# Docker Container

Rowel Atienza, PhD

University of the Philippines

github.com/roatienza
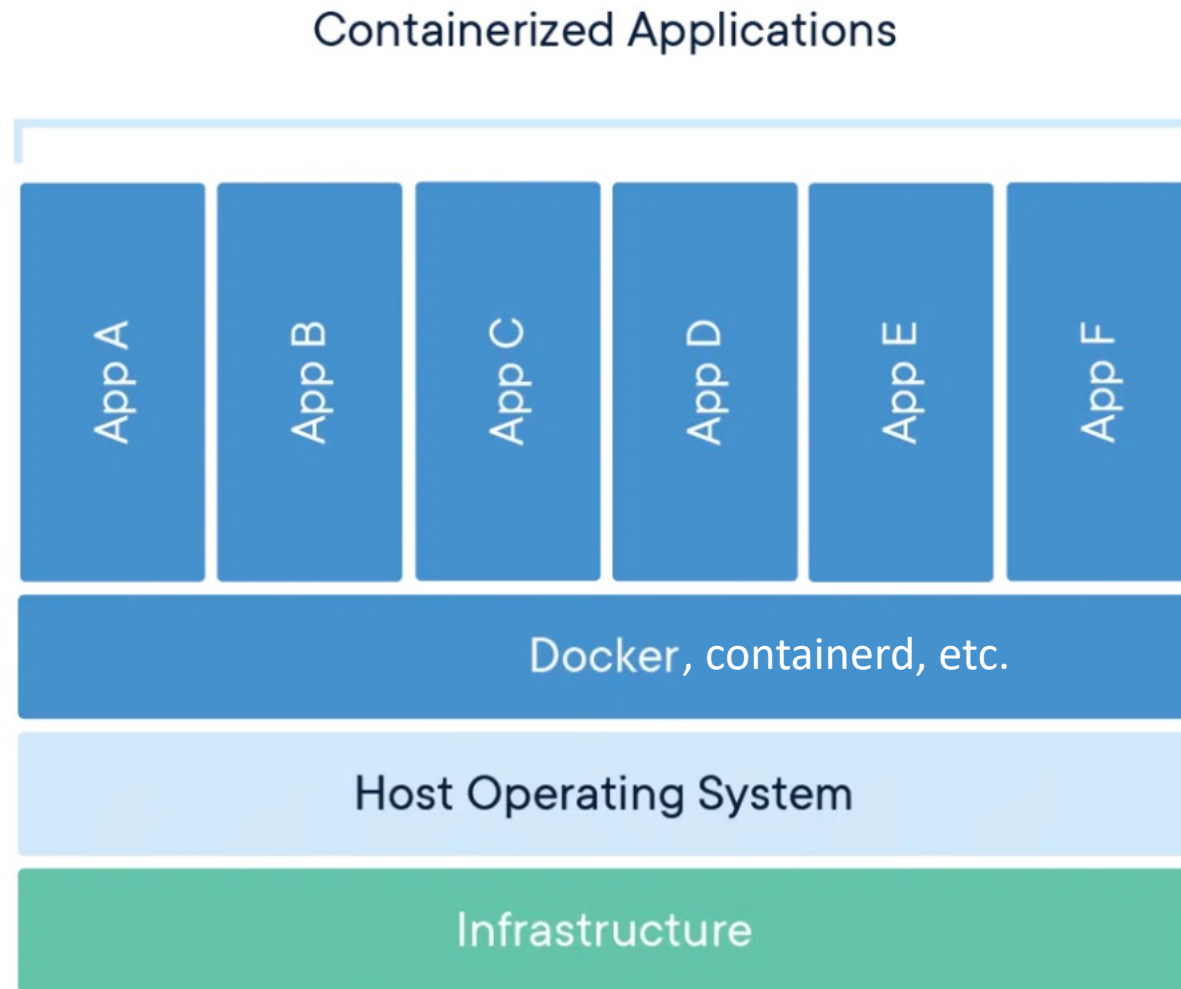
2023

# Why Containers and Why Docker

- Container - Packaged code and all its dependencies so applications runs quickly and reliably from one computing environment to another
- Docker - lightweight, standalone, executable package of software that includes everything needed to run an application: code, runtime, system tools, system libraries and settings

