

A Far-Reaching Data Set and Experimental Learn About Docker Images Throw on Docker Hub

Prof. Monika Rokade¹ Mr. Pandurang R. Shinde² Prof. Sunil Khatal³

^{1,3}Assistant Professor ²ME Student

^{1,2,3}Department of Computer Engineering

^{1,2,3}Sharadchandra Pawar College of Engineering, Otur, India

Abstract— Docker is a free and open application development and execution platform. Docker separates your applications from your infrastructure, allowing you to swiftly release software. You can manage your infrastructure in the same way that you control your applications using Docker. You may drastically minimise the time between writing code and executing it in production by utilising Docker's approaches for shipping, testing, and deploying code quickly. Docker is one of the most widely used containerization systems. Research looked into many parts of the Docker ecosystem, but it mostly focused on Dockerfiles from GitHub, which limited the types of questions that could be addressed and did not look into evolution. In this study, we aggregate data from Docker Hub, GitHub, and Bitbucket to construct a more recent and complete data set.

Keywords: FROM, WORKDIR, COPY, RUN, EXPOSE, and CMD

I. INTRODUCTION

Docker is a centralised open-source platform for developing, deploying, and running applications. Docker runs apps on the host's operating system using containers. Rather of creating a virtual operating system, it allows programmes to use the same Linux kernel as the host computer. Containers ensure that our app runs in every environment, including development, testing, and production.

Docker consists of several components, including a client, a server, a machine, a hub, and Docker composes. Docker containers are virtual machine lightweight alternatives. It enables developers to package an application along with all of its libraries and dependencies and distribute it as a single file. You don't need to assign any RAM or disc space for the programmes when you use a docker container. It generates storage and space automatically based on the application's needs.

A virtual machine is software that allows us to instal and use multiple operating systems on our machine at the same time (Windows, Linux, and Debian). Virtualized operating systems are the operating systems that run virtual machines. These virtualized operating systems can run applications and carry out functions that a true operating system can. Docker makes it simple to instal and execute software without having to worry about configuration or dependencies. When working on code with coworkers, developers utilise Docker to eliminate machine difficulties, such as "but code is worked on my laptop." Docker is used by operators to execute and manage apps in isolated containers in order to increase computing density. Docker is used by businesses to build secure agile software delivery pipelines that allow them to release new application features faster and more securely. Because Docker is not only utilised for

deployment but also as a development platform, we are able to effectively boost our customers' happiness.

This is a significant feature of Docker that allows us to quickly and easily customise the system. We'll be able to deploy our code in less time and with less work. Because Docker may be utilised in a wide range of contexts, the infrastructure requirements are no longer tied to the application's environment. By facilitating technical configuration and application deployment. It has undoubtedly increased productivity. Docker not only makes it possible to run an application in a controlled environment, but it also saves resources.

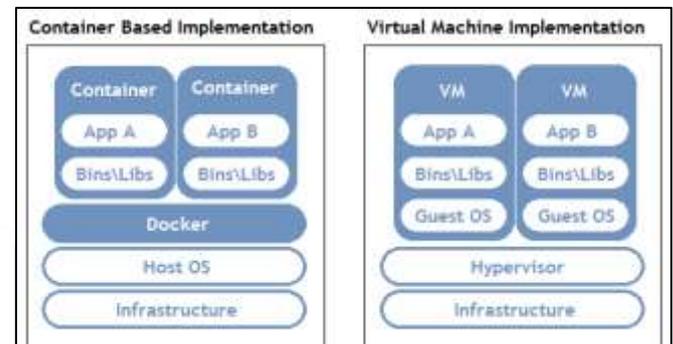


Fig. 1: System Architecture

It provides containers for running applications in a secure environment. Each container is self-contained, allowing us to run any type of programme.

