# Studying Execution Modes and Wear Leveling in Flash Memorie

G Anjaneyulu, Syed YousufUddin, Md Ismail

**Abstract:**
The impact of wear levelling on a Flash storage pack- age and its access operations' execution modes is in- vestigated. First, a simple, static logical to phys- ical mapping functions are proposed and their im- plied wear levelling is assessed for different address distributions, covering both unifrom access and hot spots, as well as the Flash chip utilisation within the whole package. Second, for the access execu- tion modes, different preemptive and non-preemptive priority schemes are considered with a range of IO arrival rates using Poisson, Erlang, Pareto and Geometric-based arrival processes. The analysis of the impact of the execution modes on the perfor- mance of the Flash memory is undertaken using a hardware simulator. The results obtained show clearly the good wear levelling obtained by the map- ping functions, even in presence of hot spots. In ad- dition, the effect of the chosen execution mode on the whole storage package for each IO workload type is clearly analysed and accurately quantified.

**Keywords**: Flash memory, Wear levelling, Priority, Preemption, Waiting time, IO performance, Chips utilisation.

## Introduction

StoragedevicesbasedonFlashmemoryarebe-comingmoreand moreprevalentinour dailylife.Thisrecenttechnologypresentsapanoplyofde-vices, continually undergoing intensive evolution inresponse to market demand for MP3 players, mobilephones,digitalcamerasusingrawFlashdevicesandforlightweightlaptopcomputers,recently evendesktopcomputersusingFlash-baseddevicesinothertermsSolidStateDrives(SSD).Infact,their use is covering both consumer and enterprisestorage products replacing Hard Drive Disks (HDD),pushing them to archiving purpose [1].Since Flashtechnology is so widely used, its performance shouldbe precisely quantified and its impact on the wholesystem, in which it is embedded, assessed relative toIO profiles and its execution mode.However, thereare currently few studies providing such information,which cannot just be deduced from the behaviouralanalysis of other storage devices such as Hard DiskDrives(HDD)andmemories(SRAM,DRAM...etc)becausetheiraccessoperationsarecompletelydifferent [2].Some studies to adapt or/and proposeaccessalgorithmsforFlash-baseddeviceswereachieved [3, 4, 5] but they still restricted to specificapplications.

The main specific features of Flash memory mak-ing it so different from the other storage devices arethe limited erase cycles and the big disparity betweentheoperations access times. The firstone associatea fixed lifetime to the Flash chip since its manufac-turing,dependingon its typeand density.In orderto maximise this Flash longevity an even erase oper-ations distribution, thus a good wear levelling shouldbeguaranteed.Inthefirstpartofthispaperandprior to giving a quantitative characterisation of thesystemunder study,weconsidernexttheeffectoftheaddress mapping chosen to provide good wear level-ling and we propose a simple static map that guar-antees data availability, and hence maximum devicelongevity. The second feature deals with the fact thatservicetimesofthethreeFlashmemoryaccessoper-ations(read/write/erase)areconstantbutpresentasignificant disparitywhateverthechiptype.Eraseoperations,beingcostlyintime,introducelongde-laysforwaitingreadorwriteoperationsperformedafterthem.

Department of civil
g.anjaneyulu143@gmail.com, yousufsyed@gmail.com, mdismail786@gmail.com
ISL Engineering College.
International Airport Road, Bandlaguda, Chandrayangutta Hyderabad - 500005 Telangana, India.