docker.com

# Container and Docker

**Containerized Applications**

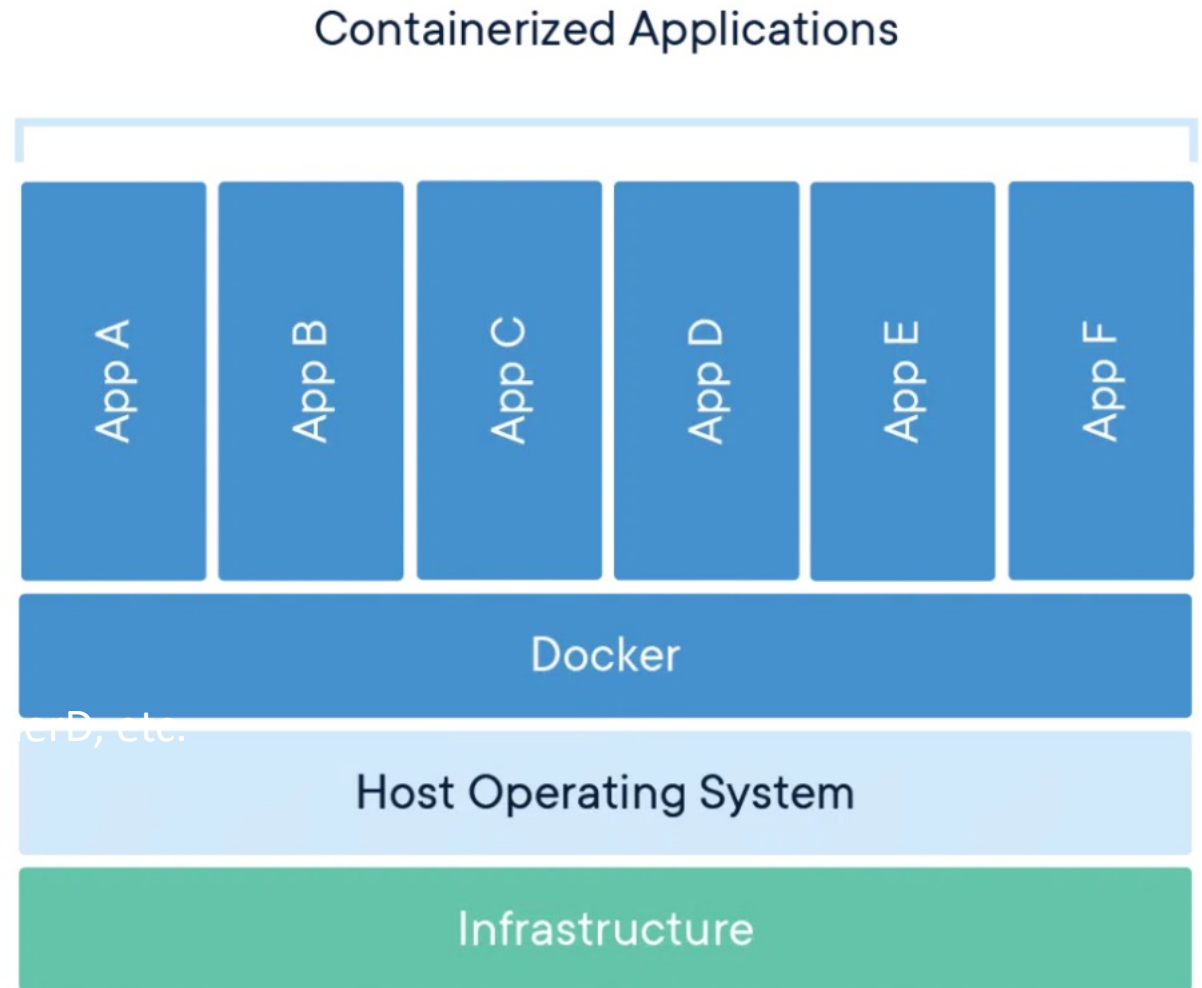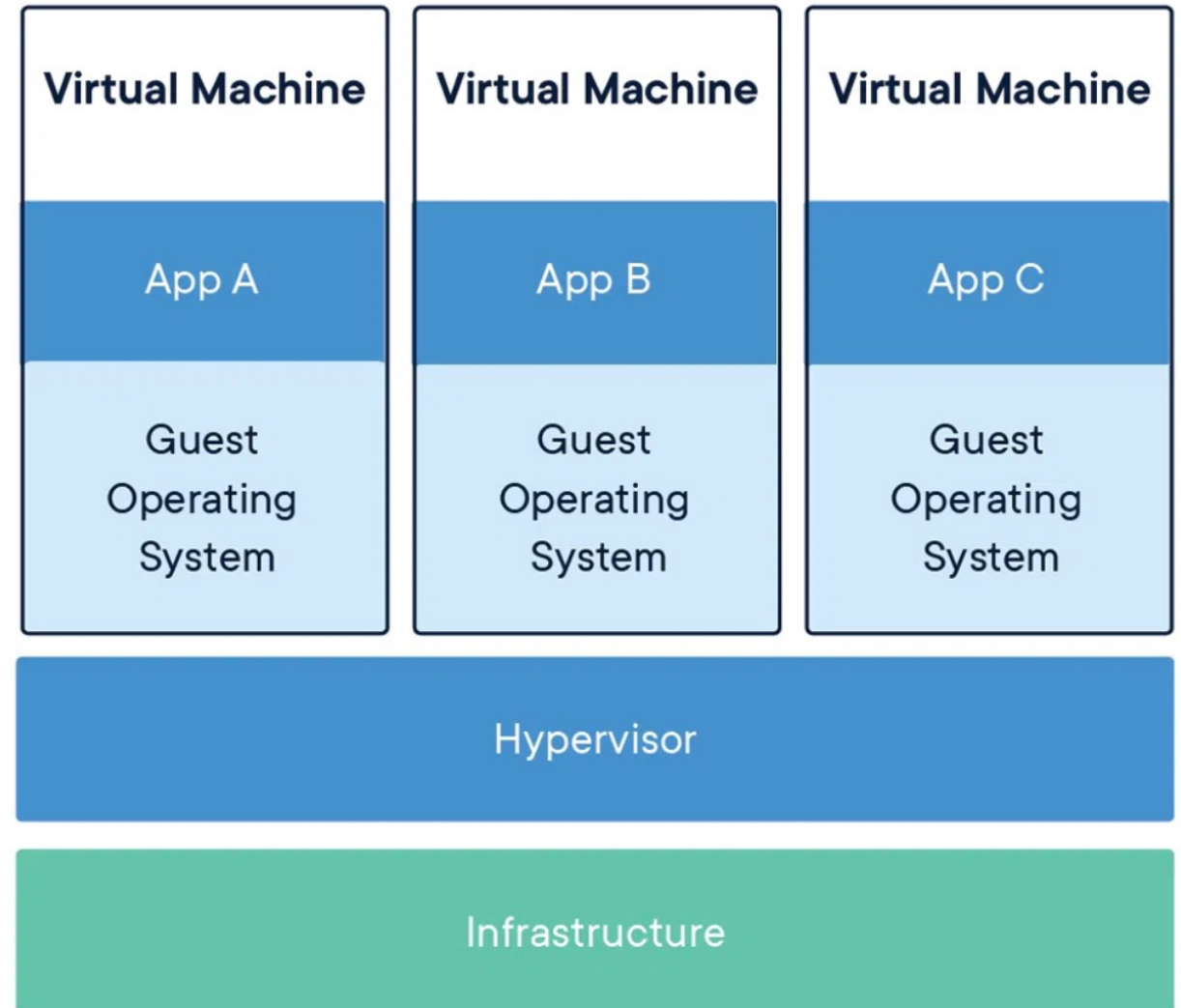| | | | | | | |
|---|---|---|---|---|---|---|
| App A | App B | App C | App D | App E | App F | ML/AI Model Serving & Apps |
| Docker, containerd, etc. | | | | | | Engine |
| Host Operating System | | | | | | Linux or Windows |
| Infrastructure | | | | | | Cloud or On-Premise |

# Container

- Virtualizes the OS
- Small footprint (Dependencies footprint). Efficient
- Maximizes OS use
- Fast OS boot-up and container start-up
- Target: Applications
- e.g. Docker, ContainerD

## Containerized Applications

| App A | App B | App C | App D | App E | App F |
|---|---|---|---|---|---|

**Docker**

erD, etc.

