

~~Воруем~~ Добываем данные из интернета
используя
Node.js

Малькевич Егор

От работы время я изучаю возможности получить разного рода данные из интернета причем таким образом что бы сервер ценой в 5\$ на Digital Ocean мог справиться с этой работой. Для меня это MVP.

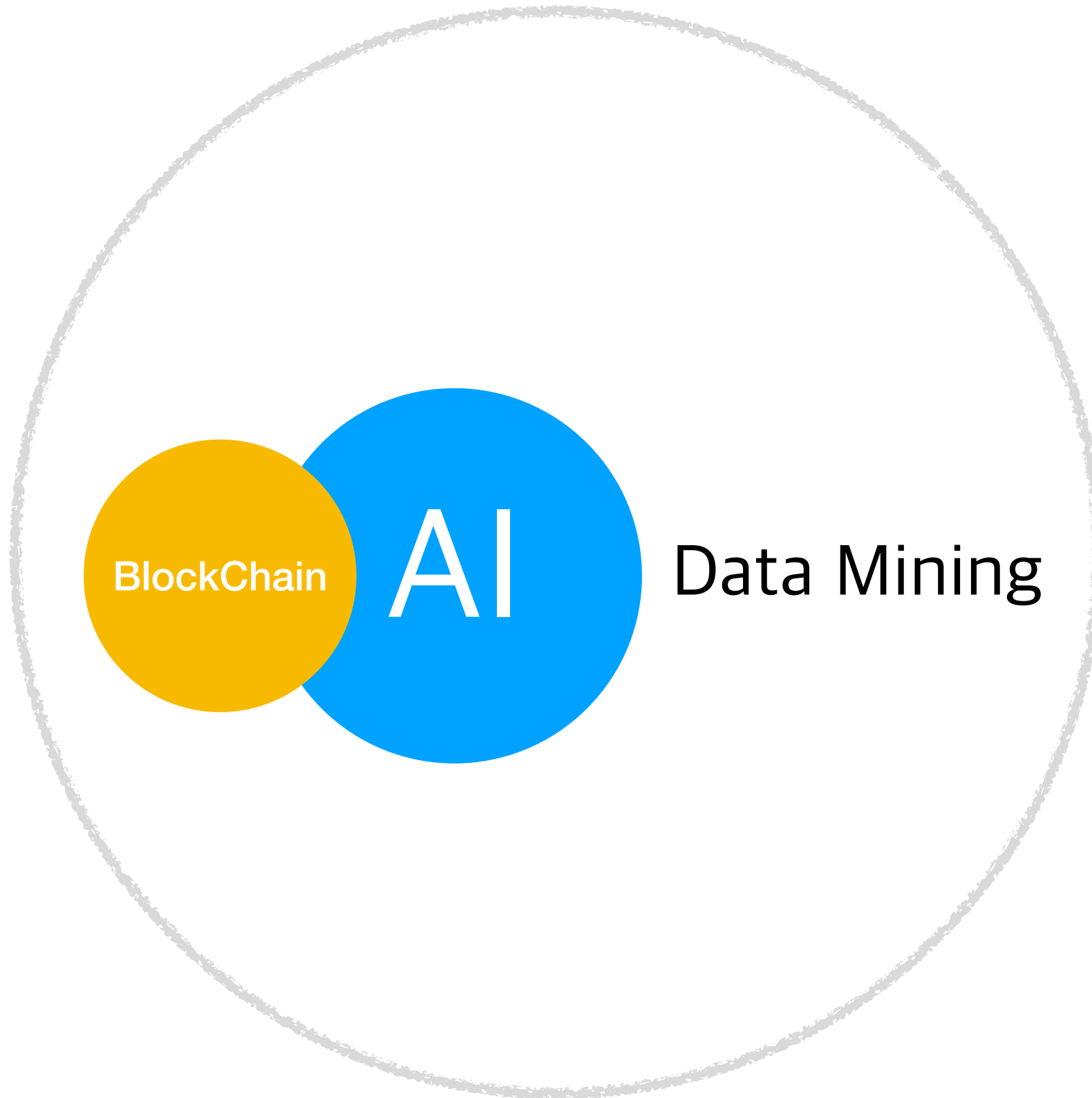
Собственно и сегодня я здесь перед вами, с целью поделиться накопленными знаниями, и показать что добывать данные нужно и нужно это всем и каждому. Пускай это будет новый голубой океан





<https://malkevich.com/>



Что сейчас Main Stream?



Парсеры

Followers	 YouTube	 facebook	 Instagram
100k - 500k	\$12,500	\$6,250	\$5,000
500k - 1m	\$25,000	\$12,500	\$10,000
1m - 3m	\$125,000	\$62,500	\$50,000
3m - 7m	\$187,500	\$93,750	\$75,000
over 7m	\$300,000	\$187,500	\$150,000

	 snapchat	 Vine	 twitter
100k - 500k	\$5,000	\$3,750	\$2,000
500k - 1m	\$10,000	\$7,500	\$4,000
1m - 3m	\$50,000	\$37,500	\$20,000
3m - 7m	\$75,000	\$56,250	\$30,000
over 7m	\$150,000	\$112,500	\$60,000

Source: Captiv8

Economist.com





Пару историй



Самолеты



Поиск дешевых авиабилетов

Лучший способ купить авиабилеты дешево

АВИАБИЛЕТЫ

ОТЕЛИ **60%**

АВТО

СТРАХОВКА

Все мы знаем такой популярный сервис как Aviasales.ru. Так вот в районе 2012 аналогичных сервисов были сотни если не тысячи. Я даже участвовал в разработке одного из них.

Основная проблема заключалась в том, что рейсы зачастую отменяют, переносят. Плюс перевозчики, повышают понижают цены. В общем данные устаревают крайне стримительно. Вы не можете один раз выкачать данные. А дальше работать со своим кэшем.


Поэтому приходилось постоянно актуализировать данные с сотен сайтов авиакомпаний, практически в режиме реального времени. Причем почти все сайты были мягко говоря кое какие, ну и разумеется никакого API.

Даже в то время мы использовали node.js 0.10 версии, и это было крайне хипстерски.

MSQ 

Город прибытия

Обратно 

1 пассажир, эконом 

Показать отели в новом окне

Найти билеты 



● **aviasales.ru**
Поисковый запрос

+ Сравнить

По всему миру ▾

2004 – настоящее вр... ▾

Все категории ▾

Веб-поиск ▾

Динамика популярности ?

100

75

50

25

1 янв. 200...

1 янв. 2009 г.

1 янв. 2014 г.

К счастью в том время Aviasales.ru рос. В него вливали деньги и он скупал всех конкурентов, кого то ради парсеров. Кого то, просто топил...

Между прочим наш сервис при 1000 еждневных пользователй в месяц зарабатывал 7 тыс \$

Примечание





Туда и обратно ▾

Один пассажир ▾

Эконом-класс ▾

Москва SVO DME VKO, Курс...

Минск MSQ

ЧТ, 11 янв. < >

ВС, 14 янв. < >



Минск – выберите рейс

Обратно: Москва, Жуковский

Сведения о путешествии

Прямой ✕

Цена ▾

Время ▾

ЮТэйр ✕

Ещё ▾

Советы по поиску рейсов



ДАТЫ

В выбранные дни билеты обычно стоят дешево

[ЕЩЁ...](#)



ЦЕНА

Минск...
меняются...

Если вдруг вы решите делать подобный сервис сегодня, обратите внимание на google.ru/flights

Потому что гугл уже парсит и индексирует миллионы подобных сервисов за вас.

АЭРОПОРТЫ

Минск: билеты на рейсы в аэропорты неподалеку могут стоить дешевле.

[ЕЩЁ...](#)



СОВЕТЫ

Минск: спланируйте путешествие

[ЕЩЁ...](#)

Рейсы вылета

Указана итоговая цена за один взрослый билет с учетом сборов за провоз багажа и другие комиссии.

[сборы за провоз](#)

Упорядочить по: ↑↓



10:00 – 11:20
ЮТэйр

1 ч. 20 мин.
VKO–MSQ

Прямой рейс

4 408 Р
туда и обратно



14:40 – 16:00
ЮТэйр

1 ч. 20 мин.
VKO–MSQ

Прямой рейс

4 408 Р
туда и обратно

Машины





Все же сталкивались с перекупами? Которые скупают тачки подешевле, шаманят и продают по дороже. Так вот, до появления таможенного союза, эти мастера продаж, могли зарабатывать влегкую по несколько тыс \$ с машины.

На сегодняшний день, ихняя маржа упала. Скажем до 300\$ с машины в среднем.

Сайтов по продаже авто стало меньше. Найти и купить самостоятельно авто - проще.

И вот забавно, но факт. Перекупы эволюционировали и нашли возможности в «высокочастотном трейдинге».

Любой регион

Все параметры

Показать 503 576 объявлений

BMW	22209	Hyundai	32863	Mercedes-Benz	26016	Renault	16523
Chevrolet	18195	Kia	21515	Mitsubishi	15596	Toyota	35500
Ford	24023	LADA (BA3)	82341	Nissan	27204	Volkswagen	27991



MAZDA6
в кредит от
285 рублей в месяц

[УЗНАТЬ БОЛЬШЕ](#)

Продажа автомобилей

По актуальности

За все время

Пр

Р



Peugeot 308

1.6 AT (120 л.с.)
привод, хэтчб

Москва

00 Р

2009

69 800 км



Suzuki Grand

2.0 AT (140 л.с.)
привод, внедо

Подольск

00 Р

2005

193 095 км

Схема такая, вы выкидываете машину на продажу, перекупы сами или используя ботов сканируют доски объявлений, и стараются первыми посмотреть тачку. И иногда успевают купить и продать в тот же день. Посути как работая как брокеры на бирже.

У автору есть платный сервис, который позволяет вам видеть объявления раньше других. Но зачастую дешевле написать свой маленький парсер, что многие и делают.

К примеру, в среднем в Минске выкидывают 30-60 новых объявлений о продаже машины в день.

То-есть если добавить туда маржу в 300\$ - максимальный суточный доход может достигать 90000\$

Видео и сериалы



фильтр

Фильтр сериалов

Любая страна

Все Зарубежные Русские

Оригинальные названия

Любой жанр

созтировать

По названию

Дополнительная фильтрация

Показать найденные: 10358

Скоро на нашем сайте

Лучшие сериалы



алфавит

Сериалы на буквы

#	А	Б	В	Г	Д
Е	Ж	З	И	Й	К
Л	М	Н	О	П	Р
С	Т	У	Ф	Х	Ц
Ч	Ш	Щ	Э	Ю	Я

+100500

+100500 на ТВ

Никто так не парсит контент, как пиратские сайты с видео.

Они пишут парсеры

- на группы вконтакте,
- сканируют торренты,
- сканируют своих конкурентов,
- сканируют сайты студий озвучки
- парсят еще кучу сомнительных сайтов
- они анализируют официальную ленту релизов фильмов и выгинеривают фэйковые страницы с ожидаемыми фильмами и сериалами, что бы попадать в высокочастотные запросы

Я не говорю про Seasonvar.ru, а говорю в целом про нишу пиратских видео.

Вы бы видели какие гавносайты являются, донорами этих агрегаторов.

Между прочим, даже средний такой пиратский сайт может приносить +-3k\$ при наличии 10k пользователей, на рекламе или пожертвованиях.

Мир наизнанку (8 сезон)

Морская Полиция: Лос Анджелес (9 сезон)

5 серия (GoldTeam)

6 серия (ColdFilm)

3-4 серия

3-4 серия (VirusProject)

9 серия

12 серия (GoldTeam)

5-6 серия (ColdFilm)

5 серия (ColdFilm)

4 серия (NewStudio)

30 серия

сериал полностью

сериал полностью

25 серия

6 серия (Ozz)

3-4 серия (Субтитры VP)

5 серия (ColdFilm)

5 серия (BaibaKo)

7 серия

4 серия (Sunshine)

1 серия (NewStudio)

6 серия

3 серия (Субтитры VP, VirusProject)

9-10 серия

19 серия

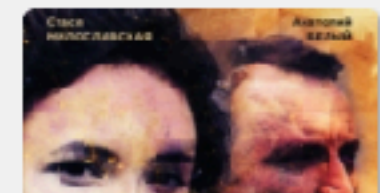
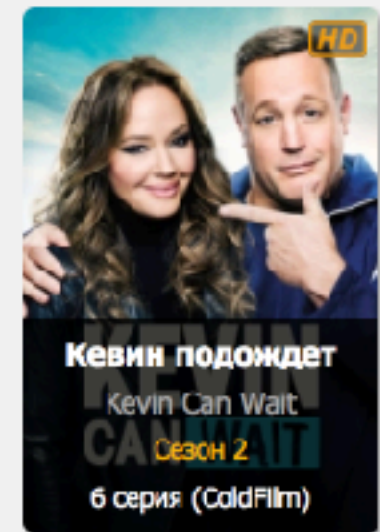
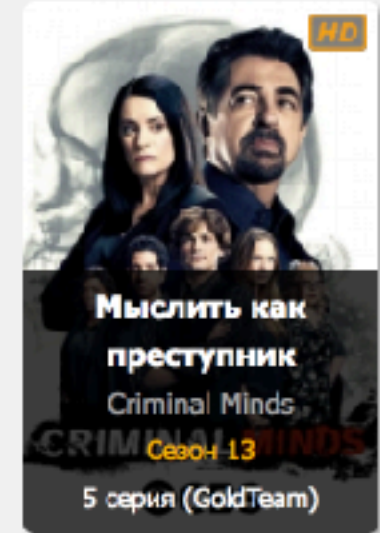
7 серия

5 серия (BaibaKo)

Обновления

Топ

Новинки





Мы с ребятами даже создали сервис, Gendalf TV. В первую очередь провести ряд экспериментов по парсингу видео.

Во вторую дать людям возможность стримить видео с любого сайта на Smart TV, без необходимости держать бук, или телефон залоченым.

Собственно в момент разработки этого сервиса я в серьез увлекся разработкой парсеров. Много экспериментировал в поиске дешевого решения, для разработки собственных парсеров.

В общем об этом опыте я и хочу с вами поделиться



<https://gendalf.tv>



API которым делаются и
которое нужно



http://seasonvar.ru/api.txt

Ответ возвращается в формате JSON.

Обработчик требует наличия двух обязательных параметров:

1. `key` - ключ авторизации, уникальный для каждого пользователя
2. `command` - название команды

Список команд (method):

1. `getSerialList` - список всех сериалов.

1. `country` (массив)

Список стран.

Пример: `array('сша', 'россия')`

2. `genre` (массив)

Список жанров.

Пример: `array('анимационные', 'комедия')`

3. `locale` (строка).

Два вида - `domestic`(отечественные), `foreign`(иностранные)

4. `sort` (массив)

Принимает два параметра: `order` и `method`

`-order` (массив)

`-kinopoisk` - сортировка по оценкам кинопоиска

`-imdb` - сортировка по оценкам imdb

`-popular` - сортировка по популярным сериалам

`-year` - сортировка по году

`-method`(строка)

`-desc` -выставляет DESC (по умолчанию ASC)

Пример:

`array('order' => 'popular', 'method' => 'DESC')`

`-lastSeasonInfo` (bool)

Дополнительная информация к последнему сезону

По умолчанию - `false`.

5. `letter` (string) - по букве (или части слова, начиная с начала)

1. `getSeason` - информацию по указанному сезону. Требует наличия дополнительного параметра `season_id` (id запрашиваемого сезона).

2. `search` - поиск. Требует наличия идентификатора пользователя (личного ключа).

<https://habrahabr.ru/post/302150/>

Хабрахабр

Публикации

Пользователи

Хабы

Компании

Песочница



ra1ym 30 мая 2016 в 12:03 [Разработка](#)

Что делать если Instagram не дал доступ к API?

PHP, Open source, API

Из песочницы

1 июня 2016 года Instagram [отключит](#) от своего API все приложения, которые не прошли модерацию. Что делать если вы в их числе?

Предыстория

Мы делаем сервис для постинга в Instagram по расписанию и используем API для получения информации об аккаунтах. Самим постингом занимаются телефоны в автоматическом режиме. Нам отказали в доступе к API после 1 июня (пробовали пройти модерацию два раза) поэтому было решено найти замену.

Сначала расскажу как мы использовали официальный API:

1. При добавлении аккаунта забираем из Instagram информацию об аккаунте: имя, фото профиля, количество постов, подписчиков, подписок.
2. Перед тем как опубликовать фото/видео мы запрашиваем количество постов, и тоже самое после публикации, если число постов увеличилось считаем публикацию успешной.
3. Если публикация прошла успешно забираем ссылку на последнее фото в профайле

https://vk.com/dev.php?method=audio_api

Продукты Документация Мои приложения Поддержка Поиск

Roadmap > Отключение публичного API для аудио

На этой странице мы рассказываем об одном из будущих изменений платформы API ВКонтакте. Чтобы посмотреть список всех недавних и грядущих изменений, перейдите к [этой странице](#).

1. Как это работает сейчас
2. Что изменится
3. Как подготовиться к изменениям

Ожидаемая дата изменений: **16 декабря 2016 года**.

Аудиозаписи — один из популярных сервисов ВКонтакте. Возможность слушать любимую музыку очень важна для наших пользователей. Мы хотим сохранить эту возможность и продолжить развитие аудиораздела.

С 2011 года мы начали переход к легальному использованию контента правообладателей. Сотрудничество с крупнейшими студиями позволяет ВКонтакте размещать эксклюзивные композиции, анонсировать новые альбомы исполнителей и улучшать сервис в целом.

При этом мы вынуждены брать на себя определённые обязательства в том, что касается борьбы с нарушением авторских прав. Это означает, что ВКонтакте не может сохранить публичный API аудиозаписей для приложений сторонних разработчиков в его текущем виде.

Это изменение не затронет пользователей веб-версии vk.com и официальных мобильных приложений ВКонтакте. Более того, в ближайшем будущем мы планируем активно совершенствовать аудиораздел в собственных продуктах.

Coding Time



What will we need?

- [request-promise](#)—Request is a simple HTTP client that allows us to make quick and easy HTTP calls.
- [cheerio](#)—jQuery for Node.js. Cheerio makes it easy to select, edit, and view DOM elements.

```
1 import request from 'request-promise'; 1.1M (gzipped: 320.5K)
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   uri: "https://www.google.com",
6   transform: (body) => cheerio.load(body)
7 };
8
9 const scrape = async (options) => {
10   return await request(options);
11 };
12
13 export default (params) => scrape({...options, ...params });
14
```

```
1 import request from 'request-promise'; 1.1M (gzipped: 320.5K)
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   uri: "https://www.google.com",
6   transform: (body) => cheerio.load(body)
7 };
8
9 const scrape = async (options) => {
10   return await request(options);
11 };
12
13 export default (params) => scrape({...options, ...params });
14
```

Example 1

```
1 import request from 'request-promise'; 1.1M (gzipped: 320.5K)
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   | uri: "https://www.google.com",
6   | transform: (body) => cheerio.load(body)
7   };
8
9 const scrape = async (options) => {
10  | return await request(options);
11  };
12
13 export default (params) => scrape({...options, ...params });
14
```

Example 1

```
1 import request from 'request-promise'; 1.1M (gzipped: 320.5K)
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   uri: "https://www.google.com",
6   transform: (body) => cheerio.load(body)
7 };
8
9 const scrape = async (options) => {
10   return await request(options);
11 };
12
13 export default (params) => scrape({...options, ...params });
14
```

Example 1

```
1 import request from 'request-promise'; 1.1M (gzipped: 320.5K)
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   uri: "https://www.google.com",
6   transform: (body) => cheerio.load(body)
7 };
8
9 const scrape = async (options) => {
10   return await request(options);
11 };
12
13 export default (params) => scrape({...options, ...params });
14
```

Example 1


```
1 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
2
3 const $ = cheerio.load(`
4 <ul id="cities">
5   <li class="large">New York</li>
6   <li id="c-medium">Portland</li>
7   <li class="small">Salem</li>
8 </ul>
9
10 <ul id="towns">
11   <li class="large">Bend</li>
12   <li id="t-medium">Hood River</li>
13   <li class="small">Madras</li>
14 </ul>`);
15
16 export default $;
```

```
1 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
2
3 const $ = cheerio.load(`
4 <ul id="cities">
5   <li class="large">New York</li>
6   <li id="c-medium">Portland</li>
7   <li class="small">Salem</li>
8 </ul>
9
10 <ul id="towns">
11   <li class="large">Bend</li>
12   <li id="t-medium">Hood River</li>
13   <li class="small">Madras</li>
14 </ul>`);
15
16 export default $;
```

Example 2

```
1 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
2
3 const $ = cheerio.load(`
4 <ul id="cities">
5     <li class="large">New York</li>
6     <li id="c-medium">Portland</li>
7     <li class="small">Salem</li>
8 </ul>
9
10 <ul id="towns">
11     <li class="large">Bend</li>
12     <li id="t-medium">Hood River</li>
13     <li class="small">Madras</li>
14 </ul>`);
15
16 export default $;
```

```
$('#cities').find('.small').text() // Salem
$('#towns').find('.small').text() // Madras

// or all

$('.small').map((index, el) => $(el).text()); // ["Salem",
"Madras"]
```

Limitations

MOST websites modify the DOM using JavaScript.

