendar_month_desc in ('MAR-2020', 'APR-2020', 'MAY-2020')
```

Plan hash value: 2475767165

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				2000K(100)			
1	SORT AGGREGATE		1	28				
* 2	HASH JOIN		24M	646M	2000K(100)	06:40:01		
* 3	TABLE ACCESS FULL	RUNS	2	32	43 (8)	00:00:01		
4	PARTITION RANGE SUBQUERY		10G	111G	1166K(100)	03:53:21	KEY(SQ)	KEY(SQ)
5	TABLE ACCESS FULL	EVENTS	10G	111G	1166K(100)	03:53:21	KEY(SQ)	KEY(SQ)

Predicate Information (identified by operation id):

- 2 - access ("S"."TIME_ID"="T"."TIME_ID")
- 3 - filter(("T"."CALENDAR_MONTH_DESC"='MAR-2020' OR "T"."CALENDAR_MONTH_DESC"='APR-2020' OR "T"."CALENDAR_MONTH_DESC"='MAY-2020'))

26 rows selected.



Dynamic Partition Pruning

Sample Plan - Bloom filter pruning

```
SELECT avg(luminosity)
FROM events s, runs t
WHERE s.time_id = t.time_id
AND t.calendar_month_desc in ('MAR-2020', 'APR-2020', 'MAY-2020')
```

Plan hash value: 365741303

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				19 (100)			
1	SORT AGGREGATE		1	28				
* 2	HASH JOIN		2	56	19 (100)	00:00:01		
3	PART JOIN FILTER CREATE	:BF0000	2	32	13 (8)	00:00:01		
* 4	TABLE ACCESS FULL	RUNS	2	32	13 (8)	00:00:01		
5	PARTITION RANGE JOIN-FILTER		960	11520	5 (0)	00:00:01	:BF0000	:BF0000
6	TABLE ACCESS FULL	EVENTS	960	11520	5 (0)	00:00:01	:BF0000	:BF0000

Predicate Information (identified by operation id):

```
2 - access("S"."TIME_ID"="T"."TIME_ID")
4 - filter( ("T"."CALENDAR_MONTH_DESC"='MAR-2020'
OR "T"."CALENDAR_MONTH_DESC"='APR-2020'
OR "T"."CALENDAR_MONTH_DESC"='MAY-2020'))
```

27 rows selected.



“AND” Pruning

```
FROM events s, runs t ...  
WHERE s.time id = t.time id ..  
AND t.cal year in (2020,2021)  
AND s.time id  
  between TO_DATE('01-JAN-2020','DD-MON-YYYY')  
  and      TO_DATE('01-JAN-2021','DD-MON-YYYY')
```

Dynamic pruning

Static pruning

- All predicates on partition key will be used for pruning
 - Dynamic and static predicates will now be used combined
- Example:
 - Star transformation with pruning predicate on both the FACT table and a dimension



“AND” Pruning

Sample Plan

Plan hash value: 552669211

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT		1	24	17 (12)	00:00:01		
1	SORT AGGREGATE		1	24				
* 2	HASH JOIN		204	4896	17 (12)	00:00:01		
3	PART JOIN FILTER CREATE	:BF0000	185	2220	13 (8)	00:00:01		
* 4	TABLE ACCESS FULL	RUNS	185	2220	13 (8)	00:00:01		
5	PARTITION RANGE AND		313	3756	3 (0)	00:00:01	KEY (AP)	KEY (AP)
* 6	TABLE ACCESS FULL	EVENTS	313	3756	3 (0)	00:00:01	KEY (AP)	KEY (AP)

Predicate Information (identified by operation id):

```
2 - access("S"."TIME_ID"="T"."TIME_ID")
4 - filter("T"."TIME_ID"<=TO_DATE(' 2021-01-01 00:00:00', 'syyy-mm-dd hh24:mi:ss') AND
          ("T"."CAL_YEAR"=2020 OR "T"."CAL_YEAR"=2021) AND
          "T"."TIME_ID">=TO_DATE(' 2020-01-01 00:00:00', 'syyy-mm-dd hh24:mi:ss'))
6 - filter("S"."TIME_ID"<=TO_DATE(' 2015-01-01 00:00:00', 'syyy-mm-dd hh24:mi:ss'))
```

22 rows selected.



Ensuring Partition Pruning

Don't use functions on partition key filter predicates

```
SELECT avg(luminosity)
FROM atlas.events s, atlas.runs t
WHERE s.time_id = t.time_id
AND TO_CHAR(s.time_id, 'YYYYMMDD') between '20200101' and '20210101'
```

Plan hash value: 672559287

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				6 (100)			
1	SORT AGGREGATE		1	12				
2	PARTITION RANGE ALL		2	24	6 (17)	00:00:01	1	16
* 3	TABLE ACCESS FULL	EVENTS	2	24	6 (17)	00:00:01	1	16

Predicate Information (identified by operation id):

```
3 - filter((TO_CHAR(INTERNAL_FUNCTION("S"."TIME_ID"), 'YYYYMMDD') >= '20200101' AND
            TO_CHAR(INTERNAL_FUNCTION("S"."TIME_ID"), 'YYYYMMDD') <= '20210101'))
```

23 rows selected.



Ensuring Partition Pruning

Don't use functions on partition key filter predicates

```
SELECT avg(luminosity)
FROM atlas.events s, atlas.runs t
WHERE s.time_id = t.time_id
AND TO_CHAR(s.time_id, 'YYYYMMDD') between '20200101' and '20210101'
```



```
SELECT avg(luminosity)
FROM atlas.events s, atlas.run t
WHERE s.time_id = t.time_id
AND s.time_id between TO_DATE('20200101','YYYYMMDD') and TO_DATE('20210101','YYYYMMDD')
```



Plan hash value: 2025449199

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT				3 (100)			
1	SORT AGGREGATE		1	12				
2	PARTITION RANGE ITERATOR		313	3756	3 (0)	00:00:01	9	13
* 3	TABLE ACCESS FULL	EVENTS	313	3756	3 (0)	00:00:01	9	13

Predicate Information (identified by operation id):

```
3 - filter("S"."TIME_ID"<=TO_DATE(' 2021-01-01 00:00:00', 'syyyy-mm-dd hh24:mi:ss'))
```

22 rows selected.

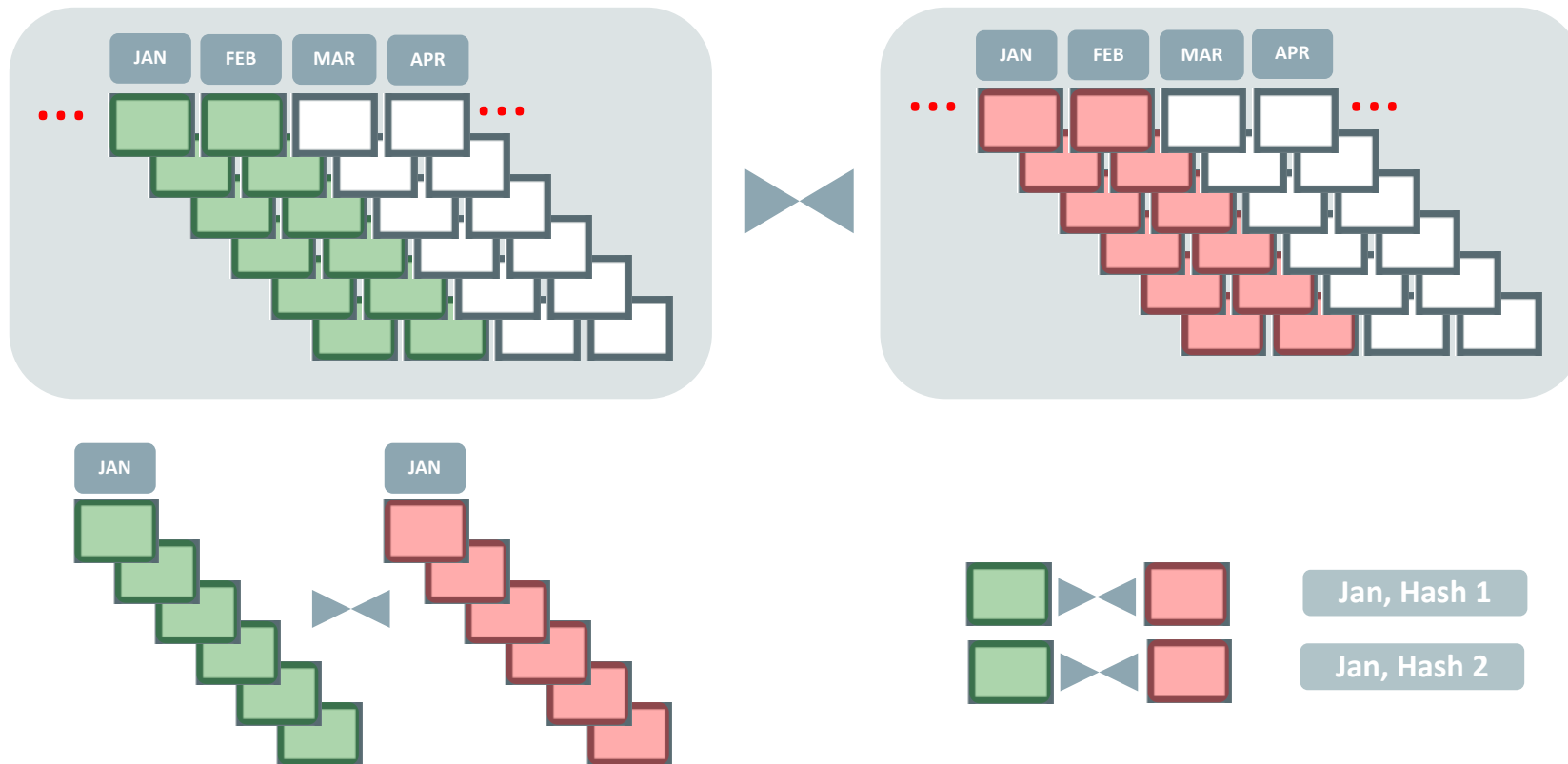
Pstart	Pstop
1	16
1	16

