



**bigdata**  
republic

# Docker for Data Scientists

---

Packaging your models & best practices



# After this webinar, you will know...

---

Why you need Docker

How Docker works

How to use Docker in your workflow



# bigdata republic



Data Scientists



Data Engineers

# Presenters and contributors

---

**Tom de Ruijter**

*Data Scientist*

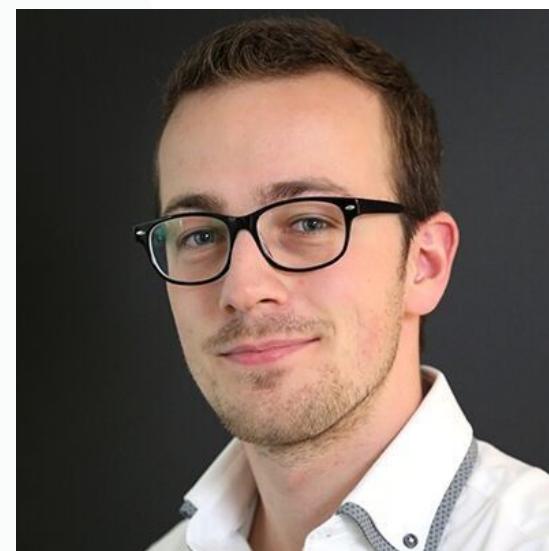
[/in/tomderuijter/](#)



**Joris Bukala**

*Data Scientist*

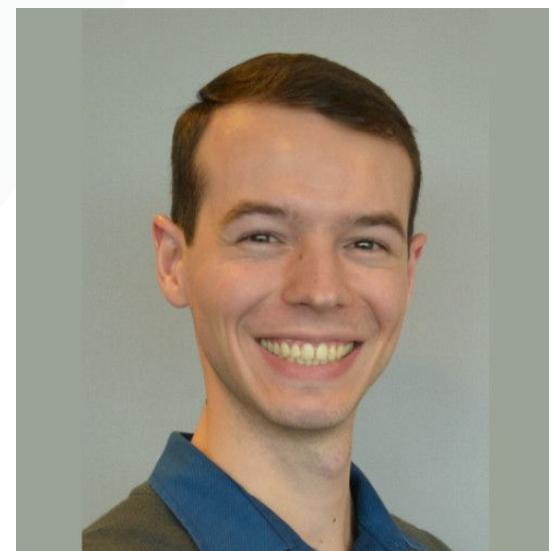
[/in/joris-bukala/](#)



**Jacobus Herman**

*Data Engineer*

[/in/jacobus-herman/](#)



