

國立清華大學

資訊工程研究所

碩士論文

容器在Swarm叢集中的遷移

Checkpoint and restore containers in Docker Swarm



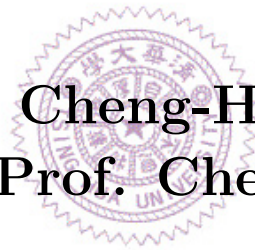
研究生：103062622 黃晟豪 (Cheng-Hao Huang)

指導教授：李哲榮 教授 (Prof. Che-Rung Lee)

中華民國一零五年七月

Design and Implementation of Container's Checkpoint-and-Restore in Docker Swarm

Student: Cheng-Hao Huang
Advisor: Prof. Che-Rung Lee



Department of Computer Science
National Tsing Hua University
Hsinchu, Taiwan, 30013, R.O.C.

July 2016

中文摘要

容器技術自 2013 年 Docker 發表後在全世界迅速竄紅，Container 解決了維護人員在伺服器進行大量部屬時的痛點，使得環境部屬只需要建立完容器映像檔後就可以進行大量部屬，並對每一個容器環境進行隔離。

在這篇論文中，我們提出了在 Docker swarm 叢集中，將容器在多個節點中相互搬移。另外，可以針對特定的容器定期設定 checkpoint 儲存至雲端儲存空間，若叢集中的節點遇到不正常的離線時，可以及時回復最近的容器狀態到健康的節點上。

Abstract



Contents

Chinese Abstract	i
Abstract	i
Contents	iv
List of Figures	vi
List of Tables	vii
List of Algorithms	viii
1 Introduction	1
2 Background	2
2.1 Docker	2
2.1.1 Docker Client	2
2.1.2 Docker Daemon	3
2.1.3 runC	3
2.2 Docker Swarm	3
2.2.1 Discovery services	4
2.2.2 Scheduler	4
2.2.3 High availability of Swarm Manager	6
2.2.4 High availability of Swarm containers	6
2.3 CRIU	6
3 Design and Implementation	7
3.1 Docker	7
3.1.1 Docker Client	7



3.1.2	Docker Daemon	8
3.2	Docker Swarm configuration	8
3.3	Docker containers migration in Docker Swarm	8
3.4	Docker Swarm checkpoint restore rescheduling policy	10
3.4.1	Docker Swarm container checkpoint tickers	10
3.4.2	Docker Swarm restore rescheduling policy	11
3.4.3	High availability of Swarm Manager in Docker Swarm check- point restore rescheduling policy	12
4	Experiments	13
5	Related Work	14
6	Conclusion	15



List of Figures

2.1	Single node Docker	2
2.2	Docker Swarm architecture	4
3.1	Docker Swarm with remote storage server	9
3.2	Containers checkpoint version in remote storage server	11



List of Tables



List of Algorithms



Chapter 1

Introduction

Traditionally, people rely single supercomputer to calculate data or deploy their companies's application. Nowadays, cloud computing has been use wildly. People would like to use a lot of normal computers that gather them into a cluster to replace a supercomputer. More and more companies build their data centres or use cloud platforms to construct their business application.

Virtualization is a popular technology that is used wildly on cloud computing, including virtual operating system, computer hardware platforms, storage device and computer network resources. Virtual machine is a technology to emulate a particular computer like a real computer. It can partition a physical machine resource, such as CPU, memory, storage, and network. A hypervisor uses execution to manage and share host physical machine hardware, it allows many different virtual machines isolated from each other.

Container is operating system level virtualization which runs as an isolated process in userspace on the host operating system and shares the same operation system kernel with other containers. It provides kernel namespace such as PID,IPC, network, mount, and user namespace to isolate each container environment to host operation system. In order to control hardware resource like CPU, memory, network, and disk I/O, container uses cgroups to limit each container resources. Container doesn't have hypervisor to isolate with host operating system, therefore, it can offer better performance than virtual machine.

