# Logging

"We therefore consider this bottleneck as the most dangerous to future scalability"

Canonical recovery algorithm
 Decouples data writing from log writing

### ARIES

Ø Write ahead logging
Ø Restoring database to pre-crash
Ø Undoing uncommitted txns

## 3 Components

 Create an additional log while undoing transactions to recover multiple crashes

Multiple Crashes

### Aether

Ø OLTP txns are small and frequent
Ø More cores → more contention
Ø Other bottlenecks disappearing

### Motivation



٢	ଷ	I/O delay
	ଷ	Lock contention
	ଷ	Schedule overhead
L	ଷ	Log buffer contentior

### Convoluted Drawing

### Tradeoffs suck.

## Speed + security rocks.

Ø Critical path
Ø Decouple to reduce contention
Ø Commit ≠ return



Ø Critical path
Ø Decouple to reduce contention
Ø Commit ≠ return



## I/O Delay and Scheduling

Release locks on commit
Track dependencies
Abort and rollback if necessary

### Early Lock Release





 Ø Decouples transaction commit from scheduling
 Ø Keeps threads busy

## Flush Pipelining

Allowing other transactions to proceed speculatively + providing the threads to actually execute them rivals asynchronous commit, but SAFELY

## Combining the Two



## Scalable Log Buffer



Log mgr. contention contention Log mgr. Other work

High core count + Ś high load  $\rightarrow$ bottleneck in the log buffer

## Log Buffer Bottleneck



ø To write to the log, threads ର୍ଷ Acquire space in the buffer ର୍ଷ Fill space ର୍ଷ Release buffer space for writing

### The Problem

Ø Group logs together to write to the buffer at once
Ø Use a consolidation array as a backoff

structure

### Solution C

### Observation

#### **linearizable stack**



Yes!

Push(<sup>O</sup>)

Pop()



23









Art of Multiprocessor Programming 25



### **Elimination-Backoff Stack**

- Lock-free stack + elimination array
- Access Lock-free stack,
  - If uncontended, apply operation
  - if contended, back off to elimination array and attempt elimination





Consolidation array combines updates rather than eliminating them
Group leader does all the work
Great in theory as well as practice

### Consolidation Array



Ø Buffer fill is not inherently serial
Ø Must release regions in proper order

## Decoupled Buffer Fill



#### я Hybrid approach я Consolidate and fill buffer in parallel

### Two > One

### Quick Summary

### Does it actually work?



 Left shows record size average held to 120B constant

Right shows thread count constant (64)

## Experimental Results

How good?



 Left shows record size average held to 120B constant

Right shows thread count constant (64)

## Experimental Results



 Flush pipelining + ELR is most important
 Log buffer contention become increasingly important as corecounts grow

### **Overall** Picture

Delegation can prevent the problem with varying record size
Threads wake the next in line
More robust, but performance penalty in the normal case

## Further Optimization

In summary...

Ø Distributed logging
Ø Higher level (txn + parameters)
Ø Different (serial excn per thread)
<u>Ø Txns strongly ordered</u>

### H-Store WAL

### H-Store WAL

## Checkpointing

R Natural points of consistently
R At least once per second
R Consistent (not fuzzy)

Frequently Consistent

### Requirements

Synchronously copies state
Asynchronously writes to disk
No checkpoint-specific work

## Naïve-Snapshot



and AS<sub>1</sub> contain the same information

5	6	1	1
9	9	0	1
7	1	1	1
2	9	1	1
4	4	0	1
3	3	0	1
AS0	AS1	MR	MW

3	3	0	1	
4	4	0	1	
2	9	1	0	
7	1	1	0	
9	9	0	1	
5	6	1	0	

3	6	0	0
9	8	1	1
7	1	1	0
2	9	1	0
4	4	0	1
3	3	0	1
AS0	451	MR	MW

(a) At the beginning of time, AS<sub>0</sub> (b) During the first checkpoint (c) The state right after switching period, some updates from the to the second checkpoint period Mutator are applied

(d) In the second checkpoint period, the Mutator applies additional updates

Figure 1: Wait-Free Zigzag Example

*π* Maintain untouched copy of every word for duration of checkpoint *π* Bits determine which copy to use

## Wait-Free Zigzag



Maintains extra version of app state
Swing pointer rather than reset all
Cache-friendly interleaved version

## Wait-Free Ping-Pong

#### Uniform overhead



#### Low overhead



#### High throughput



### Results