



Εθνικό Μετσόβιο Πολυτεχνείο

Σχολή Ηλεκτρολόγων Μηχανικών  
και Μηχανικών Υπολογιστών

Τομέας Τεχνολογίας Πληροφορικής  
και Υπολογιστών

## **Thesis subject**

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

**ΟΝΟΜΑ ΦΟΙΤΗΤΗ**

**Επιβλέπων :** Υπεύθυνος Διπλωματικής

Τίτλος Υπευθύνου

Αθήνα, Σεπτέμβριος 9999





Εθνικό Μετσόβιο Πολυτεχνείο

Σχολή Ηλεκτρολόγων Μηχανικών  
και Μηχανικών Υπολογιστών

Τομέας Τεχνολογίας Πληροφορικής  
και Υπολογιστών

## **Thesis subject**

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

**ΟΝΟΜΑ ΦΟΙΤΗΤΗ**

**Επιβλέπων :** Υπεύθυνος Διπλωματικής

Τίτλος Υπευθύνου

Εγκρίθηκε από την τριμελή εξεταστική επιτροπή την 9η Σεπτεμβρίου 9999.

.....  
Πρώτο μέλος επιτροπής  
Τίτλος μέλους

.....  
Δεύτερο μέλος επιτροπής  
Τίτλος μέλους

.....  
Τρίτο μέλος επιτροπής  
Τίτλος μέλους

Αθήνα, Σεπτέμβριος 9999

.....  
**Όνομα Φοιτητή**

Διπλωματούχος Ηλεκτρολόγος Μηχανικός και Μηχανικός Υπολογιστών Ε.Μ.Π.

Copyright © Όνομα Φοιτητή, 9999.

Με επιφύλαξη παντός δικαιώματος. All rights reserved.

Απαγορεύεται η αντιγραφή, αποθήκευση και διανομή της παρούσας εργασίας, εξ ολοκλήρου ή τμήματος αυτής, για εμπορικό σκοπό. Επιτρέπεται η ανατύπωση, αποθήκευση και διανομή για σκοπό μη κερδοσκοπικό, εκπαιδευτικής ή ερευνητικής φύσης, υπό την προϋπόθεση να αναφέρεται η πηγή προέλευσης και να διατηρείται το παρόν μήνυμα. Ερωτήματα που αφορούν τη χρήση της εργασίας για κερδοσκοπικό σκοπό πρέπει να απευθύνονται προς τον συγγραφέα.

Οι απόψεις και τα συμπεράσματα που περιέχονται σε αυτό το έγγραφο εκφράζουν τον συγγραφέα και δεν πρέπει να ερμηνευθεί ότι αντιπροσωπεύουν τις επίσημες θέσεις του Εθνικού Μετσόβιου Πολυτεχνείου.

## **Περίληψη**

Περίληψη της διπλωματικής.

### **Λέξεις κλειδιά**

Λέξη-κλειδί 1, λέξη-κλειδί 2, λέξη-κλειδί 3



## **Abstract**

Abstract of diploma thesis.

## **Key words**

Key-word 1, Key-word 2, Key-word 3





# Ευχαριστίες

Ευχαριστίες.

Όνομα Φοιτητή,  
Αθήνα, 9η Σεπτεμβρίου 9999

Η εργασία αυτή είναι επίσης διαθέσιμη ως Τεχνική Αναφορά CSD-SW-TR-\*-\* , Εθνικό Μετσόβιο Πολυτεχνείο, Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών, Τομέας Τεχνολογίας Πληροφορικής και Υπολογιστών, Εργαστήριο Τεχνολογίας Λογισμικού, Σεπτέμβριος 9999.

URL: <http://www.softlab.ntua.gr/techrep/>  
FTP: <ftp://ftp.softlab.ntua.gr/pub/techrep/>



# Contents

<b>Περίληψη</b>	5
<b>Abstract</b>	7
<b>Ευχαριστίες</b>	9
<b>Contents</b>	11
<b>List of Figures</b>	13
<b>List of Tables</b>	15
<b>1. Introduction</b>	19
1.1 Introduction/Motivation	19
1.2 Thesis structure	19
<b>2. Chapter 2</b>	21
2.1 Section 1	21
2.1.1 Subsection 1	21
2.2 Section 2	21
2.2.1 Subsection 1	21
2.2.2 Subsection 2	21
<b>3. Chapter 3</b>	23
3.1 Necessary theoretical background	23
3.1.1 Multi-threaded programming	23
3.1.2 Typical IPC	23
3.2 Archipelago	23
3.3 XSEG	24
3.3.1 Drivers	24
3.3.2 Libraries	24
3.3.3 Xtypes	24
3.3.4 Peers	24
3.3.5 Archipelago IPC	25
3.4 Request flow example	25
3.4.1 Get request	25
<b>4. Tiering</b>	27
4.1 Theoretical Background	27
4.1.1 Caching	27
4.2 Existing storage tiers	28
4.2.1 Bcache	28
4.2.2 Memcached	28

4.2.3	Blabla	28
4.2.4	Summary	28
<b>5.</b>	<b>Design of cached</b>	<b>29</b>
5.1	Design overview	29
5.1.1	Cached components	30
5.2	The xcache xtype	30
5.2.1	Entry Preallocation	31
5.2.2	Entry indexing	31
5.2.3	Entry eviction	31
5.2.4	Concurrency control	32
5.2.5	Re-insertion	33
5.2.6	Event hooks	33
5.2.7	xcache flow	33
5.3	The xworkq xtype	34
5.4	The xwaitq xtype	34
5.5	Cached internals	34
5.5.1	Buckets	35
5.5.2	Object states	35
5.6	Bucket pool	35
5.6.1	Per-object peer requests	35
5.6.2	Bucket/Object states	35
5.6.3	Write policy	35
5.7	Cached Operation	36
5.7.1	Polling for new requests	36
5.7.2	Write-through mode	36
5.7.3	Write-back mode	36
<b>6.</b>	<b>Implementation of cached</b>	<b>37</b>
6.1	Implementation of xcache	37
6.1.1	Entry Preallocation	37
6.1.2	Entry Indexing	38
6.1.3	Entry eviction	39
6.1.4	Concurrency control	40
6.1.5	Event hooks	41
6.2	Implementation of cached	41
6.3	Bucket pool	42
6.3.1	Bucket/Object states	42
<b>7.</b>	<b>Performance evaluation of cached</b>	<b>45</b>
7.1	Benchmark methodology	45
7.2	Specifications of test-bed	46
7.3	Performance comparison between cached and sosd	47
	<b>Bibliography</b>	<b>51</b>

**List of Figures**

7.1 Comparison of bandwidth performance for writes . . . . . 47

7.2 Comparison of bandwidth performance for reads . . . . . 48

7.3 Comparison of latency performance for writes . . . . . 49

7.4 Comparison of latency performance for reads . . . . . 50



**List of Tables**

6.1    [Reference counting of xcache](#) . . . . . 41

7.1    [dev100 specs](#) . . . . . 46





## List of Listings

2.1	Sample code . . . . .	21
6.1	Main <code>xcache</code> struct . . . . .	37
6.2	<code>xcache</code> struct fields for preallocated entries . . . . .	37
6.3	<code>xcache</code> entry struct . . . . .	38
6.4	<code>xcache</code> struct fields for entry indexing . . . . .	38
6.5	Indexing functions . . . . .	39
6.6	<code>xcache</code> struct fields for eviction . . . . .	39
6.7	Doubly-linked LRU list . . . . .	40
6.8	Main cached struct . . . . .	41
6.9	Cached entry struct . . . . .	42
6.10	Bucket implementation . . . . .	42



# Chapter 1

## Introduction

### 1.1 Introduction/Motivation

Bla-bla...

