

DFTL: A Flash Translation Layer Employing Demand-based Selective Caching of Page-level Address Mappings

Aayush Gupta, Youngjae Kim, Bhuvan Urgaonkar
Pennsylvania State University

Speaker: Sehwan Lee

Contents

- ④ Introduction
- ④ Background and Related Work
- ④ Design of DFTL: Our Demand-based Page-mapping FTL
- ④ The FlashSim Simulator
- ④ Experimental Results
- ④ Concluding Remarks

Introduction

Using Flash Memory in Enterprise-scale Storage

- Reliability issue
- Performance issue

The Flash Translation Layer

- Out-of-place updates
- Page/Block-level FTL
- Hybrid FTL
 - Poor garbage collection behavior
 - Too many tunable parameters
 - Not exploiting the temporal locality in access

Introduction

Research Contributions

- ④ DFTL - purely page-mapped FTL
- ④ FlashSim to evaluate the efficacy of DFTL
- ④ Using a number of realistic enterprise-scale workloads

Background and Related Work

④ Basics of Flash Memory Technology

Flash Type	Data Unit Size			Access Time		
	Page (Bytes)		Block	Page	Page	Block
	Data	OOB	(Bytes)	READ (us)	WRITE (us)	ERASE (ms)
Small Block	512	16	(16K+512)	41.75	226.75	2
Large Block	2048	64	(128K+4K)	130.9	405.9	2

④ 10K-1M erase operations -> wear-leveling issue

Background and Related Work

④ Page-level and Block-level FTL Schemes

④ Page-level FTL scheme

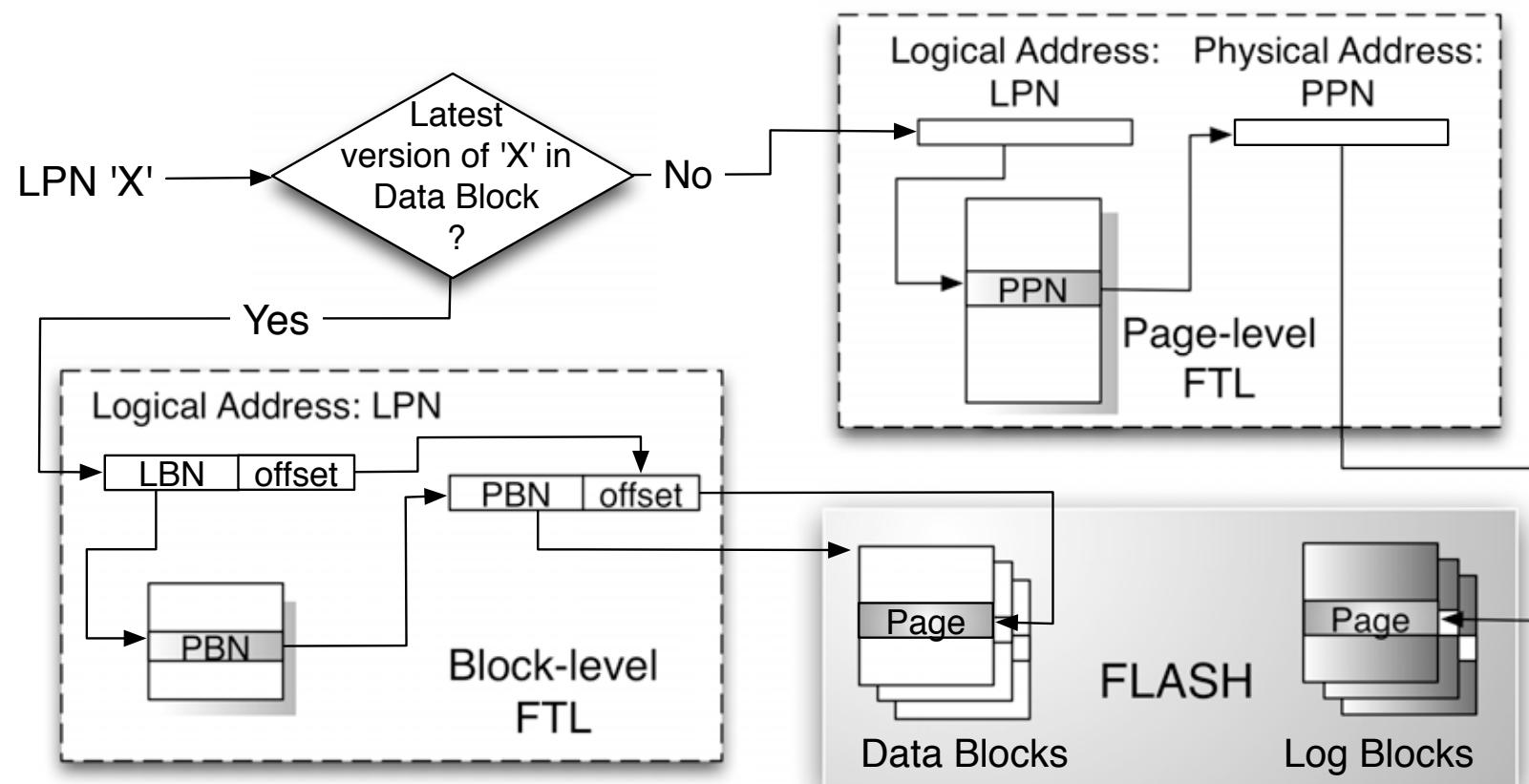
- ④ Merit: compact and efficient utilization of blocks
- ④ Demerit: large translation table

④ Block-level FTL scheme

- ④ Merit: small translation table
- ④ Demerit: garbage collection overheads + sequential programming constraint in a large block flash memory