II. BACKGROUND AND TERMINOLOGY

A. Docker:

Docker is one of the most widely used containerization systems today. A Docker image is used to distribute a Docker-encapsulated programme. A Dockerfile is used to define the Docker image, which contains all of the instructions and environment variables needed to run the application. When the image-encapsulated target application is deployed, an instance of the image is produced [14], which is referred to as the Docker container. As a result, a Docker image is similar to a snapshot of the target application and enables for reproducible container creation. One of the benefits of the Docker ecosystem is that creating a new Docker image does not necessitate writing a completely new Dockerfile. Using the FROM command, a Docker image can inherit image definitions from another base image, similar to how inheritance works in object-oriented programming. In this situation, the new image inherits all of the characteristics and files contained in the base image.

Extending and using any Docker image as a base image is possible. A Dockerfile defining a Docker image for executing a Python script is shown in Listing 1. In this case, the image developer just needs to use the FROM

python:3.7.3-stretch statement to identify the base image for the target Python run-time environment. All properties and files contained in the official Python 3.7.3 language run-time image are inherited as a result. The remaining instructions in the Dockerfile are required to execute the application, which will run in Python 3.7.3. Docker builds an image by executing Dockerfile statements and creating a layer for each instruction.

```
# syntax=docker/dockerfile:1
FROM ubuntu:18.04
COPY ./app
RUN make /app
CMD python /app/app.py
```

Fig. 2: .Example of Dockerfile

B. Docker Hub:

A Docker image registry is a location where Docker images can be hosted, indexed, and managed. Docker Hub [23], which debuted in 2014, is the world's largest, free, and public Docker image registry. As of May 3, 2020, it had over 3.4 million public Docker images, and the number is continually growing. Official photos are those published by a few certified firms and organisations, whereas non-official images are those shared by community developers. The majority of Docker images use an official image as their foundation image (for example, Ubuntu). Docker images work with a variety of hardware architectures, including ARM and x86. Docker Hub keeps track of the number of times an image has been retrieved from Docker Hub for each image.

Official images are divided into three categories: operating system, language run-time, and application. OS Docker images are Docker images that merely encapsulate a basic operating system; examples include Ubuntu and Debian. Language run-time images, such as Python and Golang, are Docker images that offer the run-time environment for a specific programming language. Application images, such as Nginx and PostgreSQL, are Docker images that encapsulate a ready-to-use application. Docker Hub is a continuous integration (CI) service that connects the Dockerfile and other Docker image files to a GitHub or Bitbucket source repository. Docker Hub automatically triggers an image build when new code is submitted into the associated source repository. Docker images that use such a service have a link to their GitHub/Bitbucket source on Docker Hub.

On their summary page, they provide extra image information, such as the Dockerfile of the most recently successfully produced image and the build status.

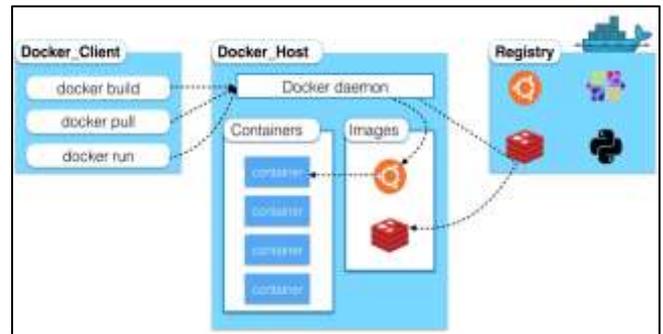


Fig. 3: far-Reaching Overview Docker Architecture

Docker is built on a client-server model. The Docker client communicates with the Docker daemon, which handles the construction, execution, and distribution of your Docker containers. You can execute the Docker client and daemon on the same machine or you can link a Docker client to a Docker daemon that is located elsewhere. A REST API, UNIX sockets, or a network interface is used by the Docker client and daemon to communicate. Docker Compose is another Docker client that allows you to interact with applications made up of many containers.

C. Docker Daemon:

The Docker daemon (dockerd) handles Docker objects such as images, containers, networks, and volumes by listening for Docker API requests. To manage Docker services, a daemon can communicate with other daemons.

D. Docker Client:

Many Docker users interact with Docker primarily through the Docker client (docker). When you run commands like docker run, the client transmits them to dockerd, which executes them. The Docker API is used by the docker command. The Docker client has the ability to communicate with many daemons.