Unfortunately Cheerio doesn't resolve parsing a modified DOM. Dynamically generated content from procedures leveraging AJAX, client-side logic, and other async procedures are not available to Cheerio.

At the end of 2017

[PhantomJS](#)' development has [stopped](#), Headless Chrome is in the spotlight—and people love it

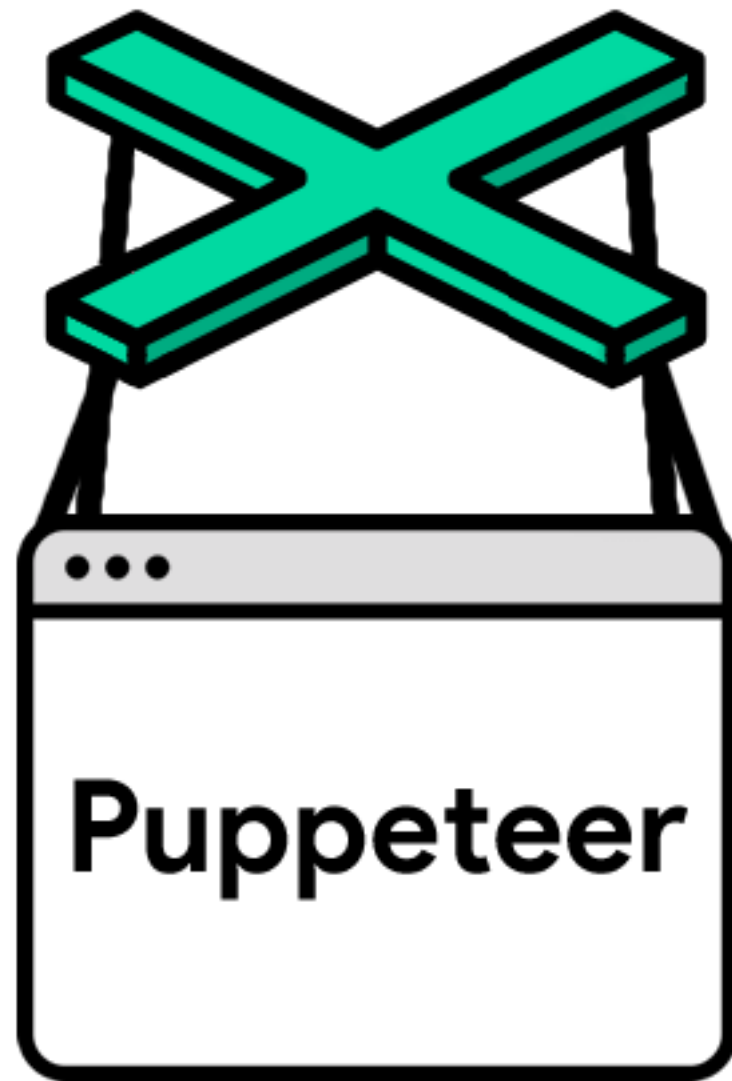


- There are a lot of libraries for controlling Chrome, pick the one you like;
- Web scraping with Headless Chrome is easy, even more so when you're aware of these tips & tricks;
- Headless browser visitors can be detected but nobody does it.

You'll need to have Node 8+ installed on your computer.

NPM package	chrome-remote-interface	chromeless	puppeteer	nickjs
GitHub stars	~2k	~10k	~11k	~0.03k (yeah!)
Maintainer	Andrea Cardaci	Graphcool	Google's DevTools team	Phantombuster
API	Generic low level mapping between NodeJS and the DevTools protocol.	Mostly complete API with command chaining (imitates NightmareJS). Allows for an easy AWS Lambda integration.	Google's take on a generic NodeJS mapping of the DevTools protocol. Complete API that tries to simplify Chrome automation.	Very simple API that gets out of your way as much as possible. Is specifically made with scraping in mind.

To do this, we'll use Puppeteer. [Puppeteer](#) is a Node library API that allows us to control headless Chrome



```
1 import puppeteer from 'puppeteer';
2
3 const scrape = async() => {
4     const browser = await puppeteer.launch({
5         headless: false
6     });
7     const page = await browser.newPage();
8     await page.setViewport({
9         width: 1240,
10        height: 680
11    });
12    await page.goto('https://google.com');
13    await page.screenshot({
14        path: 'google.png'
15    });
16    await browser.close();
17 }
18
19 export default () => scrape();
20
```

Example 3

```
1 import puppeteer from 'puppeteer';
2
3 const scrape = async() => {
4     const browser = await puppeteer.launch({
5         headless: false
6     });
7     const page = await browser.newPage();
8     await page.setViewport({
9         width: 1240,
10        height: 680
11    });
12    await page.goto('https://google.com');
13    await page.screenshot({
14        path: 'google.png'
15    });
16    await browser.close();
17 }
18
19 export default () => scrape();
20
```

Example 3


```
1 import puppeteer from 'puppeteer';
2
3 const scrape = async() => {
4     const browser = await puppeteer.launch({
5         headless: false
6     });
7     const page = await browser.newPage();
8     await page.setViewport({
9         width: 1240,
10        height: 680
11    });
12    await page.goto('https://google.com');
13    await page.screenshot({
14        path: 'google.png'
15    });
16    await browser.close();
17 }
18
19 export default () => scrape();
20
```

Example 3


```
1 import puppeteer from 'puppeteer';
2
3 const scrape = async() => {
4     const browser = await puppeteer.launch({
5         headless: false
6     });
7     const page = await browser.newPage();
8     await page.setViewport({
9         width: 1240,
10        height: 680
11    });
12    await page.goto('https://google.com');
13    await page.screenshot({
14        path: 'google.png'
15    });
16    await browser.close();
17 }
18
19 export default () => scrape();
20
```

Example 3

```
1 import puppeteer from 'puppeteer';
2
3 const scrape = async() => {
4     const browser = await puppeteer.launch({
5         headless: false
6     });
7     const page = await browser.newPage();
8     await page.setViewport({
9         width: 1240,
10        height: 680
11    });
12    await page.goto('https://google.com');
13    await page.screenshot({
14        path: 'google.png'
15    });
16    await browser.close();
17 }
18
19 export default () => scrape();
20
```

Example 3

```
1 import puppeteer from 'puppeteer';
2
3 const scrape = async() => {
4     const browser = await puppeteer.launch({
5         headless: false
6     });
7     const page = await browser.newPage();
8     await page.setViewport({
9         width: 1240,
10        height: 680
11    });
12    await page.goto('https://google.com');
13    await page.screenshot({
14        path: 'google.png'
15    });
16    await browser.close();
17 }
18
19 export default () => scrape();
20
```

Example 3



Anything that has real and lasting value is always a gift from within.

(Franz Kafka)

все, что имеет реальную и прочную ценность, всегда является подарком изнутри

Fallback to the first example

```
1 import request from 'request-promise'; 1.1M (gzipped: 320.5K)
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   | uri: "https://www.google.com",
6   | transform: (body) => cheerio.load(body)
7   | };
8
9 const scrape = async (options) => {
10  | return await request(options);
11  | };
12
13 export default (params) => scrape({...options, ...params });
14
```

```
1 import puppeteer from 'puppeteer';
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   uri: `http://books.toscrape.com/` ,
6   transform: (body) => cheerio.load(body)
7 };
8
9 const scrape = async (options) => {
10   const browser = await puppeteer.launch(options);
11   const page = await browser.newPage();
12   await page.goto(options.uri);
13   let content = await page.content();
14   let $ = options.transform(content);
15   browser.close();
16   return $;
17 };
18
19 export default (params) => scrape({...options, ...params });
20
```

Example 4


```
1 import puppeteer from 'puppeteer';
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   uri: `http://books.toscrape.com/` ,
6   transform: (body) => cheerio.load(body)
7 };
8
9 const scrape = async (options) => {
10   const browser = await puppeteer.launch(options);
11   const page = await browser.newPage();
12   await page.goto(options.uri);
13   let content = await page.content();
14   let $ = options.transform(content);
15   browser.close();
16   return $;
17 };
18
19 export default (params) => scrape({...options, ...params });
20
```

Example 4


```
1 import puppeteer from 'puppeteer';
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   uri: `http://books.toscrape.com/` ,
6   transform: (body) => cheerio.load(body)
7 };
8
9 const scrape = async (options) => {
10   const browser = await puppeteer.launch(options);
11   const page = await browser.newPage();
12   await page.goto(options.uri);
13   let content = await page.content();
14   let $ = options.transform(content);
15   browser.close();
16   return $;
17 };
18
19 export default (params) => scrape({...options, ...params });
20
```

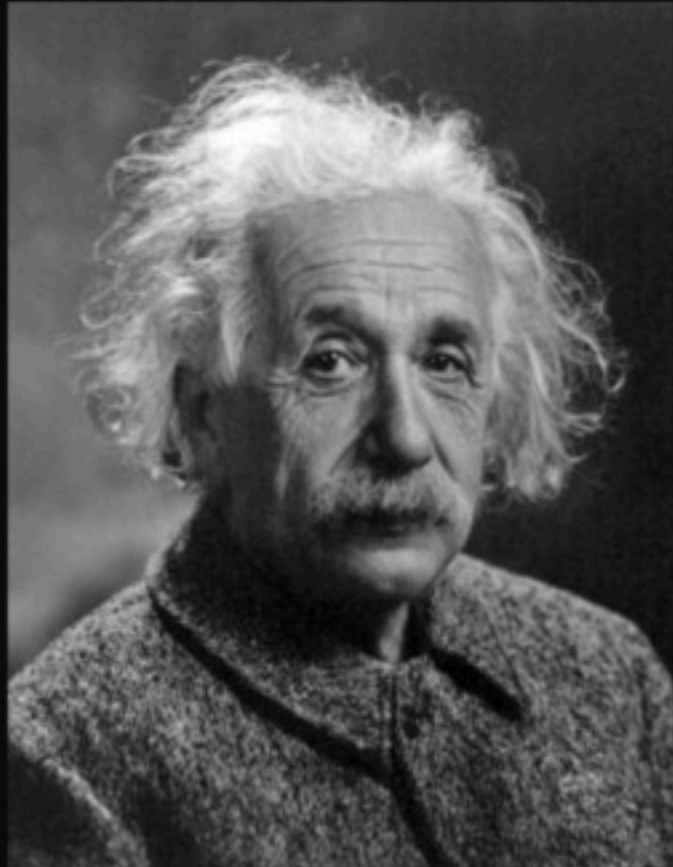
Example 4

```
1 import puppeteer from 'puppeteer';
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   uri: `http://books.toscrape.com/`,
6   transform: (body) => cheerio.load(body)
7 };
8
9 const scrape = async (options) => {
10   const browser = await puppeteer.launch(options);
11   const page = await browser.newPage();
12   await page.goto(options.uri);
13   let content = await page.content();
14   let $ = options.transform(content);
15   browser.close();
16   return $;
17 };
18
19 export default (params) => scrape({...options, ...params });
20
```

Example 4

```
1 import puppeteer from 'puppeteer';
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   uri: `http://books.toscrape.com/`,
6   transform: (body) => cheerio.load(body)
7 };
8
9 const scrape = async (options) => {
10   const browser = await puppeteer.launch(options);
11   const page = await browser.newPage();
12   await page.goto(options.uri);
13   let content = await page.content();
14   let $ = options.transform(content);
15   browser.close();
16   return $;
17 };
18
19 export default (params) => scrape({...options, ...params });
20
```

Example 4



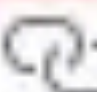
Everything should be made as simple as possible, but not one bit simpler.

(Albert Einstein)

Block or report user

 Poland / Norway / Sweden

 john@sundell.co

 <https://twitter.com/johnsundell>

1,455

Organizations



M

W

Fr



john

Pull requests Issues Marketplace Explore



Repositories 22K

Code 40M

Commits 160M

Issues 190K

Wikis 28K

Users 74K

Advanced search

74,511 users

Sort: Best match ▾



john John McGrath

I'm a solutions architect at AWS and prior to that co-founded Entelo. Interested in renewable energy, media, and democracy.

San Francisco, CA

Follow



JohnSundell John Sundell

I build apps, games and Swift developer tools! Passionate about open source & developer productivity. You can follow me on Twitter @johnsundell.

Poland / Norway / Sweden ✉ john@sundell.co

Follow



jrOcket John Stevenson

Speaker, author, conference organiser & community obsessed developer. Loves Clojure, Emacs, Cats, Cycling & Agile development.

London, UK ✉ john@jrOcket.co.uk

Follow



jdegoes John A. De Goes

Boulder, CO ✉ john@degoes.net

Follow

Languages

JavaScript	5,645
Java	4,034
Python	3,082
HTML	2,202
Ruby	1,925
PHP	1,787
C#	1,497
CSS	1,349
C++	1,342
C	1,170



Sign in to GitHub

Username or email address

Password

[Forgot password?](#)

Sign in

New to GitHub? [Create an account.](#)

[Terms](#) [Privacy](#) [Security](#) [Contact GitHub](#)

Elements Console Sources Network

```
<div id="js-pjax-container" data-pjax-container>
  <div class="auth-form px-3" id="login">
    <!-- '' -->
    <!-- </textarea></xmp> -->
    <form accept-charset="UTF-8" action="/session" method="post">
      <div style="margin:0;padding:0;display:inline">...</div>
      <div class="auth-form-header p-0">...</div>
      <div id="js-flash-container">
      </div>
      <div class="auth-form-body mt-3">
        <label for="login_field">
          Username or email address
        </label>
        <input autocapitalize="off" autocorrect="off" autofocus="autofocus" class="form-control input-block" id="login_field" name="login" tabindex="1" type="text"> == $0
        <label for="password">...</label>
        <input class="form-control form-control input-block" id="password" name="password" tabindex="2" type="password">
        <input class="btn btn-primary btn-block" data-disable-with="Signing in..." name="commit" tabindex="3" type="submit" value="Sign in">
      </div>
    </form>
    <p class="create-account-callout mt-3">...</p>
  </div>
  <div class="modal-backdrop js-touch-events"></div>
  </div>
  <div class="footer container-lg p-responsive py-6 mt-6 f6" role="contentinfo">...</div>
  <div id="ajax-error-message" class="ajax-error-message flash flash-...>...</div>
  <script crossorigin="anonymous" integrity="sha256-F7VsAbjYEuEdAvwD0jCP6snWeqx5tkyuQYm7fPXA5w0=" src="https://assets-cdn.github.com/assets/frameworks-17b56c0...js"></script>
  <script src="..." crossorigin="anonymous" integrity="sha256-...></script>
... div #js-pjax-container #login form div input#login_field.form-control.input-
```

Styles Event Listeners DOM Breakpoints Properties Accessibility Properties Aureli

Filter :hov .cls +

```
element.style {
}
.auth-form-body .input-block {
  margin-top: 5px;
  margin-bottom: 15px;
}
.input-block {
  display: block;
  width: 100%;
}
```



```
102 const scrape = async (options) => {
103     const browser = await puppeteer.launch(options.launch);
104     const page = await browser.newPage();
105
106     await page.setViewport(options.setViewport);
107     await login(page, options);
108     await search(page, options);
109
110     const numPages = await getNumPages(page, options);
111
112     let users = await numPages.times(async (index) => {
113         return await receiveUserMeta(index + 1, page, options);
114     });
115
116     browser.close();
117
118     return [].concat(...users);
119 }
```

Example 5

```
102 const scrape = async (options) => {
103     const browser = await puppeteer.launch(options.launch);
104     const page = await browser.newPage();
105
106     await page.setViewport(options.setViewport);
107     await login(page, options);
108     await search(page, options);
109
110     const numPages = await getNumPages(page, options);
111
112     let users = await numPages.times(async (index) => {
113         return await receiveUserMeta(index + 1, page, options);
114     });
115
116     browser.close();
117
118     return [].concat(...users);
119 }
```

Example 5

```
13 const login = async(page, options) => {
14     const CREDENTIALS = {
15         username: options.username,
16         password: options.password
17     };
18
19     const USERNAME_SELECTOR = '#login_field';
20     const PASSWORD_SELECTOR = '#password';
21     const BUTTON_SELECTOR = '#login > form > div.auth-form-body.mt-3 >
22
23     await page.goto('https://github.com/login');
24     // dom element selectors
25     await page.click(USERNAME_SELECTOR);
26     await page.keyboard.type(CREDENTIALS.username);
27     await page.click(PASSWORD_SELECTOR);
28     await page.keyboard.type(CREDENTIALS.password);
29     await page.click(BUTTON_SELECTOR);
30     await page.waitForNavigation();
31 }
32
```

Example 5


```
13  const login = async(page, options) => {
14      const CREDENTIALS = {
15          username: options.username,
16          password: options.password
17      };
18
19      const USERNAME_SELECTOR = '#login_field';
20      const PASSWORD_SELECTOR = '#password';
21      const BUTTON_SELECTOR = '#login > form > div.auth-form-body.mt-3 >
22
23      await page.goto('https://github.com/login');
24      // dom element selectors
25      await page.click(USERNAME_SELECTOR);
26      await page.keyboard.type(CREDENTIALS.username);
27      await page.click(PASSWORD_SELECTOR);
28      await page.keyboard.type(CREDENTIALS.password);
29      await page.click(BUTTON_SELECTOR);
30      await page.waitForNavigation();
31  }
32
```

Example 5

```
13 const login = async(page, options) => {
14     const CREDENTIALS = {
15         username: options.username,
16         password: options.password
17     };
18
19     const USERNAME_SELECTOR = '#login_field';
20     const PASSWORD_SELECTOR = '#password';
21     const BUTTON_SELECTOR = '#login > form > div.auth-form-body.mt-3 >
22
23     await page.goto('https://github.com/login');
24     // dom element selectors
25     await page.click(USERNAME_SELECTOR);
26     await page.keyboard.type(CREDENTIALS.username);
27     await page.click(PASSWORD_SELECTOR);
28     await page.keyboard.type(CREDENTIALS.password);
29     await page.click(BUTTON_SELECTOR);
30     await page.waitForNavigation();
31 }
32
```

Example 5

```
13 const login = async(page, options) => {
14     const CREDENTIALS = {
15         username: options.username,
16         password: options.password
17     };
18
19     const USERNAME_SELECTOR = '#login_field';
20     const PASSWORD_SELECTOR = '#password';
21     const BUTTON_SELECTOR = '#login > form > div.auth-form-body.mt-3 >
22
23     await page.goto('https://github.com/login');
24     // dom element selectors
25     await page.click(USERNAME_SELECTOR);
26     await page.keyboard.type(CREDENTIALS.username);
27     await page.click(PASSWORD_SELECTOR);
28     await page.keyboard.type(CREDENTIALS.password);
29     await page.click(BUTTON_SELECTOR);
30     await page.waitForNavigation();
31 }
32
```

Example 5

```
13 const login = async(page, options) => {
14     const CREDENTIALS = {
15         username: options.username,
16         password: options.password
17     };
18
19     const USERNAME_SELECTOR = '#login_field';
20     const PASSWORD_SELECTOR = '#password';
21     const BUTTON_SELECTOR = '#login > form > div.auth-form-body.mt-3 >
22
23     await page.goto('https://github.com/login');
24     // dom element selectors
25     await page.click(USERNAME_SELECTOR);
26     await page.keyboard.type(CREDENTIALS.username);
27     await page.click(PASSWORD_SELECTOR);
28     await page.keyboard.type(CREDENTIALS.password);
29     await page.click(BUTTON_SELECTOR);
30     await page.waitForNavigation();
31 }
32
```

Example 5

```
13 const login = async(page, options) => {
14     const CREDENTIALS = {
15         username: options.username,
16         password: options.password
17     };
18
19     const USERNAME_SELECTOR = '#login_field';
20     const PASSWORD_SELECTOR = '#password';
21     const BUTTON_SELECTOR = '#login > form > div.auth-form-body.mt-3 >
22
23     await page.goto('https://github.com/login');
24     // dom element selectors
25     await page.click(USERNAME_SELECTOR);
26     await page.keyboard.type(CREDENTIALS.username);
27     await page.click(PASSWORD_SELECTOR);
28     await page.keyboard.type(CREDENTIALS.password);
29     await page.click(BUTTON_SELECTOR);
30     await page.waitForNavigation();
31 }
32
```