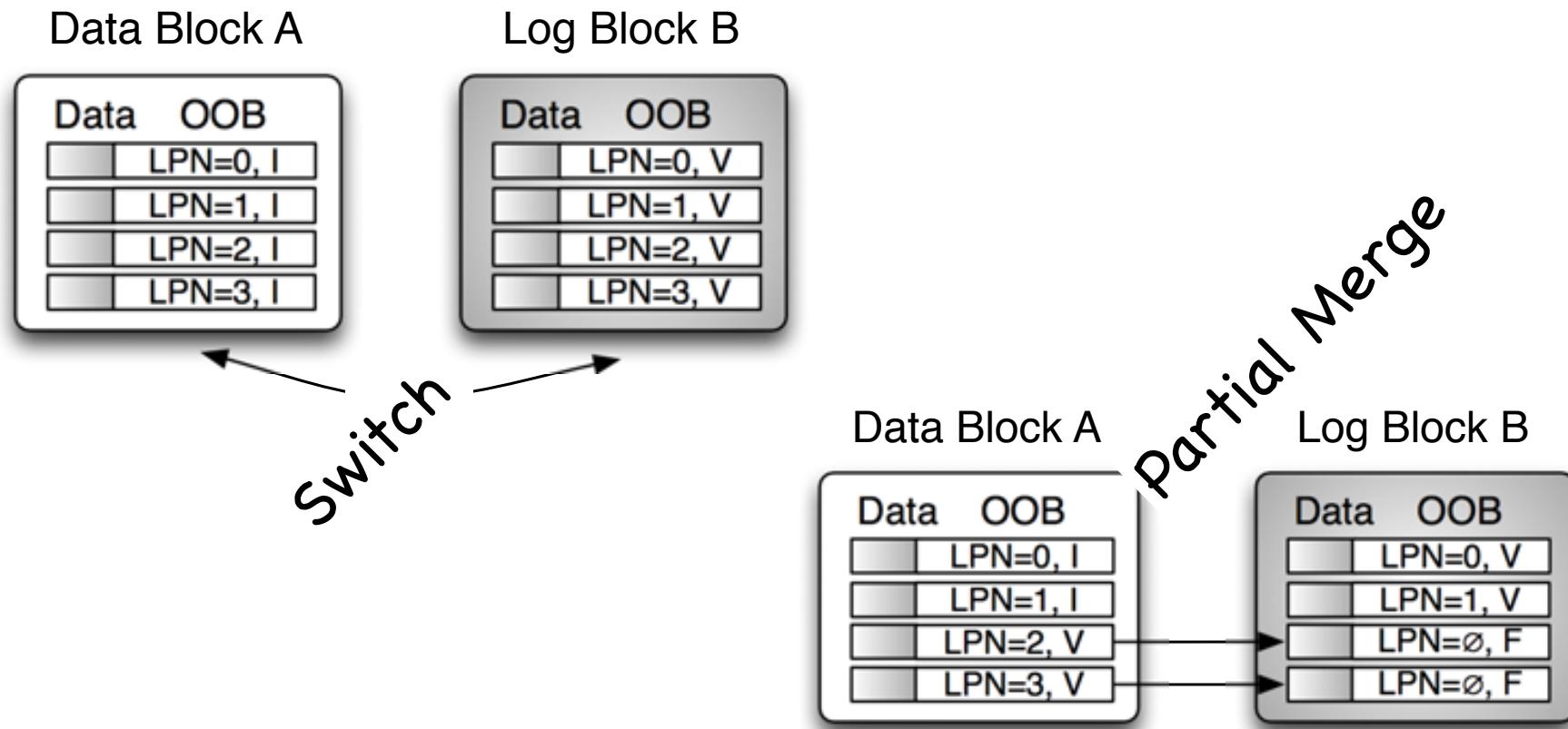
Background and Related Work

A Generic Description of Hybrid FTL Scheme



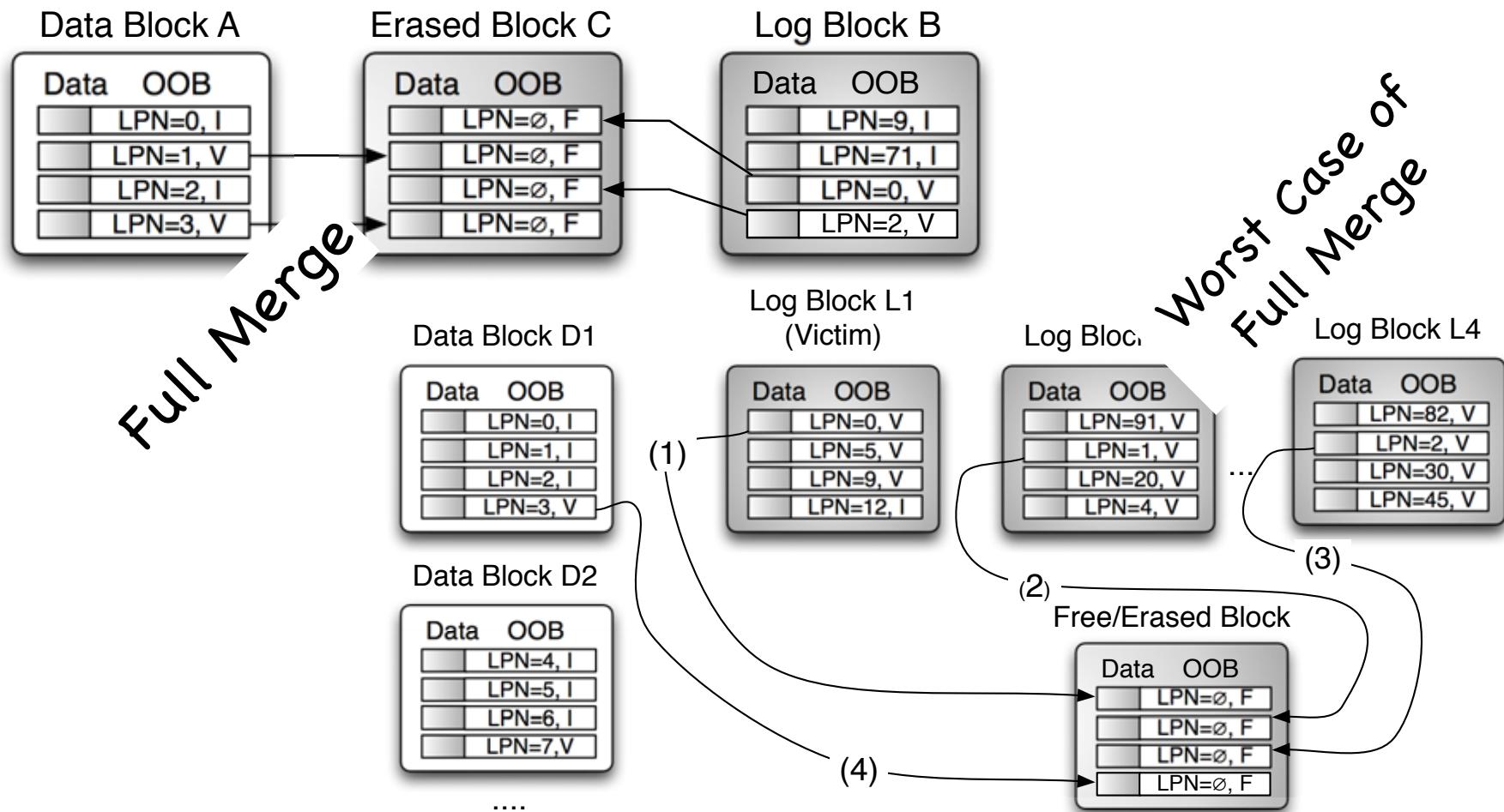
Background and Related Work

Garbage Collection in Hybrid FTLs



Background and Related Work

Garbage Collection in Hybrid FTLs



Background and Related Work

State-of-The-Art FTLs

Block Associative Sector Translation (BAST)

- data block:log block=1:m
- “log block thrashing” when many small random writes

Fully Associative Sector Translation (FAST)

