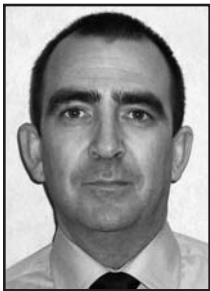




Tuning Parallel Execution



By Doug Burns

The first time I heard of Oracle's Parallel Query Option was in 1993, when my boss returned from the IOUG conference. He was an experienced database guy and I still remember his comment that "this could be as big a step forward as B*tree indexing."



I heard nothing about it for another two or three years, working at various sites as a contractor, and the facility wasn't released until Oracle 7.1. Then I worked at a site where one of the DBAs decided we should give it a try to improve the performance of our overnight batch schedule. The results were disastrous, so it was switched off again because no one had the time to investigate what went wrong. Periodically I would hear about someone trying this wonderful feature and that it had disastrous effects so it was never used as much as Oracle must have hoped.

Parallel Execution Architecture

I think that one of the reasons why Parallel Execution is not used at many sites is that Oracle's basic Instance Architecture works so well. One of the foundations of the architecture is multiple processes running in parallel that use multi-CPU servers effectively. In a typical Production OLTP system there will be a large number of concurrent sessions. In most configurations, each will have its own dedicated server process to communicate with the instance. When all of these server connection processes are considered, the reality is likely to be that you are using your available CPU resource quite effectively, particularly when you consider the various background processes as well!

However, what about those situations when you have a resource intensive job to run, whether it be a long running report query, large data load, or perhaps an index creation? That's where parallel execution can prove extremely useful. It allows the server to take a single large task, break into separate streams of work and pass those streams to parallel execution (PX) slave processes for completion. Because the PX slaves are separate processes (or threads in a Windows environment), the operating system is able to schedule them and provide timely CPU resource in the same way that it would schedule individual user sessions. In fact, each PX slave is just like

a normal dedicated server connection process in many ways, so it's like setting four or eight normal user sessions to work on one problem. Of course, those "users" need to behave in a co-ordinated manner (probably unlike most real users!).

Sensibly, Oracle limits those tasks that can use Parallel Execution to those that are likely to benefit from it. Unless a task is likely to run for a fair amount of time (minutes or hours, rather than seconds) there's not much point in splitting it into parallel streams. The process initialization, synchronization and messaging effort might take longer than the original single-threaded operation! According to the 9.2.0.5 documentation, Oracle supports parallel execution of the following operations:

Access methods

For example, table scans, index full scans, and partitioned index range scans.

Join methods

For example, nested loop, sort merge, hash, and star transformation.

DDL statements

CREATE TABLE AS SELECT, CREATE INDEX, REBUILD INDEX, REBUILD INDEX PARTITION, and MOVE SPLIT COALESCE PARTITION

DML statements

For example, INSERT AS SELECT, updates, deletes, and MERGE operations.

Miscellaneous SQL operations

For example, GROUP BY, NOT IN, SELECT DISTINCT, UNION, UNION ALL, CUBE, and ROLLUP, as well as aggregate and table functions.

In practice, you can simplify the possibilities for the SELECT statements that I'm going to focus on to those that include one of the following:

- Full Table Scans
- Partition Scans (Table or Indexes)

Or, to put it another way, those operations that are likely to process significant volumes of data. There wouldn't be any significant benefits to having multiple processes retrieving a handful of rows via a selective index.

First, let's look at the default single-threaded architecture.

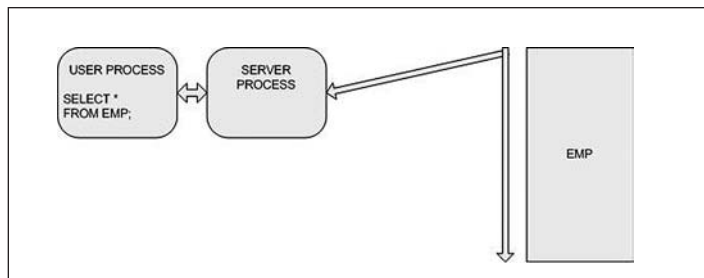


Figure 1: Standard single-threaded architecture using dedicated server process

This should be very familiar. The User Process (on the client or server) submits a SELECT statement that requires a full table scan of the EMP table and the Dedicated Server Process is responsible for retrieving the results and returning them to the User Process.

