

May 6-10, 2007
San Jose Convention Center
San Jose, California, USA

Session: A04

Troubleshooting Tuning for Dynamic SQL

IDUG® 2007
North America

Thomas Baumann
Swiss Mobiliar

May 07, 2007 04:20 p.m. – 05:20 p.m.

Platform: DB2 for z/OS



GoFurther



This presentation describes a methodical approach to database tuning, which has been developed, tested and used for emergency performance tuning after an application's rollout, but has also been shown to be quite helpful for routine database monitoring. We will discuss answers to the following three questions derived from commands, metrics and tools available in each shop:

- a) Are the application's critical queries being served in the most effective manner?
- b) Is DB2 making optimal use of resources so as to reach the desired performance levels?
- c) Are there enough primary resources available for DB2's consumption and are they configured adequately, given the current workload?

Objectives

- Learn to apply a methodical tuning approach
- Hear about new performance indicators other than cache hit ratios etc.
- Find out what is specific for *dynamic* SQL tuning
- Look at new ways to trigger runstats based on your queries' performance numbers
- Detect the dynamic SQL statement cache's secrets

GoFurther

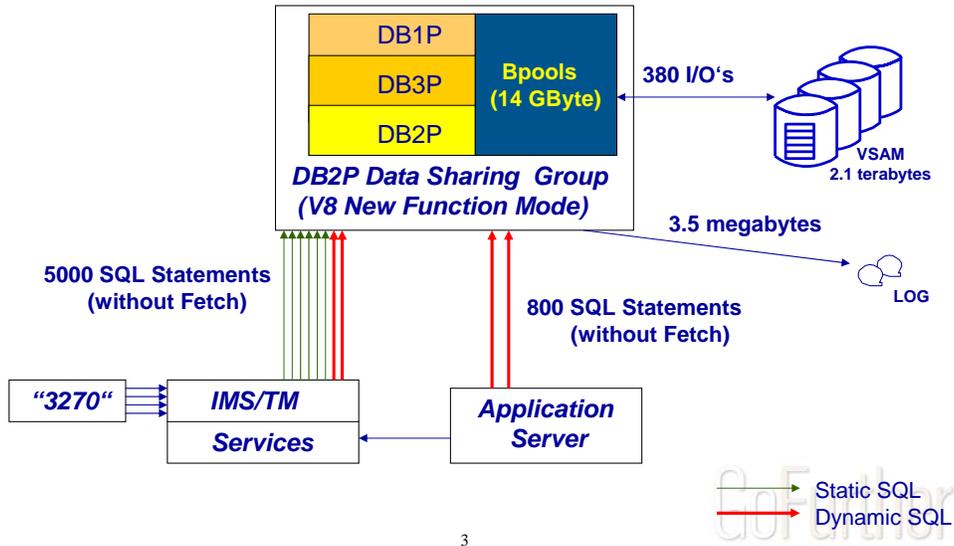
2

First, we will discuss how to identify those 'queries from hell' which could potentially monopolize your whole system. Then – with a finer granularity – we will discuss metrics and techniques to further optimize both your dynamic queries as your system as a whole.

Finally, we will discuss the question of having enough system resources or if there is a need to increase your system's overall size.

One second in the life of Mobilair's DB2

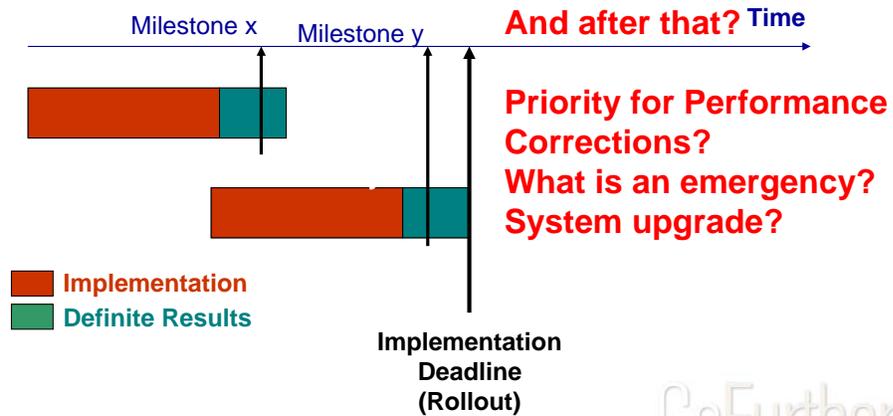
(2006-07-18 @ 10:30 AM)



This slide shows the workload at Swiss Mobilair's OLTP environment as measured 07/18/2006.

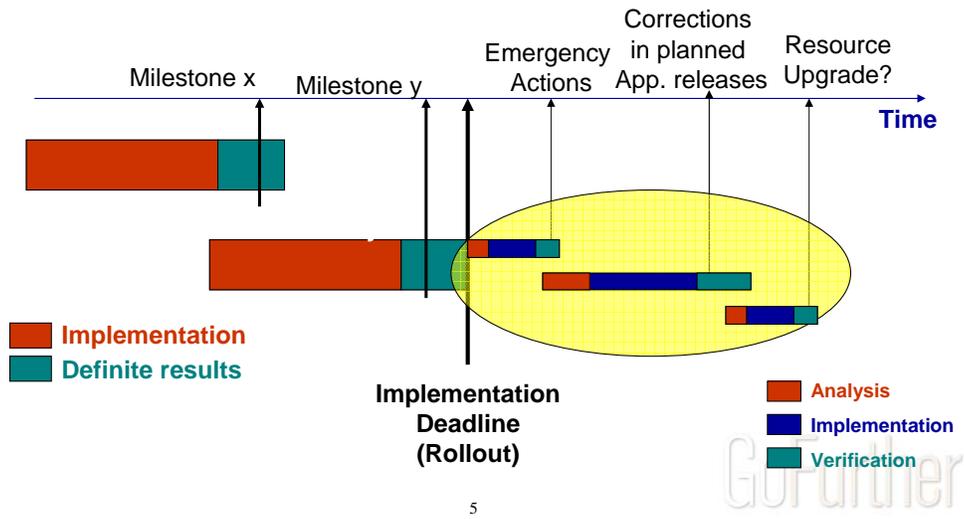
Methodology for Project Consolidation

MvM (Mobilair Action Methodology):



After the application's rollout, the structured project methodology ends. During project consolidation, it is often hard to get sufficient resources to fix performance problems, because priority is often given to correct functional deficiencies only.

Methodology for Project Consolidation



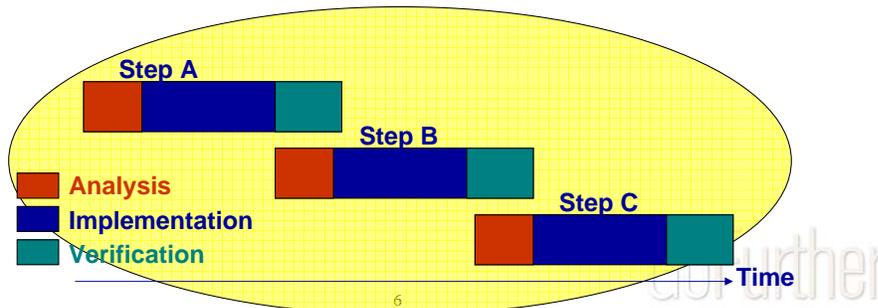
5

The methodology presented in this paper starts here: It defines milestones with performance metrics that need to be met. Milestones with a scope ranging from a single SQL statement to the overall performance of the application and the database subsystem. If the the criterions of the first step were missed, an emergency corrective action is automatically scheduled.

