



We will be starting soon



Docker for Public Health Bioinformatics

Week 01 - Introduction to Docker and Containerization

PRESENTED BY:

Inês Mendes, PhD



Course Introduction

Training Workshop Resources

Training Information, Communication, and Support

- **GitHub Repository** created to host training resources and information:
 - <https://github.com/theiagen/Mid-Atlantic-Docker4PH-2025>
- **Support contact:**
 - support@theiagen.com

Meet the Team



Inês Mendes, PhD

Senior Bioinformatics Scientist with
Theiagen Genomics



Andrew Hale, BSc

Bioinformatics Developer with
Theiagen Genomics

Training Objectives

Docker for Public Health Bioinformatics

- Articulate the **utility of Docker and software containerization** and their applications in public health bioinformatics;
- Describe the **advantages and challenges of using containerized software** for interoperable and reproducible pathogen genomics;
- Gain experience with the **steps involved in developing and testing containerized software** for use in public health bioinformatics;
- Gain familiarity with **community-driven resources** for public health bioinformatics, such as StaPH-B's docker container repository.

Course Agenda

Docker for Public Health Bioinformatics

Week 1 - April 01/03, 2025

- Introduction to the use of Docker and software containerization for public health bioinformatics
- Hands-on Exercise: Accessing the Course Repository, VS Code and Docker





Goals by End of Week 1

- Understand the concept of software containerization
- Understand the differences between Virtual Machines (VMs), virtual environments, and containers
- Understand the advantages and challenges of containers
- Learn the applications of containers in public health
- Learn how to utilize a container on the command line

OBJECTIVE



Why use Containers?

Why use Containers?

Ease of software install

- Avoid this →

How to get a bioinformatics headache

1. See tweet about new published tool
2. Read abstract - sounds awesome!
3. Fail to find link to source code - eventually Google it
4. Attempt to compile and install it
5. Google for 30 min for fixes
6. Finally get it built
7. Run it on tiny data set
8. Get a vague error
9. Delete and never revisit it again



Why use Containers?

One command to download and “install” a program onto your computer:

```
$ docker pull staphb/spades
```

```
Using default tag: latest
```

```
latest: Pulling from staphb/spades
```

```
b549f31133a9: Already exists
```

```
bf4358dc43e4: Pull complete
```

```
74ff0d5990f9: Pull complete
```

```
4f4fb700ef54: Pull complete
```

```
Digest:
```

```
sha256:b33f57d65cb63d631c6e3ba9b2a1c5a11ff4351475f38a1108ec61a5bf430077
```

```
Status: Downloaded newer image for staphb/spades:latest
```

```
docker.io/staphb/spades:latest
```

```
$ docker run staphb/spades spades.py --help
```

```
SPAdes genome assembler v3.15.5
```

```
Usage: spades.py [options] -o <output_dir>
```

Why use Containers?

Reproducibility

- A program may behave differently depending on how it and its dependencies are installed
 - example: SPAdes – requires python
 - Which python version is installed? 2.7, 3.4, 3.5 ?
- **Docker images are static**, minimizing the chance that the program will run differently when installed on different computers
 - (latest version) StaPH-B SPAdes docker image has python 3.8.10
- **Goal = have the software run the same exact way, every time**

Why use Containers?

Reproducibility

- Containers allow for pinning installations of specific:
 - database versions
 - software versions
- **Can help with validating assays**
 - Example: E. coli/STEC serotyping via SerotypeFinder run via a docker container

Why use Containers?

Portability

- Run software on almost ANY computer/cluster/server/HPC
- Containers are easy to share
- Solves the issue of the all too familiar phrase:

“Well it works on my computer...”

Why use Containers?