- a sequential log block + random log blocks

SuperBlock FTL

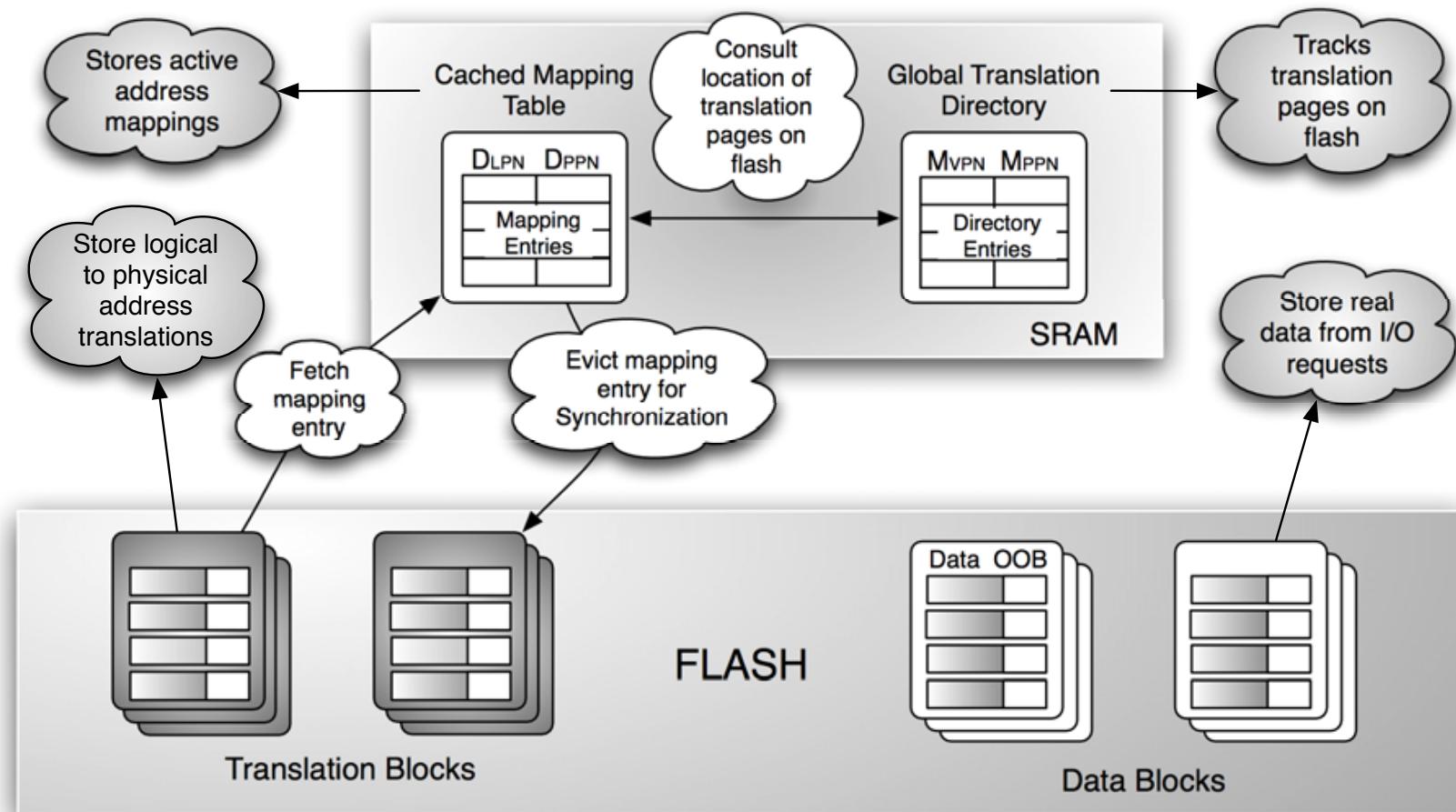
- using spatial locality in workloads

Locality-Aware Sector Translation (LAST)

- FAST + hot/cold identification to reducing full merge costs

Design of DFTL

DFTL Architecture Data Pages and Translation Pages



Design of DFTL

- ④ Logical to Physical Address Translation
 - ④ Global Mapping Table (GMT) and Global Translation Directory (GTD)
 - ④ GMT: the entire logical-to-physical address translation set in the flash memory
 - ④ CMT: active mapping present in SRAM
 - ④ GTD: translation set which keeps track of all these translation pages
 - ④ Overhead in DFTL Address Translation
 - ④ The worst case: two translation page reads + one translation page write

Design of DFTL

Logical to Physical Address Translation DFTL Address Translation Process

```
Input: Request's Logical Page Number ( $\text{request}_{lpn}$ ), Request's Size  
          ( $\text{request}_{size}$ )  
Output: NULL  
while  $\text{request}_{size} \neq 0$  do  
    if  $\text{request}_{lpn}$  miss in Cached Mapping Table then  
        if Cached Mapping Table is full then  
            /* Select entry for eviction using segmented LRU replacement  
            algorithm */  
        end  
         $\text{Translation\_Page}_{\text{request}} \leftarrow$   
        consult\_GTD( $\text{request}_{lpn}$ )  
        /* Load map entry of the request from flash into Cached Mapping  
        Table */  
        load_entry( $\text{Translation\_Page}_{\text{request}}$ )  
    end  
     $\text{request}_{type} \leftarrow$  Data Block  
     $\text{request}_{ppn} \leftarrow \text{CMT\_lookup}(\text{request}_{lpn})$   
    DFTL_Service_Request( $\text{request}$ )  
     $\text{request}_{size}--$   
end
```

