Level up your Twilio API skills in **TwilioQuest**, an
educational game for Mac, Windows, and Linux.

Download
Now

## BLOG

DOCS      LOG IN      SIGN UP      ↗ TWILIO

Build the future of
communications.

START BUILDING FOR FREE

BY **SAM AGNEW** ▪ 2020-02-07

TWITTER          FACEBOOK          LINKEDIN

# Web Scraping and Parsing HTML with Node.js and Cheerio



The internet has a wide variety of information for human consumption. But this data is often
difficult to access programmatically if it doesn't come in the form of a dedicated REST API.
With Node.js tools like Cheerio, you can scrape and parse this data directly from web pages
to use for your projects and applications.

Let's use the example of scraping MIDI data to train a neural network that can generate classic Nintendo-sounding music. In order to do this, we'll need a set of music from old Nintendo games. Using Cheerio we can scrape this data from the Video Game Music Archive.

## Getting started and setting up dependencies

Before moving on, you will need to make sure you have an up to date version of Node.js and npm installed.

Navigate to the directory where you want this code to live and run the following command in your terminal to create a package for this project:

```
1  npm init --yes
```

The `--yes` argument runs through all of the prompts that you would otherwise have to fill out or skip. Now that we have a package.json for our app

For making HTTP requests to get data from the web page we will use the Got library, and for parsing through the HTML we'll use Cheerio.

Run the following command in your terminal to install these libraries:

```
1  npm install got@10.4.0 cheerio@1.0.0-rc.3
```

Cheerio implements a subset of core jQuery, making it a familiar tool to use for lots of JavaScript developers. Let's dive into how to use it.

## Using Got to retrieve data to use with Cheerio

First let's write some code to grab the HTML from the web page, and look at how we can start parsing through it. The following code will send a `GET` request to the web page we want, and will create a Cheerio object with the HTML from that page. We'll name it `$` following the infamous jQuery convention:

```
1  const fs = require('fs');
2  const cheerio = require('cheerio');
3  const got = require('got');
```

```
 4
 5   const vgmUrl= 'https://www.vgmusic.com/music/console/nintendo/nes';
 6
 7   got(vgmUrl).then(response => {
 8     const $ = cheerio.load(response.body);
 9     console.log($('title')[0]);
10   }).catch(err => {
11     console.log(err);
12   });
```

With this `$` object, you can navigate through the HTML and retrieve DOM elements for the data you want, in the same way that you can with jQuery. For example, `$('title')` will get you an array of objects corresponding to every `<title>` tag on the page. There's typically only one `title` element, so this will be an array with one object. If you run this code with the command `node index.js`, it will log the structure of this object to the console.

## Getting familiar with Cheerio

When you have an object corresponding to an element in the HTML you're parsing through, you can do things like navigate through its children, parent and sibling elements. The child of this `<title>` element is the text within the tags. So `console.log($('title')[0].children[0].data);` will log the title of the web page.

If you want to get more specific in your query, there are a variety of selectors you can use to parse through the HTML. Two of the most common ones are to search for elements by class or ID. If you wanted to get a div with the ID of "menu" you would run `$('#menu')` and if you wanted all of the columns in the table of VGM MIDIs with the "header" class, you'd do `$('td.header')`

What we want on this page are the hyperlinks to all of the MIDI files we need to download. We can start by getting every link on the page using `$('a')`. Add the following to your code in `index.js`:
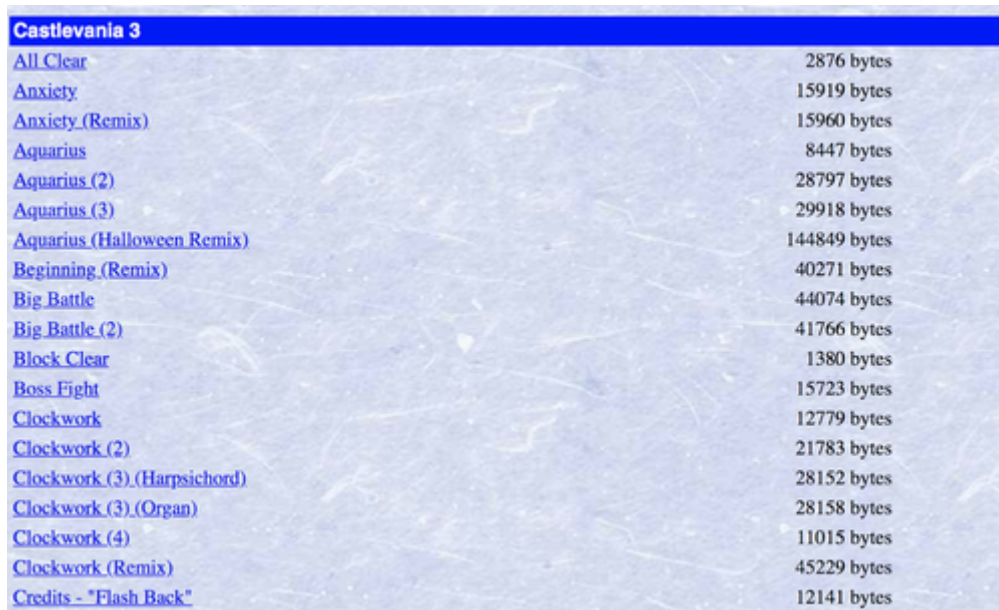
```
1   got(vgmUrl).then(response => {
2     const $ = cheerio.load(response.body);
3
4     $('a').each((i, link) => {
5       const href = link.attribs.href;
6       console.log(href);
7     });
8   }).catch(err => {
```

```
 9      console.log(err);
10   });
```

This code logs the URL of every link on the page. Notice that we're able to look through all elements from a given selector using the `.each()` function. Iterating through every link on the page is great, but we're going to need to get a little more specific than that if we want to download all of the MIDI files.