Example 5

```
13 const login = async(page, options) => {
14     const CREDENTIALS = {
15         username: options.username,
16         password: options.password
17     };
18
19     const USERNAME_SELECTOR = '#login_field';
20     const PASSWORD_SELECTOR = '#password';
21     const BUTTON_SELECTOR = '#login > form > div.auth-form-body.mt-3 >
22
23     await page.goto('https://github.com/login');
24     // dom element selectors
25     await page.click(USERNAME_SELECTOR);
26     await page.keyboard.type(CREDENTIALS.username);
27     await page.click(PASSWORD_SELECTOR);
28     await page.keyboard.type(CREDENTIALS.password);
29     await page.click(BUTTON_SELECTOR);
30     await page.waitForNavigation();
31 }
32
```

Example 5


```
33  const getSearchUrl = (options) => {  
34      const userToSearch = options.search || 'john';  
35      return `https://github.com/search?q=${userToSearch}&type=Users`  
36  }  
37
```

```
33  const getSearchUrl = (options) => {  
34      const userToSearch = options.search || 'john';  
35      return `https://github.com/search?q=${userToSearch}&type=Users`  
36  }  
37
```

```

42  async function getNumPages(page, options) {
43      const NUM_USER_SELECTOR = '#js-pjax-container > div.container > c
44
45      let content = await page.content();
46      let $ = options.transform(content);
47
48      // format is: "69,803 users"
49      let inner = ($(NUM_USER_SELECTOR)
50                  .html() || '')
51                  .replace(',', '')
52                  .replace('users', '')
53                  .trim();
54      const numUsers = parseInt(inner);
55      /**
56       * GitHub shows 10 results per page, so
57       */
58      let num = Math.ceil(numUsers / 10);
59      return num > options.maxPage ? options.maxPage || 100 : num;
60  }

```

Example 5

```

62  const receiveUserMeta = async(index, page, options) => {
63      const LENGTH_SELECTOR_CLASS = '.user-list-item';
64
65      const LIST_USERNAME_SELECTOR = '.user-list-info.ml-2 > a';
66      const LIST_EMAIL_SELECTOR = '.octicon.octicon-mail + a.muted-link';
67
68      let pageUrl = getSearchUrl(options) + '&p=' + index;
69
70      await page.goto(pageUrl, {waitUntil: ['domcontentloaded']});
71
72      let content = await page.content();
73      let $ = options.transform(content);
74
75      return $(LENGTH_SELECTOR_CLASS).map(function() {
76          return {
77              username: $(this).find(LIST_USERNAME_SELECTOR).eq(0).text(),
78              email: $(this).find(LIST_EMAIL_SELECTOR).eq(0).text()
79          };
80      }).get().filter(e => e.email);
81  }
82

```

Example 5

```
62  const receiveUserMeta = async(index, page, options) => {
63      const LENGTH_SELECTOR_CLASS = '.user-list-item';
64
65      const LIST_USERNAME_SELECTOR = '.user-list-info.ml-2 > a';
66      const LIST_EMAIL_SELECTOR = '.octicon.octicon-mail + a.muted-link';
67
68      let pageUrl = getSearchUrl(options) + '&p=' + index;
69
70      await page.goto(pageUrl, {waitUntil: ['domcontentloaded']});
71
72      let content = await page.content();
73      let $ = options.transform(content);
74
75      return $(LENGTH_SELECTOR_CLASS).map(function() {
76          return {
77              username: $(this).find(LIST_USERNAME_SELECTOR).eq(0).text(),
78              email: $(this).find(LIST_EMAIL_SELECTOR).eq(0).text()
79          };
80      }).get().filter(e => e.email);
81  }
82
```

Example 5


```
62  const receiveUserMeta = async(index, page, options) => {
63      const LENGTH_SELECTOR_CLASS = '.user-list-item';
64
65      const LIST_USERNAME_SELECTOR = '.user-list-info.ml-2 > a';
66      const LIST_EMAIL_SELECTOR = '.octicon.octicon-mail + a.muted-link';
67
68      let pageUrl = getSearchUrl(options) + '&p=' + index;
69
70      await page.goto(pageUrl, {waitUntil: ['domcontentloaded']});
71
72      let content = await page.content();
73      let $ = options.transform(content);
74
75      return $(LENGTH_SELECTOR_CLASS).map(function() {
76          return {
77              username: $(this).find(LIST_USERNAME_SELECTOR).eq(0).text(),
78              email: $(this).find(LIST_EMAIL_SELECTOR).eq(0).text()
79          };
80      }).get().filter(e => e.email);
81  }
82
```

Example 5

```

62  const receiveUserMeta = async(index, page, options) => {
63      const LENGTH_SELECTOR_CLASS = '.user-list-item';
64
65      const LIST_USERNAME_SELECTOR = '.user-list-info.ml-2 > a';
66      const LIST_EMAIL_SELECTOR = '.octicon.octicon-mail + a.muted-link';
67
68      let pageUrl = getSearchUrl(options) + '&p=' + index;
69
70      await page.goto(pageUrl, {waitUntil: ['domcontentloaded']});
71
72      let content = await page.content();
73      let $ = options.transform(content);
74
75      return $(LENGTH_SELECTOR_CLASS).map(function() {
76          return {
77              username: $(this).find(LIST_USERNAME_SELECTOR).eq(0).text(),
78              email: $(this).find(LIST_EMAIL_SELECTOR).eq(0).text()
79          };
80      }).get().filter(e => e.email);
81  }
82

```

Example 5

```

62  const receiveUserMeta = async(index, page, options) => {
63      const LENGTH_SELECTOR_CLASS = '.user-list-item';
64
65      const LIST_USERNAME_SELECTOR = '.user-list-info.ml-2 > a';
66      const LIST_EMAIL_SELECTOR = '.octicon.octicon-mail + a.muted-link';
67
68      let pageUrl = getSearchUrl(options) + '&p=' + index;
69
70      await page.goto(pageUrl, {waitUntil: ['domcontentloaded']});
71
72      let content = await page.content();
73      let $ = options.transform(content);
74
75      return $(LENGTH_SELECTOR_CLASS).map(function() {
76          return {
77              username: $(this).find(LIST_USERNAME_SELECTOR).eq(0).text(),
78              email: $(this).find(LIST_EMAIL_SELECTOR).eq(0).text()
79          };
80      }).get().filter(e => e.email);
81  }
82

```

Example 5

```

62  const receiveUserMeta = async(index, page, options) => {
63      const LENGTH_SELECTOR_CLASS = '.user-list-item';
64
65      const LIST_USERNAME_SELECTOR = '.user-list-info.ml-2 > a';
66      const LIST_EMAIL_SELECTOR = '.octicon.octicon-mail + a.muted-link';
67
68      let pageUrl = getSearchUrl(options) + '&p=' + index;
69
70      await page.goto(pageUrl, {waitUntil: ['domcontentloaded']});
71
72      let content = await page.content();
73      let $ = options.transform(content);
74
75      return $(LENGTH_SELECTOR_CLASS).map(function() {
76          return {
77              username: $(this).find(LIST_USERNAME_SELECTOR).eq(0).text(),
78              email: $(this).find(LIST_EMAIL_SELECTOR).eq(0).text()
79          };
80      }).get().filter(e => e.email);
81  }
82

```

Example 5

```

1 import puppeteer from 'puppeteer';
2 import cheerio from 'cheerio';
3
4 const timesFunction = function(callback) {
5   // ...
6 };
7
8 String.prototype.times = timesFunction;
9 Number.prototype.times = timesFunction;
10
11 const options = {
12   url: 'https://www.github.com',
13   transform: (body) => cheerio.load(body)
14 };
15
16 const login = async (page) => {
17   const CREDENTIALS = {
18     username: 'eeee',
19     password: 'eeee'
20   };
21
22   const USERNAME_SELECTOR = '#login_field';
23   const PASSWORD_SELECTOR = '#password';
24   const LOGIN_BUTTON_SELECTOR = '#login > form > div.auth-form-body.st-3 > input.btn.btn-primary.btn-block';
25
26   await page.goto('https://github.com/login');
27   // dom element selectors
28   await page.click(USERNAME_SELECTOR);
29   await page.keyboard.type(CREDENTIALS.username);
30   await page.click(PASSWORD_SELECTOR);
31   await page.keyboard.type(CREDENTIALS.password);
32   await page.click(LOGIN_BUTTON_SELECTOR);
33   await page.waitForNavigation();
34 }
35
36 const search = async (page) => {
37   const userToSearch = 'john';
38   const searchUrl = `https://github.com/search?q=${userToSearch}&type=Users&utf8=%E2%9C%93`;
39   await page.goto(searchUrl);
40   await page.waitFor(2 * 1000);
41 }
42
43 async function getNumPages(page, options) {
44   const NUM_USER_SELECTOR = '#js-pjax-container > div.container > div > div.col.col-three-fourth.codesearch-results.pr-6 > div.d-flex.flex-justify-between.border-bottom.pb-3 > h3';
45   let content = await page.content();
46   let $ = options.transform(content);
47
48   // format is "69,883 users"
49   let inner = $(NUM_USER_SELECTOR)
50     .html()
51     .replace(/,/g, '')
52     .replace('users', '')
53     .trim();
54
55   const numUsers = parseInt(inner);
56   console.log('numUsers: ', numUsers);
57   // github shows 18 results per page, so
58   //
59   return Math.ceil(numUsers / 18);
60 }
61
62 const receiveUserData = async (index, page, options) => {
63   const LIST_USERNAME_SELECTOR = '#user_search_results > div.user-list > div:nth-child(INDEX) > div.d-flex > div > a';
64   const LIST_EMAIL_SELECTOR = '#user_search_results > div.user-list > div:nth-child(INDEX) > div.d-flex > div > ul > li:nth-child(2) > a';
65   const LENGTH_SELECTOR_CLASS = 'user-list-item';
66
67   let pageUrl = searchUrl + '&p=' + index;
68   await page.goto(pageUrl);
69   let content = await page.content();
70   let $ = options.transform(content);
71
72   let listLength = $el.length;
73
74   return listLength.times((index) => {
75     let usernameSelector = LIST_USERNAME_SELECTOR.replace("INDEX", index + 1);
76     let emailSelector = LIST_EMAIL_SELECTOR.replace("INDEX", index + 1);
77
78     let username = $(usernameSelector).attr('href').replace('/', '');
79     let email = $(emailSelector).html();
80     // not all users have email visible
81     if (!email)
82       return null;
83
84     return {
85       username: username,
86       email: email
87     };
88   });
89 }
90
91 const scrape = async () => {
92   const browser = await puppeteer.launch();
93   const page = await browser.newPage();
94
95   await login(page);
96   await search(page);
97
98   const numPages = await getNumPages(page, options);
99   let users = await Promise.all(numPages.times((index) => {
100     return receiveUserData(index + 1, page, options);
101   }));
102
103   browser.close();
104
105   return users;
106 }
107
108 scrape(options);
109
110
111
112
113
114
115
116
117
118
119
120
121

```

All together 121 line of code


```
1 import scrape from './index';
2 import config from '../config';
3 jasmine.DEFAULT_TIMEOUT_INTERVAL = 999999;
4 describe('puppeteer login github', () => {
  Inconclusive | Debug
5   test('should return array', async () => {
6     let users = await scrape({
7       username: config.github.username,
8       password: config.github.password,
9       search: "bitcoin",
10      maxPage: 20
11    });
12
13    const resp = {
14      username: expect.stringMatching(/.**/),
15      email: expect.stringMatching(/@/),
16    };
17
18    expect(Array.isArray(users)).toBe(true)
19    expect(users.length).toBeGreaterThan(0);
20
21    return expect(users[0]).toMatchObject(resp);
22  });
23 });
24
```

```
# macbook at abc in ~/sandbox/presentation/presentation-one on git:master x [1:45:12]
→ npm test ./src/puppeteer-login-github
```


index.js .../puppeteer-login-github

JS index.test.js .../puppeteer-login-github x

1: node

```

1  import scrape from './index';
2  import config from '../config';
3  jasmine.DEFAULT_TIMEOUT_INTERVAL = 999999;
   1 Inconclusive
4  describe('puppeteer login github', () => {
   Inconclusive | Debug
5    test('should return array', async () => {
6      let users = await scrape({
7        username: config.github.username,
8        password: config.github.password,
9        search: "bitcoin",
10       maxPage: 20
11     });
12
13     const resp = {
14       username: expect.stringMatching(/.**/),
15       email: expect.stringMatching(/@/),
16     };
17
18     expect(Array.isArray(users)).toBe(true)
19     expect(users.length).toBeGreaterThan(0);
20
21     console.log(users);
22     return expect(users[0]).toMatchObject(resp);
23   });
24 });
25

```

```

# macbook at abc in ~/sandbox/presentation/presentation-one on git:master * [1:45:12]
+ npm test ./src/puppeteer-login-github

```

```

> crawler-tests@1.0.0 test /Users/macbook/sandbox/presentation/presentation-one
> jest "./src/puppeteer-login-github"

```

```

PASS src/puppeteer-login-github/index.test.js (20.401s)
  puppeteer login github
    ✓ should return array (19915ms)

```

```

Test Suites: 1 passed, 1 total
Tests:       1 passed, 1 total
Snapshots:  0 total
Time:       20.924s

```

```

Ran all test suites matching /.\src\puppeteer-login-github/i.

```

```

# macbook at abc in ~/sandbox/presentation/presentation-one on git:master * [1:46:31]
+ npm test ./src/puppeteer-login-github

```

```

> crawler-tests@1.0.0 test /Users/macbook/sandbox/presentation/presentation-one
> jest "./src/puppeteer-login-github"

```

```

RUNS ...puppeteer-login-github/index.test.js

```



A famous example of that is LinkedIn. Setting the **li_atcookie** will guarantee your scraper bot access to their social network (please note: we encourage you to respect your target website ToS).

<https://phantombuster.com/>

```
import puppeteer from 'puppeteer';

const scrape = async () => {
  const browser = await puppeteer.launch();
  const page = await browser.newPage();

  await page.setCookie({
    name: "li_at",
    value: "a session cookie value copied from your DevTools",
    domain: "www.linkedin.com"
  });
  // make something awesome here
}
```

page.setCookie(...cookies)

- ...cookies <...Object>
 - name <string> required
 - value <string> required
 - url <string>
 - domain <string>
 - path <string>
 - expires <number> Unix time in seconds.
 - httpOnly <boolean>
 - secure <boolean>
 - sameSite <string> "Strict" or "Lax".
- returns: <Promise>

<https://phantombuster.com/>

The screenshot shows the Chrome DevTools Network tab. The top panel displays a list of network requests with a time scale from 5000ms to 70000ms. The selected request is a PDF file download. The details panel on the right shows the following information:

- General:**
 - Request URL: `https://www.linkedin.com/dms/35540F-a3Y2_a9zHg/profile-profilePdf/G/m-156486453&e-1503999959&v-alpha&t-vfr1aXg_ot6TEkcYGS_ZP5CUmOHZsgD30qnmBBlayj0`
 - Request Method: `GFT`
 - Status Code: `200`
 - Remote Address: `105.63.144.1:440`
 - Referrer Policy: `no-referrer-when-downgrade`
- Response Headers:**
 - `cache-control`: `no-cache, no-store`
 - `content-disposition`: `attachment; filename="ClémentCombesProfile.pdf"`
 - `content-type`: `application/octet-stream`
 - `date`: `Mon, 28 Aug 2017 09:45:59 GMT`
 - `expires`: `Thu, 01 Jan 1970 00:00:00 GMT`
 - `pragma`: `no-cache`
 - `status`: `200`
 - `strict-transport-security`: `max-age=2592000`
 - `x-ambry-blob-size`: `12797`
 - `x-content-type-options`: `nosniff`
 - `x-frame-options`: `sameorigin`
 - `x-li-fabric`: `prod-lor1`

Because we're using the DevTools API, the code we write has the equivalent power of a human using Chrome's DevTools. That means your bot can intercept, examine and even modify or abort any network request.

```
10 //         cookie: config.linkedIn.token,
11 //     });
12
13 //     return expect($(".[data-control-name='identity_welcome_me
14 //     });
15 // });
16
17
```

1 Inconclusive

```
18 describe('puppeteer github cookie', () => {
  Inconclusive | Debug
19   test('should login using mine cookie', async () => {
20     let $ = await scrape({
21       uri: "https://github.com/search?q=bitcoin&type=Users&ut
22       cookie: config.github.token,
23       setCookie: {
24         name: "user_session",
25         domain: "github.com"
26       },
27     });
28
29     let nickname = $(".css-truncate-target").text();
30     console.log(nickname);
31
32     return expect(nickname).toMatch(/wegorich/);
33   });
34 });
35
```

```
# macbook at abc in ~/sandbox/presentation
n/presentation-one on git:master * [2:04:
37]
```

```
→ █
```




AdBlock

In the same vein, we can speed up our scraping by blocking unnecessary requests. Analytics, ads and images are typical targets. However, you have to keep in mind that it will make your bot less human-like (for example LinkedIn will not serve their pages properly if you block all images—we're not sure if it's deliberate or not).

```
1 import puppeteer from 'puppeteer';
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   launch: { headless: true },
6   blacklist: [
7     /.*collector\.githubapp.*/,
8     /.*fsitouchzoom\.js/,
9     /.*api\.github\.com\/_private\/browser.*/,
10    /.*google.*/,
11  ],
12  setViewport: { width: 1240, height: 680 },
13  setCookie: {
14    name: "li_at",
15    domain: "www.linkedin.com"
16  },
17  transform: (body) => cheerio.load(body)
18 };
```

```
1 import puppeteer from 'puppeteer';
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3
4 const options = {
5   launch: { headless: true },
6   blacklist: [
7     /.*/collector\.githubapp.*\/,
8     /.*/fsitouchzoom\.js\/,
9     /.*/api\.github\.com\/_private\/browser.*\/,
10    /.*/google.*\/
11  ],
12   setViewport: { width: 1240, height: 680 },
13   setCookie: {
14     name: "li_at",
15     domain: "www.linkedin.com"
16   },
17   transform: (body) => cheerio.load(body)
18 };
```

Example 8

```
20  const writeFileInterceptor = ({blacklist}) => (e) => {
21      if (blacklist.find(item => item.test(e.url))) {
22          e.abort();
23      } else {
24          e.continue();
25      }
26  }
```

```
28  const scrape = async (options) => {
29      const browser = await puppeteer.launch(options.launch);
30      const page = await browser.newPage();
31
32      await page.setViewport(options.setViewport);
33      await page.setRequestInterception(true);
34      page.on('request', writeFileInterceptor(options));
35      // "a session cookie value copied from your DevTools",
36      options.setCookie.value = options.cookie || options.setCookie.value;
37      await page.setCookie(options.setCookie);
38      await page.goto(options.uri, {waitUntil: ['domcontentloaded']});
39
40      let content = await page.content();
41      let $ = options.transform(content);
42      browser.close();
43      return $;
44      // make something awesome here
45  };
```

Example 8


```
26   let start = null;
27
28   // first item loads 100ms longer usually
29   // and 1.2s longer first run
30
31   start = now();
32   await scrape(Object.assign({}, options, {blacklist: []}));
33   let speedDefault = now() - start;
34
35   start = now();
36   await scrape(options);
37   let speedAdBlock = now() - start;
38
39   console.log(`AdBlocked speed - ${speedAdBlock}, by default - ${speedDefault}`);
40   return expect(speedAdBlock).toBeLessThanOrEqual(speedDefault);
```

```
26   let start = null;
27
28   // first item loads 100ms longer usually
29   // and 1.2s longer first run
30
31   start = now();
32   await scrape(Object.assign({}, options, {blacklist: []}));
33   let speedDefault = now() - start;
34
35   start = now();
36   await scrape(options);
37   let speedAdBlock = now() - start;
38
39   console.log(`AdBlocked speed - ${speedAdBlock}, by default - ${speedDefault}`);
40   return expect(speedAdBlock).toBeLessThanOrEqual(speedDefault);
```

Example 8

```
26     let start = null;
27
28     // first item loads 100ms longer usually
29     // and 1.2s longer first run
30
31     start = now();
32     await scrape(Object.assign({}, options, {blacklist: []}));
33     let speedDefault = now() - start;
34
35     start = now();
36     await scrape(options);
37     let speedAdBlock = now() - start;
38
39     console.log(`AdBlocked speed - ${speedAdBlock}, by default - ${speedDefault}`);
40     return expect(speedAdBlock).toBeLessThanOrEqual(speedDefault);
```

Example 8

```

1  import scrape from './index';
2  import config from '../config';
3  import now from 'performance-now'; 2.7K (gzipped: 1.1K)
4
5  jasmine.DEFAULT_TIMEOUT_INTERVAL = 999999;
6
7  describe('puppeteer adblock', () => {
8      Inconclusive | Debug
9      test('adblock should be faster', async () => {
10         let options = {
11             uri: "https://github.com/search?q=bitcoin&type=Users&ut",
12             cookie: config.github.token,
13             setCookie: {
14                 name: "user_session",
15                 domain: "github.com"
16             },
17             blacklist: [
18                 /.assets-cdn.github./,
19                 /.githubusercontent./,
20                 /.collector.githubapp./,
21                 /.fsitouchzoom.js/,
22                 /.api.github.com/_private/browser./,
23                 /.google./
24             ]
25         };
26         let start = null;

```

```

# macbook at abc in ~/sandbox/presentation/presentation-one on git:master * [2:22:05]

```

```

→ npm test ./src/puppeteer-adblock-demo

```

```

}

```

EasyList

```
[Adblock Plus 2.0]
! Version: 201712102233
! Title: EasyList
! Last modified: 10 Dec 2017 22:33 UTC
! Expires: 4 days (update frequency)
! Homepage: https://easylist.to/
! Licence: https://easylist.to/pages/licence.html
!
! Please report any unblocked adverts or problems
! in the forums (https://forums.lanik.us/)
! or via e-mail (easylist.subscription@gmail.com).
! GitHub issues: https://github.com/easylist/easylist/issues
! GitHub pull requests: https://github.com/easylist/easylist/pulls
!
! -----General advert blocking filters-----
! *** easylist:easylist/easylist_general_block.txt ***
&act=ads_
&ad.vid=$~xmlhttprequest
&ad_box_
&ad_channel=
&ad_classid= https://easylist.to/easylist/easylist.txt
&ad_height=
&ad_ids=
&ad_keyword=
```


abp-filter-parser

JavaScript Adblock Plus filter parser for lists like EasyList.