### 1.2 Thesis structure

**Chapter 2:** We define what "cloud" means and mention some of the most notable examples. Then, we give a brief overview of the synnefo implementation, its key characteristics and why it can have a place in the current cloud world.

**Chapter 3:** We present the architecture of Archipelago and provide the necessary theoretical background (mmap, IPC) the reader needs to understand its basic concepts. Then, we thoroughly explain how Archipelago handles I/O requests. Finally, we mention what are the current storage mechanisms for Archipelago and evaluate their performance.

**Chapter 4:** We explain why tiering is important and what is the state of tiered storage at the moment (bcache, flashcache, memcached, ramcloud, couchbase). Then, we provide the related theoretical background for cached (hash-tables, LRUs). Finally, we defend why we chose to roll out our own implementation.

**Chapter 5:** We explain the design of cached, the building blocks that is consisted of (xcache, xworkq, xwaitq). Then, we give some examples that illustrate the operation under different scenarios

**Chapter 6:** We present the cached implementation, the structures that have been created and the functions that have been used.

**Chapter 7:** We explain how cached was evaluated and present benchmark results.

**Chapter ??:** It connects brain parts. And its tale must be told.

**Chapter ??:** We draw some concluding remarks and propose some future work.



## Chapter 2

## Chapter 2

### 2.1 Section 1

This section has an important citation[[aeal99](#)]

#### 2.1.1 Subsection 1

This subsection has code in Haskell:

```
1 foo [] = []  
2 foo h:t = 9: foo t
```

**Listing 2.1:** Sample code

It also has a list:

**Item 1** First item

**Item 2** Second item and a footnote<sup>1</sup>.

**Item 3** Third item and text in *italics*.

And an enumerated list:

1. First item.
2. Second item and text in **bold**

### 2.2 Section 2

#### 2.2.1 Subsection 1

This subsection has a link to the block of code [2.1](#) in Section 1.

#### 2.2.2 Subsection 2

This subsection has a FIXME comment, visible only to the author.

---

<sup>1</sup> Footnote description.



## **Chapter 3**

## **Chapter 3**

### **3.1 Necessary theoretical background**

#### **3.1.1 Multi-threaded programming**

Multi-threading programming is good and is bad and here are some challenges:

1. Concurrency control
2. Challenge 2
3. Challenge 3

#### **Concurrency control**

**Locking** Three concepts for locking:

1. Lock overhead
2. Lock contention
3. Deadlocking

#### **3.1.2 Typical IPC**

Below we can see some IPC methods:

1. mmap()
2. Semaphores
3. Sockets

### **3.2 Archipelago**

Archipelago consists of the following:

1. XSEG 2. 3.

## 3.3 XSEG

XSEG is the segment on which the IPC...

There are some XSEG stuff such as:

1. Drivers 2. Libraries 3. Xtypes 4. Peers

### 3.3.1 Drivers

### 3.3.2 Libraries

### 3.3.3 Xtypes

The rationale behind xtypes is:

- Abstraction(?) layers: Creating inner abstractions layers for software is not a new concept but it's very easy to miss, especially when you start small and end up big.

In a nutshell, when writing code for a new software (in our case a peer for Archipelago but this can apply to most software that surpass the 1000 LOC<sup>1</sup> mark) it is wrong practice to create from scratch a monolithic implementation with indistinguishable parts. There is a main reason for this:

Monolithic implementations usually derive from lack of code architecture and planning. Although it is feasible for a programmer to create fully-functional code that meets the necessary requirements, albeit with a lot more effort and concentration, this approach will backfire when the programmer needs to add new features. Since there is no explicit code architecture and the fragile inner correlations are between lines of code and not separate entities, stored precariously in the developer's mind, the result will eventually be constant code refactorization.

One might think that new features happen once in a while in the development cycle but that would be wrong. This happens more often than you might think and is actually the common case in iteration and test-driven development.

The right practice instead is to...

- Re-usability:...
- User-space / Kernel-space agnosticity: (I doubt that such a word even exists...)

### 3.3.4 Peers

Peers are Archipelago components that are responsible for accepting, processing and sending of the I/O requests. They are essential for the modular nature of Archipelago since each of them can be considered as a separate entity. They do their own logging, signal handling and processing.

The main Archipelago peers can be seen in Figure ?. As we can see from this figure, peers are processes that are attached to an XSEG segment. In the previous chapter, we have mentioned that XSEG segments facilitate the IPC between different Archipelago components by offering a shared space where process can read and write to very fast. This however barely scratches the surface of IPC in Archipelago. In the following section, we will discuss more in-length the details behind Archipelago IPC

---

<sup>1</sup> Lines Of Code



### 3.3.5 Archipelago IPC

First of all, we must clarify that in Archipelago, IPC is done strictly between peers in the **same** memory segment. The reason is that we have crafted our own methods for IPC and the processes that need it must attain to a certain architecture, which is the peer architecture.

The entrance point for IPC is the peer port. When a peer is registered in the segment, it attaches itself to a port range. Peer ports are completely different to common ports (which are these ports?). When a peer wants to send a request to another peer, it must first "get" the registered port on the segment. The xseg port is a structure that holds the necessary information as to where to send the request. Every port has three different queues; reply, request, free queue.

Request queues are typically a stack that can be addressed from different peers in the same segment. For this reason, they are designed as xtypes. For speed reasons, they are pre-allocated to a certain length and re-allocated on-line, if there is need

Also, ports are designed to be considered as paths. That is, when a request is sent from one port to another...

## 3.4 Request flow example

We have bench xseg which works like so:

1. Get request
2. Prepare request
3. Create chunk
4. Allocate peer request
5. Set request (xhash)
6. Submit request

### 3.4.1 Get request

Explain here about xq or in xtypes?



## Chapter 4

# Tiering

### 4.1 Theoretical Background

#### 4.1.1 Caching

In caching, there are usually the following two policies:

- Write-through: This policy bla bla bla
- Write-back: This policy blu blu blu

#### Eviction

Caching generally means that you project a large address space of a slow medium to the smaller address space of a faster medium. That means that not everything can be cached as there is no 1:1 mapping. So, when a cache reaches its maximum capacity, it must evict one of its entries

And the big question now arises: which entry?

This is a very old and well documented problem that still troubles the research community. It was first faced when creating hardware caches (the L1, L2 CPU caches we are familiar with). In 1966, Lazlo Belady proved that the best strategy is to evict the entry that is going to be used more later on in the future[Bela66]. However, the clairvoyance needed for this strategy was a little difficult to implement, so we had to resort to one of the following, well-known strategies:

- **Random:** Simply, a randomly chosen entry is evicted. This strategy, although it seems simplistic at first, is sometimes chosen due to the ease and speed of each. It is preferred in random workloads where getting fast free space for an entry is more important than the entry that will be evicted.
- **FIFO (First-In-First-Out):** The entry that was first inserted will also be the first to evict. This is also a very simplistic approach as well as easy and fast. Interestingly, although it would seem to produce better results than Random eviction, it is rarely used though, since it assumes that cache entries are used only once, which is not common in real-life situations.
- **LRU (Least-Recently-Used)**
- **LFU (Least-Frequently-Used)**

Choosing the LRU strategy is usually a no-brainer. Not only does it *seem* more optimal than the other algorithms, but it has also been proven, using a Bayesian statistic model, that no other algorithm that tracks the last K references to an entry can be more optimal.