These delays significantly change the delivered performance and make workload scheduling crucial. On the other hand, applications require good organisation of their workload to consistently realise faster access without having to customise in every specific context for each access profile. Whilst true for any storage device technology, it is particularly important for Flash because its access time is largely location-independent, especially for reads. It is complicated and probably unnecessary to include this kind of workload access optimisation in the application layers or to add a software layer to schedule the workload generated by applications according to the Flash operation modes. An easy way to meet good performance without any scheduling of requests, is to choose the most appropriate execution mode from the most basic provided ones (priority, preemption) according to the input IO profile. In the second part of this paper, we consider the effect of the workload on Flash delivered performance using different IO in-ter arrival time distributions. We focus on the three main execution modes: without priority among the three access operations (mode1), giving priority to reads in a non-preemptive policy (mode2) and finally giving priority to reads with preemption (mode3). In the rest of the paper, section 2 gives a succinct presentation of Flash technology background. Section 3 is dedicated to the wear levelling study. It describes the used tools, presents the proposed mapping functions and details their validation. Section 4 is dedicated to the analysis of the IO profiles and the operations execution modes impact on the performance. Section 4.3 presents the obtained numerical results and discussion their significance. Finally, section 5 summarises our main conclusions and suggests directions for future work.

## Background

Flash memory is considered a major, non-volatile mass-storage component due to its shock resistance,

vibration tolerance, light weight and low energy consumption; not to mention its high capacity. It is al-ready used in a wide range of applications and environments, from daily entertainment, e.g. in MP3 players, through personal computers, for web servers machines [6] and critical systems such as satellite systems [7].

There are two types of Flash memory, labelled according to its construction: NOR and NAND. The former has lower density and higher cost but provides fast random access and can be easily re-programmed, making it most suitable for storing code. Another ad-vantage of NOR is its lower susceptibility to corruption than NAND, partly because of the bad blocks that exist in the latter from the time of manufacture. NAND Flash, on the other hand, has a very large storage capacity and provides fast data access for large read/write requests, making it most suitable for storing data [8]. This is the most widespread and the one we consider. In fact the density is about 8 times more for NAND [9], at a cost that is 4 to 8 times cheaper than NOR. Although erases are signif-icantly faster on NOR, they can be pre-scheduled in NAND, essentially running in a garbage collector. In addition, $MLC_n$ (MultiLevelCell) technology multiplies the storage capacity of the Flash memory chip by having $n$-bit information per cell. $MLC_2$ becomes a classical device, present in most mobile components such as cameras and smart-phones. The first $MLC_3$ was developed by Hynix Semiconductor in 2008, followed by Samsung in 2009 to produce initially mi-croSD cards and to support more competitive high density consumer electronics storage solutions in the near future.

A NAND Flash memory chip is composed of a fixed number of blocks, each of which is partitioned into a fixed number of pages. Every page consists of two areas: a data area for native (user) data and a spare area for data status information (figure1). A block is the erase operation's unit of storage, whilst a page is the read and write operation's unit. No 'in-place' updates are allowed in NAND Flash. When data is modified, the new version must be written to an available page – called the *live page*. The page containing the old version is considered a dead page and is invalidated. As time passes, ficient garbage collection process, using a relatively reduced mounting time. Recently, comes UBIFS [17] for Unsorted Block Images File System, designed by Nokia for Flash-based devices as Solid State Drives (SSD). It provides faster mounting time and good wear levelling comparing to the JFFS2 random one. The second class of file systems are designed to work under any file system. We can cite YAFFS (Yet An-other Flash File System) which is the first file system designed specifically for NAND devices, considering the number of dead pages increases and the system reclaims them, in order to perform further write operations, by running a garbage collection process. However, the erase/write unit mismatch generates additional copying of remaining live pages from a block, when erasing it, to another one. Another limitation of the NAND Flash technology is that the number of erase operations is limited to about $10^5$ [10] for SLC (Single Level Cell) and to $10^4$ for $MLC_2$ (Multiple Level Cell) [11]. As any recycling of dead pages introduces block erasing, an even erase-count distribution over the Flash memory blocks cannot be achieved, which results in the "wear-levelling" problem. This has a significant negative impact on the longevity of the memory chips. Much like the 'small write problem' in traditional RAID 5 systems [12].

