

How to Deploy Hugging Face Models in a Docker Container

In this short tutorial, we will explore how [Hugging Face](https://huggingface.co/) models can be deployed in a [Docker](https://docker.com) Container and exposed as a web service endpoint.

The service it exposes is a translation service from English to French and French to English.

Why someone would like to do that? Other than to learn about those specific technologies, it is a very convenient way to try and test the thousands of models that exists on Hugging Face, in a clean and isolated environment that can easily be replicated, shared or deployed elsewhere than on your local computer.

In this tutorial, you will learn how to use `Docker` to create a container with all the necessary code and artifacts to load Hugging Face models and to expose them as web service endpoints using `Flask`.

All code and configurations used to write this blog post are available in this [GitHub Repository](https://github.com/fgiasson/en-fr-translation-service) (<https://github.com/fgiasson/en-fr-translation-service>). You simply have to clone it and to run the commands listed in this tutorial to replicate the service on your local machine.

Installing Docker

The first step is to install Docker. The easiest way is by simply installing [Docker Desktop](https://www.docker.com/products/docker-desktop/) (<https://www.docker.com/products/docker-desktop/>) which is available on MacOS, Windows and Linux.

Creating the Dockerfile

The next step is to create a new `Git` repository where you will create a `Dockerfile`. The Dockerfile is where all instructions are written that tells Docker how to create the container.

I would also strongly encourage you to install and use [hadolint](https://github.com/hadolint/hadolint) (<https://github.com/hadolint/hadolint>), which is a really good Docker linter that helps people to follow Docker best practices. There is also a [plugin for VS Code](https://marketplace.visualstudio.com/items?itemName=exiasr.hadolint) (<https://marketplace.visualstudio.com/items?itemName=exiasr.hadolint>) if this is what you use as your development IDE.

Base image and key installs

The first thing you define in a Dockerfile is the base image to use to initialize the container. For this

tutorial, we will use Ubuntu's latest LTS:

```
1 | # Use Ubuntu's current LTS
2 | FROM ubuntu:jammy-20230804
```

Since we are working to create a Python web service that expose the predictions of a ML model, the next step is to add they key pieces required for the Python service. Let's make sure that you only include what is necessary to minimize the size, and complexity, of the container as much as possible:

```
1 | # Make sure to not install recommends and to clean the
2 | # install to minimize the size of the container as much as possible.
3 | RUN apt-get update && \
4 |     apt-get install --no-install-recommends -y python3=3.10.6-1~22.04 && \
5 |     apt-get install --no-install-recommends -y python3-pip=22.0.2+dfsg-1ubuntu0.3 && \
6 |     apt-get install --no-install-recommends -y python3-venv=3.10.6-1~22.04 && \
7 |     apt-get clean && \
8 |     rm -rf /var/lib/apt/lists/*
```

This instruct Docker to install `Python3`, `pip` and `venv`. It also ensures that apt get cleaned of cached files, that nothing more is installed and that we define the exact version of the package we want to install. That is to ensure that we minimize the size of the container, while making sure that the container can easily be reproduced, with the exact same codebase, any time in the future.

Another thing to note: we run multiple commands with a single `RUN` instruction by piping them together with `&&`. This is to minimize the number of layers created by Docker for the container, and this is a best practice to follow when creating containers. If you don't do this and run `hadolint`, then you will get warning suggesting you to refactor your `Dockerfile` accordingly.

Copy required files

Now that the base operating system is installed, the next step is to install all the requirements of the Python project we want to deploy in the container:

```
1 | # Set the working directory within the container
2 | WORKDIR /app
3 |
4 | # Copy necessary files to the container
5 | COPY requirements.txt .
6 | COPY main.py .
7 | COPY download_models.py .
```

First we define the working directory with the `WORKDIR` instruction. From now on, every other instruction will run from that directory in the container. We copy the local files: `requirements.txt`, `main.py` and `download_models.py` to the working directory.