# What does Docker do?

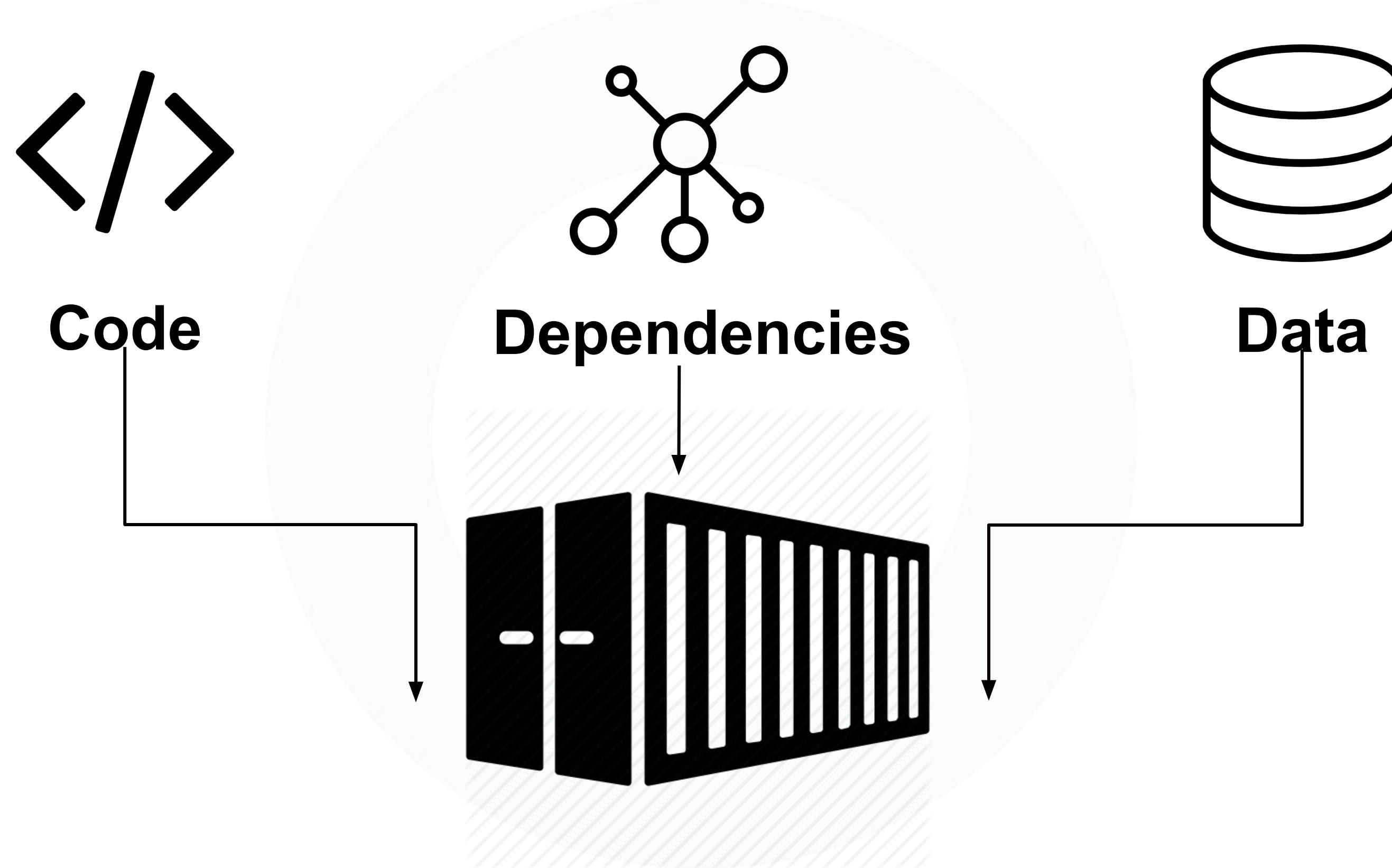
---

Docker allows us to package and run applications in an isolated environment



# Docker Containers

---



# Problems Docker solves

---

# Reproducing models/environments:

---

*Did I need this library in this branch of my code?*



# Reproducing models/environments:

---

*Did I need this library in this branch of my code?*

*What environment variables did I have to set when I started this project?*



# Reproducing models/environments:

---

*Did I need this library in this branch of my code?*

*What environment variables did I have to set when I started this project?*

*Did I use that confidential data when training this model?*



# Sharing your work

---

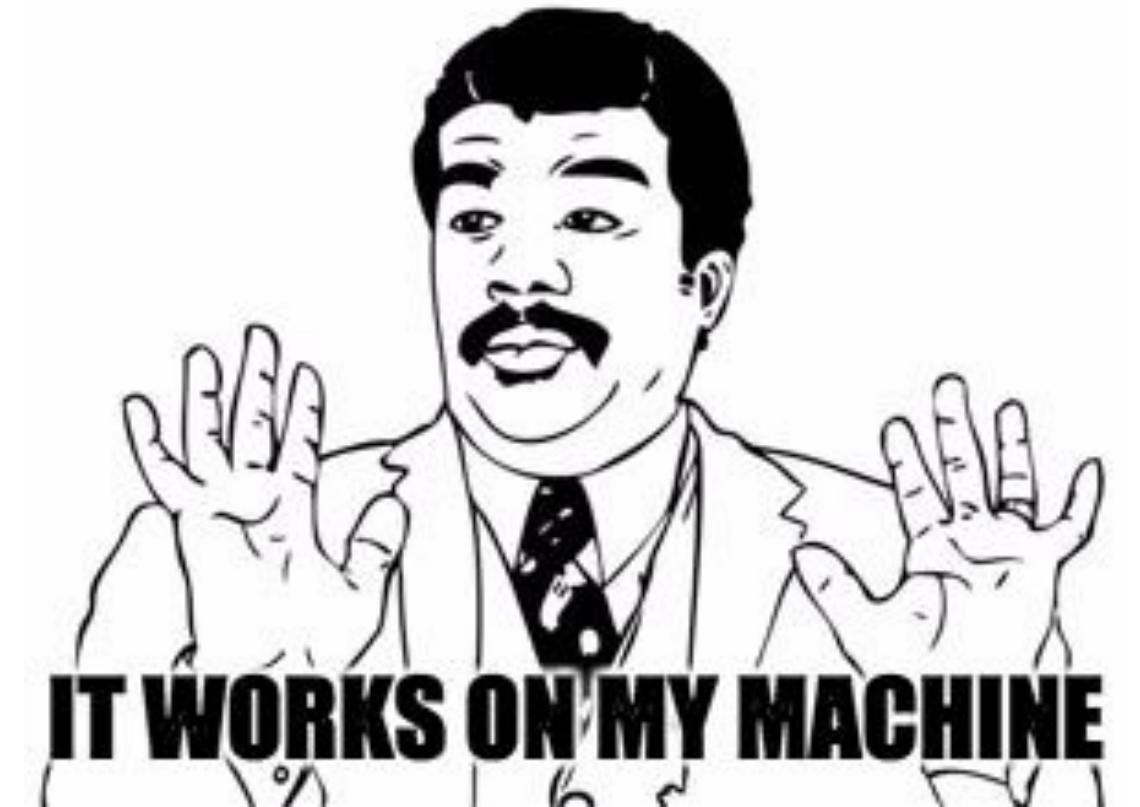
- *Import error: no module named XXXX*



# Sharing your work

---

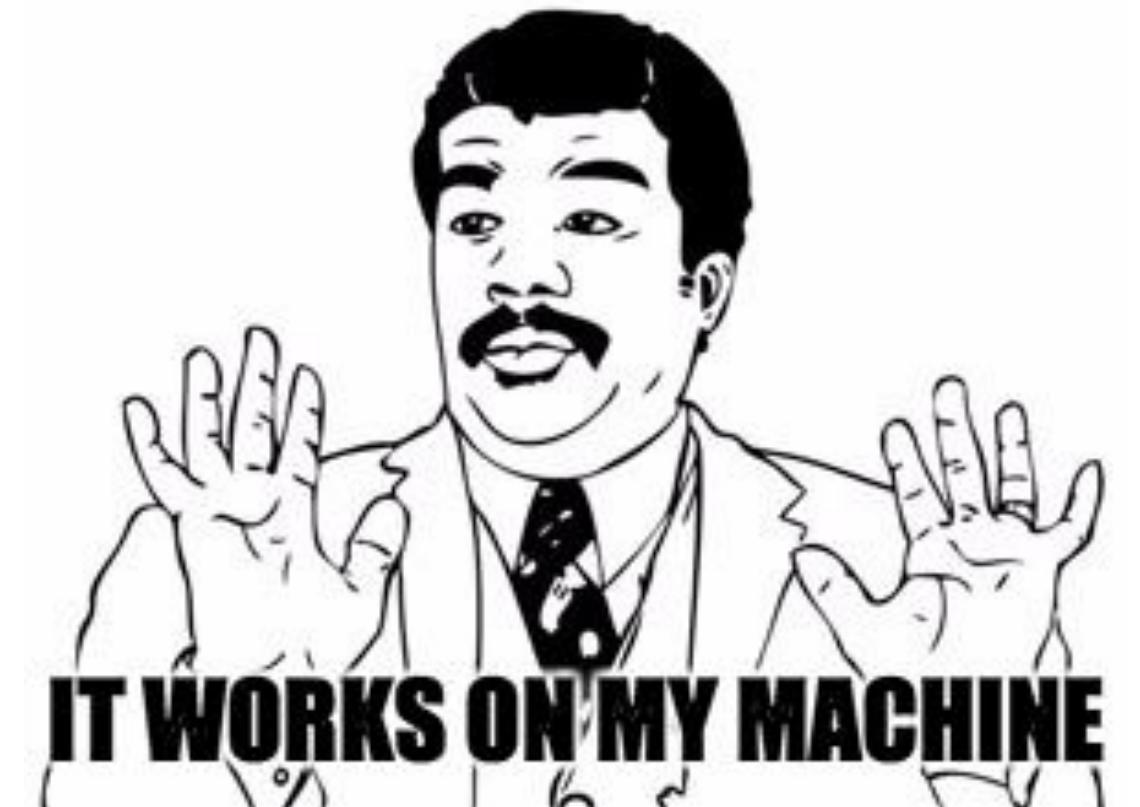
- *Import error: no module named XXXX*
- *“Oh no all my devices only run Arch Linux...”*



# Sharing your work

---

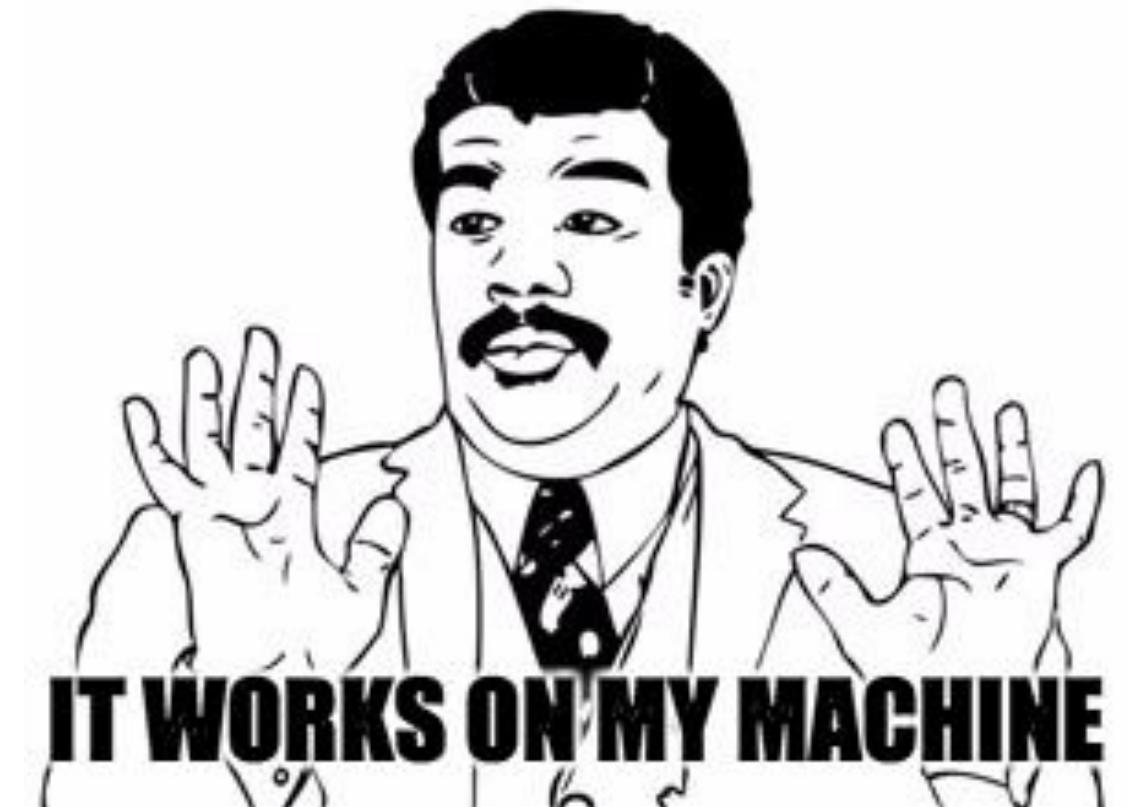
- *Import error: no module named XXXX*
- *“Oh no all my devices only run Arch Linux...”*
- *“Here is the Git repo, now you have to install X, Y, and Z. Oh right I forgot about W. Why don’t you have an up to date COBOL compiler?!”*



# Sharing your work

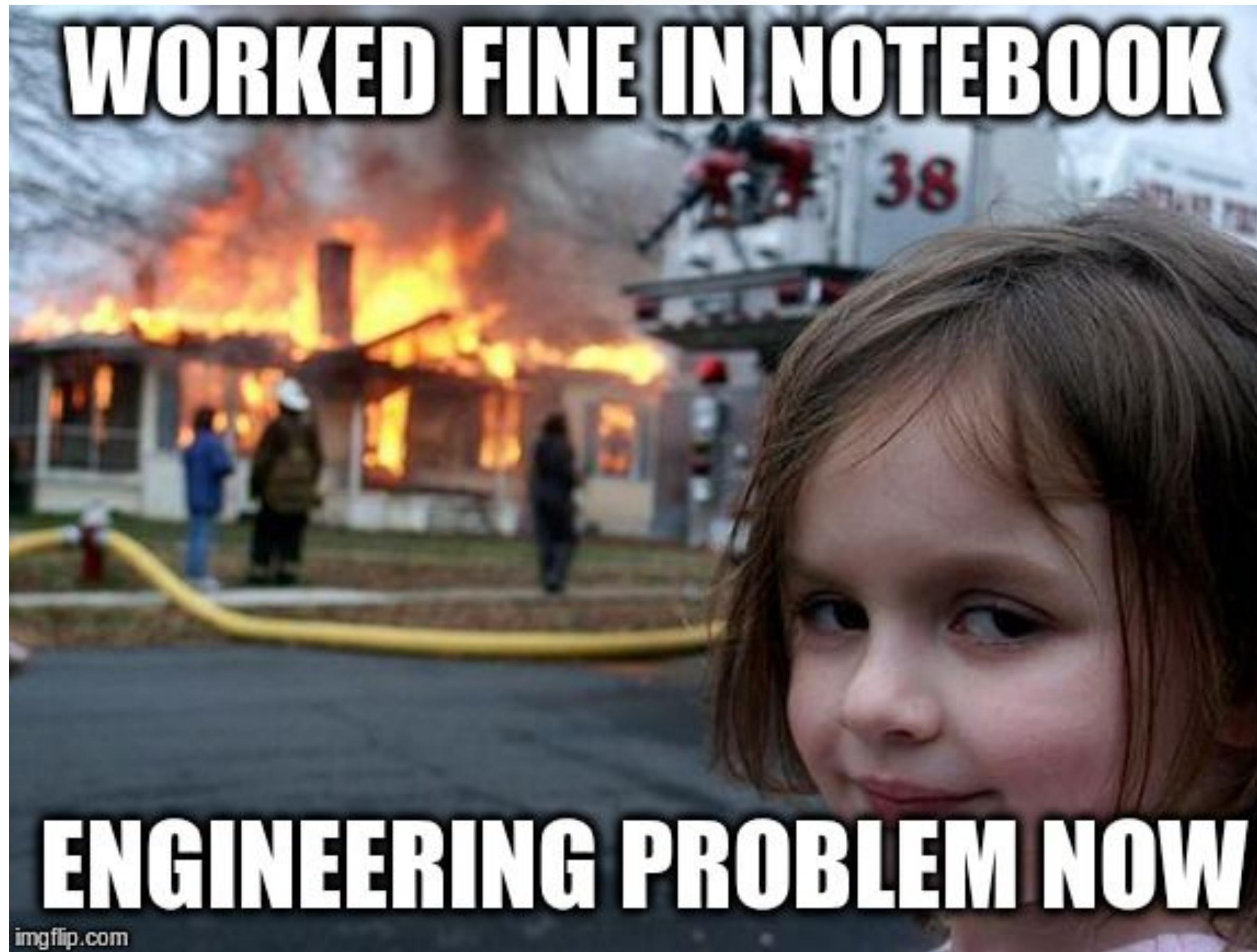
---

- *Import error: no module named XXXX*
  - *“Oh no all my devices only run Arch Linux...”*
  - *“Here is the Git repo, now you have to install X, Y, and Z. Oh right I forgot about W. Why don’t you have an up to date COBOL compiler?!”*
- Share code repo with a small text file (Dockerfile) and be done



