# How to work around Docker

Using GitHub Codespaces

# 1) Fork the repository

# 2) Go to GitHub Codespaces

# 3) Create new Codespace with forked repo



- I recommend choosing **4-core cpu**, as it is slow otherwise
- You only have limited time in free plan, but should be enough: ~120 Core hours per month

# 4) Docker-compose up in exercise05

# 5) Open a new terminal and continue

# 5.2) Open a new terminal and continue

# 6) Copy the dataset and get started



Now we need to populate both tables with data. We will use the ImportTsv utility of HBase. Populate the table `wiki_small` by running the following (keep in mind that you should run this command in the container's bash):

```
hbase org.apache.hadoop.hbase.mapreduce.ImportTsv -Dimporttsv.separator=, -Dimporttsv.columns="HBASE_ROW_KEY,page:page_title,page:page_ns,page:revision_id,author:timestamp,author:contributor_id,author:con wiki_small enwiki-20200920-pages-articles-multistream_small.csv
```

We need to specify which column in the csv maps to which column in the HBase table. Note that we make `page_id` into the `HBASE_ROW_KEY` and how we specify the mappings between the **.csv columns** and the **family:column** in the HBase table.

These commands print a lot of messages, but they are mostly informational with occasional non-critical warnings; unless something goes wrong, of course :). The commands will also report some "Bad Lines", but you can safely ignore this -- some lines may contain illegal characters and be dropped, but most of the data is in good shape.

You can count how many rows there are using this command from your head node's shell:
`hbase org.apache.hadoop.hbase.mapreduce.RowCounter 'wiki_small'`
If everything goes right, you should see `ROWS=887784` in the output.

Now let's go into HBase shell again (by running `hbase shell`) and run some queries against the `wiki_small` table. We will look at some of the filters listed by HBase if you run `show_filters` in an HBase shell, e.g., `PrefixFilter()`, `ValueFilter()`, `SingleColumnValueFilter()`.

## Task 2.1: Indexing

PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    PORTS  7    COMMENTS

```
1000158                         column=page:bytes, timestamp=1730124523047, value=24
1000158                         column=page:page_ns, timestamp=1730124523047, value=0
1000158                         column=page:page_title, timestamp=1730124523047, value=|The Space Race|
1000158                         column=page:revision_id, timestamp=1730124523047, value=19128253
1000159                         column=author:contributor_id, timestamp=1730124523047, value=55783
1000159                         column=author:contributor_name, timestamp=1730124523047, value=|CryptoDerk|
1000159                         column=author:timestamp, timestamp=1730124523047, value=2004-09-20T00:04:00Z
1000159                         column=page:bytes, timestamp=1730124523047, value=40
1000159                         column=page:page_ns, timestamp=1730124523047, value=0
1000159                         column=page:page_title, timestamp=1730124523047, value=|Waterloo College|
1000159                         column=page:revision_id, timestamp=1730124523047, value=16782323
7 row(s)
Took 0.6184 seconds
hbase(main):002:0>
```

# Additional notes

- Even if you have an ARM architecture this all will happen within the Codespace so you do not have to touch the docker-compose files

- Every time you create a new Codespace you will have to start from scratch

- Be mindful of the instructions in the exercise: when you have to be in the shell and when you instead cannot