

2FAST 2FURIOUS

designing
developing for



speed

correctness

concurrency

MARK BROADBENT

Contact...



contactme@sturmovik.net



@retracement



tenbulls.co.uk

Likes...



Guilty pleasures...



Badges...



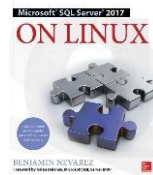
Master: SQL Server



Community...



Tech ed/ reviewer...

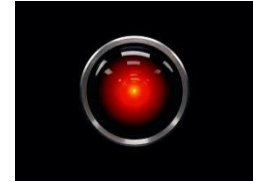


2nd ed
(Coming soon!)

Agenda



Query Tuning and
Indexing



Isolation and Locking



Mechanics



Improvements

Speed, Concurrency, Correctness?



Speed



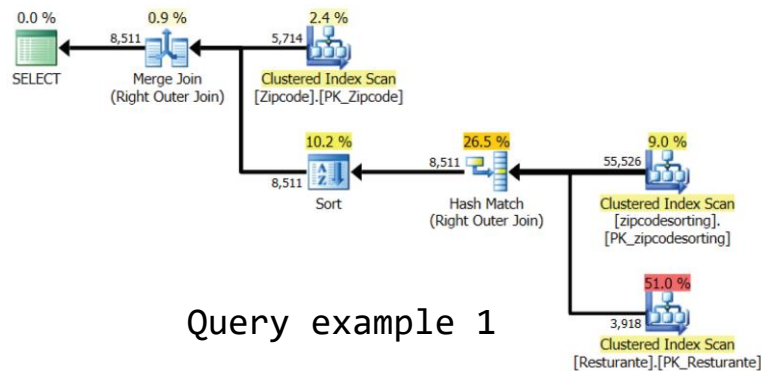
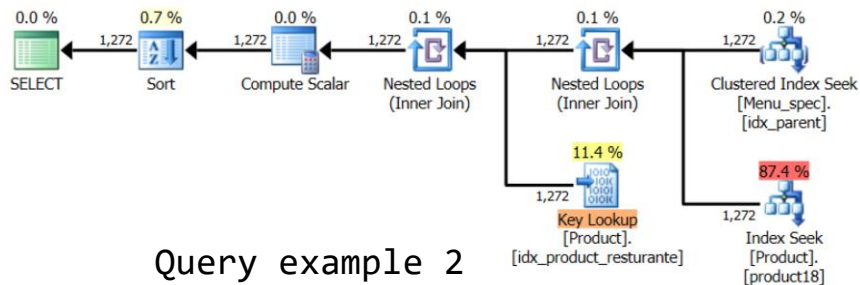
Concurrency

A large grid of data, likely a race results table, showing multiple columns of numbers and text. The data is organized in a structured format, possibly representing race times, positions, and driver names. The grid is composed of many small cells, each containing a small amount of text or numbers.

Correctness

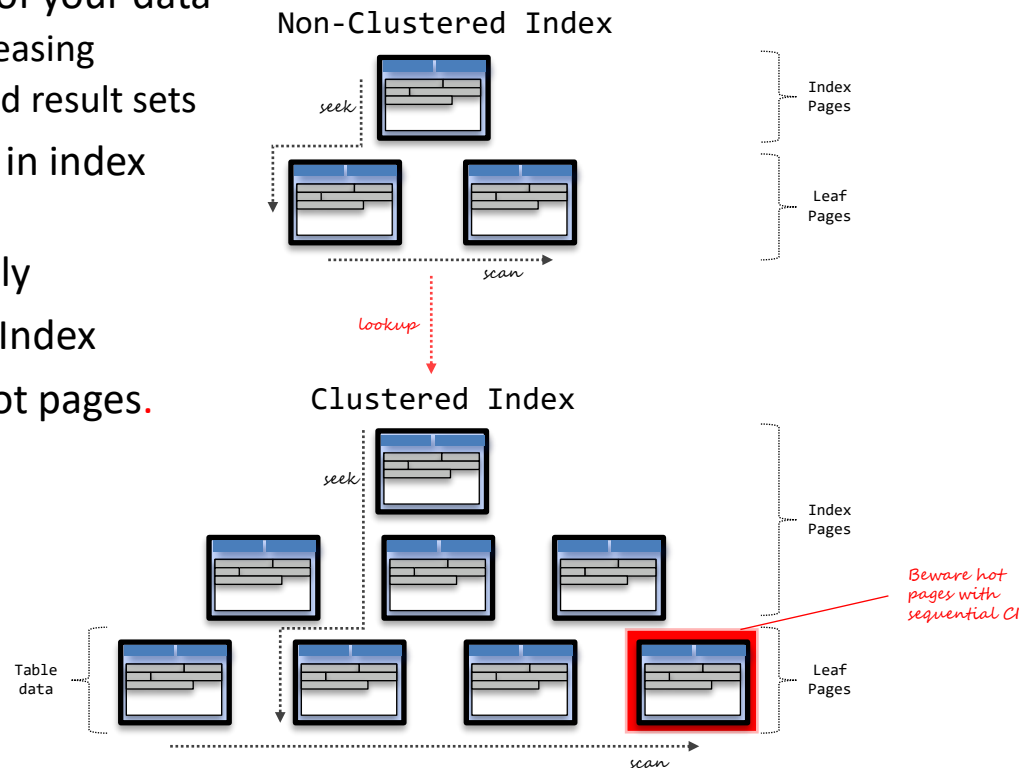
Is my query a problem?