```
' AND  
('))
```



Partition-wise Joins

Partition pruning and PWJ's "at work"

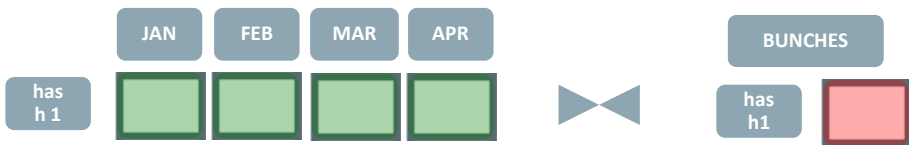
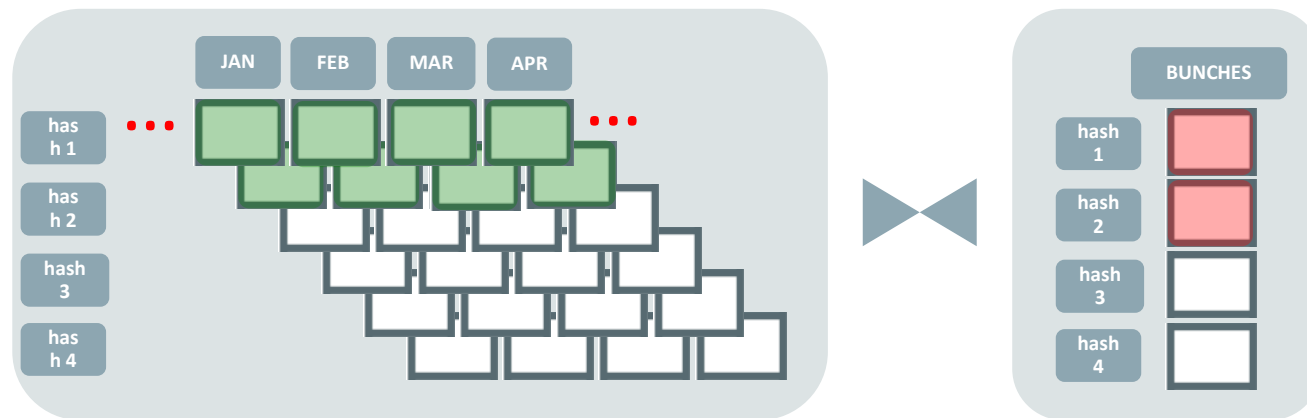


A large join is divided into multiple smaller joins, executed in parallel

- # of partitions to join must be a multiple of DOP.
- Both tables must be partitioned the same way on the join column.

Partition-wise Joins

Partition pruning and PWJ's "at work"



A large join is divided into multiple smaller joins, executed in parallel

- # of partitions to join must be a multiple of Degree of parallelism
- Both tables must be partitioned the same way on the join column

Lets start

The example shown here has been derived from a customers use case.

They collect operational database information from ADDM for all their databases (containers, pluggables).



Data management: partitioning challenges

- **CH 1:** partition pruning when using **multiple columns** for the partitioning key.
- **CH 2:** partition pruning when using **virtual columns** for the partitioning key.
- **CH 3: Minimizing effort** to create/maintain the partitions you need.
- **CH 4:** Overcoming the 1 million partitions limitation.



Challenge 1:

partition pruning when using multiple columns for the partitioning key



Examples for partition pruning

partitioning-key(host_name, instance_name, dbid, snap_id)

- Consider having a table defined like this one:

```
create table test_sysstat_XL
( host_name      varchar2(128)  not null,
  instance_name  varchar2(128)  not null,
  dbid           number(16)     not null,
  snap_id       number(16)     not null,
  ... )

(storage clause)
partition by list ( host_name, instance_name, dbid, snap_id )
```

- The partitioning key helps Oracle to prune (avoid work) in general.
- Well „in general“ but not always as shown in the previous section.



Examples for partition pruning

partitioning-key (host_name, instance_name, dbid, snap_id)

```
select ...  
from test_sysstat_XL  
where host_name = 'blade-1672' ;
```



```
where host_name = 'blade-1672'  
and instance_name = 'sbl4prd' ;
```



```
where host_name = 'blade-1672'  
and instance_name = 'sbl4prd' ;  
and dbid between 10000000000762  
and 10000000001245 ;
```



```
where host_name = 'blade-1672'  
and instance_name = 'sbl4prd' ;  
and dbid between 10000000000762  
and 10000000001245  
and snap_id between 1 and 750 ;
```



Why does Oracle prune?

- Predicates are used, that refer to the leftmost columns of the partitioning key.

What can you expect?

- You can use almost any predicate on these columns to prune effectively.
- Joins are the best examples for predicates, since tables have relationships.



Examples for partition pruning limitations

partitioning-key(host_name, instance_name, dbid, snap_id)

```
select ...  
from test_sysstat_XL  
where instance_name = 'sb14prd' ;
```



host_name is missing ...

```
where dbid between 1000000000762  
and 1000000001245 ;
```



host_name is missing ...
Instance_name is missing ...

```
where hostname = 'blade-1672'  
and dbid between 1000000000762  
and 1000000001245  
and snap_id between 1 and 750 ;
```



instance_name is missing ...

Overcoming partition pruning limitations

partitioning-key(host_name, instance_name, dbid, snap_id)

- If you don't have indexes and storage indexes or don't want to rely on them only, then you need a little help from an object containing the missing attribute values.
- Create a **partition elimination table** a **PET** as follows:

```
create table test_sysstat_PET
( host_name      varchar2(128)  not null,
  instance_name  varchar2(128)  not null,
  dbid           number(16)     not null,
  snap_id       number(16)     not null
)
(storage clause)
as select distinct host_name, instance_name, dbid, snap_id
   from test_sysstat_XL
```

- You don't perform the **select distinct** like that, it takes too long! But you really want it.



Overcoming partition pruning limitations

partitioning-key(host_name, instance_name, dbid, snap_id)

- Create a partition elimination table a PET as follows:

```
create table test_sysstat_PET
( host_name      varchar2(128)  not null,
  instance_name  varchar2(128)  not null,
  dbid           number(16)     not null,
  snap_id       number(16)     not null
)
(storage clause);
```

- Create a simple PL/SQL program to obtain the information from **DBA_TAB_PARTITIONS**
- Convert the **HIGHVALUE** to the information you need and insert it into the PET.
- You have then a table with less than 1 million rows.
- You may want to index, set it in-memory, create a global temp table ... That is your choice.



Overcoming partition pruning limitations

partitioning-key(host_name, instance_name, dbid, snap_id)

- Always join your base table and the PET using the attributes of the partitioning key of the base table. Use that for any query or DML.

```
select /*+ USE_HASH( XL PET ) */
      ( columns_you_want_from_base_table )
from   TEST_SYSSTAT_XL  XL,
      TEST_SYSSTAT_PET PET
where  XL.host_name      = PET.host_name
and    XL.instance_name = PET.instance_name
and    XL.dbid           = PET.dbid
and    XL.snap_id        = PET.snap_id
and    PET.dbid between 1000000000762 and 10000000001245;
```

- Add any predicates including functions on columns in the PET to prune partitions
- Add any predicates on the base table you need.



Overcoming partition pruning limitations

partitioning-key(host_name, instance_name, dbid, snap_id)

- Alternative query using a WITH clause (you'll like the extra benefit