Create virtual environment

Before doing anything with those files, we are better creating a virtual environment where to install all those dependencies. Some people may wonder why we create an environment within an environment? It is further isolation between the container and the Python application to make sure that there is no possibility of dependencies clashes. This is a good best practice to adopt.

```
1 | # Create a virtual environment in the container
2 | RUN python3 -m venv .venv
3 |
4 | # Activate the virtual environment
5 | ENV PATH="/app/.venv/bin:$PATH"
```

Install application requirements

Once the virtual environment is created and activated in the container, the next step is to install all the required dependencies in that new environment:

```
1 | # Install Python dependencies from the requirements file
2 | RUN pip install --no-cache-dir -r requirements.txt && \
3 | # Get the models from Hugging Face to bake into the container
4 | python3 download_models.py
```

It runs `pip install` to install all the dependencies listed in `requirements.txt`. The dependencies are:

```
1 transformers==4.30.2
2 flask==2.3.3
3 torch==2.0.1
4 sacremoses==0.0.53
5 sentencepiece==0.1.99
```

Just like the Ubuntu package version, we should (have to!) pin (specify) the exact version of each dependency. This is the best way to ensure that we can reproduce this environment any time in the future and to prevent unexpected crashes because code changed in some downstream dependencies that causes issues with the code.

Downloading all models in the container

As you can see in the previous `RUN` command, the next step is to download all models and tokenizers in the working directory such that we bake the model's artifacts directly in the container. That will ensure that we minimize the time it takes to initialize a container. We spend the time to download all those artifacts at build time instead of run time. The downside is that the containers will be much bigger depending on the models that are required.

The `download_models.py` file is a utility file used to download the Hugging Face models used by the service directly into the container. The code simply download the models and tokenizer files from Hugging Face and save them locally (in the working directory of the container):

```

1  from transformers import AutoTokenizer, AutoModelForSeq2SeqLM
2  import os
3
4  def download_model(model_path, model_name):
5      """Download a Hugging Face model and tokenizer to the specified directory"""
6      # Check if the directory already exists
7      if not os.path.exists(model_path):
8          # Create the directory
9          os.makedirs(model_path)
10
11     tokenizer = AutoTokenizer.from_pretrained(model_name)
12     model = AutoModelForSeq2SeqLM.from_pretrained(model_name)
13
14     # Save the model and tokenizer to the specified directory
15     model.save_pretrained(model_path)
16     tokenizer.save_pretrained(model_path)
17
18     # For this demo, download the English-French and French-English models
19     download_model('models/en_fr/', 'Helsinki-NLP/opus-mt-en-fr')
20     download_model('models/fr_en/', 'Helsinki-NLP/opus-mt-fr-en')

```

Creating the Flask translation web service endpoint

The last thing we have to do with the Dockerfile is to expose the port where the web service will be available and to tell the container what to run when it starts:

```

1  # Make port 6000 available to the world outside this container
2  EXPOSE 6000
3
4  ENTRYPOINT [ "python3" ]
5
6  # Run main.py when the container launches
7  CMD [ "main.py" ]

