

# Adaptive memory management framework for derivative clouds

## Master's Thesis Report

A thesis  
submitted in partial fulfillment of the  
requirements for the degree of

## Master of Technology

by

**Prashanth**

Roll No: 153050095

under the guidance of

**Prof. Purushottam Kulkarni**



Department of Computer Science and Engineering  
Indian Institute of Technology, Bombay  
Mumbai

# Abstract

Cloud computing has emerged as one of the hot topics in the computing community today. Most servers these days are either already running on cloud, or are in the virtue of shifting base to cloud. Cloud providers traditionally multiplex a set of compute resources, to group of isolated clients using hardware level virtualization techniques that make use of Virtual Machines (VM) to deploy isolated Virtual Environments (VE).

Although VMs provide a very effective methodology in provisioning compute over the cloud, they incur heavy overheads there by degrading efficiency while provisioning. Lately, there has been a new direction in the flow of research in virtualization, i.e OS-level virtualization in which compute resources in a system are virtualized at an OS-level to provision light weight isolated VEs called Containers. Containers provide similar features to that as VMs but incur much lesser overheads [1] [2] . A recent work [3] has also tried to take it a step ahead, by provisioning compute resources to clients using an nested approach in which repackages and resells resources purchased from native Infrastructure-as-a-service (IaaS) cloud provider. This approach is coined as derivative cloud.

In this work, we have made an initial attempt to understand memory management between containers. We started off with purposing hypotheses based on theoretical evidences. We performed analysis to verify the correctness of our purposed hypotheses and understand parts of memory management for which hypotheses couldn't be drawn. We then tried to extrapolate its implications on real world applications running inside a derivative cloud environment running VMs on the host machine and containers in the guest machine. These implications strongly suggested that existing memory management techniques may impact higher provisioned containers negatively. We conclude by purposing the requirements of a new desired policy that provides this notion of a differentiated reclamation to enforce deterministic allocation when the system is under memory pressure. The end goal of our work is to provide an adaptive deterministic resource provisioning framework for container based services.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Memory management in clouds . . . . .	1
1.1.1	Issues in native container environment . . . . .	1
1.1.2	Amplification of issues in derivative cloud environment . . . . .	1
1.2	Caching in the cloud . . . . .	1
1.2.1	Drawbacks of caching in native (VM) cloud setups . . . . .	1
1.2.2	Hypervisor managed caching . . . . .	1
1.2.3	Issues of caching frameworks in derivative clouds . . . . .	1
1.2.3.1	Lack of framework support in derivative clouds . . . . .	1
1.2.3.2	Dual layers of isolated control . . . . .	1
1.2.4	Application cache sensitivity is unaccounted . . . . .	2
1.3	Problem description . . . . .	2
1.3.1	Phase-1 . . . . .	2
1.3.2	Phase-2 . . . . .	2
<b>2</b>	<b>Background</b>	<b>3</b>
2.1	Memory management between processes in Linux . . . . .	3
2.1.1	Memory pages used by a process . . . . .	3
2.1.2	Memory allocation . . . . .	3
2.1.3	Memory reclamation without container support . . . . .	3
2.2	Containers . . . . .	4
2.2.1	Control groups (Cgroups) . . . . .	4
2.2.1.1	Memory Cgroup . . . . .	4
2.3	Caching . . . . .	4
2.3.1	Hypervisor managed caching . . . . .	4
2.3.1.1	T-MEM cache . . . . .	4
2.3.2	Multilevel caches . . . . .	4
2.3.3	Application specific cache partitoning . . . . .	4
2.3.3.1	MRC construction . . . . .	4
2.3.4	Double decker: Second chance cache for derivataive clouds . . . . .	4
<b>3</b>	<b>Differentiated memory management controller for containers</b>	<b>5</b>
3.1	Drawbacks of existing memory management for containers . . . . .	6