```
WITH wpl AS ( select /*+ MATERIALIZE */
               host_name,
               instance_name,
               dbid,
               snap_id
               from TEST_SYSSTAT_PET )
select /*+ USE_HASH( XL PET ) */
       ( columns_you_want_from_base_table )
from   TEST_SYSSTAT_XL XL,
       WPL PET
where  XL.host_name      = PET.host_name
and    XL.instance_name = PET.instance_name
and    XL.dbid          = PET.dbid
and    XL.snap_id       = PET.snap_id
and    PET.dbid between 1000000000762 and 10000000001245;
```



Overcoming partition pruning limitations

partitioning-key(host_name, instance_name, dbid, snap_id)

access method	host_name	instance_name	dbid	snap_id	returned rows	block read	returned rows / logical reads	elapsed time in sec	M rows returned per second	MB/s scan rate (internal SSD!)	observation made	speed-up	I/O reduction
simple table access	single value				330 M	1.8 M	180.53	38	8.76	379.00	pruning works	1.00	1.00
join with PET-table	single value				330 M	1.8 M	180.43	44	7.51	325.34	pruning works	0.86	1.00
simple table access		single value			212 M	59.0 M	3.59	1201	0.18	384.64	full table scan	1.00	1.00
join with PET-table		single value			212 M	3.5 M	59.77	85	2.49	325.10	pruning works	14.06	16.63
simple table access			single value		6.4 M	58.0 M	0.11	1215	0.01	373.48	full table scan	1.00	1.00
join with PET-table			single value		6.4 M	2.5 M	2.61	58	0.11	332.72	pruning works	20.92	23.48
simple table access				single value	6.7 M	58.0 M	0.12	1167	0.01	388.63	full table scan	1.00	1.00
join with PET-table				single value	6.7 M	2.4 M	2.71	59	0.11	328.85	pruning works	19.87	23.48

* Example TEST_SYSSTAT_XL has 10'482'424'855 records in 203'415 partitions.
No indexes used.

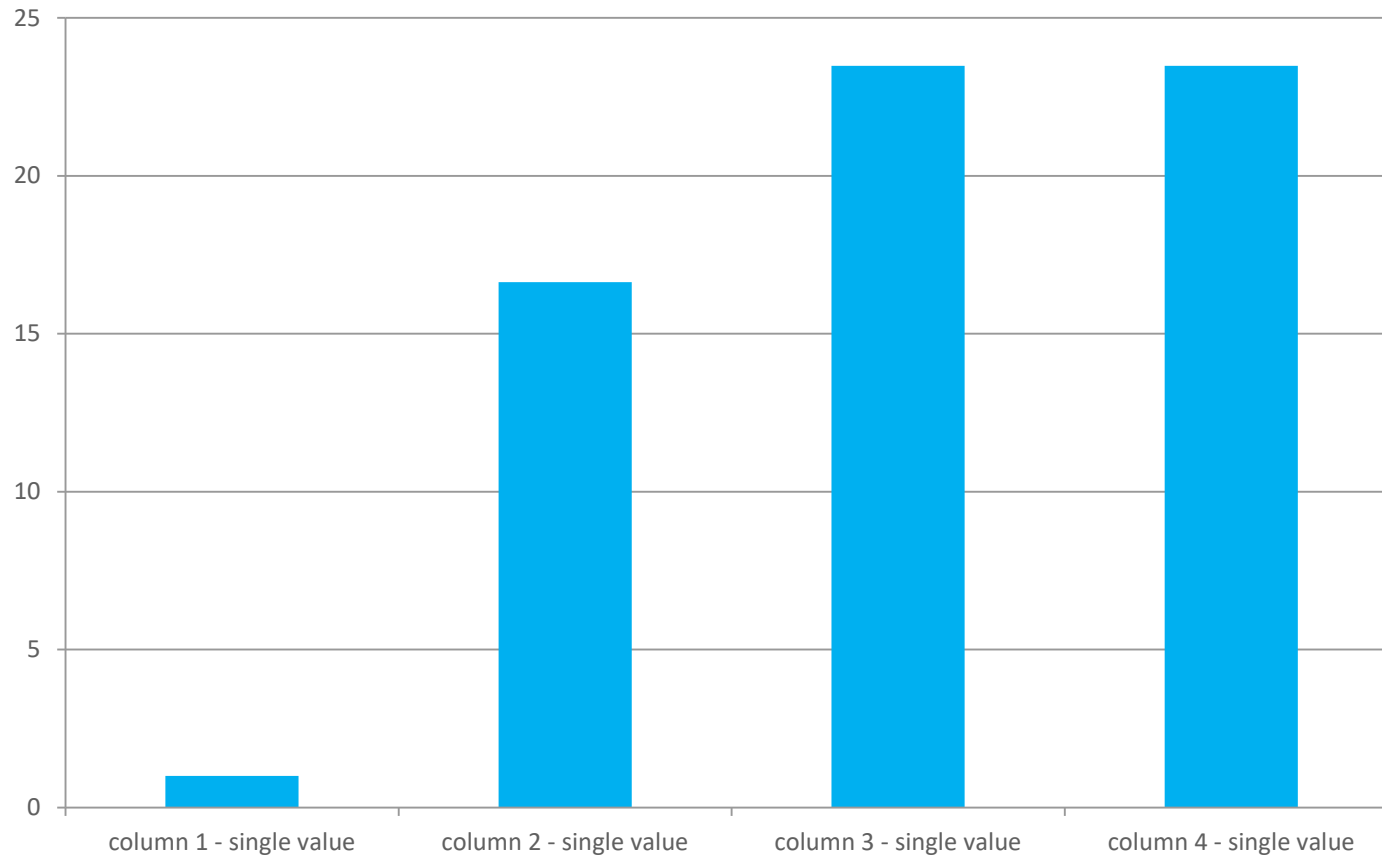
** Test performance using 12.2 on i5-quad-core with 8 GB RAM and 1 TB SSD for data 500 GB SSD for TEMP.



PET I/O reduction – 1 partitioning column – 1 value

partitioning-key(host_name, instance_name, dbid, snap_id)

PET: I/O reduction



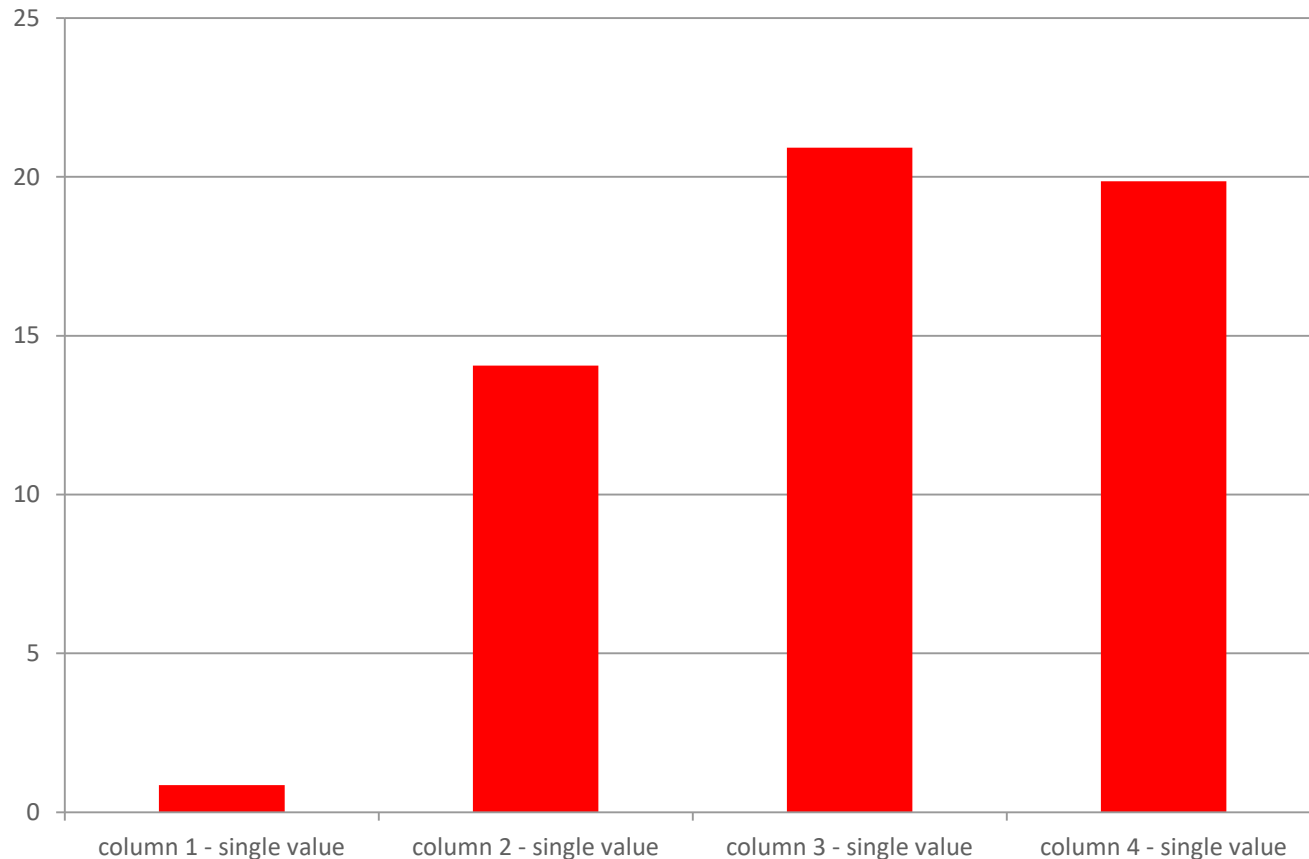
- I/O reduction is given by the distribution of values for the queried column only.
(Example: if you query for a column value accounting for 5% of the records, then generate 20 times less I/O.)
- The left-most column is used to eliminate I/O also without the PET. Thus I/O remains unchanged.



PET runtime reduction – 1 partitioning column – 1 value

partitioning-key(host_name, instance_name, dbid, snap_id)

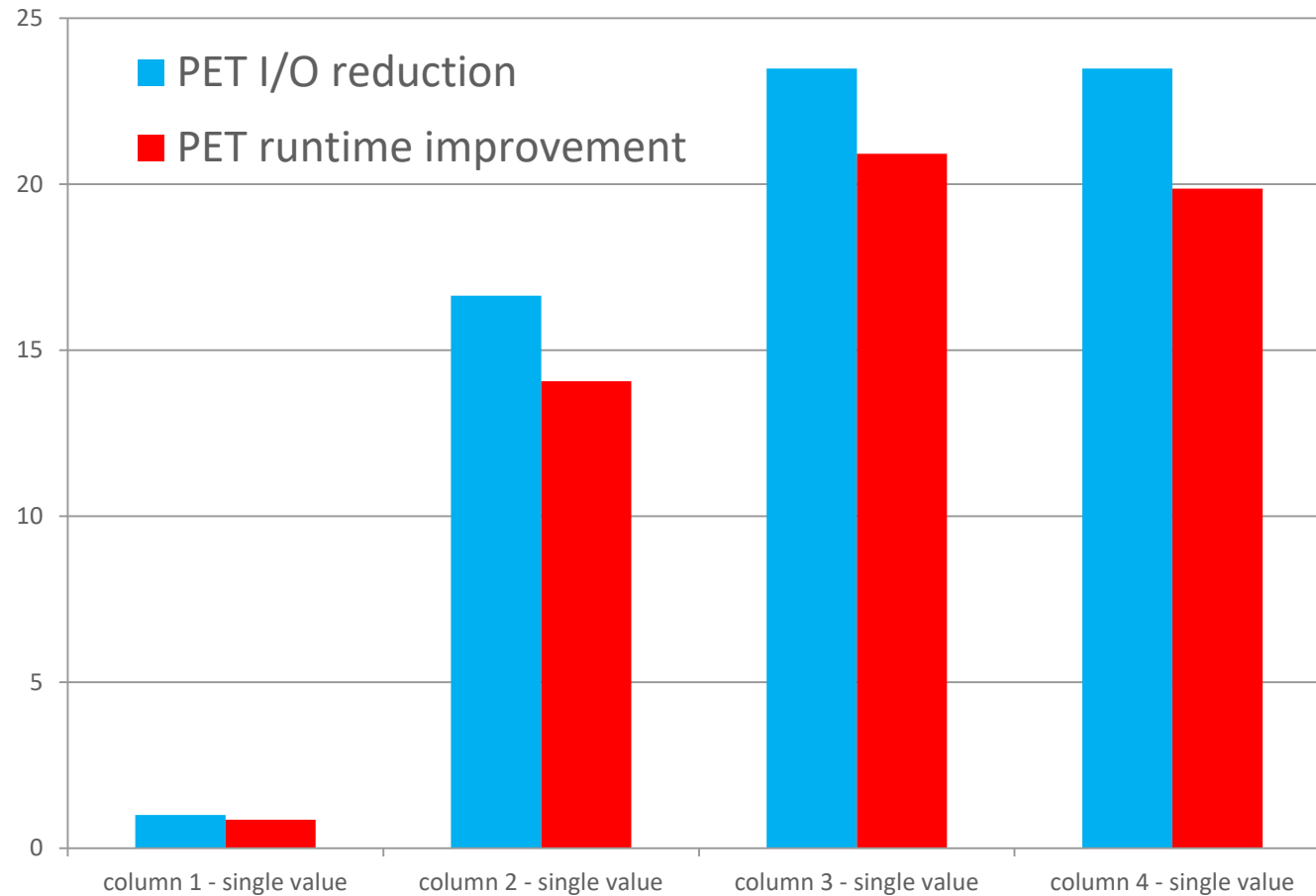
PET: runtime improvement



- Runtime improvement is aligned to I/O reduction.
- Deviations indicate an infrastructure bottleneck in the test-system.
- Using the PET for the leftmost column only results in 10-15% higher runtime in the test-system.



Result: PET I/O reduction leads to runtime improvement



- Runtime improves almost as good as I/O is reduced.
- The observed gap indicates clearly: the test-system has relatively low performance OR the impact of the join with the PET is too high.
- Conclusion: expect the runtime benefit being lower than the I/O reduction.



Overcoming partition pruning limitations – more examples

partitioning-key(host_name, instance_name, dbid, snap_id)

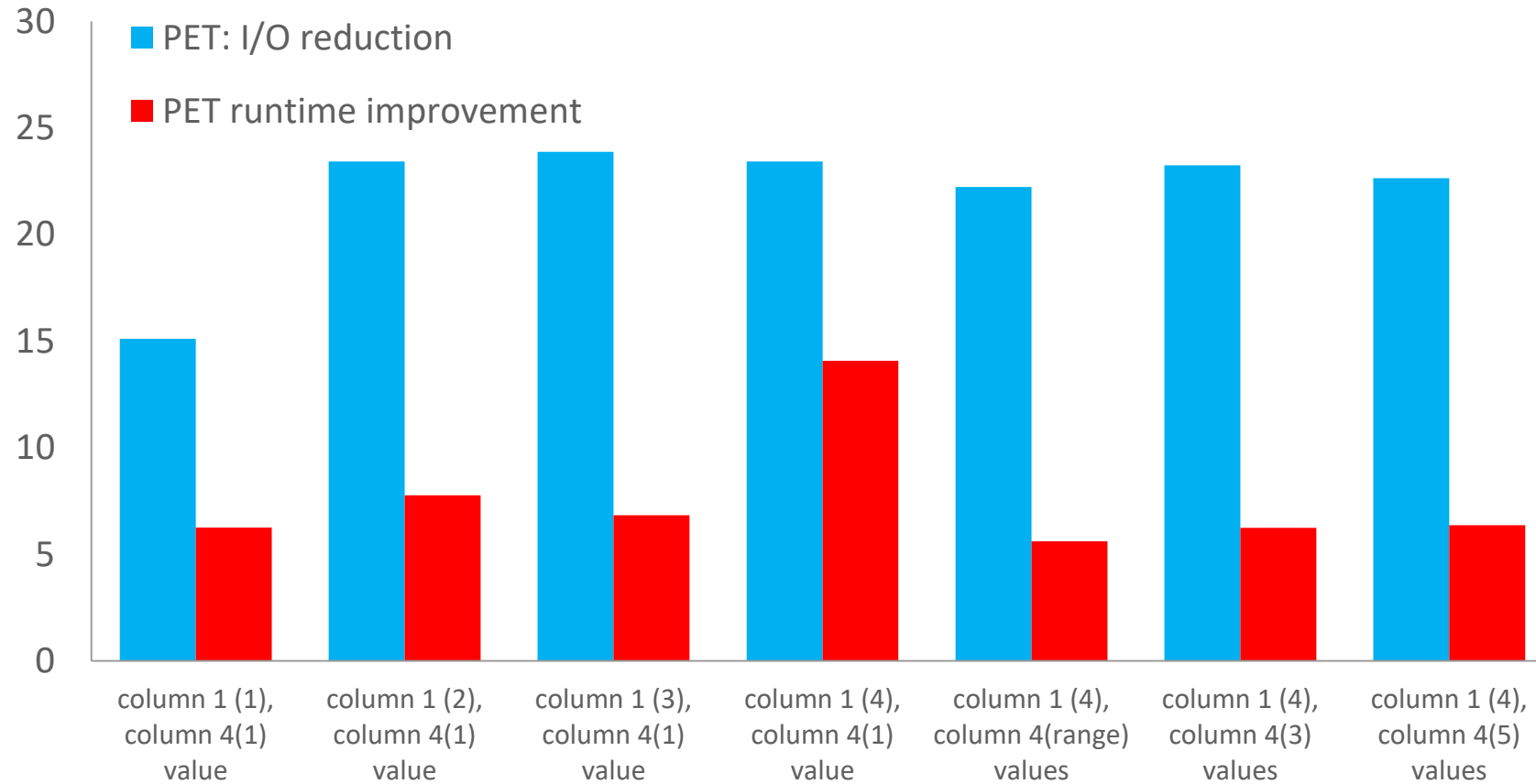
access method	host_name	instance_name	dbid	snap_id	returned rows	block read	returned rows / logical reads	elapsed time in sec	M rows returned per second	observation made	speed-up	I/O reduction
simple table access	single value		single value		6.4 M	1.8 M	3.52	41	0.16	pruning works	1.00	1.00
join with PET-table	single value		single value		6.4 M	0.1 M	53.23	7	0.99	pruning works	6.25	15.10