Checkpoint and restoration container migration

High availability

Chapter 2

Background

2.1 Docker

Docker [2] is an open-source project container engine. It provides an additional layer of abstraction and automation of operating-system-level virtualization in Linux. Besides, Docker has extra image management and layered file system to reduce disk space. Docker has two parts, including Docker Client and Docker Daemon.

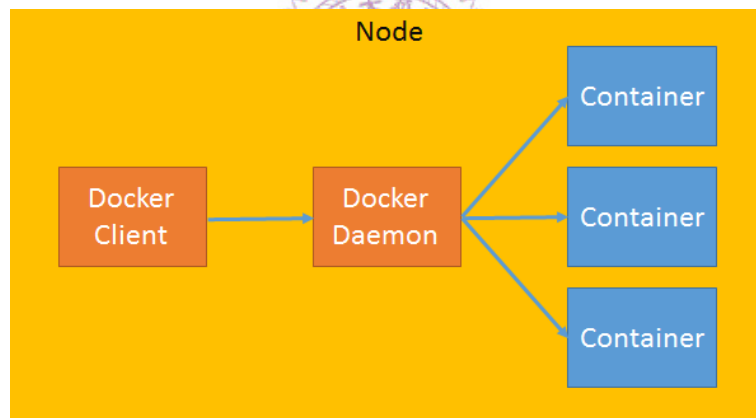


Figure 2.1: Single node Docker

2.1.1 Docker Client

Docker is a typical client/server architecture application. Docker Client is a CLI (Client Line Interface) in Docker. Docker uses Docker Client to send and receive requests to Docker Daemon. Also, Docker supports remote RESTful API to

send and receive HTTP requests to Docker Daemon. In addition, it has been implemented by more than 10 programming languages.

2.1.2 Docker Daemon

Docker Daemon is a daemon that runs as system service. It has two the most importance features:

- Receive and handle Docker Client's requests.
- Manage containers.

When Docker Daemon is running, it will run a server that receives requests from Docker Clients or remote RESTful API. After receiving requests, server will pass the requests by router to find the handler to handle the requests. By default, Docker Daemon listens UNIX socket requests, it serves root permission, or docker group membership. Whenever user wants Docker Daemon to listen remote RESTful API or Docker Swarm requests, it has to be enabled the TCP socket.

2.1.3 runC

runC is a CLI tool for running containers according to the OCI(Open Container Initiative) specification. It doesn't need any dependency from the operating system, it can control Linux Kernel include namespace, cgroups, apparmor, netlink, capabilities, firewall, etc. runC provides a standard interface to support the containers management that Docker can use it to control the containers.

2.2 Docker Swarm

Docker Swarm [3] is a native clustering for Docker. It gathers several docker engines together into one virtual docker engine. Docker Swarm serves standard Docker API, so it can be connected by Dokku, Docker Machine, Docker Compose, Jenkins, DockerUI, Drone, etc. It also support Docker Client of course.

In Docker Swarm, it has two components which are Swarm master and Swarm node. Swarm master is the manager which handles Docker Client and RESTful API requests and manages multiple Docker nodes resources. Docker Node is an agent

which sends heartbeat to Discovery Service to ensure Docker Daemon is alive in the cluster.