Let's look at how things change when we enable Parallel Execution.

continued on page 6

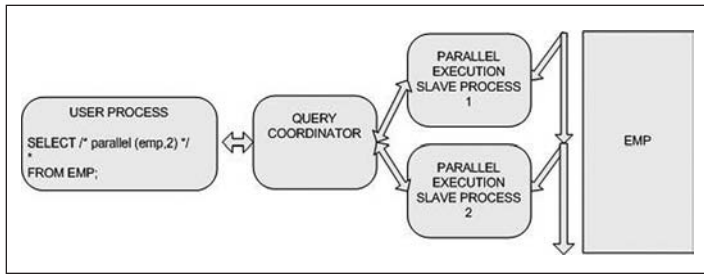


Figure 2: Parallel full table scan – Degree 2

This time, the server is going to process the query in parallel as a result of the optimizer hint. When the server sees that the requested Degree of Parallelism (DOP) for the emp table is two the dedicated server process becomes the Query Coordinator. It makes a request for two PX slaves and, if it's able to acquire them, it will divide all of the blocks that it would have had to scan in the emp table into two equal ranges. Then it will send a SQL statement similar to the following to each of the slave processes.

```
SELECT /*+ Q1000 NO_EXPAND ROWID(A1) */ A1."EMPNO",A1."ENAME",A1."JOB",
A1."MGR",A1."HIREDATE",A1."SAL",A1."COMM",A1."DEPTNO"
FROM
"OPS$ORACLE"."EMP" PX_GRANULE(0, BLOCK_RANGE, DYNAMIC) A1
```

As the data is retrieved from the emp table, it will be returned to the query coordinator which will, in turn, return the data to the user process. The way that all of the data is moved between the processes is using areas of memory called parallel execution message buffers or table queues. These can be stored in either the Shared Pool or Large Pool.

Now let's look at another parallel query.

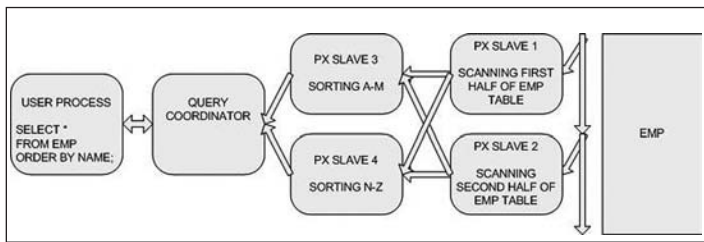


Figure 3: Parallel full table scan with sort – Degree 2

The first thing to note is that there is no PARALLEL hint in the query and yet Oracle chooses to use Parallel Execution to process it. The reason is that the emp table itself also has a parallel DEGREE setting which allows us to specify whether Oracle should attempt to parallelize operations against that table and in this case, it's been set to two:

```
ALTER TABLE emp PARALLEL 2;

SELECT table_name, degree
FROM user_tables
WHERE table_name = 'EMP';
```

TABLE_NAME	DEGREE
EMP	2

Hold on a minute! We requested a DOP of two and yet there are four PX slaves being used to process our request. This is because Oracle will sometimes use two sets of PX slaves if there are intermediate steps in the query plan that could benefit from this. So, in this case, Oracle can see that we are going to have to perform a sort because the NAME column isn't indexed, so it requests four PX slaves. The first set is responsible for scanning the EMP table and the second set for sorting the data. As you can see, though, the type of operation defines the way the workload is distributed between the slaves in each set. For the full table scan, it's based on block ranges. For the sort, each slave takes half of the possible range of values.

```
SELECT A1.C0 C0,A1.C1 C1,A1.C2 C2,A1.C3 C3,A1.C4 C4,A1.C5 C5,A1.C6 C6,A1.C7
C7
FROM
:Q1000 A1 ORDER BY A1.C0
```

There are a couple of very important things to note:

- It is essential that each PX slave in a given set must be able to communicate with all of the slaves in the other set. Even with a DOP of two, you can see this means four inter-slave connections. As the DOP increases, the number of connections will increase exponentially.
- The maximum number of processes required for a given query is 2 x DOP plus the Query Coordinator. If there are multiple step operations in the plan, then the two sets of slaves will be reused. However, this is complicated by the fact that a sub-query might also be executed in parallel.

Finally, it's worth considering what happens when we use Parallel Execution against a Partitioned Table, because Oracle behaves differently.

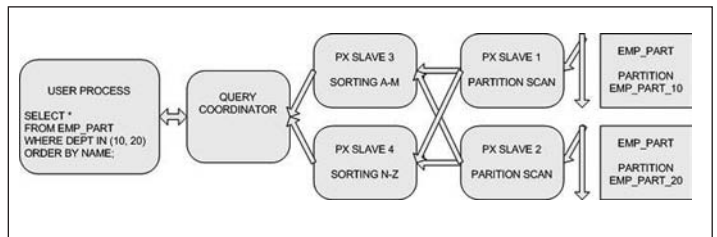


Figure 4: Parallel partition scan with sort – Degree 2

In this case, the maximum DOP that Oracle can use is the number of separate partitions that will be scanned because, instead of using block ranges, Oracle divides the work up by partitions. It's worth noting that this could lead to less than impressive results if you don't have a fairly even distribution of rows in each partition. For example, if EMP_PART_10 contained only 3,000 rows, but EMP_PART_20 contained 35,000, Oracle would still have to wait for PX SLAVE 2 to complete its work before the sort operation could be completed, so PX SLAVE 1 might be idle for much of the query processing. This isn't necessarily a problem, but our query would be processed more efficiently if each partition had 19,000 rows.