# Deployment

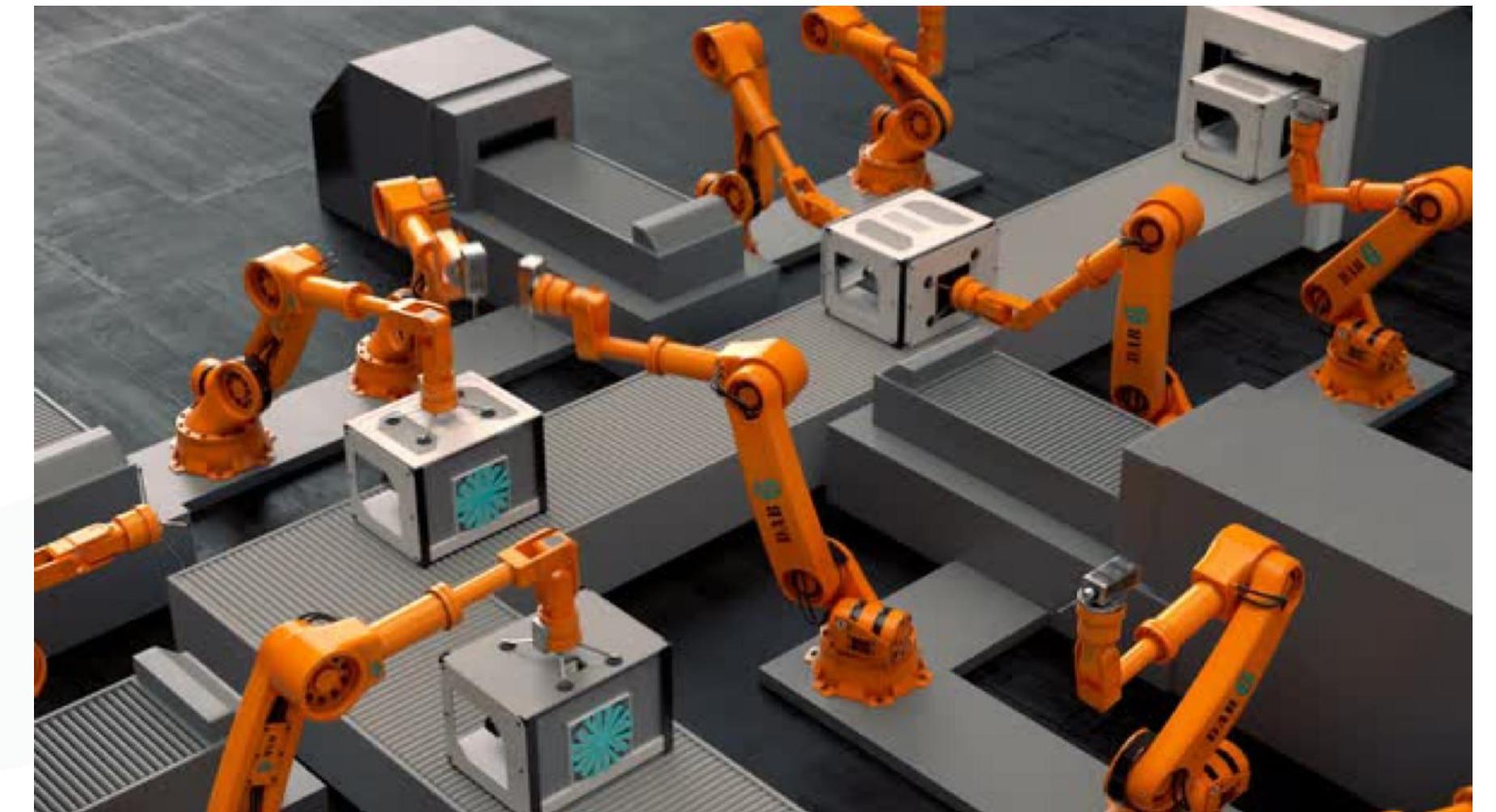
---



# Deployment

---

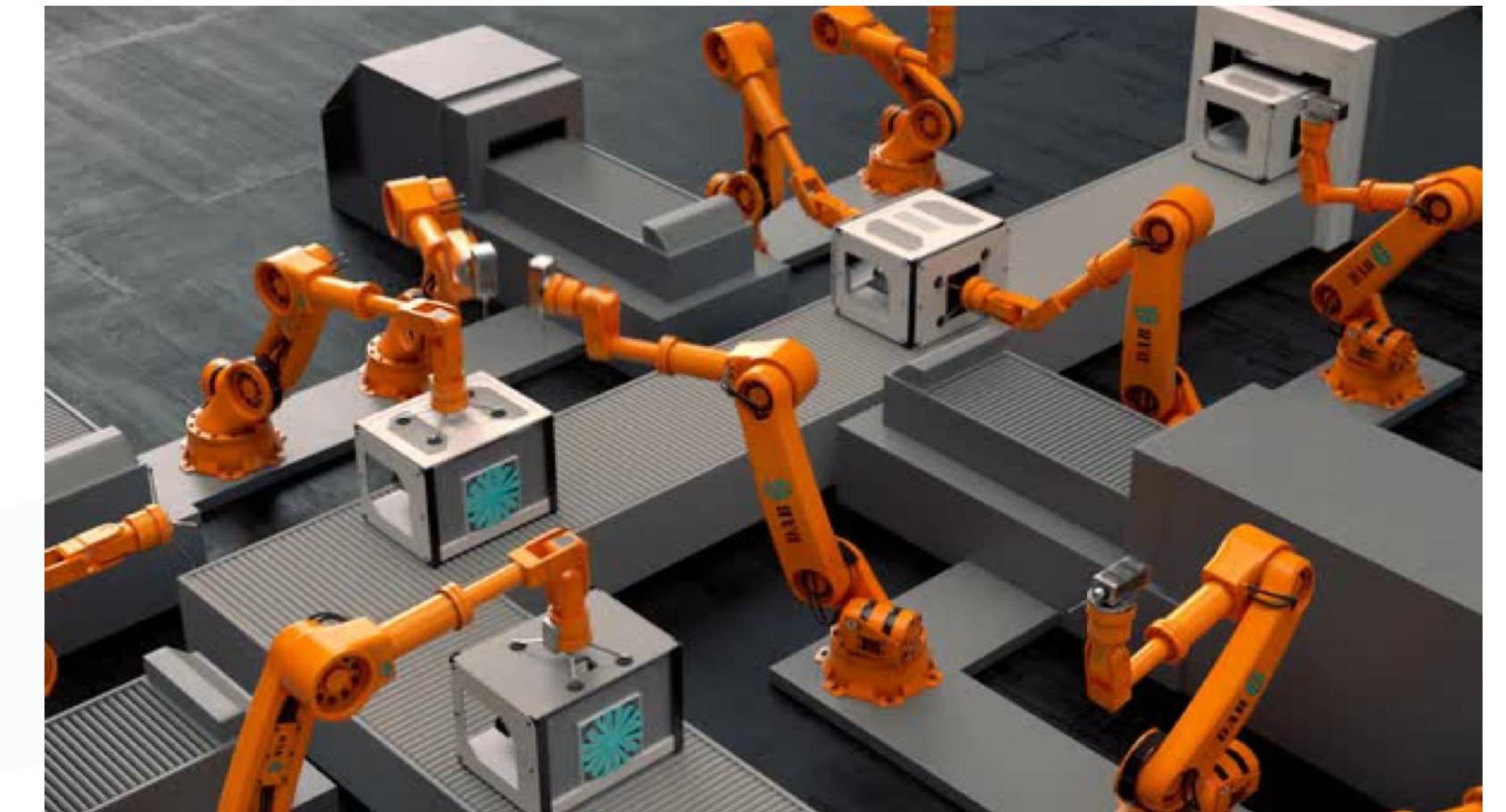
- Independent of underlying hardware



# Deployment

---

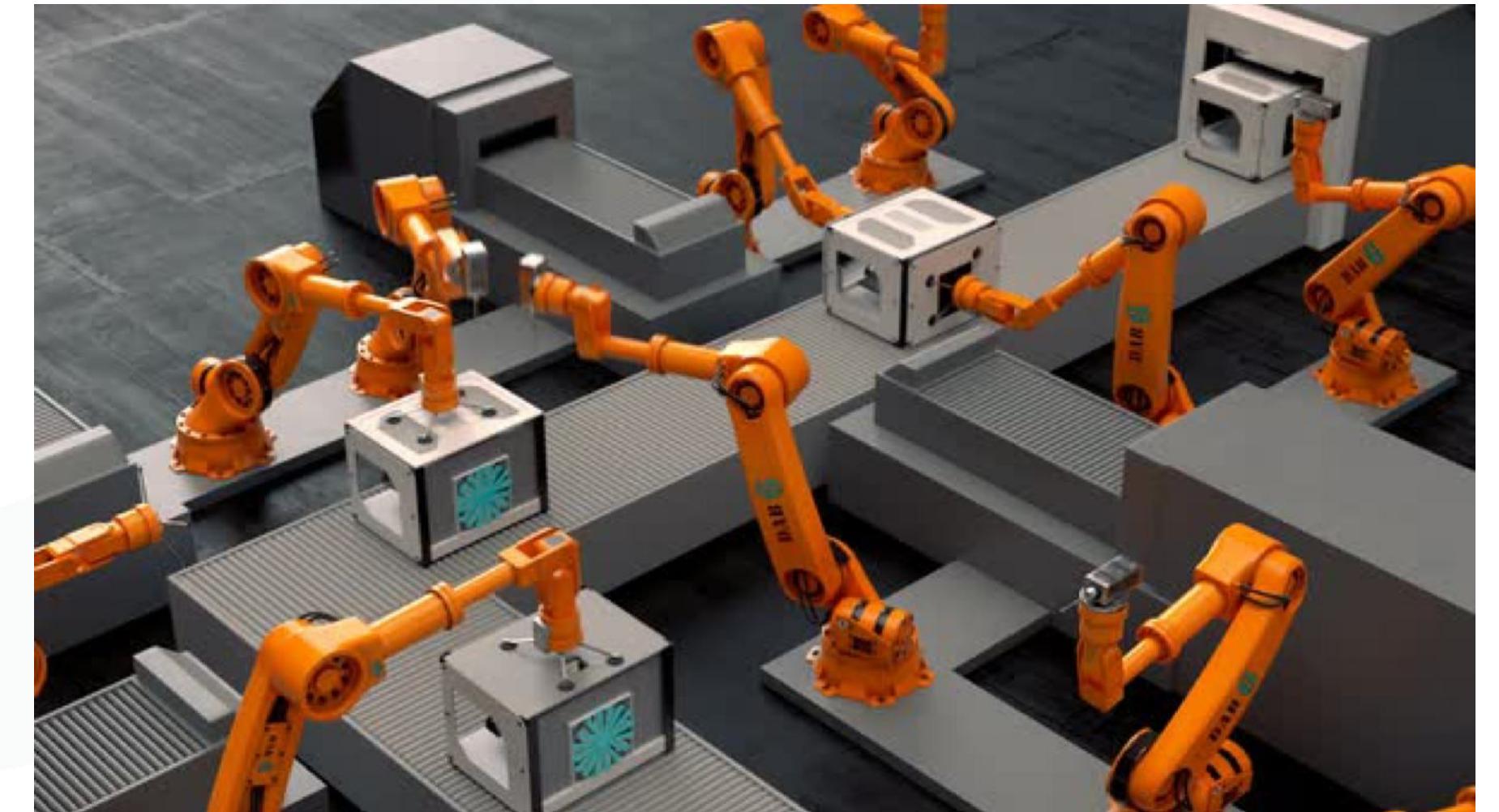
- Independent of underlying hardware
- Easy to put anywhere, in any cloud environment