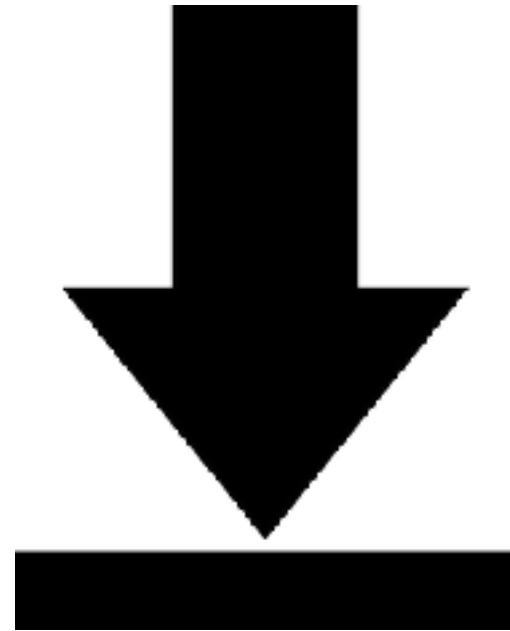
Parses filter rules as per:

<https://adblockplus.org/en/filters>

<https://adblockplus.org/en/filter-cheatsheet>

`npm install abp-filter-parser`

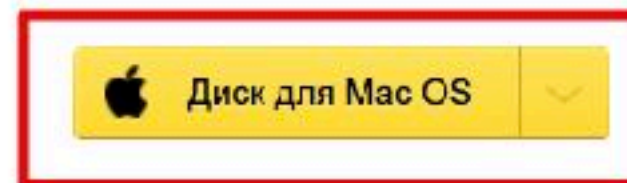
Download File





Установите Диск для Mac

Папка Яндекс.Диска выглядит так же, как обычная папка на компьютере, но хранит ваши файлы ещё и на сервере Яндекса. Они доступны только вам и тем, кого вы сами выберете. А добраться до них можно с любого устройства с интернетом.



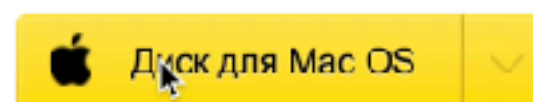
[Узнайте о бета-версии экспериментального Яндекс.Диска 2.0](#)

```
15  const writeFileInterceptor = ({file})=> (e) => {
16    |   if (file.regex.test(e.url)) {
17    |     |   e.buffer().then(buffer => {
18    |     |     |   fs.writeFileSync(file.path, Buffer.from(buffer, 'base64'));
19    |     |     |   });
20    |     |   }
21  |   }
```



Установите Диск для Mac

Папка Яндекс.Диска выглядит так же, как обычная папка на компьютере, но хранит ваши файлы ещё и на сервере Яндекса. Они доступны только вам и тем, кого вы сами выберете. А добраться до них можно с любого устройства с интернетом.



[Узнайте о бета-версии экспериментального Яндекс.Диска 2.0](#)

Request.js

+

Puppeteer.js

page.setContent(html)

- `html` `<string>` HTML markup to assign to the page.
- returns: `<Promise>`

page.goto(url, options)

- `url` `<string>` URL to navigate page to. The url should include scheme, e.g. `https://`.
- `options` `<Object>` Navigation parameters which might have the following properties:
 - `timeout` `<number>` Maximum navigation time in milliseconds, defaults to 30 seconds
 - `waitFor` `<string|Array<string>>` When to consider navigation succeeded, default is `load`. If `waitFor` is a string, navigation is considered to be successful after all events have been fired. Events are:
 - `load` - consider navigation to be finished when the `load` event is fired.
 - `domcontentloaded` - consider navigation to be finished when the `DOMContentLoaded` event is fired.
 - `networkidle0` - consider navigation to be finished when there are no more than 0 network connections. Default is 0.
 - `networkidle2` - consider navigation to be finished when there are no more than 2 network connections.

```
1 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
2 import { URL } from 'url';
3 import requestDriver from './request-driver';
4 import puppeteerDriver from './puppeteer-driver';
5
6 const options = {
7   uri: 'https://ab.onliner.by/',
8   selector: '.autoba-table-imp h2 strong',
9   launch: { headless: true },
10  blacklist: [
11    /*.*css/,
12    /*.*google.*/
13  ],
14  setViewport: { width: 1240, height: 680 },
15  setCookie: {
16    name: "li_at",
17    value: "10",
18    domain: "www.linkedin.com"
19  },
20  isDynamic: false,
21  transform: (body) => cheerio.load(body)
22 };
23
```

```
1 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
2 import { URL } from 'url';
3 import requestDriver from './request-driver';
4 import puppeteerDriver from './puppeteer-driver';
5
6 const options = {
7   uri: 'https://ab.onliner.by/',
8   selector: '.autoba-table-imp h2 strong',
9   launch: { headless: true },
10  blacklist: [
11    /*.*css/,
12    /*.*google.*/
13  ],
14  setViewport: { width: 1240, height: 680 },
15  setCookie: {
16    name: "li_at",
17    value: "10",
18    domain: "www.linkedin.com"
19  },
20  isDynamic: false,
21  transform: (body) => cheerio.load(body)
22 };
23
```

```
1 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
2 import { URL } from 'url';
3 import requestDriver from './request-driver';
4 import puppeteerDriver from './puppeteer-driver';
5
6 const options = {
7   uri: 'https://ab.onliner.by/',
8   selector: '.autoba-table-imp h2 strong',
9   launch: { headless: true },
10  blacklist: [
11    /*.*css/,
12    /*.*google.*/
13  ],
14  setViewport: { width: 1240, height: 680 },
15  setCookie: {
16    name: "li_at",
17    value: "10",
18    domain: "www.linkedin.com"
19  },
20  isDynamic: false,
21  transform: (body) => cheerio.load(body)
22 };
23
```



```
24  const isDynamicUrls = new Map();
25
26  const scrape = async (options) => {
27      let origin = (new URL(options.uri)).origin;
28      let $ = null;
29      if (options.isDynamic) {
30          isDynamicUrls.set(origin, true);
31      }
32
33      if (!isDynamicUrls.has(origin)) {
34          $ = await requestDriver(options);
35          let items = $(options.selector);
36          if (items.length) {
37              return items.map(function() {return $(this).text();}).get();
38          }
39
40          isDynamicUrls.set(origin, true);
41      }
42
43      $ = await puppeteerDriver(options);
44      return $(options.selector).map(function() {return $(this).text();}).get();
45  };
```

```
24  const isDynamicUrls = new Map();
25
26  const scrape = async (options) => {
27      let origin = (new URL(options.uri)).origin;
28      let $ = null;
29      if (options.isDynamic) {
30          isDynamicUrls.set(origin, true);
31      }
32
33      if (!isDynamicUrls.has(origin)) {
34          $ = await requestDriver(options);
35          let items = $(options.selector);
36          if (items.length) {
37              return items.map(function() {return $(this).text();}).get();
38          }
39
40          isDynamicUrls.set(origin, true);
41      }
42
43      $ = await puppeteerDriver(options);
44      return $(options.selector).map(function() {return $(this).text();}).get();
45  };
```

```
24  const isDynamicUrls = new Map();
25
26  const scrape = async (options) => {
27      let origin = (new URL(options.uri)).origin;
28      let $ = null;
29      if (options.isDynamic) {
30          isDynamicUrls.set(origin, true);
31      }
32
33      if (!isDynamicUrls.has(origin)) {
34          $ = await requestDriver(options);
35          let items = $(options.selector);
36          if (items.length) {
37              return items.map(function() {return $(this).text();}).get();
38          }
39
40          isDynamicUrls.set(origin, true);
41      }
42
43      $ = await puppeteerDriver(options);
44      return $(options.selector).map(function() {return $(this).text();}).get();
45  };
```

```
24  const isDynamicUrls = new Map();
25
26  const scrape = async (options) => {
27      let origin = (new URL(options.uri)).origin;
28      let $ = null;
29      if (options.isDynamic) {
30          isDynamicUrls.set(origin, true);
31      }
32
33      if (!isDynamicUrls.has(origin)) {
34          $ = await requestDriver(options);
35          let items = $(options.selector);
36          if (items.length) {
37              return items.map(function() {return $(this).text();}).get();
38          }
39
40          isDynamicUrls.set(origin, true);
41      }
42
43      $ = await puppeteerDriver(options);
44      return $(options.selector).map(function() {return $(this).text();}).get();
45  };
```

```
24  const isDynamicUrls = new Map();
25
26  const scrape = async (options) => {
27      let origin = (new URL(options.uri)).origin;
28      let $ = null;
29      if (options.isDynamic) {
30          isDynamicUrls.set(origin, true);
31      }
32
33      if (!isDynamicUrls.has(origin)) {
34          $ = await requestDriver(options);
35          let items = $(options.selector);
36          if (items.length) {
37              return items.map(function() {return $(this).text();}).get();
38          }
39
40          isDynamicUrls.set(origin, true);
41      }
42
43      $ = await puppeteerDriver(options);
44      return $(options.selector).map(function() {return $(this).text();}).get();
45  };
```




Автобарахолка



[Пройти диагностику](#)

Разместить объявление
всего за 90 секунд!

Актуальные

Последние добавленные

Год

Пробег

Цена, р.

Цена \$ € р.

Любая Любая

Обмен

С диагностикой **NEW**

Местонахождение

Все страны

Все области

Все города

Растаможен

Марка

Все марки

Все модели

[Добавить марку](#)

Тип кузова

Седан

Универсал

Хетчбек

Минивэн

АВТОРАЗДЕЛЫ В УСЛУГАХ

[Ремонт и обслуживание авто](#)

[Перетяжка и ремонт салона](#)

[Шиномонтаж](#)

[Автомойка](#)

[Химчистка](#)

[Другие автоуслуги](#)

[Все разделы Услуг](#)

Все расчеты на территории Беларуси между гражданами Республики Беларусь осуществляются исключительно в белорусских рублях. Цена в иностранной валюте (в случае ее наличия в объявлении граждан) указана справочно.



Автобарахолка

[Пройти диагностику](#)

Разместить объявление всего за 90 секунд!

32193 объявления

Актуальные [Последние добавленные](#)

Год Пробег Цена, р.

Цена \$ € р.

Любая Любая

- Обмен 14613
- С диагностикой 11 NEW

Местонахождение

- Все страны
- Все области
- Все города

Растаможен 31824

Марка

- Все марки
- Все модели

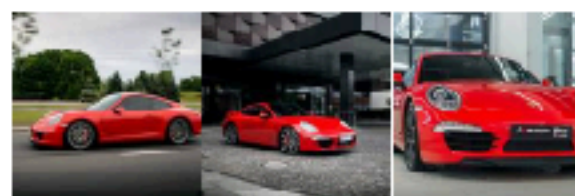
[Добавить марку](#)

Тип кузова

- Седан 11392
- Универсал 4233
- Хэтчбек 6493
- Минивэн 2411

★ **Porsche 911 Carrera S**

2013 72 000 км 159 881



купе, бензин 3.8 л, автомат, задний привод, кондиционер, кожаный салон, легкосплавные диски, ксенон, парктроник, подогрев сидений, система контроля стабилизации, навигация, громкая связь.

📷 17 UP! 2 часа назад

78 500 \$
66 810 €

★ **BMW X6 XDRIVE 30D**

2015 26 307 км 119 962



внедорожник, дизель 3 л, автомат, полный привод, кондиционер, кожаный салон, легкосплавные диски, ксенон, парктроник, подогрев сидений, система контроля стабилизации, навигация, громкая связь.

📷 8 UP! 4 часа назад

58 900 \$
50 129 €

ОБМЕН

★ **Mini Cooper**

2016 22 000 км 52 529



хэтчбек, бензин 1.5 л, автомат, передний привод, кондиционер, легкосплавные диски, ксенон, парктроник, подогрев сидений, система контроля стабилизации.

📷 8 UP! 4 часа назад

25 792 \$
21 950 €

ОБМЕН

★ **BMW X5 50i**

2010 118 992 км 47 659

23 400 \$

АВТОРАЗДЕЛЫ В УСЛУГАХ

- Ремонт и обслуживание авто
- Перетяжка и ремонт салона
- Шиномонтаж
- Автомойка
- Химчистка
- Другие автоуслуги

[Все разделы Услуг](#)

Все расчеты на территории Беларуси между гражданами Республики Беларусь осуществляются исключительно в белорусских рублях. Цена в иностранной валюте (в случае ее наличия в объявлении граждан) указана справочно.

```
14 test('the code should detect dynamic pages, the cost should be permanent',
15     let urls = [
16         { uri: "https://ab.onliner.by/" },
17         { uri: 'https://ab.onliner.by/' },
18         { uri: 'https://ab.onliner.by/', isDynamic: true },
19     ];
20
21     let otherResults = await mapSeries(urls, async (e) => {
22         return await getPerformace(() => scrape(e));
23     });
24
25     console.log(otherResults.map(e => e.speed));
26
```

Inconclusive | Debug

```
36 test('static sites should be superfast', async () => {
37   let urls = [
38     { uri: "https://ab.onliner.by/", selector: '.project-navigation__sign' },
39     { uri: 'https://ab.onliner.by/', selector: '.project-navigation__sign' },
40     { uri: 'https://ab.onliner.by/', selector: '.project-navigation__sign', isDynamic: true },
41   ];
42
43   let otherResults = await mapSeries(urls, async (e) => {
44     return await getPerfomance(() => scrape(e));
45   });
46
47   console.log(otherResults.map(e => e.speed));
48 }
```



```


6
7  const getPerformace = async (fn)=> {
8      let start = now()
9      let result = await fn();
10     return { result, speed: now() - start};
11 }
12
13 describe('puppeteer should work with request js', () => {
14     Inconclusive | Debug
15     test('the code should detect dynamic pages, the cost should be perm
16         let urls = [
17             { uri: "https://ab.onliner.by/" },
18             { uri: 'https://ab.onliner.by/' },
19             { uri: 'https://ab.onliner.by/', isDynamic: true },
20         ];
21
22         let otherResults = await mapSeries(urls, async (e)=> {
23             return await getPerformace(()=> scrape(e));
24         });
25
26         console.log(otherResults.map(e => e.speed));
27
28         otherResults.map(e=> e.speed - otherResults[0].speed).forEach(e
29             expect(e).toBeLessThanOrEqual(1500);
30         });
31
32         otherResults.map(e=> e.result).forEach(result=> {
33             expect(result.length).toBeGreaterThanOrEqual(1);

```

```

# macbook at abc in ~/sandbox/presen
tation/presentation-one on git:maste
r * [3:00:33]
→ npm test ./src/request-puppeteer-t
ogether-demo
█

```

Speed Up Puppeteer.js

```
11 let browser = null;
12 let browserTimeout = null;
13
14 const createBrowser = async ({launch})=> {
15     clearTimeout(browserTimeout);
16     if (!browser) {
17         browser = await puppeteer.launch(launch);
18     }
19
20     return browser;
21 }
22
23 const closeBrowser = ()=> {
24     clearTimeout(browserTimeout);
25     browserTimeout = setTimeout(() => {
26         browser.close();
27         browser = null;
28     }, 15000);
29 }
```


Cost

Note: When you install Puppeteer, it downloads a recent version of Chromium (~71Mb Mac, ~90Mb Linux, ~110Mb Win) that is guaranteed to work with the API.

Limitations

Encodings - Your text editor, browser, word processor or whatever else that's trying to read the document is assuming the wrong encoding (windows-1251).

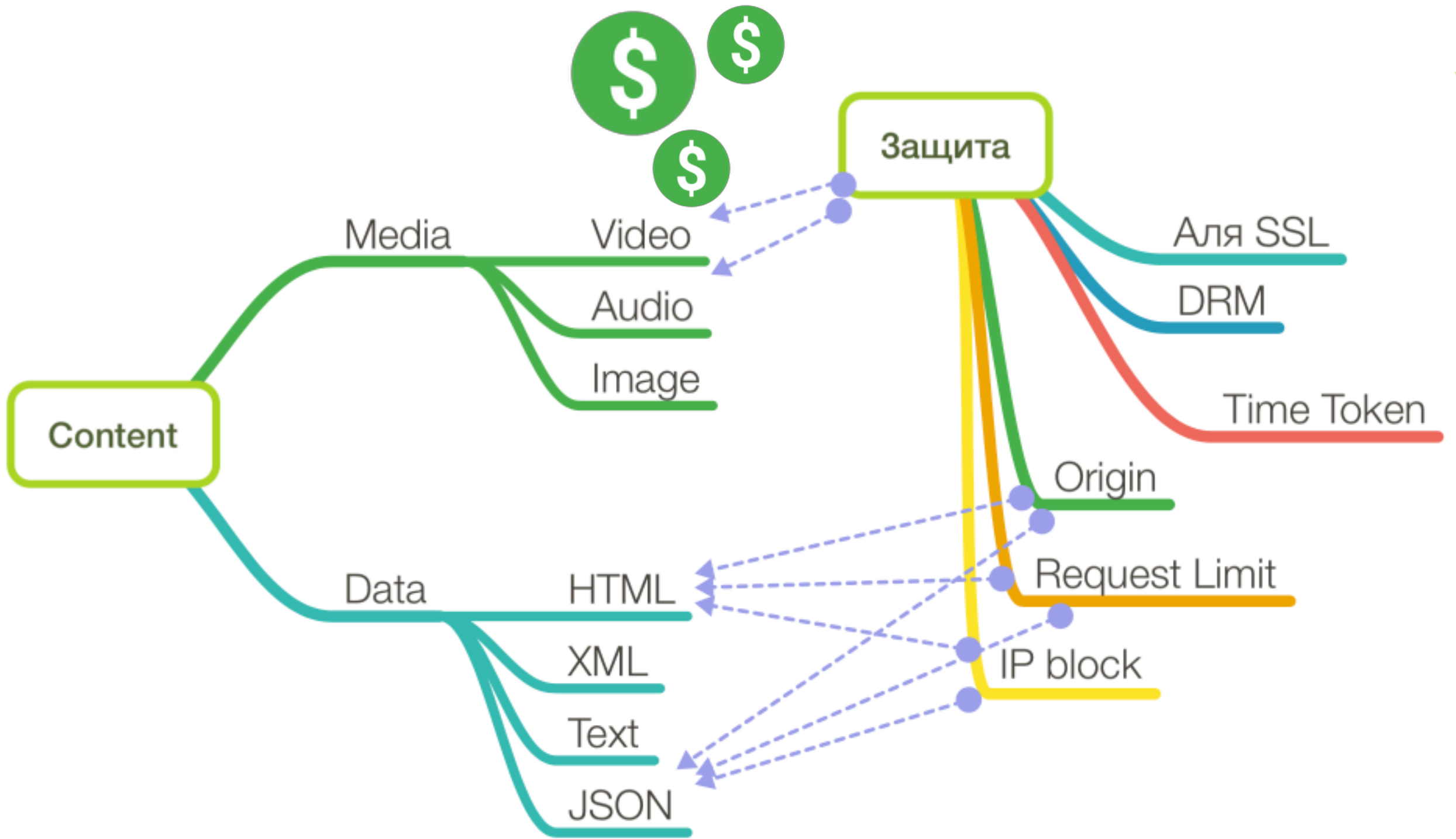
Query limits:

Rule 1: Every request per user increments an internal count. When the count exceeds the limit, the requests are denied with a HTTP 429 Too Many Requests

Rule 2: The only way for count to go away, is for an internal expiration time to expire, called the expiry, and is measured in seconds

Content





Text



Это никому не интересно, ну серьезно

Images



You could simply use this expression to match an img tag as in the example :

```
<img([\w\W]+?)/>
```

Audio



Connect on SoundCloud

Discover, stream, and share a constantly expanding mix of music from emerging and major artists around the world.