## **4.2 Existing storage tiers**

### **4.2.1 Bcache**

### **4.2.2 Memcached**

### **4.2.3 Blabla**

### **4.2.4 Summary**

## Chapter 5

### Design of cached

In the previous chapters, we have addressed the need for tiering in terms of scalability as well as performance.

We have also evaluated current caching solutions and described why they couldn't be used as a cache tier in Archipelago.

With the results of chapter ? in mind, we can provide some more strict requirements that our solution must have:

1. **Nativity:** Our solution must be native to Archipelago i.e. not need any translation layers to communicate with it.
2. **Pluggability:** Our solution must be able to provide a caching layer between peers that are already in operating mode without restarting Archipelago.
3. **Object awareness:** Our solution must be able to operate on object level, which are the data entities Archipelago understands.
4. **In-memory:** Our solution must cache requests in RAM, since the next fastest tier, SSDs, are already being used in RADOS as a journal.

For the following chapters, we will drop the "*solution*" moniker and we will use instead the proper name of our implementation, "cached", which simply means **cache daemon**).

The following two chapters are the main bulk of this thesis and they present our own implementation that aims to fill the above requirements.

More specifically, this chapter provides an in-depth description of the design of cached. Section ? provides a general overview of cached. Sections ? - ? present the building blocks of cached and their design. Section ? presents the interaction of cached and its building blocks. Finally, in Section ? we illustrate the flow of requests for cached.

#### 5.1 Design overview

First of all, cached has been designed as an Archipelago user-space peer (see Section 3.3.4 about Archipelago peers). This design decision covers the nativity requirement we posed at the beginning of this chapter.

Also, cached's purpose is to provide a caching layer between the vlmc and blocker. Due to the fact that cached is a peer, this is easily achievable even under normal operation because:

1. As we have mentioned in Section 3.3.5, XSEG ports can be registered on-line.
2. During normal operation, the administrator can add the cached port to the request path between vlmc and blocker, and all requests will seamlessly be intercepted by cached. This follows the same principle with bcache, which plugs its own request\_fn() function to the virtual device it creates. Unlike bcache however, cached can be plugged on and off at any time.

Thus, the pluggability requirement has also been covered.

Furthermore, one of the most important aspects of cached's design is the index mechanism. Specifically, we have utilized an in-memory hash table to index our cached objects.

Moreover, cached has been designed to be object aware. This means that we search in the hash table with the object name as key.

Finally, an important aspect of cached is its multi-threading support. Specifically, cached can work with multiple threads that can accept requests from cached's request queue and serve them concurrently with the other threads. Of course, multi-threading can be very tricky, especially when we are dealing with I/O requests and simultaneous accesses to the same object blocks. So, in order to achieve a balance between safety and speed, we use a fine-grained locking scheme in critical sections that can be seen is discussed in detail in Section 5.3.

### 5.1.1 Cached components

Let's see now the design of cached in detail. The cached peer consists of a number of building blocks. Per Archipelago policy, most of these building blocks have been written in the xtypes fashion (see Section ?? about xtypes). Also, we must note that these components were not created by the author and the contribution to them was limited.

The components of cached can be seen below:

- xcache, an xtype that provides indexing support, amongst many other things
- xworkq, an xtype that guarantees atomicity for execution of jobs on the same object
- xwaitq, an xtype that allows conditional execution of jobs

and their design will be discussed in-depth in the following sections.

Also, we must note that the above components predate our cached implementation and are not a contribution of this thesis<sup>1</sup>. They are presented however in this thesis for clarity reasons.

## 5.2 The xcache xtype

xcache is the main component of cached. It is responsible for several key aspects of caching such as:

- entry indexing,
- entry eviction, and

---

<sup>1</sup> xcache is an exception since we have extended its functionalities for our purposes

- concurrency control
- event hooks

We have to note here that xcache pre- our contribution to xcache

Below we can see a design overview of xcache:

As we can see above, xcache utilizes two hash tables. One hash table is responsible for indexing entries (or more generally speaking "cache entries") that are active in cache. The other hash table is responsible for indexing evicted cache entries that have pending jobs. Again, more generally speaking, evicted cache entries are entries whose refcount has not dropped to zero yet.

### 5.2.1 Entry Preallocation

Since xcache has a bounded number of entries that will allocate, there is no need to allocate them on-the fly using malloc/free. Considering that we are caching at RAM level and not at SSD level, the system call overhead will have a considerable impact on performance.

Thus, the best thing to do in our case would be to pre-allocate the necessary space.

### 5.2.2 Entry indexing

In order to index the cached entries, xcache relies on another xtype, xhash, which is a hash table. Moreover, it's actually the C implementation of the dictionary used in Python.

We have chosen to use a hash table as our index because:

Finally, the xhash xtype gives provides us with the basic hash table functions, namely:

- Insertion
- Look-up
- Deletion

### 5.2.3 Entry eviction

As we can see in figure ?, xcache has been designed to index a pre-defined number of entries. That means that when xcache reaches its maximum capacity and is requested to index a new entry, it has to resort to the eviction of a previously cached entry. We have chosen the LRU strategy

Also, an added bonus is that we won't need to sacrifice speed over optimality, since that, our hash table approach allows us to create an  $O(1)$  LRU algorithm which you can see in the following figure:

In a nutshell, our LRU implementation uses a doubly linked list and utilize the hash table to jump to the element (instead of traversing the list linearly). This design allows us to do all of the following action in constant time:

- Insert a new entry to the LRU list
- Evict the LRU entry

- Update an entry's access time (i.e. mark it as MRU)
- Remove an arbitrary entry

Another interesting feature of xcache is that evictions occur implicitly and not explicitly. The user doesn't need to interact with the LRU queue.

For example, when a user tries to insert a new entry to an already full cache, the insertion will succeed and the user will not be prompted to evict an entry manually. Also, the user will be notified via specific event hook that is triggered upon eviction that an entry has been evicted.

More about hooks can be seen in the following subsection.

### 5.2.4 Concurrency control

The concept of concurrency control has been discussed in chapter ?. The goal of xcache is to handle safely - and preferably fast - simultaneous accesses to shared memory.

In order to do so, we must first identify which are the critical sections of xcache, that is the sections where a thread can modify a shared structure. These sections are the following

- All xhash operations: Two of the three xhash operations (inserts and removals) can modify the hash table (e.g. they can resize it, add more entries or delete existing ones). This means that the third one (lookups) must not run concurrently with the other.
- Cache node claiming: Before an entry is inserted, it must acquire one of the pre-allocated nodes and we must ensure that this can happen from all threads.
- Entry migration: An entry can migrate from one hash table to the other e.g. on cache eviction. This migration involves a series of xhash operations; removal from one hash table and subsequent insertion to the other. This a scenario that must be handled properly.
- Reference counting: Every entry must have a reference counter. Reference counters provide a simple way to determine when an entry can be safely removed. You can see more about reference counting in chapter ?
- LRU updates: Most actions that involve cache entries must subsequently update the LRU queue. Being a doubly linked list, if two threads update the LRU simultaneously, we can lead to seg-faults.