Build complex & modular workflows

- Spend less time installing software, spend more time doing science
- Nearly all workflow managers utilize & prefer containers:
 - Nextflow
 - WDL
 - Snakemake

Example WDL task, utilizes docker image for genome assembly

```
1 version 1.0
2
3 task shovill_pe {
4   input {
5     File read1_cleaned
6     File read2_cleaned
7     String sample_name
8     String docker = "quay.io/staphb/shovill:1.1.0"
9     Int min_contig_length = 200
10  }
11  command <<<
12    shovill --version | head -1 | tee VERSION
13    shovill \
14      --outdir out \
15      --R1 ~{read1_cleaned} \
16      --R2 ~{read2_cleaned} \
17      --minlen ~{min_contig_length}
18    mv out/contigs.fa out/~{sample_name}_contigs.fasta
19    mv out/contigs.gfa out/~{sample_name}_contigs.gfa
20  >>>
21  output {
22    File assembly_fasta = "out/~{sample_name}_contigs.fasta"
23    File contigs_gfa = "out/~{sample_name}_contigs.gfa"
24    String shovill_version = read_string("VERSION")
25  }
26  runtime {
27    docker: "~{docker}"
28    memory: "16 GB"
29    cpu: 4
30    disks: "local-disk 100 SSD"
31    preemptible: 0
32  }
33 }
```

Why use Containers?

Build complex & modular workflows

- Example workflows you may be familiar with:
 - **WDL**
 - [TheiaCov](#), [TheiaProk](#), [TheiaEuk](#), any WDL workflows available on [Terra.bio](#)
 - **Nextflow**
 - [Bactopia](#), [Cecret](#), [Donut Falls](#), [MycoSNP-nf](#), [PHoeNix](#), [nf-core/viralrecon](#)

```
1 version 1.0
2
3 task shovill_pe {
4   input {
5     File read1_cleaned
6     File read2_cleaned
7     String sample_name
8     String docker = "quay.io/staphb/shovill:1.1.0"
9     Int min_contig_length = 200
10  }
11  command <<<
12    shovill --version | head -1 | tee VERSION
13    shovill \
14      --outdir out \
15      --R1 ~{read1_cleaned} \
16      --R2 ~{read2_cleaned} \
17      --minlen ~{min_contig_length}
18    mv out/contigs.fa out/~{sample_name}_contigs.fasta
19    mv out/contigs.gfa out/~{sample_name}_contigs.gfa
20  >>>
21  output {
22    File assembly_fasta = "out/~{sample_name}_contigs.fasta"
23    File contigs_gfa = "out/~{sample_name}_contigs.gfa"
24    String shovill_version = read_string("VERSION")
25  }
26  runtime {
27    docker: "~{docker}"
28    memory: "16 GB"
29  }
30 }
```

All of these use containers

Why use Containers?

Summary:

- Ease of installing bioinformatics software
- Reproducibility
- Portability
- Build complex and modular bioinformatics workflows



What is a Container?

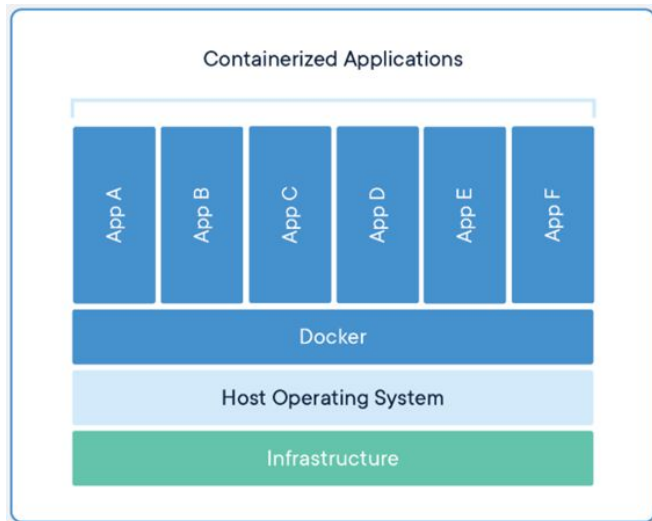
Containers

- Software container
- A standard unit of software that **packages up code and all dependencies** so the application runs quickly and **reliably from one computing environment to another**
 - <https://www.docker.com/resources/what-container/>
- Leading platforms for containers running containers
 - [Docker](#)
 - [Apptainer](#)
 - [podman](#)

