

Response Time Distribution of Flash Memory Accesses

Peter G. Harrison¹ Naresh M. Patel² and Soraya Zertal³

¹ Imperial College London, South Kensington Campus, London SW7 2AZ, UK
`pgh@doc.ic.ac.uk`

² NetApp, Inc., 495 East Java Drive, Sunnyvale, CA 94089, USA
`naresh@netapp.com`

³ PRiSM, Université de Versailles, 45, Av. des Etats-Unis, 78000 Versailles, France
`Zertal@prism.uvsq.fr`

Abstract

Flash memory is becoming an increasingly important storage component among non-volatile storage devices. Its cost is decreasing dramatically, which makes it a serious competitor of disks and a candidate for enterprise-tier storage device of the future. Consequently, there is an urgent need for models and tools to analyse its behaviour and evaluate its effects on a system's performance. We propose a fluid model with priority classes to investigate the response time characteristics of Flash memory accesses. This model can represent well the Flash access operations, respecting the erase/write/read relative priorities.

keywords: Fluid model, Flash memory, Response time distribution, Preemptive mode.

1 Introduction

Storage devices based on Flash memory are becoming more and more prevalent in our daily life. This recent technology presents a panoply of devices, subject to continuous and intensive evolution – as commonly observed in the semiconductor manufacturing industry but especially prominent here in response to the market demand of mp3 players, mobile phones, digital cameras, solid state drives and other consumer and enterprise storage products.

However, there is currently a lack of tools capable of analysing the technology's quantitative behaviour and its effect on the whole system of which it forms a part. Moreover, there are no methodologies available to evaluate the performance it can deliver and the associated cost of its very specific operating characteristics. Such tools and methodologies must be representative of the system under scrutiny, accurate at the scale of the semiconductor and quick to develop and use, in order to remain in step with the wide variety of devices and their rate of manufacture. The mathematical modelling techniques we present address these requirements with respect to both their ease of representation and accuracy. This requires careful model design, parameterisation and validation.

Up to our knowledge, a considerable body of work is done on the usage/access [3, 4, 13] and reliability [20, 22, 23] aspects of Flash with a number of specific file systems [17, 19, 16, 14, 12, 15, 21]. On the other hand, there is no analytical model of Flash memory. The only one, which is still in progress is done by Butterfield et al. [1, 2]. In this existing work, fragments of Z models of NAND Flash are presented and a sketch of main ideas are published, nevertheless no complete formalisation or performance evaluation are included. This previous work is based on the “open NAND Flash Interface (ONFI)” consortium [6] and intend to give a formal functional description of NAND Flash memory without any evaluation of its performance or analysis of its behaviour when it is submitted to different degrees of utilisation.

In our work, we propose a fluid model for a NAND Flash memory. Fluid models represent the behaviour of Flash memory well because they can take into account correlation between access streams – of erase, write, and read operations – as well as their relative priorities, and can handle different execution modes: preemptive or non-preemptive. Preemption of erase and write operations is possible but uses up erase cycles and hence causes early wear-out. We consider, as in [5], the three different sources of requests and a Flash component with an intermediate (fluid) buffer. Then, a Markov chain is used to describe the state of the access operations’ sources, in particular their input/output rates. Although the tools supporting fluid models are not yet as well developed as those for discrete state Markov models [8], there is a wide range of applications that can be modelled effectively using them; Flash memory access is one such.

In the rest of the paper, the next section describes the background to Flash technology – its physical characteristics, its access operations and its impact in the enterprise storage system domain. Section 3 is dedicated to the fluid model and how it is applied to represent our particular Flash memory system. We analyse both the low and high priority storage access classes, as well as the corresponding no-priority system, in terms of their respective response times, using busy period analysis. As a case study, we consider an online transaction processing (OLTP) system in section 4, which clarifies how the proposed, generic fluid model can be customised to represent a Flash storage system supporting this specific type of workload. Numerical results are given in section 5. Finally, we conclude in section 6, suggesting desirable model enhancements, further ideas for analysis and designs that may emerge in response to predictions arising from the modelling methodology proposed for this recent technology.

2 Background on Flash memory

The impact of Flash memory as a major, non-volatile storage component stems from its shock resistance, vibration tolerance, light weight and low energy consumption, as well as its high capacity. It already has a wide range of applications, from its use in personal computers as a solid-state drive (SSD) [13] to critical environments such as satellite systems [23].

There are two types of Flash memory, named according to their manufacturing: NOR and NAND. The former has lower density and higher cost but provides fast random access and can be easily re-programmed; this makes it most suitable for storing codes. Another advantage 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.

On the other hand, NAND Flash has a very large storage capacity and provides high performance for large read/write requests, making it more suitable for storing data [14]. In fact the density is about 8 times more for NAND [17], at a cost that

is 4 to 8 times cheaper than NOR. Further, sequential read and write accesses are faster on NAND Flash. Although erases are significantly faster on NOR, they can be pre-scheduled in NAND, essentially running in a garbage collector.

All of the previous characteristics make NAND the technology of choice for many data centres and this is what we model here, using fluid methods. Further, MLC_n (Multi Level Cell) technology multiplies the storage capacity of the Flash memory chip by having n -bit information per cell, where $n = 2$ or 4 in contemporary devices. We model essentially the basic SLC (Single Level Cell) technology with only one bit per cell, i.e. effectively MLC_1 , but it is straightforward to adapt the techniques used to arbitrary n in the mapping from discrete storage requests onto volumes of fluid – see section 4.

2.1 NAND Flash memory access

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 (figure 1). A block is the erase operation’s unit of storage, whilst a page is the

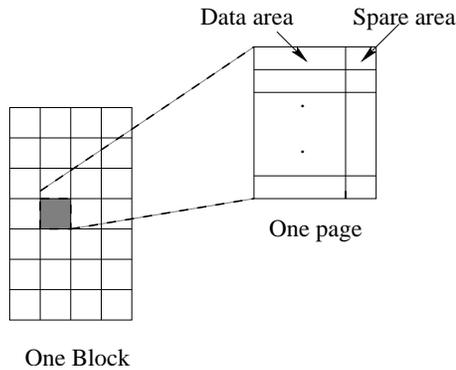


Figure 1: A NAND Flash block structure

read and write operation’s unit. No ‘in-place’ updates are allowed in NAND Flash and pages have to be written in order within a block. So, 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, the number of dead pages increases and the system reclaims them, in order to perform further write operations, by running a garbage collection process. This is possible only when these pages are erased. However, this erase/write unit mismatch generates additional copying of 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 [4] for SLC (Single Level Cell) and to 10^4 for MLC_2 (Multiple Level Cell) [22]. 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 chip.

Much like the small write problem in traditional RAID5 systems, random writes to a Flash device require a read-modify-write sequence to overcome the lack of in-place updates: read the old block, modify updated pages, and write the new block. Some form of mapping function is needed to hide this from the user and the next subsection outlines popular approaches.