Many studies from the literature were dedicated to Flash memory, especially to associated file sys-tems and more recently but less significant to pro-vide formal Flash models, as in [13, 14]. In fact, several file systems have been developed to manage data on Flash memories. We can separate them into two classes: Native Flash file systems and non-natine file systems which can be used with any operating systems. The former class is used for raw

Flashmemories, directly integrated in embedded systemsas JFFS (Journal Flash File System) which is a log-structured file system for the NOR Flash device [15].Its second version (JFFS2) [16] supports NAND deviceswithasequentialI/Ointerfaceandamoreefdataintegrityasapriority[18].IttakesintoaccounttheFlashconstraintsandexploititsfeaturestomaximisetheperformanceandtherobustness.Itssecondversion(YAFFS2)accommodatesanewerchipwith larger pages. More recently, LogFS [19] supportssnapshotsandismorespecifictolargedevicesduetoits reducedmountingtimeanditsefficientgarbagecollectionprocess[20].Finally,wecite[21,22]forhybridarchitectureshandling both FlashandRAM.AlloftheseFlashfilesystems have an FTL(FlashTranslationLayer)composedessentiallyoftwoparts :anallocatorprocessforthelogicaltophysicalspacemappingan dacleanerprocessforthegarbagecol-

lection.Themappingbetweenthelogicallocationandthephysic aloneisperformedusingmetadatainthepages'spare areas, mountedat

theinitialisationphasebeforeanyI/Ooperationtakesplace.Gar bagecollectionisperformedinthebackgroundtomake freespaceforwriteoperations.

# Wearlevellinganalysis

The reliability aspect is capital in storage systems. InNAND Flash based systems, either with SLC or MLCtechnology,thereare4typesofreliabilityproblems: Wearoutblocks,

Informationretentionloss,

Writedistrurbphenomenon,

Readdistrurbphenomenon.

The first class of problems is the most importantone as the chip blocks cannot be used anymore andthecontaineddataislostfortheuser.Toavoidthis

situation, a good wear levelling should be guaranteedto exploit the longevity of the chip at its maximum.There are many algorithms to implement the wearlevelling.Mainly,theycanbesplitintotwocategories:*static*and*dynamic*algorithms.Theformerissimple touse,itscostisnegligibleandprovidesan average wear leveling quality for all IO profiles asitiscompletelyindependentfromtheapplicationsIOrequests accessing thestored data.Thesecondclass is a bit more complex to implement as it builtsstatistics first and keeps maintaing them over time,thenadaptsthewriterequestsdistributionusingthese statistics according to the blocks use indicator.Thisiscostlyinbothcomputingtimeandstoragespace butprovidesawearlevelling"alacarte"whichisconsideredopti mal.

In this work, we propose a very simple static map-ping functions and assess their impact on the deviceuse at different levels as well as on its real longevityduration.

In this section, we present the tools used to im-plementourmappingfunctionsandstudythere-sulted wear-levelling algorithm in subsection 3.1, wedescribe our proposition in subsection 3.2 and finallywevalidateitinsubsection3.3.

## Architecture'sconfigurationandtools

In this study, we consider a specific architecture butthiscaneasily beextendedto alternate Flashconfigurations by using their description files, availablein hardware libraries. In the present case, the targetarchitecture is a package composed of 16 chips of theK9KAG08U0MNAND-Flashof2GB[23],connectedby a single 40MBps channel. This Flash storage pack-age is seen as a single address space, a logical storagepool.

Wewrotea customised event driven simulatorin Cto represent the Flash storage system and its associ-ated modes of operation. In this section, we use onlythe mapping module where the proposed static wearlevelling algorithm is implemented.For all the per-formed simulations, we considered traces of 2M reqseach.