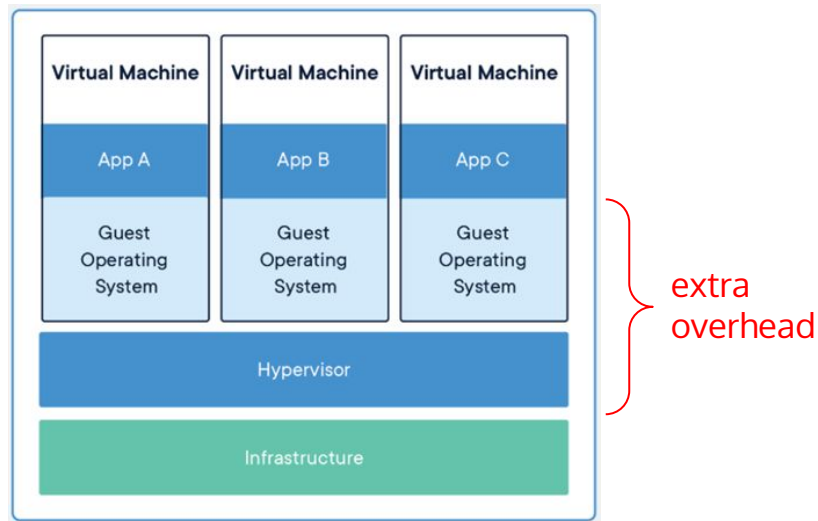


podman

Containers vs Virtual Machines



VS



- Run on existing operating system (OS) & hardware
- Can include an OS and dependencies, but in a smaller format (<1GB usually)
- Boot very fast

- VMs use more compute resources to run than containers
- Have OS, dependencies, but in a much larger format (10+ GB)
- Slow to boot

Containers vs Virtual Machines

	Containers	Virtual environments (like conda or venv)
architecture	Includes complete filesystem, and dependencies required to run software	Creates directories only containing required dependencies and python executable

Containers vs Virtual Machines

	Containers	Virtual environments (like conda or venv)
architecture	Includes complete filesystem, and dependencies required to run software	Creates directories only containing required dependencies and python executable
portability	Highly portable to any compute environment that can run docker	Usually specific to host operating system, but offers some portability with conda recipe files

Containers vs Virtual Machines

	Containers	Virtual environments (like conda or venv)
architecture	Includes complete filesystem, and dependencies required to run software	Creates directories only containing required dependencies and python executable
portability	Highly portable to any compute environment that can run docker	Usually specific to host operating system, but offers some portability with conda recipe files
isolation	Complete isolation from host system	Limited isolation from host system, host system dependencies can affect environment

Containers vs Virtual Machines

	Containers	Virtual environments (like conda or venv)
architecture	Includes complete filesystem, and dependencies required to run software	Creates directories only containing required dependencies and python executable
portability	Highly portable to any compute environment that can run docker	Usually specific to host operating system, but offers some portability with conda recipe files
isolation	Complete isolation from host system	Limited isolation from host system, host system dependencies can affect environment
resource requirements	Requires significant resources to operate (cpu, disk space, and memory). Still less than a traditional VM	Requires fewer resources and relies on host

Containers vs Virtual Machines

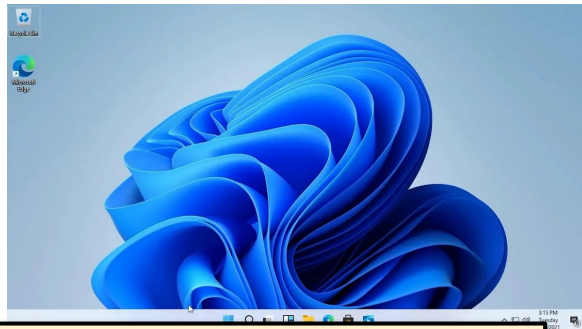
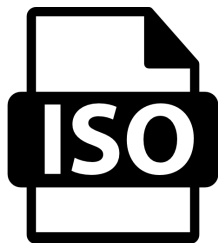
	Containers	Virtual environments (like conda or venv)
architecture	Includes complete filesystem, and dependencies required to run software	Creates directories only containing required dependencies and python executable
portability	Highly portable to any compute environment that can run docker	Usually specific to host operating system, but offers some portability with conda recipe files
isolation	Complete isolation from host system	Limited isolation from host system, host system dependencies can affect environment
resource requirements	Requires significant resources to operate (cpu, disk space, and memory). Still less than a traditional VM	Requires fewer resources and relies on host
deployment	Often used in complex, production environments	Often used for development, testing, and analysis