Let's see what guarantees we provide for each of the above scenarios:

- xhash operations: We provide a lock for each hash table
- Cache node claiming: The free node queue is protected by a fast lock
- Entry migration: We always take fist the lock for entries and then for rm\_entries
- Reference counting: Another important guarantee is the reference counting of entrys. xcache uses atomic gets and puts to update the reference count of an entry.
- LRU updates: Since all LRU operations take place for entries in "entries" hash table and LRU updates are blazing fast we can secure our LRU with the cache->lock.



### 5.2.5 Re-insertion

We have previously mentioned that in xcache, there can be data migration between hash tables. This is easy to see why in case of evictions: an entry that previously was in "entries" must now be migrated to "rm\_entries" until its reference count falls to zero and can be freed.

However, what happens when xcache receives a request for an evicted entry?

there is a concept called "re-insertion". In order for an entry to be re-inserted to the primary hash table (which will be called "entries" from now on) it must first reside in the hash table that indexes the evicted cache entries (which will be called "rm\_entries" from now on). As mentioned above, an entry that is in rm\_entries has probably pending jobs that delay its removal.

So, what happens if a lookup arrives for that entry while on this stage? In this case, we re-insert it to entries and increase its refcount by 2, since there is one reference by the hash table and one reference by the one who requested the lookup.

### 5.2.6 Event hooks

We have mentioned that xcache is not a standalone program but aims to provide core caching functionalities for other peers to use. Moreover, an important design feature is to not bother the top peer with the xcache internals, such as refcounting or which hash table its resides.

There is an issue that arise with this bla.

Different peers have different things to cache. For example, cached needs to cache objects but xcache was originally created for another peer that needed to cache file descriptors. So, xcache may pose no restriction to what is cached but that also means that it has no control over what happens to the peer's contents on insertion, eviction, removal etc. This is fair, but shouldn't the peer be alarmed for events such as eviction of an entry?

To this end, we have crafted event hooks in xcache onto which the peer can plug its own function and take brief control of the event. For example, when cached is on writeback mode, it plugs a special function on the entry eviction hook to check if this entry is dirty.

### 5.2.7 xcache flow

Below we will see three important scenarios

#### **Insertion**

Figure

#### **Lookup**

Figure

#### **Put**

Figure

## 5.3 The xworkq xtype

The xworkq xtype is a useful abstraction for concurrency control. It's purpose is to enqueue "jobs" and ensure that only one thread will execute them at a time. There is no distinction as to which thread this will be, as well as no execution condition. The winner is simply the one that acquires first the lock.

This means that xworkq is generally used when multiple threads want atomic access to the same memory segment. xworkq is also generic by nature, since the "jobs" are simply a set of functions and their input data.

In cached context, every object has a workq. Whenever a new request is accepted/received for an object, it is enqueued in the workq and we are sure that only one thread at a time can have access to the object's data and metadata.

Let's see the design of the xworkq xtype.

It consists of a queue where jobs (e.g. read from block, write to block) are enqueued. The thread that enqueues a job can attempt to execute it to, by acquiring a lock for the workq. If the lock is free, the thread will be able to execute the enqueued job. Also, other threads can enqueue their jobs, so the thread that has the lock can do those too. There is an xworkq for every object.

## 5.4 The xwaitq xtype

The xwaitq xtype bears some similarities to the xworkq xtype. Like xworkq, it is also an abstraction where "jobs" are enqueued and dequeued later on to be executed. Unlike xworkq though, in xwaitq the "jobs" that are being enqueued are considered to be thread-safe and can be executed concurrently by any thread. Moreover, the execution scheme is different. Specifically, jobs are executed only when a specific condition is true. This way, xwaitq provides an abstraction to...

In cached context...

Let's see the design of the xwaitq xtype.

## 5.5 Cached internals

The components that have been discussed in the previous sections provide core functionalities for cached. Also, besides bridging these components, cached also has some internal bla bla that are crucial for cached's operation.

Namely, these are:

1. Bucket pool
2. Per-object request queues
3. Bucket/Object/cio states
4. Write policy

### 5.5.1 Buckets

### 5.5.2 Object states

## 5.6 Bucket pool

There is however a problem when operating solely on object level. Objects have typically 4MB of size. What would happen if a user requested e.g. a 16KB chunk of an object?

In this case, we would need to read and cache the whole object just to reply to the user's request. If the user then requests a chunk from another object, we would have to cache that object too and in the end, we would thrash our cache<sup>2</sup>.

The solution to this is to further divide objects to the next and final logical entity, buckets (typically 4KB of size). Each bucket consists of its data and metadata and cannot be half-empty, or half-allocated. This way, we can also know which parts of the cached object are actually written, or are in the process of being read etc.

The buckets are pre-allocated, which means two things:

1. We don't need to care about memory fragmentation and system call overhead
2. We cannot index single buckets. <FILLME>

### 5.6.1 Per-object peer requests

Reads and writes to objects are practically read/write request from other peers, for which a peer request has been allocated. There are cases though when an object has to allocate its own peer request e.g. due to a flushing of its dirty buckets. Since this must be fast, there are pre-allocated requests hard-coded in the struct of each object which can be used in such cases.

### 5.6.2 Bucket/Object states

In order to know on how to operate on an object/bucket/cio, we must have some sort of book-keeping. The book-keeping we use is to check the state of the above. (arggh, silly)

### 5.6.3 Write policy

The user must define beforehand what is the write policy of cache. There are two options: write-through and write-back. These policies aren't new and have been discussed extensively in chapter ?, but let's see what these policies translate to in cached context.

- In **write-back** mode, cached caches writes, immediately serves the request back and marks the data as dirty. When a read arrives, it either serves the request with the dirty data (read-hit) or forwards the request to the storage peer and caches the answer (read-miss).

This policy is used when we want to improve read and write speed and can sacrifice data safety.

---

<sup>2</sup> cache thrashing occurs when we aggressively cache data that is only used once and effectively leads to a snowball of evictions

- In **write-through** mode, cached forwards writes to blocker, servers the request when blocker replies, caches the data and marks them as valid. When a read arrives, it either serves the request with the valid data (read-hit) or forwards the request to the storage peer and caches the answer (read-miss).

This policy is used when we want to improve read speed and want to make sure that no data will be lost.

These policies are specified once during cached's deployment and cannot be switched on/off later on.

## 5.7 Cached Operation

### 5.7.1 Polling for new requests

We have explained in Section ? that a peer can send an I/O request to another peer by submitting to its port. However, what is are

### 5.7.2 Write-through mode

Here we will see how cached operates in write-through mode.

#### Write

This is the flow for the write path:

#### Read

This is the flow for the read path:

### 5.7.3 Write-back mode

Here we will see how cached operates in write-back mode.

#### Write

This is the flow for the write path:

#### Read

This is the flow for the read path:

## Chapter 6

# Implementation of cached

In the previous chapter, we presented a design overview for cached and its components. In this chapter we will blabla how the above design has been implemented and explain in depth the structures and functions that have been created for this purpose.

More specifically, sections ? - ? provide implementation information for the components of cached, as described in Chapter ?. Next, section ? presents the actual initialization and blabla operations using excerpts from the code.