Sign up for free

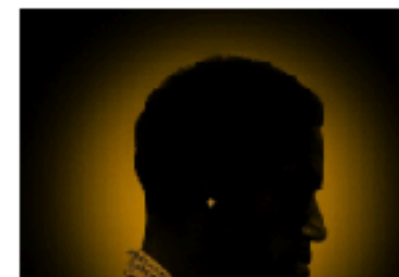
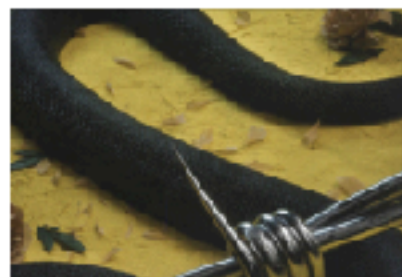
Search for artists, bands, tracks, podcasts



or

Upload your own

Hear what's trending for free in the SoundCloud community





JavaScript is disabled

You need to enable JavaScript to use SoundCloud

Show me how to enable it

rything

Audio - using soundcloud.com as example

`page.waitFor(selectorOrFunctionOrTimeout[, options])`

- `selectorOrFunctionOrTimeout` <string|number|function> A `selector`, predicate or timeout to `wait` for
- `options` <Object> Optional `waiting` parameters
- returns: <Promise>

This method behaves differently with respect to the type of the first parameter:

- if `selectorOrFunctionOrTimeout` is a `string`, then the first argument is treated as a `selector` to `wait` for and the method is a shortcut for `page.waitForSelector`
- if `selectorOrFunctionOrTimeout` is a `function`, then the first argument is treated as a predicate to `wait` for and the method is a shortcut for `page.waitForFunction()`.
- if `selectorOrFunctionOrTimeout` is a `number`, then the first argument is treated as a timeout in milliseconds and the method returns a promise which resolves after the timeout
- otherwise, an exception is thrown

Shortcut for `page.mainFrame().waitFor(selectorOrFunctionOrTimeout[, options])`.

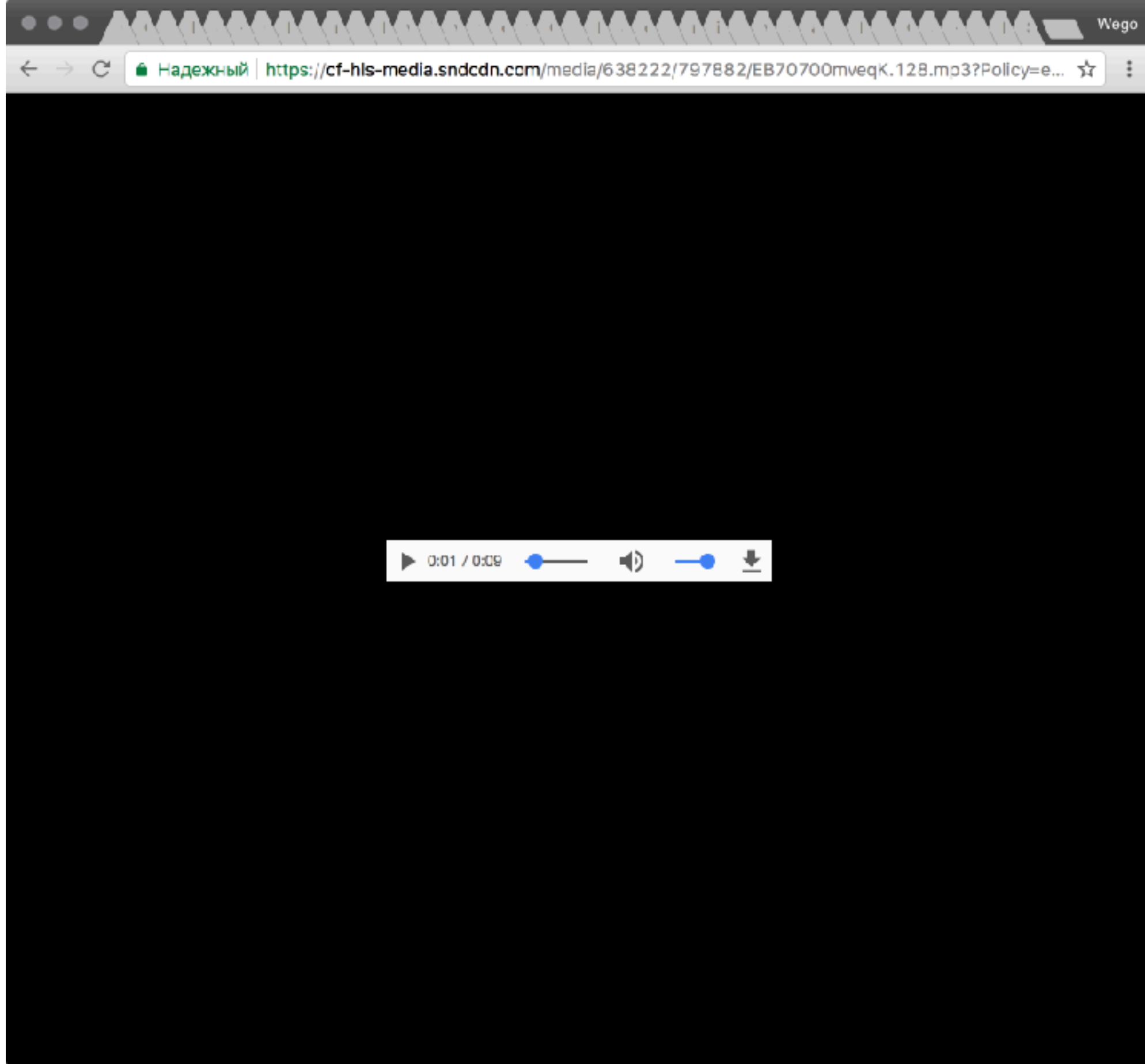
Audio - using [soundcloud.com](https://www.soundcloud.com) as example


```

▼ <article itemscope itemtype="http://schema.org/MusicRecording">
  ▼ <header>
    ▶ <h1 itemprop="name">...</h1>
      "
        published on "
        <time pubdate>2015-07-27T11:54:04Z</time>
        <meta itemprop="duration" content="PT00H05M41S">
        <meta itemprop="genre" content="Remix">
        <meta itemprop="interactionCount" content="UserLikes:263437">
        <meta itemprop="interactionCount" content="UserDownloads:0">
        <meta itemprop="interactionCount" content="UserComments:1604">
    ▼ <div itemscope itemprop="audio" itemtype="http://schema.org/
      AudioObject">
        <meta itemprop="embedUrl" content="https://w.soundcloud.com/
          player/?
          url=https%3A%2F%2Fapi.soundcloud.com%2Ftracks%2F216542711&auto_p
          lay=false&show_artwork=true&visual=true&origin=schema.org"> == $
        <meta itemprop="height" content="400px">
      </div>
    ▶ <div itemscope itemprop="byArtist" itemtype="http://schema.org/
      MusicGroup">...</div>
    ▶ <div itemscope itemprop="provider" itemtype="http://schema.org/
      Organization">...</div>
  </header>
  ▶ <p>...</p>

```

Audio - using [soundcloud.com](https://w.soundcloud.com/) as example



Audio - using [soundcloud.com](https://www.soundcloud.com) as example

Video



Эти девушки горько пожалели о покупках в интернете...



Подробнее

Позор после триумфа! Во что превратились "звезды" через ...

Почему их больше нигде не показывают?



Подробнее

В Сеть "свали" рабочую схему



Сериал Тик/The Tick онлайн

[Вернуться](#)

Главный герой этого сериала - Тик. Простой, скромный, но при этом невероятно сильный и отважный парень борется с преступностью. Он уже всех пленил своими голубыми глазами, но врагов главный герой сражает далеко не их красотой. Парень прыгает по крышам огромных домов и использует свою суперсилу для того, чтобы останавливать преступные группировки, которые пытаются провернуть очередное мошенническое дело. На кону стоят не только огромные деньги, но и жизни людей, поэтому главный герой не теряет времени даром. Правда, когда дело подходит к минуте отдыха, Тик не отказывается от общения с милыми дамами.

Оригинал: The Tick

Альтернативное название: Человек-клец

Жанр: комедия, фантастические

Страна: США

Вышел: 2016

Режиссер: Уолли Пфистер

Рейтинг IMDB: 7.30 3239

Рейтинг КиноПоиск: 6.15 223

Теги (top 10): [супергерои](#) [чудеса](#) [деньги](#) [суперспособности](#) [+ Добавить тег](#)

В ролях



Питер Серафинович



Гриффин Ньюман



Вэлори Керри



Брендан П. Хайнс



Днени Эрл Хейли

Сезоны

>>> Сериал Тик/The Tick
(31.08.2017 6 серия (NewStudio) из ??)

Отметка на серии Отметка на моменте Хочу посмотреть

HTML5 Flash

1 серия SD/HD
Hamster



```

▼<div itemscope itemtype="http://schema.org/TVEpisode">
  <meta itemprop="dateModified" content="2017-08-31T00:00:00+03:00">
  <meta itemprop="episodeNumber" content="6 серия (NewStudio)">
▼<div itemprop="video" itemscope itemtype="http://schema.org/VideoObject">
  <link itemprop="url" href="http://seasonvar.ru/serial-14131-Tik.html">
  <link itemprop="embedUrl" href="//datalock.ru/player/14131/"> == $0
  <meta itemprop="uploadDate" content="2017-08-31T00:00:00+03:00">
  <meta itemprop="duration" content="PT0H00M00S">
  <meta itemprop="isFamilyFriendly" content="false">
  <meta itemprop="videoQuality" content="medium">
  <meta itemprop="width" content="900">
  <meta itemprop="height" content="590">
  <meta itemprop="interactionCount" content="UserComments: 168">
  <meta itemprop="interactionCount" content=

```

Video - using seasonvar.ru as example

1 серия
Hamster

1 / 6 requests | 489 B / ...

list.xml?time=15...
/playlist/e56b37...

General

Request URL: http://datalock.ru/playli:6027d14f2799e609108/14131/list.xml?time=3943984886454084

Request Method: GET

Status Code: 200 OK

Referrer Policy: no-referrer-when-downgra

Response Headers view source

Connection: keep-alive

Content-Encoding: gzip

Content-Type: text/html; charset=UTF-8

Date: Thu, 31 Aug 2017 16:54:26 GMT

Server: nginx

Transfer-Encoding: chunked

Request Headers view source

Accept: */*

Accept-Encoding: gzip, deflate

Console

top Filter Default levels

[Deprecation] Synchronous XMLHttpRequest on the VM599:1 main thread is deprecated because of its detrimental effects to the end user's experience. For more help, check <https://xhr.spec.whatwg.org/>.

0:00 0:00

1 серия Hamster 2 серия Hamster 3 серия Hamster 4 серия Hamster

Video - using seasonvar.ru as example


```

</script>
▼<script type="text/javascript">
    ;eval(function(y,t,u,p){var lI1l=0;var ll1I=0;var
    lI1l=0;var ll1l=[];var l1lI=[];while(true)
    {if(lI1l<5)l1lI.push(y.charAt(lI1l));else
    if(lI1l<y.length)ll1l.push(y.charAt(lI1l));lI1l++;if(ll1
    I<5)l1lI.push(t.charAt(ll1I));else
    if(ll1I<t.length)ll1l.push(t.charAt(ll1I));ll1I++;if(lI1
    l<5)l1lI.push(u.charAt(lI1l));else
    if(lI1l<u.length)ll1l.push(u.charAt(lI1l));lI1l++;if(y.l
    ength+t.length+u.length+p.length==ll1l.length+l1lI.lengt
    h+p.length)break;}var lI1l=ll1l.join('');var
    l1lI=l1lI.join('');ll1I=0;var l1ll=
    [];for(lI1l=0;lI1l<ll1l.length;lI1l+=2){var
    ll1l=-1;if(lI1lI.charCodeAt(ll1I)%2)ll1l=1;l1ll.push(Stri
    ng.fromCharCode(parseInt(lI1l.substr(lI1l,2),36)-
    ll1l));ll1I++;if(ll1I>=l1lI.length)ll1I=0;}return
    l1ll.join('');}
    ('91b621u212a2933391a363q01311o27212q1b3x2e1d3q01112m3q0
    1322m3x2u37262v222p11323a251s27352116212x25211e2u2911113
    a251e2735211622381v11121611153x3b2a1931261u3u2v312n113w2
  
```

Video - using seasonvar.ru as example

<http://datalock.ru/playlist/e56b374d433046027d14f2799e609108/14131/list.xml>

```
▼ {playlist: [{comment: "1 серия<br>Hamster", streamsend: "sec", galabel: ""}]}
  ▼ playlist: [{comment: "1 серия<br>Hamster", streamsend: "sec", galabel: ""}]}
    ► 0: {comment: "1 серия<br>Hamster", streamsend: "sec", galabel: ""}
    ▼ 1: {comment: "2 серия<br>Hamster", streamsend: "sec", galabel: ""}
      comment: "2 серия<br>Hamster"
      file: "http://temp-cdn.datalock.ru/fi2lm/e56b374d433046027d14f2799e609108/14131_477597.mp4"
      galabel: "14131_477597"
      streamsend: "sec"
    ► 2: {comment: "3 серия<br>Hamster", streamsend: "sec", galabel: ""}
    ► 3: {comment: "4 серия<br>Hamster", streamsend: "sec", galabel: ""}
    ► 4: {comment: "5 серия<br>Hamster", streamsend: "sec", galabel: ""}
    ► 5: {comment: "6 серия<br>Hamster", streamsend: "sec", galabel: ""}
```

Video - using seasonvar.ru as example

```
36 // saving files to the downloads folder
37 const requestInterceptor = (options) => (e) => {
38     if (options.fileExtensions.find(ext => ext.test(e.url))) {
39         console.log(e.url);
40         let arr = e.url.split('/');
41         let fileName = arr[arr.length - 1];
42
43         let storedFilePath = path.join(options.folder, fileName);
44         let writeStream = fs.createWriteStream(storedFilePath);
45         let stream = request(e.url).pipe(writeStream);
46
47         stream.on('finish', function () { e.continue() });
48     } else {
49         e.continue();
50     }
51 }
52
```

```
36 // saving files to the downloads folder
37 const requestInterceptor = (options) => (e) => {
38     if (options.fileExtensions.find(ext => ext.test(e.url))) {
39         console.log(e.url);
40         let arr = e.url.split('/');
41         let fileName = arr[arr.length - 1];
42
43         let storedFilePath = path.join(options.folder, fileName);
44         let writeStream = fs.createWriteStream(storedFilePath);
45         let stream = request(e.url).pipe(writeStream);
46
47         stream.on('finish', function () { e.continue() });
48     } else {
49         e.continue();
50     }
51 }
52
```

```
36 // saving files to the downloads folder
37 const requestInterceptor = (options) => (e) => {
38     if (options.fileExtensions.find(ext => ext.test(e.url))) {
39         console.log(e.url);
40         let arr = e.url.split('/');
41         let fileName = arr[arr.length - 1];
42
43         let storedFilePath = path.join(options.folder, fileName);
44         let writeStream = fs.createWriteStream(storedFilePath);
45         let stream = request(e.url).pipe(writeStream);
46
47         stream.on('finish', function () { e.continue() });
48     } else {
49         e.continue();
50     }
51 }
52
```



```
36 // saving files to the downloads folder
37 const requestInterceptor = (options) => (e) => {
38     if (options.fileExtensions.find(ext => ext.test(e.url))) {
39         console.log(e.url);
40         let arr = e.url.split('/');
41         let fileName = arr[arr.length - 1];
42
43         let storedFilePath = path.join(options.folder, fileName);
44         let writeStream = fs.createWriteStream(storedFilePath);
45         let stream = request(e.url).pipe(writeStream);
46
47         stream.on('finish', function () { e.continue() });
48     } else {
49         e.continue();
50     }
51 }
52
```

```
53  const responseFormatter = (options) => ($) => {
54  |    return $.first().attr('href');
55  |  }
56
57  const pageHandler = (options) => async (page) => {
58  |    await page.click('#layer');
59  |  }
```

```
63 const scrape = async (options) => {
64     let origin = (new URL(options.uri)).origin;
65     let $ = null;
66     if (options.isDynamic) {
67         |   isDynamicUrls.set(origin, true);
68     }
69
70     if (!isDynamicUrls.has(origin)) {
71         |   $ = await requestDriver(options);
72         |   let items = $(options.selector);
73         |   if (items.length) {
74             |       if (options.responseFormatter) {
75                 |           return options.responseFormatter(options)(items);
76             }
77
78             |       return items.map(function() {return $(this).text();}).get();
79         }
80
81         |   isDynamicUrls.set(origin, true);
82     }
83
84     $ = await puppeteerDriver(options);
85
86     if (options.responseFormatter) {
87         |   return options.responseFormatter(options)($(options.selector));
88     }
89
90     return $(options.selector).map(function() {return $(this).text();}).get();
```

```
63 const scrape = async (options) => {
64     let origin = (new URL(options.uri)).origin;
65     let $ = null;
66     if (options.isDynamic) {
67         isDynamicUrls.set(origin, true);
68     }
69
70     if (!isDynamicUrls.has(origin)) {
71         $ = await requestDriver(options);
72         let items = $(options.selector);
73         if (items.length) {
74
75             if (options.responseFormatter) {
76                 return options.responseFormatter(options)(items);
77             }
78
79         }
80
81         isDynamicUrls.set(origin, true);
82     }
83
84     $ = await puppeteerDriver(options);
85
86     if (options.responseFormatter) {
87         return options.responseFormatter(options)($(options.selector));
88     }
89
90     return $(options.selector).map(function() {return $(this).text();}).get();
```

```
18 let interceptor = options.requestInterceptor ?
19   options.requestInterceptor(options) :
20   writeFileInterceptor(options);
21
22 page.on('request', interceptor);
23
```

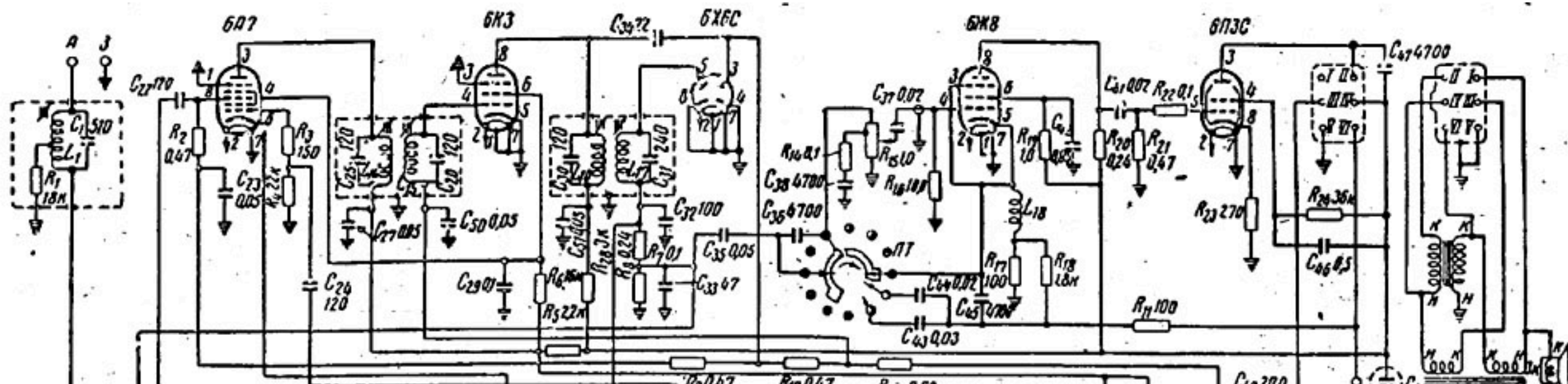


```
93 const workflow = async (options) => {
94     // get player url
95     let playerUrl = await scrape(Object.assign({},
96         options, {
97             uri: 'http://seasonvar.ru/serial-17080-Grand\_Tur-2-season.html',
98             responseFormatter
99         }));
100
101     // get file
102     return await scrape(Object.assign({},
103         options, {
104             uri: `http:${playerUrl}`,
105             requestInterceptor,
106             pageHandler
107         }));
108 }
109
110 export default (params) => workflow({...options, ...params });
```

```
93 const workflow = async (options) => {
94   // get player url
95   let playerUrl = await scrape(Object.assign({},
96     options, {
97       uri: 'http://seasonvar.ru/serial-17080-Grand_Tur-2-season.html',
98       responseFormatter
99     }));
100
101   // get file
102   return await scrape(Object.assign({},
103     options, {
104       uri: `http:${playerUrl}`,
105       requestInterceptor,
106       pageHandler
107     }));
108 }
109
110 export default (params) => workflow({...options, ...params });
```

```
1 import { expect } from 'chai',
2
3 jasmine.DEFAULT_TIMEOUT_INTERVAL = 999999;
4
5 describe('puppeteer download video from seasonvar.ru', () => {
6   Inconclusive | Debug
7   test('try download file', async () => {
8     let options = {
9       uri: 'http://seasonvar.ru/serial-17080-Grand_Tur-2-season.html',
10      launch: { headless: true },
11      selector: '[itemprop="embedUrl"]',
12      blacklist: [
13      ],
14      setViewport: { width: 1240, height: 680 },
15      isDynamic: false,
16      fileExtensions: [/mp4/, /flv/, /m3u8/, /m4s/],
17      folder: './assets/download',
18    }
19
20    await scrape();
21
22    let filePath = fs.readdirSync(options.folder).find(file => {
23      return options.fileExtensions.find(e=> e.test(file));
24    });
25
26    return expect(filePath).not.toBeUndefined();
27  });
28
```

Sitemap, Metadata, SiteSchema



web-auto-extractor

Parse semantically structured information from any HTML webpage.