# Deployment

---

- Independent of underlying hardware
- Easy to put anywhere, in any cloud environment
- Can be used in a cluster



# How does Docker work

---

# Definitions

---

Container

Image

Dockerfile

Volume

# Definitions

---

Container: “Lightweight, standalone, executable software package; includes everything: code, runtime, system tools, system libraries and settings.”

Image

Dockerfile

Volume

# Definitions

---

Container: “Lightweight, standalone, executable software package; includes everything: code, runtime, system tools, system libraries and settings.”

Image: A frozen snapshot of a container.

Dockerfile

Volume

# Definitions

---

Container: “Lightweight, standalone, executable software package; includes everything: code, runtime, system tools, system libraries and settings.”

Image: A frozen snapshot of a container.

Dockerfile: An (often) short text-file detailing how to make an image from another image.

Volume

# Definitions

---

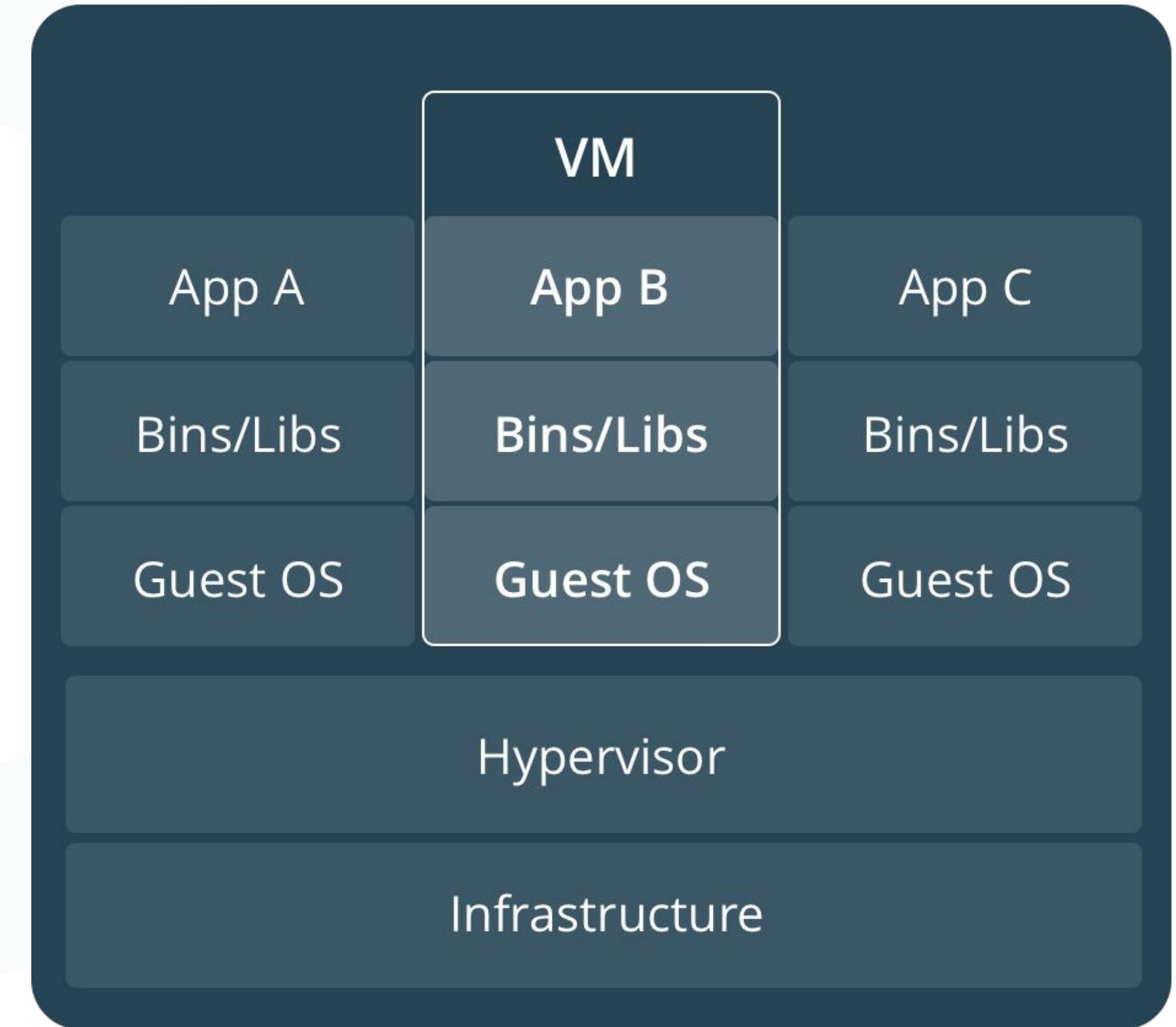
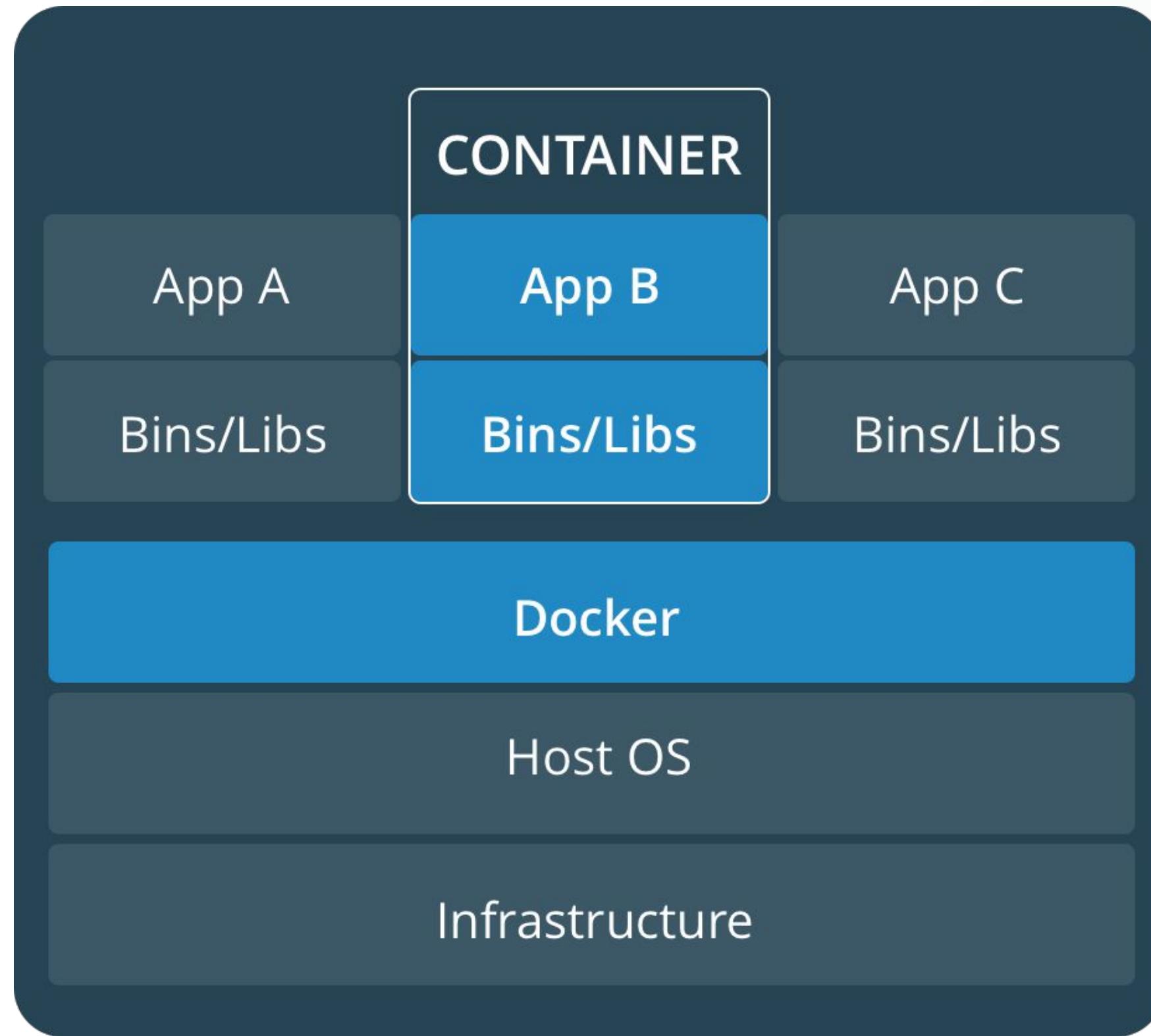
Container: “Lightweight, standalone, executable software package; includes everything: code, runtime, system tools, system libraries and settings.”