## Proposition

Wegivebelowthestaticmappingfunctionsweproposetoimple mentthewearleveling:

$chipID = ladr \% nbchip$

$blocID = (chipID/nbchip)\%chipsize$ $pageID = ladr/(chipsize * nbchip)$

*ladr*is the logical address, *nb chip* is the number ofFlash chips in the package and *chip size* is the num-ber of pages per Flash chip. The % denotes the mod-ulooperationand(*chipID, blocID, pageID*)denotethe physical address of any page (address unit) in thepackage. We observed the wear levelling achieved atthree different levels: the chip, the block and the pagelevels and we considered three parameters: the num-berofaccessestoeveryFlashchipamongthe16in

total, the number of erases of every block among the$16 \times 8$ ones and the number of writes for every pageamongthe16 $\times 8 \times 64$ones.

## Staticwearlevelingqualityvalida-tion

The validation of the efficiency of our mapping func-tions in providing a good wear levelling is performedtrough3phases: Validation of a uniform utilisation of Flash chipswithinthepackage.

Figure 2 shows the mean chip utilisation ($C_{use}$)and the package utilisation ($P_{use}$) as percentagesof time, against the arrival rate.Both are inde-pedentoftheexecution mode.Thisisrelatedto the constant service time, whether the serviceis delayed or not, interrupted or not.They arelinear, increasing with the arrival rate.The $C_{use}$represents the mean chip utilisation but can beconsidered as the chip utilisation because all thechips are almost equally used, as shown in fig-ure3,duetothegoodwearlevelling.

Validation of a good wear levelling at three hier-archicallevels:chip,blockandpage.

Figures4,5and 6represent showthegoodwear

focusing on the emulation of the hardware functionsand the second generates suitable IO traces for theconductedtests.

Weconsiderthearchitectureconfigurationde-scribed in 3.1, where data is manipulated using threecommands – read, write and erase – one at a time.The service times of these commands were estimatedas constants by measurement on

real chips as $130\mu s, 305\mu s, 1.5ms$ respectively, including the bus trans-fer time. Therefore, our performance study focuses on the waiting time, unique parameter affecting the response time.

We use our event driven simulator, representing the Flash storage system and its associated modes of operation. Its execution module processes events as they occur – for example the arrival of a request of a given type or the completion of an access – by main-taining a standard event diary. In addition to the FIFO execution mode, both priority and preemption policies can be accommodated. The simulator en-sures correct implementation of the relative priorities between two classes: a high priority class composed of reads only and a low priority class composed of writes and erases, which have equal priority and are processed in order of arrival. In the case of the pre-emptive priority mode, a write or erase is interrupted as soon as a read enters the system. The interrupted operation being resumed as soon as there are again no reads outstanding, but subject to further interrup-tion. In non-preemptive mode, a write or erase, once started, is allowed to finish, any new read arrivals being queued.

For the investigation of the behaviour of a Flash storage package serving various IO profiles, devel-oped IO generator handles different probability dis-tributions for both logical adresses and IO interar-rival times [24].

## Experimentations

Simulations were run using different synthetic IO workloads generated using our generator. The re-quests' type is consistent with standardised OLTP criteria, e.g. a fixed read:write ratio of 3:1 with ar-rival rates ranging from 50req/s to 5000req/s, and each address trace had a total of 500,000 requests. Various interarrival times were considered using dif- ferent distributions: Poisson for the 'typical user' case2, Erlang for the multi-source case and Pareto for the heavy tailed case, in an attempt to be rep- resentative of read environments. The Erlang traces are chosen to see the effect of reduced variance in the interarrival times. In fact, every exponential stage in the Erlang-n random variable is n times as fast as the Poisson process. This gives a variance of $1/(n\lambda)$, lower than the variance of the interarrival time in the corresponding Poisson process with rate $\lambda$, by a fac- tor of n. Conversely, the Pareto traces were chosen to examine the effect of increased variance, again keep- ing the mean interarrival time the same at $1/\lambda$. We used a Pareto random variable with range [1, [ and probability distribution function F (x) = 1 x$-\alpha$. Direct integration shows that this has nth moment $\alpha/(\alpha \quad n)$, which exists if and only if $\alpha$ > n; thus, $1/2 < \lambda < 1$ for the variance to exist. To achieve this, we first scale a range of actual arrival rates, so as to satisfy the inequality, then generate the Pareto- based input traces, and then scale these

back again by multiplying by the same scaling factor.