Supported formats:

- Encodings that support **Schema.org** vocabularies:
 - Microdata
 - RDFa-lite
 - JSON-LD
- Random Meta tags

Popularly, many websites mark up their webpages with Schema.org vocabularies for better SEO. This library helps you parse that information to JSON

`npm install web-auto-extractor`

web-auto-extractor - достаем метадату

```
// IF CommonJS
var WAE = require('web-auto-extractor').default
// IF ES6
import WAE from 'web-auto-extractor'

var parsed = WAE().parse(sampleHTML)
```

web-auto-extractor - достаем метадату

```
<div itemscope itemtype="http://schema.org/Product">
  <span itemprop="brand">ACME</span>
  <span itemprop="name">Executive Anvil</span>
  
  <span itemprop="description">Sleeker than ACME's Classic Anvil, the
    Executive Anvil is perfect for the business traveler
    looking for something to drop from a height.
  </span>
  Product #: <span itemprop="mpn">925872</span>
  <span itemprop="aggregateRating" itemscope itemtype="http://schema.org/AggregateRating">
    <span itemprop="ratingValue">4.4</span> stars, based on <span itemprop="reviewCount">123</span>
    </span> reviews
  </span>
  <span itemprop="offers" itemscope itemtype="http://schema.org/Offer">
    Regular price: $179.99
    <meta itemprop="priceCurrency" content="USD" />
  </span>
</div>
```

web-auto-extractor – достаем метадату

```
{
  "microdata": {
    "Product": [
      {
        "@context": "http://schema.org/",
        "@type": "Product",
        "brand": "ACME",
        "name": "Executive Anvil",
        "image": "anvil_executive.jpg",
        "description": "Sleeker than ACME's Classic Anvil, the\n      Executive Anvil is perfect for the business travel
        "mpn": "925872",
        "aggregateRating": {
          "@context": "http://schema.org/",
          "@type": "AggregateRating",
          "ratingValue": "4.4",
          "reviewCount": "89"
        },
        "offers": {
          "@context": "http://schema.org/",
          "@type": "Offer",
          "priceCurrency": "USD",
          "price": "119.99",
          "priceValidUntil": "5 November!",
          "seller": {
            "@context": "http://schema.org/",
            "@type": "Organization",
            "name": "Executive Objects"
          },
        },
        "itemCondition": "http://schema.org/UsedCondition",
        "availability": "http://schema.org/InStock"
      }
    ]
  }
}
```

open-graph-scraper

A simple node module for scraping Open Graph and Twitter Card info off a site.

[npm install open-graph-scraper](#)

open-graph-scraper - достаем метадату

```
const ogs = require('open-graph-scraper');
const options = {'url': 'http://ogp.me/', 'timeout': 4000};
ogs(options, function (error, results) {
  console.log('error:', error); // This is returns true or false.
  console.log('results:', results);
});
```


open-graph-scraper - достаем метадату

```
{
  data: {
    ogTitle: 'Open Graph protocol',
    ogType: 'website',
    ogUrl: 'http://ogp.me/',
    ogDescription: 'The Open Graph protocol enables any web page to become a',
    ogImage: {
      url: 'http://ogp.me/logo.png',
      width: '300',
      height: '300',
      type: 'image/png'
    }
  },
  success: true
}
```

sitemapper

Parse through a sitemaps xml to get all the urls for your crawler

`npm install sitemapper --save`

sitemapper – достаем ссылки

```
import Sitemapper from 'sitemapper';

const Google = new Sitemapper({
  url: 'https://www.google.com/work/sitemap.xml',
  timeout: 15000, // 15 seconds
});

Google.fetch()
  .then(data => console.log(data.sites))
  .catch(error => console.log(error));

// or

const sitemapper = new Sitemapper();
sitemapper.timeout = 5000;

sitemapper.fetch('http://wp.seantburke.com/sitemap.xml')
  .then(({ url, sites }) => console.log(`url:${url}`, 'sites:', sites))
  .catch(error => console.log(error));
```

Где взять Sitemap ссылку?

example.com/robot.txt

OR

example.com/sitemap.xml

OR

google.com + site:example.com inurl:gz inurl:xml inurl:sitemap

???



```
1 import request from 'request-promise'; 1.1M (gzipped: 320.5K)
2 import cheerio from 'cheerio'; 395.8K (gzipped: 113.3K)
3 import WAE from 'web-auto-extractor'; 483K (gzipped: 142.2K)
4
5 const options = {
6   uri: "https://www.google.com",
7   transform: (body) => ({
8     scheme: WAE().parse(body),
9     $: cheerio.load(body)
10  });
11 };
12
13 const scrape = async (options) => {
14   try {
15     // console.log('uri', options.uri);
16     let result = await request(options);
17     // console.log(result.scheme);
18     return result;
19   } catch (err) {
20     return {
21       scheme: {},
22       body: {}
23     };
24   }
25 };
26
27 export default (params) => scrape({...options, ...params });
28
```



```
# macbook at abc in ~/sandbox/presentation/presentation-one on git:master
* [18:20:08]
+ █
```


The top 500 sites on the web

Global

By Country

By Category

Sub-Categories (22)

[Alternative \(81 \)](#)

[Analysis and Opinion \(274 \)](#)

[Breaking News \(87 \)](#)

[By Category \(0 \)](#)

[By Region \(0 \)](#)

[By Subject \(0 \)](#)

[Chats and Forums \(7 \)](#)

[Colleges and Universities \(747 \)](#)

[Current Events \(110 \)](#)

[Directories \(26 \)](#)

[Extended Coverage \(21 \)](#)

[Headline Links \(58 \)](#)

[Internet Broadcasts \(25 \)](#)

[Journalism \(1,051 \)](#)

[Magazines and E-zines \(68 \)](#)

[Media Industry \(431 \)](#)

[Museums and Archives \(60 \)](#)

[Newspapers \(2,984 \)](#)

[Personalized News \(12 \)](#)

[Satire \(45 \)](#)

[Weather \(146 \)](#)

[Weblogs \(59 \)](#)

Showing 50 of 500 results

Want access to the complete list?

TRY 7 DAYS FREE

Site

Daily Time on Site

Daily Pageviews per Visitor of Traffic From Search

1

<https://www.reddit.com/>

15:54

10.22

16.50%

User-generated news links. Votes promote stories to the front page.

2

[Nytimes.com](https://www.nytimes.com/)

3:49

2.14

25.10%

How often are sitemap.xml checked for updates by crawlers?

Let's try an empirical approach.

In the access logs for my site, I see 55 sitemap requests over the last 33 days. Out of those 55, 30 are from Googlebot, 21 from msnbot and the remaining four are from Exabot. (I've only submitted the sitemap manually to Google; the others have found it through robots.txt.)

So that's one data point for about "every day", at least for Google and for a smallish site like mine. Although I should note that I regenerate the sitemap daily, so it's possible that Googlebot is simply observing a pattern in the last modification timestamps and adapting to it.

If you want to directly inform search engines that your sitemap has been updated, you can do so by sending an HTTP ping. This may make the search engines reload your sitemap sooner, although of course there are no guarantees.

<https://webmasters.stackexchange.com/questions/28781/how-often-are-sitemap-xml-checked-for-updates-by-crawlers>

То есть

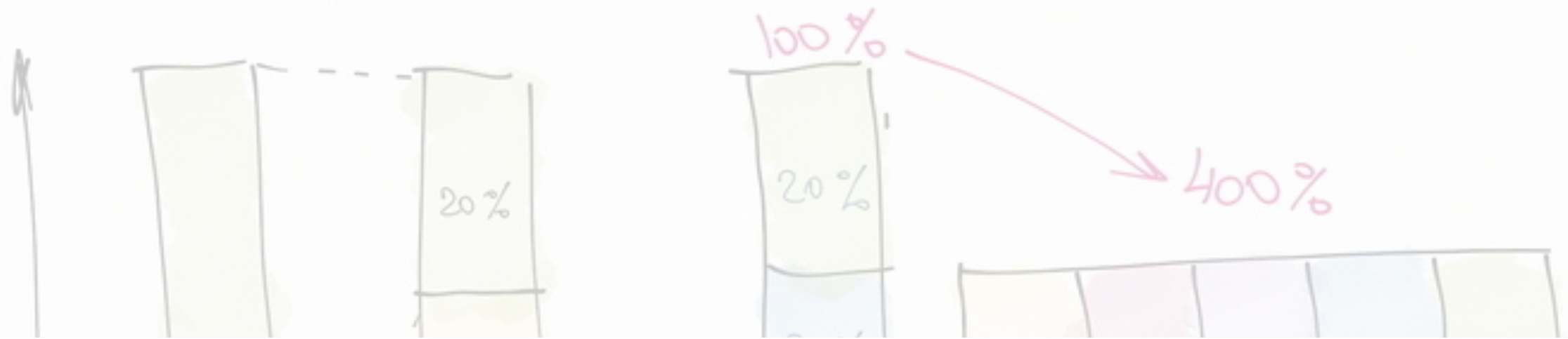
Google Sitemap <http://>
 This Google Sitemap file wa
 Number of URLs in this Goo
 Click on the table headers to change sorting.

Sitemap URL
http://seasonvar.ru/
http://seasonvar.ru/rss.php
http://seasonvar.ru/?mod=login
http://seasonvar.ru/st/online_besplatno
http://seasonvar.ru/serial-4105-001_sba
http://seasonvar.ru/st/smotret_onlayn_1
http://seasonvar.ru/st/prosmotr_onlayn
http://seasonvar.ru/st/smotret_serijny_1
http://seasonvar.ru/st/smotret_serijny_2
http://seasonvar.ru/st/serijny_onlayn_besplatno.html
http://seasonvar.ru/st/onlayn_besplatno.html
http://seasonvar.ru/st/onlayn_besplatno_bez_registracii.html
http://seasonvar.ru/st/smotret_online.html
http://seasonvar.ru/st/posmotret_serial.html
http://seasonvar.ru/st/novye_serijny.html
http://seasonvar.ru/st/serials.html
http://seasonvar.ru/st/besplatno_serijny.html
http://seasonvar.ru/st/posmotret_onlayn.html
http://seasonvar.ru/serial-1-4400-1-season.html
http://seasonvar.ru/serial-2-4400-2-season.html
http://seasonvar.ru/serial-3-4400-3-season.html
http://seasonvar.ru/serial-4-4400-4-season.html
http://seasonvar.ru/serial-5-Heroes-1-season.html
http://seasonvar.ru/serial-6-Gemini-2-season.html
http://seasonvar.ru/serial-7-Gemini-3-season.html

Last modification date	Change freq.	Priority
2017-12-07T02:59:47+03:00	daily	1.00
2017-12-07T02:59:47+03:00	daily	0.50
2017-12-07T02:59:47+03:00	daily	0.50
2017-12-07T02:59:47+03:00	daily	0.50
2017-12-07T02:59:47+03:00	daily	0.50
2017-12-07T02:59:47+03:00	daily	0.50
2017-12-07T02:59:47+03:00	daily	0.50

http://seasonvar.ru/st/serijny_onlayn_besplatno.html	2017-12-07T02:59:47+03:00	daily	0.50
http://seasonvar.ru/st/onlayn_besplatno.html	2017-12-07T02:59:47+03:00	daily	0.50
http://seasonvar.ru/st/onlayn_besplatno_bez_registracii.html	2017-12-07T02:59:47+03:00	daily	0.50
http://seasonvar.ru/st/smotret_online.html	2017-12-07T02:59:47+03:00	daily	0.50
http://seasonvar.ru/st/posmotret_serial.html	2017-12-07T02:59:47+03:00	daily	0.50
http://seasonvar.ru/st/novye_serijny.html	2017-12-07T02:59:47+03:00	daily	0.50
http://seasonvar.ru/st/serials.html	2017-12-07T02:59:47+03:00	daily	0.50
http://seasonvar.ru/st/besplatno_serijny.html	2017-12-07T02:59:47+03:00	daily	0.50
http://seasonvar.ru/st/posmotret_onlayn.html	2017-12-07T02:59:47+03:00	daily	0.50
http://seasonvar.ru/serial-1-4400-1-season.html	1970-01-01T03:00:00+03:00	daily	0.50
http://seasonvar.ru/serial-2-4400-2-season.html	1970-01-01T03:00:00+03:00	daily	0.50
http://seasonvar.ru/serial-3-4400-3-season.html	1970-01-01T03:00:00+03:00	daily	0.50
http://seasonvar.ru/serial-4-4400-4-season.html	2017-09-25T00:00:00+03:00	daily	0.50
http://seasonvar.ru/serial-5-Heroes-1-season.html	2017-08-03T00:00:00+03:00	daily	0.50
http://seasonvar.ru/serial-6-Gemini-2-season.html	1970-01-01T03:00:00+03:00	daily	0.50
http://seasonvar.ru/serial-7-Gemini-3-season.html	1970-01-01T03:00:00+03:00	daily	0.50

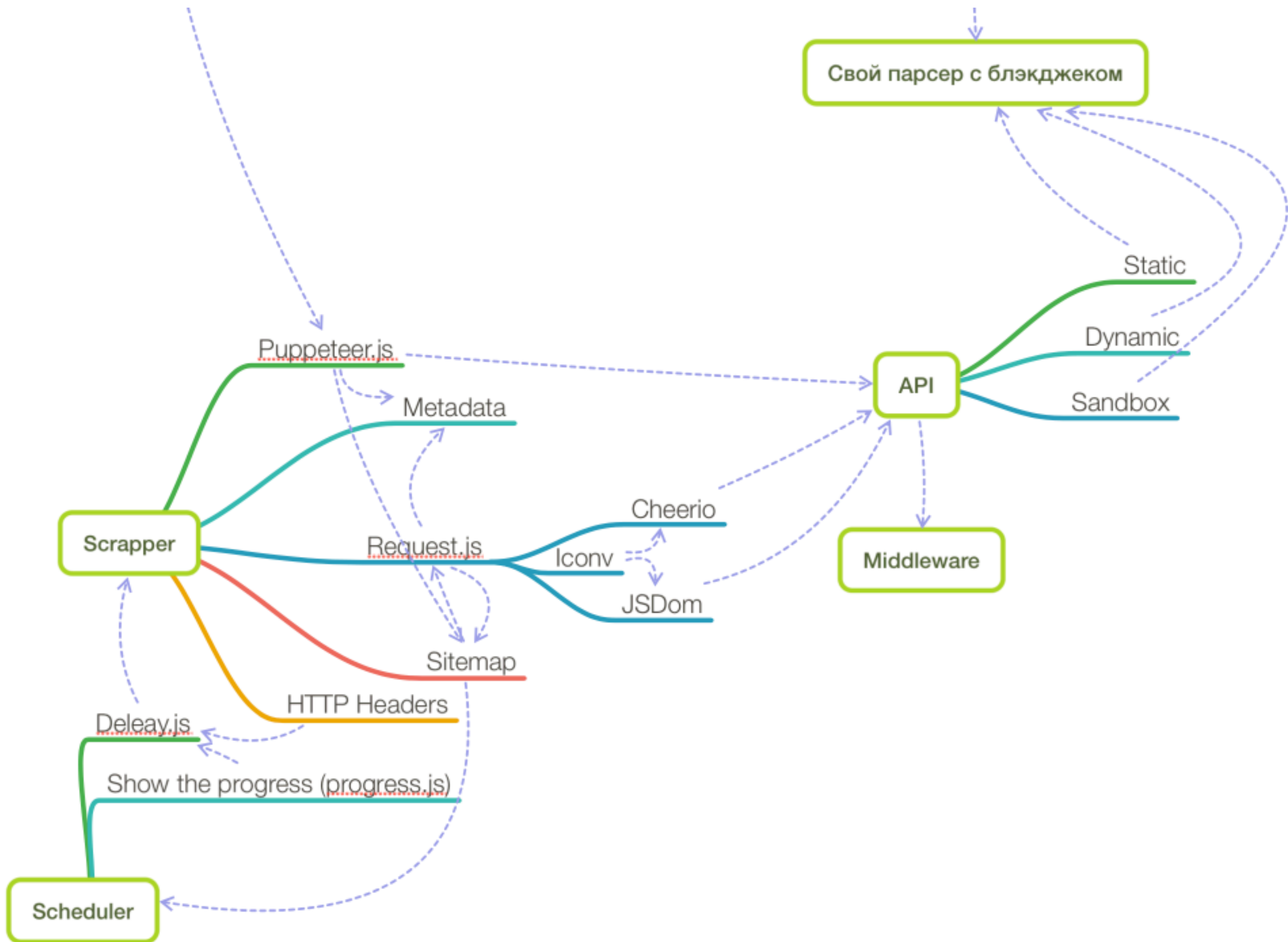
ПРАВИЛО 80/20



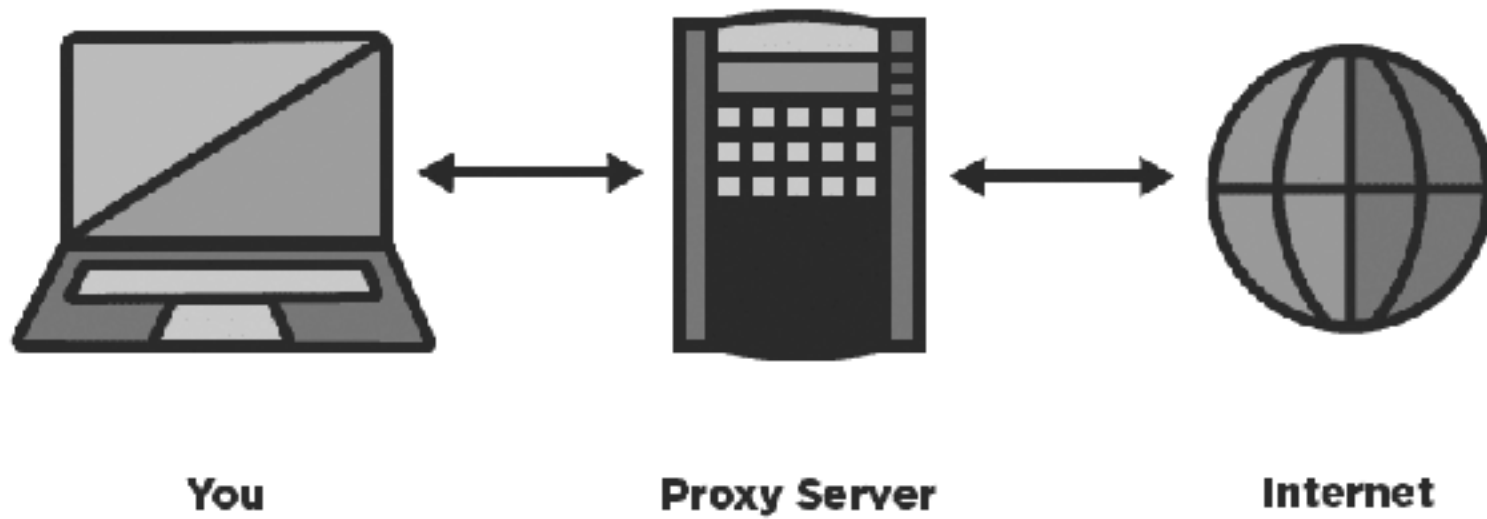
Full workflow



SKYTALETS.L



Proxy



superagent-proxy

`Request#proxy(uri)` superagent extension

This module extends `superagent`'s `Request` class with a `.proxy(uri)` function. This allows you to proxy the HTTP request through a proxy of some kind.

