



APACHE MAVIBOT

A NEW BACKEND FOR APACHE DIRECTORY SERVER

Emmanuel Lécharny <u>elecharny@apache.org</u> http://symas.com







"BLUE BOAT" IN TURKISH (THANKS TO ERSIN ER FOR THE NAME) ALSO BUILT FROM MVBT (MULTI-VERSIONS B-TREE)







- A NODE HAS AT LEAST N/Z KEYS AND UP TO N KEYS (EXCEPT THE ROOT NODE)
- LEAVES HAVE KEYS AND VALUES
- THE B+TREE IS BALANCED (IE ALL LEAVES ARE AT THE SAME LEVEL) 3 5
- OPERATIONS :
 - Add
 - DELETE
 - BROWSE





- LOG N OPERATIONS
- LOCKS ARE USED TO GUARANTEE CONSISTENCY
- ON AVERAGE, THE B-TREE NODES ARE 3/4 FULL
- TRANSACTIONS ARE COMPLEX TO IMPLEMENT ...



- USED EVERYWHERE
 - DATABASE (INDEXES)
 - FILE SYSTEM
 - HSF+
 - NTF5
 - BTRF5
 - EXTY
 - REISER Y (B*TREE)



- STANDS FOR "MULTI-VERSION CONCURRENCY CONTROL"
- INVENTED BACK IN 1978 (DAVID P. REED)
- ALLOWS CONSISTENT CONCURRENT READS AND WRITES
- LOCKS CAN BE AVOIDED WHILE READING
- TRANSACTIONS CAN EASILY BE IMPLEMENTED

THINK ABOUT VERSION CONTROL SYSTEMS ...



APACHE DIRECTORY SERVER

- ENTRIES STORED IN A B-TREE
- INDEX STORED IN B-TREES
- IN-MEMORY B-TREES NEEDED
- TRANSACTION SUPPORT CROSS B-TREES
- MULTI-VALUE SUPPORT



• USING JDBM

- NO CROSS B-TREE TRANSACTIONS
- · NO LOCKS
- NO CRASH RECOVERY SYSTEM
- · NO BULK LOAD
- AVL TREES IN MEMORY
- LOCKS ALL OVER THE SERVER TO GUARANTEE CONSISTENCY



• BDB JE

- INCOMPATIBLE LICENSE ...
- · COMPLEX CONFIGURATION
- **HEAVY...**
- LMDB
 - C CODE BASE (IE, JNI ...)
 - FIXED SIZE INITIALIZATION
 - NEED A JAVA WRAPPER









• COUCHDB

- 2006 CONFERENCE
- SADLY, ERLANG CODE BASE
- HOWEVER, A PERFECT FIT
- JDBM LIMITS REACHED
 - ZOII : ADDING LOCKS ALL OVER THE SERVER
 - ZOIZ : IMPLEMENTING A CRASH RECOVERY SYSTEM RESULTING IN AWFUL PERFORMANCES
 - NO TRANSACTION (WELL, NO CROSS B-TREE TRANSACTIONS)
 - · NO BULK LOAD TOOL
 - DATABASE CORRUPTION



- MVCC
- CROSS B-TREES TRANSACTIONS
- CRASH RESISTANCE
- BULK-LOAD
- IN-MEMORY AND PERSISTED
- MULTI-VALUE SUPPORT
- FAST























• READ

- CURSOR BROWSE()
- CURSOR BROWSEFROM(K)
- · BOOLEAN CONTAINS(K, V)
- · V GET(K)
- · BOOLEAN HASKEY(K)

- WRITE
 - DELETE(K)
 - · INSERT(K, V)



• NAVIGATION

- AFTERLAST()
- BEFOREFIRST()
- NEXT()
- PREV()

ACTION

- CLOSE()
- **GET()**

- STATE
 - AVAILABLE()
 - HASNEXT()
 - · HASPREV()
 - ISCLOSED()



- ADDING A B-TREE :
 - CREATE ITS DEFINITION
 - INSERT IT IN THE RECORD MANAGER
 - CREATES THE B-TREE
 - STORES ITS FIRST REVISION IN THE BOB
 - ADDS IT TO THE LIST OF MANAGED B-TREES
- NO DELETION (ATM)



```
RecordManager recordManager = new RecordManager( "/tmp/mavibot.db" );
```

```
// Create a new BTree
try ( WriteTransaction transaction = recordManager.beginWriteTransaction() )
{
    btree = recordManager.addBTree(
        transaction,
        "test",
        LongSerializer.INSTANCE,
        StringSerializer.INSTANCE );
}
```



```
long[] values = new long[]
{
    14, 7, 43, 37, 49, 3, 20, 26, 17, 29,
    40, 33, 21, 18, 9, 30, 45, 36, 12, 8
};
try (WriteTransaction writeTxn = recordManager.beginWriteTransaction() )
{
    for ( long value : values )
    ł
        BTree<Long, String> btree = writeTxn.getBTree( "test" );
        btree.insert( writeTxn, value, Long.toString( value ) );
    }
}
```



```
try ( Transaction readTxn = recordManager.beginReadTransaction() )
{
  BTree<Long, String> btree = readTxn.getBTree( "test" );
  TupleCursor<Long, String> cursor = btree.browse( readTxn );
  while ( cursor.hasNext() )
  {
    cursor.next();
    Tuple<Long, String> tuple = cursor.get();
    System.out.println( "<" + tuple.getKey() + "," + tuple.getValue() + ">" );
  }
}
```



