

# Container

Thomas Juettemann, EMBL-EBI  
aqua\_faang\_course@ebi.ac.uk

May 10-12, 2021

## 1 Introduction

The aim of this tutorial is to help you understand how container and docker works. While bioinformatic tools are used for some examples, they are not the main focus. Two versions of the document exist, `container_solutions` & `container_no_solutions`. Use the solutions at any time you are stuck, or of course to check if your answer is correct. Finishing everything in the time given is not important, focus on understanding each step, and take the time play around!

## 2 Docker Hub

Docker Hub is the world's largest repository of container images, at the time of this course it hosts 8.3 million repositories. In this section we explore a couple of repositories and run our first container.

### 2.1 Hello World

Because that is how every IT adventure begins.

#### Task

1. On Dockerhub, find the official **hello-world** image and click on it.
2. In the description page, find the **docker run** command.
3. Execute in a terminal window.

### Solution

## 2.2 Versions

Dockerhub offers to host several versions of a software. These versions are called tags Here we take a quick look at the similarity and differences using PostgreSQL as an example.

### Task

1. On Dockerhub find the official **postgres** image and click on it.
2. Which is the latest version?
3. What happens when you click on *latest*?
4. What is the difference between the first two tags in the first line?
5. What is the major difference (hint: *FROM*) between the first tag in line 1 and the first tag in line 2?

### Solution

## 3 BioContainers

BioContainers is a community-driven project that provides the infrastructure and basic guidelines to create, manage and distribute bioinformatics packages (e.g conda) and containers (e.g docker, singularity). Many tools used by computational biologists on a daily basis are readily available. A mayor difference to Docker Hub is that users can request container to be created.

### 3.1 Exploring

In this section we will search for a specific tool on BioContainers, find the newest version, download the image and run it locally in different ways. The tool of choice is BLAST, and we will use it to create a blast database and then run an alignment.

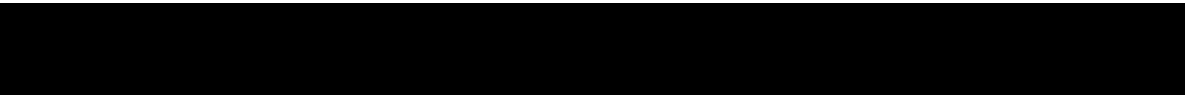
#### 3.1.1 Finding an image

As a community driven project, BioContainers does not have the notion of an "official image", and therefore no such filter exists. Those images are usually used (*=pulled*) much more often than others. We exploit that fact to find the "official" NCBI version of BLAST.

#### Task

1. Got to <https://biocontainers.pro/>
2. Click on Registry (wait for the search interface to load)
3. Search for **blast**
4. Below the search bar, change the sorting from "**Default**" to "**Pull No**" and the sort order from "**Asc**" to "**Desc**".
5. The NCBI image is the first or second hit, click on it

#### Solution



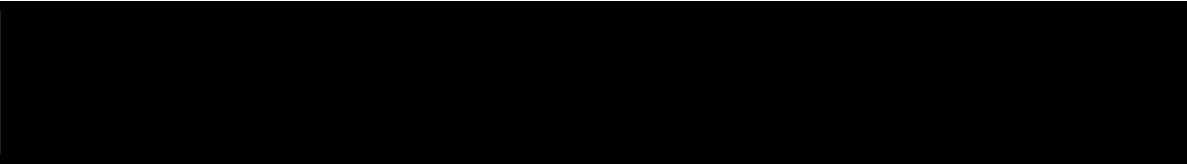
#### 3.1.2 Pulling a specific version

The landing page contains generic information about the tool and different way of utilising it: Docker, Singularity and Conda. The version used in the various commands is not necessarily the most recent one.

### Task

1. Click on "Packages and Containers"
2. Find the version **quay.io/biocontainers/blast:2.11.0-pl526he19e7b1\_0**
3. Copy the **docker pull** command
4. Open a terminal
5. Paste and execute the cached **docker pull** command
6. Store (**=export**) the image name in an environment variable called *blast\_c*
7. List all images

### Solution



### 3.1.3 Inspecting an image

Time to explore the image! Docker provides low-level information (by default in JSON format) for all objects. Finding out what happened to a failed container or getting information about an image is done by inspecting the object.

### Task

1. *Inspect* the image
2. In the **Config** section, what are the values for
  - (a) **WorkingDir**
  - (b) **Entrypoint**
3. Which command (**Cmd**) is executed when the container is run?