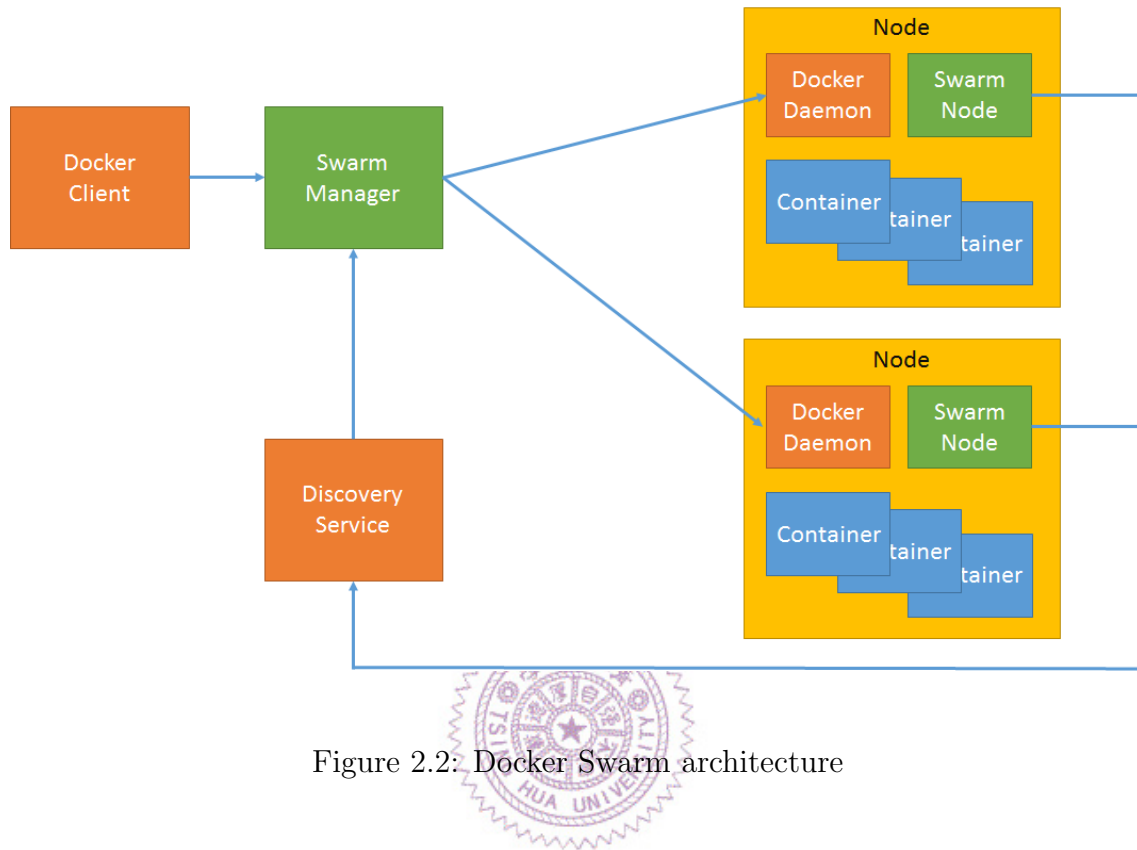


Figure 2.2: Docker Swarm architecture

2.2.1 Discovery services

Docker Swarm provides multiple Discovery Services backends. They are used to discover nodes in the cluster. There are:

- Using a distributed key/value store, like Consul, Etcd and Zookeeper.
- A static file or list of nodes.
- Docker Hub as a hosted discovery service

Otherwise, it also supports any modules which satisfy Discovery API interface.

2.2.2 Scheduler

Docker Swarm scheduler decides which nodes to use when creating and running a container. It has two steps. First, It follows user's filters to decide which nodes are conform. Second, It passes through strategies to select the best node in the cluster.

Filter

Filters are divided into two categories, node filters and container configuration filters. Node filters operate on characteristics of the Docker host or on the configuration of the Docker Daemon. Container configuration filters operate on characteristics of containers, or on the availability of images on a host. The node filters are:

- Constraint
- Container slots
- Health filter

The container configuration filters are:

- Affinity
- Dependency
- Port filter



Strategies

The Docker Swarm scheduler features multiple strategies for ranking nodes. Docker Swarm currently supports these values:

- Spread
- Binpack
- Random

Spread and Binpack strategies compute rank according to a nodes available CPU, its RAM, and the number of containers it has. It selects a node at random. Under the

Spread strategy, Swarm optimizes for the node with the least number of containers. The Binpack strategy causes Swarm to optimize for the node which is most packed. The Random strategy uses no computation, chooses nodes at random regardless of their available CPU or RAM.

