

NoFTL: Database Systems on FTL-less Flash Storage

Sergej Hardock
Databases and Distributed
Systems Group
TU-Darmstadt, Germany
hardock@dvs.tu-
darmstadt.de

Iliia Petrov
Data Management Lab
Reutlingen University,
Germany
ilia.petrov@reutlingen-
university.de

Robert Gottstein
Databases and Distributed
Systems Group
TU-Darmstadt, Germany
gottstein@dvs.tu-
darmstadt.de

Alejandro Buchmann
Databases and Distributed
Systems Group
TU-Darmstadt, Germany
buchmann@dvs.tu-
darmstadt.de

ABSTRACT

The database architecture and workhorse algorithms have been designed to compensate for hard disk properties. The I/O characteristics of Flash memories have significant impact on database systems and many algorithms and approaches taking advantage of those have been proposed recently. Nonetheless on system level Flash storage devices are still treated as HDD compatible block devices, black boxes and fast HDD replacements. This backwards compatibility (both software and hardware) masks the native behaviour, incurs significant complexity and decreases I/O performance, making it non-robust and unpredictable. Database systems have a long tradition of operating directly on RAW storage natively, utilising the physical characteristics of storage media to improve performance.

In this paper we demonstrate an approach called *NoFTL* that goes a step further. We show that allowing for native Flash access and integrating parts of the FTL functionality into the database system yields significant performance increase and simplification of the I/O stack. We created a real-time data-driven Flash emulator and integrated it accordingly into Shore-MT. We demonstrate a performance improvement of up to 3.7x compared to Shore-MT on RAW block-device Flash storage under various TPC workloads.

1. INTRODUCTION

Many key database architectural principles and workhorse algorithms have been designed to leverage the properties of HDD. Flash memories are a new technology crucial to database systems, which comes with a set of different I/O characteristics. A large body of algorithmic approaches has

been proposed over the last years to natively address Flash properties. Nonetheless on system level, Flash devices still support the same block level interface as HDD. On the one hand, the block device compatibility favours adoption by making replacement seamless. On the other hand, as a legacy interface, it is a major source of unpredictability, non-robustness. The negative performance impact ultimately precludes any Flash relevant optimisations.

The *Flash Translation Layer (FTL)* is an on-device layer that ensures low-level block interface compatibility, masking physical characteristics, and making a Flash device behave like a hard drive [6, 4]. Some of the negative FTL aspects are: (i) Unpredictable and state-dependent performance due to background processes [6, 5]; (ii) adverse performance impact due to limited on-device computational resources [5, 12]; (iii) redundant functionality also present at different layers along the critical I/O path [5, 12].

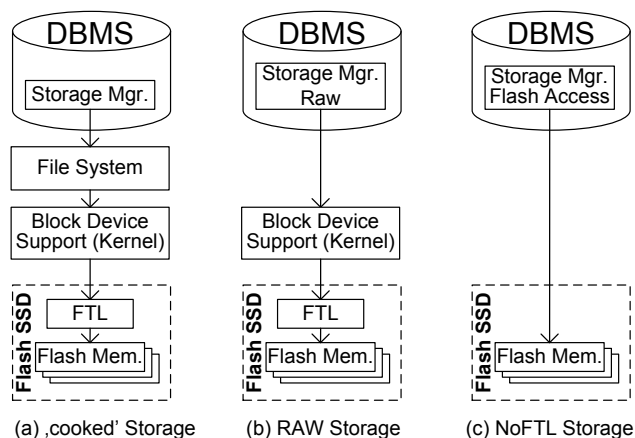


Figure 1: DBMS storage alternatives: (a) Traditional 'cooked' DBMS storage; (b) RAW DBMS storage; (c) NoFTL

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Articles from this volume were invited to present their results at The 39th International Conference on Very Large Data Bases, August 26th - 30th 2013, Riva del Garda, Trento, Italy. *Proceedings of the VLDB Endowment*, Vol. 6, No. 12. Copyright 2013 VLDB Endowment 2150-8097/13/10... \$ 10.00.

block devices (Figure 1.a). Database systems on raw storage (Figure 1.b) eliminate file system overhead, enable raw storage access and direct physical data placement, achieving better performance [17]. Newer approaches propose departing from block device interfaces, achieving: atomic writes, computational efficiency and parallelism [19], stripped down FTL and a native interface to host [5, 12]. With *NoFTL* (Figure 1.c) we stop treating Flash devices as a closed system, consider native Flash access, and explore FTL integration in the DBMS.