### Solution

#### 3.1.4 From image to container

A docker image is build from the instructions in a dockerfile. The image is then used as a template (or base), which we can copy and use to run an application. The application needs an isolated environment in which to run: a container. From the BLAST image that we downloaded, we will create a container, and using that container we will create a BLAST database, and then query this database using blastp.

makeblastdb takes 2 arguments:

1. **-in** *path/to/input/fasta*
2. **-dbtype** *prot or nucl*
3. Documentation: <https://ncbi.github.io/magicblast/cook/blastdb.html>

Keep in mind that docker can not access any files in our local filesystem, unless it is mounted. Remember that mount takes 3 arguments:

1. type (usually bind)
2. source (absolute path to the directory outside the container)
3. target (absolute path to the directory inside the container )
4. Documentation (bind mounts): <https://docs.docker.com/storage/bind-mounts/>

### Task

1. Using the BLAST container (`=blast_c`), create a BLAST database (**makeblastdb**).
2. Sequence in FASTA format is available at **\$HOME/train-aquafaang-bioinf/container/data/zebrafish.1.protein.faa**

### Solution

## Container tutorial

---

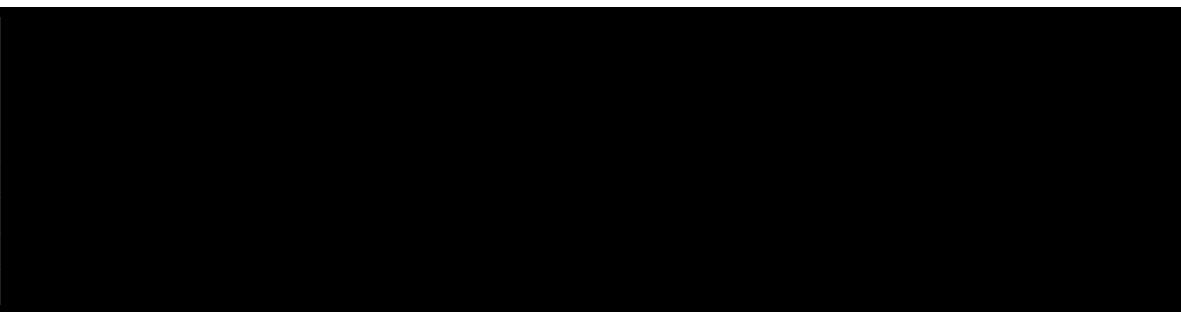
### 3.1.5 Life inside a container

To gain a better understanding of a container, we will run the actual alignment inside the container. We also want to know which user you are inside a container, and where all executables are located. If you are not sure about the parameters, remember that docker offers help for each sub command (e.g. *docker run --help* ).

#### Task

1. **run** the BLAST container **interactive**, using a **pseudo tty** and **bash** as **command**.
2. Inside the container, run **whoami** and **which blastp**
3. List the contents of /mnt

#### Solution



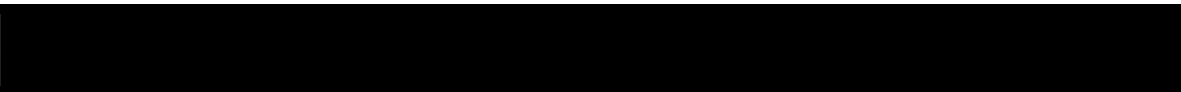
#### Mounting

Mounting local storage into a container is essential for any type of data analysis and will be used frequently in this course.

#### Task

1. Run the same command as before, but also **mount** the **container/data** directory to **/mnt**

#### Solution





## Container tutorial

---

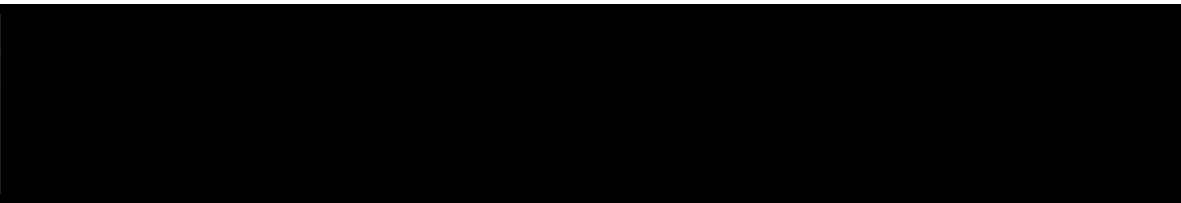
### Explore a running container

Another way of understanding container is exploring a running container. Here we have a look at the state and the inside of it.