- Review query cost and memory grant
- Compare #rows returned against #rows accessed by the child operators
- Problem operators: scans, lookups, sorts, joins – why are they used?
- Predicate logic (e.g. using OR requires a separate op, but cannot return duplicates so will often require a sort and temp table to eliminate)
- Beware of joins or filters on different predicate types –implicit conversion.



Indexes

- Clustered Index defines logical order of your data
 - Keep unique, narrow, static, ever increasing
 - good for range scans, ordered, covered result sets
- # Key columns determine space used in index AND leaf pages
- Included columns in NC leaf pages only
- Cluster key present in Non-Clustered Index
- Sequential Cluster key can result in hot pages.



Query Tuning and Indexing review

- Limit query columns (and if all necessary) INCLUDE in index (for covering)
 - Index key columns may be used for the join or filter predicates for efficient seeks
 - INCLUDE(d) columns avoid expensive lookups
- Indexes add overhead – HEAPS ONLY really only good for loading (*caveat on RID size*)
- Keep index key columns few and small
 - But index on multiple columns to achieve high SARGability
- Index creates statistics object. Temp stats could indicate missing indexes
- Leading wildcards on predicates will not use index, use 'M%' rather than '%ark'
- Calculated predicates will scan (e.g WHERE Sales * Qty < 100).

Demo

Query “Speed”

(Demo 1)

Isolation and Locking

Isolation Level	“Bad Dependencies”	TX Dependencies	Concurrency
READ UNCOMMITTED	Dirty Read, Non-Repeatable Read, Phantoms, Lost Updates	WRITE → WRITE	GOOD: prone to allocation order scan failures
READ COMMITTED	Non-Repeatable Read, Phantoms , Lost Updates	WRITE → WRITE READ → WRITE WRITE → READ	OK: wait on both only writers held to EOT
READ COMMITTED (with Snapshot)	Non-Repeatable Read, Phantoms , Lost Updates	WRITE → WRITE	GOOD: no wait on readers only writers held to EOT
REPEATABLE READ	Phantoms	WRITE → WRITE READ → WRITE WRITE → READ	LOW: read and write locks held to EOT
SERIALIZABLE	None	WRITE → WRITE WRITE → READ READ → WRITE	LOWEST: read and write locks held to EOT
SNAPSHOT	<i>None (though Causal consistency concerns, lost update prevention and other behaviours)</i>	WRITE → WRITE	GOOD: only wait on writes held to EOT

Isolation levels with In-Memory OLTP are either restricted and/ or non-blocking

Isolation review

- Isolation Levels provide balance between consistency and correctness
- Set at session level, transaction level and statement
- Only RC or RCSI can be set as a default (on disk)
- Both on-disk SI (on disk) and RCSI still block on writes – but remove need for NOLOCK
- Use In-Memory OLTP is functionally better than either.

Locking review

- Locks ONLY memory structures (lock blocks), can consume a lot of memory
 - Try to avoid the need for escalation!
- Explicitly use table locks **ONLY** when it makes sense (bulk loading/ reporting)
- NOLOCK is prone to failures for long running queries (and dirty reads) instead use:
 - Use SNAPSHOT ISOLATION for consistency
 - Use READ COMMITTED SNAPSHOT for concurrency (and as new default)
- Lock duration depends upon ISOLATION LEVEL
 - Write locks last for the entire transaction!
 - Read locks (by default) are taken and released for each row read.