Instance Configuration

Switching Parallel Execution on at the Instance Level is surprisingly easy and only requires a few parameters to be set. However, you need to think carefully about the overall effects on the server, so it's worth reading this document and as many of the references listed at the end of this paper as you can before deciding on the values that you think are best for your particular application.

parallel_automatic_tuning

Default Value	FALSE
Recommended Value	TRUE

This parameter was first introduced in Oracle 8i and its very presence is instructive! I remember in Version 6 being able to tune the various areas of the dictionary (or row) cache using the dc_ parameters. When 7 came along, the facility was taken away, and that particular cache became self-tuning. Oracle has attempted much the same at various stages in the development history of the server, with varying degrees of success! Other examples include automatic PGA management and System Managed Undo. To me, this parameter is a sign that users have experienced difficulty in configuring PX themselves so Oracle is trying to make the job easier. In this case, I think they probably have a point. When parallel_automatic_tuning=true, it has several effects.

- The message buffers are stored in the Large Pool rather than the Shared Pool to improve efficiency. However, you need to be aware of this and set large_pool_size accordingly. The calculation for this is in the documentation (see the list of references at the end)
- It sets various parameters to more sensible values than their defaults (e.g. parallel_execution_message_size is increased from 2Kb to 4KB)
- Importantly, it enables the parallel adaptive multi-user algorithm (see next section)

According to Oracle, this parameter is deprecated in 10g as the default values for parallel configuration parameters are optimal. However, when I've tried changing the value, apart from generating a warning, it seems to me that the behavior is the same as previous versions! Perhaps this is one of those situations where Oracle will be making some far-reaching changes in approach in future versions and this is an early sign.

parallel_adaptive_multi_user

Default Value	FALSE
Automatic Tuning Default	TRUE
Recommended Value	Depends

Parallel_adaptive_multi_user is one of the most significant PX configuration parameters and you need to be fully aware of its effects. Imagine a situation where perhaps we have an Oracle Discoverer report running against a Data Warehouse that takes eight minutes to run. So we modify the DOP setting on the relevant tables to see if using PX will improve performance. We find that a DOP of four gives us a run-time of 90 seconds and the users are *extremely* happy. However, to achieve this, the report is using a total of nine server processes. We decide to stress test the change to make sure that it's going to work, don't we? Or maybe we just decide to release this to production because it's such a fantastic improvement! The only difference is going to be whether we have a disastrous test (bearable) or dozens of unhappy users (unbearable).

The problem is that we've just multiplied the user population by nine and, while this worked fantastically well with just one report running, it probably won't scale to large user populations unless you have some extremely powerful hardware. The likelihood is that it won't take long before Oracle manages to suck the very last millisecond of CPU away and the effect on the overall server performance will be very noticeable!

continued on page 8

To give you an example, one day we were testing our overnight batch run that used multiple job streams running in parallel (using the scheduling tool) each of which was parallelized using PX. In effect, we had over a hundred server processes running on a server with a handful of CPUs. It ran quickly enough for our purposes but, while it was running, it was difficult to do anything else on the server at all. It was so slow that it appeared to be dead.

Clearly, the last thing we want is for our server to grind to a halt. The users wouldn't be getting their reports back *at all*, never mind in eight minutes! So Oracle introduced the Adaptive Multi-User facility to address this problem. Effectively, Oracle will judge how busy the server is when it's doling out the PX slaves and, if it decides the machine is too busy, it will give you a smaller degree of parallelism than requested.

Initially this seems like an excellent idea and, on balance, I think it probably is because if the server is absolutely saturated, everyone loses out. However, let's question and be clear on the impact of this. To help me do this, I'm going to use a quote:

"This provides the best of both worlds *and* what users expect from a system. They know that when it is busy, it will run slower."

Effective Oracle by Design. Thomas Kyte

Tom Kyte is a man that I admire very much for all the work he's done for the Oracle community and for his considerable technical and communication skills. It's difficult for me to find anything he has to say about PX that I'd disagree with. All I'm interested here is in different opinions and perspectives, not technical detail.

Ask yourself this question: "Do my users expect the same report to run four or eight times more slowly depending on what else is going on on the server?" I'm not talking about 90 seconds versus 100 seconds, more like 90 seconds against eight minutes, at unpredictable moments (from the point of view of the user). In my opinion, the statement doesn't really reflect what a lot of users are like at all. The one thing they don't want is unpredictable performance. In fact, I just finished working at a site where the managers were very particular about the fact that they wanted a reasonable but, more important, *consistent* level of performance.