Befehle, Metriken, Werkzeuge for SQL tuning

Methodology for SQL Troubleshooting Tuning

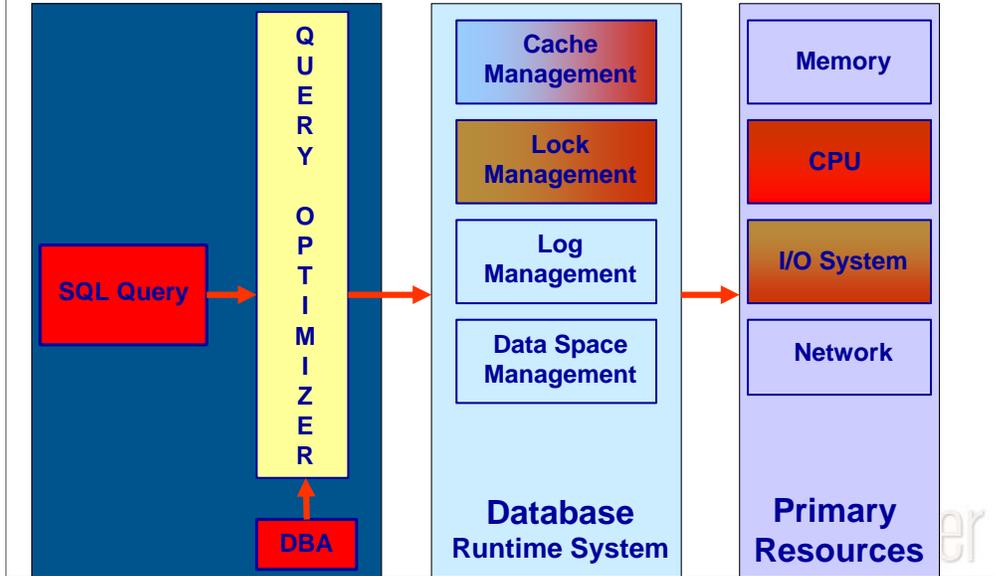
- **Step A: Are there extremely critical SQL queries?**
- Step B: Is the workload being adequately supported by the database system?
- Step C: Are enough system resources (CPU, Memory, ...) available?



This BMW troubleshooting (or emergency database performance) tuning methodology is divided in three steps: In each of these steps, there is one single question which is analyzed in full detail, and the problems detected during this step must be resolved before entering into the next step.

Step 1 detects extremely inefficient queries (*,queries from hell‘*). Once they are corrected, the next steps is initiated: Different thresholds answer the question *‘are my queries served the optimal way?’*. This might concern both SQL queries and database runtime system features and parameters. Only in the third step, after all the problems detected in the first two steps are fully resolved, the question *‘Do my queries get sufficient resources ?‘* will be asked.

BMW step A:
Are there extremely critical SQL queries?



In a system with single very inefficient SQL statements, these might be the typical symptoms: A very high CPU usage rate, an I/O system close at the limits, bufferpools which do not reach their performance goals, and some locking problems.

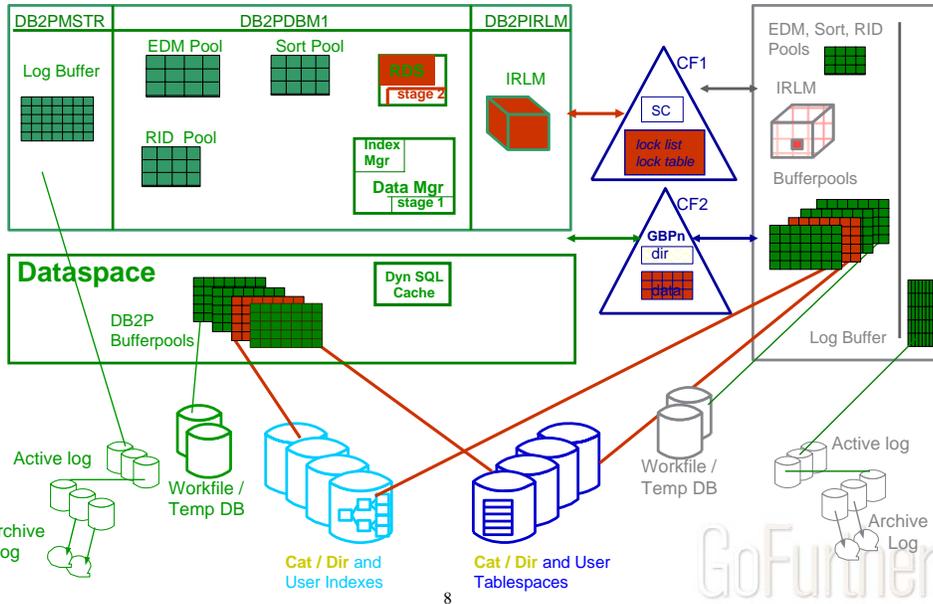
All these 'trouble areas' are marked in red colour – the more red the more urgent they want your attention.

And all caused by a few inefficient SQL statements. That is why we start investigating this system from left to right, from SQL queries to basic resources and not vice versa.

The Big Picture

DB2P

DB3P



This slide contains exactly the same information as the one before – but is more understandable for the programmer, whereas the slide before was the ‘management summary’ of performance problems.

But once again – our basic problem is still the inefficient dynamic SQL statement still unknown and not yet detected.

A1: Queries From Hell

Measurement criterion A1: Extremely high CPU usage

Identify the CPU resource-intensive Queries with

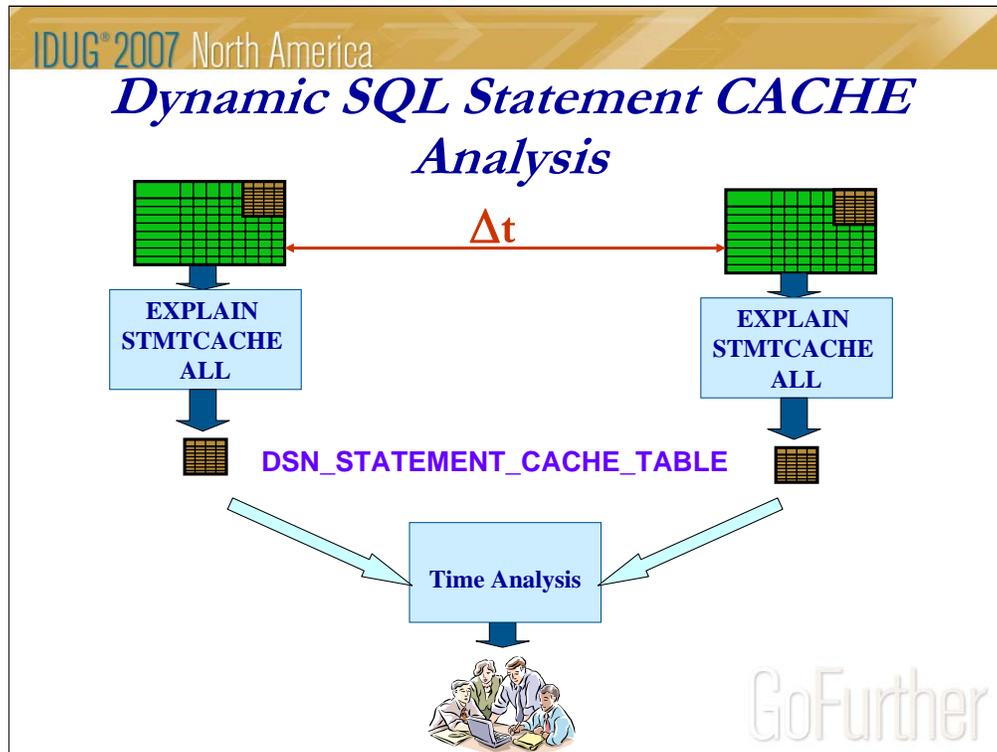
$$CPU / SQL > 0.5 \text{ sec} \ \& \ \sum_h CPU / SQL > 5 \text{ min}$$

With static SQL, there are more stringent limits (for all measurement criteria):
The upper boundary of permissible resource usage for dynamic SQL is 12 times
higher than that which is permissible for static SQL:

A1: CPU/SQL > 0.5 sec & \sum CPU/SQL > 10 min

GoFurther

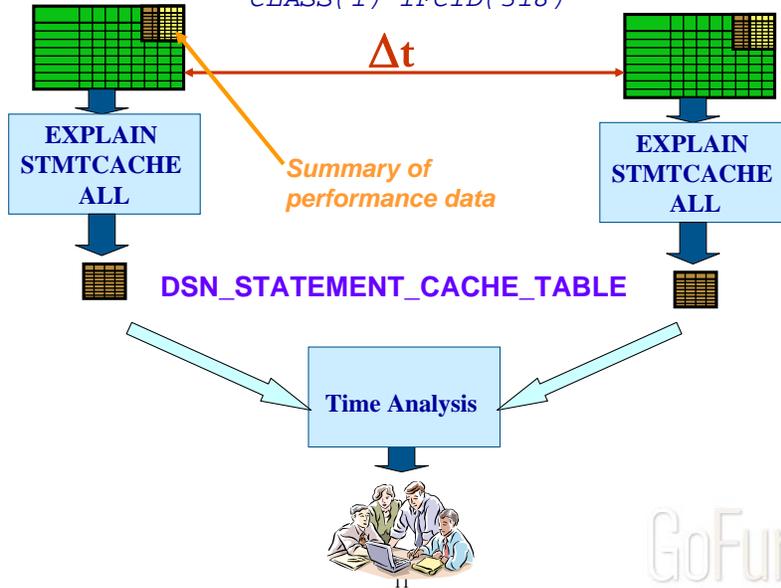
This very first criterion identifies the most inefficient dynamic SQL statements the ones which should be attacked at first – the *queries from hell*.



Analysis of the global dynamic statement cache has become easy in DB2 for z/OS V8: Simply define the `dsn_statement_cache_table`, and populate it by executing `EXPLAIN STMTCACHE ALL`.

If you repeat this command a certain time later you can subtract the performance numbers from the statements with the same `stmt_id`, thus producing a result of all queries executed between these two measurements. Make sure to have this time interval small enough in order to capture all statements.

Dynamic Stmt cache with `-STA TRACE(MON)` `CLASS(1) IFCID(318)`



“The collection of statistics for statements in the dynamic statement cache can increase the processing cost for those statements. When IFCID 0318 is inactive, DB2 tracks the statements in the dynamic statement cache, but does not accumulate the statistics as those statements are used. When you are not actively monitoring the cache, you should turn off the trace for IFCID 0318.”

A2: High Number of Timeouts/Deadlocks

Measurement criterion A2: Timeout / Deadlock

DSNMSTR address space: Identify tables with

$$\sum_h (\text{timeouts or deadlocks}) > 5$$

Analyze queries against these tables with DSPC

GoFurther

Another objective of emergency tuning is the detection of statements which cause timeouts and deadlocks.

A3: Missing Parameter Markers

Measurement criterion A3: Efficient Use of Dynamic Statement Cache

$$\text{Query Execution / Cache Entry} < 2 \quad \& \quad \sum_h (\text{Query Execution}) > 10000$$

GoFurther

13

One of the largest problems in dynamic SQL tuning is the usage of parameter markers: If you don't use them, every single instance of a dynamic statement is considered to be a separate statement, and needs a complete bind process. And will be placed in the dynamic statement cache even if it will never be rereferenced.

On the other hand, if you use parameter markers for all kind of values, the optimizer can not benefit from the frequency values statistics.

In general, you should use parameter markers for all numbers and names which are really variable, such as client numbers, customer names, addresses, etc.

Dynamic Statement Cache Time Analysis

```

DSPC 1.4.2 Aggregated SQL Overview Row 1 from 1197
ID: DB3P0214 10301045 Extract from: 2006.02.14 10:30:09
to : 2006.02.14 10:45:09

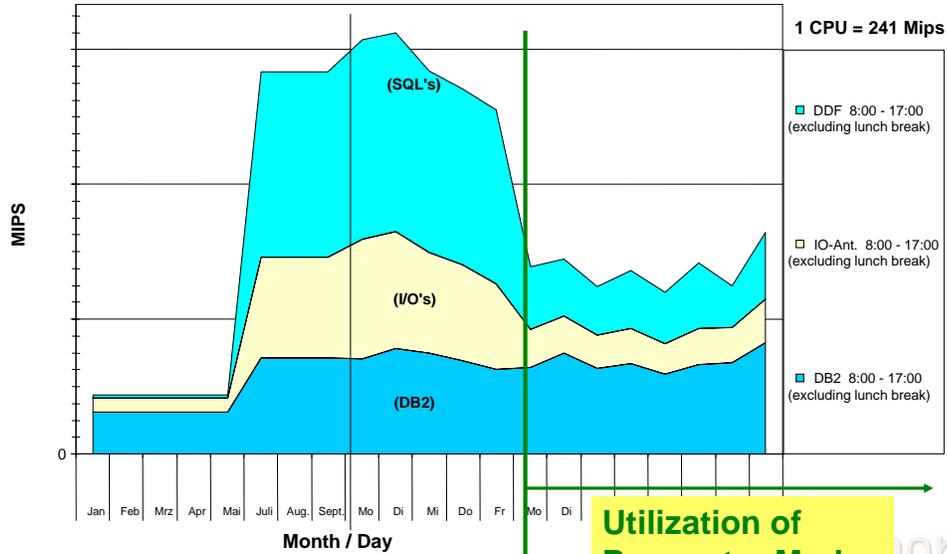
Line Commands: PI Performance Information T Show SQL Statement Text
CE Cache Entries TE Text Editor

LC          Total CPU      Total Elapsed      Executed      CE SQL Type
-----
00:00:00.762423 00:00:00.973919      3.770      2.177 SELECT
00:00:00.478855 00:00:02.072450      2.917      1.640 SELECT
    
```

Parameter markers should be implemented to efficiently use the dynamic sql statement cache: This is the case for really variable values (like customer numbers or account numbers). *Constants* should be used for rather static values such as a status code, a flag or a currency code (columns with a low cardinality and/or a very skewed data distribution).

Mips Usage of DB2, DDF + I/O portion: Production

Average over 8 hours

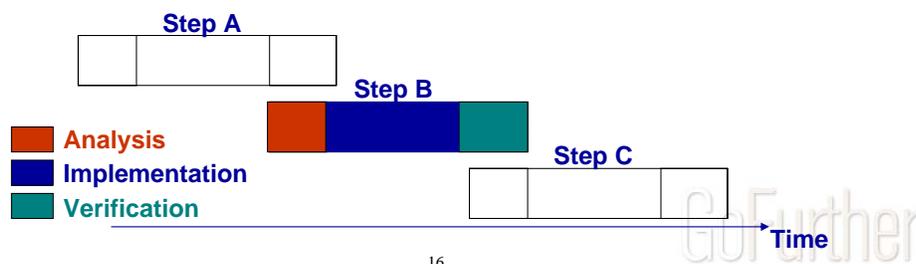


Utilization of Parameter Marker

This slide shows the effect of introducing parameter markers in one of Swiss Mobiliar's key applications. The overall CPU consumption (not only the CPU consumption of dynamic SQL statements) was divided by two!

BMW for SQL Tuning: Step B

- ✓ Step A: Are there extremely critical SQL queries?
- **Step B: Is the workload being adequately supported by the database system?**
- Step C: Are enough system resources (CPU, Memory, ...) available?