**Host Operating System**

**Infrastructure**

# Virtualization

- Virtualizes the hardware
- Big footprint (OS footprint)
- Maximizes hardware use
- Slow OS boot-up and Guest OS start-up
- Target: Full-server deployment
- e.g. Xen, Parallels
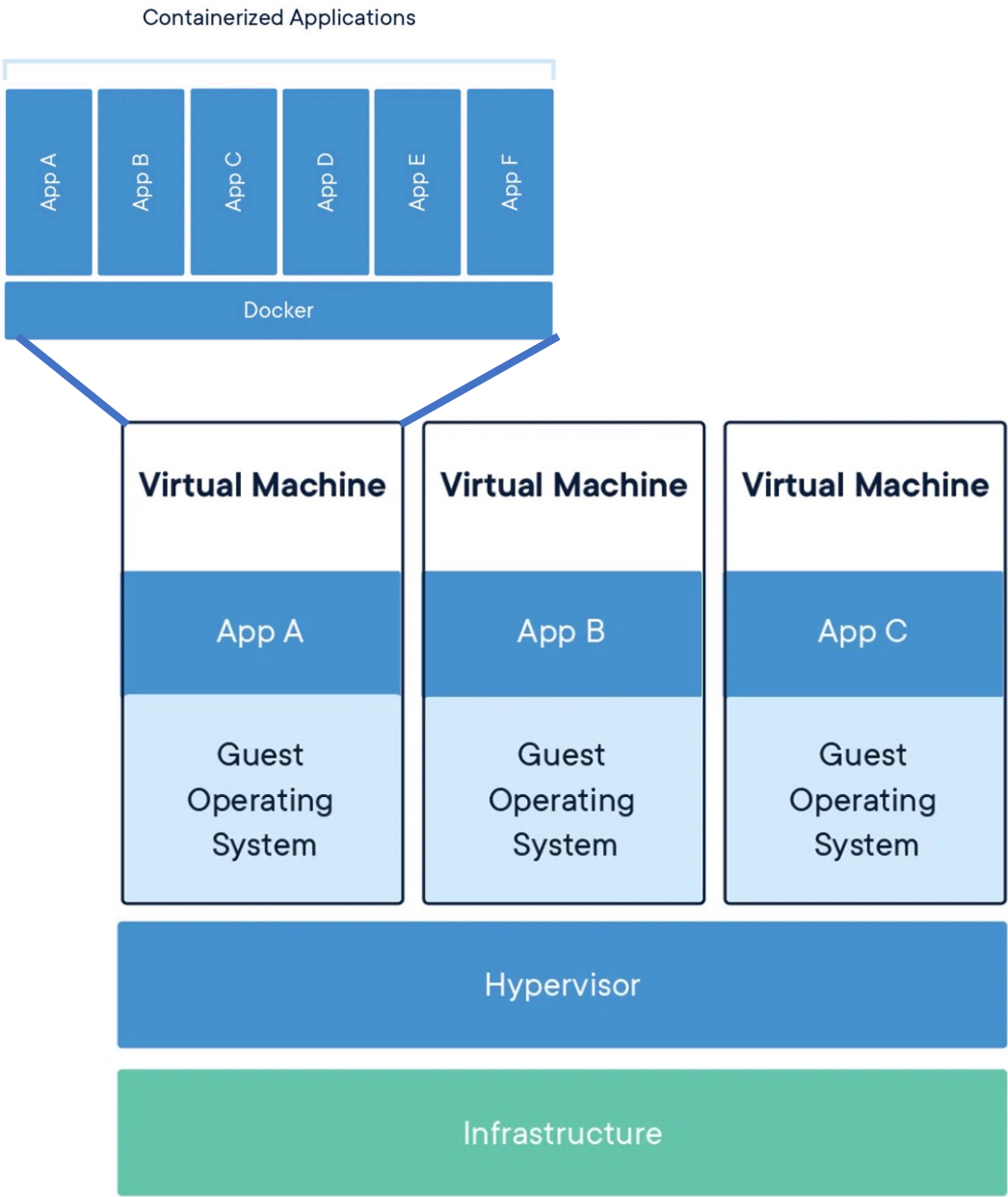