## Resultsanddiscussion

The focus is placed on the queueing time for each of the three operation classes (read, write and erase) because its variability has a significant effect on overall Flash performance due to its fixed service time.

Queueing time is considered the main performance metric for every type of operation because of the con-stant service (or access) time for all three access types (read/write/erase). Thus queueing time is a pur e performance metric, being entirely an overhead that depends only on the execution mode. We investi-gated the three main modes: no priority among the operations (mode 1); non-preemptive priority to read operations (mode 2); and finally preemptive read pri-ority (mode 3).

## Conclusion

In this paper, we have proposed simple static map- ping functions that ensure good wear levelling for uni- formly distributed accesses to the storage space, as well as for accesses with hot spots, even in quite ex- treme case with 1% of the storage space accessed 100 times more frequently than the other 99%. Such a static scheme avoids the complex implementation and the frequent statistics extraction and management routines called under dynamic mappings. Consider- ing the queueing time, we confirmed using the Pois- son distribution which provides a standard against which to assess other workload types, that (of course) it increases as the arrival rate increases, more rapidly as the system approaches instability. We observed similar qualitative behaviour for the Erlang case but the queueing times are smaller, due to the lower vari- ance in the interarrival times, while the queueing times for the Pareto distribution, where the variance is larger, increase at certain arrival rates. We showed that the chips within the package are equally under- used even when the arrival rate is high, and similarly for the corresponding queueing time. This suggests that, rather than using the package as a sole unit, it is better to exploit the parallelism available among the Flash chips for improved usage of the hardware com- ponents and a reduced queueing time. This should be achieved in the short term, taking into account the concurrent access management to chips, the re- quests' scheduling policies and the shared bus alloca- tion. Further, it should be extended to the Flash bi- dimensional vector configuration. In the longer term, we plan to extend our fluid model for the Flash pack- age storage system [14] to handle Flash chips operat- ing in parallel.

# References

J. Gray, "Tape is dead,disk is tape, flash is disk, ram lo- cality is king," 2007, pres. at the CIDR Gong Show, Asilomar,CA,USA.gineering Complex Computer Systems, 2007.

P. G. Harrison, N. M. Patel, and S. Zertal, "Response time distribution of flash memory accesses," in Valuetools, 2008.

D. Woodhouse, "Jffs : The journalling flash file system," in Ottawa Linux Symposium, 2001.

——, "Journaling flash file system (jffs) and journaling flash filesystem 2 (jffs2)," website, http://sources.redhat.com/jffs2/jffs2.

A. Schierl, G. Schellhorn, D. Haneberg, and W. Reif, Abstract Specification of the UBIFS File System for Flash Memory. Springer Berlin/Heidelb erg, 2009, vol. 5850/2009, pp. 190– 206.

A. O. Limited, "Yaffs2 readme-linux v1.1," website, 2007, http://aleph1.co.uk.

LogFS, "Logfs home page," website, http://logfs.org/logfs/LogFS.

J. Engel, D. Botle, and R. Mertens, "Garbage collection in logfs," in Linyx conf., 2007.
Y. Park, S.-H. Lim, C. Lee, and K. H. Park, "Pffs: A scalable flash file system for the hybrid architecture of phase-change ram and nand flash," in 2008 ACM Symposium onApplied Computing, 2008, pp. 1498–14 503.

Y. Park, S. Lim, C. Lee, and K. Park, "Pffs: a scalable flash memory file system for the hybrid architecture of phase-change ram and nand flash," in ACM Symposium on Applied Computing, 2008.

Samsung, "K9kag08u0m nand-flash," website, http://www.samsung.com/global/system/business/ semi-conductor/product/.

R. Jain, The Art of Computer Systems Performance Analysis:Techniques for Experimental Design, Measurement, Simulation, and Modeling. Wiley Interscience, 1991.