2.2.3 High availability of Swarm Manager

In Docker Swarm, Swarm Manager responses the cluster and manages the resources of multiple Docker nodes at scale. If Swarm master dies, we have to create a new one and deal with the interruption of service.

The High availability feature allows Docker Swarm has multiple Swarm Manager instances. We can create a primary manager and multiple replica instances. Whenever we send requests to replica instances, it will be automatically proxied to the primary manager. In addition, if the primary manager fails, the others replica instances will lead a new primary manager.

2.2.4 High availability of Swarm containers

In Docker Swarm, it has rescheduling policy. As we set the reschedule policy when we start a container, whenever Swarm nodes fail, Swarm Manager will restart all of the containers which on the fail nodes to another alive Swarm Nodes.

2.3 CRIU

CRIU [1] (Checkpoint/Restore in Userspace) stands for Checkpoint and Restore in User Space, creates a complete snapshot of the state of a process, including things like memory contents, file descriptors, and even open TCP connections. It can be used for suspending and resuming processes, or live migrating them from one machine to another.

Chapter 3

Design and Implementation

3.1 Docker

In native Docker, It has two part, Docker Client and Docker Daemon. Docker Daemon has many components include server, engine, registry, graph, driver and runC. To support dump checkpoint and restore request, some of these steps should be implemented.

3.1.1 Docker Client

We implement 3 Docker commands in Docker Client, including checkpoint, restore, migrate. In checkpoint command should have these configurations:

- image directory - Dump checkpoint image directory.
- work directory - Dump checkpoint image log directory.
- leave running - After dumping checkpoint image, keeping running container or not.
- pre-dump - Pre-dump checkpoint memory image to minimize frozen time.
- pre image directory - Define which version image to compare.
- track memory - Track memory to pre image directory image to minimize disk space.

In restore command should have these configurations:

- image directory - Checkpoint image directory to restore from.
- work directory - Directory for restore log.
- force - Force restoring container from image directory whether container is running or not.

In migrate command, it focus on Docker Swarm Scheduler filter configurations. In run command, we can do this with the environment variable or the label. Therefore, we implement environment variable and the label configurations in migrate command.

3.1.2 Docker Daemon

In native Docker Daemon, it doesn't support checkpoint and restore command. Fortunately, it is already implemented in runC, so we have to add some code to proxy Docker Client's requests through Docker Daemon to runC.

3.2 Docker Swarm configuration

As Figure 3.1, we prepare a remote storage server for saving Docker containers dump checkpoint images. It should have fault tolerant to avoid service shout down. For these reasons, We choose glusterFS to be our experiments remote storage server.

After setting up glusterFS, we mount it to every Docker nodes in the same folder. To avoid mount it as absolute path in every Docker nodes, we mount it to Docker root which we can change configuration in Docker Daemon.

3.3 Docker containers migration in Docker Swarm

Docker Swarm creates containers through Swarm scheduler to dispatch Docker nodes. If we want to specific assign which node we want to create containers, we have to set filters like constraint, affinity or dependency. To migrate containers in Docker Swarm, we must avoid the containers which we want to migrate that migrate to an other nodes, instead of the same node.

Step 1. Check Docker Swarm cluster has at least two Swarm nodes.