2.2 Flash management systems

Several file systems have been developed to manage data on Flash memories. JFFS (Journal Flash File System) is a log-structured file system for the NOR Flash device [19]. Its second version (JFFS2) supports NAND devices with a sequential I/O interface and a more efficient garbage collection process, using a relatively reduced mounting time. YAFFS (Yet Another Flash File System) is the first file system designed specifically for NAND devices, considering data integrity as a priority [15], and its second version (YAFFS 2) accommodates a newer chip with larger pages. More recently, LogFS [16] supports snapshots and is more specific to large devices due to its reduced mounting time and its efficient garbage collection process [12]. Finally, we cite [21] for hybrid architectures handling both Flash and RAM.

All of these Flash file systems have an FTL (Flash Translation Layer) composed essentially of two parts: an allocator process for the logical to physical space mapping and a cleaner process for the garbage collection. The mapping between the logical location and the physical location is performed using metadata in the pages' spare areas, mounted at the initialisation phase before any I/O operation takes place. Garbage collection is performed in the background to make free space for write operations. The block size is an important design parameter which has a significant impact on the amount of metadata and so mounting time. The flash controller therefore has many design choices, which affect the output stream of storage requests it delivers to the Flash devices. Hence we consider a fairly general input stream, using a Markov modulated arrival process to model the input to a Flash memory array. More specific, high-level models, based on specific flash controller designs, could then be developed by appropriately parameterising the base model we describe here.

2.3 Flash for enterprise storage

As more consumer devices incorporate Flash and new process technology is implemented, the cost of Flash is expected to decrease rapidly in the coming years. This trend may allow Flash to compete with high-speed enterprise disks in the future. Most enterprise storage systems already provide fast writes using Non-Volatile RAM but random read operations from disk incur long latencies (typically 5 to 10 milliseconds). On the other hand, sequential read operations from disk are very fast and so the randomness of the access pattern significantly impacts the observed performance. The key advantage of Flash is that the random reads and sequential reads have the same low latency. On the downside, Flash implementations have to deal with the complexities of the write constraints outlined above.

While 4KB page read operations from Flash have fairly low average latency (say 25 microseconds for access and 105 microseconds for channel transfer), the distribution of the latency may have a long tail because of contention with write and erase operations. Page write operations take about 105 microseconds on the channel and 200 to 700 microseconds for access (often an increasing function of erase count). Erasing a block (say 256 KB, or 64 pages) can take more than 1.5 milliseconds. These numbers vary considerably among the Flash device manufacturers, but often the latency concern stems from the possibility that small read operations have to wait behind relatively long write and erase operations. Hence we would like to model the distribution of latency for read operations (in particular) under various workloads. Fluid modelling techniques provide a promising way of improving our intuition about response time distributions for given Flash device characteristics, as well as providing quantitative estimates of quantiles.

3 Fluid model

We use a fluid queue, with input defined by a four-state (or four-phase) continuous time semi-Markov chain (CTSMC), to model the behaviour of a Flash device. Such a queue is appropriate to model correlated streams of incoming data requests of different types, especially (but not exclusively) at moderate to high utilisation. In fact, in high utilisation scenarios, the erase operations and their related page copying can be generated frequently due to the high rate of native write/update operations. The four states, numbered 0 to 3, correspond to *off* (no input), *read requests*, *write requests* and *erases*. Class 1 fluid ‘particles’ (i.e. those arriving in state 1) have priority over states 2 and 3 to account for the priority actually given to reads, which are usually more critical in allowing a process to continue.

The queue is parameterised, for an n -state CTSMC so as to allow for model extensions with more than four states, as follows:

- the CTSMC has probability transition matrix $P = (p_{ij} \mid 0 \leq i, j \leq n - 1)$, defined at state transition instants;
- the state holding time in state i , given that the next state is j , has probability distribution $H_{ij}(t)$, for $0 \leq i, j \leq n - 1$.
- the fluid arrival rate in state i is the constant λ_i volume-units of fluid per unit time, for $0 \leq i \leq n - 1$, so $\lambda_0 = 0$;
- the rate at which the server outputs fluid when its buffer is non-empty is the constant μ volume-units of fluid per unit time;
- the net input *rate matrix* $R = \text{diag}(r_0, \dots, r_{n-1})$, where $r_i = \lambda_i - \mu$ for $0 \leq i \leq n - 1$, and the rate vector $\vec{r} = (r_0, \dots, r_{n-1})$;
- the vector $e = (1, 1, \dots, 1)$.

In fact we will require that all states other than the high-priority state 1 have exponential holding times with parameter that does not depend on the next state transitioned to; these are characterised by a *rate* γ_i for state $i = 0, 2, 3$, so that $H_{ij}(t) = 1 - e^{-\gamma_i t}$ for all states $j \neq i$. In this paper, we further assume that state 1 has exponential holding time, so we have a modulating continuous time Markov chain (CTMC), specified by a generator matrix, Q say¹.

The response time distribution (or, rather, its Laplace-Stieltjes transform (LST)) of each class in this fluid queue can be calculated using the following observations:

1. high priority fluid (arriving in state 1) is processed as if the other classes had fluid arrival rates 0, which effectively constitute additional off-states.
2. response time for a state-1 particle is therefore the time taken to serve the state-1 fluid in the buffer present at the arrival instant of that particle – that is, simply the volume of state-1 fluid present divided by μ , given by the reduced model described in the previous point.
3. response time for particles arriving in lower priority states is the sum of the time to serve all fluid (of both classes) already present on arrival – of volume F say – added to all additional fluid of class 1 that arrives during that time period. The amount of this additional fluid is the product of λ_1 and the sum of the *busy periods* generated by each of these state-1 arrivals. When

¹This assumption is not necessary, however, since it is only used to determine the equilibrium fluid level distribution, which can be derived fairly simply as a functional equation for an on-off process in which the off-period is exponential.

there is only one off-state (i.e. in an on-off process), the number of these busy periods, when conditioned on the time period F/μ , is a Poisson random variable. However, this is not the case here and a more complex analysis is given in section 3.4. Note that the fluid queue model for the period considered is that described in the first point since lower class fluid waits behind class 2 and so can be ignored here. However, F is given by the steady state fluid level in the original 4-state fluid queue *with no priority classes*, because of the service rate being the same for all classes. The probability distribution of this fluid level is well known for any number of states with any fluid input rates.

We therefore have to calculate:

- the steady state probability distribution of the fluid level F in the no-priority queue – a routine task;
- the busy period for a class 1 fluid arrival in the reduced model – this is determined in section 3.2 using results from [9];
- the Laplace transform of the response time distributions – this is done in section 3.4 in terms of conditional expectations on the number of high priority arrivals during the processing of the F units of fluid.