Docker Architecture

Docker image vs Docker container

- A docker image is a read-only template with instructions for creating a Docker container
- Similar to the file type used to install an operating system - ISO
- A docker container is the runnable instance of an image
- Containers are ephemeral (i.e. are temporary and usually deleted after use)



docker **image** is used to create the docker **container**

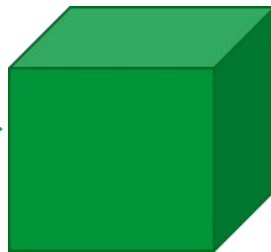
Central Dogma of Containers

Dockerfile

```
1 FROM ubuntu:xenial
2
3 # metadata
4 LABEL base.image="ubuntu:xenial"
5 LABEL version="1"
6 LABEL software="SPAdes"
7 LABEL software.version="3.13.0"
8 LABEL description="de novo DBG genome assembler"
9 LABEL website="http://cab.spbu.ru/files/release3.13.0/manual.html"
10
11 # Maintainer
12 MAINTAINER Curtis Kapsak <curtis.kapsak@state.co.us>
13
14 RUN apt-get update && apt-get install -y python \
15     wget
16
17 RUN wget http://cab.spbu.ru/files/release3.13.0/SPAdes-3.13.0-Linux.tar.gz && \
18     tar -xzf SPAdes-3.13.0-Linux.tar.gz && \
19     rm -r SPAdes-3.13.0-Linux.tar.gz && \
20     mkdir /data
21
22 ENV PATH="${PATH}:/SPAdes-3.13.0-Linux/bin"
23 WORKDIR /data
```

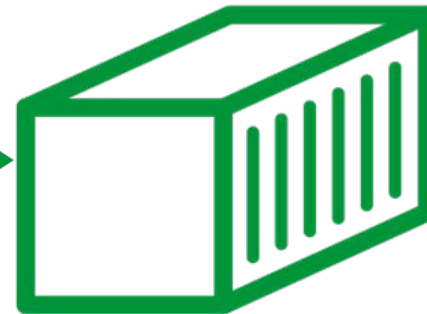
docker build

Dockerfile image

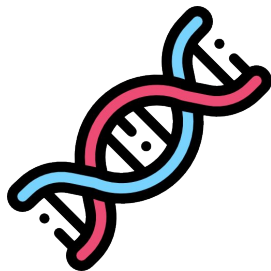


docker run

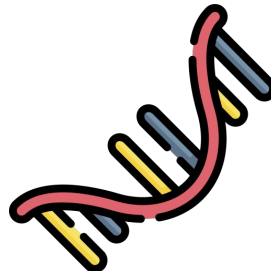
Docker container



DNA



RNA



Protein



neigen[®]
GENOMICS

The Dockerfile

- In order to build a docker image, you need at a minimum one file: **the Dockerfile**
 - Dockerfile = set of instructions used to build a docker image
 - Similar to an installation script or a .yaml file used for making/sharing conda environments

```
1 FROM ubuntu:xenial
2
3 # metadata
4 LABEL base.image="ubuntu:xenial"
5 LABEL version="1"
6 LABEL software="SPAdes"
7 LABEL software.version="3.13.0"
8 LABEL description="de novo DBG genome assembler"
9 LABEL website="http://cab.spbu.ru/files/release3.13.0/manual.html"
10
11 # Maintainer
12 MAINTAINER Curtis Kapsak <curtis.kapsak@state.co.us>
13
14 RUN apt-get update && apt-get install -y python \
15     wget
16
17 RUN wget http://cab.spbu.ru/files/release3.13.0/SPAdes-3.13.0-Linux.tar.gz && \
18     tar -xzf SPAdes-3.13.0-Linux.tar.gz && \
19     rm -r SPAdes-3.13.0-Linux.tar.gz && \
20     mkdir /data
21
22 ENV PATH="${PATH}:/SPAdes-3.13.0-Linux/bin"
23 WORKDIR /data
```