simple table access	two values		single value		6.4 M	57.9 M	0.11	1300	0.00	full table scan	1.00	1.00
join with PET-table	two values		single value		6.4 M	2.5 M	2.61	167	0.04	pruning works	7.77	23.43
simple table access	three values		single value		6.4 M	59.0 M	0.11	1202	0.01	full table scan	1.00	1.00
join with PET-table	three values		single value		6.4 M	2.5 M	2.61	176	0.04	pruning works	6.83	23.87
simple table access	four values		single value		6.4 M	57.9 M	0.11	1999	0.00	full table scan	1.00	1.00
join with PET-table	four values		single value		6.4 M	2.5 M	2.61	142	0.05	pruning works	14.08	23.43
simple table access	four values		range		32.1 M	57.9 M	0.55	1190	0.03	full table scan	1.00	1.00
join with PET-table	four values		range		32.1 M	2.6 M	12.32	212	0.15	pruning works	5.61	22.23
simple table access	four values		three values		19.3 M	59.0 M	0.33	1215	0.02	full table scan	1.00	1.00
join with PET-table	four values		three values		19.3 M	2.5 M	7.60	195	0.10	pruning works	6.24	23.24
simple table access	four values		five values		32.1 M	59.0 M	0.54	1210	0.03	full table scan	1.00	1.00
join with PET-table	four values		five values		32.1 M	2.6 M	12.32	190	0.17	pruning works	6.37	22.65

* Example TEST_SYSSTAT_XL has 10'482'424'855 records in 203'415 partitions.
No indexes used.

** Test performance using 12.2 on i5-quad-core with 8 GB RAM and 1 TB SSD for data 500 GB SSD for TEMP.



Result: PET I/O reduction leads to runtime improvement



- I/O reduction works fine, if you query multiple values per column
- runtime reduction is lower than I/O reduction, because the test-system is memory bound.



Overcoming partition pruning limitations

using a partition elimination table (PET)

- You can achieve pruning as shown before for the example using LIST partitioning for a multiple column partitioning-key.
- You are not limited in the selection of predicates on the partitioning-key.
- Do you really need a UNIQUE INDEX on your large fact table?
If you don't need then: save storage and index maintenance cost.
(What if you just define constraints on dimension tables and join to them while performing DML?)
- Do you really need a STORAGE INDEX on you large fact table?
No you don't need one. So you can achieve all this on any platform.
If you have STORAGE INDEXES, then this will help always.
- Do you need IN-MEMORY to improve performance significantly?
Yes if you have a subset of partitions, that require lowest response time.
No if the storage system is sufficiently good.
Never forget: The amount of information in RAM and price/performance is most relevant.



Challenge 2:

partition pruning when using virtual columns for the partitioning key



Virtual columns used for the partitioning key

- A virtual column allows you to group records into partitions.
- There is no limitation using virtual columns in the definition of the partitioning key.
- Our previous example had a limitation: We can have only 1 million combinations of (`host_name`, `instance_name`, `dbid`, `snap_id`)
- You might have for a given DBID many SNAP_Ids. Consider the standard settings in Oracle: every 30 minutes you generate a new SNAP_ID. If you want to retain a history of months or years of SNAPS, then you group them.
- You might have additional columns, that qualify for a partitioning key very well. So you might think: can I combine columns into one key column?
That leads to snow-flakes ... 😊



Using virtual column(s) in your base table

partitioning-key(host_name, instance_name, dbid, grp_snap_id)

- Consider changing our previously defined table into this one:

```
create table test_sysstat_XL
( host_name          varchar2(128)    not null,
  instance_name      varchar2(128)    not null,
  dbid               number(16)       not null,
  snap_id            number(16)       not null,
  grp_snap_id        as ( round( snap_id / 30*48, 0 ) ),
  ... )
```

(storage clause)

partition by list (host_name, instance_name, dbid, grp_snap_id)

- The **grp_snap_id** groups 30 days having 48 snapshots = 1 months of SNAP details. Even, if you have 10k combinations of (host_name, instance_name, dbid), you can store on average 100 groups (= 100 months of SNAP details = 8 years and 4 month)
- Limitation: Partition pruning on **grp_snap_id** does not work. You need again a PET.



Using virtual column(s) to your PET table

partitioning-key(host_name, instance_name, dbid, grp_snap_id)

- Create a different partition elimination table a PET as follows:

```
create table test_sysstat_PET
( host_name          varchar2(128)    not null,
  instance_name      varchar2(128)    not null,
  dbid               number(16)       not null,
  grp_snap_id        number(16)       not null    -- do NOT include SNAP_ID
)
(storage clause);
```

- Create a simple PL/SQL program to obtain the information from **DBA_TAB_PARTITIONS**
- Convert the **HIGHVALUE** to the information you need and insert it into the PET.
- You have then a table with less then 1 million rows.
- You may want to index, set it in-memory, create a global temp table ... That is your choice.



Overcoming partition pruning limitations

partitioning-key(host_name, instance_name, dbid, grp_snap_id)

- Always join your base table and the PET using the attributes of the partitioning key of the base table. Use that for any query or DML.