```

We expose the port `6000` to the outside world, and we tell Docker to run the `python3` command with `main.py`. The `main.py` file is a very simple file that register the web service's path using Flask, and that makes the predictions (translations in this case):

```

1  from flask import Flask, request, jsonify
2  from transformers import AutoTokenizer, AutoModelForSeq2SeqLM
3
4  def get_model(model_path):
5      """Load a Hugging Face model and tokenizer from the specified directory"""
6      tokenizer = AutoTokenizer.from_pretrained(model_path)
7      model = AutoModelForSeq2SeqLM.from_pretrained(model_path)
8      return model, tokenizer
9
10 # Load the models and tokenizers for each supported language
11 en_fr_model, en_fr_tokenizer = get_model('models/en_fr/')
12 fr_en_model, fr_en_tokenizer = get_model('models/fr_en/')
13
14 app = Flask(__name__)
15
16 def is_translation_supported(from_lang, to_lang):
17     """Check if the specified translation is supported"""
18     supported_translations = ['en_fr', 'fr_en']
19     return f'{from_lang}_{to_lang}' in supported_translations
20
21 @app.route('/translate/<from_lang>/<to_lang>', methods=['POST'])
22 def translate_endpoint(from_lang, to_lang):
23     """Translate text from one language to another. This function is
24     called when a POST request is sent to /translate/<from_lang>/<to_lang>"""
25     if not is_translation_supported(from_lang, to_lang):
26         return jsonify({'error': 'Translation not supported'}), 400
27
28     data = request.get_json()
29     from_text = data.get(f'{from_lang}_text', '')
30
31     if from_text:
32         model = None
33         tokenizer = None
34
35         match from_lang:
36             case 'en':
37                 model = en_fr_model
38                 tokenizer = en_fr_tokenizer
39             case 'fr':
40                 model = fr_en_model
41                 tokenizer = fr_en_tokenizer
42
43         to_text = tokenizer.decode(model.generate(tokenizer.encode(from_text, return_tensors='|
44
45         return jsonify({'f'{to_lang}_text': to_text})
46     else:
47         return jsonify({'error': 'Text to translate not provided'}), 400
48
49 if __name__ == '__main__':
50     app.run(host='0.0.0.0', port=6000, debug=True)

```

Building the container

Now that the `Dockerfile` is completed, the next step is to use it to have Docker to build the actual image of the container. This is done using this command in the terminal:

Bash

```
1 | docker build -t localbuild:en_fr_translation_service .
```

Note that we specified a tag to make it easier to manage it in between all the other images that may exists in the environment. The output of the terminal will show every step defined in the `Dockerfile`, and the processing for each of those step. The final output looks like:

```
docker build -t localbuild:en_fr_translation_service .
[+] Building 0.8s (13/13) FINISHED
=> [internal] load build definition from Dockerfile
=> => transferring dockerfile: 1.27kB
=> [internal] load .dockerignore
=> => transferring context: 2B
=> [internal] load metadata for docker.io/library/ubuntu:jammy-20230804
=> [1/8] FROM docker.io/library/ubuntu:jammy-20230804@sha256:ec050c32e4a6085b423d36ecd025c0d3ff00c38ab93a3d71a460ff1c44fa6d77
=> [internal] load build context
=> => transferring context: 213B
=> CACHED [2/8] RUN apt-get update && apt-get install --no-install-recommends -y python3=3.10.6-1~22.04 && apt-get install
=> CACHED [3/8] WORKDIR /app
=> CACHED [4/8] COPY requirements.txt .
=> CACHED [5/8] COPY main.py .
=> CACHED [6/8] COPY download_models.py .
=> CACHED [7/8] RUN python3 -m venv .venv
=> CACHED [8/8] RUN pip install --no-cache-dir -r requirements.txt && python3 download_models.py
=> exporting to image
=> => exporting layers
=> => writing image sha256:50cbc55e8070cf83ba056cbdb248b74fc4a3f0e3955a07712f48eb4909a9b2b8
=> => naming to docker.io/library/localbuild:en_fr_translation_service

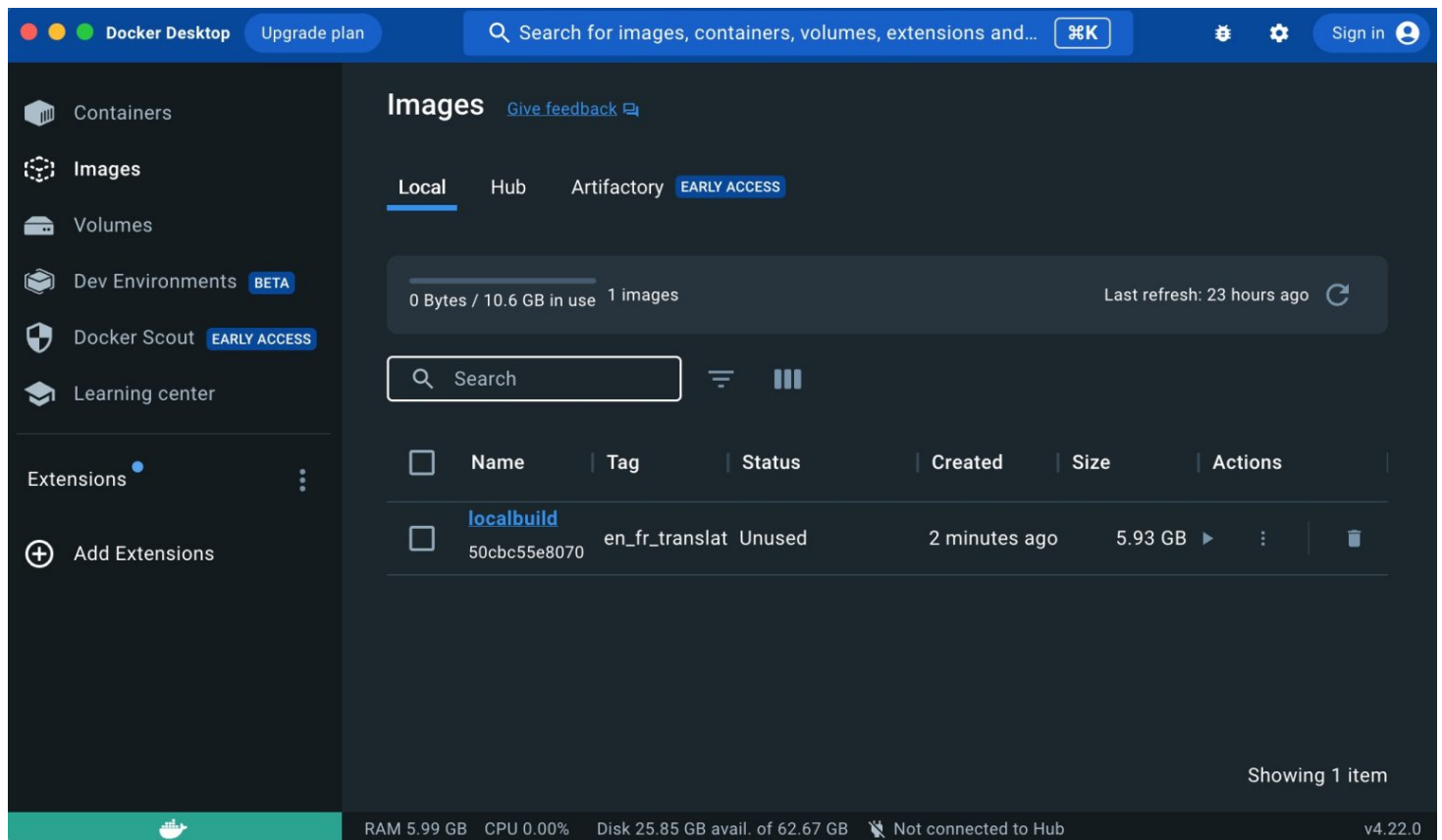
What's Next?
View summary of image vulnerabilities and recommendations → docker scout quickview
```

(https://fgiasson.com/blog/wp-content/uploads/2023/08/build_output.jpg)

Running and Querying the service

Now that we have a brand new image, the next step is to test it. In this section, I will use Docker Desktop's user interface to show how we can easily do this, but all those step can easily be done (and automated) using the `docker` command line application.