## 6.1 Implementation of xcache

In this section, we describe how we implemented the design concept of section 5.2. The main xcache structure is the following:

```
1 struct xcache {
2     struct xlock lock;           /* Main xcache lock */
3     uint32_t size;               /* Upper limit of entries */
4     uint32_t nr_nodes;           /* Shadow entries */
5     struct xq free_nodes;        /* Unclaimed (?) entries */
6     xhash_t *entries;            /* Hash-table for valid entries */
7     xhash_t *rm_entries;         /* Hash-table for evicted entries */
8     struct xlock rm_lock;        /* Lock for rm_entries */
9     struct xcache_entry *nodes; /* Data segment */
10    struct xcache_entry *lru;     /* O(1) lru implementation-specific */
11    struct xcache_entry *mru;     /* O(1) lru implementation-specific */
12    struct xcache_ops ops;        /* Hooks */
13    uint32_t flags;               /* Flags */
14    void *priv;                   /* Pointer to peer struct */
15 };
```

**Listing 6.1:** Main xcache struct

Each of the above xcache struct fields serves a design purpose. Let's see which fields help in what:

### 6.1.1 Entry Preallocation

The relevant fields for this purpose can be seen in the following code listing:

```
1 struct xcache {
2     ...
3     uint32_t size;               /* Upper limit of entries */
4     uint32_t nr_nodes;           /* Shadow entries */
5     ...
6 }
```

```

5     struct xq free_nodes;           /* Unclaimed (?) entries */
6     ...
7     struct xcache_entry *nodes; /* Data segment */
8     ...
9 };

```

**Listing 6.2:** xcache struct fields for preallocated entries

and the definition of the xcache entry struct which shows up in xcache struct can be seen below:

```

1 struct xcache_entry {
2     struct xlock lock;           /* Entry lock */
3     volatile uint32_t parallel_puts; /* Concurrency control */
4     volatile uint32_t ref;        /* Reference counter */
5     uint32_t state;              /* Evicted or active state */
6     char name[XSEG_MAX_TARGETLEN + 1]; /* Entry name */
7     xbinheap_handler h;          /* Index in data segment */
8     struct xcache_entry *older;  /* Less(?) recent entry in LRU
9     queue */
10    struct xcache_entry *younger; /* More(?) recent entry in LRU
11    queue */
12    void *priv;                  /* Pointer to data contents */
13 };

```

**Listing 6.3:** xcache entry struct

Let's start by listing what xcache entry consists of. First of all, it must have a name. Since we preallocate the entries and cannot know in runtime their length, we must allocate as much space as possible. The char name[XSEG\_MAX\_TARGETLEN + 1] field, which is 256 characters long, is long enough to hold the target's name. Also, as we have mentioned in Section 5.2.1, xcache must be agnostic of the cache contents. To this end, we use the generic void \*priv field as a pointer to the actual entry content. The rest of the fields will be explained in the following chapters.

Let's continue now with the fields of Listing 6.2. Since we preallocate the entries using malloc, they take up a contiguous space in memory. The start of this space is the where the \*nodes field points to. The free\_nodes field works similarly to the free\_entries field in Section 3.4.1 i.e. it is a stack where indexes to unused nodes are pushed. These indexes will be seen in various code excerpts in this chapter and have a specific name, xcache\_handler<sup>1</sup>.

## 6.1.2 Entry Indexing

The relevant fields for this purpose can be seen in the following code listing:

```

1 struct xcache {
2     ...
3     xhash_t *entries;           /* Hash-table for valid entries */
4     xhash_t *rm_entries;        /* Hash-table for evicted entries */
5     ...
6 };

```

**Listing 6.4:** xcache struct fields for entry indexing

As we have mentioned in Section 5.2.2, we utilize two hash tables, one for the cached entries and one for the former cached entries (or evicted entries or removed entries). These hash tables can be accessed from the xcache struct and are \*entries and \*rm\_entries respectively.

<sup>1</sup> #define xcache\_handler uint64\_t

These are the functions which are related to indexing and xcache exposes to the peer function:

```

1 xcache_handler xcache_lookup(struct xcache *cache, char *name);
2 xcache_handler xcache_insert(struct xcache *cache, xcache_handler h);
3 int xcache_remove(struct xcache *cache, xcache_handler h);
4 int xcache_invalidate(struct xcache *cache, char *name);

```

**Listing 6.5:** Indexing functions

All of these function need a pointer to the xcache struct. Here's a brief description of them:

**xcache\_lookup:** Takes the target's name as an argument and searches for it in cache.

Returns on failure: NoEntry<sup>2</sup>

Returns on success: the requested handler.

**Note:** Looks only in entries.

**xcache\_insert:** Takes the handler of an allocated entry as an argument and uses it to index that entry.

Returns on failure: NoEntry. Returns on success:  $\alpha$ () the same handler or  $\beta$ () another one, if this entry already exists in cache. **Note:** It looks up first if the entry exist in entries or rm\_entries. The later case can lead to re-insertions.

**xcache\_remove:** Takes the handler of an allocated entry as an argument and uses it to remove that entry.

Returns on failure: -1. Returns on success: 0. **Note:** Removes entries only from entries hash table.

**xcache\_invalidate:** An xcache\_remove spin-off. Takes the name of the entry as an argument, looks it up and then removes it Returns on failure: -1. Returns on success: 0. **Note:** Unlike remove, entries can either be on entries or rm\_entries hash table.

### 6.1.3 Entry eviction

The relevant fields for this purpose can be seen in the following code listing:

```

1 struct xcache {
2     ...
3     struct xq free_nodes;           /* Unclaimed (?) entries */
4     xhash_t *entries;               /* Hash-table for valid entries */
5     xhash_t *rm_entries;            /* Hash-table for evicted entries */
6     struct xlock rm_lock;           /* Lock for rm_entries */
7     ...
8     struct xcache_entry *lru;       /* 0(1) lru implementation-specific */
9     struct xcache_entry *mru;       /* 0(1) lru implementation-specific */
10    ...
11 };

```

**Listing 6.6:** xcache struct fields for eviction

As we have mentioned in Section ??, we resort to eviction when the cache is full and new entries can't be inserted. This entry is the Least Recently Used entry. The doubly-linked list we maintain for this end can be seen below:

<sup>2</sup> #define NoEntry (xcache\_handler)-1

```

1 struct xcache {
2     ...
3     struct xcache_entry *lru;    /* O(1) lru implementation-specific */
4     struct xcache_entry *mru;    /* O(1) lru implementation-specific */
5     ...
6 };
7
8 struct xcache_entry {
9     ...
10    struct xcache_entry *older;    /* Less(?) recent entry in LRU
queue */
11    struct xcache_entry *younger;  /* More(?) recent entry in LRU
queue */
12    ...
13 };

```

**Listing 6.7:** Doubly-linked LRU list

Lets explain these fields a bit:

**lru:** Obviously, it's the least recently used entry. It can be considered as the one end of the doubly linked list.

**mru:** The entry that has just been used. It can be considered as the other end of the doubly-linked list

**younger:** This entry-specific field points to an entry used right after our entry was used.

**older:** Same as "younger", it points to the entry that has been used right before our entry was used.

Finally, as we have explained in Section ??, the eviction internals should normally not bother the user. However, if the user wants to, xcache provides the exposes the following functions:

**xcache\_evict\_lru:** The name says it all, it evicts the recently used item.

**xcache\_peek\_and\_get\_lru:** This function allows the user to atomically take a peek on the Least Recently Used entry and also update its refcount.