The Dockerfile

- Dockerfile instructions (**FROM**, **RUN**, **COPY**, **ENV**, etc.) will add a “layer” to the docker image
- Images are multi-layered and different images may share layers like the base image
 - **FROM** `ubuntu:focal`

Spades Dockerfile

- <https://github.com/StaPH-B/docker-builds/blob/master/spades/3.15.5/Dockerfile>

Dockerfile Cheat Sheet

- https://kapeli.com/cheat_sheets/Dockerfile.docset/Contents/Resources/Documents/index

```
1 FROM ubuntu:focal as app
2
3 # to make it easier to upgrade for new versions; ARG variables only persist during docker image build time
4 ARG SPADES_VER="3.15.5"
5
6 LABEL base.image="ubuntu:focal"
7 LABEL dockerfile.version="2"
8 LABEL software="SPAdes"
9 LABEL software.version="${SPADES_VER}"
10 LABEL description="de novo DBG genome assembler"
11 LABEL website="https://github.com/ablab/spades"
12 LABEL license="https://github.com/ablab/spades/blob/v3.15.3/assembler/LICENSE"
13 LABEL maintainer="Curtis Kapsak"
14 LABEL maintainer.email="kapsakcj@gmail.com"
15
16 # install dependencies; cleanup apt garbage
17 # python v3.8.10 is installed here; point 'python' to python3
18 RUN apt-get update && apt-get install --no-install-recommends -y python3 \
19     python3-distutils \
20     wget \
21     pigz \
22     ca-certificates && \
23     apt-get autoclean && rm -rf /var/lib/apt/lists/* && \
24     update-alternatives --install /usr/bin/python python /usr/bin/python3 10
25
26 # install SPAdes binary; make /data
27 RUN wget http://cab.spbu.ru/files/release${SPADES_VER}/SPAdes-${SPADES_VER}-Linux.tar.gz && \
28     tar -xzf SPAdes-${SPADES_VER}-Linux.tar.gz && \
29     rm -r SPAdes-${SPADES_VER}-Linux.tar.gz && \
30     mkdir /data
31
32 # set PATH and locale settings for singularity
33 ENV LC_ALL=C.UTF-8 \
34     PATH="${PATH}:/SPAdes-${SPADES_VER}-Linux/bin"
35
36 WORKDIR /data
37
38 # test layer
39 FROM app as test
40
41 # print version and run the supplied test flag
42 RUN spades.py --version && spades.py --test && spades.py --help
```

The Container

Docker Images can be built locally **or** pre-built images can be downloaded from public repositories like:

- **Docker hub:** <https://hub.docker.com/>
- **Quay.io:** <https://quay.io/>
- **GitHub container registry (GHCR):** <https://ghcr.io/>
- Cloud provider container registries:
 - GCP Artifact Registry
 - Amazon Elastic Container Registry
 - Microsoft Azure Container Registry
- Private registries are an (paid) option

Container Registries

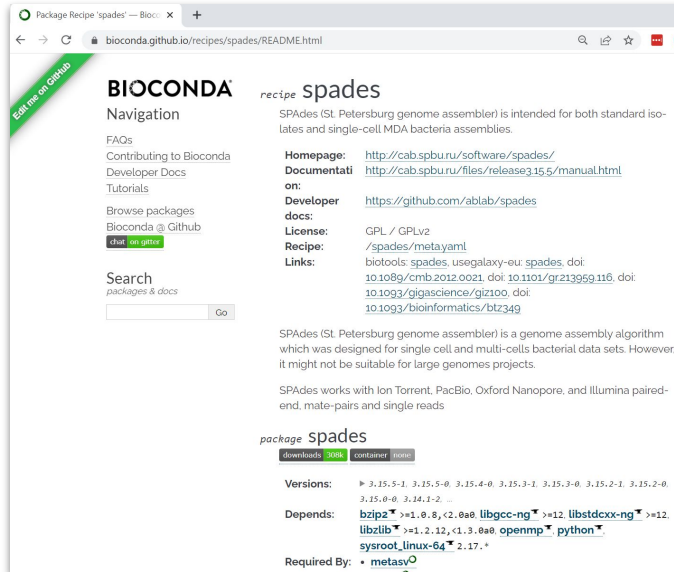
- Docker Hub: <https://hub.docker.com>
- Quay.io: <https://quay.io/>



Red Hat
Quay.io

Container Registries

- All bioconda packages are available as docker images on quay.io:
 - <https://bioconda.github.io/recipes/spades/README.html>



BIOCONDA

Navigation

- FAQs
- Contributing to Bioconda
- Developer Docs
- Tutorials
- Browse packages
- Bioconda on GitHub
- [View on GitHub](#)

Search

packages & docs

Go

recipe spades

SPAdes (St. Petersburg genome assembler) is intended for both standard isolates and single-cell MDA bacteria assemblies.