Of course, the performance will only vary because you're asking the server to do more than it's capable of so I think Oracle's solution is very sensible and I recommend it. But you must remember the implications and be able to articulate them to your users!

parallel_max_servers and parallel_min_servers

Default Value	Derived from the values of CPU_COUNT, PARALLEL_AUTOMATIC_TUNING and PARALLEL_ADAPTIVE_MULTI_USER
Recommended Value	Completely dependant on number of CPUs – use your initiative?

As Oracle uses PX for user requests, it needs to allocate PX slaves and it does this from a pool of slaves. These two parameters allow you to control the size of the pool and are very straightforward in use. The most difficult thing is to decide on the maximum number of slaves that you think is sensible for your server. I've seen people running dozens or hundreds of slaves on a six CPU server. Clearly that means that each CPU could be trying to cope with 10-20 or more processes and this probably isn't a good idea. Don't forget, though, that this is precisely what your home PC is doing all of the time! However if your disk subsystem is extremely slow, it may be that a number of

slaves per CPU is beneficial because most of your processes are spending most of their time waiting on disk i/o rather than actually doing anything. However, that needs to be balanced against the extra work that the operating system is going to have to do managing the run queue.

A sensible range of values is perhaps two-to-10 times the number of CPUs. The most important thing is to perform some initial stress testing and to monitor CPU and disk usage and the server's run queue carefully.

parallel_threads_per_cpu

Default Value	OS Dependent, but usually 2
Automatic Tuning Default	2
Recommended Value	Increase if i/o-bound, decrease if CPU-bound

This is closely related to the previous parameter. Although it may be worth increasing this from the default of 2, you should be careful that you don't overload the CPUs as a result.

parallel_execution_message_size

Default Value	2Kb
Automatic Tuning Default	4Kb
Recommended Value	4-8Kb

This parameter controls the size of the buffers used to pass messages between the various slaves and the query coordinator. If a message is larger than this size, then it will be passed in multiple pieces, which may have a slight impact on performance. Tellingly, parallel_automatic_tuning increases the size from the default of 2Kb to 4Kb so this is probably a useful starting point, but it may be worth increasing to 8Kb or even larger. Bear in mind, though, that increasing this value will also increase the amount of memory in the Large or Shared Pool, so you should check the sizing calculations in the documentation and increase the relevant parameter appropriately

Other Significant Parameters

In addition to the parallel_ parameters, you should also think about the effect that all of the additional PX slaves will have on your server. For example, each is going to require a process and a session and each is going to be using a sub-task SQL statement which will need to exist in the Shared SQL area. Then we need to think about all of the additional sort areas. The documentation is very good in this area, though, so I'll refer you to that.

Data Dictionary Views

The easiest approach to high-level real-time performance monitoring is to use data dictionary views. There is some information in the standard views, such as V\$SYSSTAT, which I'll come to later. First, though, let's take a look at the PX specific views. These begin with either V\$PQ or V\$PX, reflecting the change in Oracle's terminology over time. Typically, the V\$PX_ views are the more recent and Oracle change the views that are available reasonably frequently so it's always worth using the query below to find out what views are available on the version that you're using.

```
SELECT table_name
FROM dict
WHERE table_name LIKE 'V%PQ%'
OR table_name like 'V%PX%';
```

```
TABLE_NAME
-----
V$PQ_SESSTAT
V$PQ_SYSSTAT
V$PQ_SLAVE
V$PQ_TQSTAT
V$PX_BUFFER_ADVICE
V$PX_SESSION
V$PX_SESSTAT
V$PX_PROCESS
V$PX_PROCESS_SYSSTAT
```

V\$PQ_SESSTAT

V\$PQ_SESSTAT shows you PX statistics for your *current* session.

```
SELECT * FROM v$pq_sesstat;
```

STATISTIC	LAST_QUERY	SESSION_TOTAL
Queries Parallelized	1	2
DML Parallelized	0	0
DDL Parallelized	0	0
DFO Trees	1	2
Server Threads	7	0
Allocation Height	7	0
Allocation Width	1	0
Local Msgs Sent	491	983
Distr Msgs Sent	0	0
Local Msgs Recv'd	491	983
Distr Msgs Recv'd	0	0

It's a nice easy way to confirm that your queries are being parallelized and also gives you a taste of the amount of messaging activity that's required even for a fairly straightforward task.

V\$PQ_SYSSTAT

This view is useful for getting an instance-wide overview of how PX slaves are being used and is particularly helpful in determining possible changes to `parallel_max_servers` and `parallel_min_servers`. For example if "Servers Started" and "Servers Shutdown" were constantly changing, maybe it would be worth increasing `parallel_min_servers` to reduce this activity.