```
select /*+ USE_HASH( XL PET ) */  
      ( columns_you_want_from_base_table )  
from   TEST_SYSSTAT_XL  XL,  
       TEST_SYSSTAT_PET PET  
where  XL.host_name      = PET.host_name  
and    XL.instance_name = PET.instance_name  
and    XL.dbid           = PET.dbid  
and    XL.grp_snap_id    = PET.grp_snap_id -- think twice: is this needed?  
and    PET.dbid between 1000000000762 and 10000000001245;
```

- Add any predicates including functions on columns in the PET to prune partitions
- Add any predicates on the base table you need.
If we add a predicate on XL.SNAP_ID, then we are not that much selective.



Challenge 3:

Minimizing the effort to create the partitions you need



Managing partitions

- In Oracle 11g and 12c you need to define the values for LIST partitions.
- Oracle 12.2 allows you to **automate** the creation of partitions as follows:

```
create table test_sysstat_XL
( host_name          varchar2(128)    not null,
  instance_name      varchar2(128)    not null,
  dbid               number(16)       not null,
  snap_id            number(16)       not null,
  grp_snap_id        as ( round( snap_id / 30*48, 0 ) ),
  stmt_id            number(16)       not null,
  grp_stmt_id        as ( round( stmt_id / 1000, 0 ) ),
  ... )

(storage clause)
partition by list ( host_name, instance_name, dbid,
                   grp_snap_id, grp_stmt_id ) automatic;
```

- Partition names are generated. For example **SYS_P1627471** is the 1'627'471 th automatically generated partition in your instance.



Challenge 4:

Overcoming the 1 million partitions limitation *

* This is an implementation limit.



1 million partitions

Naming them	<ul style="list-style-type: none">• You need to name the partitions using a program or allow Oracle doing it automatically for you.
Getting close to the limit	<ul style="list-style-type: none">• You might feel the panic already: what can I do, once I hit the limit?
Storing them	<ul style="list-style-type: none">• Example: you design partitions to scan roughly for 1 second on your machine. 1 second means today several GB. Thus: a PB scale table is in reach.• Check the size of partitions: compress data and resize segments on schedule.
Using them	<ul style="list-style-type: none">• Check your PGA size (check for PGA_AGGREGATE_TARGET alerts)• Manageability: work selectively for indexing and optimizer statistics.• Using data lifecycle management is strongly advised



Resolution for challenge 4

- Consider you have a table **TEST_SYSSTAT_XL** having many partitions.
- Reorganise your data into tables named **TEST_SYSSTAT_XL_nnn**, where **nnn** is a number.
- Create a view to work with your table **V_TEST_SYSSTAT_XL_UA** using

```
create or replace view V_TEST_SYSSTAT_XL_UA as
select ... from TEST_SYSSTAT_XL_01
UNION ALL
...
UNION ALL
select ... from TEST_SYSSTAT_XL_32;
```

- Flexibility: you define the storage-clause as per table as needed.
- Limitation: DML using the view - you can manage this as we do it for sharding.
- Limitation: a unique key across all these tables – you can manage having additional constraints.
- Limitation: application changes – you can manage this using a synonym instead.



Resolution for challenge 4

- Design rule: define a common partitioning key for these tables.

Example (you might have some application, that collects all your DB performance information):

```
create table test_sysstat_XL_01
( host_name          varchar2(128)  not null,
  instance_name     varchar2(128)  not null,
  dbid              number(16)      not null,
  snap_id           number(16)      not null,
  grp_snap_id       as ( round( snap_id / 30 * 48, 0 ) ),
  stmt_id           number(16)      not null,
  grp_stmt_id       as ( round( stmt_id / 1000, 0 ) ),
  ... )

(storage clause)
partition by list ( host_name, instance_name, dbid,
                   grp_snap_id, grp_stmt_id )
```

- Add further partition key columns, if deemed necessary (it works!).



Example: create detail tables (covering all the partitions)

- Design rule: use virtual column definitions consistently for **master** and any detail tables.

Example:

TEST_SYSSTAT_XL is the master table for the **TEST_SYSSTAT_XL_nnn** tables.

TEST_SYSSTAT_PET is naturally the PET for the **TEST_SYSSTAT_XL_nnn** tables.

```
create table test_ SYSSTAT_XL
( host_name          varchar2(128)  not null,
  instance_name     varchar2(128)  not null,
  dbid              number(16)      not null,
  snap_id           number(16)      not null,
  grp_snap_id       as ( round( snap_id / 30 * 48, 0 ) ),
  stmt_id           number(16)      not null,
  grp_stmt_id       as ( round( stmt_id / 1000, 0 ) ),
  ... )

(storage clause)
partition by list ( host_name, instance_name, dbid,
                   grp_snap_id, grp_stmt_id )
```



Example: access view definition – a local PET per detail table

```
-- UNION ALL access view to simplify all : using the local PET tables for the partitioning keys INCLUDING the origin for the virtual column

create or replace view V_TEST_SYSSSTAT_PET_LCL as
select XL... -- the columns you need
from ( select myxl.*
      from TEST_SYSSSTAT_XL_01 myxl,
           TEST_SYSSSTAT_PET_01 mypel
      where myxl.host_name
            = mypel.host_name
            and myxl.instance_name
            = mypel.instance_name
            and myxl.dbid
            = mypel.dbid
            and myxl.grp_snap_id
            = mypel.grp_snap_id
            and myxl.grp_stmt_id
            = mypel.grp_stmt_id
      UNION ALL
      ...
      UNION ALL
      select myxl.*
      from TEST_SYSSSTAT_XL_32 myxl,
           TEST_SYSSSTAT_PET_32 mypel
      where myxl.host_name
            = mypel.host_name
            and myxl.instance_name
            = mypel.instance_name
            and myxl.dbid
            = mypel.dbid
            and myxl.grp_snap_id
            = mypel.grp_snap_id
            and myxl.grp_stmt_id
            = mypel.grp_stmt_id
      ) XL;
```

- A *local PET* is a partition elimination table containing the set of partition key combinations for a given detail table only.
- Consider having local PET tables named TEST_SYSSSTAT_PET_01 to 32.



Example: access view definition – a global PET for all detail tables

```
-- UNION ALL access view to simplify all : using the global PET table for the partitioning keys excluding the origin for the virtual column

create or replace view V_TEST_SYSSTAT_PET_ALL as
WITH PEL AS ( select /*+ MATERIALIZE */
                host_name,
                instance_name,
                dbid,
                grp_snap_id,
                grp_stmt_id
            from   TEST_SYSSTAT_PET_ALL
            )
select /*+ USE_HASH( XL EL ) */
XL.*
from   ( select *
        from   TEST_SYSSTAT_XL_01
        UNION ALL
        ...
        UNION ALL
        select *
        from   TEST_SYSSTAT_XL_32
        ) XL,
PEL EL
where  XL.host_name
      = EL.host_name
and    XL.instance_name
      = EL.instance_name
and    XL.dbid
      = EL.dbid
and    XL.grp_snap_id
      = EL.grp_snap_id
and    XL.grp_stmt_id
      = EL.grp_stmt_id;
```

- A *global PET* is a partition elimination table containing the super set of partition key combinations across all detail tables.
- Consider having local PET tables named TEST_SYSSTAT_PET_01 to 32. Think of the global PET table being populated using:

```
insert into TEST_SYSSTAT_PET_ALL
select distinct host_name, instance_name, dbid,
grp_snap_id, grp_stmt_id
from
(
select * from TEST_SYSSTAT_PET_01 union
...
select * from TEST_SYSSTAT_PET_32
);
```



Example: 32 tables, 40 billion records, 299 k partitions

MY_TABLE_NAME	NUM_PARTITIONS	NUM_RECORDS	MIN_RECORDS_PER_PART	MAX_RECORDS_PER_PART	AVG_RECORDS_PER_PART
TEST_SYSSTAT_XL_01	923	10123264	0	472064	10967
TEST_SYSSTAT_XL_02	15856	142682112	0	9216	8998
TEST_SYSSTAT_XL_03	4	18432	0	9216	4608
TEST_SYSSTAT_XL_04	4	18432	0	9216	4608
TEST_SYSSTAT_XL_05	7483	139860518	0	1024000	18690
TEST_SYSSTAT_XL_06	7262	130572432	0	18000	17980
TEST_SYSSTAT_XL_07	4	18432	0	9216	4608
TEST_SYSSTAT_XL_08	4	18432	0	9216	4608
TEST_SYSSTAT_XL_09	2248	283251551	0	968704	126001
TEST_SYSSTAT_XL_10	2785	526804230	0	1024000	189157
TEST_SYSSTAT_XL_11	2248	237494749	0	507904	105647
TEST_SYSSTAT_XL_12	1873	10597725	0	9216	5658
TEST_SYSSTAT_XL_13	2248	283042990	0	968704	125908
TEST_SYSSTAT_XL_14	2746	523097999	0	1024000	190494
TEST_SYSSTAT_XL_15	2248	243959412	0	507904	108522
TEST_SYSSTAT_XL_16	1873	10586133	0	9216	5651
TEST_SYSSTAT_XL_17	2248	283251551	0	968704	126001
TEST_SYSSTAT_XL_18	2785	526804230	0	1024000	189157
TEST_SYSSTAT_XL_19	2248	237494749	0	507904	105647
TEST_SYSSTAT_XL_20	1873	10597725	0	9216	5658
TEST_SYSSTAT_XL_21	2248	283042990	0	968704	125908
TEST_SYSSTAT_XL_22	2746	523097999	0	1024000	190494
TEST_SYSSTAT_XL_23	2248	243959412	0	507904	108522
TEST_SYSSTAT_XL_24	29808	7969194372	0	6575880	267350
TEST_SYSSTAT_XL_25	38327	5381290946	0	7240780	140404
TEST_SYSSTAT_XL_26	15520	2033793641	0	1037880	131043
TEST_SYSSTAT_XL_27	17893	1899957992	0	874944	106184
TEST_SYSSTAT_XL_28	38460	7197421036	0	11714880	187140
TEST_SYSSTAT_XL_29	17965	2264343920	0	968704	126041
TEST_SYSSTAT_XL_30	27476	1186765122	0	2357340	43192
TEST_SYSSTAT_XL_31	17941	1951675296	0	507904	108782
TEST_SYSSTAT_XL_32	30008	5608871697	0	9082046	186912
all tables	299603	40143709521	0	11714880	133989

* No indexes used.

** Test performance using 12.2 on i5-quad-core with 8 GB RAM and 1 TB SSD for data 500 GB SSD for TEMP.



Practical test-CASE: query any partition key column values

- If you query any partitioning key column for an invalid value, then you identify this quickly.
- No matter, if you use local PETs or a global PET: the join leads to an invalid partition range on the tables containing data.
- Response times on my test-system (same tables as described earlier): 0.3 to 0.5 seconds.
- **Conclusion: you can easily avoid a full-table-scan!**



Practical test-CASE: leftmost column -query in invalid value

- Query the UNION ALL view with a key-value for the left-most partitioning key that does not exist. You should see an execution plan, that shows no selection of partitions and thus just observe the time to identify, that no valid partitions exist for any of the tables used.

