Data

# Web Scraping using Puppeteer and Node

June 10, 2019

**Hrithik Jha**
Software Engineer

## Introduction

Puppeteer is a node library which provides an API to control Google Chrome and Chromium. It can be used to scrape all aspects of a Chrome (or Chromium) window including the Chrome Developer Tools. Today, we'll be scraping lobste.rs.

## Why Scrape Data?

Scraping can be used for many different purposes. It helps collect data for Machine Learning or can be used for Data Visualization. It can also be used to automate processes. An amazing example would be to scrape movie tickets and

prices of shows in theaters nearby. With the data you can sort in order of the price, show timing even with data from different websites to make decisions.

## Getting Started

For starters, you'd need a text editor (Pro Tip: MS Word with Consolas) and a system with Node. You can get Node from here or you can use `brew install node` on a Mac, or `sudo apt-get install nodejs npm` on Linux.

To install Puppeteer, simply run:

```
npm install puppeteer
```

Puppeteer is built on Node and can be run with any existing node program. One might use Requests or Selenium (more scraping and automation libraries) with Flask or Django, similarly we can use Puppeteer for programs written on Node.

## Into The Code

The first thing we do is to import the Puppeteer library to be used in the program.

```
const puppeteer = require('puppeteer');
```

The scraping logic is written inside the appropriately named function, `scrape()`. You would also notice the async keyword. That is needed as most of the Puppeteer API are asynchronous and have to be written inside an async function.

```
let scrape = async () => {
```

Inside the function, we initialize the browser and the page. It is very similar to the flow of opening a browser and website as one would, open browser > open new tab > go to the URL.

```
const browser = await puppeteer.launch({headless: false});
const page = await browser.newPage();
await page.goto('https://lobste.rs/');
```

We have a delay to allow the page load the DOM contents as our queries would not work if the element we want to scrape has not loaded.
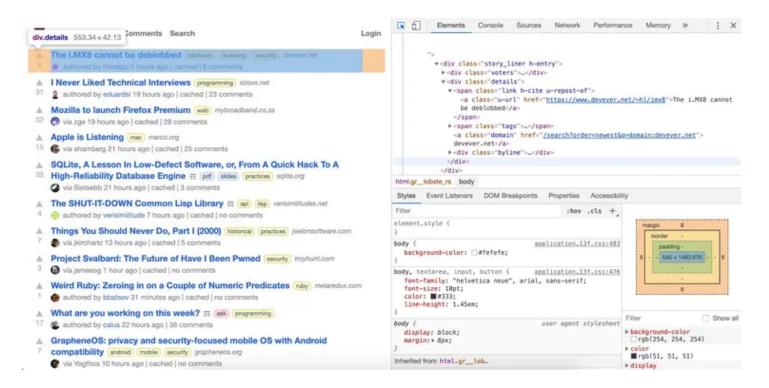
```
await page.waitFor(1000);
```

In Puppeteer, all interaction with the website's content is done using the page.evaluate function.
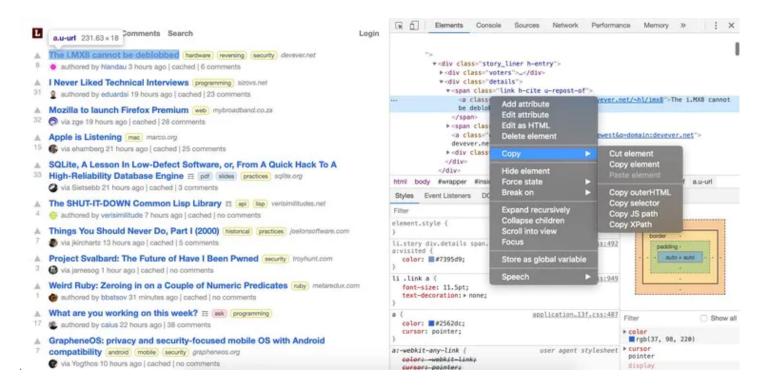
```
const result = await page.evaluate(() => {
```

Inside the function, we initialise an empty array to store the scraped data.

Here, we save the elements with the class name `u-url` and `score`. To get the class name of the element you want to scrape press Ctrl + Shift + I on Windows or Cmd + Opt + I on Mac. Here, click the top left button to select an element. Once that element is highlighted in the Elements tab, right click on it and copy the class name of the element.

In this example we use Class Names, but there are many other methods such as using the CSS Selector, the element ID or the XPath. The JavaScript Document object has support for using these methods as well.

```
let articles = []
let titles = document.getElementsByClassName('u-url');
let scores = document.getElementsByClassName('score');
```

After copying it, you can then change the parameter in
`.getElementsByClassName` to get the data in that element.

Next, we write a for loop to save all Titles and Scores in the Articles array. Before doing that we check if the number of Titles and Scores we have are the same to avoid any logical errors.

```
if(titles.length == scores.length) {
        for (i = 0; i < titles.length; i++) {
            articles.push(
            {
                'title': titles[i].textContent,
                'score': scores[i].textContent
            })
        }
    }

    return {
        articles
    }

});
```

After scraping, we close the browser and return the Articles.

```
    browser.close();
    return result;
};
```

We return Articles, to the caller of the function. Here, we have used `.then()` to print the values in console only after the function is done running.

Clubhouse!

```
        console.log(value);
});
```

# Conclusion

We finally have all the scraped data in the console. This can further be modified to save the data in an Excel sheet or a Json file. Apart from scraping Puppeteer can be used to automate processes like submitting assignments or filling forms online.

Thank you for reading the article!

**Share this Developer How-to**

**3 Comments**         **Clubhouse**                    🔴1  **Login**  ⌄

♡ **Recommend**          🐦 **Tweet**         f **Share**                        **Sort by Best** ⌄

Join the discussion…

**LOG IN WITH**

**OR SIGN UP WITH DISQUS** ⑦

Name

**Brent Engelbrecht** • 7 months ago

Thanks! Awesome, quick 'n easy tutorial.

ᐱ | �v • **Reply** • **Share ›**

**James Vreeken** • 8 months ago • edited

You could also use webtask.io if you can't install node on your
computer for whatever reason... or if you wanted to set it on a
cron job to run every so often.

Edit: doesn't seem to work as webtask is serverless and
doesn't have chrome installed.

I've used webtask.io with parsehub's webhooks to automate
and its worked pretty well so I was interested to try this out.

Question:
Instead of waiting an arbitrary 1000ms the "page.goto"
method is supposed to wait for the DomContentLoaded event:
https://github.com/GoogleCh...

There's also networkidle0 ? any thoughts on this so we don't
have to arbitrarily pick a random number of time we expect it
to take to load?

ᐱ | �v • **Reply** • **Share ›**

**Hrithik Jha** ➔ James Vreeken • 7 months ago

Hi James, I greatly apperciate your response. I mostly
use Puppeteer for automation and web scraping is a

# Join the over 25,000 developers, designers, and product managers already using Clubhouse ⚑

Get started—free forever

G  Sign up with Google

⚑

## Product

Features

Pricing

Enterprise

Integrations

Write beta

## Company

About Us

Customers

Careers

Press

Contact

Brand Guidelines

FAQ

## Developers

REST API Docs

Webhook API Docs

Community

Open Source Projects

Developer How–Tos

## Resources

Help Center

Blog

Webinars

Status

Referral Program

**The Clubhouse Blog**

Collaborate faster at scale with Group @mentions and notifications

Execute commands more efficiently with the Action Bar

Privacy Policy                    Terms of Use                    Security                    Your Cookies