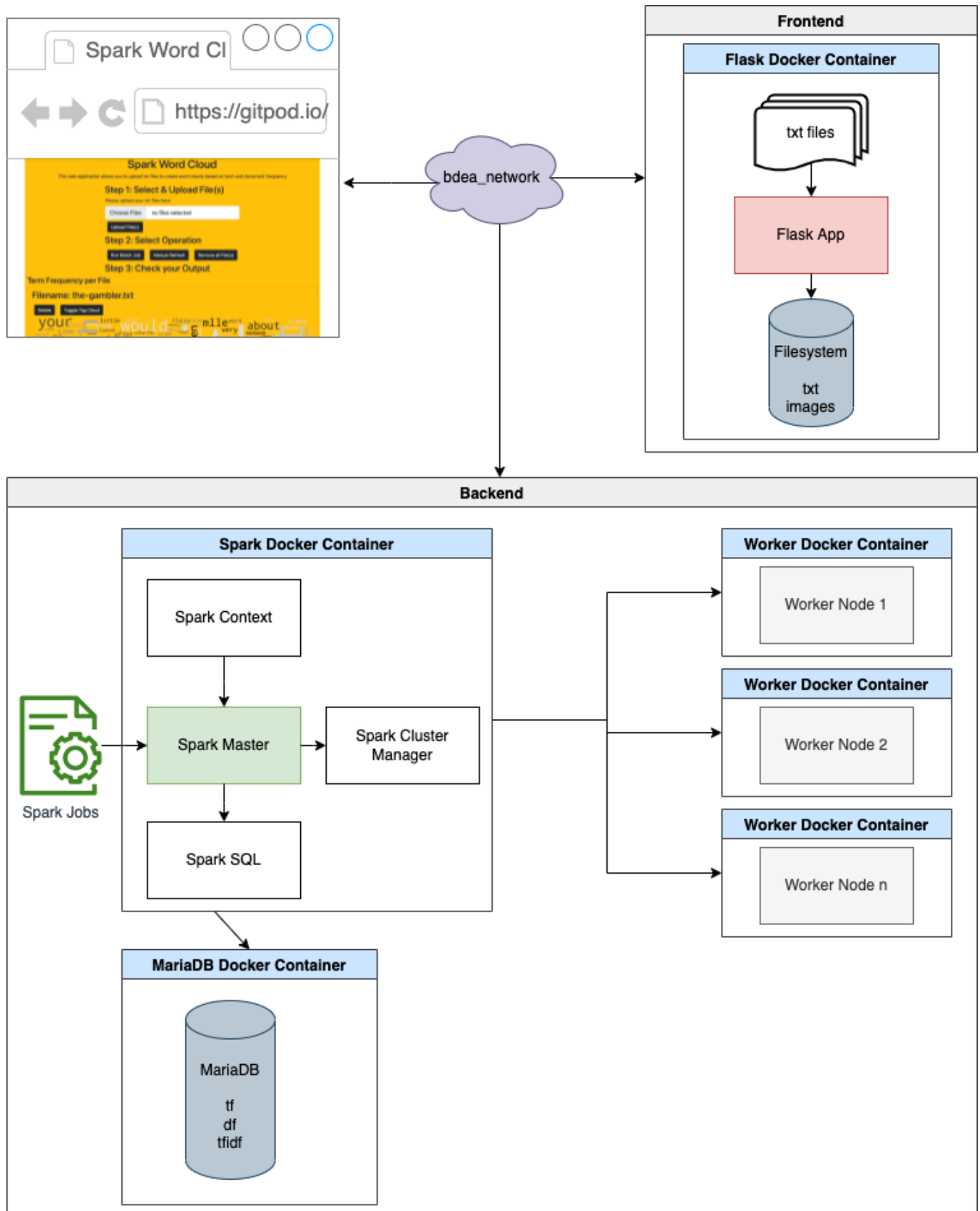


- df (stores the document frequency for all words)
- tfidf (stores the TFIDF value for each document)

Due to its setup, Spark distributed the workload between the available worker nodes (please read below on how to scale them).

Architecture

The following image is a visualization of the description above:



Challenges & Problems

- Read and write in the same spark job -> can cache the table
- Special characters are not correctly saved in database -> change char_set and collate to utf-8
- Queries in SQL are sometimes not so easy
- Bootstrap elements center -> frontend was annoying
- Setup with database and JDBC Driver -> link driver by spark-submit and in docker-image

- Docker setup with all components -> try and error until it works
- Spark Jobs depends strongly on systems hardware e. g. 4 MB .txt-file PC -> 11-20sec, laptop -> 2-3min

Development

You can easily develop this application by opening up GitPod (see above) and have the whole environment up and running. Alternatively you can clone the repo and develop locally - simply run the following commands from the root of the repository:

```
docker-compose -f webapp/docker-compose.yml up
```

Worker can be scaled with `--scale spark-worker=`:

```
docker-compose -f webapp/docker-compose.yml up --scale spark-worker=2
```

Sources

Text File Sources

- [The Grand Inquisitor by Fyodor Dostoyevsky](#)
- [The Brothers Karamazov by Fyodor Dostoyevsky](#)
- [The Gambler by Fyodor Dostoyevsky](#)
- [The Idiot by Fyodor Dostoyevsky](#)