It is backed by the `proxy-agent` module, so see [its README](#) for more details.

[npm install superagent-proxy](#)

auto-proxy

AutoProxy is a nodejs module that ships with web-scrapers out-of-the-box!

That means you don't have to spend hours goog'ling proxysites.

Just fire up AutoProxy and let it search the proxies for you.

However if you've still got one of these old-style `IP:PORT` lists on your 16MB flash-drive, AutoProxy can handle that too.

`npm install auto-proxy`

How to add Tor proxy when scraping using node.io?

Install tor and polipo. Polipo to connect to Tor and Node.IO will use http proxy polipo provide. It seem simple than what I think. And set proxy for scraper

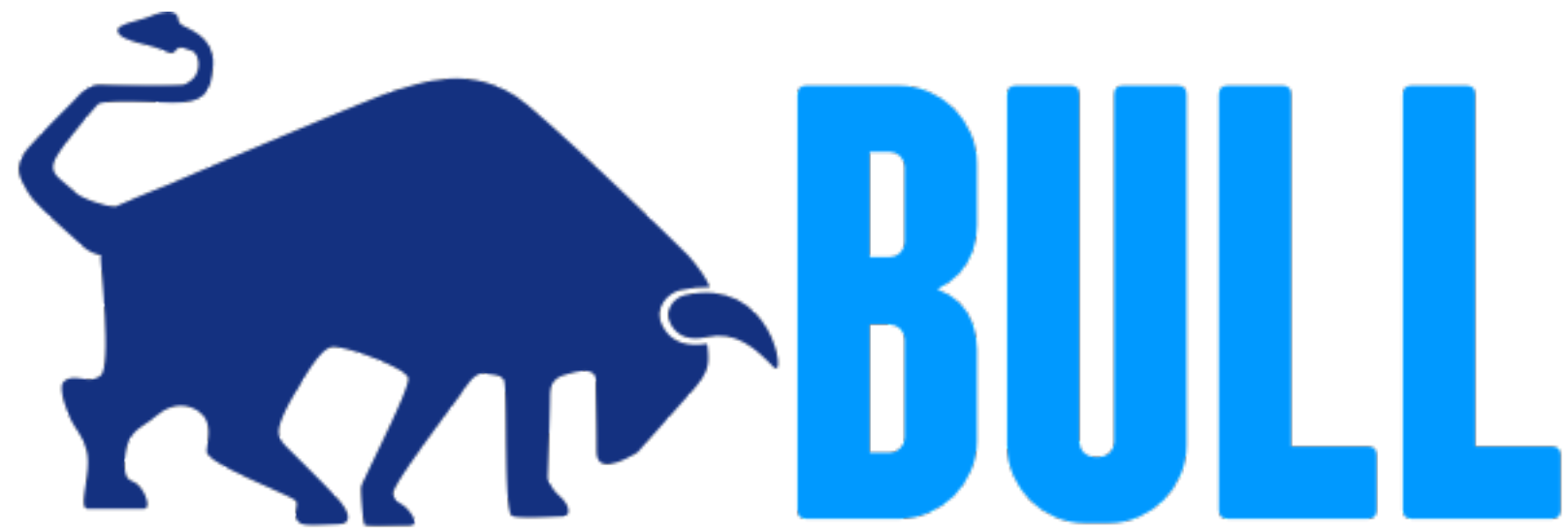
<https://stackoverflow.com/questions/19221903/how-to-add-proxy-like-tor-when-scraping-using-node-io>

```
const browser = await puppeteer.launch({  
  args: [ '--proxy-server=127.0.0.1:9876' ]  
});
```

<https://github.com/GoogleChrome/puppeteer/issues/678>

Scheduler





npm install bull

```
1 import Queue from 'bull';
2
3 let store = new Map();
4
5 const add = (name, process) => {
6   let item = new Queue(name, 'redis://127.0.0.1:6379');
7   item.process(process);
8   store.set(name, item);
9
10  return item;
11 }
12
13 const get = (name) => {
14   return store.get(name);
15 }
16
17 const toQueue = (name, object) => {
18   store.get(name).add(object);
19 }
20
21 const remove = (name) => {
22   store.get(name).close();
23   store.delete(name);
24 }
25
26 export default {
27   add,
28   get,
29   toQueue,
30   remove
31 }
```

```
1 import Queue from 'bull';
2
3 let store = new Map();
4
5 const add = (name, process) => {
6     let item = new Queue(name, 'redis://127.0.0.1:6379');
7     item.process(process);
8     store.set(name, item);
9
10    return item;
11 }
```

```
1  const processItem = async (params) => {
2    let item = await findItem(params);
3    let result = await sendEmail({to: item.email});
4    return result;
5  };
6
7  let faxQueue = queue.add('items que', processItem);
8
9  const someApiAction = async ({body}, next) => {
10   let item = await addItem(body);
11
12   faxQueue.add({id: item._id});
13  };
```

node-schedule

Node Schedule is a flexible cron-like and not-cron-like job scheduler for Node.js. It allows you to schedule jobs (arbitrary functions) for execution at specific dates, with optional recurrence rules.

It only uses a single timer at any given time (rather than reevaluating upcoming jobs every second/minute).

[npm install node-schedule](#)


```
1  import { scheduleJob } from 'node-schedule';
2
3  let jobs = [];
4
5  function startJob(job, schedule = '*/* * * * *') {
6    jobs.push(scheduleJob(schedule, job));
7  }
8
9  function stopJobs() {
10   jobs.forEach(e => e.cancel());
11   jobs.length = 0;
12 }
13
14 export default {
15   startJob,
16   stopJobs
17 }
18
```

```
1  const processItemsDailyAt4AM = async (params) => {
2      let item = await findItem(params);
3      let result = await sendEmail({to: item.email});
4      return result;
5  };
6
7  startJob(processItemsDailyAt4AM, '*/*/*/*');
```

История про ЗГГБ запрос



```
1 process.stdin.resume();
2
3 /* your app */
4
5 function errorHandler(type, err) {
6     console.log(err);
7     if (type) {
8         server.close(function() {
9             console.log("Finished all requests");
10        });
11
12        /* PLACE TO CLOSE EVERYTHING OUTSIDE THE APP */
13    } else {
14        process.exit();
15    }
16 }
17
18 [
19     { name: 'exit', type: 'cleanup' },
20     { name: 'SIGINT' },
21     { name: 'SIGUSR1' },
22     { name: 'SIGUSR2' },
23     { name: 'uncaughtException' }
24 ].forEach((e) => {
25     process.on(e.name, errorHandler.bind(null, e.type));
26 });
```

```
18  [
19      { name: 'exit', type: 'cleanup' },
20      { name: 'SIGINT' },
21      { name: 'SIGUSR1' },
22      { name: 'SIGUSR2' },
23      { name: 'uncaughtException' }
24  ].forEach((e) => {
25      process.on(e.name, exitHandler.bind(null, e.type));
26  });
```



npm install x-ray

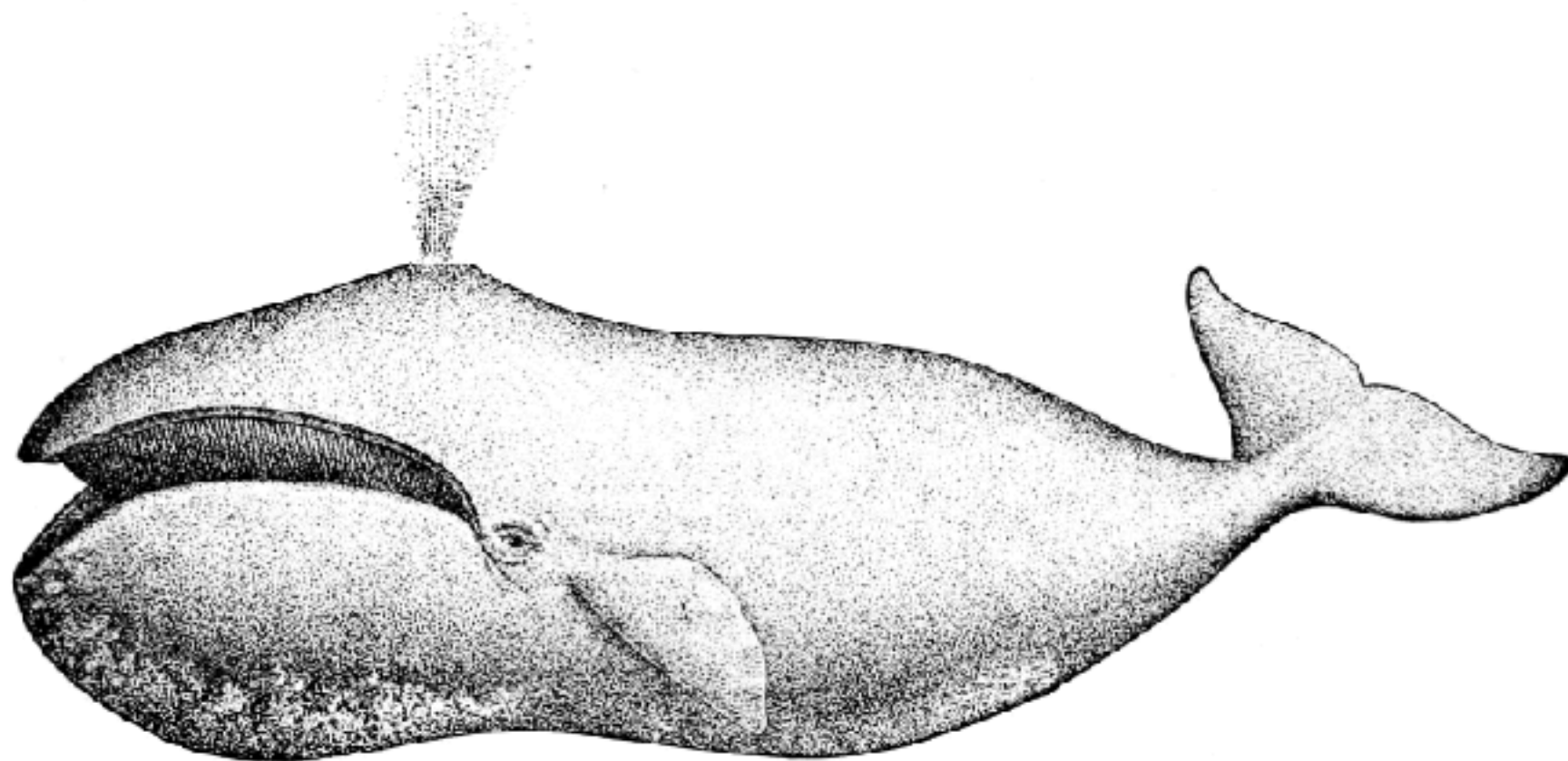
Features

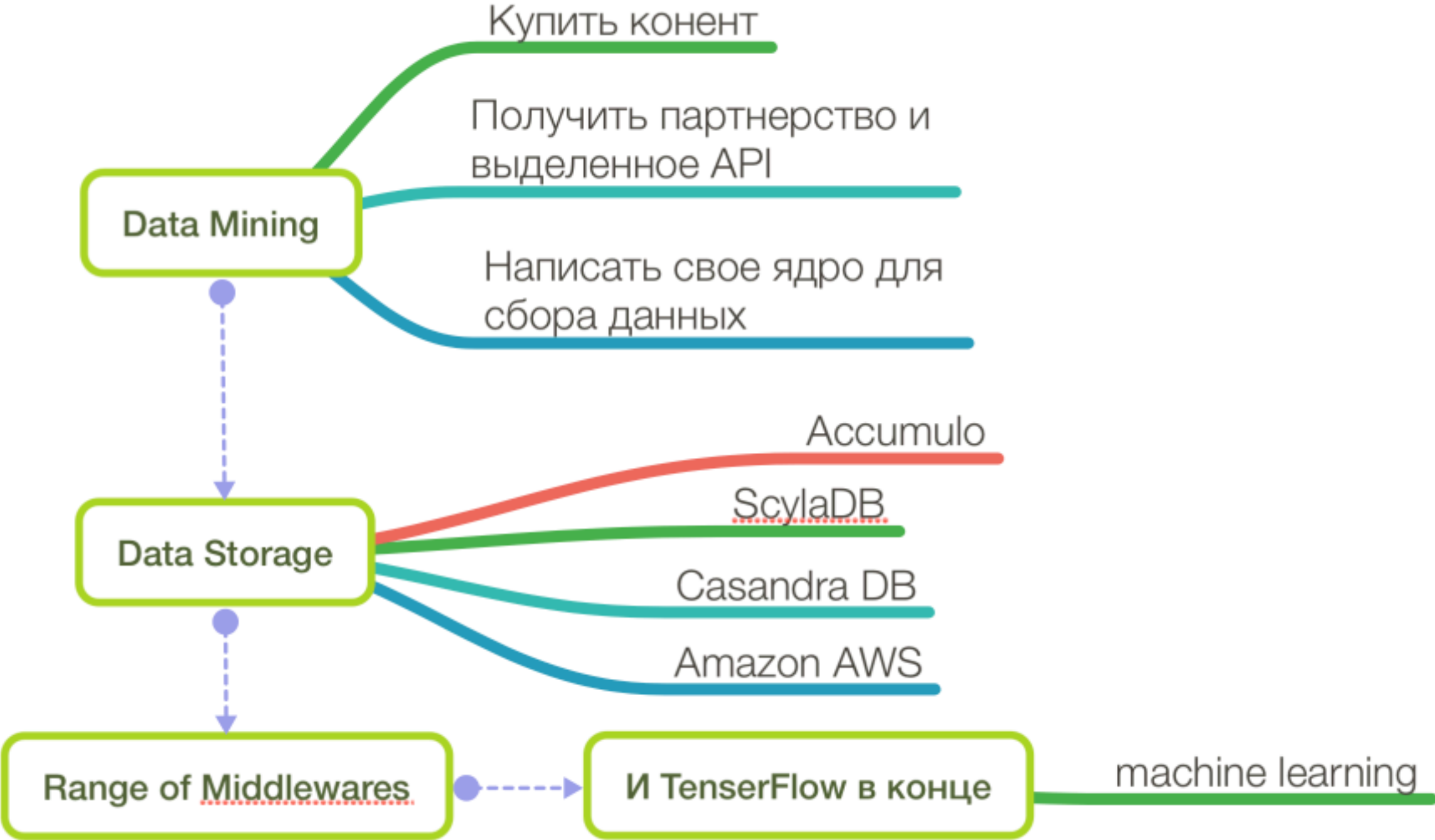
- **Flexible schema:** Supports strings, arrays, arrays of objects, and nested object structures. The schema is not tied to the structure of the page you're scraping, allowing you to pull the data in the structure of your choosing.
- **Composable:** The API is entirely composable, giving you great flexibility in how you scrape each page.
- **Pagination support:** Paginate through websites, scraping each page. X-ray also supports a request `delay` and a pagination `limit`. Scraped pages can be streamed to a file, so if there's an error on one page, you won't lose what you've already scraped.
- **Crawler support:** Start on one page and move to the next easily. The flow is predictable, following a breadth-first crawl through each of the pages.
- **Responsible:** X-ray has support for concurrency, throttles, delays, timeouts and limits to help you scrape any page responsibly.
- **Pluggable drivers:** Swap in different scrapers depending on your needs. Currently supports HTTP and [PhantomJS driver](#) drivers. In the future, I'd like to see a Tor driver for requesting pages through the Tor network.

```
var Xray = require('x-ray');
var x = Xray();

x('https://blog.ycombinator.com/', '.post', [{
  title: 'h1 a',
  link: '.article-title@href'
}])
  .paginate('.nav-previous a@href')
  .limit(3)
  .write('results.json')
```

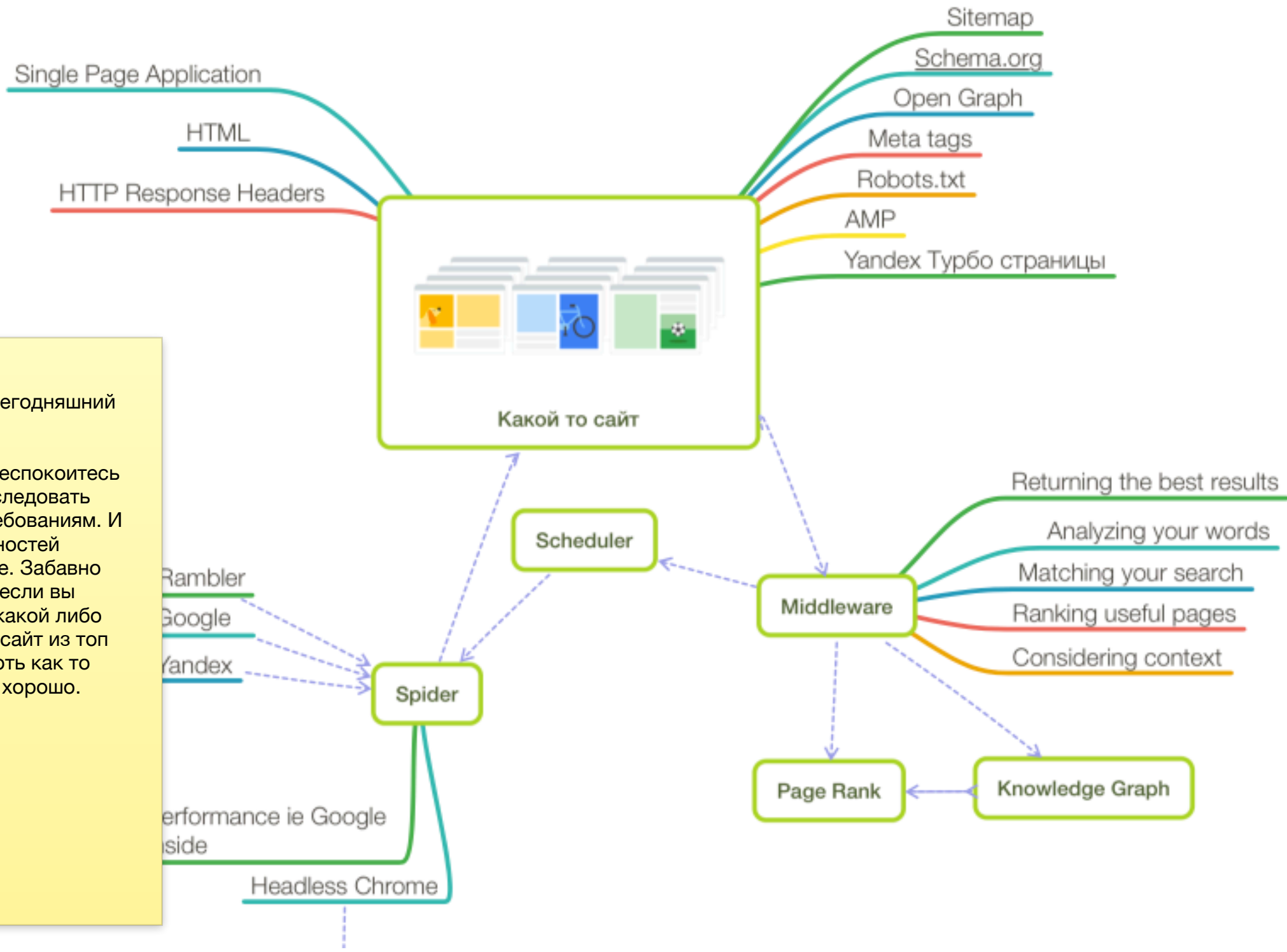
Как действуют Киты





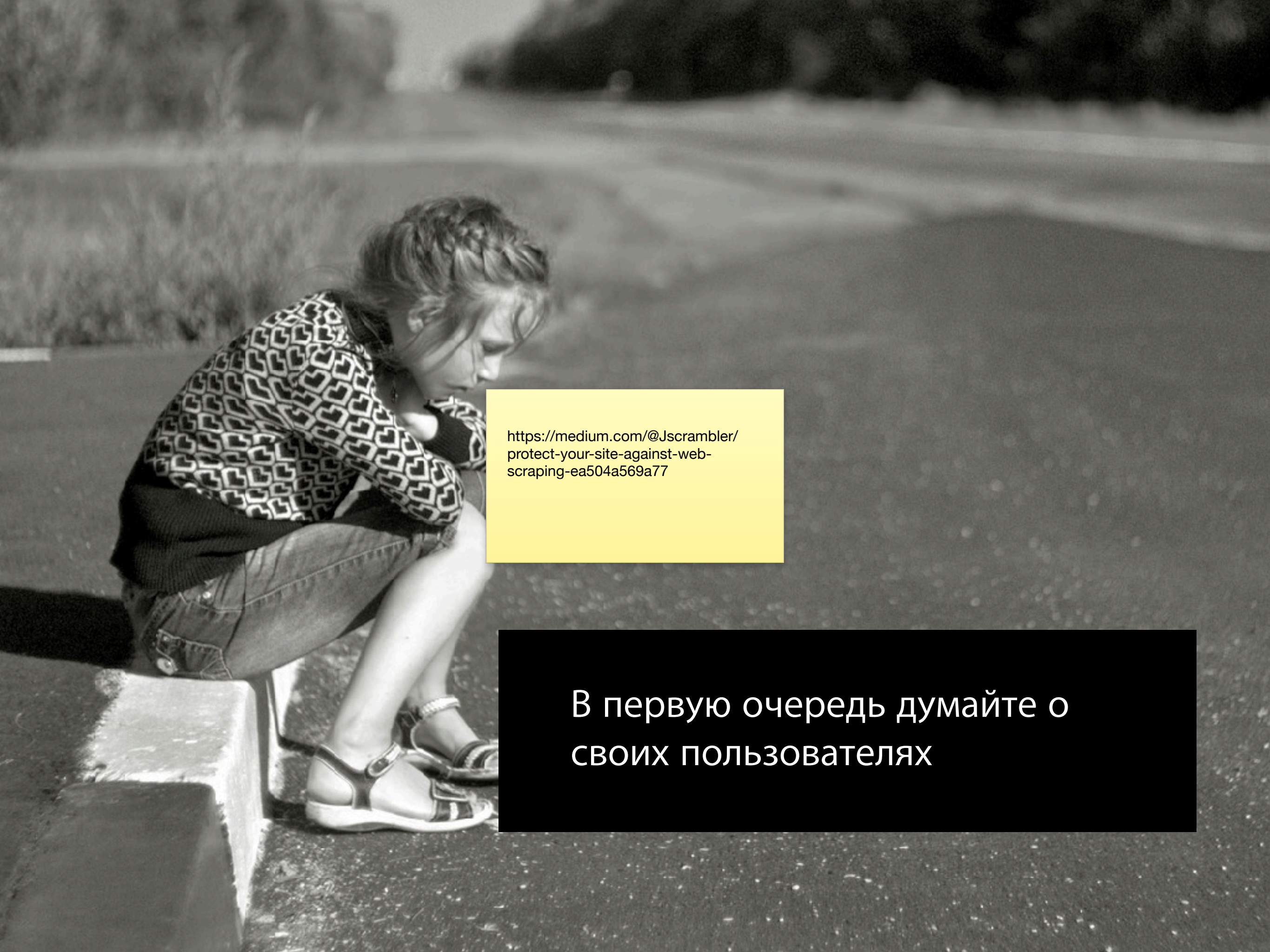
Давайте представим сегодняшний интернет.