V\$PX_PROCESS_SYSSTAT contains similar information.

```
SELECT * FROM v$pq_sysstat WHERE statistic LIKE 'Servers%';
```

STATISTIC	VALUE
Servers Busy	0
Servers Idle	0
Servers Highwater	3
Server Sessions	3
Servers Started	3
Servers Shutdown	3
Servers Cleaned Up	0

V\$PQ_SLAVE and V\$PX_PROCESS

These two views allow us to track whether individual slaves are in use or not and track down their associated session details.

```
SELECT * FROM v$px_process;
```

SERV STATUS	PID SPID	SID	SERIAL#
P001 IN USE	18 7680	144	17
P004 IN USE	20 7972	146	11
P005 IN USE	21 8040	148	25
P000 IN USE	16 7628	150	16
P006 IN USE	24 8100	151	66
P003 IN USE	19 7896	152	30
P007 AVAILABLE	25 5804		
P002 AVAILABLE	12 6772		

V\$PQ_TQSTAT

V\$PQ_TQSTAT shows you table queue statistics for the current session and you must have used parallel execution in the current session for this view to be accessible. I like the way that it shows the relationships between slaves and the query coordinator very effectively. For example, after running this query against the 25,481 row attendance table:

```
SELECT /*+ PARALLEL (attendance, 4) */ *
FROM attendance;
```

The contents of V\$PQ_SYSSTAT look like this:

```
SELECT dfo_number, tq_id, server_type, process, num_rows, bytes
FROM v$pq_tqstat
ORDER BY dfo_number DESC, tq_id, server_type, process;
```

DFO_NUMBER	TQ_ID	SERVER_TYP	PROCESS	NUM_ROWS	BYTES
1	0	Consumer	QC	25481	443612
1	0	Producer	P000	6605	114616
1	0	Producer	P001	6102	105653
1	0	Producer	P002	6251	110311
1	0	Producer	P003	6523	113032

We can see here that four slave processes have been used acting as row Producers, each processing approximately 25% of the rows, which are all consumed by the QC to return the results to the user. Whereas for the following query:

```
SELECT /*+ PARALLEL (attendance, 4) */ *
FROM attendance
ORDER BY amount_paid;
```

We'll see something more like this.

continued on page 10

```
SELECT dfo_number, tq_id, server_type, process, num_rows, bytes
FROM v$px_tqstat
ORDER BY dfo_number DESC, tq_id, server_type, process;
```

DFO_NUMBER	TQ_ID	SERVER_TYP	PROCESS	NUM_ROWS	BYTES
1	0	Consumer	P000	15351	261380
1	0	Consumer	P001	10129	182281
1	0	Consumer	P002	0	103
1	0	Consumer	P003	1	120
1	0	Producer	P004	5744	100069
1	0	Producer	P005	6304	110167
1	0	Producer	P006	6303	109696
1	0	Producer	P007	7130	124060
1	0	Ranger	QC	372	13322
1	1	Consumer	QC	25481	443612
1	1	Producer	P000	15351	261317
1	1	Producer	P001	10129	182238
1	1	Producer	P002	0	20
1	1	Producer	P003	1	37

There are a few new things going on here:

- P004, P005, P006 and P007 are scanning 25% of the blocks each.
- The QC acts as a Ranger, which works out the range of values that each PX slave should be responsible for sorting.
- P0001, P002, P003 and P004 act as Consumers of the rows being produced by P004-P007 and perform the sorting activity.
- They also act as Producers of the final sorted results, for the QC to consume.

What is a little worrying from a performance point of view is that P000 and P001 seem to be doing a lot more work than P002 and P003, which means that they will run for much longer, and we're not getting the full benefit of a degree 4 parallel sort. It's a good idea to look at the range of values contained in the sort column.

```
SELECT amount_paid, COUNT(*)
FROM attendance
GROUP BY amount_paid
ORDER BY amount_paid
/
```

AMOUNT_PAID	COUNT(*)
200	1
850	1
900	1
1000	7
1150	1
1200	15340
1995	10129
4000	1

This indicates where the problem lies. We have extremely skewed data because the vast majority of rows have one of only two values, so it's very difficult to parallelize a sort on this column.

V\$PX_BUFFER_ADVICE

```
SELECT * FROM v$px_buffer_advice;
```

STATISTIC	VALUE
Servers Highwater	8
Buffers HWM	76
Estimated Buffers HWM	72
Servers Max	20
Estimated Buffers Max	360
Buffers Current Free	51
Buffers Current Total	96

This view is new to Oracle Database 10g and is intended to help us work out the amount of memory that should be added to the large or shared pool by multiplying the estimated maximum number of buffers by the message size.