Demo

Correctness and on-disk concurrency

(Demo 2 and 3)

SQL Server IO

spid 115

```
SELECT name FROM
CarInventory
WHERE id = 10001
```

Return records
using locking
to logical table
structure

Potential
block

"On-disk" logical table(rows)

CarInventory

Access Cache pages
through Latching

Buffer Cache (memory pages)

Read pages to
cache when
not present

Write dirty
pages to disk
on checkpoint

Data File/s

spid 115

```
UPDATE
CarInventory
SET stocklvl = 2
WHERE id = 10001
```

Lock records and modify
pages in Cache

spid 162

```
BEGIN TRAN
INSERT CarInventory(id, stocklvl)
VALUES (10002,10)
INSERT CarInventory(id, stocklvl)
VALUES (10002,10)
...
COMMIT
```

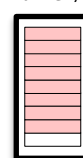
Write records on execution

spid 88

```
EXEC nsp_UpdManufacturer
'Ford', 'Michigan'
```

Write records
on commit

Log
Buffer/s



Signal to flush
on buffer full/
on commit^{*1}

Log
Writer

Persist log
record to disk

Log File

In-Memory table (all rows in memory)

products

Existing row is expired through timestamp

Chkpt file pairs + log used
to load IM table on startup



Checkpoint file pairs

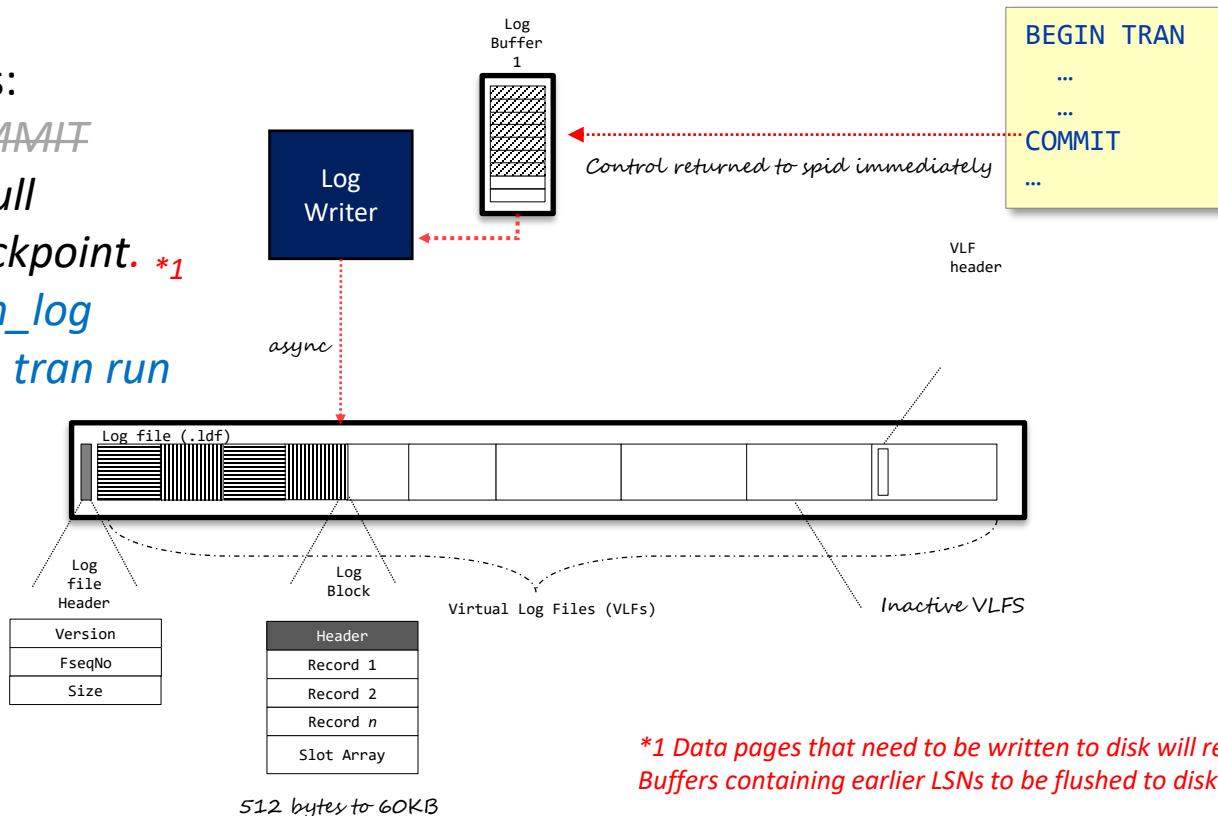
^{*1} Except when delayed
durability is in use