# Filtering through HTML elements with Cheerio

Before writing more code to parse the content that we want, let's first take a look at the HTML that's rendered by the browser. Every web page is different, and sometimes getting the right data out of them requires a bit of creativity, pattern recognition, and experimentation.



Our goal is to download a bunch of MIDI files, but there are a lot of duplicate tracks on this webpage, as well as remixes of songs. We only want one of each song, and because our ultimate goal is to use this data to train a neural network to generate accurate Nintendo music, we won't want to train it on user-created remixes.

When you're writing code to parse through a web page, it's usually helpful to use the developer tools available to you in most modern browsers. If you right-click on the element you're interested in, you can inspect the HTML behind that element to get more insight.

With Cheerio, you can write filter functions to fine-tune which data you want from your selectors. These functions loop through all elements for a given selector and return true or false based on whether they should be included in the set or not.

If you looked through the data that was logged in the previous step, you might have noticed that there are quite a few links on the page that have no `href` attribute, and therefore lead nowhere. We can be sure those are not the MIDIs we are looking for, so let's write a short function to filter those out as well as making sure that elements which do contain a `href` element lead to a `.mid` file:

```
1  const isMidi = (i, link) => {
2    // Return false if there is no href attribute.
3    if(typeof link.attribs.href === 'undefined') { return false }
4
5    return link.attribs.href.includes('.mid');
6  };
```

Now we have the problem of not wanting to download duplicates or user generated remixes. For this we can use regular expressions to make sure we are only getting links whose text has no parentheses, as only the duplicates and remixes contain parentheses:

```
1  const noParens = (i, link) => {
2    // Regular expression to determine if the text has parentheses.
3    const parensRegex = /^((?!\().)*$/;
4    return parensRegex.test(link.children[0].data);
5  };
```

Try adding these to your code in `index.js`:

```
1  got(vgmUrl).then(response => {
2    const $ = cheerio.load(response.body);
3
4    $('a').filter(isMidi).filter(noParens).each((i, link) => {
5      const href = link.attribs.href;
```

```
6      console.log(href);
7    });
8  });
```

Run this code again and it should only be printing `.mid` files.

# Downloading the MIDI files we want from the webpage

Now that we have working code to iterate through every MIDI file that we want, we have to write code to download all of them.

In the callback function for looping through all of the MIDI links, add this code to stream the MIDI download into a local file, complete with error checking:

```
1  $('a').filter(isMidi).filter(noParens).each((i, link) => {
2    const fileName = link.attribs.href;
3
4    got.stream(`${vgmUrl}/${fileName}`)
5      .on('error', err => { console.log(err); console.log(`Error on $
6      .pipe(fs.createWriteStream(`MIDIs/${fileName}`))
7      .on('error', err => { console.log(err); console.log(`Error on $
8      .on('finish', () => console.log(`Finished ${fileName}`));
9  });
```

Run this code from a directory where you want to save all of the MIDI files, and watch your terminal screen display all 2230 MIDI files that you downloaded (at the time of writing this). With that, we should be finished scraping all of the MIDI files we need.

```
Downloaded: Levels 6 and 9
Downloaded: Cave
Downloaded: Great Palace Boss
Downloaded: House
Downloaded: Overworld
Downloaded: Palace
Downloaded: Palace Boss
Downloaded: Princess Zelda Awakes
Downloaded: Title Screen
Downloaded: Town
Downloaded: Whistle
```

Go through and listen to them and enjoy some Nintendo music!

## The vast expanse of the World Wide Web

Now that you can programmatically grab things from web pages, you have access to a huge source of data for whatever your projects need. One thing to keep in mind is that changes to a web page's HTML might break your code, so make sure to keep everything up to date if you're building applications on top of this.

If you're looking for something to do with the data you just grabbed from the Video Game Music Archive, you can try using Python libraries like Magenta to train a neural network with it.

I'm looking forward to seeing what you build. Feel free to reach out and share your experiences or ask any questions.

- Email: sagnew@twilio.com

- Twitter: @Sagnewshreds

- Github: Sagnew

- Twitch (streaming live code): Sagnewshreds

AUTHORS | Sam Agnew

## Search

Build the future of communications. Start today with Twilio's APIs and services.

START BUILDING FOR FREE

POSTS BY STACK

| JAVA | .NET | RUBY | PHP | PYTHON | SWIFT | ARDUINO | JAVASCRIPT |

POSTS BY PRODUCT

| SMS | AUTHY | VOICE | TWILIO CLIENT | MMS | VIDEO | TASK ROUTER | FLEX | SIP | IOT |

PROGRAMMABLE CHAT   STUDIO

CATEGORIES

Code, Tutorials and Hacks

Customer Highlights

Developers Drawing The Owl

News

Stories From The Road

The Owl's Nest: Inside Twilio

TWITTER                                                        FACEBOOK

# Developer stories to your inbox.

Subscribe to the Developer Digest, a monthly dose of all things code.

Enter your email...

You may unsubscribe at any time using the unsubscribe link in the digest email. See our privacy policy for more information.

NEW!

## Tutorials

Sample applications that cover common use cases in a variety of languages. Download, test drive, and tweak them yourself.

Get started

SIGN UP AND START BUILDING

Not ready yet? Talk to an expert.

ABOUT
LEGAL
COPYRIGHT © 2020 TWILIO INC.
ALL RIGHTS RESERVED.
PROTECTED BY RECAPTCHA – PRIVACY– TERMS