

# **Special Topics in Web Scraping with Python**

# Hi, I'm Isaac Vidas

Software Engineer @ HERE Technologies

 [Linkedin: in/isaac-vidas](https://www.linkedin.com/in/isaac-vidas)

 [Twitter: @kazuarous](https://twitter.com/@kazuarous)

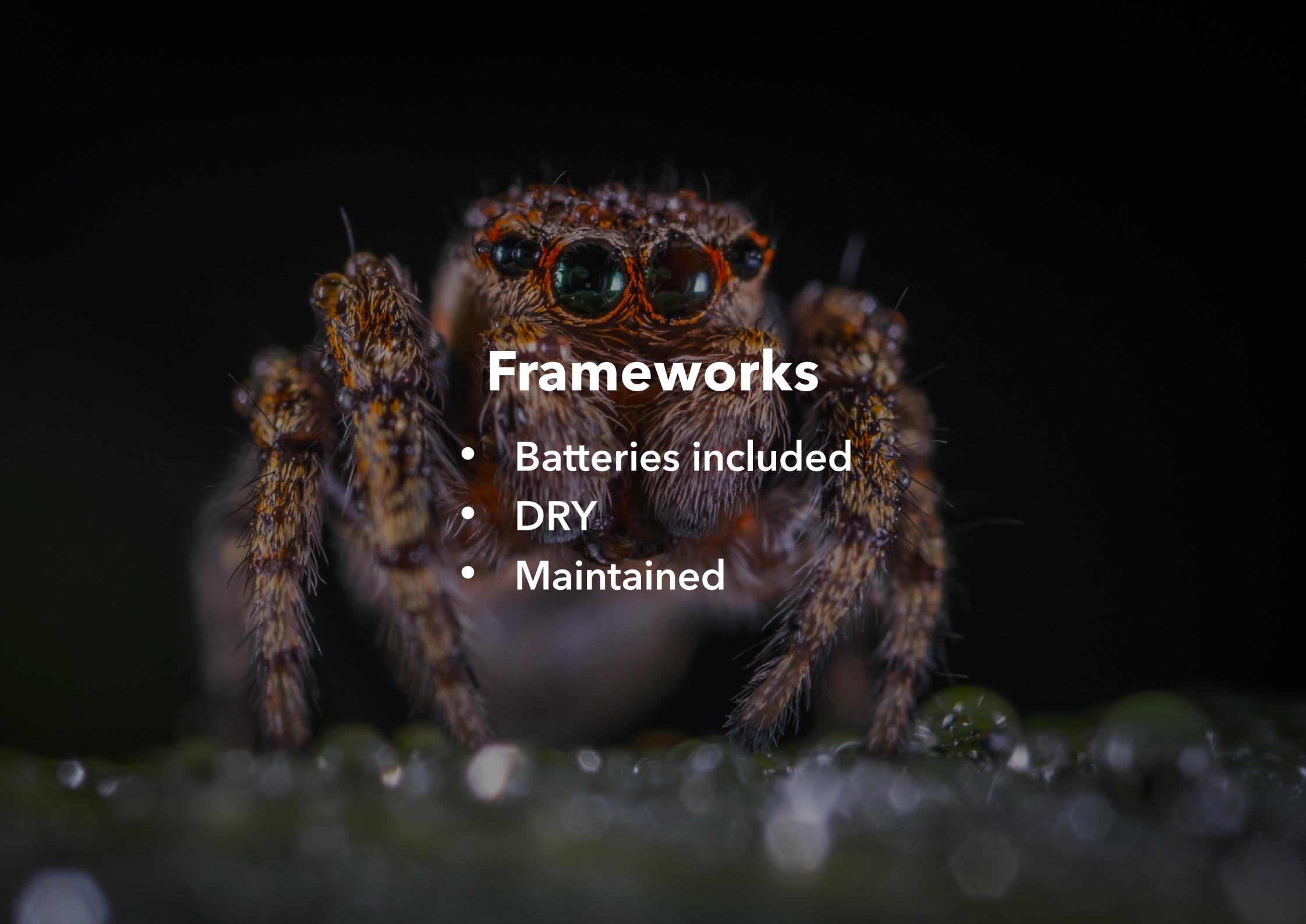
 [Blog: http://vid.as](http://vid.as)