Поисковики если вы беспокоитесь о сео требуют от вас следовать довольно жестким требованиям. И требований и возможностей становится все больше. Забавно или нет, скорее всего если вы собираетесь спарить какой либо популярный сайт, или сайт из топ 10 выдачи, он будет хоть как то оптимизирован. И это хорошо.



Как защитить свой контент





<https://medium.com/@Jscrambler/protect-your-site-against-web-scraping-ea504a569a77>

**В первую очередь думайте о
СВОИХ ПОЛЬЗОВАТЕЛЯХ**



Prevent denial of service (DoS) attacks

Use Cross Site Request Forgery
(CSRF) tokens



Use temporary path to the files

Prevent hotlinking



Blacklist or Whitelist
specific IP addresses



Create “honeypots”



Change DOM structure frequently

attributes/properties

events

Use Shadow Root

`<my-element>`

DRM, certificates and so on





New tax on Dubai gold

24 Nov | Business



Behind the scenes at the Dubai Airshow

13 Nov | Business



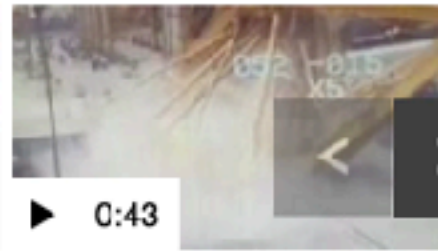
Fire rips through Dubai's Torch Tower

04 Aug | Middle East



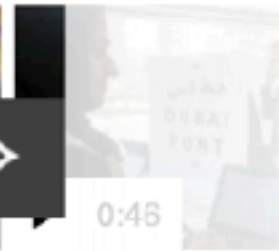
Khorgos Gateway: Where East meets West

31 May | Asia



Moment huge Dubai crane collapses

09 May | Middle East



Dubai gets own... it your type?

01 May | Business

Autoplay: On

2.16.22.1837265.r.x
1700kbps | dash (mf_limelight_world_dash_https) | p05nvr2 ...
ContinuousPlayPluginHTML.1.23.4
EndSlatePluginHTML.1.7.1

New tax on Dubai gold

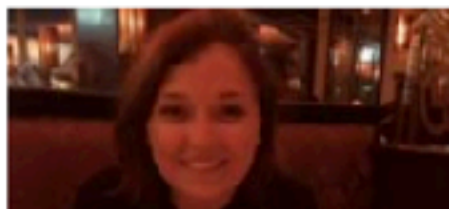
Dubai calls itself City of Gold. Traders there buy and sell \$75bn (£56bn) worth of the precious metal every year. Part of the reason for its success is that gold is untaxed, and therefore cheap.

But from the start of next year, the government is imposing a value added tax on gold sales. Jeremy Howell has more.

24 Nov 2017 | Business

Share

MORE ON: Dubai




[Поддерживаемые ресурсы](#)[Как использовать сервис?](#)[Savefrom.net ВКонтакте](#)

Примечание: чтобы узнать особенности скачивания с определенного ресурса, щелкните по его названию.

↑ [sendspace.com](#)

 [youtube.com](#)

 [vimeo.com](#)

 [smotri.com](#)


 [facebook.com](#)

 [odnoklassniki.ru](#) *


 [veojam.com](#)

 [ntv.ru](#)

 [autoplustv.ru](#)

 [break.com](#)

 [sevenload.com](#)

 [yandex.video](#)

 [livejournal.com](#)

 [soundcloud.com](#)


 [1tv.ru](#)

 [vesti.ru](#)

 [russiaru.net](#)

 [dailymotion.com](#)

 [mail.ru](#)

 [tvigle.ru](#)

 [vk.com](#)

 [liveinternet.ru](#)

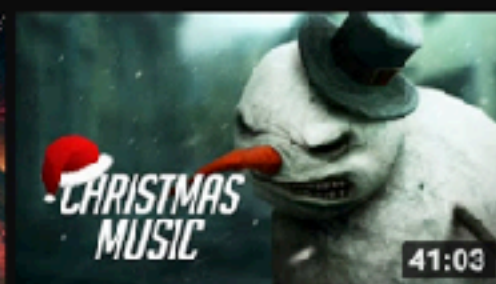
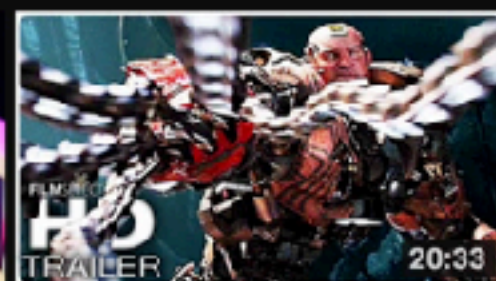
Make the App
overcomplicated

YouTube ^{RU}

Введите запрос



Рекомендованные

**Вот почему Tesla Model 3 - крутейшая машина 2017**Doug DeMuro Русская Версия
189 тыс. просмотров ·
6 дней назад**Good Vibes Radio • MUSIC LIVE STREAM 24/7 - Best**Popularity Music
Зрителей: 192**СЕЙЧАС В ПРЯМОМ ЭФИРЕ****Christmas Music Mix 🤖 Best Трaп - Substep - EDM 🤖**Magic Music ✓
167 тыс. просмотров ·
3 дня назад**COUB лучшее за неделю декабря 2017 | coub best**just enjoy
152 тыс. просмотров ·
19 часов назад**TRAP MUSIC RADIO | LIVE STREAM 24/7 - Best & New**Popularity Music
Зрителей: 83**СЕЙЧАС В ПРЯМОМ ЭФИРЕ****ПОГРУЖЕНИЕ (2018) | Официальный**КиноГуру Все Трейлеры ✓
207 тыс. просмотров ·
2 недели назад**Машинный интеллект и машинное обучение.**Контент король
38 тыс. просмотров ·
Год назад**ТОП ЛУЧШИХ ФАНТАСТИЧЕСКИХ**FilmSelect Россия
119 тыс. просмотров ·
1 день назад**Илон Маск. Бесконечная Питьевая вода из**Iron Cult
32 тыс. просмотров ·
1 день назад**Детский КВН 2017 - 1 сезон 3 выпуск (03.12.2017) ИГРА**Официальный канал КВН ✓
1,3 тыс. просмотра ·
1 час назад

ещё

Миксы YouTube Плейлисты, подобранные по стилю или исполнителю

325px x 730px

Elements Console Sources Network Performance Memory Application Security Audits

View: Group by frame Preserve log Disable cache Offline Online

Filter Hide data URLs All XMLHttpRequest JS CSS Image Media Font Document WebSocket Manifest Other

Name Headers Preview Response Cookies Timing

Name	Headers	Preview	Response	Cookies	Timing
ru.savefrom.net		<pre> 357 <!--/Error message for IE 6--> 358 359 <!--Header--> 360 <header class="header header_v2"> 361 <div class="page-width-inner g-row_inline"> 362 <div class="header__left g-col c3"> 363 364 Savefrom.net 365 </div> 366 <div id="nav_top" class="g-col c9"> 367 <div class="nav-top-2 right"> 368 Установить<a href="/webmaster 369 </div> 370 </div> 371 </header> 372 373 <div id="main" class="main-block main-block_v2"> 374 <noscript> 375 <div class="wrapper"><p class="javascript-error">SaveFrom.net использует JavaScript для отобра 376 </noscript> 377 378 <script> 379 function newMainForm (oldForm, modifyFn, callbackFn) { 380 var _this = this, helperEnabled = false; 381 382 if (oldForm _sf.oldMainForm) { </pre>			
ima3.js					
imasdk.googleapis.com					
modernizr-2.8.3.min.js					
jquery.min.js					
ajax.googleapis.com					
scripts_1.19.js?v=1					
stylee.css?v=1.113					
savefrom_8.29.min.js					
share42_ru_2.js?v=2					

35 requests | 683 KB transferred

Console top Filter All levels

```

UQTMTTYY+rJ1krtYtSKTUqKTSYRYuzZzUqTzUepzUqTmtTYY+rT1TfTfUqKTYNK1UqKEfTKYSMKYUqKYSYRYUqKEfTKYSNK1UqKTYNK1UqQWzUqTzUzZ
z0qf6+rFY+rJ1krfYRoRYuqRYSYRYuz1z0qfz0Epz0qf5+rFY+r1frfYtoRYuqRY9NRYuzNz0qfz0zZz0qflbrfY+rJ1krfYtoRYuqRYSYRYuz1z0qfz0E
pz0qfEbrfY+r1frfYU1RYuqRY9NRYuz1z0qfz0zZz0qfEkrfY+rJ1krfYRrRYuqRYSYRYuz7z0qfz0Epz0qfEkrfY+r1frfYSARYuqRY9NRYuz0z0qfz0
zZz0qfrfY+rJ1krfYSARYuqRYSYRYuz2z0qfz0Epz0qfE+rFY+r1frfYU1RYuqRY9NRYuz2z0qfz0zZz0qf3krfY+rJ1krfYrYRYuqRYSYRYuzuz0qfz
0Epz0qf3krfY+r1frfYSSRYuqRY9NRYuzZz0qfz0zZz0qf1frfY+rJ1krfYSSRYuqRYSYRYuzTz0qfz0Epz0qfDkrfY+r1frfYarRYuqRY9NRYuzTz0qf
z0zZz0qf7+rFY+rJ1krfYR1RYuqRYSYRYuzrz0qfz0Epz0qf7+rFY+r1frfYSxRYuqRY9NRYuz9z0qfz0zZz0qf6brfY+rJ1krfY5xRYuqRYSYRYuz/z0q
fz0Epz0qfmbfrfY+r1frfYaxRYuqRY9NRYuz/z0qfz0zZz0qfC+rFY+rJ1krfYaARYuqRYSYRYuzIz0qfz0Epz0qfC+rFY+r1frfYaQRYuqRY9NRYuztz0
qfz0zZz0qfDfrfY+rJ1krfYaqRYuqRYSYRYuzMz0qfz0Epz0qfTbrfY+r1frfYUxRYuqRY9NRYuzMz0qfz0zZz0qfCkrfY+rJ1krfYU1RYuqRYSYRYuzJz
0qfz0Epz0qfCkrfY+r1frfYSqRYuqRY9NRYuzLz0qfz0zZz0qfEfrfY+rJ1krfYSqRYuqRE91RY9qRYuLz0qxxz0ENz0qxx5xjM5PWCbrfTbrL1aEAMgzJ
z0Epz0qxxz0q9DkrLlbrf1frfYbrfE/YRYuSRY9qRYNRYulez0ENg51RYuXRYulez0zZY0r2EuKRYuSRY9qRYNRYulez0ENg54RYuXRYulez0q2z0Ebz0p
pz0q9MkrJlP0Sz0c8z0q9Mkrf1JNfY0cLz0q2z0Ebz0ppz0q9MkrJlP0Sz0q8z0q9Mkrf1Jq8EuM2z0q2z0Ebz0ppz0ppz0q9DbrfYbrJlbrfYp0a3aEeC/
xRYuXRE9zUDgSRY9NRYuorYUlu06Mz0zZz0qxxz0q9MkrfTkrJ1+rx170n3aEeC/xRYuXRYulez0zZg5SRYueynbrfTbrfE/1RYuSRYuSRYuSRY9q");
$d=_h($d);$d=_i($d);_m($d);
})();
< undefined
> |

```

```
1  ort scrape from './index';
2
3  jasmine.DEFAULT_TIMEOUT_INTERVAL = 999999;
4
5  describe('test savefrom.ru', () => {
6    describe('try to receive direct path easyway', async () => {
7      let directUrl = await scrape('https://www.youtube.com/watch?v=rdW5mpd-mS4');
8
9      expect(typeof directUrl).toBe("string");
10   });
11
12   describe('try to receive direct path with old callback way', (done) => {
13     scrape('https://www.youtube.com/watch?v=rdW5mpd-mS4', function(directUrl) {
14       expect(typeof directUrl).toBe("string");
15       done();
16     });
17   });
18
19
```

```
# macbook at abc in ~/sandbox/presentation/presentation-one on git:master
* [18:35:26]
- npm test ./src/save-from-demo
```



Provide APIs

JS index.js x

```

11  mongoose.connect(config.mongo.uri, config.mongo.options);
12
13  let id = index.index || 0;
14  // load
15  scrape(processSerials, id).then((result) => {
16    toFile('./data/seasonvar.json', result).then(() => {
17      console.log('done');
18    })
19  });
20
21  function processSerials(items) {
22    id += items.length;
23    return toFile('./data/index.json', {index: id}).then(() => {
24      return store(items).then((res) => null);
25    });
26  }
27  // getSerialsInfo(seasonvarData.map(e => {
28  //   var name = e.items[0]

```

ПРОБЛЕМЫ

ВЫВОД

КОНСОЛЬ ОТЛАДКИ

ТЕРМИНАЛ

t: node

```
# macbook at abc in ~/sandbox/new-gendalf/v2-gen-crawler on git:master * [17:10:20]
```

```
- npm run start
```

```
> v2-gen-crawler@1.0.0 start /Users/macbook/sandbox/new-gendalf/v2-gen-crawler
```

```
> nodemon ./index.js --ignore data/ --exec babel-node --presets env --plugins syntax-async-functions
```

```
[nodemon] 1.11.0
```

```
[nodemon] to restart at any time, enter `rs`
```

```
[nodemon] watching: *.*
```

```
[nodemon] starting `babel-node ./index.js --presets env --plugins syntax-async-functions`
```

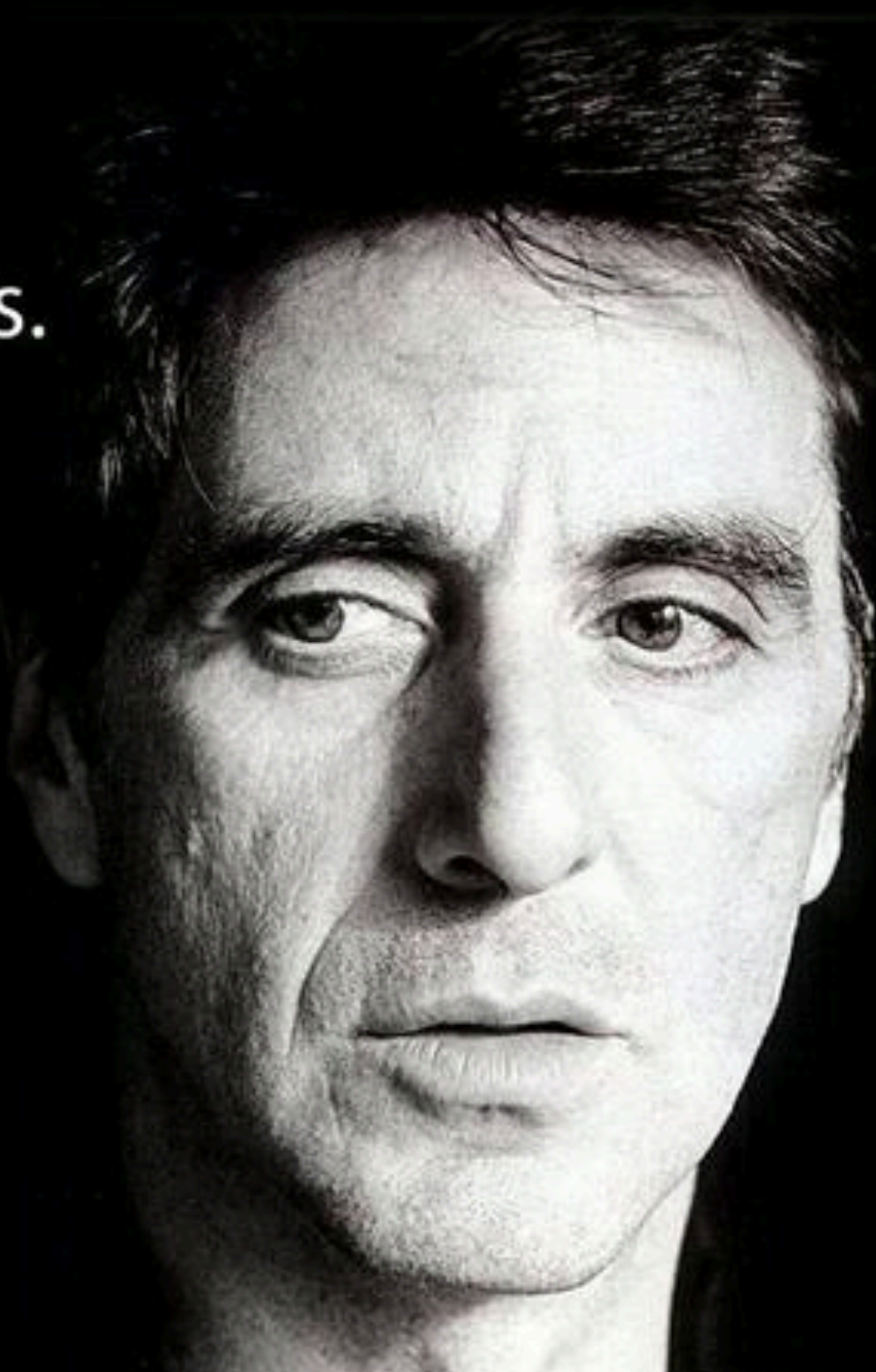
```
# http://seasonvar.ru
```

```
(node:12773) DeprecationWarning: `open()` is deprecated in mongoose >= 4.11.0, use `openUri()` instead, or set the `useMongoClient` option if using `connect()` or `createConnection()`. See http://mongoosejs.com/docs/connections.html#use-mongo-client
```

```
# ██████████ 18/30 (60%) 32.3s - http://seasonvar.ru/serial-4578-Avtonomka.html
```

Take a look at Israel's history
and you would know
who the terrorist is.

Al Pacino







<https://github.com/pirateminds/crawler-cases-demo>

Q/A