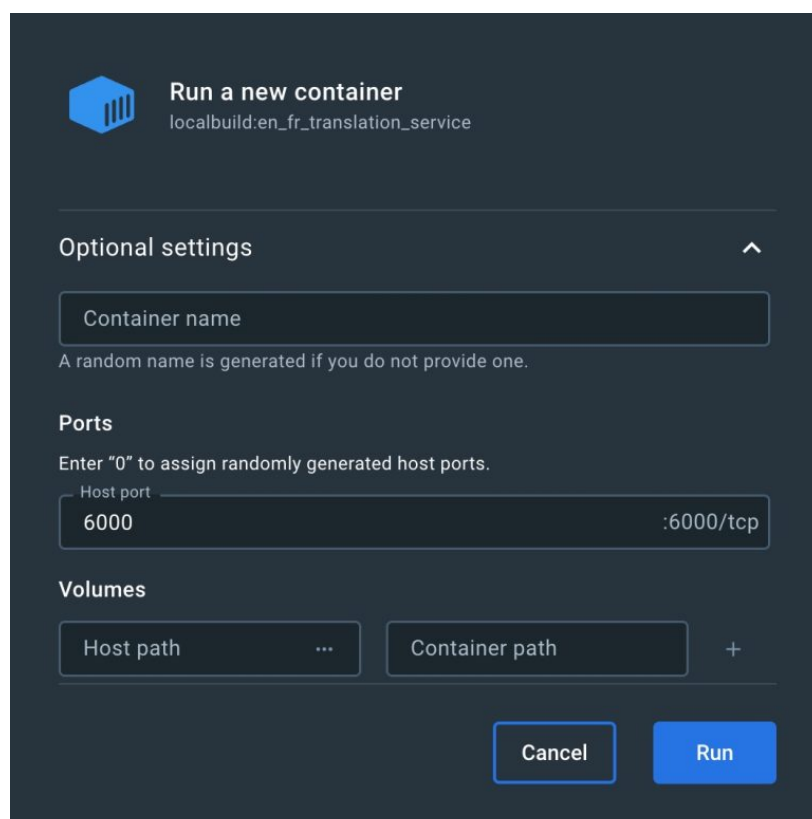
After you built the image, it will automatically appear in the `images` section of Docker Desktop:



(<https://fgiasson.com/blog/wp-content/uploads/2023/08/images.jpg>).

You can see the tag of the image, its size, when it was created, etc. To start the container from that image, we simply have to click the **play arrow** in the **Actions** column. That will start running a new container using that image.