Image: A frozen snapshot of a container.

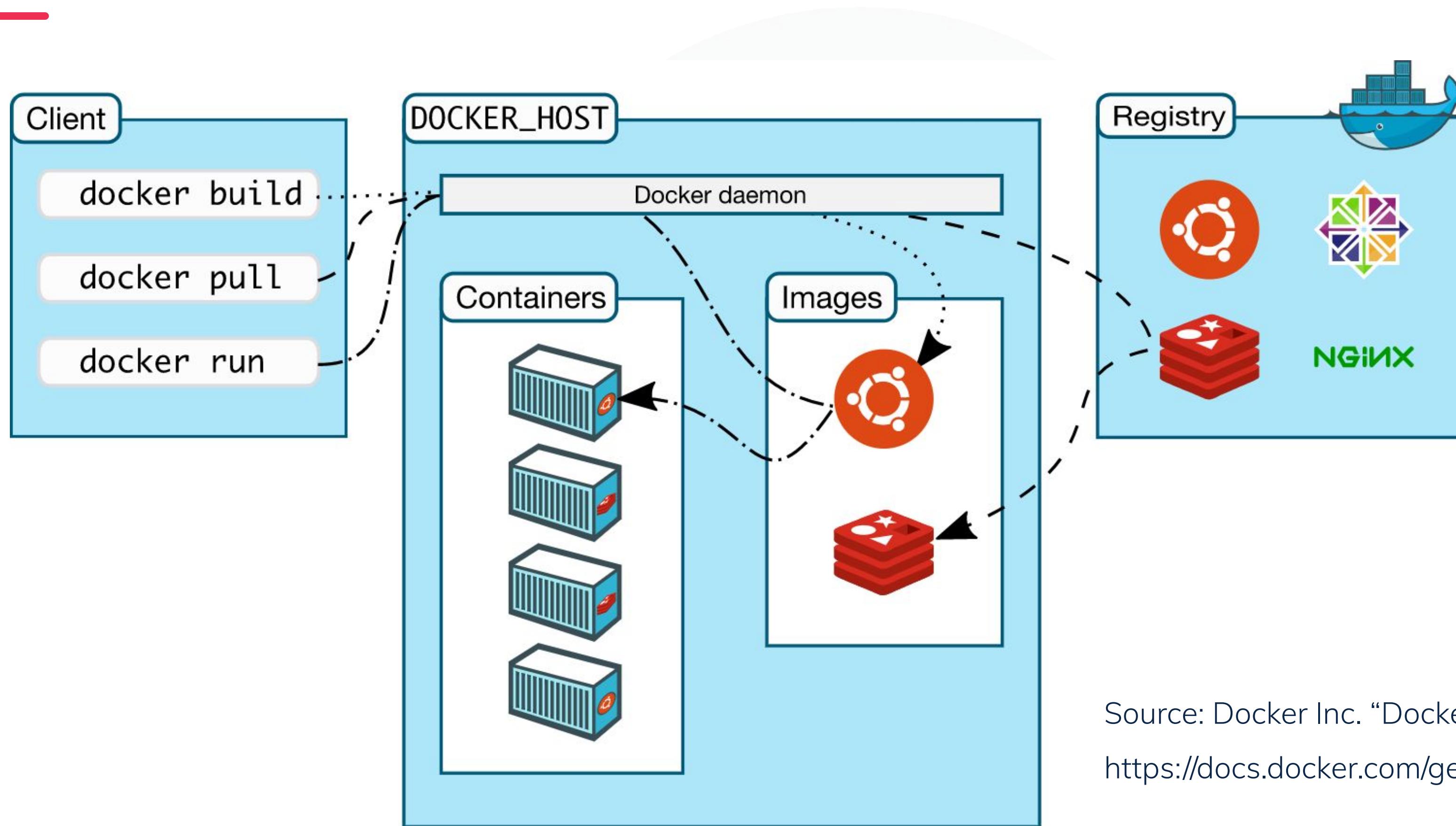
Dockerfile: An (often) short text-file detailing how to make an image from another image.

Volume: A chunk of filesystem (“folder”) that can be accessed by both containers as well as the host computer, and can persist.

# Definitions: Containers vs Virtual Machines



# How it actually works



Source: Docker Inc. “Docker overview” located at  
<https://docs.docker.com/get-started/overview/>



Search

OFFICIAL REPOSITORY

# python

Last pushed: 5 hours ago

Repo Info

Tags

## Short Description

Python is an interpreted, interactive, object-oriented, open-source programming language.

## Full Description

Supported tags and respective [Dockerfile](#) links

## Simple Tags

- `3.7.0b3-stretch`, `3.7-rc-stretch`, `rc-stretch` ([3.7-rc/stretch/Dockerfile](#))
- `3.7.0b3-slim-stretch`, `3.7-rc-slim-stretch`, `rc-slim-stretch`, `3.7.0b3-slim`, `3.7-rc-slim`, `rc-slim` ([3.7-rc/stretch/slim/Dockerfile](#))
- `3.7.0b3-alpine3.7`, `3.7-rc-alpine3.7`, `rc-alpine3.7`, `3.7.0b3-alpine`, `3.7-rc-alpine`, `rc-alpine` ([3.7-rc/alpine3.7/Dockerfile](#))

# Building our first container

---

Using Docker

# Dockerfile basics

---

```
1  FROM ubuntu:18.04
2  RUN apt-get update && apt-get install -y python python-pip
3  RUN pip install numpy
4  COPY ./calculate_pi.py /
5  ENTRYPOINT [ "python", "/calculate_pi.py" ]
6  CMD ["1000"]
7  |
```

[https://github.com/BigDataRepublic/docker\\_webinar](https://github.com/BigDataRepublic/docker_webinar)

# Dockerfile basics

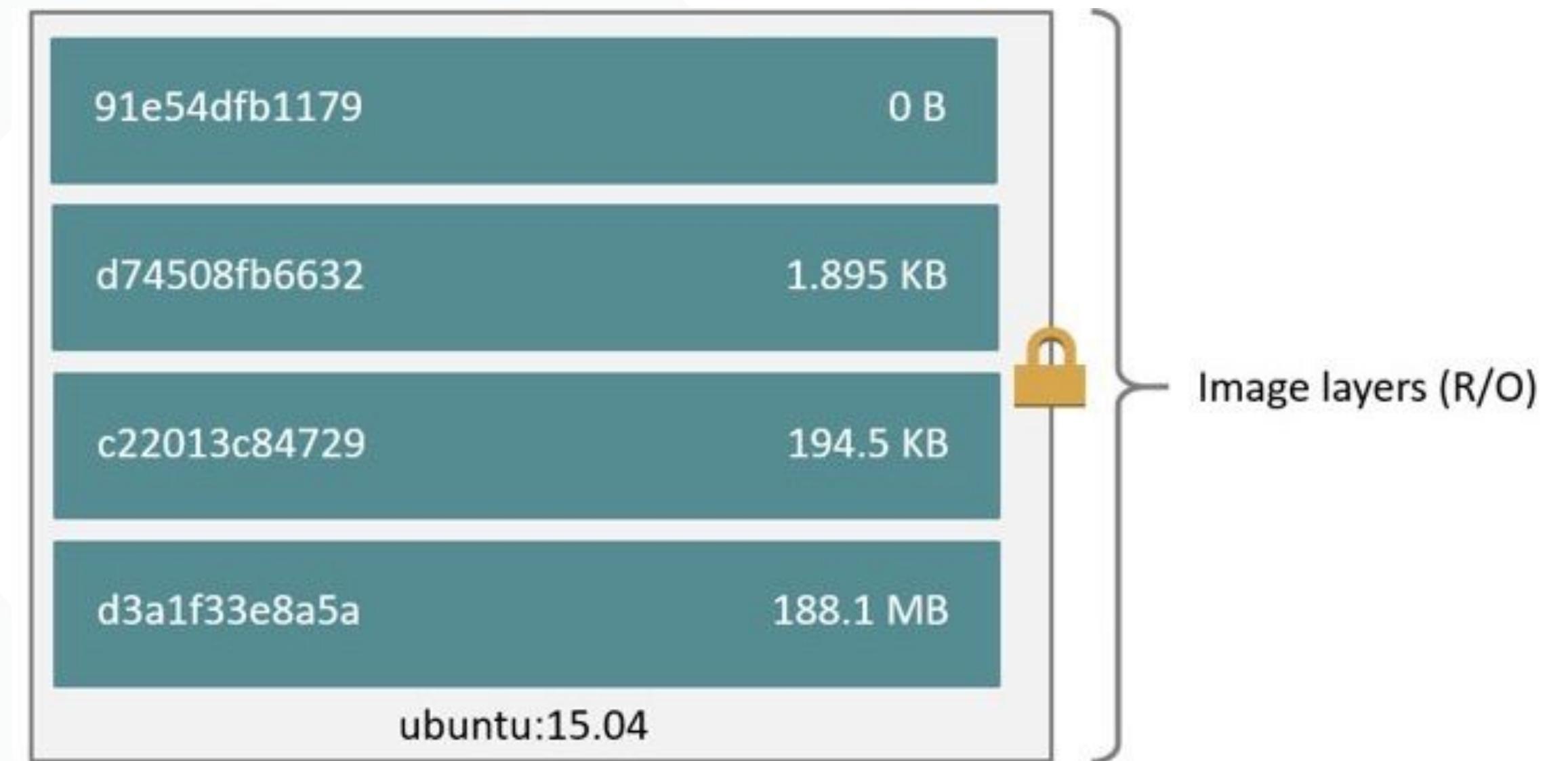
---

- Constructed as layers
- Each instruction ~ 1 layer
- Start with base image - **FROM**
- Copy files to image - **COPY** or **ADD**
- Execute shell commands - **RUN**
- Run model - **ENTRYPOINT** with **CMD**

