

# proof\_platform

Scalable Web Scrapping

Installation guide

***Tomáš Kocman, Libor Polčák***



TARZAN project VI20172020062 document  
Faculty of Information Technology, Brno University of Technology

Last change: March 2, 2020



# proof\_platform — Installation guide

Tomáš Kocman, Libor Polčák

Faculty of Information Technology, Brno University of Technology, e-mail:  
`polcak@fit.vutbr.cz`

This project was created as part of the TARZAN project and the diploma thesis of Tomáš Kocman<sup>1</sup>. This project allows scalable web scrapping, i.e. archiving of web pages.

## 1 Getting Started

Please instal Docker as the first step<sup>2</sup>. This project is only supported as a Docker image. If you want to run Docker-less, you will need to build `scrapit` and `lemit` (see the subdirectories). For Docker-less install, you will also need PostgreSQL RDBMS<sup>3</sup>, Redis<sup>4</sup>, pgAdmin<sup>5</sup>.

Once you install Docker, edit the `docker-compose.yml` volume section. Select a directory where to store the archives.

Then, installation of the applications is as simple as:

```
docker-compose build
docker-compose up -d
```

---

<sup>1</sup> <https://www.fit.vut.cz/study/thesis/21459/>

<sup>2</sup> <https://docs.docker.com/install/>

<sup>3</sup> <https://www.postgresql.org/>

<sup>4</sup> <https://redis.io/>

<sup>5</sup> <https://www.pgadmin.org/>