```
select count(*) my_num_rows
from   V_TEST_SYSSTAT_XL...
where  host_name = 'this is not a host' ;
```

- If a UNION ALL view is used, what is the best advise for the PET?

You will definitely check the JOIN-cost of the PET with the individual tables used in the UNION-ALL view versus the UNION-ALL view as such.

Then you decide having local PET tables or a global PET table.
(Hint: you can have both at a time ...)



Practical test-CASE: query any partition key column values

Typical problems, when testing on huge data sets. Right predicates? Existing values used?

```
select count(*)  num_rows
from    V_TEST_SYSSTAT_PET_...
where   ( ( host_name      = 'my.host.20' ) OR
          ( host_name      = 'my.host.21' ) OR
          ( host_name      = 'my.host.22' ) OR
          ( host_name      = 'my.host.23' )      )
and     ( ( dbid           = 200000000024 ) OR
          ( dbid           = 200000000025 ) OR
          ( dbid           = 200000000026 ) OR
          ( dbid           = 200000000027 ) OR
          ( dbid           = 200000000028 )      ) ;
```

Again: no full table scan although we skipped the 2nd partition key column i.e. Instance_name



... and the execution plan is:

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time	Pstart	Pstop	TQ	IN-OUT	PQ Distrib
0	SELECT STATEMENT		1		19M (2)	00:12:48					
1	SORT AGGREGATE		1								
2	TEMP TABLE TRANSFORMATION										
3	LOAD AS SELECT (CURSOR DURATION MEMORY)	SYS_TEMP_0FD9D661A_2317DE6F									
4	TABLE ACCESS FULL	TEST_SYSSTAT_PET_ALL	238K	6512K	138 (2)	00:00:01					
5	PX COORDINATOR										
6	PX SEND QC (RANDOM)	:TQ10002	1						Q1,02	P->S	QC (RAND)
7	SORT AGGREGATE		1						Q1,02	PCWP	
8	VIEW	V_TEST_SYSSTAT_PET_ALL	32		19M (2)	00:12:48			Q1,02	PCWP	
* 9	HASH JOIN		32	10944	19M (2)	00:12:48			Q1,02	PCWP	
10	PX RECEIVE		32	5472	19M (2)	00:12:48			Q1,02	PCWP	
11	PX SEND HASH	:TQ10001	32	5472	19M (2)	00:12:48			Q1,01	P->P	HASH
12	VIEW		32	5472	19M (2)	00:12:48			Q1,01	PCWP	
13	UNION-ALL								Q1,01	PCWP	
14	PX PARTITION LIST INLIST		1	30	2732 (2)	00:00:01	KEY(I)	KEY(I)	Q1,01	PCWC	
* 15	TABLE ACCESS FULL	TEST_SYSSTAT_XL_01	1	30	2732 (2)	00:00:01	KEY(I)	KEY(I)	Q1,01	PCWP	
...											
77	PX PARTITION LIST INLIST		1	31	1174K (3)	00:00:46	KEY(I)	KEY(I)	Q1,01	PCWC	
* 78	TABLE ACCESS FULL	TEST_SYSSTAT_XL_32	1	31	1174K (3)	00:00:46	KEY(I)	KEY(I)	Q1,01	PCWP	
79	BUFFER SORT								Q1,02	PCWC	
80	PX RECEIVE		238K	38M	250 (2)	00:00:01			Q1,02	PCWP	
81	PX SEND HASH	:TQ10000	238K	38M	250 (2)	00:00:01				S->P	HASH
* 82	VIEW		238K	38M	250 (2)	00:00:01					
83	TABLE ACCESS FULL	SYS_TEMP_0FD9D661A_2317DE6F	238K	6512K	250 (2)	00:00:01					

Looks great and takes advantage of predicate push-down.



... and predicate push-down is:

Predicate Information (identified by operation id):

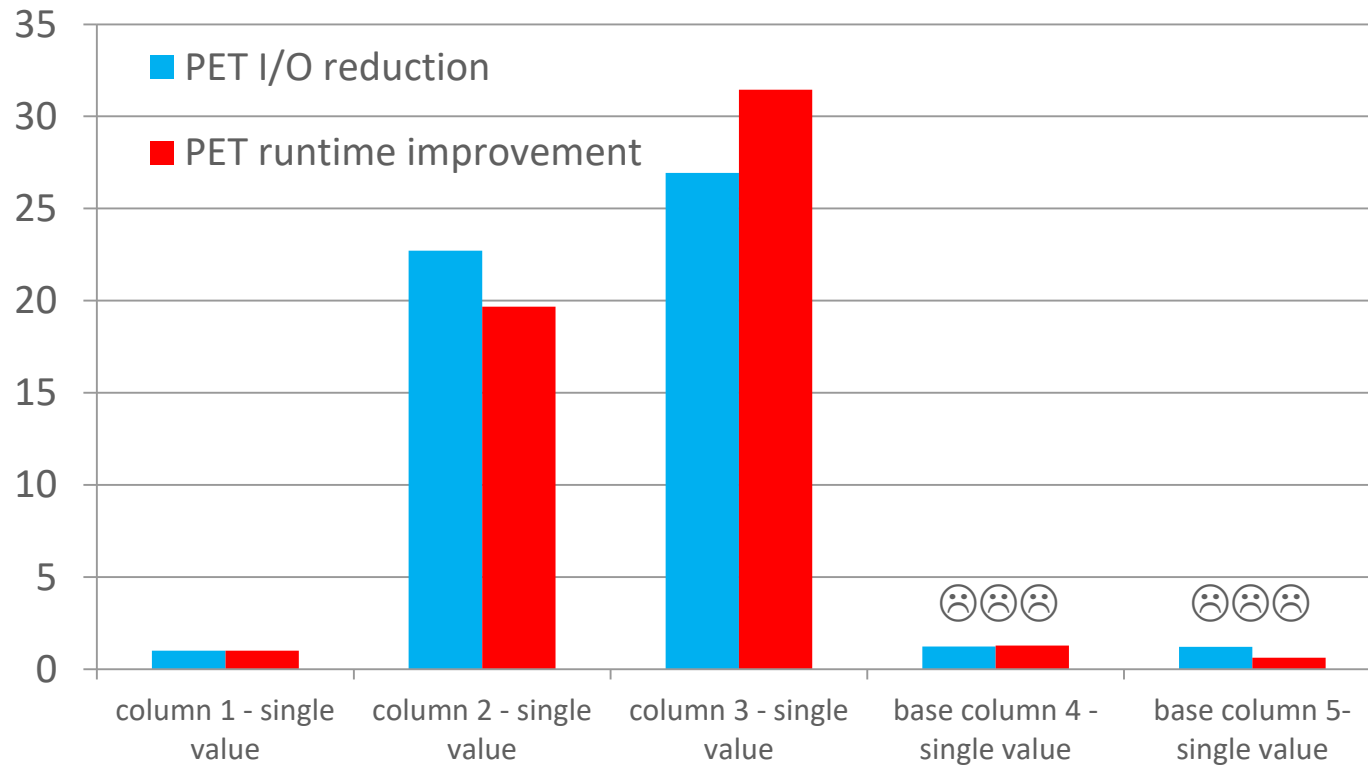
```
-----  
9 - access("XL"."HOST_NAME"="EL"."HOST_NAME" AND "XL"."INSTANCE_NAME"="EL"."INSTANCE_NAME" AND "XL"."DBID"="EL"."DBID" AND  
"XL"."GRP_SNAP_ID"="EL"."GRP_SNAP_ID" AND "XL"."GRP_STMT_ID"="EL"."GRP_STMT_ID")  
15 - filter(("TEST_SYSSTAT_XL_01"."HOST_NAME"='my.host.20' OR "TEST_SYSSTAT_XL_01"."HOST_NAME"='my.host.21' OR  
"TEST_SYSSTAT_XL_01"."HOST_NAME"='my.host.22' OR "TEST_SYSSTAT_XL_01"."HOST_NAME"='my.host.23') AND ("TEST_SYSSTAT_XL_01"."DBID"=20000000024 OR  
"TEST_SYSSTAT_XL_01"."DBID"=20000000025 OR "TEST_SYSSTAT_XL_01"."DBID"=20000000026 OR "TEST_SYSSTAT_XL_01"."DBID"=20000000027 OR  
"TEST_SYSSTAT_XL_01"."DBID"=20000000028))  
...  
78 - filter(("TEST_SYSSTAT_XL_32"."HOST_NAME"='my.host.20' OR "TEST_SYSSTAT_XL_32"."HOST_NAME"='my.host.21' OR  
"TEST_SYSSTAT_XL_32"."HOST_NAME"='my.host.22' OR "TEST_SYSSTAT_XL_32"."HOST_NAME"='my.host.23') AND ("TEST_SYSSTAT_XL_32"."DBID"=20000000024 OR  
"TEST_SYSSTAT_XL_32"."DBID"=20000000025 OR "TEST_SYSSTAT_XL_32"."DBID"=20000000026 OR "TEST_SYSSTAT_XL_32"."DBID"=20000000027 OR  
"TEST_SYSSTAT_XL_32"."DBID"=20000000028))  
82 - filter(("EL"."HOST_NAME"='my.host.20' OR "EL"."HOST_NAME"='my.host.21' OR "EL"."HOST_NAME"='my.host.22' OR "EL"."HOST_NAME"='my.host.23') AND  
("EL"."DBID"=20000000024 OR "EL"."DBID"=20000000025 OR "EL"."DBID"=20000000026 OR "EL"."DBID"=20000000027 OR "EL"."DBID"=20000000028))
```

Looks great because the predicates are pushed down to each of the branches of the union all view.



Overcoming partition pruning limitations – part 1

partitioning-key(host_name, instance_name, dbid, grp_snap_id , grp_stmt_id)



- I/O reduction and runtime are significantly reduced for the UNION ALL view. (predicates are pushed to each branch)
- Predicates on base columns for virtual columns DO NOT lead to partition elimination ☹️☹️☹️

* 32 tables, 40 G records, 300 k partitions, local PETs used. No indexes used.

** Test performance using 12.2 on i5-quad-core with 8 GB RAM and 3 TB HDD for data, 500 GB SSD for TEMP.



Overcoming partition pruning limitations – part 2

partitioning-key(host_name, instance_name, dbid, grp_snap_id , grp_stmt_id)

- 1st attempt failed with:

Add the following predicate to trigger elimination on partition key column grp_snap_id
and grp_snap_id = round(snap_id / 2000, 0);

