

Realization of DBS	External Storage Management	
Mapping of files & blocks File system Mapping of blocks Enhancing fault tolerance Mapping of tegments & pages Shadow page	 Mapping of files and blocks 2-level storage hierarchy General tasks of storage management File system Operations, addressing, and mapping Methods for block allocation Static Dynamic extent allocation Dynamic block allocation Measures to enhance fault tolerance Reads and writes of blocks 	
Differential files	 Use of storage redundancy Mapping of segments and pages Segment concept Deferred propagation strategies Shadow page mechanism Differential file Evaluation of page allocation methods 	3-2



Realization of DBS	Mapping of Files and Blocks – File Services	
	Example: call of an operation at the upper interface	
Mapping of files & blocks	PAM UPAMDAT, RD, FECB=ESTBLOCK, HP=1	
File system	UPAMDAT FCB FCBTYPE=PAM, LINK=EIN, IOAREA1=EINBER	
Mapping of blocks	Mapping function	
Enhancing fault tolerance Mapping of segments & pages	block number ——— cyl-#, track-#, rec-#	
Shadow page mechanism Differential files	■ I/O execution: Arm positioning: EXCP \rightarrow CCW Track selection: EXCP \rightarrow CCW Record transmission: EXCP \rightarrow CCW	
	 Properties of the upper interface Set of numbered blocks within files Read- and write access via block numbers Block accesses can result in read errors, blocks can contain old or invalid data 	
© 2011 AG DBIS	3-	4















Realization of DBS	Block Allocation – Optimization				
	Log-structured files				
Mapping of files & blocks	Proposal rather appropriate for applications having a multitude of small files, mostly with a short life time				
File system	 Avoidance of random reads using large main-memory buffers 				
Mapping of blocks	 Write optimization of modified blocks embodies the key idea. All modified blocks are sequentially written at the end of file processing (transaction), in a single batch, to the current end of the file which is organized like a log file and is overwritten in a cyclic way. Aged blocks are then released 				
Enhancing fault tolerance	Complex management of blocks and their metadata:				
Mapping of segments & pages;	concept cannot be easily transferred to DB files				
Shadow page mechanism	Catalog Catalog				
Differential files	Files: A B C D E F G H A B C D E F G H				
COLLAR DELE	Disadvantages (http://en.wikipedia.org/wiki/Log-structured_file_system): The design rationale for log-structured file systems assumes that most reads will be optimized away by ever-enlarging memory caches. This assumption does not always hold: On magnetic media—where seeks are relatively expensive—the log structure may actually make reads much slower, since it fragments files that conventional file systems normally keep contiguous with in-place writes. On flash memory—where seek times are usually negligible—the log structure may not confer a worthwhile performance gain, because write fragmentation has much less of an impact on write throughput. However many flash based devices can only write a complete block at a time because they must first perform a (slow) erase cycle before being able to write. 3-12				



Realization of DBS	Measures to Enhance Fault Tolerance					
	MTTF of different disk errors*					
Mapping of files & blocks	Type of Error	MTTF	Recovery	Consequences		
Mapping of	Soft data read error	1 hour	Retry or ECC (error correcting code)	None		
Enhancing fault tolerance	Recoverable seek error	6 hours	Retry	None		
Mapping of segments & pages	Maskable hard data read error	3 days	ECC	Remap to new sector and rewrite good data		
Shadow page mechanism	Unrecoverable data read error	1 year	None	Remap to new sector, old data lost		
Differential files	Device needs repair	5 years	Repair	Data unavailable		
	Miscorrected data read error	10 ⁷ years	None	Read wrong data		
	* Gray, J., Reuter, A.: Transactior	Processing – Co	oncepts and Techniques, Mo	organ Kaufmann, 1993.		
© 2011 AG DBIS				3-14		

Realization of DBS	Measures to Enhance Fault Tolerance (2)	
Mapping of files & blocks File system	 Simple read Read can be interrupted by transient errors Additional measure: unsuccessful read is repeated n times (to protect against transient errors) 	
Mapping of blocks Enhancing fault tolerance Mapping of segments & pages Shadow page mechanism	 Simple write Block write as an atomic action (either the entire block is correctly transferred to the specified slot or the slot remains unchanged) cannot be guaranteed Write actions can lead to wrong results because of transient and permanent errors Write error in the catalog? 	
Differential files	 After a simple write, the block is immediately read and compared with the original block Operation sequence is repeated until the block is successfully written Write is protected against transient errors. Permanent errors can, however, lead to wrong results 	
© 2011 AG DBIS		3-15