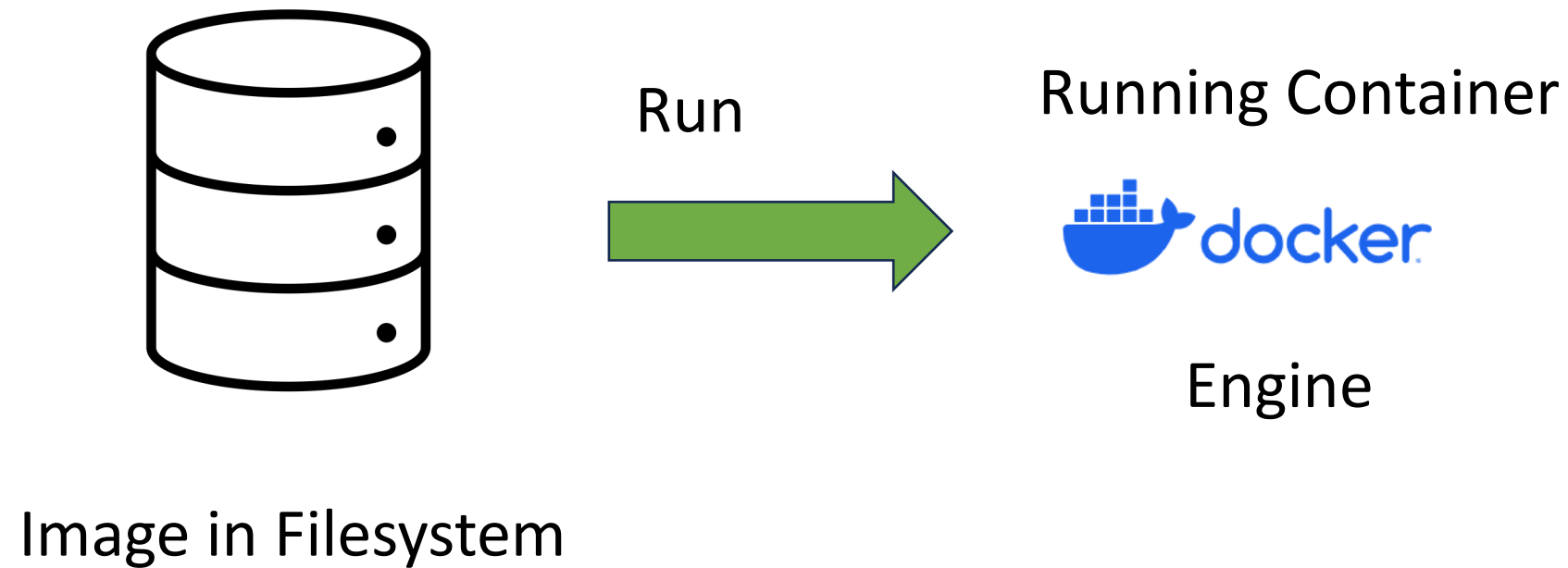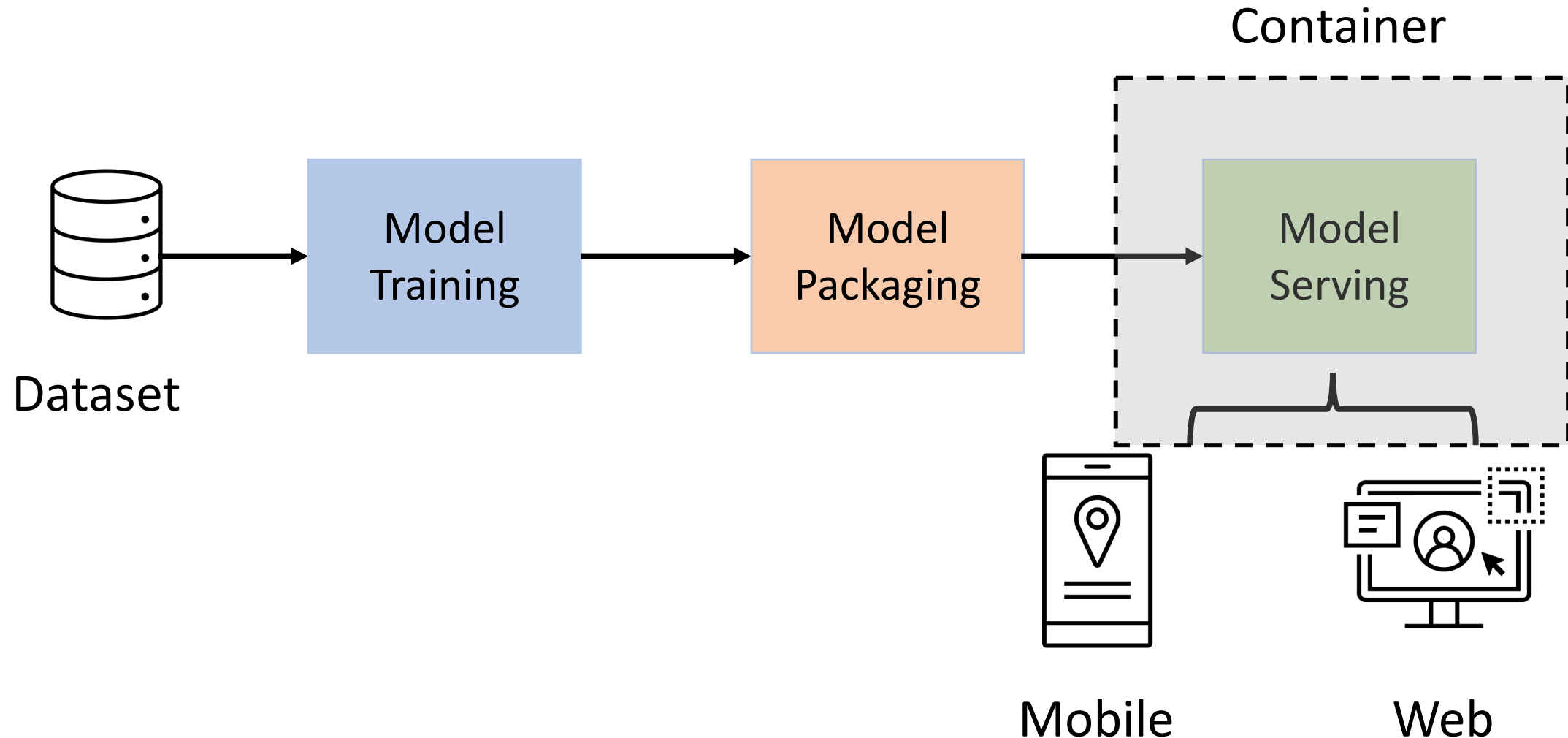
# Container + Virtualization