V\$PX_SESSTAT

This view is a bit like V\$SESSTAT but also includes information about which QC and which Slave Set each session belongs to, which allows us to see a given statistic (e.g., Physical Reads) for all steps of an operation.

```
SELECT stat.qcsid, stat.server_set, stat.server#, nam.name, stat.value
FROM v$px_sesstat stat, v$statname nam
WHERE stat.statistic# = nam.statistic#
AND nam.name LIKE 'physical reads%'
ORDER BY 1,2,3
```

QCSID	SERVER_SET	SERVER#	NAME	VALUE
145	1	1	physical reads	0
145	1	2	physical reads	0
145	1	3	physical reads	0
145	2	1	physical reads	63
145	2	2	physical reads	56
145	2	3	physical reads	61

Monitoring the Parallel Adaptive Multi-User Algorithm

If you are using the Parallel Adaptive Multi-User algorithm, it's vital that you are able to check whether any particular operations have been severely downgraded because the server is too busy. There are additional statistics in V\$SYSSTAT that show this information.

```
SELECT name, value FROM v$sysstat WHERE name LIKE 'Parallel%'
```

NAME	VALUE
Parallel operations not downgraded	546353
Parallel operations downgraded to serial	432
Parallel operations downgraded 75 to 99 pct	790
Parallel operations downgraded 50 to 75 pct	1454
Parallel operations downgraded 25 to 50 pct	7654
Parallel operations downgraded 1 to 25 pct	11873

Clearly, you should be most concerned about any operations that have been downgraded to serial as these may be running many times more slowly than you'd expect. This information is also available in a STATSPACK report, so it's easy to get a view over a period of time. Unfortunately the name column is truncated in the report, which makes it a little difficult to read, but you soon get used to this.

Monitoring the SQL Being Executed by Slaves

As with most dictionary views, we can write queries that combine them to show us interesting or useful information. To offer just one small example, this query will show us the SQL statements that are being executed by active PX slaves. (N.B. The slave must be active, otherwise the SID and SERIAL# it was previously associated with is not contained in the v\$px_process vie.)

```
set pages 0
column sql_test format a60

select p.server_name,
       sql.sql_text
from v$px_process p, v$sql sql, v$session s
WHERE p.sid = s.sid
and p.serial# = s.serial#
and s.sql_address = sql.address
and s.sql_hash_value = sql.hash_value
/
```

Even more interestingly, you'll see completely different results if you run this query on Oracle Database 10g than on previous versions. First some example results from Oracle 9.2:

```
P001 SELECT A1.C0 C0,A1.C1 C1,A1.C2 C2,A1.C3 C3,A1.C4 C4,A1.C5 C5,
A1.C6 C6,A1.C7 C7 FROM :Q3000 A1 ORDER BY A1.C0
```

Whereas on 10g the results look like this:

```
P000 SELECT /*+ PARALLEL (attendance, 2) */ * FROM attendance ORD
ER BY amount_paid

P003 SELECT /*+ PARALLEL (attendance, 2) */ * FROM attendance ORD
ER BY amount_paid

P002 SELECT /*+ PARALLEL (attendance, 2) */ * FROM attendance ORD
ER BY amount_paid

P001 SELECT /*+ PARALLEL (attendance, 2) */ * FROM attendance ORD
ER BY amount_paid
```

This is an example of a more general change in 10g. When tracing or monitoring the PX slaves, the originating SQL statement is returned, rather than a block range query as shown earlier in this document. I think this makes it much easier to see at a glance what a particular long-running slave is really doing, rather than having to tie it back to the QC as on previous versions.

Session Tracing and Wait Events

Tracing an application that uses parallel execution is a little more complicated than tracing non-parallel statements in a few ways.

- A trace file will be generated for each slave process as well as for the query coordinator.
- As I've just mentioned, it's time consuming prior to 10g to identify precisely what application operation a PX slave is involved in processing.
- The trace files for the slave processes may be created in background_dump_dest, rather than the standard user_dump_dest. This is version dependant and the trace file for the query coordinator will be in user_dump_dest in any case
- As a result of all the synchronisation and message passing that occurs between the different processes, there a number of additional wait events.

You won't find too much specific information about tracing Parallel Execution because it's based on exactly the same principles as standard tracing, with the few differences mentioned above. The biggest problem tends to be the large number of trace files that you have to analyze!

Parallel-specific Wait Events

The first thing to get used to when monitoring PX, whether it be using STATSPACK at the high level or event tracing at the low level, is that you are going to see a lot more wait events, including types that you won't have seen before. Here are some of the Parallel-specific wait events.

Events indicating Consumers re waiting for data from Producers:

- PX Deq: Execute Reply
- PX Deq: Table Q Normal

Oracle's documentation states that these are idle events because they indicate the normal behavior of a process waiting for another process to do its work, so it's easy to ignore them. However, if you have excessive wait times on these events it could indicate a problem in the slaves. To give you a real-world example, here is the top timed events section of a STATSPACK report from a production system I worked on.