3.1.1	Issues in native environment . . . . .	6
3.1.1.1	Reclamation above soft limits . . . . .	6
3.1.1.2	Reclamation below soft limits . . . . .	6
3.1.2	Amplification of issue in derivative clouds . . . . .	6
3.2	Requirements for a new memory management controller . . . . .	6
3.3	Proposed memory management controller . . . . .	6
3.3.1	Controller architecture . . . . .	6
3.3.2	Policies supported by the controller . . . . .	6
3.4	Modifications made to Linux memory Cgroup . . . . .	6
3.4.1	Per container configurable weights . . . . .	6
3.4.2	Deterministic reclamation . . . . .	6
3.4.3	Flexible reclamation size . . . . .	6
3.5	Empirical evaluation of our controller . . . . .	6
3.5.1	Effectiveness of our controller . . . . .	6
3.5.2	Differential QOS containers . . . . .	6
3.5.3	Impact of reclamation chunk size . . . . .	6
<b>4</b>	<b>Double decker: A memory management framework for derivative clouds</b>	<b>7</b>
4.1	Application cache sensitivity . . . . .	8
4.1.1	Provisioning of caches at different levels based on application sensitivity . . . . .	8
4.1.2	Inability of cache partitioning framework to support anonymous memory applications	8
4.2	Rethink of existing design . . . . .	8
4.2.1	Decentralized memory management framework . . . . .	8
4.2.1.1	Native provider cache partitioning framework . . . . .	8
4.2.1.2	Derivative provider memory management framework . . . . .	8
4.2.2	Hybrid cache . . . . .	8
4.2.2.1	Multilevel configurable caches . . . . .	8
4.2.2.2	Movement of cache objects . . . . .	8
4.3	Implementation details . . . . .	8
4.3.1	Existing implementation status . . . . .	8
4.3.2	Hybrid cache . . . . .	8
4.3.2.1	Pools to accommodate both memory and SSD objects . . . . .	8
4.3.2.2	Asynchronous kernel threads for movement of objects . . . . .	8
4.3.2.3	Multilevel stats . . . . .	8
4.4	Correctness of implementation . . . . .	9
4.4.1	Experimental setup . . . . .	9
4.4.2	Arithmetic validation of stats . . . . .	10
4.4.3	Movement of objects between both levels of cache . . . . .	11
4.4.3.1	Memory to SSD cache . . . . .	11
4.4.3.2	SSD to memory cache . . . . .	12
4.5	Evaluation of Double Decker . . . . .	13

4.5.1	Experimental setup . . . . .	13
4.5.2	Provisioning for anonymous and file backed workloads . . . . .	13
4.5.3	Hybrid cache provisioning . . . . .	13
<b>5</b>	<b>Conclusions</b>	<b>14</b>
<b>6</b>	<b>Future Extensions</b>	<b>15</b>

# List of Tables

4.1	Comparison between expected and actual values . . . . .	11
-----	---	----

# List of Figures

4.1	Experimental testbed for checking correctness . . . . .	10
-----	---	----

# Introduction

Mention about derivative clouds in intro itself

## Memory management in clouds

Issues in native container environment

Amplification of issues in derivative cloud environment

## Caching in the cloud

Drawbacks of caching in native (VM) cloud setups

Hypervisor managed caching

Issues of caching frameworks in derivative clouds

Lack of framework support in derivative clouds

Dual layers of isolated control

Derivative provider has no control over cache partitioning    Cache-level sensitivity

Native provider has no control over application memory allocations    Anonymous memory sensitivity



Application cache sensitivity is unaccounted

## Problem description

Phase-1

Phase-2

# Background

## Memory management between processes in Linux

Memory is allocated/deallocated in terms of pages in any operating system. Memory management in Linux is done using techniques like virtual memory, demand paging, swapping caching etc. They separate between the memory needed by a process and the memory physically allocated on the RAM. The OS creates a large virtual address space for each process. In this section we focus on how memory is managed between processes or a group of processes. We mainly focus on how memory is assigned and reclaimed between them.

### Memory pages used by a process

Memory used by processes are divided into 2 types of pages

1. Anonymous Pages: Pages those which are not associated with any files on disk. They are process memory pages.
2. Page cache pages: Are an in-memory representation of a part files on the disks.
3. Mapped pages: File page with VA mappings

### Memory allocation

When the process needs memory to be allocated, Linux decides the how this memory is going to be allocated physically on the RAM. The process/ application does not see in physical RAM addresses. It only sees virtual addresses from the virtual space assigned to each process. The OS uses a page file located on the disk to assist with memory requests in addition to the RAM. Less RAM means more pressure on the Page file. When the OS tries to find a piece of memory that's not in the RAM, it will try to find in the page file, and in this case they call it a page miss. The actual physical memory allocated (RSS) to a process depends on how much free memory is available in the system. On free memory becoming freshly available in the system, the OS tries to equally distribute the available memory to all processes that are demanding for more memory.