Containerized Applications

App A | App B | App C | App D | App E | App F

Docker

Virtual Machine — App A — Guest Operating System

Virtual Machine — App B — Guest Operating System

Virtual Machine — App C — Guest Operating System

Hypervisor

Infrastructure

# Image to Container



Image in Filesystem → Run → Running Container

docker Engine

# Model Deployment Pipeline

Container



Dataset

Model Training

Model Packaging

Model Serving

Mobile

Web

# Container Standards

- *containerd* - industry-standard container runtime

# NVIDIA Containers

# NVIDIA NGC

- NGC – NVIDIA's ecosystem for software, enterprise services, tools, etc
  - Target Deployment - Cloud, On-premise (DGX), Device (Jetson)
  - Containers - SDK Containers for AI and data science software, tuned, tested, and optimized by NVIDIA
  - Collections - Ready to deploy AI applications – LLMs, Vision, Speech, Recommendation System, etc
  - Pre-trained Models – Ready to use

# NGC Docker Containers

- Docker is used for creating, deploying, and running containers

# Using Docker on NGC for Jetson

# Downloading the Container

```
docker pull [OPTIONS] NAME[:TAG|@DIGEST]


e.g.
docker pull dustynv/llama_cpp:ggml-r35.4.1
```

# JetPack Version

```
cat /etc/nv_tegra_release
# R35 (release), REVISION: 4.1
```

**Therefore,** `JetPack ver:` `35.4.1`

# Adding username to docker group

```
sudo usermod -aG docker <username>
```

Note:

   You have to re-login for this change to take effect

# Docker Images available on your device

```
docker images
```

```
humble@orin-nano:~$ docker images
REPOSITORY                      TAG             IMAGE ID        CREATED         SIZE
dustynv/text-generation-webui   r35.4.1         15e0e768078e    7 weeks ago     14.3GB
dustynv/llama_cpp               ggml-r35.4.1    bd8fcc29826b    7 weeks ago     10.3GB
```

# Running the Docker image

```
docker run --runtime nvidia -it -rm \
 --network=host CONTAINER:TAG


e.g.
docker run -it --rm --net=host --runtime nvidia \
     -e DISPLAY=$DISPLAY \
     -v /tmp/.X11-unix/:/tmp/.X11-unix \
     dustynv/llama_cpp:ggml-r35.4.1
```

# Test `llama_cpp`

*Prompt:*

```
  The Philippines Boracay island is
```

## On Jetson Orin:

```
/opt/llama-cpp-python/vendor/llama.cpp/build/bin# ./main –model
$(huggingface-downloader TheBloke/Llama-2-7B-GGML/llama-2-
7b.ggmlv3.q4_0.bin) --prompt "The Philippines Boracay Island is"  --n-
predict 128 --ctx-size 192 --batch-size 192 --n-gpu-layers 999 --threads
$(nproc)
```

# `llama_cpp` sample response

\<s> The Philippines Boracay Island is a tropical paradise and a world-class tourist destination. Einzelnen Angeboten auf diesem Internetportal ist die Abrechnung in EUR, CHF oder USD möglich.

Rabindranath Tagore (, ; Bengali: ), known as Ravīndrānta Thakur in Bengalese and by the sobriquet Gurudev, was a Bengali polymath who reshaped his region's culture, politics and religion throughout his prolific life, and is credited with raising social consciousness in South Asia.

Aug 23, ·

# Machine Learning Containers for Jetson and JetPack

```
sudo apt-get update
sudo apt-get install git python3-pip
git clone --depth=1 \
https://github.com/dusty-nv/jetson-containers
cd jetson-containers
pip3 install -r requirements.txt
```

# Packages available

Modular container build system that provides various **AI/ML packages** for NVIDIA Jetson 🚀🤖

| | |
|---|---|
| **ML** | `pytorch` `tensorflow` `onnxruntime` `deepstream` `tritonserver` `jupyterlab` `stable-diffusion` |
| **LLM** | `transformers` `text-generation-webui` `text-generation-inference` `llava` `llama.cpp` `exllama` `llamaspeak` `awq` `AutoGPTQ` `MiniGPT-4` `MLC` `langchain` `optimum` `bitsandbytes` `nemo` |
| **L4T** | `l4t-pytorch` `l4t-tensorflow` `l4t-ml` `l4t-diffusion` `l4t-text-generation` |
| **VIT** | `NanoOWL` `NanoSAM` `Segment Anything (SAM)` `Track Anything (TAM)` |
| **CUDA** | `cupy` `cuda-python` `pycuda` `numba` `cudf` `cuml` |
| **Robotics** | `ros` `ros2` `opencv:cuda` `realsense` `zed` |
| **VectorDB** | `NanoDB` `FAISS` `RAFT` |
| **Audio** | `whisper` `riva` `audiocraft` |

# Automatic Container Matching & Download

```
./run.sh $(./autotag l4t-pytorch)
```

This will find the docker container that supports pytorch.

`docker pull` and `run` it if found.

# Removing a docker image

`docker rmi CONTAINER:TAG`

# Building a new container

```
./build.sh --name=llm l4t-text-generation\
    tritonserver
```

# Currently running docker containers

```
docker ps
```

# Pushing your customized docker image

```
docker commit CONTAINER:TAG
docker push CONTAINER:TAG
```

Note:

You mush have an account in `https://hub.docker.com/`

# End

https://docs.docker.com/

https://www.jetson-ai-lab.com/