Contributions and Demo. In this paper we demonstrate an approach called *NoFTL*. We argue for a significant simplification of the I/O stack; integrating Flash management in the database and using DBMS knowledge to control storage; direct access and exposure of a native Flash interface; utilisation of database server’s computational resources instead of on-device resources. The contributions of this paper are: (i) we implemented a real-time data-driven Flash emulator as a character device driver; (ii) we incorporated a DFTL implementation; (iii) we extended Shore-MT with a page mapping FTL and integrated the real-time simulator; (iv) live TPC-C, TPC-B and TPC-H tests under Shore-MT indicate a *NoFTL* performance improvement of 1.5x to 3.7x over the DFTL configuration.

2. RELATED WORK

In the past numerous designs of FTLs have been proposed (e.g. [20], [10], [14], [13], [15], [16] etc.). Such approaches can be classified as Page-, Block- or Hybrid-/Log-Block- Mapping FTLs. An evaluation and comparison of different FTLs is provided in [6] and [7]. DFTL is a page-mapping FTL and is introduced in [10]. There are multiple Flash simulation frameworks such as FlashSim [11] or DiskSim. There is further research on omitting certain on-device FTL functionalities, e.g. an approach that is not using the block I/O interface is presented in [19], [5] presents a hybrid approach which can bypass the on-device FTL. Specialized Flash Server Storage moves the FTL from a device into the driver, such as FusionIO [1]. *NoFTL* completely removes the on-device FTL, enabling the application to take full control of the Flash storage device.

3. THE NOFTL APPROACH

At the core of the *NoFTL* evaluation and demonstration (Figure. 2) is a real-time data-driven Flash emulator, simulating a Flash device according to the ONFI standard, while storing the data in a large RAM buffer. The emulated Flash device is attached to Shore-MT [3], which is a recognised storage engine supporting ACID transactions, ARIES-type logging, Indices, Buffer management. Furthermore, Shore-MT supports raw devices and standard TPC benchmark implementations.

The *NoFTL* emulator exposes two interfaces: a native Flash interface (Figure 2.a) and a block device interface (Figure 2.b). In the former case (*NoFTL*) we also extended Shore-MT with a Page-Mapping FTL, including a DatabasePage-to-PhysicalPage mapping, integrating Flash space management and wear levelling in the Shore-MT storage manager. The latter case required an FTL implementation to ensure block device functionality. DFTL [10] was chosen since it provides better results than most hybrid FTL schemes

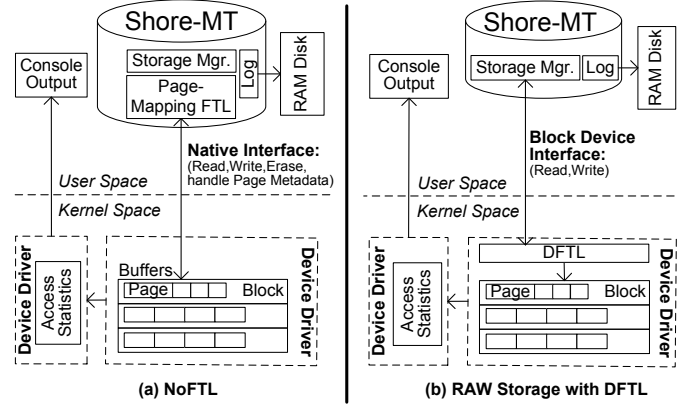


Figure 2: Architecture

proposed in the recent years. Both alternatives are implemented as device drivers. We gather reliable device statistics not interfering with the real-time simulator by using an additional device driver, in both configurations.

The Flash device emulation is based on the observation that typical Flash latencies (of $25\mu\text{s}$ and $250\mu\text{s}$ write of a 4KB page and $700\mu\text{s}$ for a block erase – Micron MT29F16G08ABABA) provide enough CPU time to perform the emulation. The device size and layout are configurable and designed according to the ONFI specification: a Flash device contains a number of chips comprising multiple planes (between 1 and 4), which in turn comprise multiple dies (typically 4). The current *NoFTL* emulator architecture and implementation allows for very accurate timings and robust performance. We validated the *NoFTL* emulator under different workloads using I/O benchmarking tools (*FIO*) and against of-line Flash Simulators (FlashSim, DiskSim).