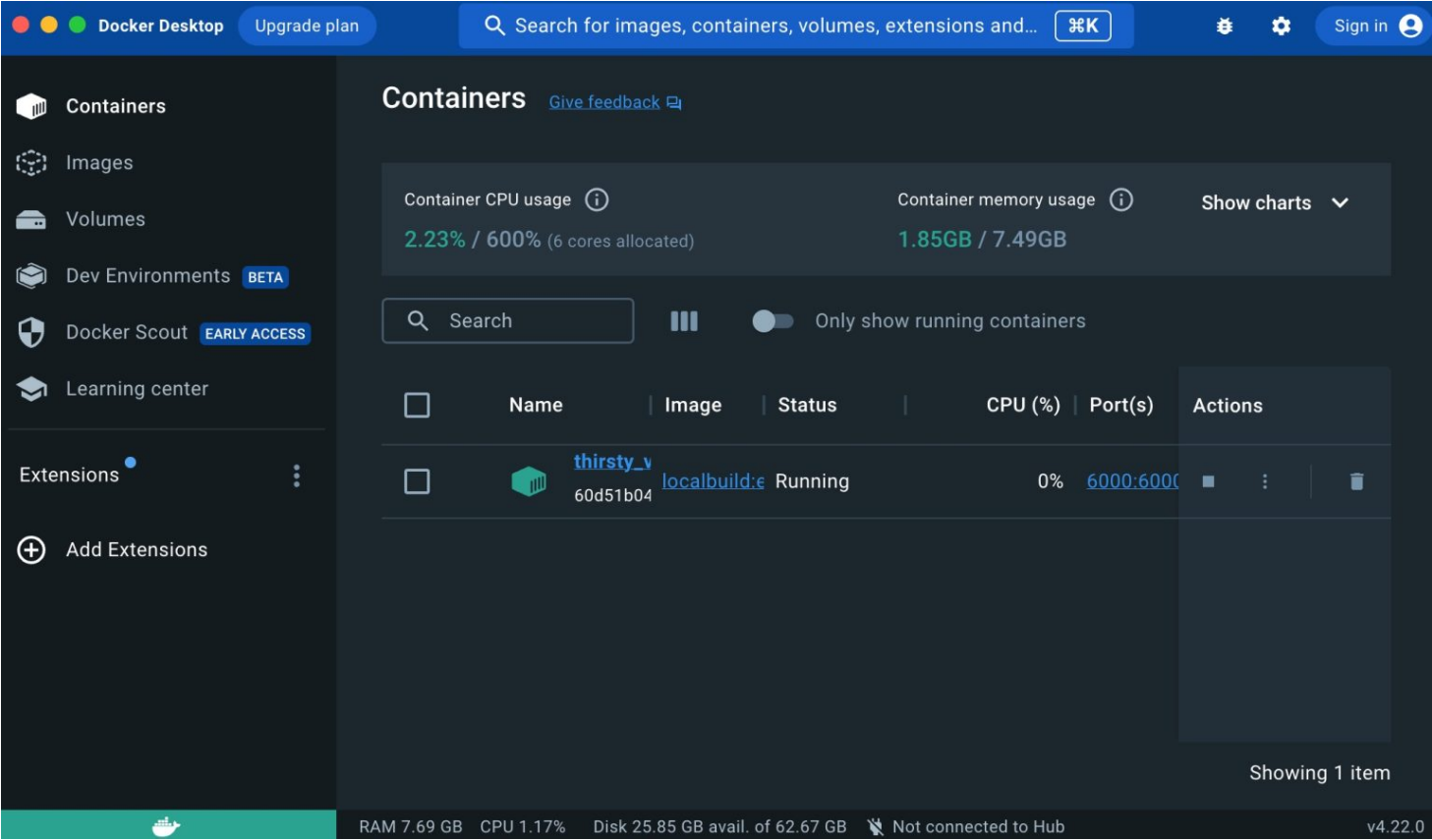
Docker Desktop will enable you to add some more parameter to start the container with the following window:



(https://fgiasson.com/blog/wp-content/uploads/2023/08/run_container.jpg).