# Image layers

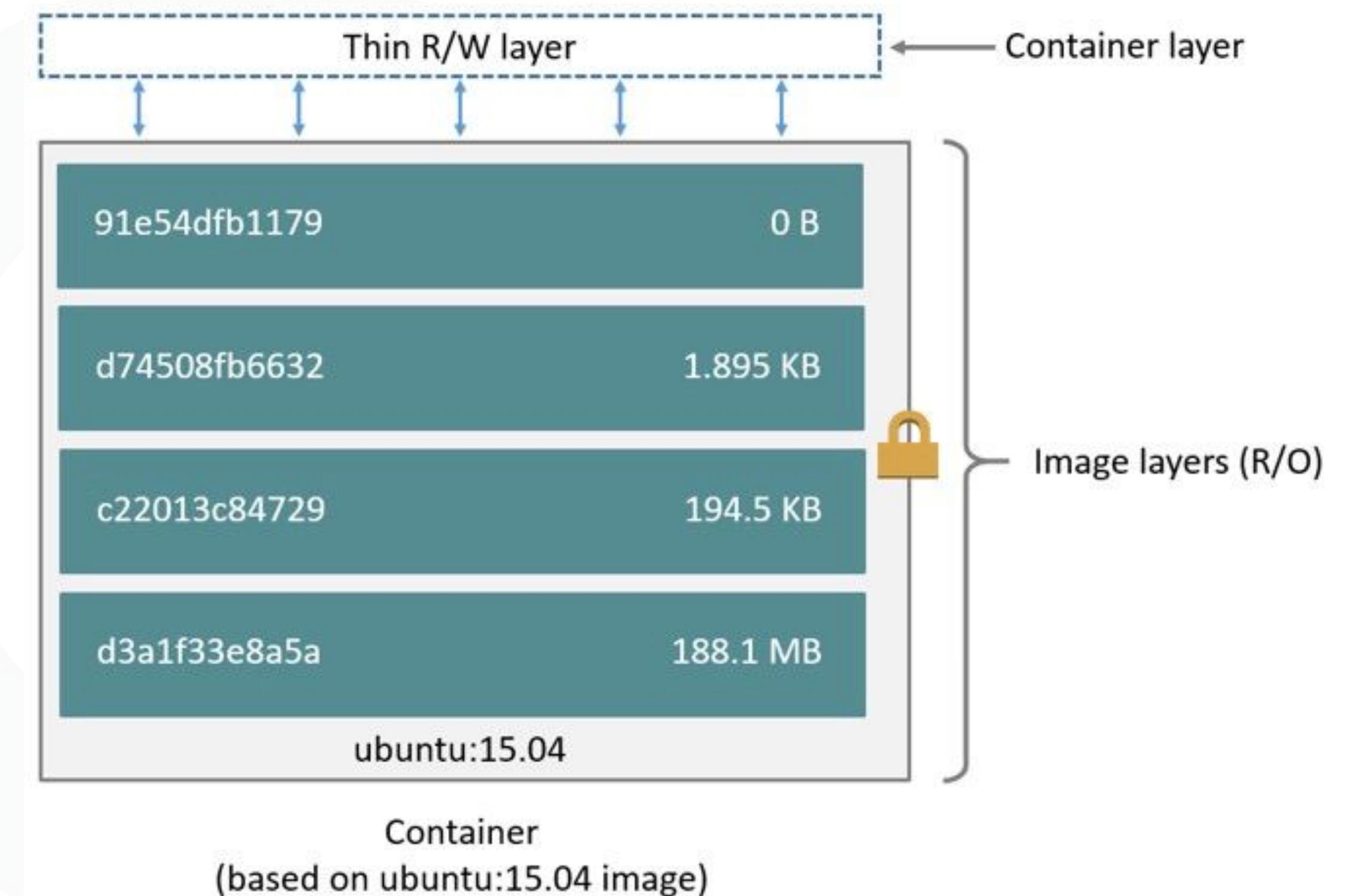
---

- Each layer read-only
- Layer stores changes from previous layer



# Image layers → container

- Each layer read-only
- Layer stores changes from previous layer
- Container > adds writeable layer on last image layer



# Running it all

---

. docker build -t approx\_pi -f

Dockerfile\_first .

. docker run approx\_pi:latest

```
docker build -t approx_pi -f Dockerfile_first .
Sending build context to Docker daemon 587.8kB
Step 1/6 : FROM ubuntu:18.04
--> 4e5021d210f6
Step 2/6 : RUN apt-get update && apt-get install -y python python-pip
--> Using cache
--> 19ccdd3989ae
Step 3/6 : RUN pip install numpy
--> Using cache
--> 4551ccf0b8a2
Step 4/6 : COPY ./calculate_pi.py /
--> Using cache
--> 5c5d4bb4a615
Step 5/6 : ENTRYPOINT [ "python", "/calculate_pi.py" ]
--> Using cache
--> 2ef60e6036a5
Step 6/6 : CMD ["1000"]
--> Using cache
--> 523087ef7bd6
Successfully built 523087ef7bd6
```

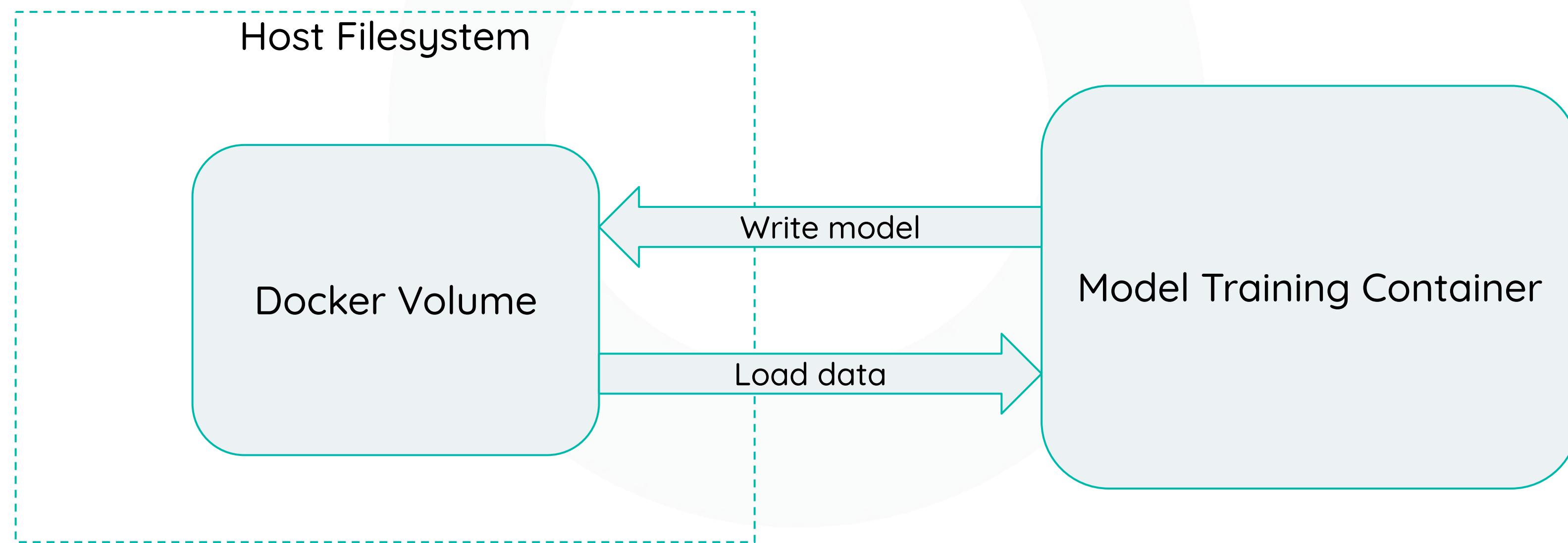
# Training a model with Docker

---

Using Docker for Data Science

# Training & saving models

- Python model training in container
- Need to retrieve our model after training
- Docker volumes: sharing storage



# Training & saving models

---

- To initialize with local filesystem folder:

```
docker run -v /foo/data:/app/data my_container
```

- docker build -t iris\_trainer -f Dockerfile\_training .
- docker run -v /full/path/to/model:/app/model iris\_trainer

# Training & saving models continued

---

## Dockerfile\_training

```
1  FROM ubuntu:18.04
2
3  LABEL version="1.0"
4  LABEL description="This container trains a random forest on the iris dataset."
5  LABEL maintainer="joris.bukala@bigdatarepublic.nl,jacobus.herman@bigdatarepublic.nl"
6
7  COPY --chown=root:root ./security/add_user_group.sh ./training_app/requirements.txt /app/
8  RUN apt-get update -y && \
9      apt-get install -y python-pip python-dev && \
10     pip install -r /app/requirements.txt && \
11     chmod 550 /app/add_user_group.sh && \
12     /app/add_user_group.sh && \
13     rm /app/add_user_group.sh && \
14     rm -rf /var/lib/apt/lists/* && \
15     mkdir -p /app/model && \
16     chown -R dockerapp:dockerapp /app/model
17
18  USER dockerapp
19
20  COPY --chown=dockerapp:dockerapp ./training_app /app
21
22  VOLUME [ "/app/model" ]
23  ENTRYPOINT [ "python" ]
24  CMD [ "/app/train.py" ]
```

# Serving predictions

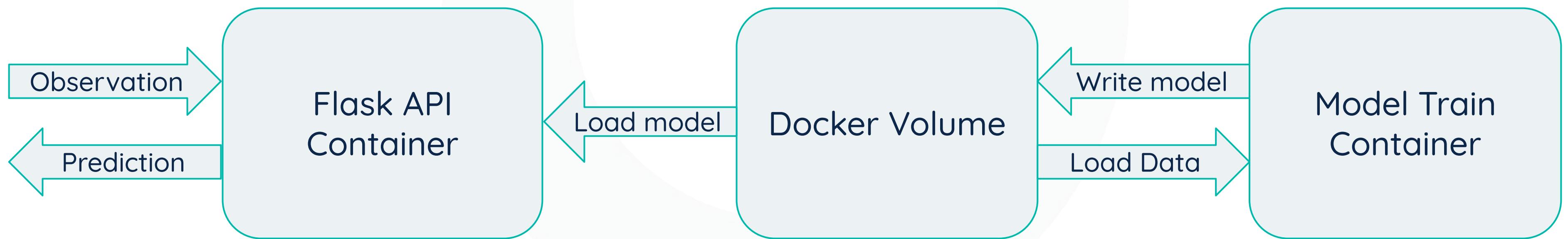
---

Using Docker for Data Science

# Serving models

---

- Use trained model with Flask API
- Share same volume
- Receives array with 2 numbers & returns number