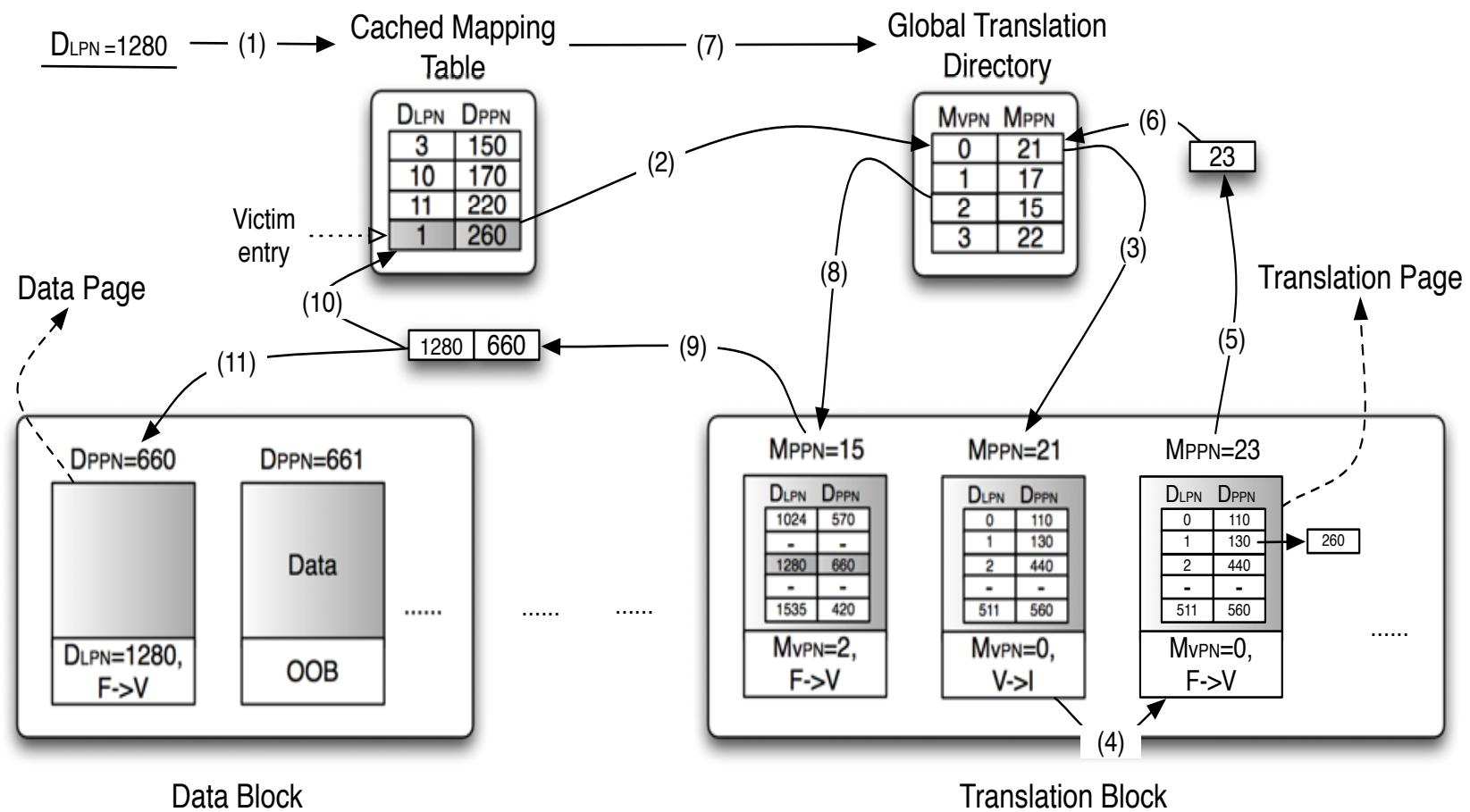
Design of DFTL

Logical to Physical Address Translation DFTL Address Translation Process

```
if Cached Mapping Table is full then
    /* Select entry for eviction using segmented LRU replacement
       algorithm */
    victimlpn ← select_victim_entry()
    if victimlast_mod_time ≠ victimload_time then
        /*victimtype : Translation or Data Block
         Translation_Pagevictim : Physical
         Translation-Page Number containing victim entry */
        Translation_Pagevictim ← consult_GTD
        (victimlpn)
        victimtype ← Translation Block
        DFTL_Service_Request(victim)
    end
    erase_entry(victimlpn)
end
```

Design of DFTL

Logical to Physical Address Translation DFTL Address Translation Process



Design of DFTL

- ④ Read/Write Operation and Garbage Collection
 - ④ Two blocks which is maintained by DFTL
 - ④ Current Data Block + Current Translation Block
 - ④ GC_{threshold}
 - ④ a high watermark which invokes a garbage collection
 - ④ victim in a translation block
 - ④ copying the valid pages to the current translation block
 - ④ updating the GTD
 - ④ victim in a data block
 - ④ copying the valid pages to the current data block
 - ④ updating all the translation pages and CMT entries associated with them
 - ④ Reducing overheads techniques
 - ④ lazy copying
 - ④ updating the CMT for those data pages whose mappings are present in it
 - ④ batch updates

Design of DFTL: Demand-based Page-mapped FTL

Comparison of Existing State-of-the-art FTLs with DFTL

	Replacement Block FTL[1]	BAST [2]	FAST [19]	SuperBlock [13]	LAST [20]	DFTL	Ideal Page FTL
FTL type	Block	Hybrid	Hybrid	Hybrid	Hybrid	Page	Page
Mapping Granularity	Block	DB-Block LB - Page	DB-Block LB-Page	SB-Block LB/Blocks within SB-Page	DB/Sequential LB - Block Random LB - Page	Page	Page
Division of Update Blocks (M)	-	-	1 Sequential + (M-1) Random	-	(m) Sequential-(M-m) (Hot and Cold)	-	-
Associativity of Blocks (Data:Update)	(1:K)	(1:M)	Random LB-(N:M-1) Sequential LB-1:1	(S:M)	Random LB-(N:M-m) Sequential LB-(1:1)	(N:N)	(N:N)
Blocks available for updates	Replacement Blocks	Log Blocks	Log Blocks	Log Blocks	Log Blocks	All Data Blocks	All Blocks
Full Merge	Yes	Yes	Yes	Yes	Yes	No	No

- Full Merge/Partial Merge
- Random Write Performance
- Block Utilization

The FlashSim Simulator

FlashSim

- ④ A simulation framework for flash-based storage systems
- ④ Modular architecture with the capability to model a holistic flash-based storage environment
- ④ Simulating multiple parts
 - ④ Device drivers+controllers+caches+flash devices+various interconnects
- ④ Implementing multiple scheme
 - ④ a block based FTL scheme (replacement-block scheme)
 - ④ FAST
 - ④ DFTL
 - ④ an idealized page-based FTL

Experimental Results

Evaluation Setup

Workloads

Workloads	Avg. Req. Size (KB)	Read (%)	Sq. (%)	Avg. Req. Inter-arrival Time (ms)
Financial [25]	4.38	9.0	2.0	133.50
Cello99 [10]	5.03	35.0	1.0	41.01
TPC-H [28]	12.82	95.0	18.0	155.56
Web Search [26]	14.86	99.0	14.0	9.97