Homepage: <http://cab.spbu.ru/software/spades/>

Documentation: <http://cab.spbu.ru/files/release3.15.5/manual.html>

Developer docs: <https://github.com/ablab/spades>

License: GPL / GPLv2

Recipe: /spades/meta.yaml

Links: biotools: spades, usegalaxy-eu: spades, doi: 10.1089/cmb.2012.0021, doi: 10.1101/gr213959.116, doi: 10.1093/gigascience/gjz100, doi: 10.1093/bioinformatics/btz349

SPAdes (St. Petersburg genome assembler) is a genome assembly algorithm which was designed for single cell and multi-cells bacterial data sets. However, it might not be suitable for large genomes projects.

SPAdes works with Ion Torrent, PacBio, Oxford Nanopore, and Illumina paired-end, mate-pairs and single reads

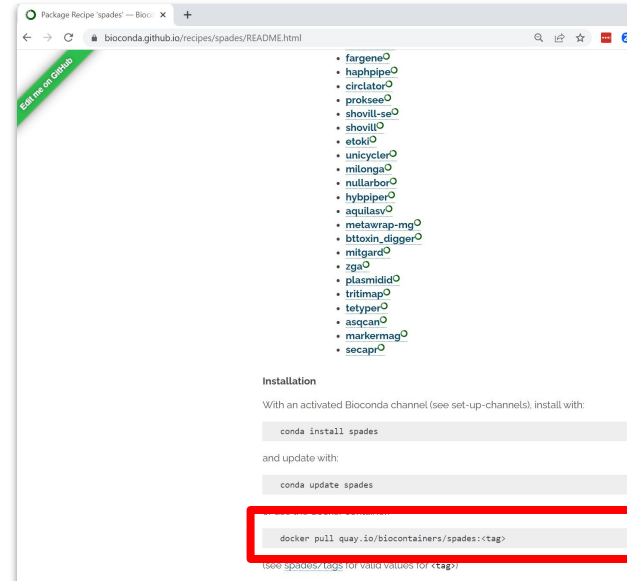
package spades

downloads 306k container recipe

Versions: 3.15.5-1, 3.15.5-0, 3.15.4-0, 3.15.3-1, 3.15.3-0, 3.15.2-1, 3.15.2-0, 3.15.0-0, 3.14.1-2, ...

Depends: `bzip2 >=1.0.8, <2.0a0 libgcc-ng >=12 libstdc++-ng >=12, libzlib >=1.2.12, <1.3.0a0 openmp*, python*, sysroot_linux-64* 2.17.*`

Required By: metaspv



Installation

With an activated Bioconda channel (see set-up-channels), install with:

```
conda install spades
```

and update with:

```
conda update spades
```

Docker

```
docker pull quay.io/biocontainers/spades:<tag>
```

(see spadesr/logs for valid values for <tag>)

Container Registries

- Sometimes bioinfo tool developers publish their own docker images
 - Bakta - <https://github.com/oschwengers/bakta>
 - NCBI AMRFinderPlus - <https://github.com/ncbi/amr/wiki/Installing-AMRFinder>

Installation

Bakta can be installed via BioConda, Docker, Singularity and Pip. However, we encourage to use [Conda](#) or [Docker/Singularity](#) to automatically install all required 3rd party dependencies.

In all cases a mandatory [database](#) must be downloaded.