3.1 Non-priority fluid queue results

Here we restate relevant results for a single fluid queue with a Markov modulated arrival process of any number of states. Consider a single, n -state Markovian fluid queue, with generator matrix $Q = (q_{ij} \mid 1 \leq i, j \leq n)$ for the input process, which has equilibrium probabilities $\vec{\pi}$ (so that $\pi Q = \vec{0}$ and $\pi \vec{e}^T = 1$).

Using the notation defined in the previous section, it can be shown that the vector density function of the fluid level at equilibrium (when this exists), denoted by $\vec{f}(x) = \frac{\partial \vec{F}}{\partial x}$ where $\vec{F}(x)$ is the corresponding probability distribution function, has Laplace transform:

$$\vec{f}^*(\theta) = \vec{F}(0)R(R - Q/\theta)^{-1} \quad (1)$$

The constant $\vec{F}(0)$ is given by the boundary conditions that the fluid level in states with fluid arrival rate greater than the service rate is positive with probability 1, and taking the limit $\theta \rightarrow 0$. Details can be found in [7], but note that this is not the only solution method; see for example [5].

For the high priority class's response time, we can consider the queue as if the other classes had zero arrival rate. To get the equilibrium probability density function of the volume of class 1 fluid in the queue, we simply solve the same problem with $\lambda_2 = \lambda_3 = 0$. Let this density function have Laplace transform denoted by $\vec{f}_1^*(\theta)$, that for the non-priority fluid level F being $\vec{f}^*(\theta)$, as above.

3.2 High priority class busy times

The model we solve here again has four states, but three of these have zero fluid arrival rate at a given time t . We therefore make a further simplification by aggregating the write and erase states, 2 and 3, into a new state 2, which we assume to have exponential holding time. This is not as unreasonable as it may appear – in fact probably no less unreasonable than assuming the read, write and erase states were exponential in the first place. Actually, the problem can be solved without this aggregation, but at greatly increased computational cost². The mean holding

²As we will see below, we would need to solve a cubic for a parameter β , rather than a quadratic, and then three simultaneous quadratic equations, rather than two.

time of the new state 2, equal to the reciprocal of the parameter of its exponential distribution, is determined as the mean time elapsed between entering either state 2 or 3 and next entering either state 0 or 1. This is a routine, first passage time calculation, which yields a mean holding time of $h_{23} = \frac{\pi_2 h_2 + \pi_3 h_3}{\pi_2 + \pi_3}$ where

$$h_2 = \frac{q_{23} - q_{33}}{q_{22}q_{33} - q_{23}q_{32}} \quad \text{and} \quad h_3 = \frac{q_{32} - q_{22}}{q_{22}q_{33} - q_{23}q_{32}}$$

and the q_{ij} are the generators of the 4-state process.

Consistent with [9], the set of draining-states is denoted by $\mathcal{E} = \{0, 2\}$ and the set of all states by $\mathcal{S} = \{0, 1, 2\}$. The number of draining-states is $n_e = |\mathcal{E}|$ and the total number of states is $n = |\mathcal{S}|$; thus in our example, $n_e = 2$ and $n = 3$.

We now define the following matrices:

- the $n \times n$ matrix $\Gamma = (\gamma_{ij})$, the generator matrix of the 3-state Markov process, where we write $\gamma_i = -\gamma_{ii}$ for the positive total transition rate out of state $i \in \mathcal{S}$.
- the $1 \times n_e$ matrix (vector) $\mathbf{V} = (V_{2j}^*(\alpha))$, the j element of which is the required Laplace transform, with real parameter $\alpha \geq 0$, of the probability density function of the busy period – i.e. of the time elapsed between entering into the filling-state 1 when the queue is empty (hence, coming from a draining-state) and the queue next becoming empty again in the draining-state j .
- the $n_e \times n$ matrix $\mathbf{U} = (u_{ij})$ by $u_{ij} = -\gamma_{ij}(1 - \delta_{ij})/r_i$ for $i \in \mathcal{E}$ and $j \in \mathcal{S}$.
- the $n_e \times n_e$ diagonal matrix $\mathbf{D} = (d_{ij})$ by $d_{ij} = [\beta - (\alpha + \gamma_i)/r_i]\delta_{ij}$ for $i, j \in \mathcal{E}$, where $\beta \in \mathcal{D}$.
- the $n \times n_e$ matrix $\mathbf{W} = (w_{ij})$ by $w_{1j} = V_{1j}^*(\alpha)$, $w_{ii} = 1$ for $i \in \mathcal{E}$ and $w_{ij} = 0$ otherwise, for $j \neq i \in \mathcal{E}$.

The $n_e \times n_e$ matrix $\mathbf{M} = \mathbf{D} - \mathbf{UW}$ has determinant, $\Delta \stackrel{\text{def}}{=} |\mathbf{M}|$ say, which is a polynomial of degree n_e in β and whose cofactors are polynomials of degree $n_e - 1$. Therefore each element in the inverse matrix $(\mathbf{D} - \mathbf{UW})^{-1}$ (assuming this exists, i.e. that the determinant $|\mathbf{M}|$ is non-zero) is a rational function of β which we may write in partial fractions as

$$\mathbf{L} = (\mathbf{D} - \mathbf{UW})^{-1} = \left(\sum_{s=1}^{n_e} \frac{a_{s;ij}(\alpha)}{\beta - b_s(\alpha)} \mid i, j \in \mathcal{E} \right)$$

assuming no degeneracies, so that $\{b_s(\alpha) \mid 1 \leq s \leq n_e\}$ are the distinct roots of the equation $\Delta = 0$ (in β) and the $a_{s;ij}(\alpha)$ are independent of β , defined by (dropping the parameter α when the meaning is clear):

$$a_{s;ij} = \left[\frac{c_{ji}}{\Delta_s} \right]_{\beta=b_s} \quad (2)$$

where $\mathbf{C} = (c_{ij} \mid i, j \in \mathcal{E})$ is the cofactor matrix of \mathbf{M} and $\Delta_s = \Delta/(\beta - b_s) \equiv \prod_{1 \leq i \neq s \leq n_e} (\beta - b_i)$. Our special case of only a single filling-state (or on-state) now gives a simplified result for the Laplace transform of the high priority class busy time probability distribution, $V^*(\theta)$ [9]:

Theorem 1 For $j \in \mathcal{E} = \{0, 2\}$,

$$V_{1j}^*(\alpha) = \sum_{k \in \mathcal{E}} \sum_{s \in \mathcal{E}} a_{s;kj}(\alpha) p_{1k} H_{1k}^*(\alpha - r_1 b_s(\alpha))$$

for general holding times in the high priority state 1. In the Markovian case (when the holding time in state 1 is exponential with parameter γ_1)

$$V_{1j}^*(\alpha) = \sum_{k \in \mathcal{E}} \sum_{s \in \mathcal{E}} \frac{a_{s;kj}(\alpha) q_{1k}}{\gamma_1 + \alpha - r_1 b_s(\alpha)}$$

Notice that the terms $a_{s;kj}(\alpha)$ are functions of the busy period Laplace transforms $V_{1i}^*(\alpha)$ ($i, j, k \in \mathcal{E}$), giving a set of mutual, non-linear equations. Of course, if $n_e = 1$, there is only one such equation, for $V_{10}^*(\alpha)$, which is the well known one for the busy period in a fluid queue with on-off source.