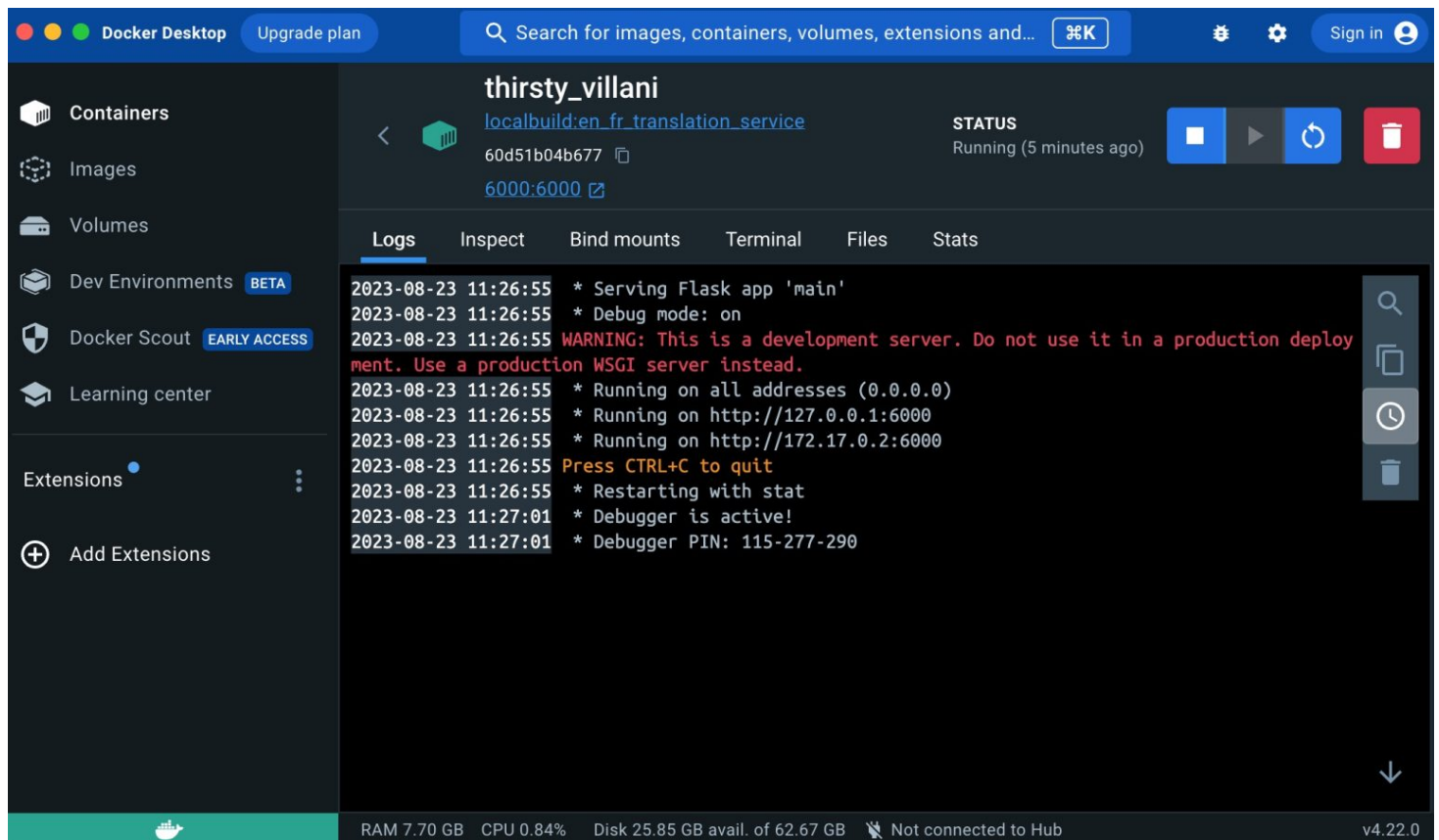
The most important thing to define here is to `Host port`. If you leave it empty, then the `port 6000` we exposed in the Docker file will become `unbound` and we won't be able to reach the service running in the container.

Once you click the `Run` button, the container will appear in the `Containers` section:



(<https://fgiasson.com/blog/wp-content/uploads/2023/08/containers.jpg>)

And if you click on it's name's link, you will have access to the internal of the container (the files it contains, the execution logs, etc.):



(https://fgiasson.com/blog/wp-content/uploads/2023/08/container_details.jpg)

Now that the container is running, we can query the endpoint like this:

```
Bash
1 | curl http://localhost:6000/translate/en/fr/ POST -H "Content-Type: application/json" -v -d '{"en_t
```

It returns:

```
JSON
1 | {
2 |   "fr_text": "Vers la certification des syst\u00e8mes distribu\u00e9s fond\u00e9s sur l'apprentis
3 | }
```

And then for the French to English translation:

```
1 | curl http://localhost:6000/translate/fr/en/ POST -H "Content-Type: application/json" -v -d '{"fr_te
```

It returns:

JSON

```
1 | {  
2 |   "en_text": "What is admirable in the happiness of others is that one believes in it."  
3 | }
```