# What is Flask?

- Micro framework for building web applications
- Allows creation of interface to something through web communication
- Similar to website: user interacts with content
- End result (web) Application Programming Interface (API)



Flask source: <https://palletsprojects.com/p/flask/>

# Serving models

---

- . docker build -t iris\_predictor -f Dockerfile\_prediction .
  - . docker run -p 5000:5000 -v ~/model\_volume:/app/model iris\_predictor
- 
- . **-p** maps host port to container port
  - . Communication to 5000 on host goes to 5000 on container (like a mirror)
  - . Ports can be different, e.g. **-p 80:5000**

# Serving models image

---

## Dockerfile\_prediction

```
1  FROM ubuntu:18.04
2
3  LABEL version="1.0"
4  LABEL description="This container exposes a fitted random forest model with Flask on the Iris dataset."
5  LABEL maintainer="joris.bukala@bigdatarepublic.nl,jacobus.herman@bigdatarepublic.nl"
6
7  COPY --chown=root:root ./security/add_user_group.sh ./prediction_app/requirements.txt /app/
8  RUN apt-get update -y && \
9      apt-get install -y python-pip python-dev && \
10     pip install -r /app/requirements.txt && \
11     chmod 550 /app/add_user_group.sh && \
12     /app/add_user_group.sh && \
13     rm /app/add_user_group.sh && \
14     rm -rf /var/lib/apt/lists/* && \
15     mkdir -p /app/model && \
16     chown -R dockerapp:dockerapp /app/model
17
18  USER dockerapp
19
20  COPY --chown=dockerapp:dockerapp ./prediction_app /app
21
22  ENV FLASK_APP=/app/main.py
23
24  VOLUME [ "/app/model" ]
25  EXPOSE 5000
26  ENTRYPOINT [ "flask" ]
27  CMD [ "run", "--host=0.0.0.0" ]
```

# Making a front-end

---

Using Docker for Data Science

# Add web UI

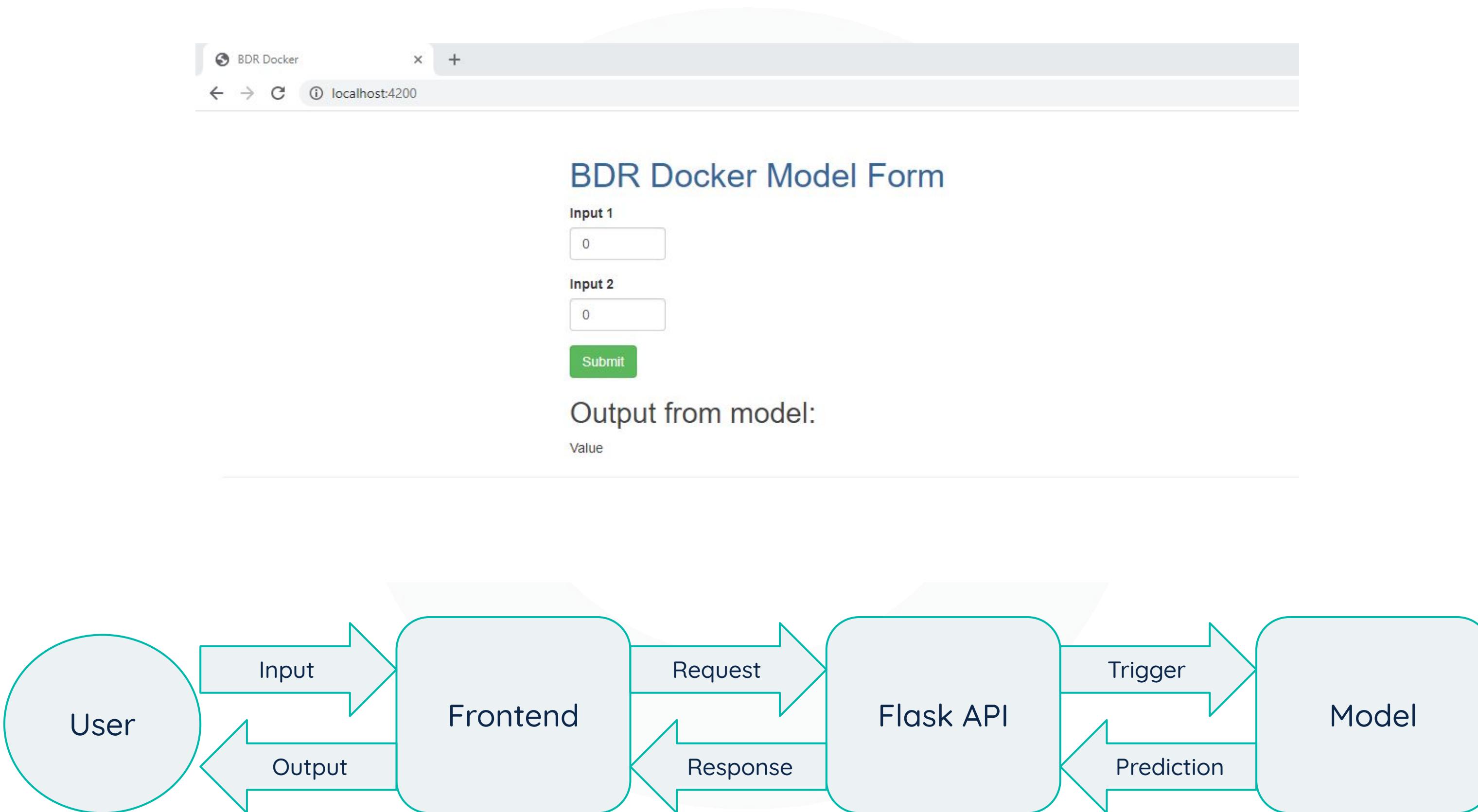
---

- . AngularJS web application for access to production model
- . Frontend container interacts with prediction container

## To run:

- . `docker build -t iris_frontend -f Dockerfile_frontend .`
- . `docker run -d -p 4200:4200 iris_frontend`

# Web UI flow



# Web UI image

---

## Dockerfile\_frontend

```
1  FROM node:12.16.2
2
3  LABEL version="1.0"
4  LABEL description="This container runs an Angular JS app that connects to a REST API on localhost:5000."
5  LABEL maintainer="joris.bukala@bigdatarepublic.nl,jacobus.herman@bigdatarepublic.nl"
6
7  WORKDIR /app
8
9  ENV PATH /app/node_modules/.bin:$PATH
10
11 COPY --chown=root:root ./security/add_user_group.sh ./frontend/model-display/package.json /app/
12 RUN npm install && \
13     chmod 550 /app/add_user_group.sh && \
14     /app/add_user_group.sh && \
15     rm /app/add_user_group.sh
16
17 USER dockerapp
18
19 COPY --chown=dockerapp:dockerapp ./frontend/model-display /app/
20
21 EXPOSE 4200
22 ENTRYPOINT [ "ng", "serve" ]
23 CMD [ "--host", "0.0.0.0", "--disableHostCheck", "true" ]
```