 [GitHub: kazuar](https://github.com/kazuar)

## Agenda

- Frameworks - why do we need them?
- Static vs. dynamic websites
- Don't be evil
- In production





# Frameworks

- Batteries included
- DRY
- Maintained



# Scrapy

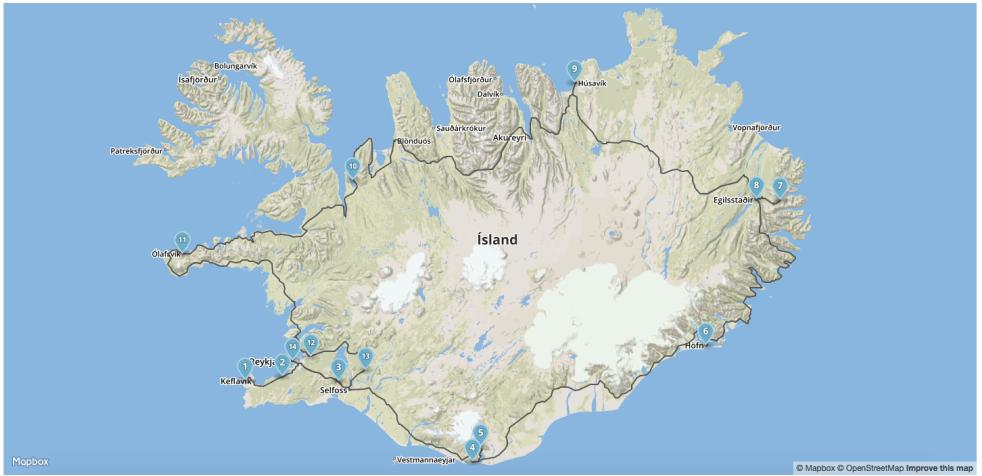
```
1 $ pip install scrapy
2 $ cat > myspider.py <<EOF
3 import scrapy
4
5 class BlogSpider(scrapy.Spider):
6     name = 'blogspider'
7     start_urls = ['https://blog.scrapinghub.com']
8
9     def parse(self, response):
10         for title in response.css('.post-header>h2'):
11             yield {'title': title.css('a ::text').get()}
12
13         for next_page in response.css('a.next-posts-link'):
14             yield response.follow(next_page, self.parse)
15 EOF
16 $ scrapy runspider myspider.py
```

# Static vs. Dynamic

Isaac Vidas  
This is not for you

Blog About

## Visualize your trip with Flask and Mapbox



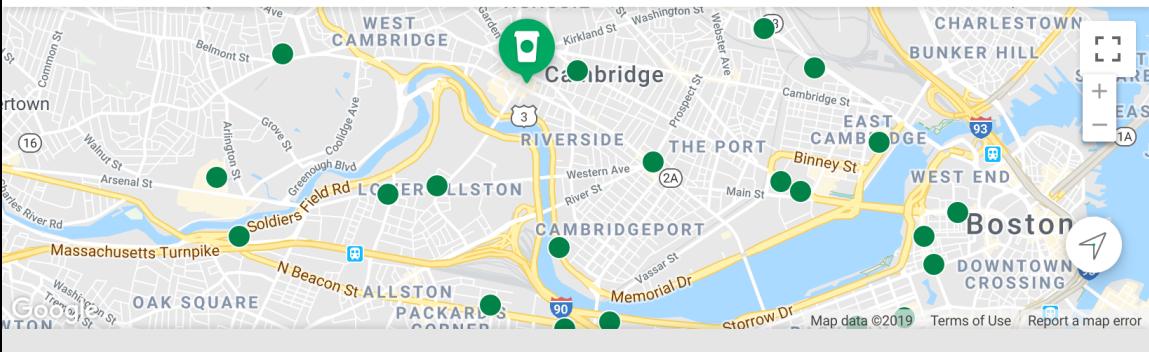
Following my recent trip to Iceland (which I highly recommend), I wanted to visualize the trip using a map, complete with the route and stop locations along the way.