We will also require the related quantity $Z_{ij}^*(x, \alpha)$ for each $i, j \in \mathcal{E}, x > 0$, which is the LST of the probability distribution of the time taken, $Z_{ij}(x)$, for a buffer containing x units of fluid in draining state i to first become empty in state j . The following result gives this:

Theorem 2 For $i, j \in \mathcal{E}$ and real $x > 0$,

$$Z_{ij}^*(x, \alpha) = \sum_{s=1}^r a_{s;ij}(\alpha) e^{b_s(\alpha)x}$$

Theorems 1 and 2 have been implemented in Mathematica[®] and our numerical results are obtained using them, together with the analysis of response times given in the next section. To illustrate the theorems' use, we first need the matrix \mathbf{M} to calculate the terms b_s and $a_{s;kj}$, for $j, k, s = 0, 2$.

The matrix $\mathbf{D} = \text{diag}(\beta + (\alpha + \gamma_0)/\mu, \beta + (\alpha + \gamma_2)/\mu)$, $u_{ij} = \gamma_{ij}/\mu$ ($i = 0, 2, j = 0, 1, 2$) and

$$\mathbf{W} = \begin{pmatrix} 1 & 0 \\ V_{10}(\alpha) & V_{12}(\alpha) \\ 0 & 1 \end{pmatrix}$$

Hence

$$\begin{aligned} \mathbf{M} &= \mathbf{D} - \mathbf{UW} \\ &= \begin{pmatrix} \beta + \frac{\alpha + \gamma_0 - \gamma_{01} V_{10}^*(\alpha)}{\mu} & -\frac{\gamma_{02} + \gamma_{01} V_{12}^*(\alpha)}{\mu} \\ -\frac{\gamma_{20} + \gamma_{21} V_{10}^*(\alpha)}{\mu} & \beta + \frac{\alpha + \gamma_2 - \gamma_{21} V_{12}^*(\alpha)}{\mu} \end{pmatrix} \end{aligned}$$

$b_1(\alpha), b_2(\alpha)$ are the roots (assumed distinct) of the equation

$$\begin{aligned} &\beta^2 + \frac{2\alpha + \gamma_0 + \gamma_2 - \gamma_{01} V_{10}^*(\alpha) - \gamma_{21} V_{12}^*(\alpha)}{\mu} \beta + \\ &\frac{(\alpha + \gamma_0 - \gamma_{01} V_{10}^*(\alpha))(\alpha + \gamma_2 - \gamma_{21} V_{12}^*(\alpha))}{\mu^2} - \frac{(\gamma_{02} + \gamma_{01} V_{12}^*(\alpha))(\gamma_{20} + \gamma_{21} V_{10}^*(\alpha))}{\mu^2} = 0 \end{aligned}$$

and $a_{1;ij}, a_{2;ij}$ are the i - j th elements of the matrices

$$\mathbf{A}_1 = \frac{1}{b_1(\alpha) - b_2(\alpha)} \begin{pmatrix} b_1(\alpha) + \frac{\alpha + \gamma_2 - \gamma_{21} V_{12}^*(\alpha)}{\mu} & \frac{\gamma_{02} + \gamma_{01} V_{12}^*(\alpha)}{\mu} \\ \frac{\gamma_{20} + \gamma_{21} V_{10}^*(\alpha)}{\mu} & b_1(\alpha) + \frac{\alpha + \gamma_0 - \gamma_{01} V_{10}^*(\alpha)}{\mu} \end{pmatrix}$$

and

$$\mathbf{A}_2 = -\frac{1}{b_1(\alpha) - b_2(\alpha)} \begin{pmatrix} b_2(\alpha) + \frac{\alpha + \gamma_2 - \gamma_{21} V_{12}^*(\alpha)}{\mu} & \frac{\gamma_{02} + \gamma_{01} V_{12}^*(\alpha)}{\mu} \\ \frac{\gamma_{20} + \gamma_{21} V_{10}^*(\alpha)}{\mu} & b_2(\alpha) + \frac{\alpha + \gamma_0 - \gamma_{01} V_{10}^*(\alpha)}{\mu} \end{pmatrix}$$

Defining the vector $\mathbf{V}(\alpha) = (V_{10}^*(\alpha), V_{12}^*(\alpha))$, we finally have

$$\mathbf{V}(\alpha) = \mathbf{g}_1(\alpha) \cdot \mathbf{A}_1(\alpha) + \mathbf{g}_2(\alpha) \cdot \mathbf{A}_2(\alpha)$$

where the vectors $\mathbf{g}_s = (\gamma_{10}, \gamma_{12})/(\gamma_1 + \alpha - b_s(\alpha)r_1)$ for $s = 1, 2$.

3.3 Response time of high priority class

Let the queueing time of a high priority fluid particle – representing the start of a read – be Q_1 . Then the LST of the distribution function of Q_1 is

$$Q_1^*(\theta) = \mathbb{E}[e^{-\theta Q_1}] = \mathbb{E}[e^{-\theta L_1/\mu}] = L_1^*(\theta/\mu) \quad (3)$$

where L_1 is the volume of high priority fluid in the queue on arrival of the particle. This is simply the total volume of fluid present at equilibrium in the same queue when the arrival rates of fluid in all states other than 1, that of the high priority class, are zero. $L_1^*(\theta)$ is therefore given by equation 1 for the four state model with fluid rate vector $\vec{r}_0 = (-\mu, \lambda_1 - \mu, -\mu, -\mu)$, as $L_1^*(\theta) = f_1^*(\theta; \vec{r}_0)/f_1^*(0; \vec{r}_0)$.