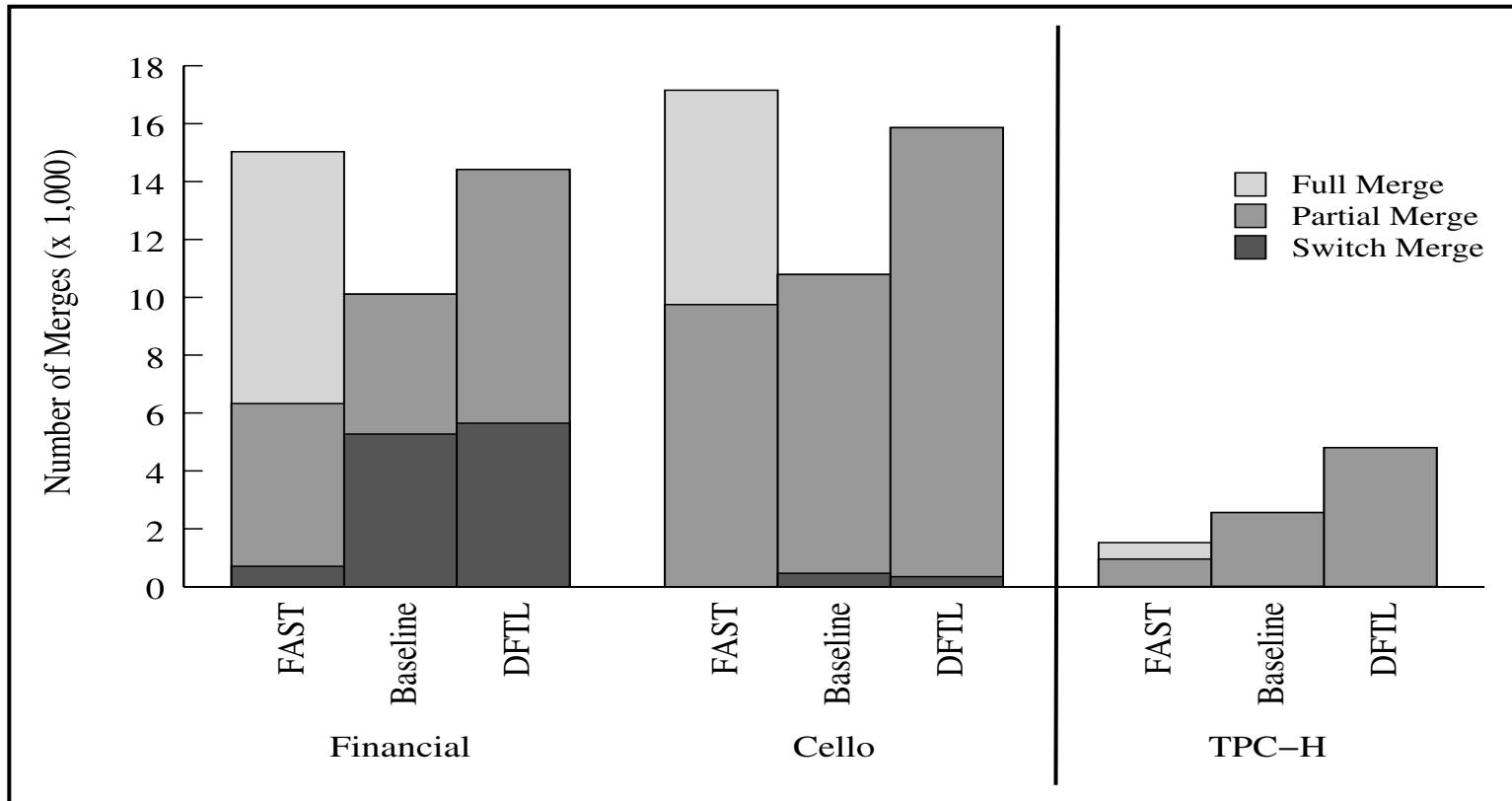
Performance Metrics

- Device service time
- indicators of the garbage collector's efficacy + response time as seen at the I/O driver

Experimental Results

Analysis of Garbage Collection and Address Translation Overheads

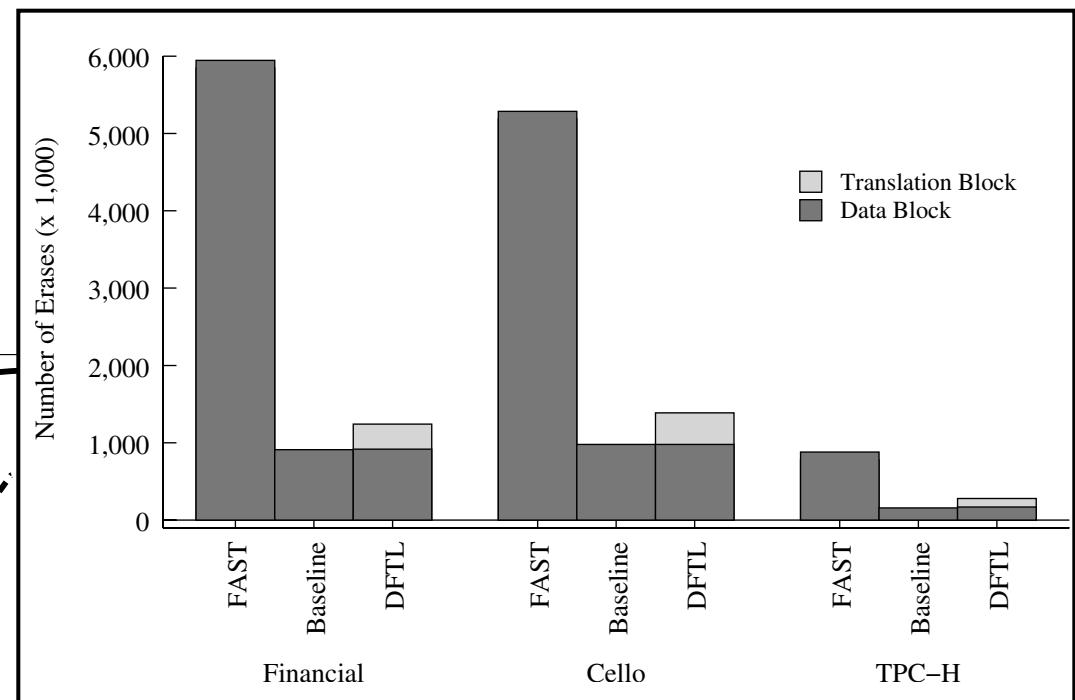
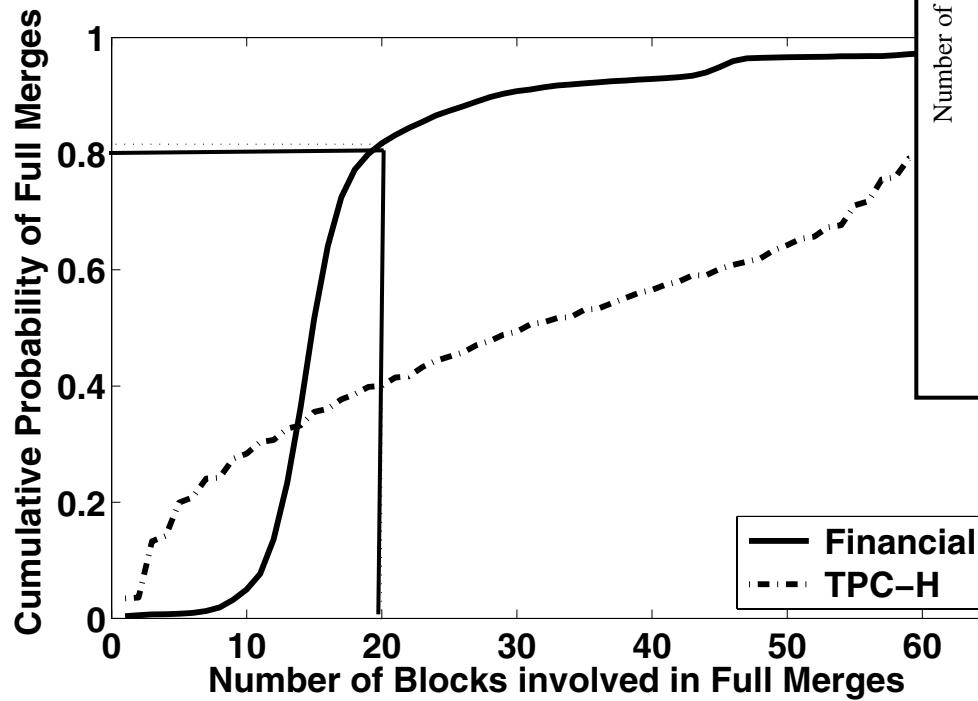
Switch Merges