#### Task

1. Open a second terminal
2. List all running container
3. Copy the container ID of the BLAST container
4. Inspect it
5. What is the **Cmd** in **Config**?
6. What do you see in the **Mounts** section?

#### Solution



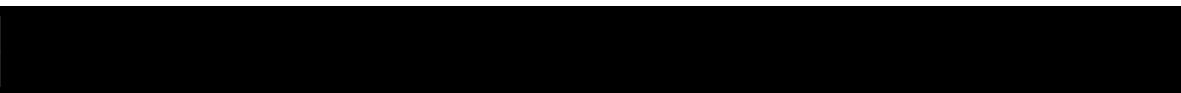
### Execute a command in a running container

The *docker exec* command runs a new command in a running container. Using *exec* on the running container, execute a `ls -l` and a `ls -l /mnt`. What do you see?

#### Task

1. Using the container ID, execute a `ls -l`
2. Using the container ID, execute a `ls -l /mnt`

#### Solution



## Container tutorial

---

### Alignment

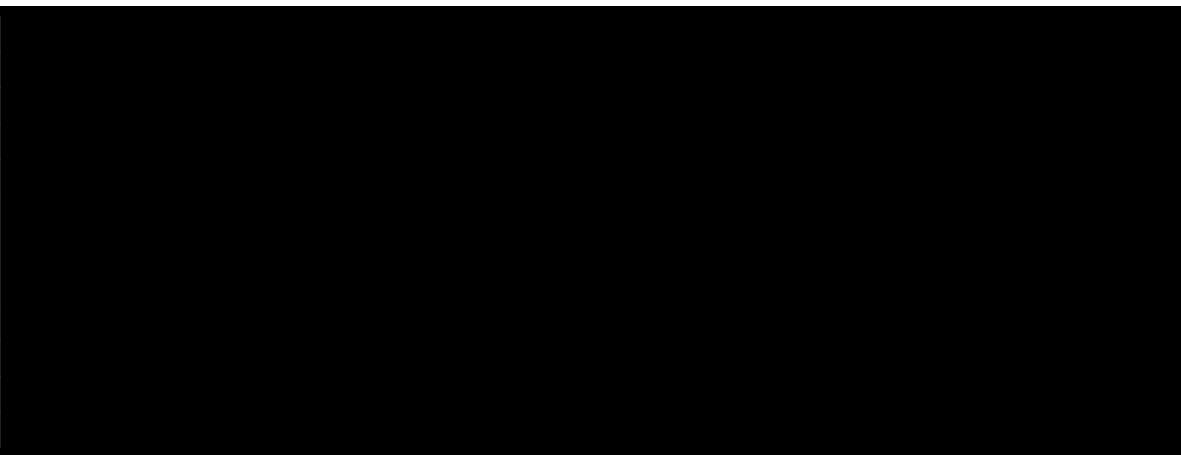
There are 3 ways of running the alignment:

1. Using the container the same way as we did with *makeblastdb* (swap *makeblastdb* for *blastp*). This is the most common solution in a workflow / pipeline or ad-hoc situation.
2. Using *docker exec* while the first container is still running. This creates a basic version of a BLAST server that can be accessible to other users on the same server.
3. Running *blastp* inside the first container. Mostly used for debugging or development.

### Task

1. Run the alignment in all 3 different ways.
2. Use different names for the out file
3. To make the mounting short, create a environment variable:  
**export docker\_mount="type=bind,source=\$HOME/train-aquafaang-bioinf/container/data,target=/mnt"**
4. Exit the running container

### Solution



## 4 Cleaning up

In the this section we will be looking at how to stop rouge containers and how to get rid of unused containers and images.





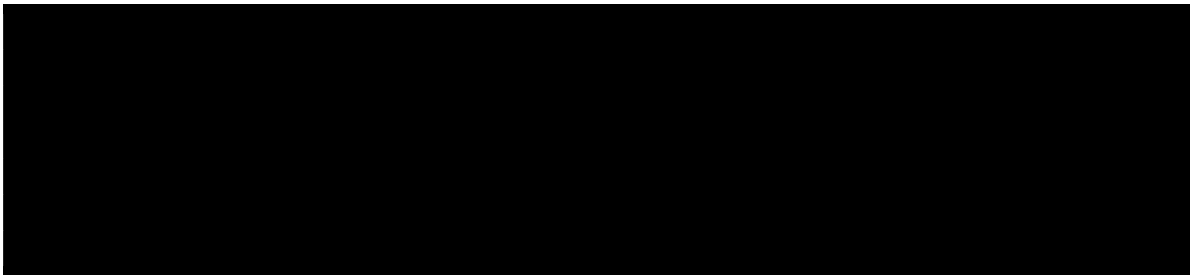
### 4.1 Stop or get killed

Containers are meant to be replaceable. If a container becomes inactive, it is easy to replace it with a new one. But how to get rid of a hanging container? Docker has two ways of getting rid of unwanted containers: stop (kill -15, then kill -9) or kill (kill -9). More information can be found at <https://docs.docker.com/engine/reference/commandline/docker/>