E. Docker Register:

Docker images are stored in a Docker registry. Docker Hub is a public registry that anybody may use, and Docker is set up by default to look for images on it. It's even possible to run your own private registry.

The relevant images are pulled from your configured registry when you use the docker pull or docker run commands. Your image gets pushed to your configured registry when you run the docker push command.

F. Docker Object:

You create and use images, containers, networks, volumes, plugins, and other things when you use Docker. This section provides a quick overview of a few of those items. A Docker image is a read-only template that contains instructions for building a Docker container. A lot of the time, an image is based on another image with some tweaking. For example, you could create an image that is based on the Ubuntu image but also includes the Apache web server and your application, as well as the configuration data required to execute your application. You can either make your own photos or rely on those generated by others and stored in a registry. You create your own image by writing a Dockerfile with a simple syntax for outlining the procedures required to create and operate the image. Each Dockerfile instruction builds a layer in the

image. Only the layers that have changed are rebuilt when you edit the Dockerfile and rebuild the image. When compared to other virtualization technologies, this is part of what makes pictures so light, tiny, and fast.

III. LITERATURE SURVEY

A literature review summarizes the work of many academics who found that Docker is far superior to Virtual Machines.

Virtual machines have no effect on CPU or memory utilization, but they do affect I/O and OS interaction. This overhead manifests itself as additional cycles for each I/O operation. This overhead raises I/O latency and lowers the number of CPU cycles available for productive activity. The throughput is likewise limited as a result of this. One of the reasons for the development and implementation of Docker in a cloud environment is because of this.

Docker has various capabilities that make it easier to use than LXC-style raw containers, such as stacked images and NAT, but these advantages come at a performance penalty. As a result, with Docker's default settings, it's possible that it's no quicker than KVM. Applications that rely on the file system or disc should use volumes instead of AUFS. Using `--net=host` to minimise NAT cost is simple, but it sacrifices the benefits of network namespaces. The paradigm of one IP address per container, according to researchers, can give flexibility and performance.

A. Containerization in the Cloud Environment:

A containerized application can be constructed by assembling individual pictures, possibly based on basic images from the repositories. Through the image layering and extension process, containers can encapsulate a variety of application components. A container can hold a variety of user applications and platform components.

Container granularity, or the amount of applications inside, varies. Some prefer the one-container-per-app approach, which allows for easy creation of new stacks (for example, changing the Web server in an application) and reuse of common components (for example, monitoring tools or a single storage service like "me cached" - either locally or predefined from a repository like the Docker Hub). Apps are simple to create, rebuild, and administer. When compared to multi-app containers, the negative is a larger number of containers with associated interface and administrative costs, however container efficiency should help [4].

IV. PROPOSED SYSTEM

Docker automates boring configuration processes and is utilised across the development lifecycle for fast, easy, and portable desktop and cloud application development. Docker's end-to-end platform comprises user interfaces, command-line interfaces, APIs, and security that are all designed to function together.

Throughout the lifecycle of an application's delivery. The full virtualization method entails virtualizing everything from the hardware to the operating system. And this type of virtualization is a combination of hardware and software that is employed very simply, making it simple to comprehend. Because of the separated OS, this strategy is simple to comprehend. It's also simple to implement their

system, network, user, and security policies, and they support a variety of operating systems. It does, however, have several flaws. In the case of VM distribution, it distributes images, including operating systems, which wastes resources. In comparison to KVM, Docker Cloud is simply deployed and distributed, without the need for a Guest OS. It can be implemented quickly and efficiently while also virtualizing and distributing with the Linux Operating System.

V. CONCLUSION

Docker is the most widely used containerization software at the moment. Given this popularity and the associated impact, it's critical to understand the Docker ecosystem's evolution patterns and quality challenges. We produced a new far-reaching and comprehensive data collection in this study, encompassing information about 3,364,529 Docker images and 378,615 source repositories behind them that were thrown on Docker Hub. We conducted a wide-ranging experimental learn based on this data set about many features of Docker images and their changes over time.