Realization of DBS	Measures to Enhance Fault Tolerance (3)
Mapping of files & blocks File system	 Duplexed write (stable write) Every block has a version number which is incremented upon each invocation of the operation The block is written in specified sequence in two different slots S_J and S_k Principle: stable storage
Mapping of blocks	block B ₁ Simple write: 1) then 2) (synchronous)
Enhancing fault tolerance Mapping of segments & pages	slots S _j 1) S _k 2) Atomicity for <u>a single</u> block Various disks,
mechanism	"far away"
Direrential files	 Assumption: block is not written to a wrong slot; otherwise, read-after-write is required Read of B_i takes place at first from slot S_j. If successful, it is assumed that the block is the youngest valid version of B_i If read of slot S_i fails, slot S_i is read
	 Because a system crash can only interrupt a single write, there is always a version of the block available at restart
© 2011 AG DBIS	 Hence, writes are protected against permanent errors. It is unexpected that both versions are not readable 3-16



Realization of DBS	Further Principles of Storage Redundancy (2)
	Detection of erroneous blocks
Mapping of files & blocks	 Protection against certain kinds of corruption: Parity bits, check sums Disk hardware can automatically figure out using parity bits, whether or not a sector was completely written
File system	→ Sufficient, if a block is fully stored in a sector (here identical to slot)
Mapping of blocks	Slot = sector
Enhancing fault tolerance	 Use of a single bit both in the first and last byte of a block Every time, identical values are assigned to both consistency bits
Mapping of segments & pages;	Multi-sector slots
Shadow page mechanism	 Erroneous block is only detected, if the block sequence (starting at block begin) is observed when the sectors are written
Differential files	 SCSI drives autonomously carry out a reordering of the sector writes to improve write performance
	 Check sums across entire blocks substantially reduce probability to miss a partial write; however, this error cannot fully be excluded (very expensive)
	 Corresponding to the n Sectors, a block is (logically) divided; from each of these fractions, a bit is taken and used as check bit. These n bits are filled with identical values and inverted upon each write
Datenbanken und Informationskysteme	\rightarrow What has to be done that these n bits do not falsify user data in the block?



Realization of DBS	Mapping of Segments and Pages	
	Segment concept	
Mapping of files & blocks	 Enables deferred propagation Allows to selectively introduce additional properties (attributes) 	
File system	e.g., to increase fault tolerance	
Mapping of	 Offers segments as units of locking, recovery, and access control 	
Enhancing	 Preserves the advantages of the file concept when appropriately mapped to files 	
fault tolerance		
Mapping of segments & pages	 DB buffer management (see chapter 4) Fix and Unfix of pages in the DB buffer 	
Shadow page mechanism	Preparation of I/O requests to the file services Outprint of mode account strate rise	
Differential files	 Optimization of replacement strategies Support of segments of different types 	
	New tasks: management of pages of variable size and of large objects	
	Division of logical DB address space in segments having visible page boundaries	
© 2011 AG DBIS	3-20	



Realization of DBS		Segr	nent Tyj	pes in a	DBMŜ			
	l	Cla	ssification	of segme	nt types			
Mapping of files & blocks		proper- ties	- segment- types	type 1	type 2	type 3	type 4	type 5
File system			usage	public	private	private	private	private
Mapping of blocks			life duration	permanent	permanent	permanent	permanent	temporary in a transaction
Enhancing fault tolerance	opening and closing		automatically by the system		explicitly by the user		e user	
Mapping of segments & pages Shadow page	recovery in case of a failure		automatic sys	ally by the tem	explicitly by no recovery mechanis the user		ery mechanism	
Differential files	 Examples characterizing the use of segment types Type 1: catalog, schema information, log, all shared DB parts Type 2: parts of the DB reserved for distinct users or user groups 							
		•	Type 3:	ocal copies of	parts of the	DB (views) fo	r specific us	ers (<i>snapshots</i>)
		•	Туре 4:	auxiliary files f	for user progr	ams		
Datenbalanskysteren Datenbalanskysteren Ø 2011 AG DBIS		•	Type 5:	temporary sto	rage, e.g., for	sort program	IS	3-22

