SQL IO review

- Large reads or writes will reduce your Buffer cache page life expectancy
 - Use page compression on a table and index pages
 - Set **max server memory** so that OS has a reservation –Approx 2GB, set to limit buffer cache
 - Implement indexing strategy and tune queries
- Transaction log will always be a key bottleneck, improve by
 - Reduce IO overhead (reduce data change)/ Consider In-Memory OLTP
 - Consider Delayed Durability
 - Improve hardware
 - Avoid logfile log growths.

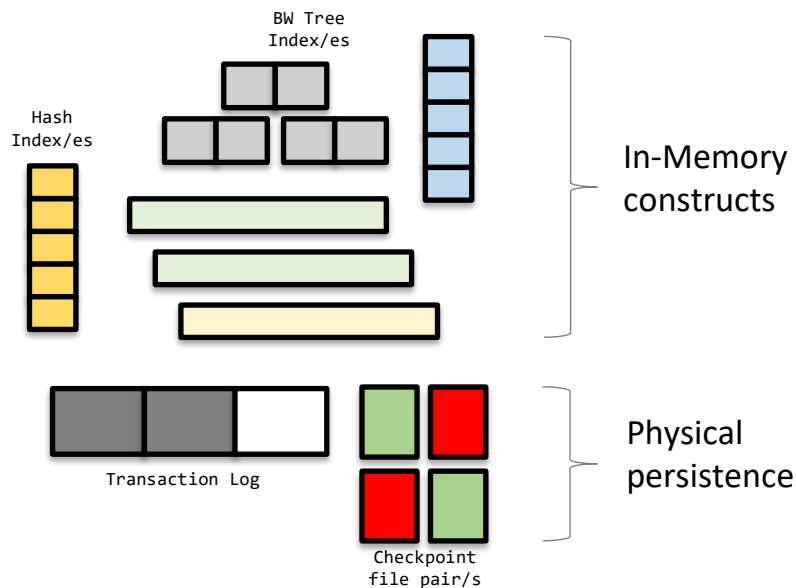
~~● On COMMIT~~

- *1 Data pages that need to be written to disk will require Buffers containing earlier LSNs to be flushed to disk*



Using In-Memory OLTP

- In-Memory data structures “optimised for memory”
- Persistence thru Transaction Log and checkpoint file pair/s
- Query thru interop or Natively Compiled Stored Procedures
- No TempDB overhead – versioning In-Memory
- Logging optimizations and improvements
- Lockless and latchless operation
- No fragmentation concerns.