16

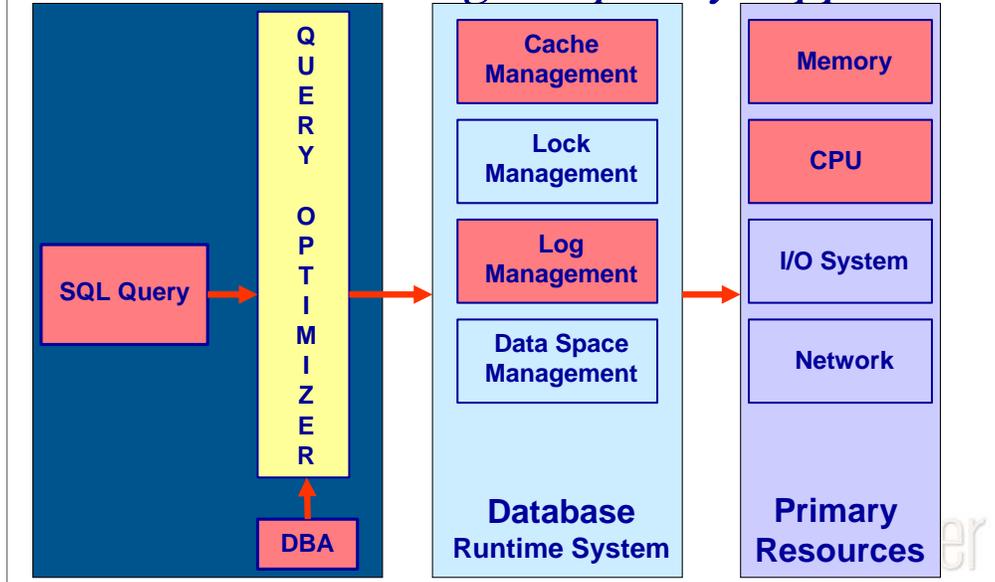
Now the first step has been done. The most critical queries are corrected. If they would still be in the system, there is a large probability that they would critically impact the results of the analysis of this second step – and leading to the wrong performance tuning actions.

If this methodology is used for routine monitoring, the first step should not identify any candidates.

Optimization opportunities detected during this second step do not require immediate (emergency) actions, but should be corrected until the next application release or system maintenance window. .

BMW step B:

Is the workload being adequately supported?



After tuning the most urgent SQL problems, there are still some areas covered in red which require your attention. So we don't want to lose too much time and start to resolve these problems immediately.

B1: Efficient WHERE predicates

Measurement criterion B1: SQL WHERE predicates / indexes

Identify and optimize SQL queries with

$$\frac{\text{Number of rows examined}}{\text{Number of rows fetched}} > 10$$

$$\& \sum_h \text{CPU / SQL} > 30 \text{ sec}$$

This criterion lists SQL statements which need a certain (too large) amount of rows to be analyzed before finding a row which corresponds to the where predicates. This might have different reasons:

- no indexes available which cover the query's predicates
- stage 2 predicates
- old or insufficient catalog statistics

etc.

Dynamic Statement Cache Time Analysis

example of inefficient index usage

```

SQL Statement Performance Information
COMMAND ==>                                SCROLL ==> CSR
ID: DB3P0424 10301045      Extract from   : 2006.04.24 10:30:53
                           to       : 2006.04.24 10:45:15
Aggr. ID : 1.857.158      in Cache since : 2006.04.22 09:03:18
SQL Type : SELECT        # Stmt execs  : 13
-----
                                Total      Average
CPU Time .....: 0:00:20.587893  0:00:01.583684
Elapsed Time .....: 0:00:23.529180  0:00:01.809936
# Synchronous Buffer reads .....: 16
# Getpage Operations .....: 727.677
# Rows examined for Statement ...: 1.211.327
# Rows returned for Statement ...: 39
# Sorts performed for Statement .: 0
# IDX Scans performed for Statement 1.709
# TS Scans performed for Statement 0
Wait Time Synchronous I/O .....: 0:00:00.356062  0:00:00.027389
    
```

These are typical values reported for queries with inefficient search attributes.

B2: Avoid using (Sort-) Workfiles

Measurement criterion B2: Do NOT use (sort-) workfiles

Identify and optimize SQL with

*(accesses to DSNWFQB Table or
number of sorts > 0)*

& $\sum_h \text{CPU} / \text{SQL} > 30 \text{ sec}$

The DSNWFQB table is a table which you or your DBA colleagues probably never have defined. It represents the materialization of a query block during query execution (a *workfile* temporarily created by DB2), and has a great potential of causing performance problems.

B3: Less Overhead for Locking

Measurement criterion B3: High degree of locking

Identify and optimize SQL Queries with

*(Wait Time Lock / SQL > 0.01 sec or
Wait Time Lock / SQL > 5% of Elapsed Time or
Wait Time Global Locks > 5% of Elapsed Time)*

& \sum_h Elapsed Time / SQL > 30 sec

Instead of looking at queries which cause deadlocks and timeouts, the queries reported here show a high amount of time spent in waiting for locks, but without necessarily reaching the timeout limits.

B4: Additional Information for the Optimizer

Measurement criterion B4: Query Optimizer problems

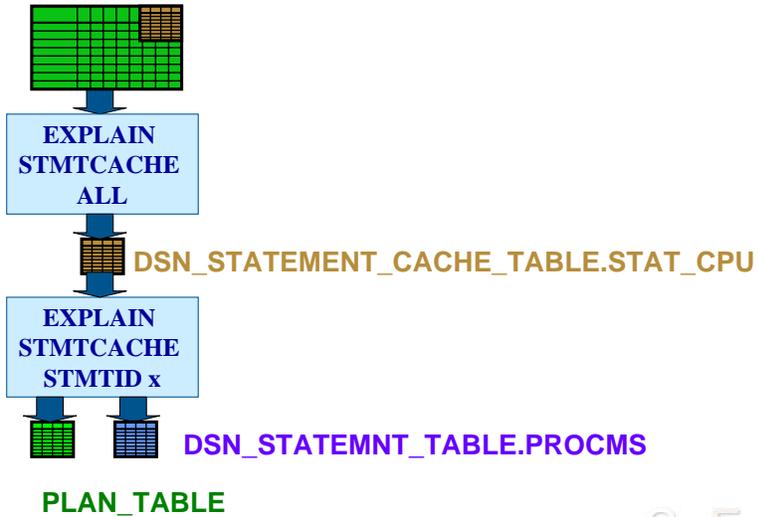
Identify, analyze and optimize SQL Queries with

$$\frac{\text{actual CPU (measured)}}{\text{estimated CPU (dsn_statemnt_table)}} > 5$$
$$\& \sum_h \text{CPU / SQL} > 20 \text{ sec}$$

If the optimizer chooses a wrong access path, the final cpu measurements are very often far away from what it initially estimated.

The ratio of measured compared to estimated CPU is a very good starting point to identify the queries which have inefficient access paths, often due to missing or wrong statistical information at the time the query was initially boud or prepared.

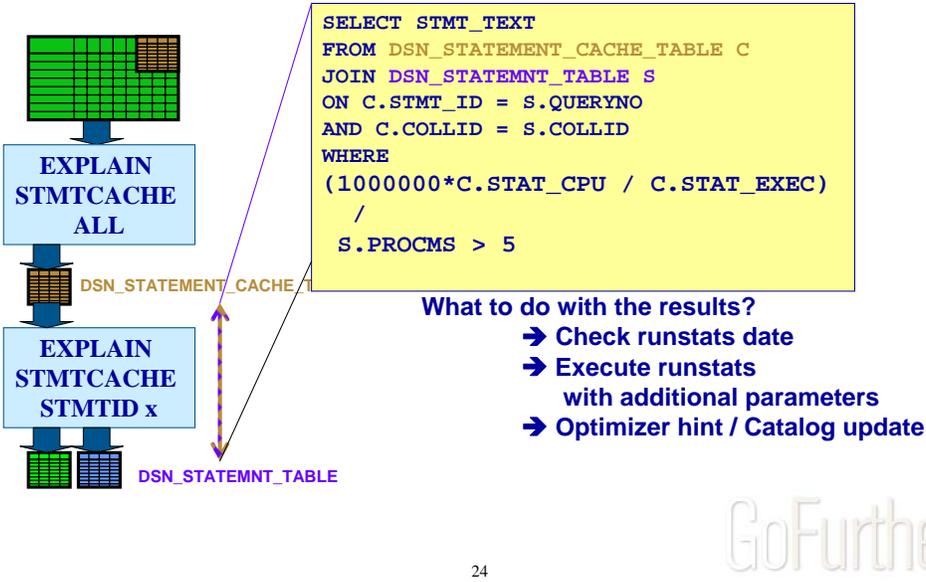
dsn_statemnt_table



GoFurther

This is how to get the information to calculate the *measured vs. estimated cpu ratio*. The EXPLAIN STMTCACHE STMTID x statement does not recalculate the access path, but does insert the originally calculated access path information into the plan_table, and the original cpu estimate into the dsn_statement_table's PROCMS column.

dsn_statemnt_table



The intensive usage of dynamic SQL changes RUNSTATS: Now, this utility has a certain potential to increase a system's stability: The new statistics are immediately active and might change access paths (in most cases, they become more efficient). In rare cases, they might become much more inefficient. Therefore, the usage of RUNSTATS should be more careful.