Event	Waits	Timeouts	Time (s)	(ms)	/txn
direct Path read	2,249,666	0	115,813	51	25.5
PX Deq: Execute Reply	553,797	22,006	75,910	137	6.3
PX qref latch	77,461	39,676	42,257	546	0.9
library cache pin	27,877	10,404	31,422	1127	0.3
db file scattered read	1,048,135	0	25,144	24	11.9

The absolute times aren't important here, just the events. First, it's worth knowing that PX slaves perform direct path reads rather than db file scattered reads. You may already be used to direct path reads because they're used with temporary segments for example. On this system, which was a European-wide Data Warehouse, we were performing long-running SELECT statements as part of the overnight batch run, so a high level of disk I/O was inevitable. (Whether an average wait time of 51 minutes is acceptable when you've spent a small fortune on a Hitachi SAN is another matter!)

The next event is PX Deq: Execute Reply, which Oracle considers to be an idle event, as I've mentioned. So we ignore that and move down to the next event. The PX qref latch event can often mean that the Producers are producing data quicker than the Consumers can consume it. On this particular system, very high degrees of parallelism were being used during an overnight batch run so a great deal of messaging was going on. Maybe we could increase parallel_execution_message_size to try to eliminate some of these waits or we might decrease the DOP.

But the real problem that we were able to solve was the next event—library cache pin. This event represents Oracle trying to load code into the Library Cache so you wouldn't normally expect to see a significant percentage of wait time for this event unless the Shared Pool is really struggling (which it was on this system).

So next we drill down and start to try session tracing to establish the source of these events. Initially I was unsuccessful in tracking them down until I realized that the PX Deq: Execute Reply was a useful hint. The fact is that many of these wait events were happening in the PX slaves and many of the PX Deq: Execute Reply events were caused by the QC waiting for the PX slaves, which were waiting for the library cache pin latch! So sometimes idle events are important.

continued on page 12

Eventually it turned out to be a pretty bad bug in earlier versions of 9.2 (fixed in 9.2.0.5) that caused some of our two-minute SQL statements to occasionally take two hours to run. (Yes, that really does say two hours.) Anyway, back to more wait events.

Events indicating producers are quicker than consumers (or QC):

- PX qref latch

I've found that PX qref latch is one of the events that a system can spend a lot of time waiting on when using Parallel Execution extensively (as you can see from the earlier STATSPACK example). Oracle suggest that you could try to increase `parallel_execution_message_size` as this might reduce the communications overhead, but this could make things worse if the consumer is just taking time to process the incoming data.

Synchronization message events:

- PX Deq Credit: need buffer
- PX Deq: Signal Ack
- PX Deq: Join Ack

Although you will see a lot of waits on these synchronization events—the slaves an QC need to communicate with each other—the time spent should not be a problem. If it is, perhaps you have an extremely busy server that is struggling to cope and reducing the Degree of Parallelism and `parallel_max_servers` would be the best approach.

Query Coordinator waiting for the slaves to parse their SQL statements:

- PX Deq: Parse Reply

Long waits on this event would tend to indicate problems with the Shared Pool as the slaves are being delayed while trying to parse their individual SQL statements. (Indeed, this was the event I would have expected to see as a result of the bug I was talking about earlier but the library cache pin waits were appearing in the Execute phase of the PX slave's work.) Again, the best approach is to examine the trace files of the PX slaves and track down the problem there.

Partial Message Event:

- PX Deq: Msg Fragment

This event indicates that `parallel_execution_message_size` may be too small. Maybe the rows that are being passed between the processes are particularly long and the messages are being broken up into multiple fragments. It's worth experimenting with message size increases to reduce or eliminate the impact of this.

Some Common Sense

One of my favorite descriptions of performance tuning, although I can't remember where I first heard it, is that it is based on "informed common sense." That really captures my own experiences of performance tuning. Yes, you need to use proper analysis techniques and often a great deal of technical knowledge, but that's all devalued if you're *completely missing the point*. So let's take a step away from the technical and consider the big picture.

- Don't even think about implementing Parallel Execution unless you are prepared to invest some time in initial testing, followed by ongoing performance monitoring. If you don't, you might one day hit performance problems either server-wide or on an individual user session that you'd never believe (until it happens to you).