[Source code](#)

Sign in

COFFEE TEA MENU COFFEEHOUSE SOCIAL IMPACT BLOG GIFT CARDS

### Harvard Square, Cambridge, MA, USA



**Harvard Yard**  
Open until 11:00 PM  
0.5 miles away

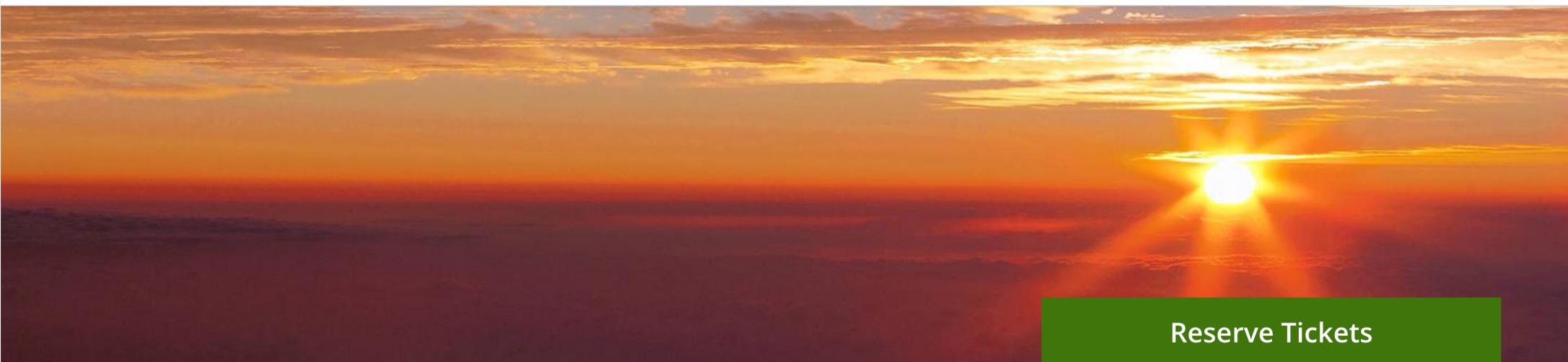
**Broadway Marketplace**  
Open until 9:00 PM  
0.6 miles away

**Shepard Post**  
Open until 10:00 PM  
0.1 miles away

**Automate your  
browser with  
Selenium**



# Too Many Guests

[Help](#)[Sign Up](#)[Log In](#)

## Reserve Tickets

[Home](#) / [Haleakala National Park](#) / [Haleakala National Park Summit Sunrise Reservations](#)

## Haleakala National Park Summit Sunrise Reservations

Part of [Haleakala National Park](#)

[Overview](#)[Need to Know](#)[Fees & Cancellations](#)[Getting Here](#)[Contact](#)

Haleakala Sunrise - Summit

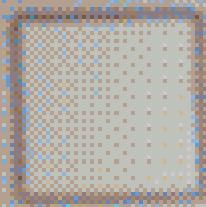
1 Vehicle Pass

08/14/2019

3:00 AM

Too Many Guests

[View Tours](#)