Add the following predicate to trigger elimination on partition key column grp_stmt_id
and grp_stmt_id = round(stmt_id / 2000, 0);

- 2nd attempt using: :your_value1 and :your_value2 explicitly

Add the following predicate to trigger elimination on partition key column grp_snap_id
and grp_snap_id = round(:your_value1 / 2000, 0);

Add the following predicate to trigger elimination on partition key column grp_stmt_id
and grp_stmt_id = round(:your_value2 / 2000, 0);

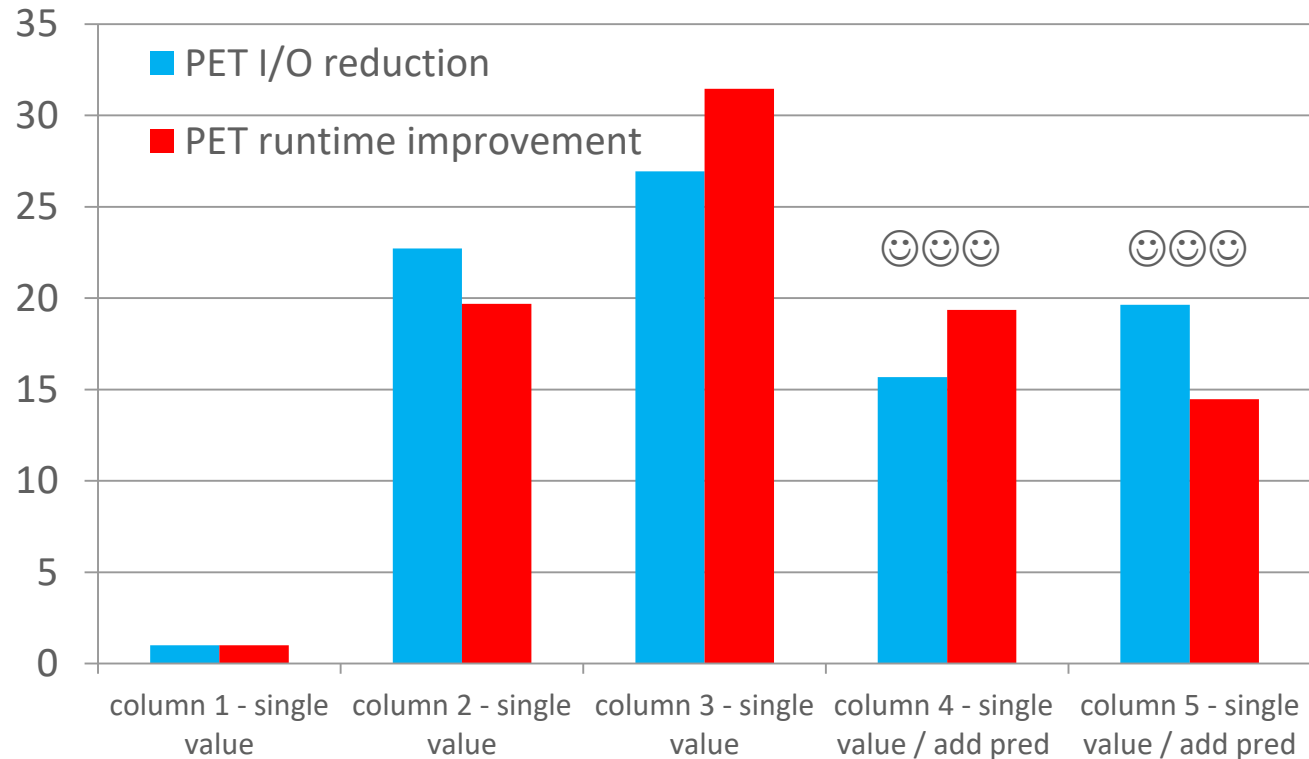
Conclusion:

using a virtual column for partition elimination requires great care, when generating the predicates.



Overcoming partition pruning limitations – part 2

partitioning-key(host_name, instance_name, dbid, grp_snap_id , grp_stmt_id)



- Adding the right predicates helps eliminating I/O and reducing runtime for virtual columns too.

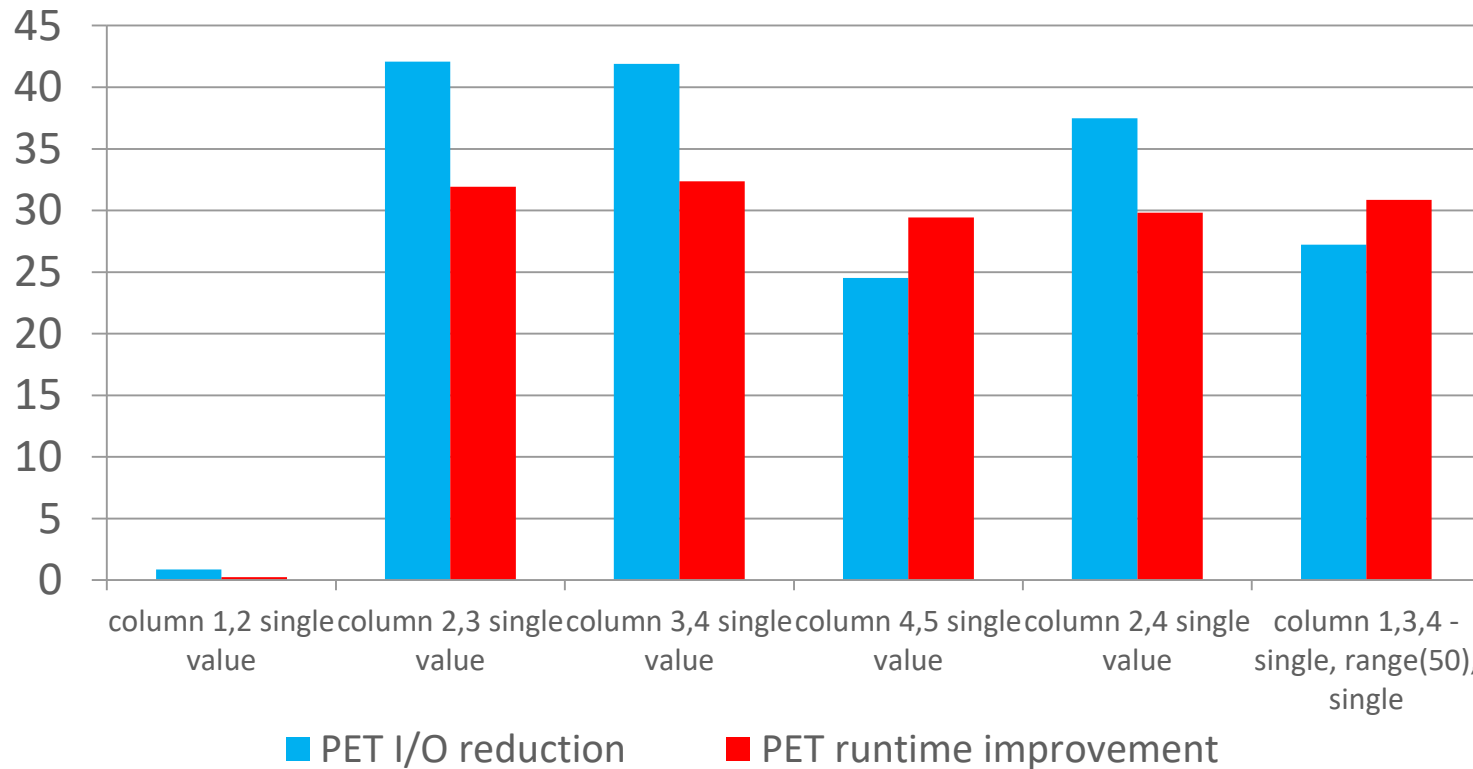
* 32 tables, 40 G records, 300 k partitions, **local PETs** used. No indexes used.

** Test performance using 12.2 on i5-quad-core with 8 GB RAM and 3 TB HDD for data, 500 GB SSD for TEMP.



Overcoming partition pruning limitations – part 3

partitioning-key(host_name, instance_name, dbid, grp_snap_id , grp_stmt_id)



- What happens, if we have predicates on two of the partitioning key columns?
(and in general on any subset thereof)
- Example: use a predicate for two partitioning key columns
- As expected: a significant advantage compared to a full table scan.
- As expected: a combination using the leftmost column 1 is already efficient.

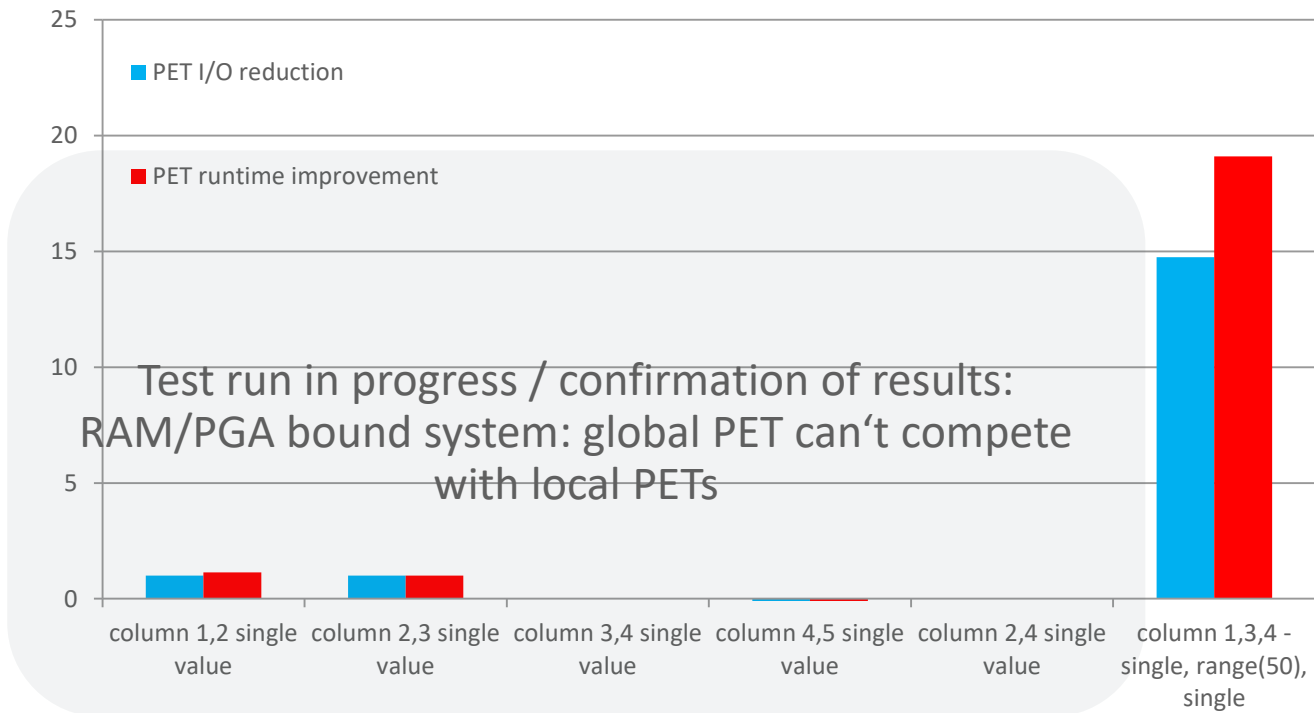
* 32 tables, 40 G records, 300 k partitions, **local PETs** used. No indexes used.

** Test performance using 12.2 on i5-quad-core with 8 GB RAM and 3 TB HDD for data, 500 GB SSD for TEMP.



Overcoming partition pruning limitations – part 3

partitioning-key(host_name, instance_name, dbid, grp_snap_id , grp_stmt_id)



- What happens, if we have predicates on two of the partitioning key columns? (and in general on any subset thereof)
- Example: use a predicate for two partitioning key columns
- As expected: a significant advantage compared to a full table scan.
- As expected: a combination using the leftmost column 1 is already efficient.
- Identifying an empty result set is quick. (column 4,5)

* 32 tables, 40 G records, 300 k partitions, **global PET** used. No indexes used.

** Test performance using 12.2 on i5-quad-core with 8 GB RAM and 3 TB HDD for data, 500 GB SSD for TEMP.



Overcoming partition pruning limitations – part 3

partitioning-key(host_name, instance_name, dbid, grp_snap_id , grp_stmt_id)

- If you have an application, that needs to check for existence of data (working assumption is: very often an empty result set), then *a global PET* can answer the question most efficiently.
- If you have an application, that needs to crawl all tables (working assumption is: always return a significant subset of partition data) then *local PETs* are more efficient.
- In practise: you have both: the local PETs and a global PET. You decide in your application, which access view to use.
- Don't forget: having many partitions works, but object caches requirements become large and PGA consumption very large.



Challenge 1 & 2 & 3 & 4:

Combine all of them



Any open questions?

- Just do it – I've tested it and it works really nicely for Oracle 12.2 using LIST partitioning.
- Just use partitioning key columns being NOT NULL.
You have dimensions having a NULL in its key value – seriously? Having a DEFAULT dimension key is OK!
- You can have up to 16 columns in your partitioning key.
Note: changing partition key column values may lead to ROW MOVEMENT between partitions.
- Keep the PET table(s) updated. Just check for changes to the partitions of data tables.
Note: you might eliminate partition key combinations from your PETs, if a partition becomes empty.
- You you will likely hit the 1 million partition limit per object, unless you have functional dependency of partition key columns (aka sparsity) or use a set of tables to address the challenge.
- Using a set of tables does NOT imply using the same partitioning key.
You can have local PETs having the adequate set of columns and it will still work.
(If you have more than 16 common dimensions: think of snow-flake models – combine dimensions!)
(Think of this being m-dimensional cubes in a much larger n-dimensional cube.)



Known limitations

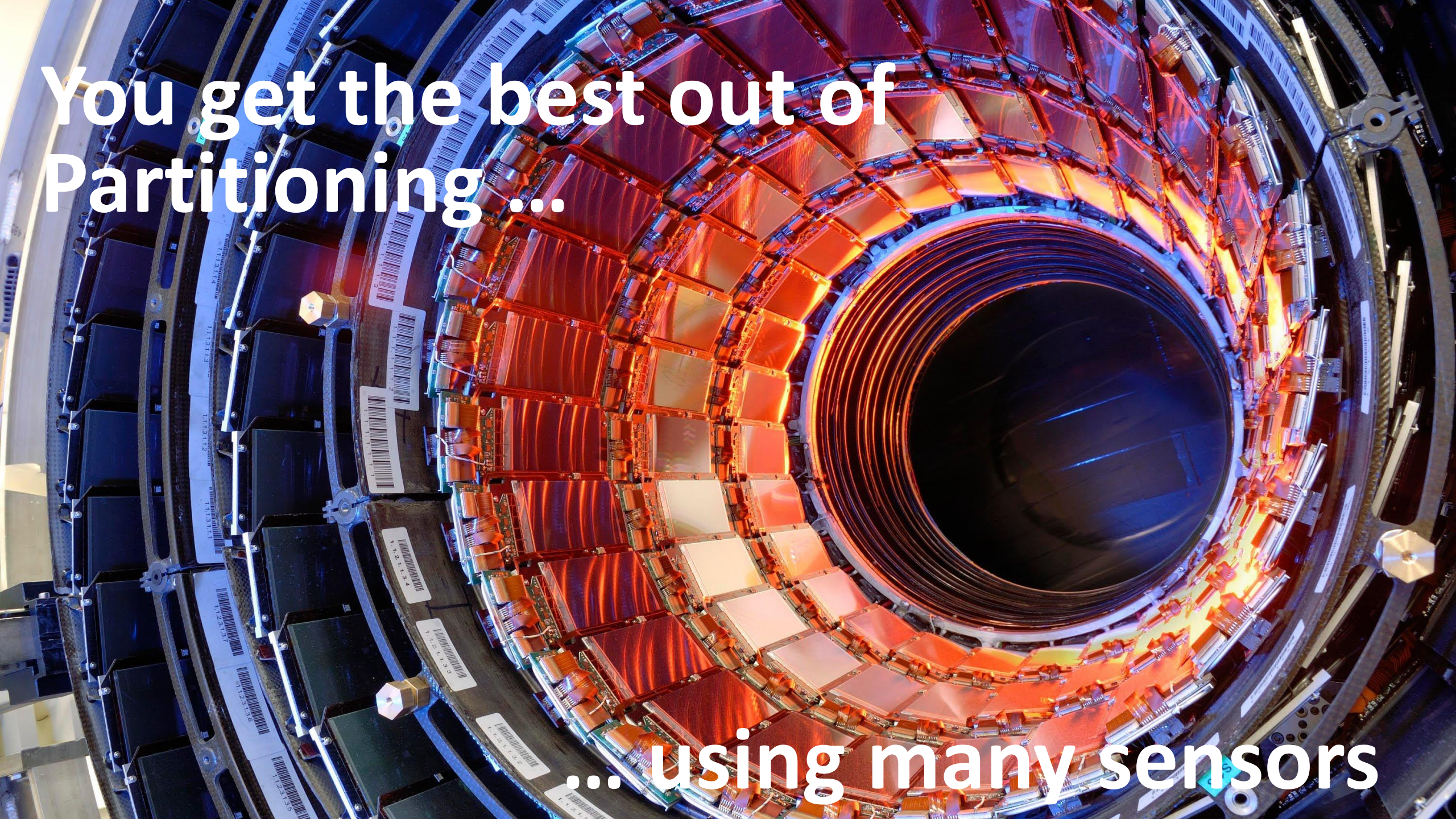
- The method can be extended working with INTERVAL partitioning but is *more complex*.
- The method can be extended working with RANGE partitioning but this is *very complex*.
- The method won't work for HASH partitioning but this *seems impossible without using the hash-function*.



Why can't Oracle perform all this automatically for me?

Let us know, what Oracle can do for you.





You get the best out of
Partitioning ...

... using many sensors

Integrated Cloud

Applications & Platform Services

Safe Harbor Statement

The preceding is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Oracle's products remains at the sole discretion of Oracle.

ORACLE®