B5: Dyn SQL Statement Cache Sizing

Measurement criterion B5: Dynamic SQL Statement Cache Size

Minimum cache duration < 10 min → too little cache

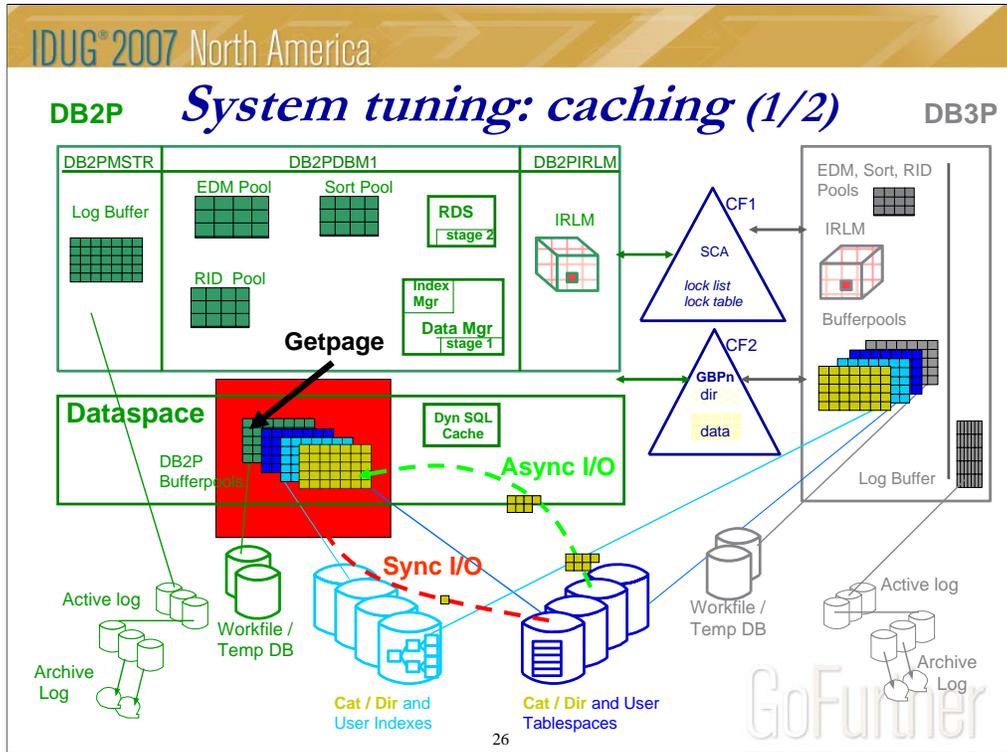
Approximation formula:

$$\text{Number of stmts in Cache} - \frac{\# \text{Stmt in Cache} * \text{hit ratio}}{\log(\# \text{Stmts in Cache})} < 600$$

$$\frac{\sum \text{new references per sec.}}{\log(\# \text{Stmts in Cache})} < 600$$

You might wonder how we found an upper limit of 10 minutes to be a good value to determine the dynamic statement cache's size. In fact we also do.

What we tested and measured was the impact of a larger dynamic statement cache (with a minimum cache residency time of up to 60 minutes). The benefit for the overall system was not much better than with a cache size according to the 10 minutes limit. However, with a smaller dynamic statement cache, the overall numbers and specifically the values measured for prepare-statements, were higher. So we decided to have 10 minutes as our performance goal.



Once again, the ,internal‘ view of where we are.

B6: Buffer Pool Configuration

Measurement criterion B6: Buffer Pool Size

Number of re-reads (I/O's) of the same page within 2 minutes > 0 → Bpool too small

Better than the Bufferpool Hit Ratio!

GoFurther

27

Bufferpool hit ratio is a very bad measurement criterion for bufferpool tuning: A single inefficient query (which retrieves for example the same few of a subquery a billion times) can lead to a bufferpool hit ratio of close to 100% - even with inefficiently organized bufferpools.

High bufferpool hit ratios can be a sign of good bufferpool tuning, but also a sign of inefficient queries being executed.

Another criterion (additionally or instead of 'No of. ReRead within 2 min > 0' is the minimal page residency time. If it is less than 2 min for a bufferpool the pool might be too small. The number of re-reads is more efficient (only really re-read pages are counted), whereas the minimal page residency time only identifies a potential for re-reads.

However, computing the number of re-reads involves executing an SQL performance trace (IFCID 101), whereas the minimal page residency time can be calculated using the same formula which was used for the dynamic statement cache. For FIFO type buffer pools, this formula is even simpler:

$$\frac{\text{Number of pages in Bpool} * (100 - \text{vpseqt})}{\text{Sync I/O per second} - 100} < 120$$

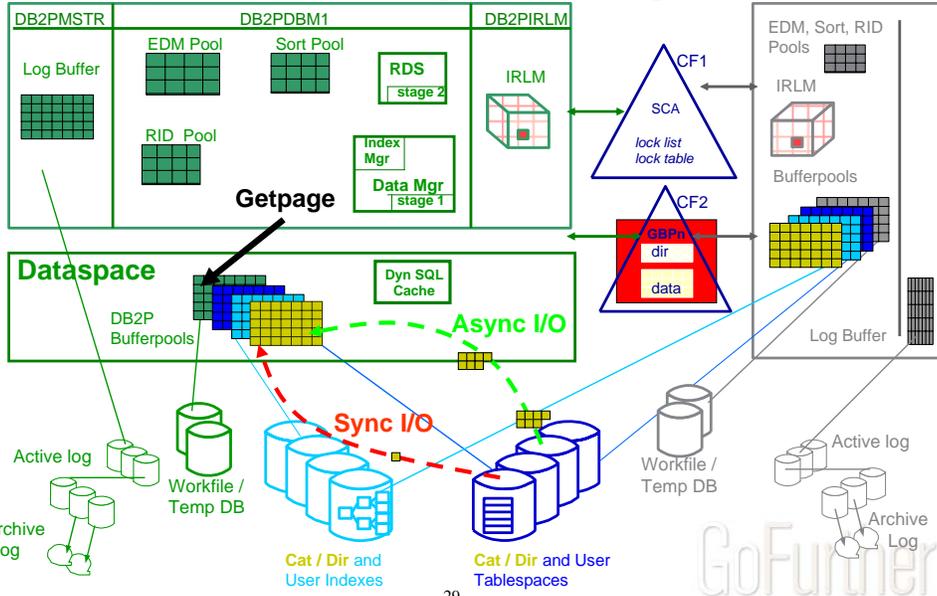
Buffer Pool Configuration DB2P

BP	Data type	Size	VPSEQT	DWT	VDWT	Method
0	Catalog	10K	80%	50%	10%	LRU
1	Data	1.0M	20%	0%	0%	LRU
2	Index	1.3M	20%	50%	10%	LRU
3	In-memory Heavy-read	40K	0%	50%	10%	FIFO
4	In-memory Heavy-upd	620K	0%	90%	90%	FIFO
5	Pipeline	5K	80%	80%	50%	FIFO
6	Workfile	20K	100%	70%	50%	LRU