Realization of DBS	Shadow Page Mechanism – Functioning Principle
Mapping of files & blocks File system	 If segment k is to be opened for updates, the following steps have to be executed: Copy V_{k0} to V_{k1} State(k) := 1 Write Master in an uninterruptible operation Create a working copy CMAP of MAP₀ in main memory
Mapping of blocks Enhancing fault tolerance Mapping of segments & pages	 If page P_i is to be updated the first time since segment is opened, the following actions have to be executed: Read page P_i from Block j = V_{k0}(i) Find a free block j' in CMAP V_{k0}(i) = j' Mark page P_i as updated in V_{k0}(i)
Shadow page mechanism	 For further updates of P₁, block j' is used Ending an update interval (creating a checkpoint): Create the bit list in MAP₁ reflecting the current storage occupancy (new blocks occupied, old ones released) Write MAP₁ (no overwrite of MAP₀) Write V_{k0}
Contractions	 Write all modified blocks State(k) := 0, Mapswitch = 1 (MAP₁ is current) Write Master in an uninterruptable operation Rollback of opened segments only requires to copy V_{k1} into V_{k0} and to set State(k) to 0 3-31



Realization of DBS	Evaluation of the Shadow Page Mechanism
	Advantages
Mapping of files & blocks	 Rollback to most recent consistent state is very simple More flexible propagation protocols for log files: buffering is possible until switching to the new state
File system	 Logical logging is possible, because an operation-consistent state is always available
Mapping of blocks	 In case of disastrous failures, a higher probability exists to reconstruct a useful DB state
fault tolerance	- Disadvantages
Mapping of	
segments & pages	 Auxiliary structures (MAP and page table V_i) become so large that block decomposition is necessary
Differential files	 Page tables V_i occupy about 0.05-0.1% of the DB size, which, in case of large DBs (≥ n TB), may lead to a large share of page faults in the DB buffer provoked by accesses to V_i (argument may be weakened in future)
	 Periodical checkpoints enforce propagation of all modified pages in the DB buffer
	• Physical clustering of logically related pages is impaired resp. annihilated
	 Additional storage space for double occupancy: therefore, long batch (update) programs are not well supported
© 2011 AG DBIS	3-33







Realization of DBS	Summary
Mapping of files & blocks	 Storage allocation structures require efficient file concept Many files are of varying, not statically fixed size Dynamic growth and shrinking is required Dermonent and temperary files
File system Mapping of blocks Enhancing fault tolerance	 Permanent and temporary files Recommended file properties: Direct and sequential block access Block size should be defined on a file basis Block allocation via dynamic extents
Mapping of segments & pages Shadow page mechanism	 Dynamic block allocation (in UNIX represented as a multi-level tree) not appropriate for large files Segment concept
Differential files	Enables the implementation of additional properties for DB processing (recovery, clustering for tables, etc.)
	 Two-level mapping Of segment/page onto file/block and these to slots on disk enables introduction of mapping redundancy through deferred propagation
© 2011 AG DBIS	3-37

Realization of DBS	Summary (2) and Literature
Mapping of files & blocks File system Mapping of blocks Enhancing fault tolerance	 Deferred propagation strategies Are more expensive than direct ones, own, however, implicit fault tolerance They burden normal processing in favor of recovery Direct propagation strategy (update-in-place) Simple implementation No additional costs for page allocation at runtime Fault tolerance only through explicit logging&recovery functions
Mapping of segments & pages Shadow page mechanism Differential files	 Lorie, R.: Physical Integrity in a Large Segmented Database. ACM Transactions on Database Systems (TODS) 2(1): 91-104 (1977) Rahm, E.: Hochleistungs-Transaktionssysteme, Vieweg, 1993, Ch. 5 and 6, http://lips.informatik.uni-leipzig.de/pub/1993-2 Rosenblum, M., Ousterhout, J.K.: The Design and Implementation of a Log-Structured File System. ACM TODS 10(1): 26-52 (1992)
© 2011 AG DBIS	3-38