LET'S GET OUR HANDS

DIRTY...



MAVIBOT DATABASE IS STORED IN A SINGLE FILE

THIS FILE IS SPLIT IN FIXED SIZE PHYSICAL PAGES THAT MAY BE LINKED TOGETHER : PAGEIO





PHYSICAL PAGES (PAGEIO) ARE HOLDING THE CONTENT OF LOGICAL PAGES :

- LEAF
- NODE
- BTREEHEADER
- BTREEINFO





MAVIBOT LAYOUT : RMH

THE RECORDMANAGER HEADER 15

THE VERY FIRST PAGE

IT STORES REFERENCES TO MAVIBOT CONTENT :

- B-TREE OF B-TREES
- COPIED PAGES B-TREE
- FREE PAGES LIST
- ID COUNTER



RecordManager Header



- ONE SINGLE PAGE, CAN BE WRITTEN IN ONE I/O
- ALWAYS UPDATED AFTER EVERY REVISION CREATION OR VERSION CLEANUP
- STORES THE CURRENT STATUS :
 - CURRENT REVISION
 - CURRENT PAGE ID
 - NUMBER OF MANAGED B-TREES
 - · PAGE-IO SIZE



MAVIBOT LAYOUT : B-TREE

A MAVIBOT B-TREE 15 :

- A B-TREE HEADER PAGE
- A B-TREE INFO PAGE
- SOME PAGES
 - NODES
 - LEAVES





- THE B-TREE ENTRY POINT
- POINTS TO THE ROOT-PAGE
- POINTS TO THE B-TREE INFO
- CONTAINS THE NUMBER OF ELEMENTS
- CONTAINS THE CURRENT
 REVISION

	B-tree Header
Page ID	
Revision : N	
NbElems : M	
RootPage[N]	
BTreeInfoOffset	1233



THE B-TREE INFORMATIONS

- B-TREE NAME
- NB ELEMENTS PER PAGE
- KEY SERIALIZER
- VALUE SERIALIZER

IMMUTABLE

	B-tree Ir
pageNbElem	
nameSize	
B-tree name	- 1
keySerializerSize	
keySerializerFQCN	
valueSeri alizerSize	
valueSerializerFQCN	



Z FLAVORS OF PAGES :

- · LEAF : UP TO N KEYS AND VALUES
- NODE : UP TO N KEYS AND N+1 OFFSET

EVERY PAGE HAS :

- · A UNIQUE ID (A 64 BITS LONG)
- A REVISION
- THE NUMBER OF ELEMENTS STORED

Р	age
con	mmon
pagelD	
revision	
nbPageElems	
Node (nbPageElems < 0)	Leaf (nbPageElems >= 0)
Page[0] offset	Value[0]
Key[0]	Key[0]
Page[1] offset	Value[1]
Key[1]	Key[1]
Page[nbPageElems - 1] offset	Value[nbElems - 1]
Key[nbPageElems - 1]	Key[nbElems – 1]
Page[nbPageElems] offset	



- THIS B-TREE STORES A REFERENCE TO ANY MANAGED B-TREE (EXCEPT ITSELF AND THE COPIED PAGES B-TREE)
- WE CAN RETRIEVE A GIVEN B-TREE FROM IT
- NEW REVISIONS ARE STORED IN THIS B-TREE

READ AT STARTUP TO LOAD EACH B-TREE IN ITS LATEST VERSION



- WHEN WE CREATE A NEW VERSION, THE COPIED PAGES ARE STORED IN THIS B-TREE
- WHEN THE REVISION IS FREED, THE ASSOCIATED COPIED PAGES CAN BE RELEASED AND MOVED TO THE FREE LIST



- WE USE A LIST OF FREE PAGEIOS
- THEY ARE RE-USED ON DEMAND
- WHEN EMPTY, WE CREATE A NEW PAGEIO AT THE END OF THE FILE
- THE FREE PAGEIOS ARE LINKED
- WE CAN GRAB MORE THAN ONE PAGEIO
- WRITES USE ONE SINGLE THREAD, NO LOCK NEEDED FOR THIS LIST







VALUES/KEYS HAS TO BE (DE)SERIALIZED (BYTE] <-> OBJECT)





NODE'S CHILDREN ARE REFERENCED USING THEIR PAGE OFFSET

PAGES ARE CACHED (LRV)

BETTER HAVE A BIG CACHE !