BioConda

```
conda install -c conda-forge -c bioconda bakta
```

Docker

```
sudo docker pull oschwengers/bakta
sudo docker run oschwengers/bakta --help
```

Installation instructions and get-started guides: [Docker docs](#)

For further convenience, we provide a shell script (`bakta-docker.sh`) handling Docker related parameters (volume mounting, user IDs, etc):

```
bakta-docker.sh --db <db-path> --output <output-path> <input>
```

Installing AMRFinder

Arjun Prasad edited this page on Jan 12 · 21 revisions

Installation

AMRFinderPlus requires HMMER, BLAST+, Linux, and perl. We provide a [bioconda](#) package, Linux binaries, and the source code is available to compile AMRFinderPlus yourself though we haven't extensively tested compiling AMRFinderPlus on other systems and aren't supporting non-Linux systems at this time.

How to install

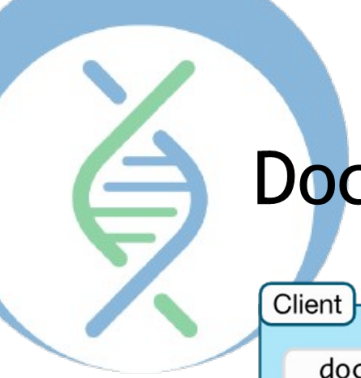
We distribute AMRFinderPlus, through [Bioconda](#), as a [linux x86 binary](#), as a [docker image](#), and as [source code](#) through [GitHub](#).

Bioconda

Here are links to instructions for three methods of installation. The simplest, and recommended method is to [install AMRFinderPlus](#) and all of the prerequisites with [bioconda](#). See [Install with bioconda](#).

Docker

A DockerHub image [NCBI/amr](#) is provided, and should be downloadable using `docker pull ncbi/amr`. The build instructions for the docker image are [available on GitHub](#). See the [docker README](#) for details and some examples of how to use it.