# Best practices

---

For writing Dockerfiles

# Docker-compose

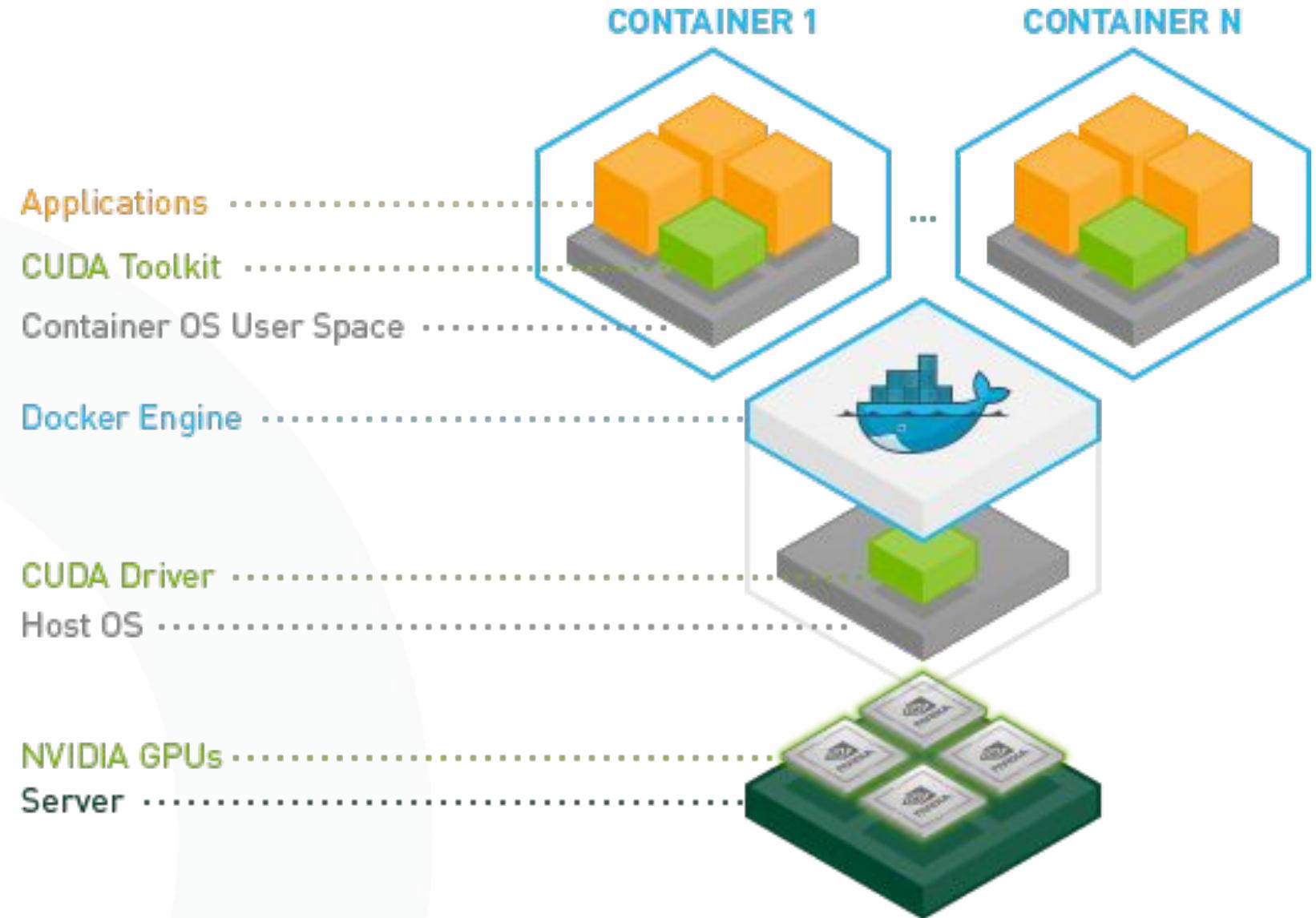
---

- All containers with arguments in one config file
- docker-compose up -d
- docker-compose down

```
1 version: "3.4"
2 services:
3   iris_frontend:
4     build:
5       context: .
6       dockerfile: Dockerfile_frontend
7     ports:
8       - "4200:4200"
9     depends_on:
10      - iris_prediction
11   networks:
12     default:
13       aliases:
14         - frontend
15   iris_prediction:
16     build:
17       context: .
18       dockerfile: Dockerfile_prediction
19     ports:
20       - "5000:5000"
21     depends_on:
22       - iris_training
23     volumes:
24       - "iris_model_vol:/app/model"
25   networks:
26     default:
27       aliases:
28         - rest-api
29   iris_training:
30     build:
31       context: .
32       dockerfile: Dockerfile_training
33     volumes:
34       - "iris_model_vol:/app/model"
35   volumes:
36     iris_model_vol:
37   networks:
38     default:
```

# Using a GPU with NVIDIA-Docker

- Training/ preparing DL model on your (Linux) laptop
- Quick deployment to Cloud service when necessary



```
#### Test nvidia-smi with the latest official CUDA image  
docker run --gpus all nvidia/cuda:10.0-base nvidia-smi
```

# Dockerfile writing

---

- Use **.dockerignore**: avoids unnecessary files, speeds up build process
- Create stateless images: able to restart without data loss
- Single purpose: decoupling interacting parts for simple containers
- Multi-stage builds: single image file to build and deploy app

```
1  **/.env
2  **/.idea
3  **/.vscode
4  **/.git
5  **/.cache
6  **/docker-compose.yml
7  **/*.md
8  **/node_modules
```

```
1  FROM golang:1.7.3
2  WORKDIR /go/src/github.com/alexellis/href-counter/
3  RUN go get -d -v golang.org/x/net/html
4  COPY app.go .
5  RUN CGO_ENABLED=0 GOOS=linux go build -a -installsuffix cgo -o app .
6
7  FROM alpine:latest
8  RUN apk --no-cache add ca-certificates
9  WORKDIR /root/
10 COPY --from=0 /go/src/github.com/alexellis/href-counter/app .
11 CMD ["./app"]
```

# Dockerfile writing continued

---

- Minimize layers: combine **RUN** statements
- Frequently changing commands bottom of Dockerfile
- Reduce container size:
  - Use Alpine Linux base image for smallest size
  - Add `rm -rf /var/lib/apt/lists/*` to apt-get install
  - Try `--no-install-recommends` flag in apt-get install
- Use trusted official base images
- Specify image versions

```
$ docker build -t mynewimage .
Sending build context to Docker daemon 3.072kB
Step 1/3 : FROM python:3.8-slim-buster
--> 3d8f801fc3db
Step 2/3 : COPY build.sh .
--> 541b65a7b417
Step 3/3 : RUN ./build.sh
--> Running in 9917e3865f96
Building...
```

# Summary & next steps

---

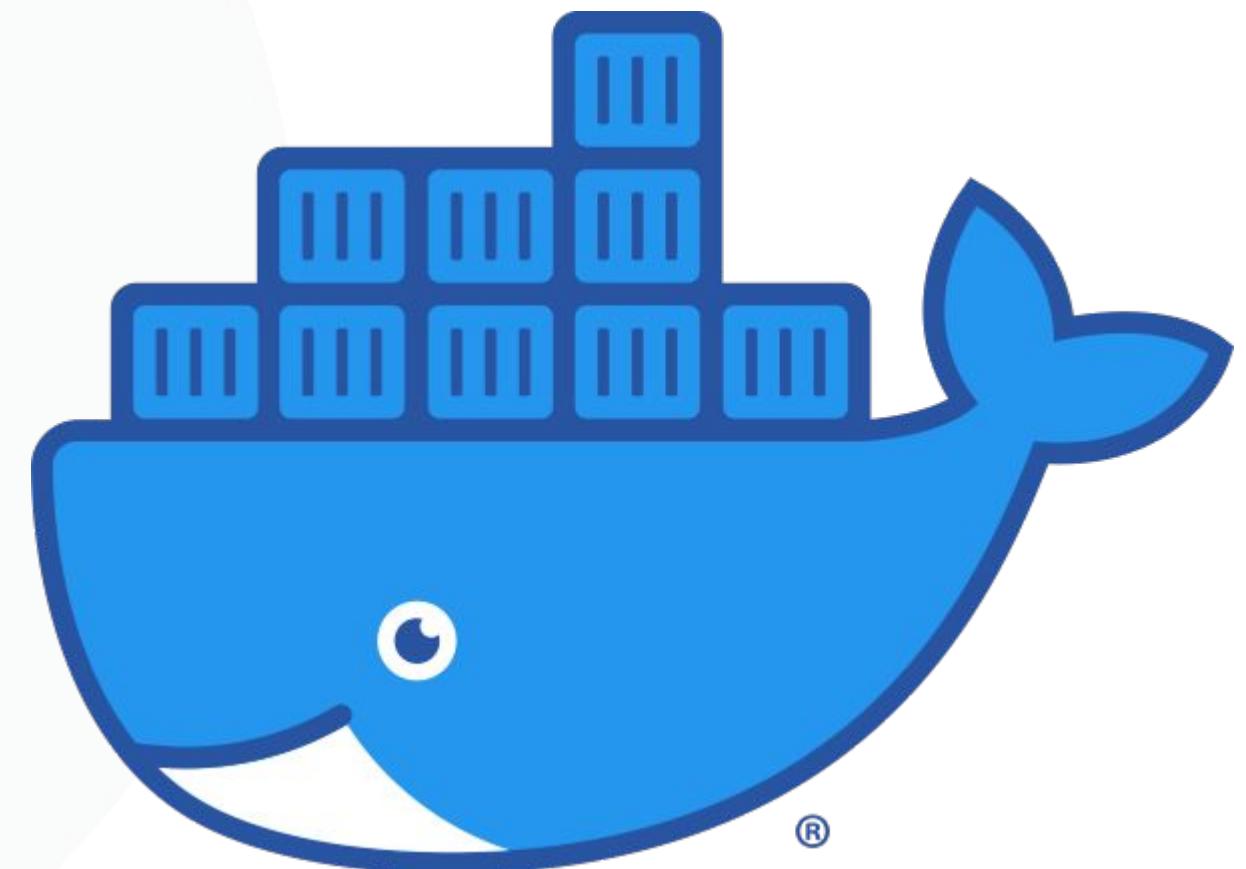
## Docker usage

# Summary

---

- Docker decouples project from environment → reproducibility
  - Eases sharing of models → better collaboration
  - Improves deployability → greater impact
- 
- Be aware of security and optimizations

*Next steps: Docker-swarm & Kubernetes*



[https://github.com/BigDataRepublic/docker\\_webinar](https://github.com/BigDataRepublic/docker_webinar)

# Presenters and contributors

---

**Tom de Ruijter**

*Data Scientist*

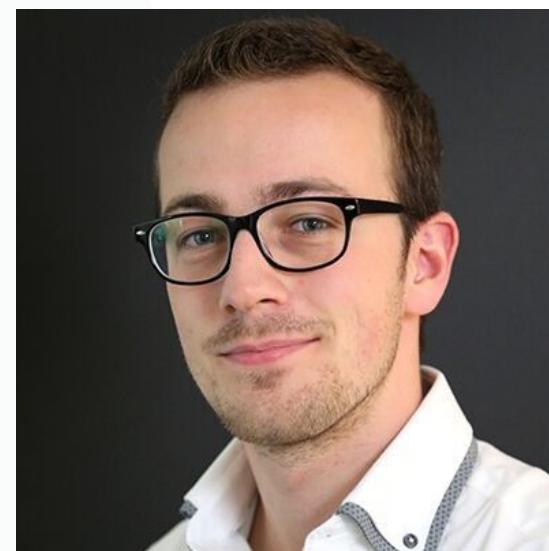
[/in/tomderuijter/](#)



**Joris Bukala**

*Data Scientist*

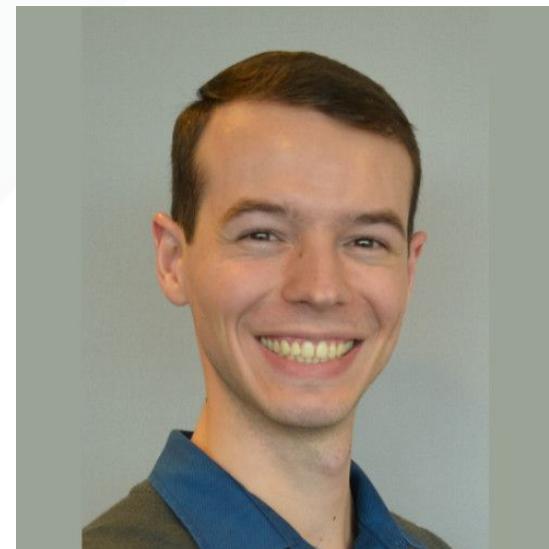
[/in/joris-bukala/](#)



**Jacobus Herman**

*Data Engineer*

[/in/jacobus-herman/](#)





0100010

**Phone** +31 (0)168 479294

**Email** [info@bigdatarepublic.nl](mailto:info@bigdatarepublic.nl)

**Address** Colbaan 4C, 3439 NG Nieuwegein, The Netherlands