WEAKREFERENCE 15 TOO SLOW ...

		r4								
Node			D		J		1			
	of	offset offset		offset		/		1		
cache			*							
			<u>r4</u>				1			
Leaf	A	В		С						
	Value ref	Value ref		Value ref		1				
	offset	of	fset	/	offset	1				
	VA									



NODE'S CHILDREN ARE REFERENCED USING THEIR PAGE OFFSET

PAGES ARE CACHED (LRU)

BETTER HAVE A BIG CACHE !

WEAKREFERENCE 15 TOO SLOW ...















- LOAD THE DATA FROM THE BOTTOM UP
- DATA MUST BE SORTED BEFOREHAND
 - · QUICKSORT IN MEMORY O(N LOG(N))
 - · MERGE SORT WHEN MEMORY IS LIMITED O(N LOG(N))
 - RESULT IS DENSE
- BIG TRANSACTIONS
 - SIMPLER BUT SLOWER
 - LESS DENSE







TRANSACTIONS



- READ TRANSACTIONS ARE CONCURRENT
- WRITE TRANSACTIONS ARE SERIALIZED
- EVERY COMMITTED WRITE TRANSACTION CREATES A NEW VERSION
- WRITE TRANSACTIONS CAN GATHER MORE THAN ONE OPERATION
- READ/WRITE TRANSACTIONS ARE CROSS-B-TREES







ONE VS MANY OPERATION/TXN





MAVIBOT + APACHEDS



- APACHEDS EMBED VERSION 1.0.0-M8
- TWICE FASTER THAN JDBM
- NO TRANSACTION SUPPORT IN 1.0.0-M8
- FREE PAGE MANAGEMENT IS NOT OPTIMAL
- MULTI-VALUE ELEMENTS SUPPORTED
- IN-MEMORY B-TREE SUPPORTED



- NO MULTI-VALUE SUPPORT : APACHEDS TO DEAL WITH IT
- CROSS-B-TREE TRANSACTION SUPPORT
- FREE PAGE MANAGEMENT ON THE FLY
- NO IN-MEMORY B-TREE
- FASTER THAN 1.0.0-M8
- CACHE



APACHEDS INDEX/B-TREES

- ENTRY
- APACHEDS INDEX
 - * PRESENT
 - * OBJECTCLASS
 - · * ENTRYUUID
 - * ENTRYCSN
 - ** RDN + REVERSE
 - ALIAS + REVERSE
 - ONEALIAS
 - SVBALIAS
 - ADMINROLE

		Obj	ectCl	ass index							
Ŗ	oerso	n	123	456, 38765	, 499,				-		
et	orgPe	rson	123	456,98805	,	_					
	_										
				cn index							
	3	jdoe 123456									
		acme 98805, 111024									
2		_						2			
		RDN index									
200	_	<dc< td=""><td>=apa</td><td>che,dc=org</td><td>, 0></td><td>1001</td><td></td><td></td><td></td><td></td><td></td></dc<>	=apa	che,dc=org	, 0>	1001					
		<	ou=p	eople, 100	1>	24512					
			<cn=< td=""><td>joe, 24512:</td><td>></td><td>123456</td><td></td><td></td><td></td><td></td><td></td></cn=<>	joe, 24512:	>	123456					
				R	ON rev	ert index					
			10	01	<dc=< td=""><td colspan="3"><dc=apache.dc=org. 0=""></dc=apache.dc=org.></td><td></td><td></td><td></td></dc=<>	<dc=apache.dc=org. 0=""></dc=apache.dc=org.>					
	-		24	512	<0	<ou=people, 1001=""></ou=people,>					
		123456			<	cn=joe, a	24512				
									7	7	
	2								123	456	
								-			
						ſ					
							cn=jo	ioe,ou=	people	,uc=apache,u	=org
						0	Object Object	Class: t Class: r	op Derson		
						C	Object	Class: i	netorg	Person	
						5	in: Joh	n Doe			
						9	n: Jol	nn Anache	- Softwa	are Foundation	
						C	. The	Apacin		ine roundation	







- 18 253 SLOCS
- LONG LASTING EFFORT (DAY JOB, FAMILY EVENTS ...)
- SOLVE A LOT OF ISSUES
- SHOULD BE COMPLETED SOON !







THANKS!