Coupling the emulator to Shore-MT is the second cornerstone of the demonstrated approach. For *NoFTL* (Figure 2.a) to work, typical FTL functionality was integrated into Shore-MT: (i) page-mapping, wear levelling, and garbage collection; (ii) integration of the native Flash interface (native use of *read*, *write* and *erase* at the respective granularity); (iii) eliminating redundant functionality.

The experimental analysis was performed on an Intel Xeon server with two quad-core Intel Xeon 5630 2.5 GHz processors (256 KB L1 cache, 1 MB L2 cache and 12 MB L3 cache) and 48 GB RAM and a QPI bus architecture. We instrumented Shore-MT for both TPC-C (scale factors: 5, 25 and 50), TPC-B (scale factors: 48, 200, 500) and TPC-H (scale factor: 1) on different Flash volume sizes (1 GB, 5GB and 10GB) - Flash volume equals DB size. We compare the *NoFTL* (Figure 2.a) against the RAW Flash with *DFTL* (Figure 2.b). DFTL was configured with different mapping cache sizes (2%, 15%, 35% of all mappings cached).

The experimental results are shown in Figure 3. *NoFTL* is 1.5x to 3.7x faster than comparable *DFTL* configurations. The speedup increases with the database size. Three major factors contribute to the speedup: (a) limited on-device computational resources; (b) FTL restrictions; (c) redundant functionality. These are discussed in detail in the paragraphs below.

Limited on-device computational resources: DFTL and the FTL in general is executed on slow on-device hardware. In the emulator we account for the physical I/O required to page in and page out the physical-to-logical address mapping

TPC-C: SF=5 TPC-B: SF=48 DB SIZE: 1GB	NoFTL		DFTL 2%			DFTL 15%			DFTL 35%		
	Tx/s	STDEV	Tx/s	STDEV	Slow-down	Tx/s	STDEV	Slow-down	Tx/s	STDEV	Slow-down
TPC-C	43.34	0.34	26.64	0.35	1.63	25.40	0.34	1.71	30.00	0.87	1.44
TPC-B	462	1.05	275.7	3.70	1.67	254	2.84	1.82	295.9	2.26	1.56

TPC-C: SF=25 TPC-B: SF=200 DB SIZE: 5GB	NoFTL		DFTL 2%			DFTL 15%			DFTL 35%		
	Tx/s	STDEV	Tx/s	STDEV	Slow-down	Tx/s	STDEV	Slow-down	Tx/s	STDEV	Slow-down
TPC-C	41.4	0.22	22.7	0.15	1.82	21.5	0.09	1.92	19.0	0.17	2.18
TPC-B	449.7	7.87	267.2	6.70	1.68	233.0	2.46	1.93	149.4	10.4	3.01

TPC-C: SF=50 TPC-B: SF=500 DB SIZE: 10GB	NoFTL		DFTL 2%			DFTL 15%			DFTL 35%		
	Tx/s	STDEV	Tx/s	STDEV	Slow-down	Tx/s	STDEV	Slow-down	Tx/s	STDEV	Slow-down
TPC-C	41.2	0.38	19.1	0.31	2.16	13.8	0.13	2.98	10.9	0.42	3.77
TPC-B	409.7	3.01	215.5	0.94	1.90	163.3	15.0	2.51	113.7	9.97	3.60

TPC-H: SF=1 DB SIZE: 1.5GB	NoFTL		DFTL 2%			DFTL 15%			DFTL 35%		
	Sec.	STDEV	Sec.	STDEV	Slow-down	Sec.	STDEV	Slow-down	Sec.	STDEV	Slow-down
Query 1	58.7	0.54	79.1	0.05	1.35	121.9	0.27	2.08	178.2	3.75	3.04
Query 6	58.4	0.42	75.2	0.05	1.29	98.6	0.11	1.69	149.5	3.77	2.56
Query 12	109.3	11.36	120.5	1.80	1.10	157.1	5.94	1.44	208.0	5.56	1.90
Query 14	98.4	2.72	113.5	0.63	1.15	142.2	2.9	1.45	194.7	3.05	1.98

Figure 3: TPC-B, TPC-C and TPC-H results NoFTL and DFTL

table as well as the computational overhead. On real devices the less powerful on-device CPU will have an even more negative impact. In *NoFTL* scenarios the FTL fully benefits from the DB server’s ample computational resources.