### Memory reclamation without container support

Two lists

## Containers

Control groups (Cgroups)

Memory Cgroup

Memory reclamation with Cgroups

## Caching

Hypervisor managed caching

T-MEM cache

Multilevel caches

Application specific cache partitioning

MRC construction

Double decker: Second chance cache for derivative clouds

# Differentiated memory management controller for containers

[Overview here](#)

## Drawbacks of existing memory management for containers

Issues in native environment

Reclamation above soft limits

Reclamation below soft limits

Amplification of issue in derivative clouds

## Requirements for a new memory management controller

### Proposed memory management controller

Controller architecture

Policies supported by the controller

## Modifications made to Linux memory Cgroup

Per container configurable weights

Deterministic reclamation

Flexible reclamation size

## Empirical evaluation of our controller

Effectiveness of our controller

Differential QOS containers

Impact of reclamation chunk size



# Double decker: A memory management framework for derivative clouds

## Application cache sensitivity

Provisioning of caches at different levels based on application sensitivity

Inability of cache partitioning framework to support anonymous memory applications

## Rethink of existing design

Decentralized memory management framework

Native provider cache partitioning framework

Derivative provider memory management framework

Hybrid cache

Multilevel configurable caches

Movement of cache objects

## Implementation details

Existing implementation status

Hybrid cache

Pools to accommodate both memory and SSD objects

Asynchronous kernel threads for movement of objects

Multilevel stats

# Correctness of implementation

## Experimental setup

The following section describes the experimental setup used to verify the correctness of our implementation.

### Experimental configurations

The set of configurations used for an analysis of memory management framework for a derivative environment must be relevant, and easy to apply. The following configurations fit this criteria, and have been used for the evaluation.

- **Memory Requirement:** Memory requirement of each container, the estimated total memory used by a container.
- **Container memory limit:** Size of memory allocated to a container at the Cgroup level (soft and hard limits).
- **Memory cache limit:** Size of memory (L1) cache assigned to a container.
- **SSD cache limit:** Size of SSD (L2) cache assigned to a container.
- **Workload:** Workload application that is running inside each of the container.
- **Number of containers:** Number of containers that are currently executing in the system.
- **Number of VMs:** Number of virtual machines that are currently executing in the system.

For the sake of simplicity in the evaluations of correctness of our setup. We have only considered a single container, single VM setup which makes use of synthetic workload to stress our system.

### Metrics of interest

The following are the metrics of interest that would help us establish the correctness of our implementation.

- **Container memory usage:** Guest memory usage of the container.
- **Memory cache usage:** Memory cache used by the container.
- **SSD cache usage:** SSD cache used by the container.
- **Demoted:** Objects moved from memory to SSD cache.
- **Promoted:** Objects moved from SSD to memory cache.

The following metrics are collected both for memory and SSD cache

- **Puts:** Number of objects successfully put into this container cache.
- **Gets:** Number of objects successfully got from this container cache.
- **Flushes:** Number of objects flushed from this container cache.
- **Evicts:** Number of objects evicted from this container cache.