#### Task

1. Download the latest alpine image
2. Run an alpine container using the -i option, nothing else
3. In a second terminal, list all running container
4. Copy the <container ID> of the alpine container
5. Stop or Kill the container using the container ID

#### Solution



### 4.2 Spring cleaning

Docker takes a conservative approach to cleaning up unused objects (often referred to as “garbage collection”), such as images, containers, volumes, and networks. These objects are generally not removed unless you explicitly ask Docker to do so. This can cause Docker to use extra disk space. For each type of object, Docker provides a prune command. In addition, you can use docker system prune to clean up multiple types of objects at once.

#### Task

1. Prune images
2. Prune all images which are not associated to at least one container
3. Prune container
4. Prune everything

## Solution



## 5 Create your own

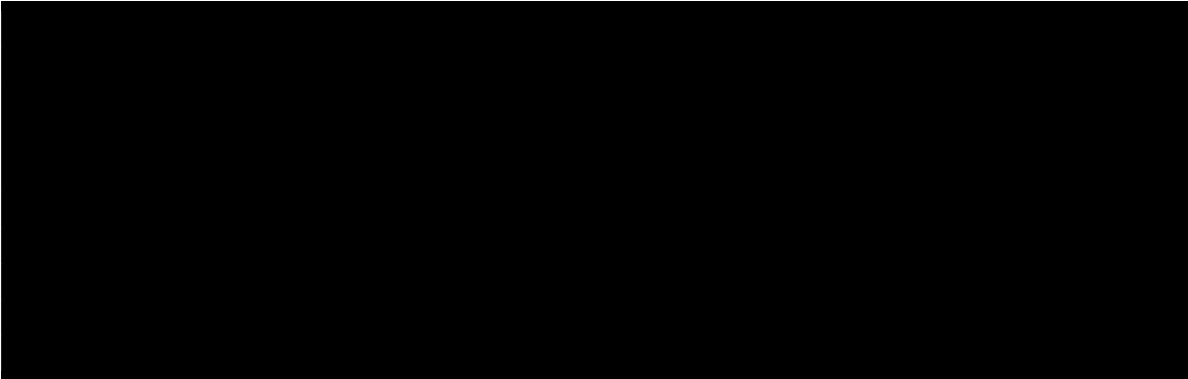
Lastly, a peak into writing your own Dockerfile. This exercise ties everything together and will require you to do some reading to understand the different components. Dockerfiles are sets of instructions, and it always in the format ***INSTRUCTION arguments***. Instructions used in this tutorial are **FROM**, **WORKDIR**, **COPY**, **ENTRYPOINT**, **CMD**. Look at <https://docs.docker.com/engine/reference/builder/> for more information about each.

Each dockerfile starts with a base image, indicated by FROM, e.g. **FROM ubuntu:latest**. **COPY** is used to copy local files into the image, e.g. **COPY ./requirements.txt .**, which copies requirements.txt to the current (root) directory of the image. **WORKDIR** is the equivalent of the unix *cd*. **ENTRYPOINT** indicates the command that should be executed everytime the starts. **CMD** is used to provide defaults for an executing container.

### Task

1. CD into **\$HOME/train-aquafaang-bioinf/container/hello\_world**
2. Using your favourite editor, create a file called **Dockerfile**
3. Use python 3 as your base image
4. Change the working directory in the image to **/usr/src/app**
5. Copy **\$HOME/train-aquafaang-bioinf/container/data/hello\_world.py** into the current working dir (".")
6. Set the entrypoint to python
7. Set the default argument to **hello\_world.py**
8. **Build** your container, using the **tag your\_name/hello\_world:latest**
9. **Run** your container
10. Documentation: Go to dockerhub and find the official python image

### Solution



### Task

You can use the Python image that you downloaded also to execute scripts that are not part of your container. This can be really helpful if you need a certain version of a software that you can install in an environment.

1. Create a container from the python 3 base image and execute the `hello_world.py`

### Solution