3.4 Response time of low priority class

Let the queueing time of a low priority, state-2 or state-3, fluid particle – representing the start of a write or erase operation – be Q_2 . This is the time required to serve *all* the fluid present on arrival, of volume L_2 units, say, as well as any high-priority fluid that arrives during this period – i.e. $Q_2 = Z_{2j}(L_2)$ when the queueing time ends in state $j \in \mathcal{E}$. Then the LST of the distribution function of Q_2 is

$$\begin{aligned} Q_2^*(\theta) &= \sum_{j \in \mathcal{E}} \mathbb{E}[e^{-\theta Z_{2j}(L_2)}] \\ &= \sum_{j \in \mathcal{E}} \mathbb{E}\left[\sum_{s=1}^{n_e} a_{s;2j} e^{b_s L_2}\right] \quad \text{by Theorem 2} \\ &= \sum_{j \in \mathcal{E}} \sum_{s=1}^{n_e} a_{s;2j} L_2^*(-b_s) \\ &= \sum_{j \in \mathcal{E}} \sum_{s=1}^{n_e} a_{s;2j} \frac{\lambda_2 f_2^*(-b_s)/f_2^*(0) + \lambda_3 f_3^*(-b_s)/f_3^*(0)}{\lambda_2 + \lambda_3} \end{aligned}$$

where $f_k^*(\theta) \equiv f_k^*(\theta; \vec{r})$, $k = 2, 3$ in equation 1; no fluid arrival rates are set to zero in the low priority case. The weighting proportionate to λ_2 and λ_3 gives the LST of the fluid level's probability distribution at the arrival instant of the considered low priority particle, i.e. for the random variable L_2 . Note that $f_2^*(\theta)$ and $f_3^*(\theta)$ are computed for the original, four-phase model whereas the terms b_s and $a_{s;2j}$ relate to the reduced, three-phase model.

In the special case that $n_e = 1, \mathcal{E} = \{2\}$, we find that

$$a_{1;22} = 1, b_1(\alpha) = -\frac{\alpha + \gamma_2(1 - V_{12}^*(\alpha))}{r_2}$$

so that

$$Q_2^*(\theta) = f_2^*((\theta + \gamma_2(1 - V_{12}^*(\theta)))/\mu)/f_2^*(2)$$

This can be obtained by letting N be the number of high priority (class 1) busy periods that punctuate the time period L_2/μ . Since N is a Poisson random variable,

abbreviating V_{12}^* by V^* , we then have

$$\begin{aligned}
Q_2^*(\theta) &= \mathbb{E}[\mathbb{E}[e^{-\theta(V_1+\dots+V_N+L_2/\mu)} \mid L_2]] \\
&= \mathbb{E}[e^{-\theta L_2/\mu} \mathbb{E}[e^{-\theta V} \mid L_2]] \\
&= \mathbb{E}[e^{-L_2(\theta+\gamma_{01}(1-V^*(\theta)))/\mu}] \\
&= L_2^*((\theta + \gamma_{01}(1 - V^*(\theta)))/\mu)
\end{aligned} \tag{4}$$

The expressions for $Q_1^*(\theta)$ and $Q_2^*(\theta)$ allow a direct comparison to be made between the latencies of the priority classes. It is straightforward to extract the moments of latency, by differentiation at $\theta = 0$, from which average performance and its variability can be compared. However, it is also possible to invert the Laplace transforms numerically, allowing information in the tails to be obtained and hence quantile-based benchmarks to be assessed.

4 Example use case and model parameterisation

For a concrete example, we consider an OLTP workload with read:write data ratio of 3:1. The random reads and random writes are both 8KB, which is the native database block size. However, we assume that the file system (such as Write Anywhere File Layout–WAFL [10]) converts all the random writes into sequential writes with negligible overhead. The Flash controller (and associated FTL) will take care of wear-leveling within Flash devices and WAFL will ensure that all Flash devices in the array are uniformly used. The OLTP application will generate traffic across the entire Flash array, but for our purposes here we consider a single Flash chip. Suppose this 32GB flash package is servicing 1600 I/O’s per second (so that the Flash access density of 50 I/O’s per GB is an order of magnitude better than high-speed disk drives). We also assume that the chip handles at most a single command (read, write, or erase) at a time so that power requirements for the chip are minimal. In practice, Flash controllers can pipeline several commands to different banks in the same package at the same time, though they will contend for the shared channel interface for data transfer. Given the read:write ratio of 3:1 and I/O size of 8KB, the Flash chip will see 2400 page read commands per second and 800 page write commands per second. In order to stay ahead of the write rate, the chip will also have to execute $800/64=12.5$ erase commands per second (assuming 64 pages per erase block).

A concept of volume, viz. *virtual bytes* (“vbytes”), helps transform real-world problems with multiple classes and fixed service rates into fluid models. Essentially, we take one of the classes and compute the bytes transferred in an operation and hence a volume service rate based on the average service time for that class. Suppose we have a large buffer for vbytes that get serviced at a fixed rate of μ vbytes per microsecond. For Flash, a 4KB read takes $105+25$ microseconds, which implies that $\mu = 31.5$ vbytes per microsecond. The classes of interest are now as follows.

- to establish a reference point for the vbyte-measure of volume, when a 4KB read request is issued, we start a flow that puts 4096 vbytes into the buffer. Recalling that the buffer is drained at 31.5 vbytes/usec, this will take 130 μ secs to complete under no contention, as required;
- a 4KB write request takes $105+200\mu\text{sec} = 305\mu\text{secs}$. Although we are writing 4096 bytes of user data, we will need to put in $4(305/130) = 9610$ virtual bytes on average to simulate the actual service time of $305\mu\text{secs}$;
- similarly, an erase takes $1500\mu\text{secs}$ of actual time, so we have to feed the buffer a total of $4(1500/130) = 47,262$ vbytes on average to simulate this service time.

In this way we can map between the number of virtual bytes in the request buffer and the volume in the fluid queue model. The constant service – or *drain* – rate is simply a transfer rate in virtual bytes per microsecond. During periods of no contention, a command starts a flow with rate that is greater than the constant drain rate (31.5 vbytes/ μ sec) and the completion of the active period for that command-mode (read, write or erase) indicates its completion with the specified duration. Notice that the actual rate at which vbytes enter the buffer is arbitrary to some degree: it must exceed the drain rate and higher rates reflect increased burstiness, as discussed below.

For the modulating CTMC, we also have to estimate the mean state holding times, along with the transition probabilities between states. These parameters are necessary to construct the generator matrix and can only be obtained by monitoring workloads typical (in this use case) for Flash devices with OLTP workloads; see [18] for example. In particular, greater burstiness can be represented by decreasing state holding times (increasing their out-transition rates) and correlation in sequences of command-mode types can be accounted for by the state transition probabilities.

5 Numerical results

5.1 Fluid model parameterisation and queueing times

We used Mathematica[®] 5.2 to implement the model described in section 3 and applied it to the use case of the previous section, as a preliminary illustration of its capability. The rates of fluid input in each command mode and drain rate, measured in vbytes per microsecond, were set as described in the previous section. The CTMC’s transition probability matrix was set, somewhat arbitrarily in the absence of a reliable OLTP workload model, to:

$$\begin{pmatrix} -0 & 0.74 & 0.25 & 0.01 \\ 0.74 & -0 & 0.25 & 0.01 \\ 0.25 & 0.74 & -0 & 0.01 \\ 0.1 & 0.148 & 0.752 & -0 \end{pmatrix}$$

The instantaneous transition rates out of each state were set to (2, 5, 2, 10)/1000. This choice reflects reasonable burstiness and plausible correlations in a typical OLTP sequence of database accesses. Moreover, it yields correctly the observed numbers of vbytes issued by each command type in one microsecond, *viz.* 9.83, 7.69, 0.59 for reads, writes and erases, respectively. These figures become 2400, 800, 12.5 when scaled by the vbyte-factors $4096 \times$ the single command relative access times. This is achieved by choosing instantaneous vbyte-arrival rates of 0, 43.0, 27.9, 220.5 for each mode (rate 0 for idle mode) and using the fact that the CTMC with the above generators has equilibrium state probabilities 0.4929, 0.2284, 0.276, 0.0027. The component-wise product of these vectors gives the observed average rates.