28

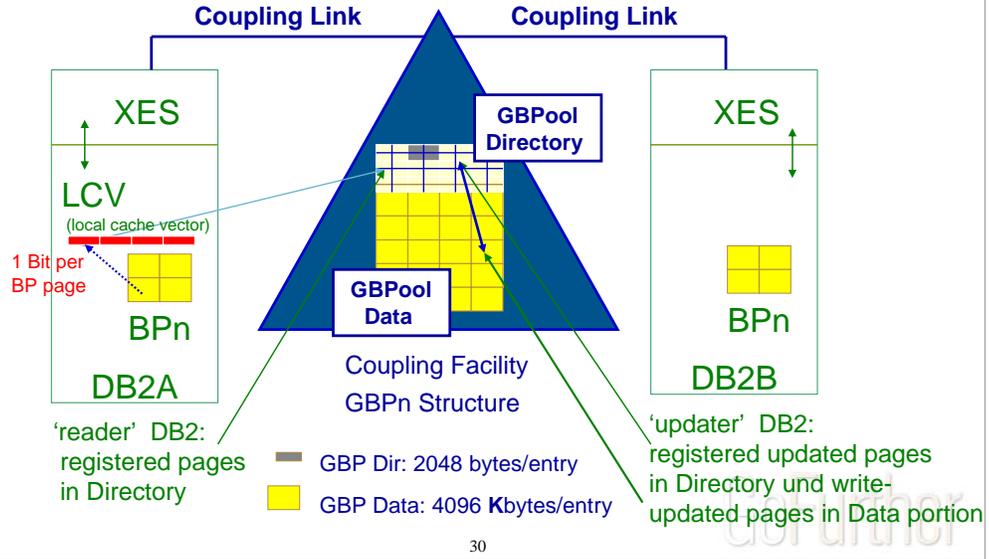
Probably, you have your own bufferpool strategy. It is not my goal to sell my bufferpool tuning strategy. This is just for information how another company has defined their bufferpools.

What is important is that you have a certain strategy, which is both easy to implement and easy to understand.



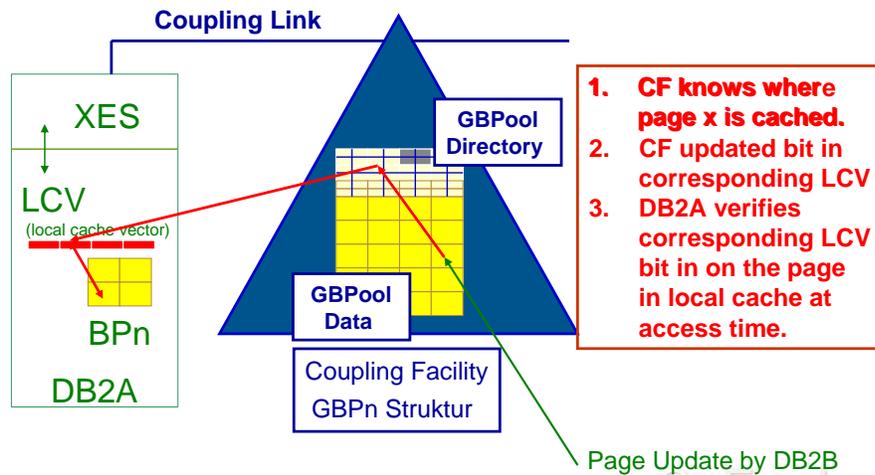
Now we look at the group buffer pool's point of view.

GBPOOL Overview



Group Bufferpools are used to make sure that all data sharing members get the most current state of a page. A positive side effect of this feature is the ability to avoid I/Os, which might already have been done by another data sharing member for the same page.

GBPOOL cross invalidation



31

GoFurther

There are two events which lead to ,cross invalidation‘:

1. A page update in data sharing member DB2B cross invalidates the local buffer pool entry for the same page in DB2A.
2. If, during the Page Registration Process, no free slot in the directory is found, the oldest (LRU) directory entry will be deleted and the corresponding page in the local bufferpool will be invalidated (via the same mechanism used during the page-update process). If this page is later re-referenced in the local bufferpool, it needs an additional disk I/O.

B7: Group Buffer Pool Checklist

Measurement criterion B7: Group Buffer Pool Tuning

- a) **Cross Invalidations due to Dir Reclaims > Cross Invalidations due to Writes & SyncRead (XI) Data not found >> SyncRead (XI) Data found**
 → GBPool Directory too small
- b) **Writes Failed due to Lack of Storage > 0**
 → GBPool Data portion too small
- c) **RMF: Sync Avg Service Time > 25 μ s**

- a) `-DIS GBPOOL(GBP1) GDETAIL(*)` and `-DIS GBPOOL(GBP1) MDETAIL(*)`
- b) `-DIS GBPOOL(GBP1) GDETAIL(*)`

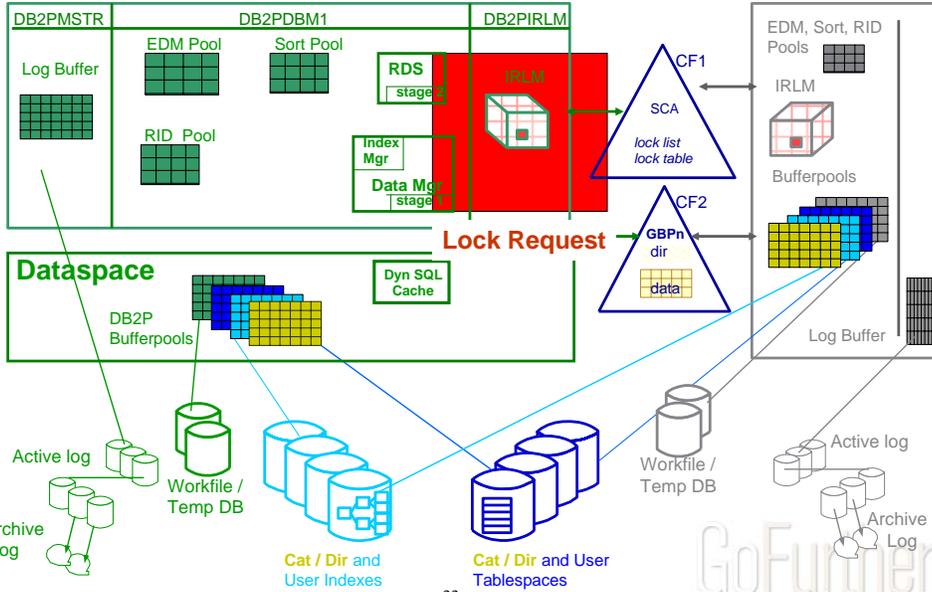
32

These are the key values which we are happy to meet. Specifically for dynamic SQL workload for DB2 for z/OS V8, look at these values for GBP8K0 and GBP16K0, since these buffer pools contain now catalog data necessary for bind (and prepare) actions. If you see writes failed due to lack of storage in these pools, make them higher (10K should be ok), and recalculate elapsed times for bind and prepare operations of complex queries.