FTL restrictions: SSD vendors offer sparse details about the implemented algorithms. Research findings converge towards hybrid FTL schemes. We opted for DFTL (a page-mapping approach), which represents an optimistic choice: DFTL wastes less paging I/Os for the mapping tabs and has more efficient garbage collection than Hybrid FTL schemes. Nonetheless, for large mapping buffers the DFTL I/O savings do not compensate for the high computational overhead incurred by the mapping table cache maintenance, clearly visible for large data (5 or 10GB and 35% cached mappings) – Figure 3.

Redundant functionality: in terms of address mapping, space management, page and block placement database storage managers, file systems and FTL schemes contain similar functionality. Eliminating some layers of abstraction and integrating functionality is where *NoFTL* has most potential. Due to the Shore-MT integration we managed to: (i) reduce the number of block erases and make them Flash-friendly while avoiding TRIM; (ii) couple garbage collection to Shore-MT space management and eliminate file system space management functionality; (iii) simplify wear-levelling with DB information about dirty block eviction; (iv) buffering and buffer management is consolidated and performed only on DBMS level (not by the DBMS, File System, OS Kernel and SSD); (v) address mapping is performed only on DBMS level (not on DBMS (Tuple-Block), File System (FileOffset-LBA) and SSD level (LBA-PBA)). Many possible further extensions result from integration with MVCC and Log-based Storage Managers [8], access paths [9], buffer management and eviction strategies [18], etc.

4. DEMONSTRATION DESCRIPTION

In this section we describe the demonstration of the main features of *NoFTL*. We also describe the main scenarios and how the audience can interact with the system.

High-level Description. The main scenario involves instrumentation and comparative testing of NoFTL versus DFTL devices. We introduce the audience to the system and explain what the expected influence of the different knobs is. We then let the audience pick a test scenario configuration and perform comparative benchmarking.

Entry-level I/O scenario. We let the audience stress the NoFTL emulator with a simple I/O benchmarking tool - *FIO*. The audience will experience the different statistics for different metrics e.g. reads, writes, overwrites, IOPS, etc. and the influence of different system parameters on the performance (see Figure 4).

```

root@timbuktu-Aspire: /home/timbuktu/simulator/flashsim
fio > .fio /home/timbuktu/simulator/flashsim/10_tests/random_write_block.fio
random_write_block: (g=0): rw=randwrite, bs=4K-4K/4K-4K, ioengine=sync, todepth=1
fio-2.0.8
Starting 1 process
Jobs: 1 (f31): [w] [100.0% done] [0K/15916K /s] [0 /3979 iops] [eta 00m:00s]
random_write_block: (groupid=0, jobs=1): err=0, pid=8712: Fri Mar 22 10:11:18 2013
write: io=990.0MB, bw=15909KB/s, iops=3977, runt= 63721msec
    clat (usec): min=249, max=316, avg=250.25, stdev= 1.70
    lat (usec): min=249, max=316, avg=250.40, stdev= 1.72
    clat percentiles (usec):
    | 1.00th=[ 249], 5.00th=[ 251], 10.00th=[ 251], 20.00th=[ 251],
    | 30.00th=[ 251], 40.00th=[ 251], 50.00th=[ 251], 60.00th=[ 251],
    | 70.00th=[ 251], 80.00th=[ 251], 90.00th=[ 251], 95.00th=[ 251],
    | 99.00th=[ 255], 99.50th=[ 255], 99.90th=[ 282], 99.95th=[ 294],
    | 99.99th=[ 306]
    bw (KB/s): min=15888, max=15920, per=100.00%, avg=15913.64, stdev= 5.99
    lat (usec): 250=3.67%, 500=96.33%
    cpu: usr=0.51%, sys=98.42%, ctx=5374, majf=0, minf=27
    IO depths : 1=100.0%, 2=0.0%, 4=0.0%, 8=0.0%, 16=0.0%, 32=0.0%, >=64=0.0%
    submit   : 0=0.0%, 4=100.0%, 8=0.0%, 16=0.0%, 32=0.0%, 64=0.0%, >=64=0.0%
    complete : 0=0.0%, 4=100.0%, 8=0.0%, 16=0.0%, 32=0.0%, 64=0.0%, >=64=0.0%
    issued   : total=r=0/w=253440/d=0, short=r=0/w=0/d=0

Run status group 0 (all jobs):
WRITE: io=990.0MB, agrb=15909KB/s, minb=15909KB/s, maxb=15909KB/s, mint=63721msec

```

Figure 4: Screenshot testrun FIO on DFTL

Shore-MT scenarios. We then let the audience to pick a TPC benchmark (TPC-B, TPC-C or TPC-H) and dataset size depending to the personal preference.