In this scenario, the flash unit is running at a moderate utilisation of 57.5% (average total arrival rate, 18.1, divided by the drain rate, 31.5) and so we can expect relatively small fluid levels and access times at equilibrium. In fact, differentiating equation 1, we find that the mean fluid level is 2364.4 (about half a read access) and its standard deviation is 7202.6 – suggesting fairly long tails. For the high priority reads, from equation 3, we obtain the queueing time probability density function shown in figure 2, with mean queueing time 82.19 μ secs. For the low priority writes/erases, we get the density function shown in figure 3, with mean queueing time 184.2 μ secs. This is broadly what we would expect at utilisations of around 50%.

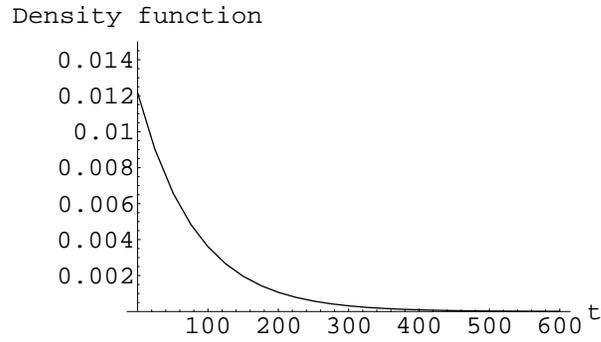


Figure 2: Queuing time density for high priority reads: moderate utilisation

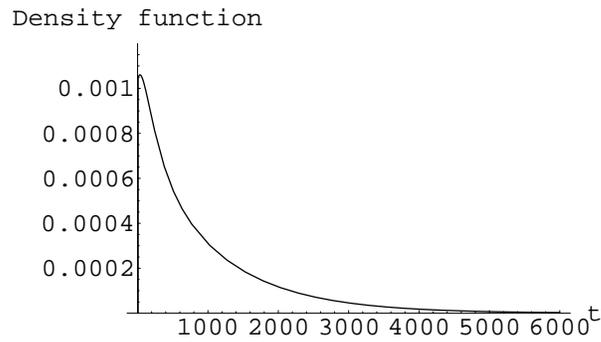


Figure 3: Queuing time density for low priority write/erases: moderate utilisation

To examine the effect of increased workload, we scaled all the arrival rates up by a factor 1.6, giving a new total average arrival rate of 28.97 vbytes/ μ sec and utilisation of 92%. The mean fluid level is now 85328.3 and its standard deviation is 96225.3. Mean queuing times grow to 365.6 and 5907.0 μ secs for the high and low priority classes respectively, and the densities are shown in figures 4 and 5.

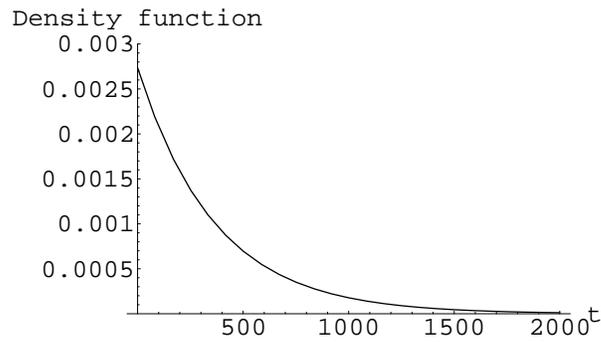


Figure 4: Queuing time density for high priority reads: high utilisation

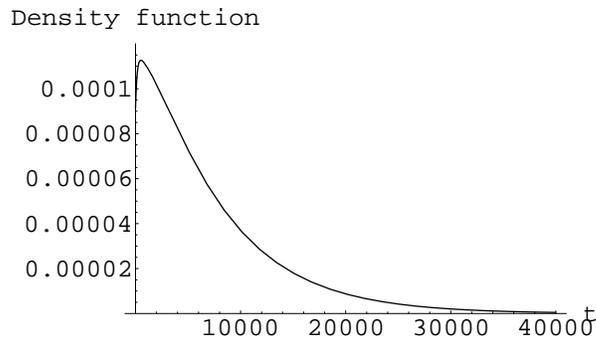


Figure 5: Queuing time density for low priority write/erases: high utilisation

5.2 Simulation

We wrote a customised event driven simulator in C to represent the Flash storage system and its associated modes of operation. The simulator handles a specific architecture for the current study but it can easily be extended to alternate Flash configurations by using their description files available in hardware libraries. In the present case, the target architecture is a package of 16 K9KAG08U0M NAND-Flash of 2GB³ each, connected by a single 40MBps channel. This Flash storage package is seen as a single address space, logical storage pool. Data is manipulated using three commands – read, write and erase – one at a time. The service times of these commands were estimated as constants by measurement on real chips. The simulator’s 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 maintaining a standard event diary. Both preemptive and non-preemptive priority mode can be accommodated, but here we consider only the preemptive variant to be consistent with the current fluid model. The simulator ensures correct implementation of the relative priorities between the two classes: a high priority class composed of reads only and a low priority class composed of writes and erases. For the preemptive mode, a read is never interrupted whereas a write or an erase (low priority event) is launched only if no reads are in the system and is interrupted as soon as a read enters the system. The interrupted write or erase is resumed as soon as there are no reads in the system, but subject to further interruption.

Simulations were run using different workloads, consistent with the stated OLTP criteria, e.g. the fixed read:write ratio of 3:1. Various request arrival rates were used, yielding package utilisation percentages ranging from 40% to 92%. We used a set of ten sets of five MMPP traces, generated for input to the simulator, so facilitating a batched means estimation of confidence bands [11]. Each trace had 100000 request arrivals. The results obtained are presented and discussed in the next subsection.

5.3 Numerical validation

A constant command service time is one of the defining characteristics of Flash and so queuing time and its variability as the package utilisation changes are quantities that have a significant impact on overall Flash performance. We therefore focus our attention here on the mean queuing time for each of the two classes (reads and writes-erases) at different arrival rates, and so at different utilisations. We also observe the evolution of the mean queuing time, as a moving average, together

³<http://www.samsung.com/global/system/business/semiconductor/product/>

with the numbers of requests in the system for both priority classes at low, medium and high package utilisations.