DB2P

System tuning: locking

DB3P



B8: High degree of Locking

Measurement criterion B8: High degree of Locking Activity

- a) #Unlock Requests / Commit > 5**
- b) #Lock Requests / #Unlock Requests < 3**
- c) #Lock Requests > #SQL Statements**
- d) Data Sharing Only:**
 - False Lock Contention < 2% → Lock Table too small**
 - RMF: Sync Avg Service Time LOCK1 > 15 μs**

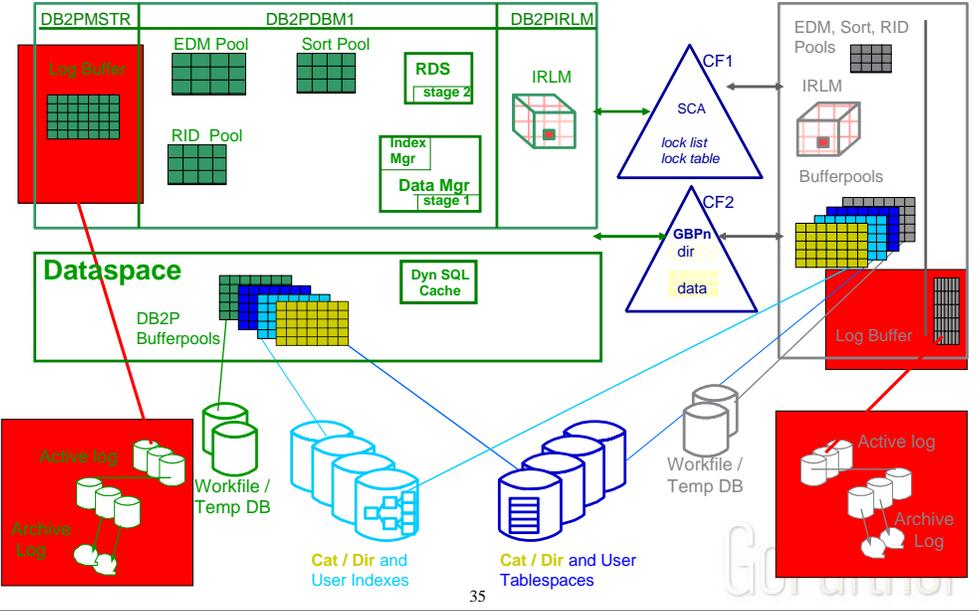
GoFurther

Unfortunately, there is no knob to reduce locking or to resolve locking problems in general. Most important is the isolation level, and best practise is to have uncommitted read as your standard isolation level for remote applications – with exceptions at the statement level.

DB2P

System tuning: logging

DB3P



B9: Unnecessary Log Waits

Measurement criterion B9: Log Tuning

*Wait for Active Log Write I/O > 1% of Elapsed Time *)*

→ enlarge the Output Log Buffer

**) applies only if Log Speed <= 2 Mbyte / sec*

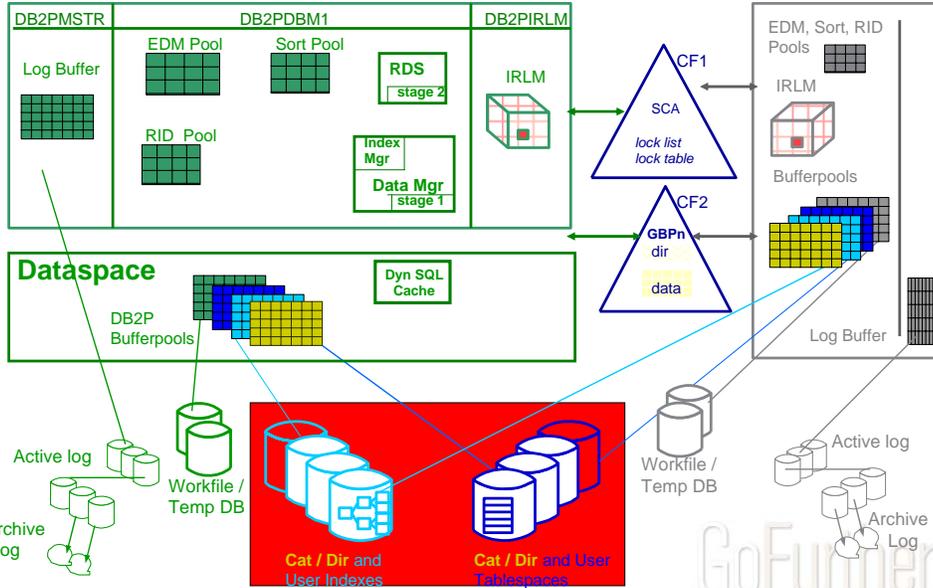
GoFurther

36

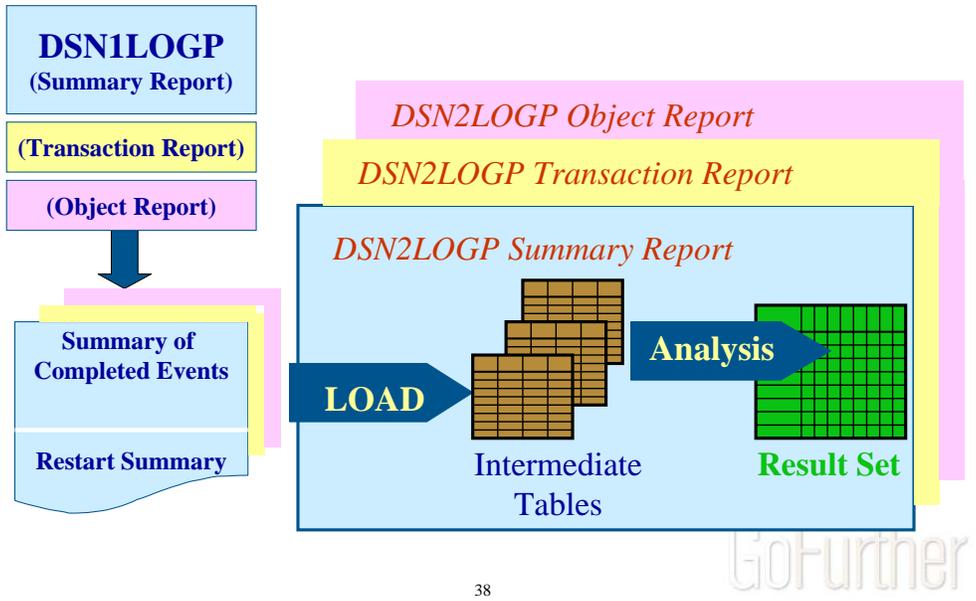
The log output buffer is situated in the DSNxMSTR address space, and not in the DSNxDBM1 address space. The good news is that there are no memory constraints, the bad news is that the additional memory which you require when increasing OUTBUFF zparm (which controls the output log buffer size) does not come for free, and is in competition with the memory demands for buffer pools, dynamic statement cache, etc.

In a nutshell: Make the output log buffer as large as necessary, but keep it as small as possible!

System tuning: Data Space Management



Log analysis: DSN2LOGP



38

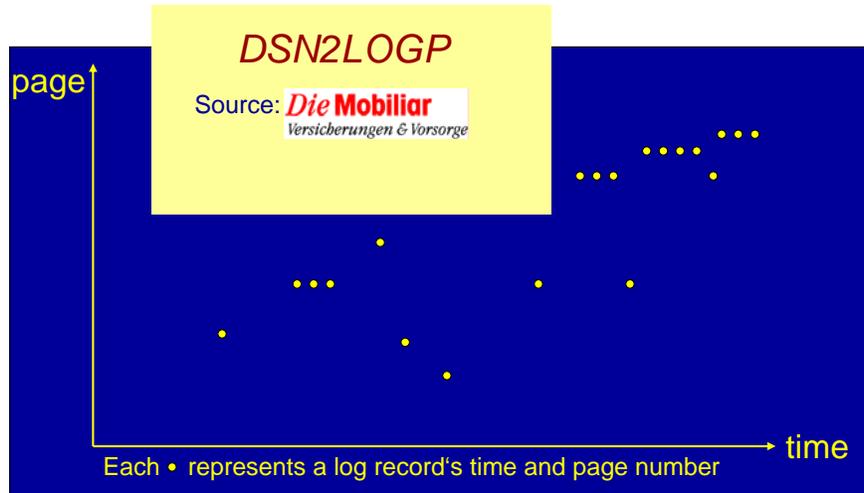
As DSN1LOGP's output is hard to interpret (especially if you want to do that using a systematic approach), we load the DSN1LOGP output into intermediate tables from which we calculate performance results, such as the page numbers accessed at given time points.