#### 6.1.4 Concurrency control

##### Locking

**Reference counting** The refcount model in xcache should be familiar to most people:

- When an entry is inserted in cache, the cache holds a reference for it (ref = 1).
- Whenever a new lookup for this cache entry succeeds, the reference is increased by 1 (ref++)
- When the request that has issued the lookup has finished with an entry, the reference is decreased by 1. (ref-)
- When a cache entry is evicted by cache, the its ref is decreased by 1. (ref-)

Some common refcount cases are:



Case	Refcount
active entry with pending jobs	ref > 1
active entry with no pending jobs	ref = 1
evicted entry with pending jobs	ref > 0
evicted entry with no pending jobs	ref = 0

**Table 6.1:** Reference counting of xcache

- active entry with pending jobs (ref > 1)
- active entry with no pending jobs (ref = 1)
- evicted entry with pending jobs (ref > 0)
- evicted entry with no pending jobs (ref = 0)

and, as always, the entry is freed only when its ref = 0.

### 6.1.5 Event hooks

The hooks that xcache provides to users are:

- on\_init: called on cache entry initialization.
- on\_put: called when the last reference to the cache entry is put
- on\_evict: called when a cache entry is evicted.
- on\_node\_init: called on initial node preparation.
- post\_evict: called after an eviction has occurred, with cache lock held.
- on\_free: called when a cache entry is freed.
- on\_finalize: called to hint the user that the cache entry's ref has dropped to zero.
- on\_reinsert: called when a cache entry has been in cache

## 6.2 Implementation of cached

The cached is the following:

```

1 struct cached {
2     struct xcache *cache;
3     uint64_t total_size; /* Total cache size (bytes) */
4     uint64_t max_objects; /* Max number of objects (plain) */
5     uint64_t max_req_size; /* Max request size to blocker (bytes) */
6     uint32_t object_size; /* Max object size (bytes) */
7     uint32_t bucket_size; /* Bucket size (bytes) */
8     uint32_t buckets_per_object; /* Max buckets per object (plain) */
9     xport bportno;
10    int write_policy;
11    struct xworkq workq;

```

```

12     struct xwaitq pending_waitq;
13     struct xwaitq bucket_waitq;
14     struct xwaitq req_waitq;
15     unsigned char *bucket_data;
16     struct xq bucket_indexes;
17     struct cached_stats stats;
18     //scheduler
19 };

```

**Listing 6.8:** Main cached struct

and the cached entries are the following

```

1 struct ce {
2     uint32_t status;           /* cache entry status */
3     uint32_t *bucket_alloc_status_counters;
4     uint32_t *bucket_data_status_counters;
5     struct bucket *buckets;
6     struct xlock lock;        /* cache entry lock */
7     struct xworkq workq;      /* workq of the cache entry */
8     struct xwaitq pending_waitq;
9     struct peer_req pr;
10 };

```

**Listing 6.9:** Cahed entry struct

## 6.3 Bucket pool

And the bucket implementation is the following:

```

1 struct bucket {
2     unsigned char *data;
3     uint32_t flags;
4 };

```

**Listing 6.10:** Bucket implementation

### 6.3.1 Bucket/Object states

Every object has a state, which is set atomically by threads. The state list is the following:

- **READY:** the object is ready to be used
- **FLUSHING:** the object is flushing its dirty buckets
- **DELETING:** there is a delete request that has been sent to the blocker for this object
- **INVALIDATED:** the object has been deleted
- **FAILED:** something went very wrong with this object

Also, object buckets have their own states too. These are divided in allocation states:

- **FREE:** the bucket has not been touched or has been flushed

- CLAIMED: the bucket has been claimed

and data states:

- INVALID: the same as empty
- LOADING: there is a pending read to blocker for this bucket
- VALID: the bucket is clean and can be read
- DIRTY: the bucket can be read but its contents have not been written to the underlying storage
- WRITING: there is a pending write to blocker for this bucket

Also, for cios we have the following:

- FAILED: at least one of the cio request(s) has failed
- ACCEPTED: normal mode
- READING: at least one of the cio request(s) is pending a read
- WRITING: at least one of the cio request(s) is pending a write
- SERVED: all requests have been served

Finally, for every object there are bucket state counters, which are increased/ decreased when a bucket state is changed. These counters give us an  $O(1)$  glimpse to the bucket states of an object.



## Chapter 7

### Performance evaluation of cached

*"There are three kinds of lies:  
lies, damned lies,  
and statistics benchmarks."*  
Mark Twain (modernized)

It may seem as an ironic statement, considering that we are about to provide benchmark results for cached, but it's actually is a valid one. What Mr. Twain tries to say here <sup>1</sup> is that the presentation of partials facts for something can be used to fabricate a plausible truth for it. In science's case, it so often happens that promising results for an experiment can seem more important to the researcher's eye than negative ones due to positive reinforcement.

In our case, we will try not to merely smear the next pages with diagrams but first explain the benchmarking methodology behind them.

The skeleton of this chapter is the following: Section 7.1 explains the methodology behind our measurements. Section ?? presents the results of the benchmarks that we have done and provides in-depth explanations about each of them. Finally, Section ? is undefined.

#### 7.1 Benchmark methodology

The benchmarks that have been executed and whose results are presented in this chapter, will be split in two categories, both of which have their own distinct goals:

The first category is the comparison between using cached on top of the sosd (sosd has been discussed here ?) and using solely the sosd as the Archipelago storage. The category's goal is to "defend" one of the core thesis arguments, that tiering is a key element that will improve the performance of Archipelago.

In order to compare effectively the performance of cached and sosd, we must consider the following:

1. The comparison of the two peers should try to focus on what is the best performance that these peers can achieve for a series of tough workloads.
2. The circumstances under which both peers will be tested need not be thorough but challenging. For example, it may be interesting to test both peers against sequential requests, but i) such patterns are rarely a nuisance for production environments ii) they do not stress the peers enough to provide something conclusive iii) they are out of the scope of this section as there can be many of them and adding them all here will impede the document's readability.

---

<sup>1</sup> and that's a phrase usually not heard in programming contexts...

3. Both peers must be tested under the same, reasonable workload, i.e a workload that can be encountered in production environments.
4. If the peer doesn't show a consistent behavior for a workload, it must be depicted in the results.

Having the above in mind, the next step is to choose a suitable workload. This choice though is fairly straight-forward; in production environments, the most troublesome workload is the stampede of small random reads/writes and is usually the most tested scenario.

One may ponder however, how many requests can be considered as a "stampede" or which block size is considered as "small". Of course, there is not only one answer to this question so, we will work with ranges. For our workload, we will use block sizes ranging from 4KB to 64KB and parallel requests ranging from 4 to 16.

The second category deals solely with the inner-workings of cached and its behavior on different operation modes or states. Its aim is not to capture the performance against a tough workload, but to explain **why** this performance is observed and how each of the options affect it. For example, we will measure things (blarg?) such as writethrough mode vs writeback mode, single-threaded vs multi-threaded etc.

Also, here is the following list is the options of cached that affect the measurements:

1. Bucket size
2. IOdepth
3. Cache size
4. Max cached objects
5. Write policy
6. ...

Finally, in the following sections for brevity reasons we will talk about comparing cached and sosd. What the reader must keep in mind however is that cached is the cache layer above sosd.

## 7.2 Specifications of test-bed

The specifications of the server on which we conducted our benchmarks is the following.