5.3.1 Model Validation

Simulations were run using MMPP-generated workload traces corresponding to specific sets of the model’s parameters. The results obtained by each model (fluid and simulation) therefore relate to the same context. Figures 6, 7, 8 and 9 show the queueing times obtained by simulation compared with the fluid model’s predictions for requests of high priority (reads) and low priority (writes and erases), respectively. We see reasonable agreement in the first two graphs, which relate to the model parameterisation used in section 5.1. For the reads, the graphs have exactly the same shape for both the simulation and the model, but they are separated by a small vertical shift. This is likely because the reads do not constitute “heavy traffic”, which is the domain in which fluid models perform accurately. In fact, the difference between these curves is less than one read service time (130 microseconds) and so could be accounted for largely by the difference in the models’ granularities.

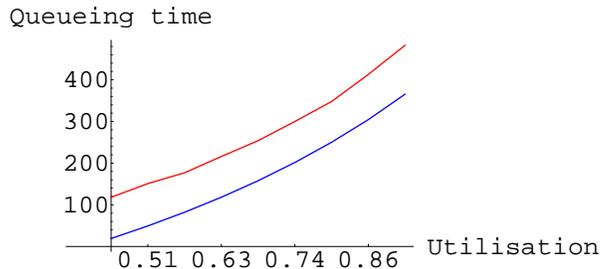


Figure 6: Queueing time for high priority class

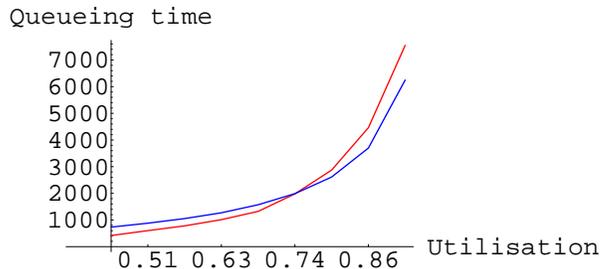


Figure 7: Queueing time for low priority class

We therefore increased the burstiness of the input traffic by modifying the state transition rate vector of the CTMC to $(1, 7, 3.5, 3)/1000$. This increased the arrival rates in each state of the CTMC considerably, it’s off-state now having equilibrium probability $\pi_0 = 0.75$. The resulting graphs showed much better agreement, although the simulation showed signs of not having converged adequately for the heavily loaded low priority class.

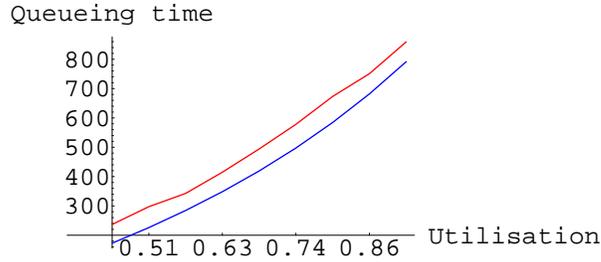


Figure 8: Queueing time for high priority class at higher burstiness

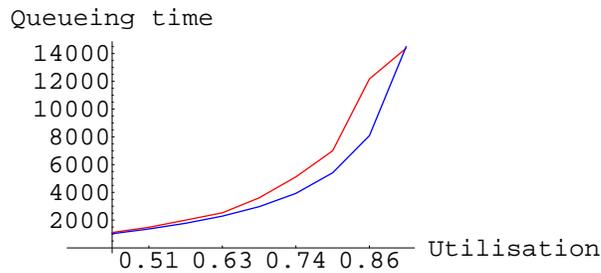


Figure 9: Queueing time for low priority class at higher burstiness

5.3.2 Queueing time vs. package utilisation

To highlight the trend of the queuing time with the growth of the package utilisation, we can see on figure 10 that it is almost constant for the high priority class (reads) but rapidly, seemingly exponentially, increasing for the low priority class (writes and erases). The sharp increase – the “knee” in the curve – begins at a

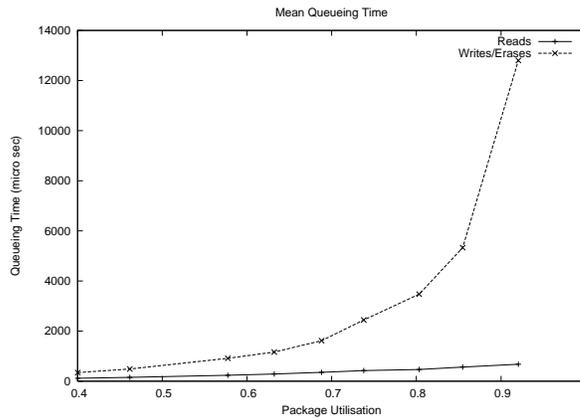


Figure 10: Queueing time Vs package utilisation

utilisation of around 70%, which is a figure that commonly arises as a threshold

in such storage and communication systems. Above this level, it is of little, and reducing, benefit to use the package.

5.3.3 Queueing time vs. number of events

The cumulative mean queueing time changes over time, increasing with the number of requests in the system before reaching a stable value. Figures 11, 12 and 13 show that this is also the case for reads, but to a much lesser extent, due to the preemptive mode of operation, regardless of the level of the Flash package’s utilisation. The same figures show clearly the evolution of the queueing time against the number of requests in the system, reaching different stable values reflecting the Flash package’s utilisation.

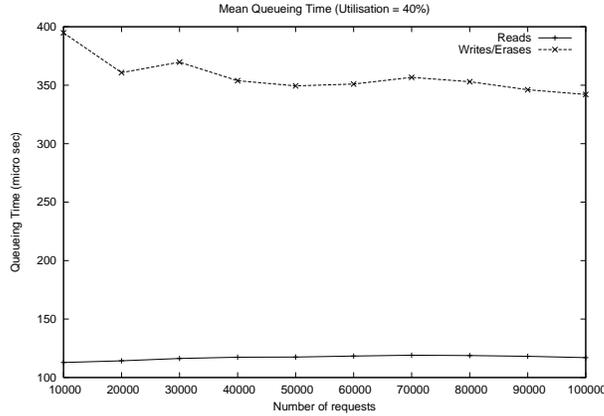


Figure 11: Queueing time Vs Requests number (Low utilisation)

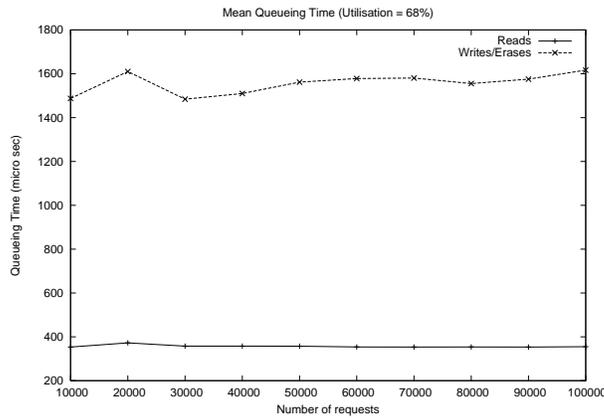


Figure 12: Queueing time Vs Requests number (Medium utilisation)