Experimental Results

Analysis of Garbage Collection and Address Translation Overheads

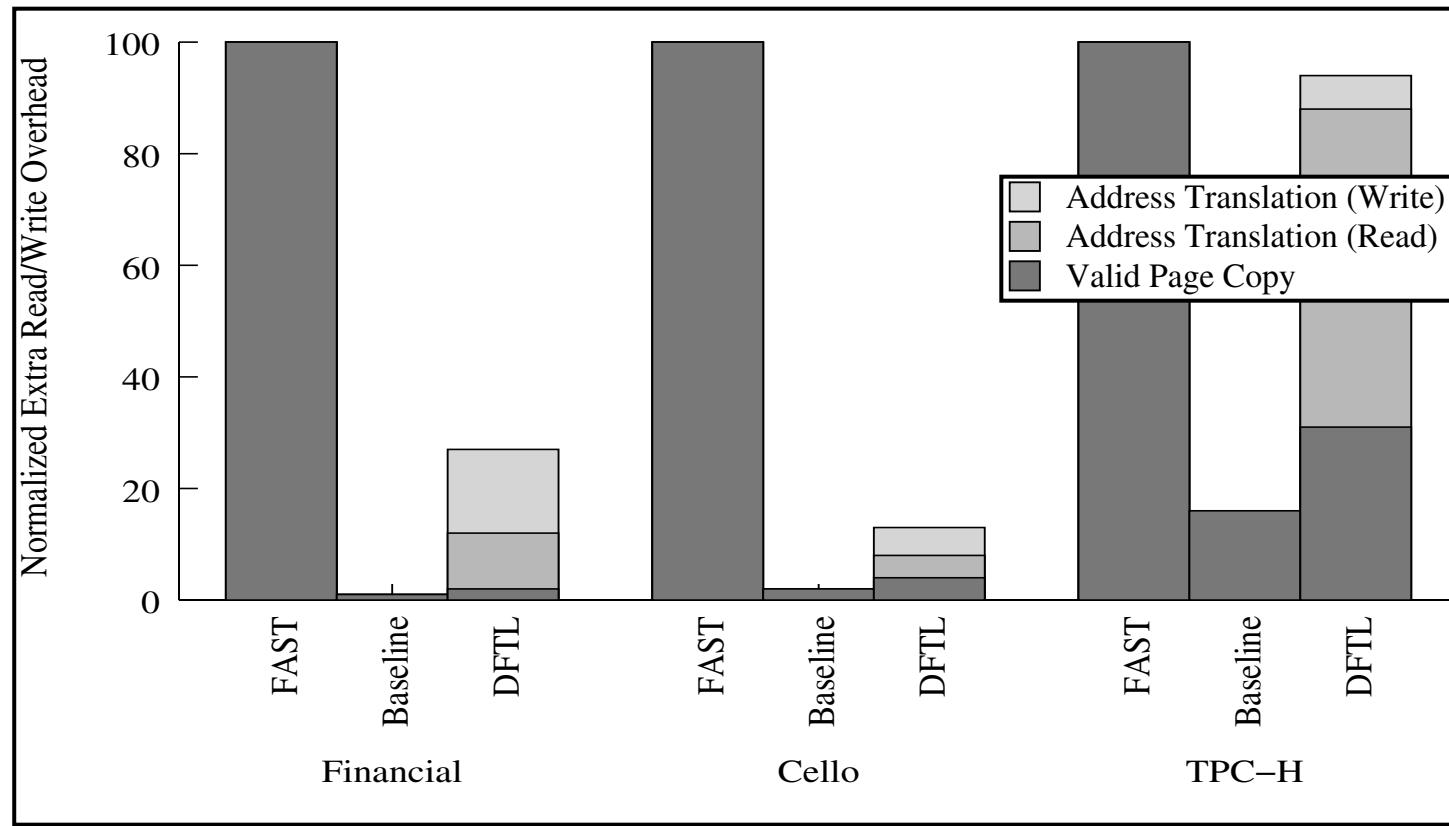
Full Merges



Experimental Results

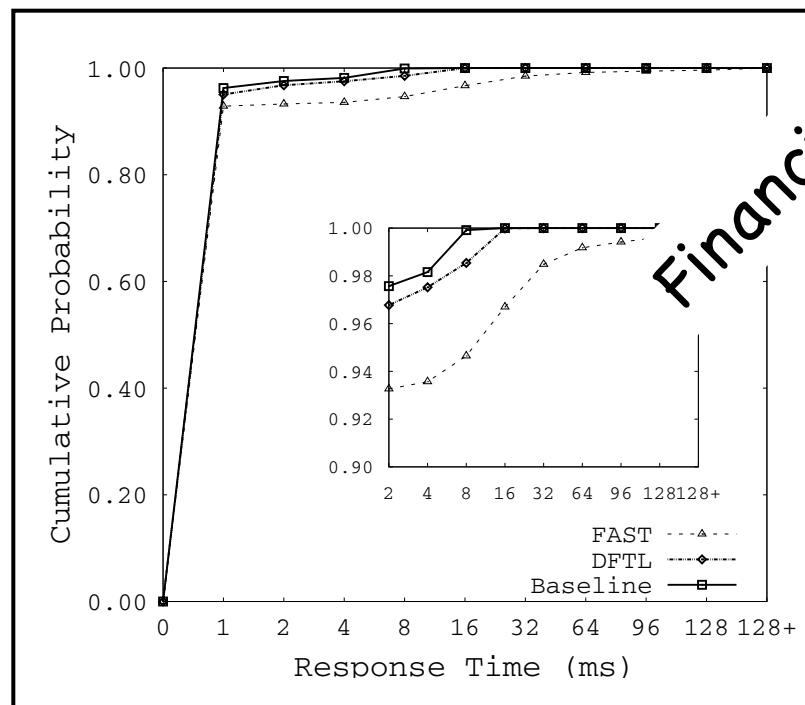
Analysis of Garbage Collection and Address Translation Overheads

