BzTree: A High-Performance Latch-free Range Index for Non-Volatile Memory

JOY ARULRAJ JUSTIN LEVANDOSKI UMAR FAROOQ MINHAS PER-AKE LARSON Microsoft Research

NON-VOLATILE MEMORY [NVM]



DURABILITY

DEVICE CHARACTERISTICS

CHARACTERISTIC	DRAM	NVM	SSD
Device Latency	1x	10x	1000x
Byte-Addressability	\checkmark	\checkmark	×
Durability	×	\checkmark	\checkmark
High Capacity	×	\checkmark	\checkmark

BWTREE: LATCH-FREE B+TREE



SINGLE-WORD COMPARE-AND-SWAP INSTRUCTION



BZTREE: NVM-CENTRIC LATCH-FREE B+TREE









BWTREE INDEX



EXPERIMENTAL RESULTS

BWTREE: SSD-CENTRIC ARCHITECTURE



BWTREE: LATCH-FREE ALGORITHMS



BWTREE: LOGGING & RECOVERY PROTOCOL



BEGIN TRANSACTION Update Stock by Stock ID COMMIT TRANSACTION



BWTREE: RECAP

- Delivers high performance on a DRAM + SSD system
 - SSD-centric architecture
 - Latch-free algorithms
 - Logging & recovery **protocol**
- Limitations
 - NVM invalidates the key design assumptions of BwTree
 - Challenging to design & extend such latch-free data structures

PROBLEM #1: ALGORITHMIC COMPLEXITY



PROBLEM #2: PROTOCOL COMPLEXITY



PROBLEM #3: ARCHITECTURAL COMPLEXITY



HOW CAN WE SIMPLIFY LATCH-FREE PROGRAMMING ON NON-VOLATILE MEMORY?







BWTREE INDEX



EXPERIMENTAL RESULTS

BZTREE: OVERVIEW

- NVM-centric design
 - Based on a new NVM-centric software primitive
 - Provides same guarantees as disk-centric BwTree
- BzTree supersedes BwTree (skipped BxTree and ByTree)
 - Because we think that it is the last index you will ever need!
- Key techniques
 - Adopt a simpler NVM-centric architecture
 - Reduce complexity using software primitive

NVM-CENTRIC SOFTWARE PRIMITIVE





VOLATILE SINGLE-WORD COMPARE-AND-SWAP

PERSISTENT MULTI-WORD COMPARE-AND-SWAP



EASY LOCK-FREE INDEXING IN NON-VOLATILE MEMORY ICDE 2018

BZTREE: NVM-CENTRIC ARCHITECTURE



BZTREE: DURABILITY & ATOMICITY







OPERATION TABLE

	LOCATION	N EXPECTED OLD VALUI	NEW VALUE	FLUSHED
•	0x100	OLD CHILD POINTER	NEW CHILD POINTER	1
NVM	0x200	OLD NODE STATUS	NEW NODE STATUS	1
	0x300	OLD PARENT POINTER	NEW PARENT POINTER	0

SOLUTION #1: ALGORITHMIC COMPLEXITY



SOLUTION #2: PROTOCOL COMPLEXITY



INDEX





LOCATION	OLD VALUE	NEW VALUE	FLUSHED
0x100	OLD CHILD POINTER	NEW CHILD POINTER	1
0x200	OLD NODE STATUS	NEW NODE STATUS	1
0x300	OLD PARENT POINTER	NEW PARENT POINTER	0

NO INDEX-SPECIFIC PROTOCOL

DURABILITY & ATOMICITY

SOLUTION #3: ARCHITECTURAL COMPLEXITY

NO MAPPING TABLE

NO DELTA RECORDS & INDIRECTION OVERHEAD



NO LOG STRUCTURED STORE











BWTREE INDEX



EXPERIMENTAL RESULTS

EVALUATION

- Index data structures: BzTree vs. BwTree index
 - Code complexity
 - Runtime performance
 - Recovery time
- Benchmark: Yahoo Cloud Serving benchmark
 - Read-mostly & Balanced workloads
- Storage device
 - Emulated Non-Volatile Memory

CODE COMPLEXITY

		CODE COMPLEXITY METRIC	BWTREE	BZTREE		
Lower is Better	CYCLOMATIC COMPLEXITY	12	7	₩	2 x	
	LINES OF CODE	750	200	₩	4 x	
1 FEWER INTERMEDIATE STAT		FEWER INTERMEDIATE STATES	NO INDEX- LOGGING P	SPECIFIC ROTOCOL	·	

RUNTIME PERFORMANCE



RECOVERY TIME

- BzTree: no recovery logic
 - Recovery is entirely handled by software primitive
 - Rolls back operations that were in progress during the crash



CONCLUSION

- NVM invalidates design assumptions in data structures
- Presented the design of a NVM-centric latch-free B+tree
- Importance of tailoring data structures for NVM

