

# Webscraper framework

---

It is not permitted to use on webpages that are not allowing webscraping. This can be checked on most sites with a /robot.txt after. An example is <https://www.youtube.com/robots.txt> No robot has permission to get all data from that is under User-agent: \* (\* = meaning all).

## Table of Contents

- [Webscraper framework](#)
  - [Table of Contents](#)
  - [Framework Structure](#)
  - [Getting Started](#)
    - [Installation](#)
      - [See Run demo of webscraper for how to use it](#)
  - [How to run the demo for the webscraper](#)
    - [As a new project](#)
    - [Use the existing test project.](#)
      - [Usage](#)
  - [Additional notes](#)
    - [Where to find the framework after downloading to your project?](#)
    - [What happens when running `npx install-webscraper`?](#)
    - [Yet to be implemented](#)

## Framework Structure

This framework has the following structure from the git repo, before installation to a existing project.

- [server.ts](#): Where all the API called happen.
- [routes/](#): This directory contains all route definitions.
- [routes/utils/helper](#) Where all functions used in routes is located.
- [hooks/](#): Where the premade connection to the server form making it easier for the developer is.
- [components/](#) - Where the demo file is.
- [components/ScraperDemo.ts](#) - Is a demonstration on how to use this framework.
- [test-next/](#) - This folder contains a next.js project that is using this framework. This could be just as a test for the framework with no need of coding, just the prerequisites mention in the [use the existing test project](#)
- [server.test.ts](#) - Where the test enviroment is located
- [install-webscraper.js](#) - Is the script that is adding this framework to the developers project.
- [README.md](#) - Is this current file/ containing getting started/ installation etc...
- [README.pdf](#) - PDF version of README.md

## Getting Started

Note: How

Installation

(ps: This is after installation of Next.js with command `npx create-next-app projectnamehere` or `.` for this folder.)

1. Install [Node.js](#) and install [Git](#). Otherwise all commands may not work. You can check if you have them already installed with `git -v` and `node -v` in a terminal.
2. Run `npm install Bass4Nation/webscraper` in the project directory for downloading and installation some of the elements used
3. Run `npx install-webscraper` this will merge package.json with webscrapers package.json. This will also copy and make folders needed for the webscraper to work.

### See Run demo of webscraper for how to use it

4. Run `npm run dev` to start the server and the website.
5. On the website, paste the URL of the website you want to scrape and click "Scrape". The scraped data will be displayed on the page/console and on the server at this moment. (This is just a test to see if the scraping works. And an example on how to use the scraper.)

## How to run the demo for the webscraper

### As a new project

This is how I would do it at a start next js project. (ps: This is just a example on where to place the demo)

1. In pages/index.tsx.
  1. `import "ScraperResult" from '@components/ScraperResult.tsx'`
  2. Add `<ScraperDemo />` somewhere in the frontend in the file. (Must be in the front end web part of the project.)
  3. In terminal `npm run dev`
  4. In the terminal it will say that start-server and start-client has started at port 3000 and 3002.
  5. Just copy a url from any pages. Ex: <https://b4n.no/> it is okay to use since it is my portfolio page and I allow the use of a scraper on that page.

### Use the existing test project.

1. Clone the project from github instead of using npm install.
  1. <https://github.com/Bass4Nation/webscraper> - Webscraper's Github repository that would be needed to for cloning
2. In a terminal navigate to `test-next` folder.
3. Run `npm install` to install the dependencies.
4. Just run `npm run dev` to start the server and the website.
5. On the website, paste the URL of the website you want to scrape and click "Scrape". The scraped data will be displayed on the page/console and on the server at this moment. (This is just a test to see if the scraping works. And an example on how to use the scraper.) Could be a good idea to use a page that is not protected by a robot.txt file. Ex: <https://b4n.no/> it is okay to use since it is my portfolio page and I allow the use of a scraper on that page.

### Usage

1. Paste the URL of the website you want to scrape in the input field.

2. Click "Scrape".
3. The scraped data will be displayed on the page/console and on the server at this moment.
4. You can also download the scraped data as a JSON file. (Not yet implemented)

Prerequisite: git, node, terminal, a code editor (I used Visual studio code to create this framework)

```
npm install Bass4Nation/webscraper npx install-webscraper npm run dev
```

Server and client on separate port. Usually server on 3002 and client is 3000.

For testing the demo that is in the code. You must import ScaperResult.tsx Can be shown in files at webscraper github how to import it and use it import ScaperResult from '@components/ScaperResult.tsx'

```
npm install user-agents
```

## Additional notes

This framework may be updated after school delivery date 1.June 2023. But it will be a waiting periode because this is a school assignment/exam, so I want it to be graded first.

Where to find the framework after downloading to your project?

This is after `npm i Bass4Nation/webscraper`. All packages/frameworks will be installed in `node_modules/` folder and will be named `webscraper`

What happens when running `npx install-webscraper`?

It will copy all files that are required for the framework from `node_modules/webscraper` to root level of your project and also create all folders needed for the existing webscraper framework. Also merging this framework's package.json with the projects current package

Yet to be implemented

1. Arrange the scraped data in a more readable format.
2. Possible to use this in other than Node.js/Next.js environments.
3. Make the framework more dynamic for the developer. Like scrapeproduct, now it is set in stone on how to to scrape, so the developer must make their store scrape function similiar to the existing but adjusting it to the store they want to scrape from.
4. Uninstalling the framework from a current project. Another way then just deleting files