Translation and Valid Page Copying Overheads

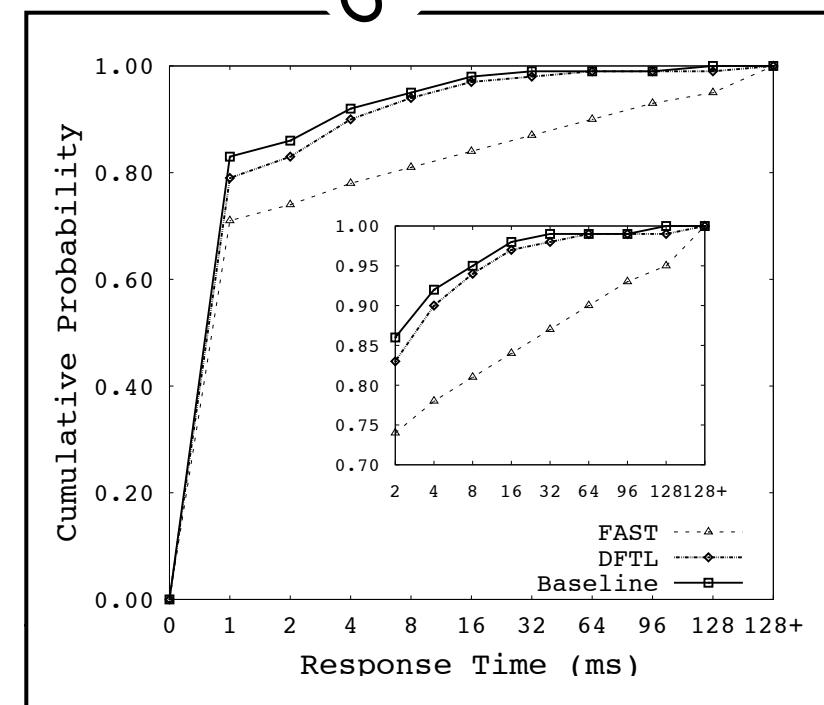


Experimental Results

Performance Analysis



Financial Trace

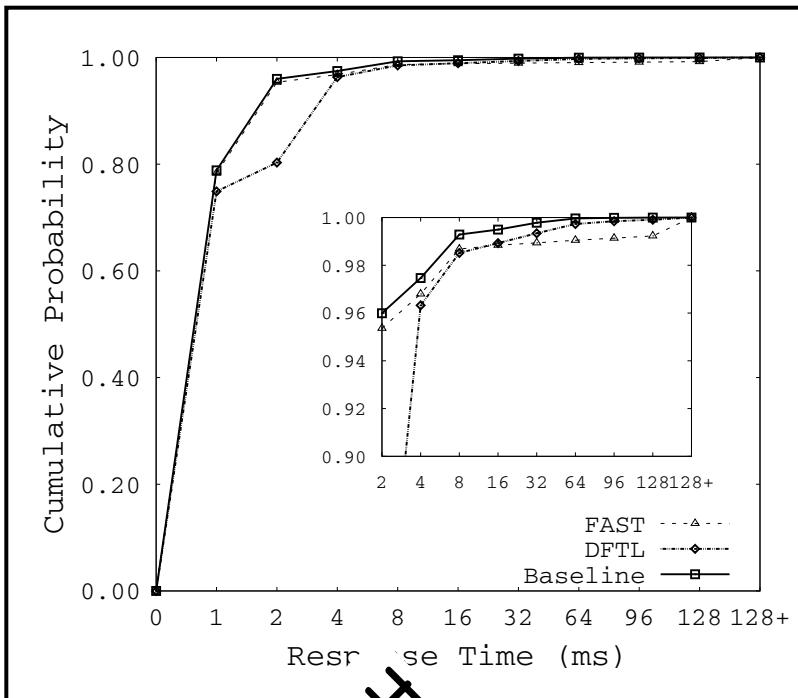


Cellogg

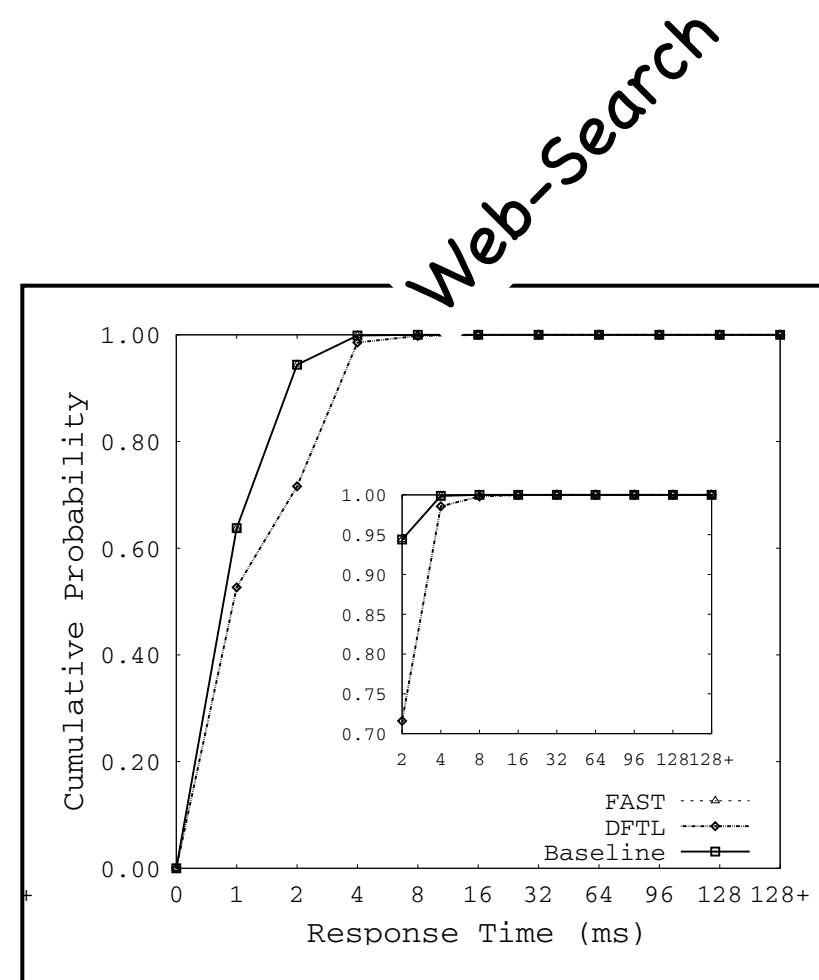
Experimental Results



Performance Analysis



TPC-H

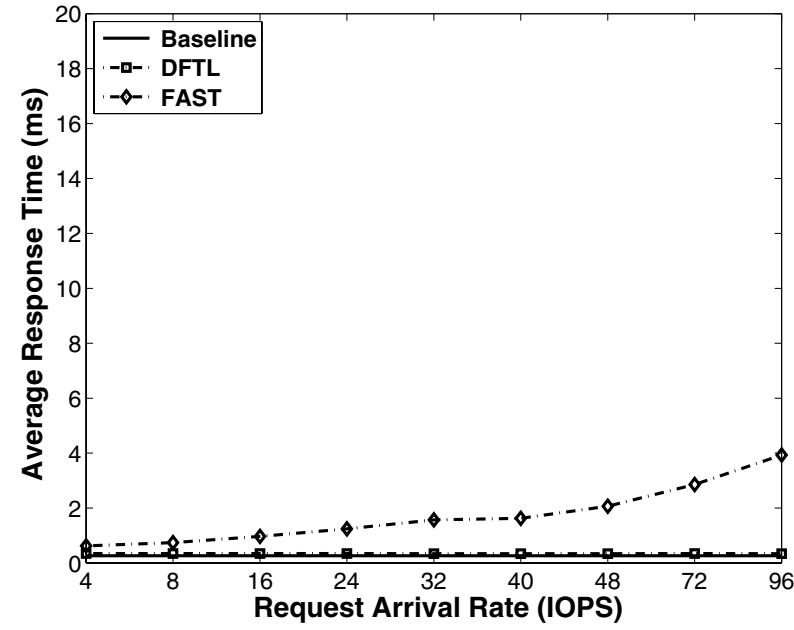
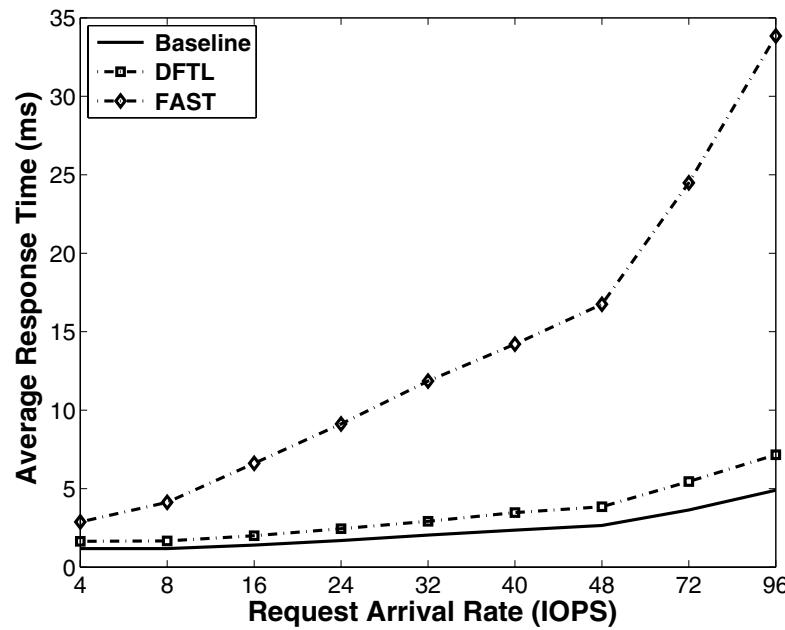