Docker architecture

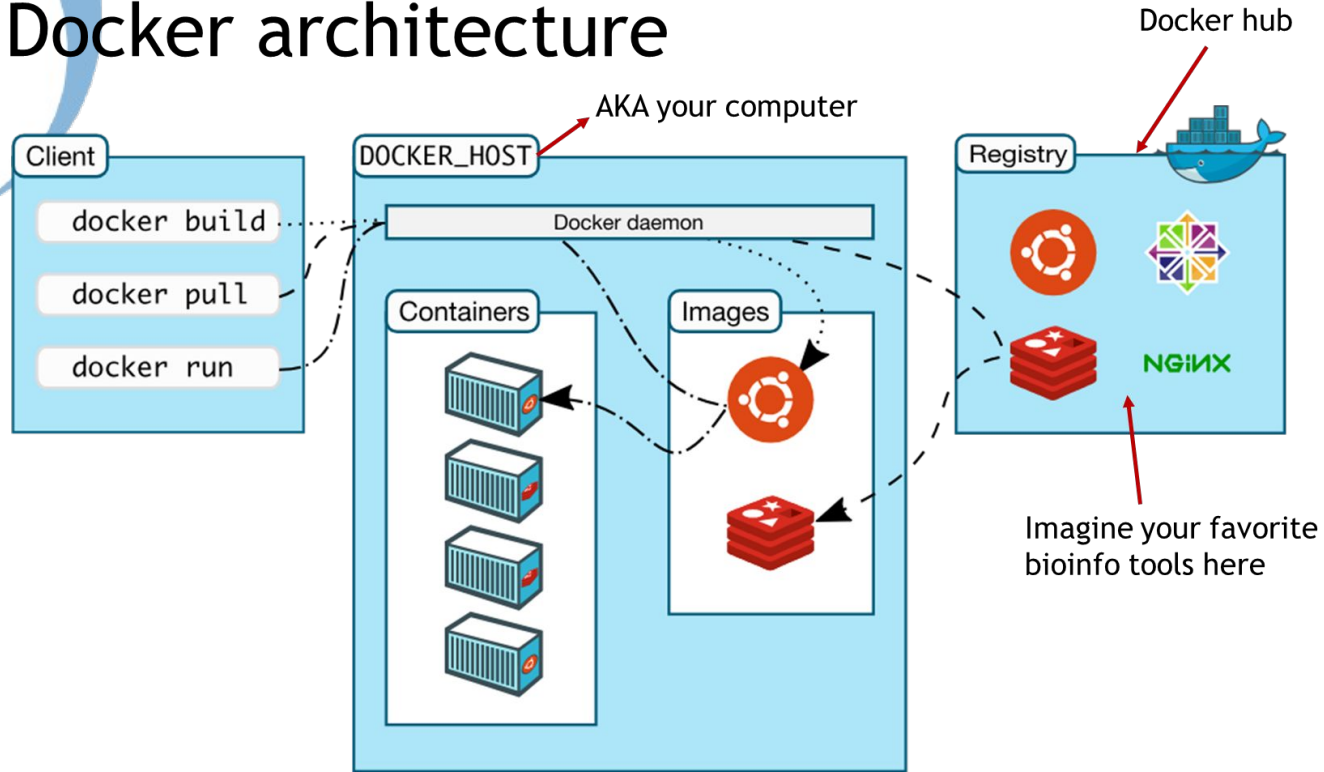


Image from: <https://docs.docker.com/engine/docker-overview/>

Docker Architecture

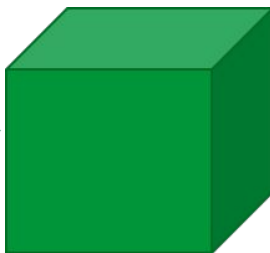
Summary:

- **Dockerfile** is used to create the docker **image**
- Docker **image** is used to create the docker **container**
 - Container is the runnable instance of an image

Dockerfile

```
1 FROM ubuntu:xenial
2
3 # metadata
4 LABEL base.image="ubuntu:xenial"
5 LABEL version="1"
6 LABEL software="SPAdes"
7 LABEL software.version="3.13.0"
8 LABEL description="de novo DBG genome assembler"
9 LABEL website="http://cab.spbu.ru/files/release3.13.0/manual.html"
10
11 # Maintainer
12 MAINTAINER Curtis Kapsak <curtis.kapsak@state.co.us>
13
14 RUN apt-get update && apt-get install -y python \
15     wget
16
17 RUN wget http://cab.spbu.ru/files/release3.13.0/SPAdes-3.13.0-linux.tar.gz && \
18     tar -xzf SPAdes-3.13.0-linux.tar.gz && \
19     rm -r SPAdes-3.13.0-linux.tar.gz && \
20     mkdir /data
21
22 ENV PATH="$PATH:/SPAdes-3.13.0-Linux/bin"
23
24 WORKDIR /data
```

Dockerfile image



docker build

docker run

Docker container



Docker Architecture

Summary:

- Where can I get docker images?
 - Container registries - docker hub, quay.io, github container registry, cloud provider registries (GCP, AWS, Azure)
- Container = A standard unit of software that **packages up code and all dependencies** so the application runs quickly and **reliably from one computing environment to another**



Docker on the Command Line

Docker on the Command Line

Common docker commands:

- **docker pull** - downloads an image from a repository
- **docker run** - run a command in a container
- **docker build** - build an image from a Dockerfile
- **docker images** - list all available images
- **docker system prune** - delete unused images and containers
- **docker rmi** - remove an image
- **docker inspect** - get information about an image

Docker CLI Cheat Sheet:

https://docs.docker.com/get-started/docker_cheatsheet.pdf

Further Reading & Resources

- Docker Documentation - a wealth of info here. Note that we use Docker Community Edition, as you have to pay for the Enterprise Edition
 - <https://docs.docker.com/>
- An awesome tutorial/workshop on docker for bioinformatics
 - <https://github.com/PawseySC/bio-workshop-18>
- Template for your Dockerfile
- <https://github.com/StaPH-B/docker-builds/blob/master/dockerfile-template/Dockerfile>
- Some best practices
 - https://staphb.org/docker-builds/make_containers/



Hands-On Exercise

Exercise 01: Accessing the Course Repository, VS Code and Docker

Exercise Goal:

- Access development environment via GitPod
- Use VS Code to most commonly used Docker commands
- Use a docker container to download a *Klebsiella pneumoniae* genome FASTA file from [NCBI](#)
- Use a docker container to run [kleborate](#) on FASTA file for subtyping, serotyping, virulence and A.M.R. prediction



Theiagen[®]
G E N O M I C S

www.theiagen.com
support@theiagen.com