- Parallel Execution is designed to use hardware as heavily as possible. If you are running on a single-CPU server with two hard disk drives and 512Mb RAM, don't expect significant performance improvements just because you switch PX on. The more CPUs, disk drives, controllers, and RAM you have installed on your server, the better the results are going to be.
- Although you may be able to use Parallel Execution to make an inefficient SQL statement run many times faster, that would be incredibly stupid. It's essential that you tune the SQL first. In the end, doing more work than you should be, but more quickly, is still doing more work than you should be! To put it another way, don't use PX as a dressing for a poorly designed application. Reduce the workload to the minimum needed to achieve the task and then start using the server facilities to make it run as quickly as possible. Seems obvious, doesn't it?
- If you try to use PX to benefit a large number of users performing online queries you may eventually bring the server to its knees. Well, maybe not if you use the Adaptive Multi-User algorithm, but then it's essential that both you and, more important, your users understand that response time is going to be very variable when the machine gets busy.
- Using PX for a query that runs in a few seconds is pointless. You're just going to use more resources on the server for very little improvement in the run time of the query. It might well run more slowly!
- The slower your I/O sub-system, the more benefit you are likely to see from PX—but shouldn't you fix the real problem?
- Consider whether PX is the correct parallel solution for overnight batch operations. It may be that you can achieve better performance using multiple streams of jobs, each single-threaded, or maybe you would be better with one stream of jobs which uses PX. It depends on your application so the only sure way to find out is to *try the different approaches*.

Conclusion

Oracle's Parallel Execution capability can improve the performance of long-running tasks significantly by breaking the tasks into smaller sub-tasks that can execute in parallel. The intent is to use as much hardware resource as possible to deliver results more quickly. However, it works best:

- On a server which has spare CPU, RAM and i/o throughput capacity;
- For tasks which run for more than a few seconds;
- For a limited number of concurrent users.

If you can meet all of these requirements then the performance improvements can be dramatic but you should consider the potential downsides carefully:

- Tracing sessions becomes more difficult, although things are supposed to become easier with 10g;
- Unless you are using the Adaptive Multi-user facility you may find your server grinding to a halt one day;
- If you are using the Adaptive Multi-user facility you may find one or more user sessions slowing down dramatically under heavy server workloads.

As with many aspects of Oracle, it's important to plan an effective implementation and test it as thoroughly as possible before inflicting it on your users but when used appropriately, parallel execution is hard to beat.

Bibliography and Resources

The best source of information on Parallel Execution is the Oracle documentation. It's amazing how often I find the (free) manuals far superior to (paid for) third-party books! Specifically, the Data Warehousing Guide contains a couple of relevant chapters:

Using Parallel Execution

http://download-west.oracle.com/docs/cd/B10501_01/server.920/a96520/tuningpe.htm#19664

Parallelism and Partitioning in Data Warehouses

http://download-west.oracle.com/docs/cd/B10501_01/server.920/a96520/parpart.htm#745

There are also a number of useful resources on Metalink, including a specific section containing the most useful notes that you can access by selecting "Top Tech Docs," "Database," "Performance and Scalability," and then "Parallel Execution" from your Metalink home page. The following documents are particularly relevant to this paper:

- 184417** – Where to track down the information on Parallel Execution in the Oracle documentation!
- 203238** – Summary of how Parallel Execution works.
- 119103** – Parallel Execution Wait Events (contains links to event-specific information).
- 201799** – init.ora parameters for Parallel Execution.
- 240762** - pqstat PL/SQL procedure to help monitor all of the PX slaves running on the instance or for one user session.
- 202219** – Script to map PX slaves to Query Coordinator (An alternative to using the procedure in note 240762.1).
- 275240** – Discusses Bug no. 2533038 in the Parallel Adaptive Multi-User algorithm (fixed in 9.2.0.3) which makes it sensitive to idle sessions, leading to an under-utilised system.
- 238680** – Investigating ORA-4031 errors caused by lack of memory for queues.
- 242374** – Discusses some of the issues around PX session tracing (but not in any great depth).
- 237328** – Direct Path Reads (brief discussion).

The following books contain some useful information about Parallel Execution.

Harrison, Guy. *Oracle SQL High Performance Tuning*. Prentice Hall.

Kyte, Thomas. *Effective Oracle by Design*. Oracle Press.

Lewis, Jonathan. *Practical Oracle 8i – Building Efficient Databases*. Addison Wesley.

Mahapatra, Tushar and Mishra, Sanjay. *Oracle Parallel Processing*. O'Reilly & Associates, Inc.

There's also an excellent Parallel Execution conference paper by Jeff Maresh. As well as covering server configuration it contains more developer-orientated information about using PX within your application than I've covered in this document.

Maresh, Jeff. *Parallel Execution Facility Configuration and Use*.
http://www.evdbt.com/PX_2003.doc

And, as usual Jonathan Lewis has some previous articles on this subject! Although these are geared to Oracle 7.3, much of the content makes sense across versions.

Lewis, Jonathan. http://www.jlcomp.demon.co.uk/ind_pqo.html

Finally, here's a nice little PX slave monitoring query

http://www.jlcomp.demon.co.uk/faq/pq_proc.html



About the Author

Doug Burns is an independent consultant who has 14 years experience working with Oracle in a range of industries and applications and has worked as a course instructor and technical editor for both Oracle UK and Learning Tree International. He can be contacted at dougburns@yahoo.com, and this document and other articles are available on his Web site at <http://doug.burns.tripod.com>.