Web-Search

Experimental Results



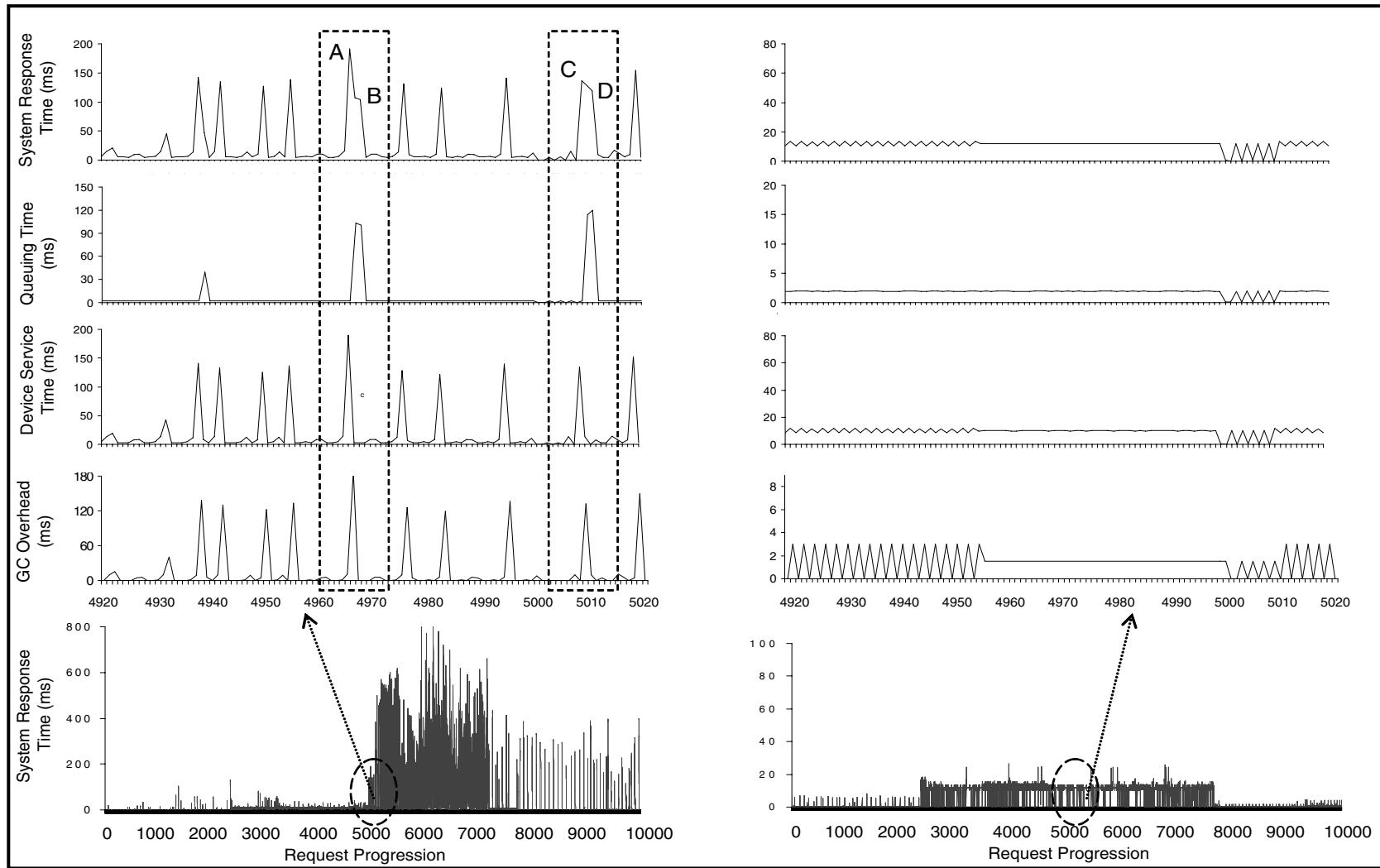
Exploring a Wider Range of Workload Characteristics



Random write dominant workload vs. Sequential write dominant workload

Experimental Results

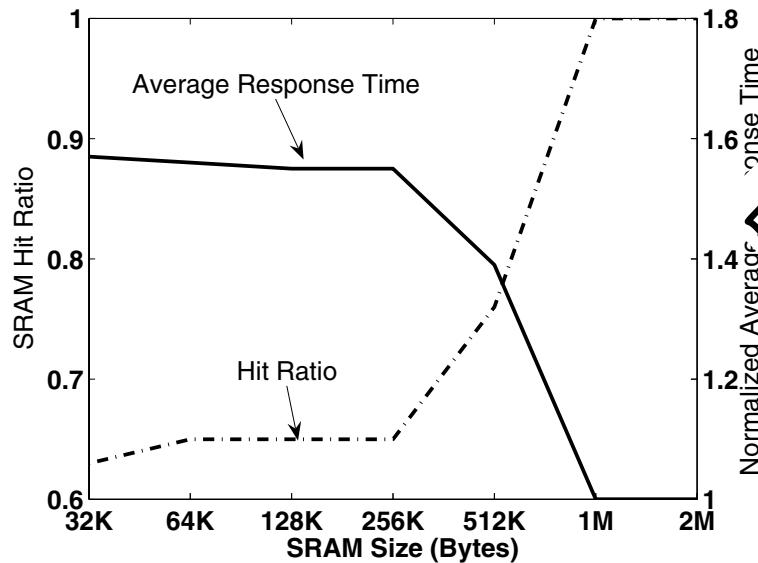
Microscopic Analysis



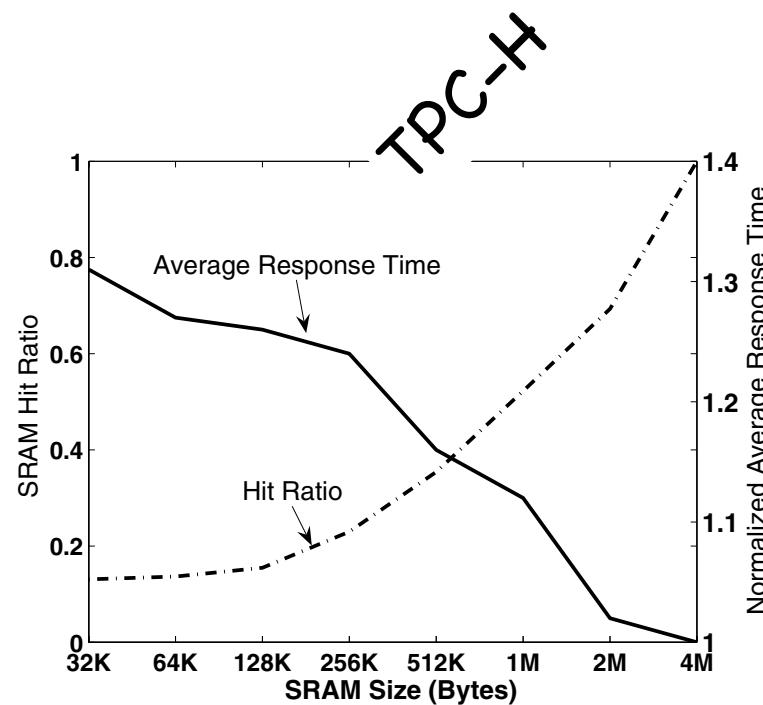
Experimental Results



Impact of SRAM size



Financial Trace



TPC-H

Concluding Remarks

- ④ We proposed a complete paradigm shift in the design of the FTL with our Demand-based Flash Translation Layer (DFTL)
- ④ DFTL is...
 - ④ improved performance
 - ④ reduced garbage collection overhead
 - ④ improved overload behavior
 - ④ free from tunable parameters