Component	Description
CPU	2 x Intel(R) Xeon(R) CPU E5645 @ 2.40GHz [e564] Each CPU has six cores with Hyper-Threading enabled, which equals to 24 threads.
RAM	2 banks x 6 DIMMs PC3-10600 Peak transfer rate: 10660 MB/s

**Table 7.1:** dev100 specs

Also, mention that we evaluated both peers by sending requests directly at their request queues

### 7.3 Performance comparison between cached and sosd

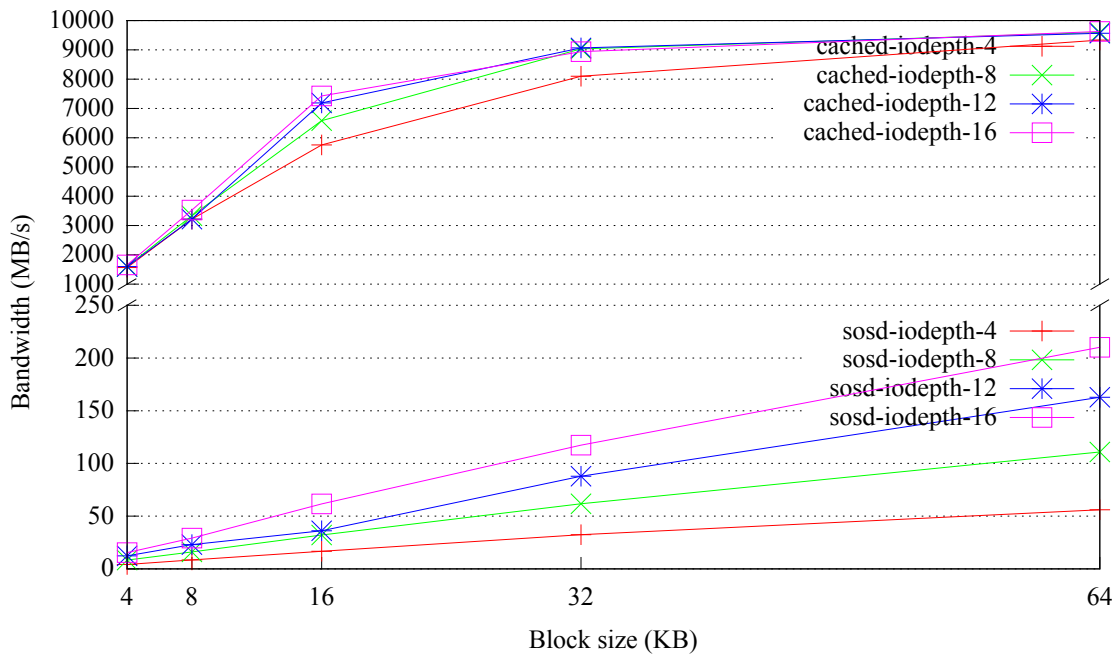
As mentioned above, for our first test, we will evaluate the read and write performance of cached and sosd for a random workload with parallel requests of small size. In order to measure accurately their performance, we will use two different metrics:

**Bandwidth:** We will measure the bandwidth for the bulk of our requests. This is a metric from the peer's perspective that reflects how much request size requests our peers can handle per second.

**Latency:** We will measure the average time the requests need to be served. Unlike bandwidth, this metric reflects the responsiveness of the implementation. For example...

**IOPS:** (Should we include them?)

Let's see now the bandwidth performance of our peers. The write performance can be seen in Figure 7.1 while the read performance can be seen in Figure 7.2.

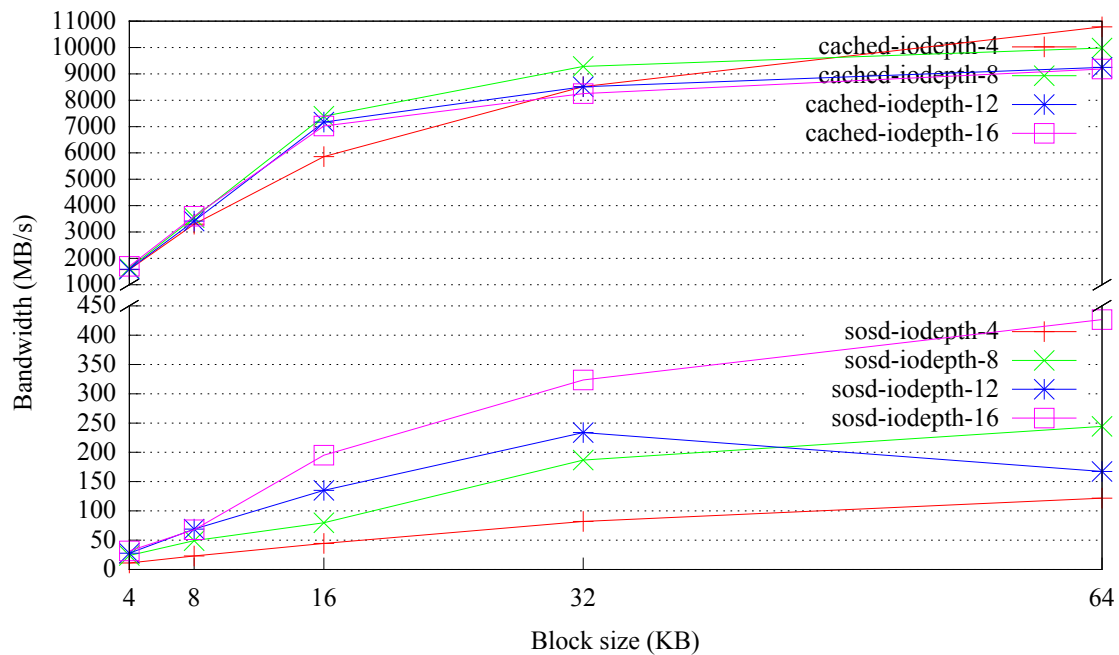


**Figure 7.1:** Comparison of bandwidth performance for writes

Before we proceed with the interpretation of the diagram results, we will briefly comment on the diagram structure. Due to the fact that the performance of the two peers differs in at least two orders of magnitude, the results would look too flat in a conventional diagram that would scale from 0 to 11000. To amend this, we have broken the y-axis of our diagrams in two parts with different scales and starting values, in order to make the comparison easier to the eye.

We will begin the diagram interpretation with the bandwidth performance of the two peers. First, let's see the speedups of write and read requests with the addition of cached. For write requests, the speedup for very small block sizes (4KB - 16KB) is approximately 100x whereas for larger ones (32KB - 64KB) it ranges from 50x to 200x. For read requests, the speedup for very small block sizes is approximately 50x, whereas for larger block sizes it ranges between 20x - 75x.

From these first two diagrams we can extract the following points:



**Figure 7.2:** Comparison of bandwidth performance for reads

### Why is there such a vast improvement?

(Explain that it's because we don't write past cache's size)

### Where is the speedup difference between read and writes attributed to?

The speedup difference is attributed to three factors:

1. The performance of cached is almost the same for writes and reads. This is expected behavior as the read and write paths for cached have many common parts (see Figure ? and ?). However, if we go one step further we will see that under closer inspection, the reads seem consistently a bit faster than writes. So, how can these happen if both paths are the same?

This is actually typical RAM behavior. Reads are reportedly faster than writes (find paper) due to the fact that the update of a bit of an SDRAM is slower than the read (arg, it's dumb).

2. Cached doesn't scale much past the 16KB block size. This is an interesting observation with an unexpected answer. It may seem implausible at first, but what happens is that we are actually hitting the bandwidth limit of the server's RAM. You can see in Table 7.1 that the bandwidth limit is 10.7GB/s. This limit is approached asymptotically as the block size increases and the CPU overhead decreases (more about the CPU overhead later on). Once more, the experienced eye may see that in reads we surpass this limit, which is logical given the fact that we have a multi channel RAM setup that reportedly increases marginally the RAM's performance (why reportedly and how marginally? Explain...)
3. On reads, sosd is benefited from the existence of caches in various levels: on OSD level, on RAID controller level and thus is faster.