Phase I. A device layout is selected and the device is being initialised (Figure 5). In addition, the benchmark data set is loaded and the system is prepared for benchmarking.

```

root@timbuktu-Aspire: /home/timbuktu/sim
#!/bin/sh
insmod flashsim.ko \
\
block_device_major_number=0 \
block_device_sector_size=4096 \
block_device_num_blocks=4096 \
block_device_sectors_per_block=64 \
block_device_delay_mode=1 \
block_device_request_statistics_mode=0 \
block_device_max_req_stat_entries=1000000 \
block_device_read_sector_latency_us=50 \
block_device_program_sector_latency_us=250 \
block_device_erase_block_latency_us=700 \
block_device_oob_size=32 \
block_device_read_oob_latency_us=0 \
block_device_program_oob_latency_us=0 \
block_device_ftl_mode=1 \
block_device_cache_level_1=4080 \
block_device_cache_level_2=1020 \
block_device_num_extra_blocks=116 \
\

```

Figure 5: Device Layout Instrumentation

Phase II. Having loaded the data into Shore-MT, the chosen TPC benchmark is run under the NoFTL and DFTL setups for 10 min. The audience tracks the execution progress both in terms of Shore-MT performance statistics (Figure 6)

but also in terms of low-level device and request statistics (Figure 7). A preview of this stage is provided in [2].

```

root@timbuktu-Aspire: /home/timbuktu
linux-mon: ./src/util/procstat.cpp:254:print_interval: (-nan) (30.0)
linux-mon: ./src/util/procstat.cpp:254:print_interval: (-nan) (22.0)
linux-mon: ./src/util/procstat.cpp:254:print_interval: (-nan) (24.0)
linux-mon: ./src/util/procstat.cpp:254:print_interval: (-nan) (31.0)
linux-mon: ./src/util/procstat.cpp:254:print_interval: (-nan) (30.0)
linux-mon: ./src/util/procstat.cpp:254:print_interval: (-nan) (29.0)
root-thread: ./src/tests/shore_kits.cpp:442:_cmd_MEASURE_impl: end mea
root-thread: ./src/workload/tpcc/shore_tpcc_xct.cpp:131:print_throughp
QuertedSF: (5.0)
Spread: (No)
Threads: (1)
Trxs Att: (1713)
Trxs Abt: (13)
Trxs Dld: (0)
NOrd Com: (785)
Secs: (60.00)
IOChars: (0.00M/s)
AvgCPUs: (-nan) (-nan%)
TPS: (28.33)
tpm-C: (785.00)

```

Figure 6: Screenshot testrun Shore-MT on DFTL

```

root@timbuktu-Aspire: /home/timbuktu/simulator/flashsim
ERASE REQUESTS...
ERASE[DATA_UNIT]: 4186
ERASE[MAP_UNIT]: 7883
ERASE: 12069
***** END of GENERAL STATISTICS *****
***** REQUEST STATISTICS *****

```

IO #	RW	LSN	SIZE	T.DURATION	READ	WRITE	OOB_R	OOB_W	ERASE	PSN	F.DURATION
0	R	0	1	53	2	0	0	0	0	0	51
1	R	1	1	52	2	0	0	0	0	1	51
2	R	3	1	52	2	0	0	0	0	3	50
3	R	254448	1	58	2	0	0	0	0	254448	53
4	R	254462	1	53	2	0	0	0	0	254462	50
5	R	0	1	4	1	0	0	0	0	0	1
6	R	1	1	3	1	0	0	0	0	1	1
7	R	254463	1	54	2	0	0	0	0	254463	51
8	R	254431	1	54	2	0	0	0	0	254431	51
9	R	254456	1	52	2	0	0	0	0	254456	50
10	R	254432	1	51	2	0	0	0	0	254432	50
11	R	254414	1	51	2	0	0	0	0	254414	50
12	R	256	1	52	2	0	0	0	0	256	50
13	R	3	1	2	1	0	0	0	0	3	1
14	R	7	1	51	2	0	0	0	0	7	50
15	R	15	1	51	2	0	0	0	0	15	50
16	R	2	1	51	2	0	0	0	0	2	50
17	R	16	1	52	2	0	0	0	0	16	50
18	R	8	1	53	2	0	0	0	0	8	50
19	R	32	1	55	2	0	0	0	0	32	51
20	R	33	1	52	2	0	0	0	0	33	50

Figure 7: Screenshot NoFTL Request Statistics

5. CONCLUSIONS

In this paper we demonstrated an approach called *NoFTL*. We argue for a significant simplification of the I/O stack; integrating Flash management in the database and using DBMS knowledge to control storage; direct access to storage and exposure of native Flash interface. We also argue that the performance gain can be maximised by even closer DBMS integration. *NoFTL* is implemented as a real-time data-driven Flash emulator. We integrated it into Shore-MT, which was also extended with a page mapping FTL. As a comparative FTL-system, we incorporated a DFTL implementation. We demonstrate live TPC tests under Shore-MT indicating a *NoFTL* performance improvement of up to 3.7x.