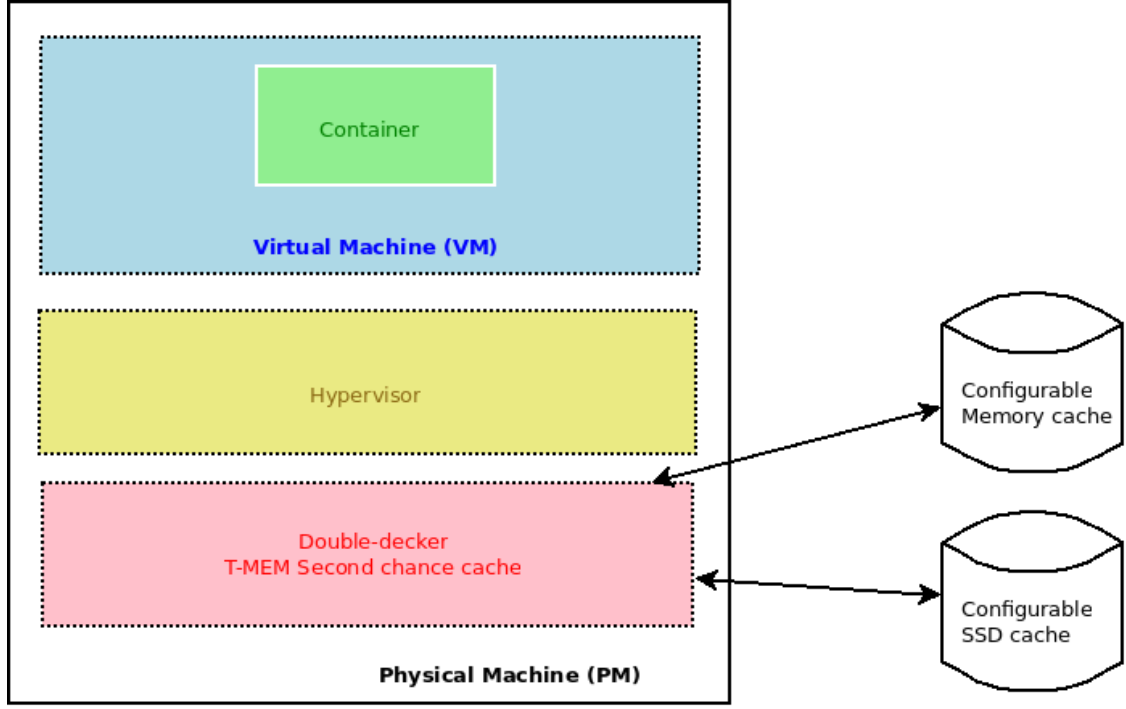


Figure 4.1: Experimental testbed for checking correctness

### Workload

For establishing the correctness of our workload, we have considered a self generated workload generated using `cat` command that outputs the content of a file onto `/dev/null`.

### Testbed

Our testbed consists of a single VM, single container running on top of our hybrid implementation of Double decker as shown in Fig 4.1. The hypervisor used is KVM, and the container manager used is LXC.

The physical machine configuration used is as described below,

1. Intel(R) Core(TM) i7-3770 CPU @ 3.40GHz
2. 4 CPU cores (with multi-threading)
3. 8 GB of physical RAM
4. 120 GB SSD disk

### Arithmetic validation of stats

#### Question

To verify the correctness in accounting of stats while accessing cache at both levels.

#### Procedure

We ran several experiments and computed the actual cache usage (memory and SSD) in our implementation. To an calculated an estimated cache usage we used the formula given below. We used the same formula for both memory and SSD cache.

$$EstimatedUsed = Puts + ObjectsMovedIn - (Gets + Flushes + ObjectsMovedOut) \quad (4.1)$$

Metric	Approx. estimated Value	Observed Value
Puts (MB)	2078	2080
Gets (MB)	0	1
Container memory usage (MB)	512	509
Memory cache usage (MB)	504	503
SSD cache usage (MB)	1008	1000
Evicts (MB)	54	64
Flushes (MB)	0	0
Cumulative usage (MB)	2078	2076
Demotions (MB)	1062	1064

Table 4.1: Comparison between expected and actual values

### Observations

The values for *EstimatedUsed* and *ActualUsed* (present in our stat counter) matched in most cases. However, over long periods of run, with quite a large number of cache operations there was a marginal difference between the two ( $<1\%$ ).

### Inference

Since the correctness of the actual value of cache used depends on the other stats that we have collected, and the matching of *EstimatedUsed* and *ActualUsed* would only mean that all the stats collected are right.

## Movement of objects between both levels of cache

To verify the correctness of our implementation empirically, we have taken our synthetic workload described in section 4.4.1 and ran a couple of simple experiments to demonstrate the expected behavior of our cache to support cache operations like puts, gets, promotions and demotions.

### Memory to SSD cache

#### Question

To verify the correctness in accounting of stats while accessing cache and moving objects from memory (L1) to SSD (L2) cache.

#### Procedure

We start of the experiment with powering on the VM, followed by the container. We assigned complete memory and SSD cache at the double-decker back-end to support the container. The container had a **memory requirement of 2078 MB**, 2048 MB workload requirement and 30 MB container requirement (container requirement was obtained by running the same experiment while having a nearly 0 MB workload). The container was allocated with 2048 MB (512 MB of container memory + 512 MB of memory cache + 1024 MB of SSD cache). Now, the workload performed a sequential read of its workload (i.e 2048 MB) once. Table 4.1 shows the list of approx. estimated values (which are based on our implementation) and observed values for the metrics at the cache at the end of the experiment.