Demo

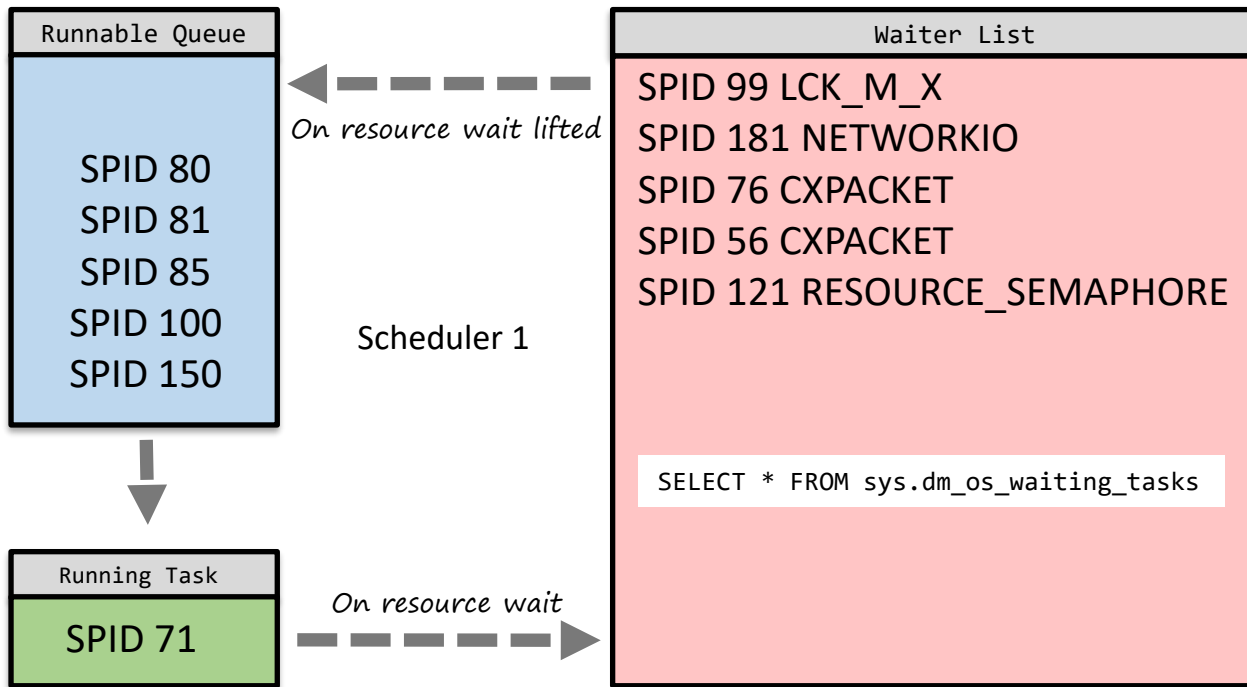
Other improvements and utilizing other technologies
(Demo 4 and 5)

CPU (SQLOS Scheduling)

- 1 scheduler per logical CPU
- **wait_time_ms** is sum of time on waiter list + time on runnable queue
- **signal_wait_time_ms** is time spent on runnable queue

THEREFORE

- High signal waits → high CPU pressure
- Large time on waiter list → high resource bottlenecks



`SELECT * FROM sys.dm_os_wait_stats`

Wait_type	waiting_tasks_count	wait_time_ms	Max wait_time_ms	signal_wait_time_ms
RESOURCE_SEMAPHORE	6	50	20	8
CXPACKET	120	50	200	18
LCK_M_S	19199	5000	780	80
LCK_M_X	3	21	7	3

SQLOS Scheduling review (common waits)

- LCK_M_* (wait for lock grant for read/write/update) → locking/ concurrency bottleneck.
- RESOURCE_SEMAPHORE (wait for query memory grant) → memory bottleneck
- CXPACKET waits always exist. Excessive CXPACKET waits → imbalance of parallelism
- SOS_SCHEDULER_YIELD (exhausted runnable queue quantum) → CPU pressure
- PAGEIOLATCH_* (latch wait for IO read/write to memory) → disk IO bottleneck
- PAGELATCH_* (latch wait for read/write to a page in memory) → system overloaded
- LATCH_* (latch wait for non-buffer resource) → resource overloaded.

<https://www.sqlskills.com/help/waits/>

[https://docs.microsoft.com/en-us/previous-versions/sql/sql-server-2005/administrator/cc966413\(v=technet.10\)](https://docs.microsoft.com/en-us/previous-versions/sql/sql-server-2005/administrator/cc966413(v=technet.10))

<https://www.sqlskills.com/blogs/paul/capturing-wait-statistics-period-time/>

Wait stats are expected and relative to your workload

SQLOS Scheduling review

- Parallel queries will utilize all schedulers for their runnable quantum (and most by default will go parallel!)
 - Set ***max degree of parallelism*** to avoid single heavy parallel query from overwhelming cores (consider NUMA configuration too)
 - Set ***cost threshold for parallelism*** to sensible level.
- Monitor waits and queues performance impact over time (rather than a snapshot)
- Consider implementing Resource Governor.

In summary

- Keep transactions short and reduce complex logic
- Use RCSI or SI isolation – especially when on disk
- Capture, review, and optimise all query plans (smallest cost, most efficient access)
- Implement new technologies to improve concurrency
 - In-Memory OLTP, Columnstore, and Delayed Durability should be considered
- **ALWAYS** Test your results in parallel to prove consistency.