This straight forward approach to a graphical interface of a pageset's activity is sometimes of great help.

DSN2LOGP Object Report: Time View

Time	Page	Transaction URID	Plan
2006-01-11-15.18.28.236272	00007F33	142071F38D34	P72S30
2006-01-11-15.18.32.291728	0000864C	142072298BF9	P72S30
2006-01-11-15.18.47.692144	00000003	0BB37A14E000	P14501
2006-01-11-15.18.47.793520	00007F33	142072C7F689	P72S30

DSN2LOGP Object Report: Time View



There is no simple measurement criterion for cases such as this. Based on pattern detection, adequate optimization techniques might be applied, always bearing in mind the tradeoff between optimal locking behaviour and minimal I/Os.

B10: Data Access Patterns

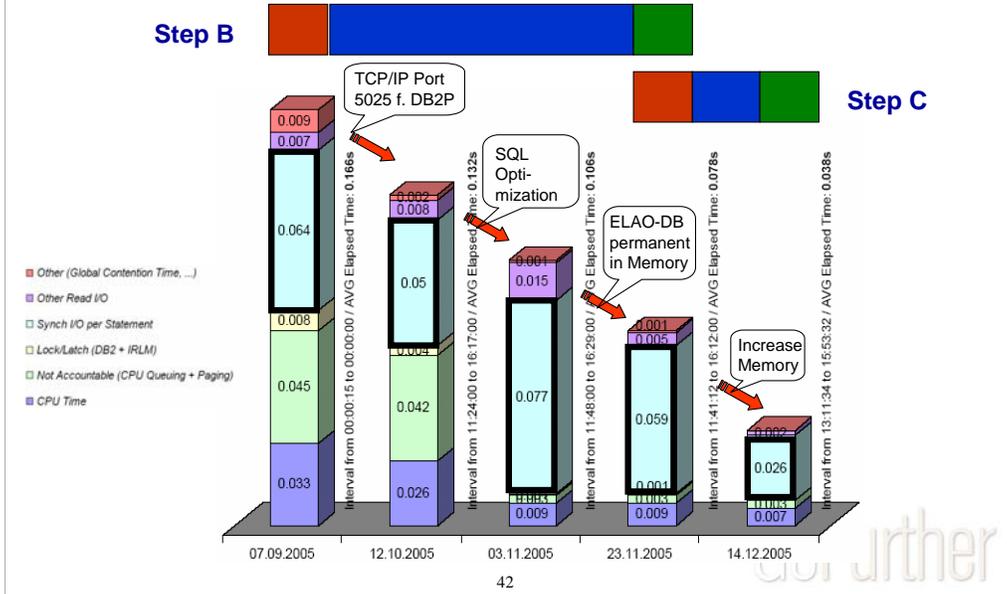
Measurement criterion B10: Page access for heavily updated tables

Are there red flags in reference to...

- a) *time-related #Getpages/SQL differences?***
- b) *time-related synchronous I/O wait times/SQL differences?***
- c) *Page updates spread over too few pages?***
 - *Hot Spot Pages for objects that are often changed?***
- d) *Page Updates spread over too many pages?***
 - *Inefficient clustering ?***

Some questions arising from DSN2LOGP object analysis.

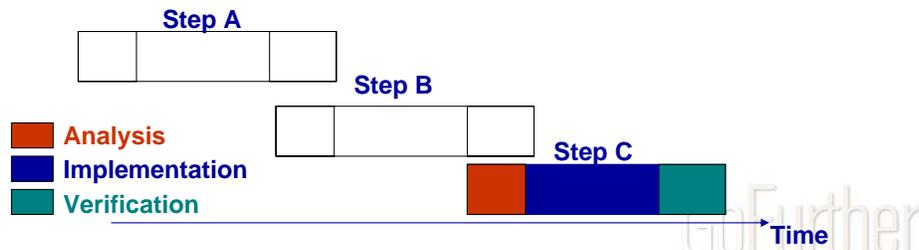
DB2 Response Time Improvement



These measurements were produced by IBM during a review of our performance activities after a major application's rollout. The numbers identify average processing times of single SQL statements of this application.

BMW for SQL tuning: Step C

- ✓ Step A: Are there extremely critical SQL queries?
- ✓ Step B: Is the workload being adequately supported?
- **Step C: Are enough system resources (CPU, memory...) available?**

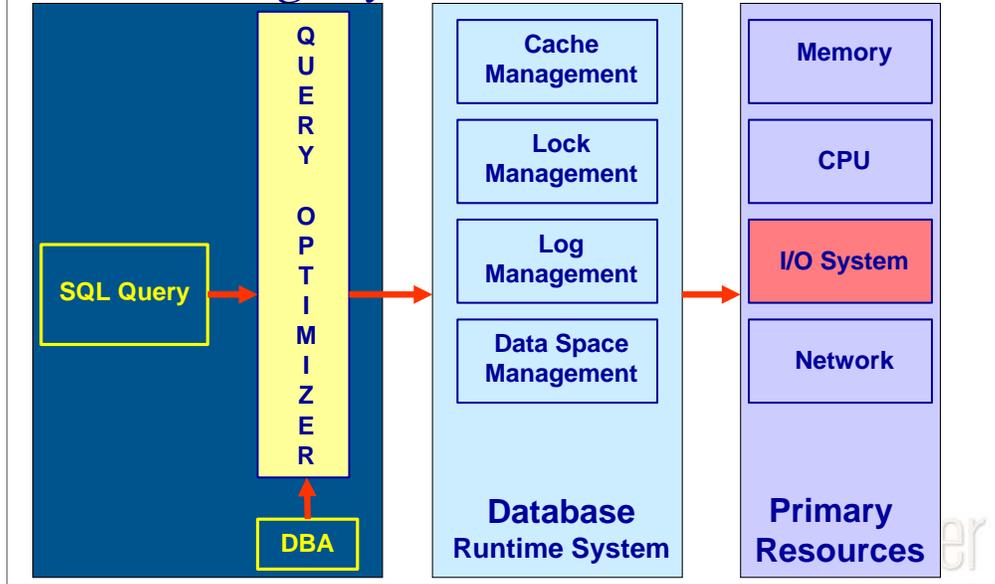


43

Now the second step of the methodology has been finished. Again – all the problems identified must be resolved before the third step can start..

BMW Step C:

Are enough system resources available?



Now it is time to concentrate on primary resources.

C1-C4: Check System Resources

Identify bottlenecks:

C1) Memory Check:

Bufferpool re-read within 2min > 0, Paging ?

C2) I/O Check:

Average Sync I/O Suspension > 8msec

C3) Network Check:

End-to-end response time/SQL processing time > 2

C4) CPU Check:

'Wait for CPU' / SQL processing time > 0.5

Some simple resource checks.

Session: A04
Troubleshooting Tuning for Dynamic SQL

Thomas Baumann

Swiss Mobiliar

thomas.baumann@mobi.ch



Thomas Baumann has been a database performance and availability professional for many years. At Swiss Mobiliar, he is, as a data architect, responsible for DB2 and is engaged in the analysis and integration of advances in the database technology with existing and new Swiss Mobiliar applications. He holds a master degree from Swiss Federal Institute of Technology (ETH) Zürich in computer science and probability theory. Thomas has also lectured on data management at the Zurich University for Applied Sciences and has conducted various workshops and delivered presentations worldwide. He is the author of the *Advanced Recovery Functions* seminar series and a member of the IDUG hall of fame, a small but distinguished group of DB2 professionals. You can reach him at thomas.baumann@mobi.ch.