I'm not a robot



Click here to generate direct links

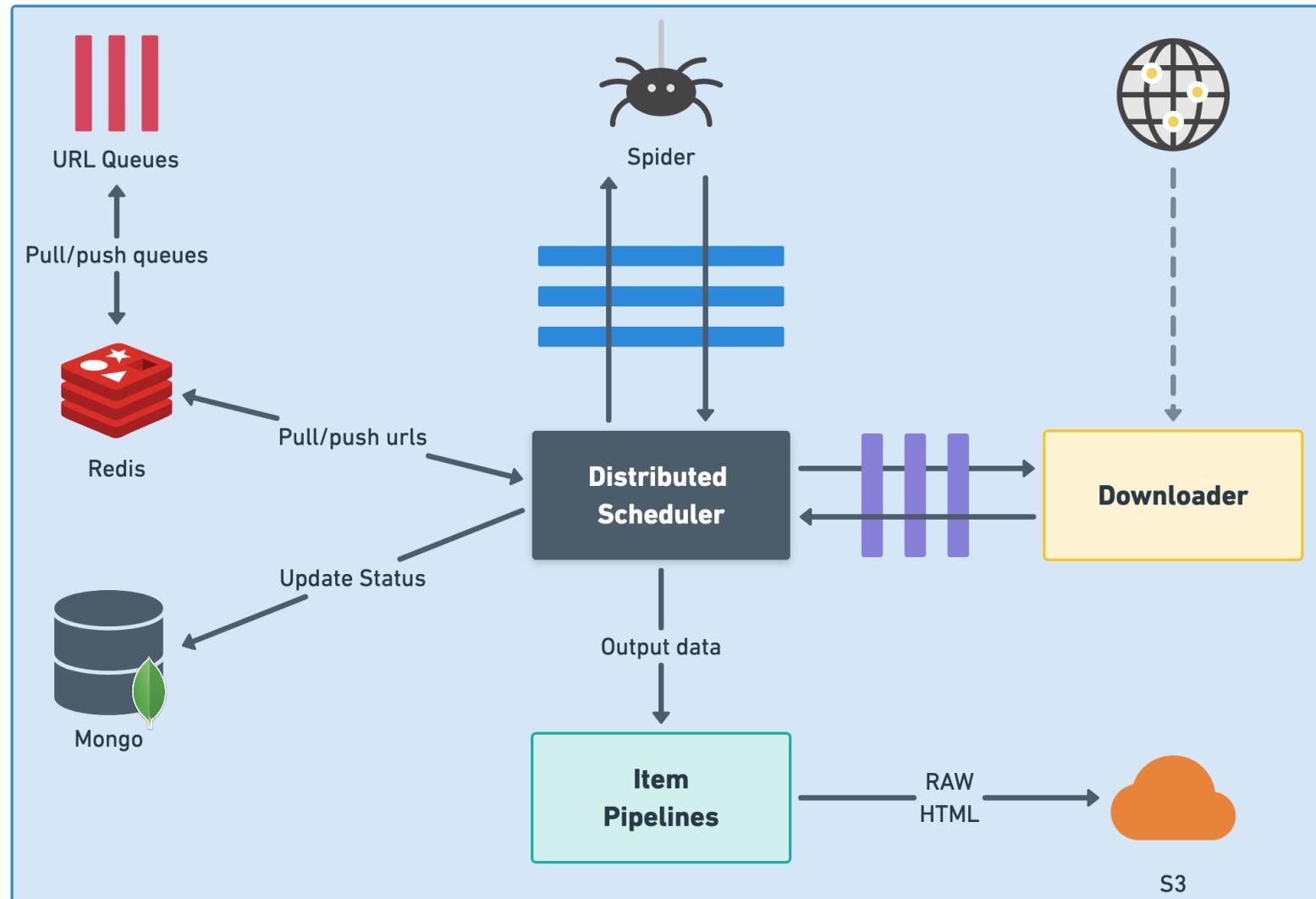
# Don't Be Evil

- Politeness
- Robots.txt
- Hits per second
- Duplicate page filter
- User agent
- Alternatives to crawling



## Web Scraping in Production

- Scaling
- Fault tolerance
- Monitoring
- Testing strategy
- Refresh policy





**Thanks!**

## Resources

- Scrapy - <https://scrapy.org/>
- Selenium - <https://www.seleniumhq.org/>
- Login with Python - <http://kazuar.github.io/scraping-tutorial/>
- Crawl politely -  
<https://blog.scrapinghub.com/2016/08/25/how-to-crawl-the-web-politely-with-scrapy>