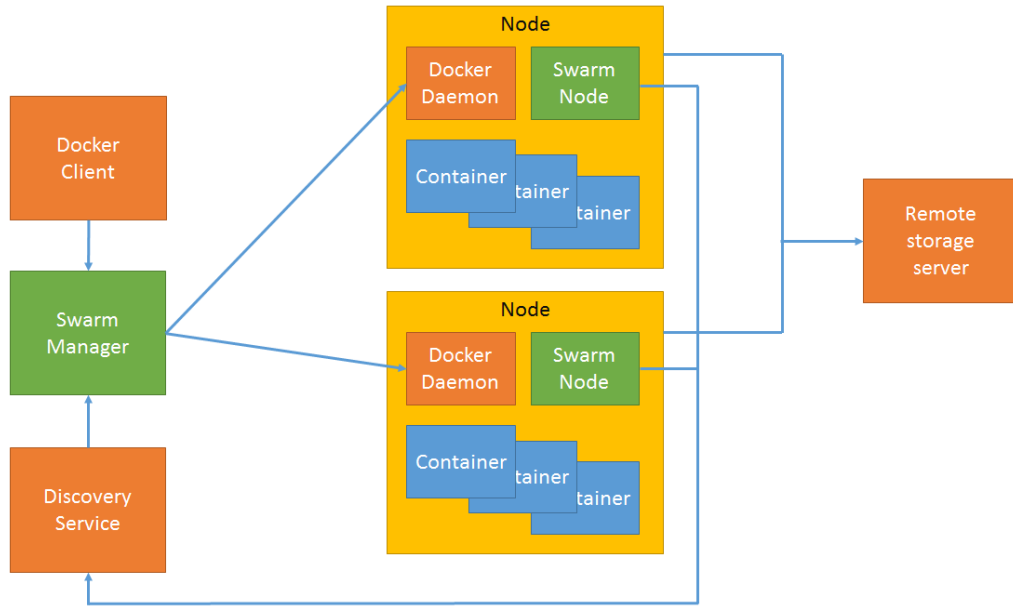


Figure 3.1: Docker Swarm with remote storage server

- Step 2. Parse Docker Client requests to analyse label and environment variables, and transform label and environment variables to Docker Swarm filters.
- Step 3. Add constraint filter to make sure the container which we want to migrate does not migrate to the same node.
- Step 4. Pre-dump the container checkpoint image which we want to migrate to decrease container frozen time.
- Step 5. Dump the container checkpoint image by tracking memory from pre-dump checkpoint image.
- Step 6. Create empty container on the Docker Swarm scheduler chooses node.
- Step 7. Restore the container to the Docker Swarm scheduler chooses node.
- Step 8. Delete the checkpoint images.
- Step 9. If the container was migrated which had setted checkpoint restore rescheduling policy, it will restart checkpoint restore rescheduling policy 3.4.

3.4 Docker Swarm checkpoint restore rescheduling policy

In Docker Swarm, it has rescheduling policy. As we set the reschedule policy when we start a container, whenever Swarm nodes fail, Swarm Manager will restart the containers which on the fail nodes to another alive Swarm Nodes.

We improve this policy that we checkpoint every containers which we want to keep checkpoint for every checkpoint ticker. Whenever Swarm Nodes fail, Swarm Manager will restore the containers which Swarm manager has dumped checkpoint. Otherwise, restore rescheduling policy provides version checkpoint by memory track. It only dump different memory page checkpoint to new version checkpoint.

In addition, it also support high availability that whenever Docker Swarm primary manager fails, the others Swarm Manager replica instances will lead a new primary manager. After replica leading a new primary manager, it will restart container checkpoint tickers.

3.4.1 Docker Swarm container checkpoint tickers

Step 1. Set checkpoint-time label and keep version label when we create the container.

Step 2. Swarm Manager analyzes checkpoint label when the container has created. If keep version label isn't setted, keep version will be setted to 5.

Step 3. After the container starting, starting the container checkpoint ticker. Container checkpoint tickers are for when users want to checkpoint container repeatedly at regular intervals. Container checkpoint tickers will do these steps:

Step a. Swarm Manager sends pre-dump checkpoint request to Swarm Node's Docker Daemon when every Pre-dump version start.

Step b. After pre-dumping the container, Swarm Manager sends dump new version checkpoint request to Swarm Node's Docker Daemon. Every new checkpoint images track memory to last checkpoint version image, just save memory difference in new checkpoint image. We

save keep-version(default 5) versions in the same directory(Figure 3.2).

Step c. Send delete checkpoint request to Swarm Node's Docker Daemon to delete last Pre-dump version directory whenever container has two Pre-dump versions.