6 Conclusion and future work

With some simplifying assumptions, we have shown that fluid modelling techniques provide a promising tool for the performance evaluation of NAND Flash memory systems. This approach has a considerable advantage over traditional queueing models in that it can account for both correlated input and priorities. Moreover,

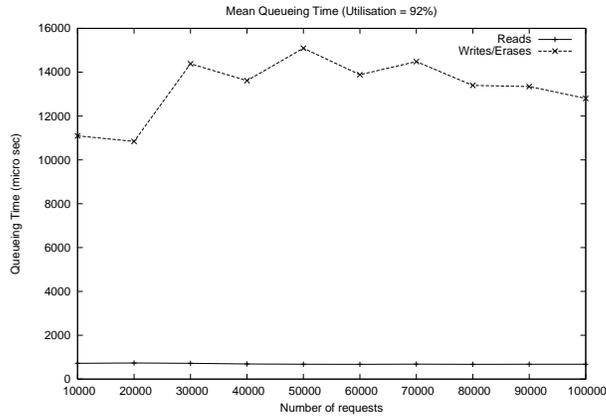


Figure 13: Queueing time Vs Requests number (High utilisation)

under heavy load, fluid models are well known to be accurate, often a stochastic limit of the analogous queue. The preliminary numerical results are encouraging, showing the right qualitative behaviour, and testing against simulation provided good quantitative validation. Further assessment against experimental observations is also crucial in the short term.

We also intend to extend the fluid model developed here, which has a preemptive priority policy, to handle non-preemptive policies that better match current Flash devices. Indeed a comparison of FCFS, preemptive resume, and non-preemptive policies would show the impact of the long erase times. As Flash devices become denser, erase times are expected to increase and even if erases are done in the background, it is possible that a long sequence of writes will use up all available free space and force erases that could delay writes. This motivates the extension of our model to investigate these effects. On the hardware side, currently, Flash memory represents the lowest level in the hierarchy. The shared channel interface and flash dye are modelled as one resource, but in fact they are separate resources that are simultaneously held during command and data transfers. Another interesting extension would be to handle command queueing inside the Flash device so that more than one memory bank could be active at any one time. Then, we can exploit the parallel access to Flash banks in order to improve the overall Flash system performance. We also intend to consider a set of Flash devices, just like disk drives, and investigate data striping and related parallel access policies, using analytical modelling and validation with respect to simulation.

Finally, we will consider alternative workloads, representing other environments taken from the wide range that use Flash devices. This will facilitate the prediction of Flash storage system behaviour under the demands of a variety of applications.

References

- [1] Leo Freitas Andrew Butterfield and Jim Woodcock. Mechanising a formal model of flash memory. *Science of Computer Programming*, 74, 2009.
- [2] Andrew Butterfield and Jim Woodcock. Formalising flash memory: First steps. In *12th IEEE International Conference on Engineering Complex Computer systems*, pages 251–260, 2007.

- [3] Li-Pin Chang and Tei-Wei Kuo. An adaptive striping architecture for flash memory storage systems of embedded systems. In *Proc. Eighteen IEEE Real Time and Embedded Technology and Application Symposium (RTAS'02)*, 2002.
- [4] Li-Pin Chang and Tei-Wei Kuo. Efficient management for large scale memory storage systems. In *Proc. ACM Symposium On Applied Computing*, 2004.
- [5] D. Mitra D. Anick and M. Sondhi. Stochastic theory of a data-handling system with multiple sources. *Bell Syst. Tech. Journal*, 8(61), 1982.
- [6] Hynix Semiconductor et al. Open nand flash interface specification. Technical report, ONFI, December 2006. www.onfi.org.
- [7] Anthony G. Field and Peter G. Harrison. An approximate compositional approach to the analysis of fluid queue network. *Performance Evaluation*, 64, 2007. Proceedings of Performance 07.
- [8] Marco Gribaudo and Miklos Telek. *Fluid Models in Performance*, volume 4486. SpringerLink, 2007.
- [9] Peter G. Harrison. A tandem network of fluid queues with on-off arrivals. Technical report, Department of Computing, Imperial College London, 2008.
- [10] Dave Hitz, James Lau, and Michael Malcolm. File system design for an nfs file server appliance. In *WTEC'94: Proceedings of the USENIX Winter 1994 Technical Conference on USENIX Winter 1994 Technical Conference*, pages 19–19, Berkeley, CA, USA, 1994. USENIX Association.
- [11] Raj Jain. *The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling*. Wiley Interscience, 1991.
- [12] Dirk Botle Jorn Engel and Robert Mertens. Garbage collection in logfs. In *Linyx conf.*, 2007.
- [13] Hyojun Kim and Seongjun Ahn. Bplru: A buffer management scheme for improving random writes in flash storage. In *Proc. 6th USENIX Conference on File and Storage Technologies*, 2008.
- [14] Seung-Ho Lim and Kyu-Ho Park. An efficient nand flash file system for flash memory storage. *IEEE Transactions On Computers*, 55(7), 2006.
- [15] Aleph One Limited. Yaffs2 readme-linux v1.1, 2007. <http://aleph1.co.uk>.
- [16] LogFS. Logfs home page. <http://logfs.org/logfs/LogFS>.
- [17] Charles Manning. Introducing yaffs, the first nand-specific flash file system, 2002. <http://www.linuxdevices.com/articles>.
- [18] Naresh M. Patel Peter G. Harrison and Soraya Zertal. Response time distribution of flash memory accesses. In *Proc. ValueTools*, 2008.
- [19] David Woodhouse. Journaling flash file system (jffs) and journaling flash filesystem 2 (jffs2). <http://sources.redhat.com/jffs2/jffs2-html>.
- [20] Chin-Hsien Wu and Tei-Wei Kuo. The design of efficient initialization and crash recovery for log-based file systems over flash memory. *ACM Transactions on Storage*, 2(4), 2006.

- [21] Chul Lee Youngwoo Park, Seung-Ho Lim and Kyu Ho Park. Pffs: A scalable flash file system for the hybrid architecture of phase-change ram and nand flash. In *2008 ACM Symposium on Applied Computing*, pages 1498–1503, 2008.
- [22] Jen-Wei Hsieh Yuan-Hao Chang and Tei-Wei Kuo. Endurance enhancement of flash memory storage systems: An efficient static wear levelling design. In *Proc. DAC*, June 2007.
- [23] Soraya Zertal. A reliability enhancing mechanism for a large flash embedded satellite storage system. In *Proc. IEEE International Conference on Systems (ICONS)*, 2008.