## Observations

The following are the observations,

1. There is a small number of gets, probably occurring due to pages used by other container applications.
2. Puts in the cache, is marginally (<2 MB) greater than expected value, and this deviation is due to the small number of Gets which are occurring.
3. The memory cache usage, is exactly as expected. However, SSD cache usage is slightly lesser than the expected value, but however the deviation seems to be an acceptable value.
4. The demotions (movement from memory to SSD) and cumulative usage values are nearly the same with a subtle deviation (<2 MB) which is an acceptable value.

## Inference

The accounting stats, are nearly as expected. This verifies the correctness of most of the stats of our implementation (except promotions).

## SSD to memory cache

### Question

To verify the correctness in accounting of stats while moving objects from SSD (L2) to memory (L1) cache.

### Procedure

We start of the experiment with powering on the VM, followed by the container. We assigned complete memory and SSD cache at the double-decker back-end to support the container. The container had a **memory requirement of 2078 MB**, 2048 MB workload requirement and 30 MB container requirement (container requirement was obtained by running the same experiment while having a nearly 0 MB workload. The container was allocated with 2560 MB (512 MB of container memory + 2048 MB of SSD cache). The workload performed a sequential read of its workload (i.e 2048 MB) once, this lead to using up of nearly 1566 MB of SSD cache.

Now, we changed the memory cache size to 256 MB while performing basic operations at the container which triggered the promotion (movement of objects from SSD to memory) of the objects to the memory cache. The promotion triggers all objects until the memory cache reaches a threshold usage - 192 MB in our case, as this threshold is calculated as,

$$MemoryCacheLowerThreshold(192MB) = MemoryLimit(256MB) - LimitSize(64MB) \quad (4.2)$$

Hence we would expect 192 MB worth of objects be promoted from SSD to memory cache in an ideal case.

### Observation

Using our stats, it was observed that the estimated promotion and the actual promotion of objects were of an exact match with the number being 192 MB.

### Inference

This verifies the correctness in the accounting stats in movement of objects from SSD to memory cache.

Hence movement of objects from both memory to SSD and vice-versa have been verified empirically.

## Evaluation of Double Decker

### Experimental setup

#### Experimental configurations

#### Metrics of interest

#### Workload

#### Testbed

#### Provisioning for anonymous and file backed workloads

#### Hybrid cache provisioning

# Conclusions

We have made an initial attempt to understand memory management in Linux containers. We started off with purposing hypotheses based on theoretical evidences. We performed empirical analysis to verify the correctness of our purposed hypotheses. We also performed a few more empirical analysis to establish parts of memory management for which hypotheses couldn't be drawn. We then tried to extrapolate its implications in the real world applications running inside a derivative cloud environment. These implications strongly suggested that existing memory management techniques may impact higher provisioned containers negatively, when the system is under memory pressure. We conclude by purposing the requirements of a new desired policy that provides this notion of a differentiated reclamation to enforce deterministic allocation when the system is under memory pressure.

# Future Extensions

The following are the list of works that are to be taken up in the near future,

1. Design and implement a new memory management policy for containers.
2. Analyze memory hierarchy in cgroups, and see how this affects containers.
3. Explore other resource controller in the container framework, identify issues and provide appropriate fixes.
4. The end goal is to provide an adaptive resource provisioning framework for containers.

# Bibliography

- [1] W. Felter, A. Ferreira, R. Rajamony, and J. Rubio, “An updated performance comparison of virtual machines and linux containers,” in *Performance Analysis of Systems and Software (ISPASS), 2015 IEEE International Symposium On*, pp. 171–172, IEEE, 2015.
- [2] R. Morabito, J. Kjallman, and M. Komu, “Hypervisors vs. lightweight virtualization: a performance comparison,” in *Cloud Engineering (IC2E), 2015 IEEE International Conference on*, pp. 386–393, IEEE, 2015.
- [3] P. Sharma, S. Lee, T. Guo, D. Irwin, and P. Shenoy, “Spotcheck: Designing a derivative iaas cloud on the spot market,” in *Proceedings of the Tenth European Conference on Computer Systems*, p. 16, ACM, 2015.