As a result, boot-time, image generation, and distribution times are all short. When compared to VM Cloud, this is a benefit of adopting Docker cloud. VM is operated in such a way that it may be described as providing a user with a new machine. As a result, system, network, user, and security policies are simple to manage and apply. In addition, regardless of the Host OS, the user can use a variety of virtualization operating systems. Docker containers' security can be improved by running them as "non-privileged" and enabling additional hardening options in the Linux kernel, such as AppArmor or SELinux. Following this report, further research could compare Docker container security to that of alternative containerization platforms or virtual machines. Such research could lead to things like a full static analysis of Docker or a larger understanding of container security in general.

REFERENCES

- [1] Changyuan Lin, Sarah Nadi and Hamzeh Khazaee, "A Large-scale Data Set and an Empirical Study of Docker Images Hosted on Docker Hub," <https://docs.docker.com/develop/develop-images/dockerfile-best-practices>, 2020, [Online; accessed May-4-2020].
- [2] N.Kratzke, "A brief history of cloud application architectures," *Applied Sciences*, vol. 8, no. 8, p. 1368, 2018.
- [3] H. Lee, K. Satyam, and G. Fox, "Evaluation of production serverless computing environments," in 2018 IEEE 11th International Conference on Cloud Computing (CLOUD). IEEE, 2018, pp. 442–450.
- [4] "Get started, part 1: Orientation and setup — dockerdokumentation," <https://docs.docker.com/get-started/>, 2020, [Online; accessed May-1-2020].
- [5] "The top programming languages 2019," <https://spectrum.ieee.org/computing/software/the-top-programming-languages-2019>, 2019, [On-line; accessed April-4-2020].
- [6] Dockerfilereference, "https://docs.docker.com/engine/reference/builder/", 2020, [Online; accessed May-9-2020].

- [7] Monika D.Rokade ,Dr.Yogesh kumar Sharma,"Deep and machine learning approaches for anomaly-based intrusion detection of imbalanced network traffic." IOSR Journal of Engineering (IOSR JEN),ISSN (e): 2250-3021, ISSN (p): 2278-8719
- [8] Monika D.Rokade ,Dr.Yogesh kumar Sharma"MLIDS: A Machine Learning Approach for Intrusion Detection for Real Time Network Dataset", 2021 International Conference on Emerging Smart Computing and Informatics (ESCI), IEEE
- [9] Monika D.Rokade, Dr. Yogesh Kumar Sharma. (2020). Identification of Malicious Activity for Network Packet using Deep Learning. International Journal of Advanced Science and Technology, 29(9s), 2324 - 2331.
- [10] Sujata D. Sumbare, Sonali B. Jadhav, Pandurang R. Shinde, Prof. Monika Rokade .(2020)" Depression Detection using Facial features with Machine LearningTechniques: An overview" International Journal of Innovative Research in Computer and Communication Engineering, ISSN(Online): 2320-9801, ISSN (Print): 2320-9798,vol.8, Issue 4, April 2020
- [11] Sunil S.Khatal ,Dr.Yogesh kumar Sharma, "Health Care Patient Monitoring using IoT and Machine Learning.", IOSR Journal of Engineering (IOSR JEN), ISSN (e): 2250-3021, ISSN (p): 2278-8719
- [12] Sunil S.Khatal ,Dr.Yogesh kumar Sharma, "Data Hiding In Audio-Video Using Anti Forensics Technique ForAuthentication ", IJSRDV4I50349, Volume : 4, Issue : 5
- [13] Pandurang R. Shinde, Sonali B. Jadhav, Sujata D. Sumbare, Prof. Monika Rokade .(2020)" Depression Detection using Machine Learning and AI Techniques" International Journal of Innovative Research in Computer and Communication Engineering, ISSN: 2319-8753,vol.9,Issue5,May2020
- [14] Sunil S.Khatal Dr. Yogesh Kumar Sharma. (2020). Analyzing the role of Heart Disease Prediction System using IoT and Machine Learning. International Journal of Advanced Science and Technology, 29(9s), 2340 - 2346.