### Why is cached's performance increased along with block size?



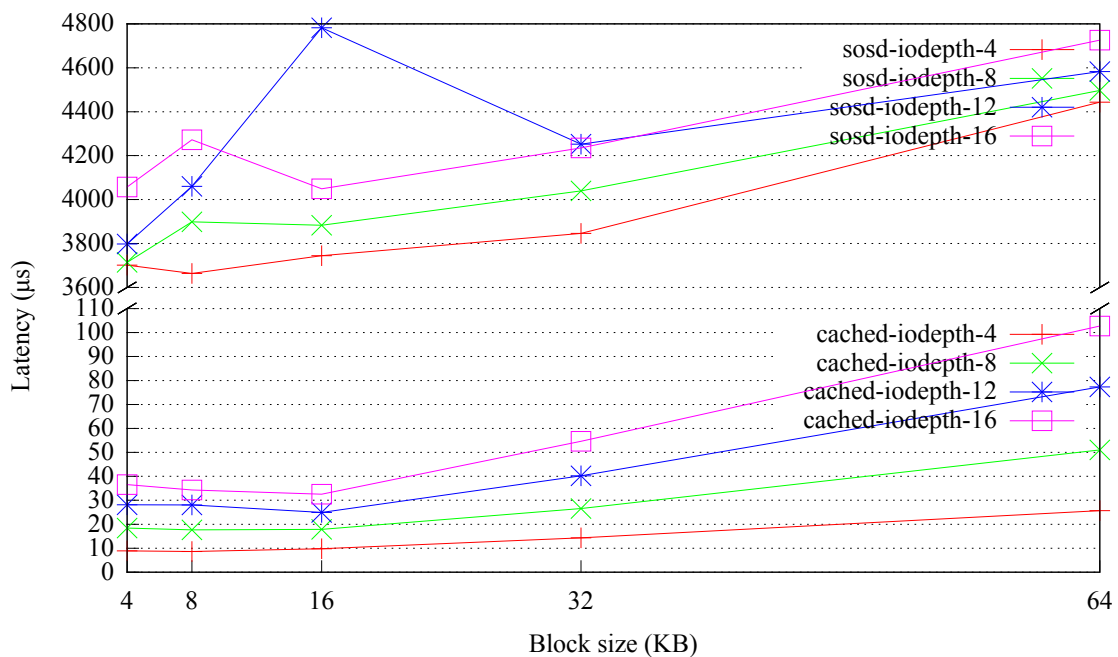
This is another interesting observation but first, why doesn't it apply on sosd? This is because sosd's primary storage are hard disk drives, which have a major drawback; their seek time is not constant and is affected by the location of the contents in the disk platters. Cached however stores data in RAM and we would expect that writing 16 x 4k blocks and 1 x 64k block to take approximately the same time.

From this observation, we can extract that the indexing-related stuff (job enqueueing, lock spinning, hash table indexing) dominate the cached's performance. We can make sure this is the case if the latency results are in microseconds instead of nanoseconds, which is typical for RAM.

### Why isn't the performance of cached improved proportionally as the parallel requests increase?

The reason why we see a minor increase in the performance of cached, even though it's multi-threaded, is because our locking scheme is not fine-grained enough. We have a single lock for our request queue, a single lock for most of the hash table accesses and this inevitably causes a lot of threads to spin. This slight improvement we see is mainly due to the fact that requests are effectively being pipelined while waiting each other to release locks. (explain better)

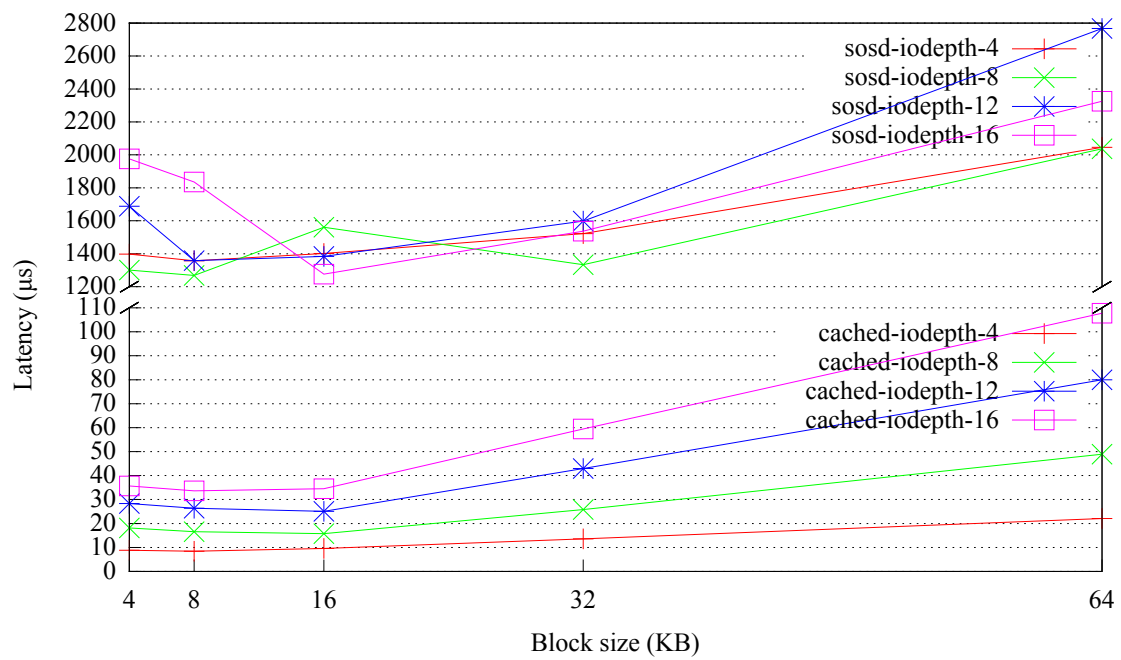
Let's proceed now to the latency results. The write performance can be seen in Figure 7.3 while the read performance can be seen in Figure 7.4.



**Figure 7.3:** Comparison of latency performance for writes

We will measure the speedup of the write and read performance of cached. It is bla-bla for writes and bla-bla for reads.

We can also see that the latency results confirm our previous observations. The latency is increased proportionally with the iodepth, which indicates that we need a more fine-grained locking scheme. Also, latency results are in the order of microseconds instead of nanoseconds which further supports the assertion that the results are dominated by index-related stuff (argh, think of something prettier)



**Figure 7.4:** Comparison of latency performance for reads

## Bibliography

- [aeal99] Some author et al., “Name of citation”, in Proceedings of the 99th ACM Symposium on Something (POPL’99), pp. 999–999, Nine, 9999.
- [Bela66] L.A. Belady, “A study of replacement algorithms for a virtual-storage computer”, IBM Systems Journal, vol. 5, no. 2, pp. 78 – 101, 1966.
- [e564] “Intel(R) Xeon(R) CPU E5645 Specifications”, [http://ark.intel.com/products/48768/Intel-Xeon-Processor-E5645-12M-Cache-2\\_40-GHz-5\\_86-GTs-Intel-QPI](http://ark.intel.com/products/48768/Intel-Xeon-Processor-E5645-12M-Cache-2_40-GHz-5_86-GTs-Intel-QPI).