Acknowledgements

We are grateful to Goetz Graefe for his comments and suggestions on earlier drafts of this paper. This work was supported by the DFG (Deutsche Forschungsgemeinschaft) project “Flashy-DB”.

6. REFERENCES

[1] Going beyond ssd: The fusionio software defined flash memory approach, 2013.
www.fusionio.com/white-papers/beyond-ssd/.

[2] Noftl demo preview.
http://dblab.reutlingen-university.de/tl_files/downloads/demo_noftl.mp4, 2013.

[3] Shoremt storage engine.
<http://diaswww.epfl.ch/shore-mt/>, 2013.

[4] N. Agrawal and e. A. Prabhakaran. Design tradeoffs for ssd performance. In *Proc. ATC'08*, pages 57–70, 2008.

[5] P. Bonnet, L. Bouganim, I. Koltsidas, and S. D. Viglas. System co-design and data management for flash devices. In *Proc. VLDB 2011*, 2011.

[6] F. Chen, D. A. Koufaty, and X. Zhang. Understanding intrinsic characteristics and system implications of flash memory based solid state drives. In *Proc. SIGMETRICS '09*, pages 181–192, 2009.

[7] T.-S. Chung, D.-J. Park, S. Park, D.-H. Lee, S.-W. Lee, and H.-J. Song. A survey of flash translation layer. *J. Syst. Archit.*, 55(56):332 – 343, 2009.

[8] R. Gottstein, I. Petrov, and A. Buchmann. Append storage in multi-version databases on flash. In *Proc. of BNCOD*, 2013.

[9] G. Graefe. Write-optimized b-trees. In *Proc. VLDB'04*, pages 672–683, 2004.

[10] A. Gupta, Y. Kim, and B. Urgaonkar. Dftl: a flash translation layer employing demand-based selective caching of page-level address mappings. In *Proc. ASPLOS*, ASPLOS XIV, pages 229–240, 2009.

[11] Y. Kim, B. Tauras, A. Gupta, and B. Urgaonkar. Flashsim: A simulator for nand flash-based solid-state drives. In *In Proc. SIMUL'09*, pages 125–131, 2009.

[12] I. Koltsidas and S. D. Viglas. Data management over flash memory. In *Proc. SIGMOD '11*, 2011.

[13] S.-W. Lee, W.-K. Choi, and D.-J. Park. Fast: An efficient flash translation layer for flash memory. In X. e. a. Zhou, editor, *Emerging Directions in Embedded and Ubiquitous Computing*, LNCS. 2006.

[14] S.-W. Lee, D.-J. Park, T.-S. Chung, D.-H. Lee, S. Park, and H.-J. Song. A log buffer-based flash translation layer using fully-associative sector translation. *ACM TECS*, 6(3), July 2007.

[15] S.-P. Lim, S.-W. Lee, and B. Moon. Faster ftl for enterprise-class flash memory ssds. In *SNAPI*, 2010.

[16] D. Ma, J. Feng, and G. Li. Lazyftl: a page-level flash translation layer optimized for nand flash memory. In *Proc. SIGMOD '11*, 2011.

[17] Oracle. A quantitative comparison between raw devices and file systems for implementing oracle databases. white paper. 2004.
www.oracle.com/technetwork/database/performance/twp-oracle-hp-files-130020.pdf.

[18] Y. Ou, J. Xu, and T. Härder. Towards an efficient flash-based mid-tier cache. In *DEXA*, pages 55–70, 2012.

[19] X. Ouyang, D. W. Nellans, R. Wipfel, and D. Flynn. Beyond block i/o: Rethinking traditional storage primitives. In *HPCA*, pages 301–311, 2011.

[20] C. Park, W. Cheon, J. Kang, K. Roh, W. Cho, and J.-S. Kim. A reconfigurable ftl architecture for nand flash-based applications. *TECS*, 7(4):38:1–38:23, 2008.