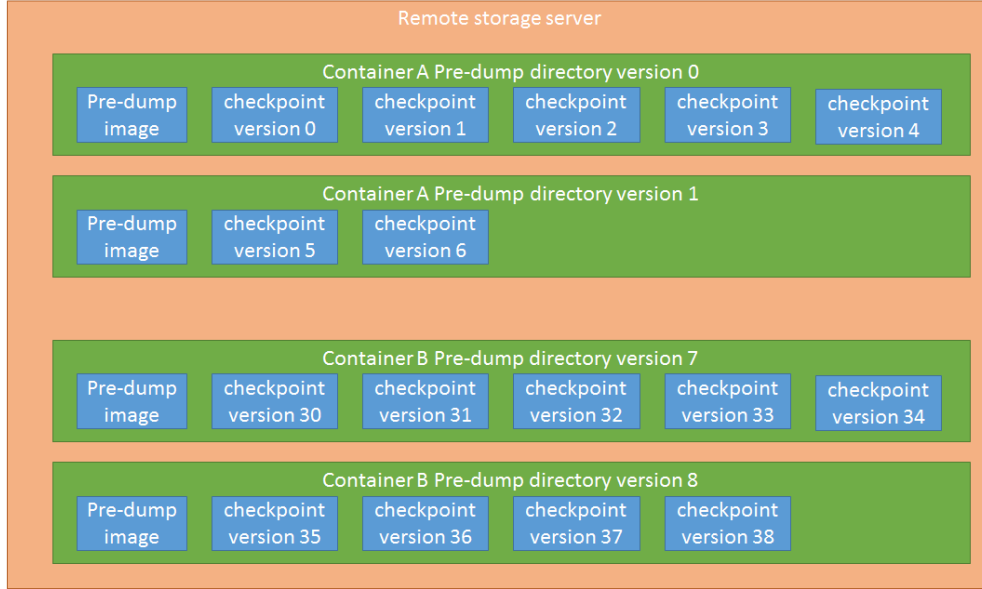


Figure 3.2: Containers checkpoint version in remote storage server

3.4.2 Docker Swarm restore rescheduling policy

Step 1. Set reschedule:restore label when we create the container.

Step 2. Swarm Manager analyzes restore label when the container start.

Step 3. Whenever Swarm Nodes fail, Swarm Manager will restore the containers which Swarm Manager has dumped checkpoint last version to another Swarm nodes.

Step 4. To avoid dumping checkpoint version failing at the same time, if restoring last version failing, Swarm Manager will retry second last version checkpoint to restore. It will retry keep-version(default 5) times.

Step 5. If Swarm Manager retries keep-version are all fail, it will create and start a new container as normal Docker Swarm rescheduling policy.

3.4.3 High availability of Swarm Manager in Docker Swarm checkpoint restore rescheduling policy

Whenever Docker Swarm primary manager fails, the others Swarm Manager replica instances will lead a new primary manager. After replica leading a new primary manager, it searching every Docker Node's containers which has checkpoint restore rescheduling policy's label. If the containers has checkpoint restore rescheduling policy's label, Docker Swarm new primary manager will restart container checkpoint tickers.



Chapter 4

Experiments



Chapter 5

Related Work

[4]



Chapter 6

Conclusion

CRIU provides the container for dumping checkpoint and restoring in the same Docker Daemon. In this thesis, we have purposed it to expand to Docker Swarm in the cluster.



Bibliography

- [1] Criu. https://criu.org/Main_Page.
- [2] Docker. <https://www.docker.com/>.
- [3] Docker swarm. <https://docs.docker.com/swarm/overview/>.
- [4] Michael R Hines and Kartik Gopalan. Post-copy based live virtual machine migration using adaptive pre-paging and dynamic self-ballooning. In *Proceedings of the 2009 ACM SIGPLAN/SIGOPS international conference on Virtual execution environments*, pages 51–60. ACM, 2009.

