

CRAWLING THE CONCRETE JUNGLE

WITH SCRAPY

OVERVIEW



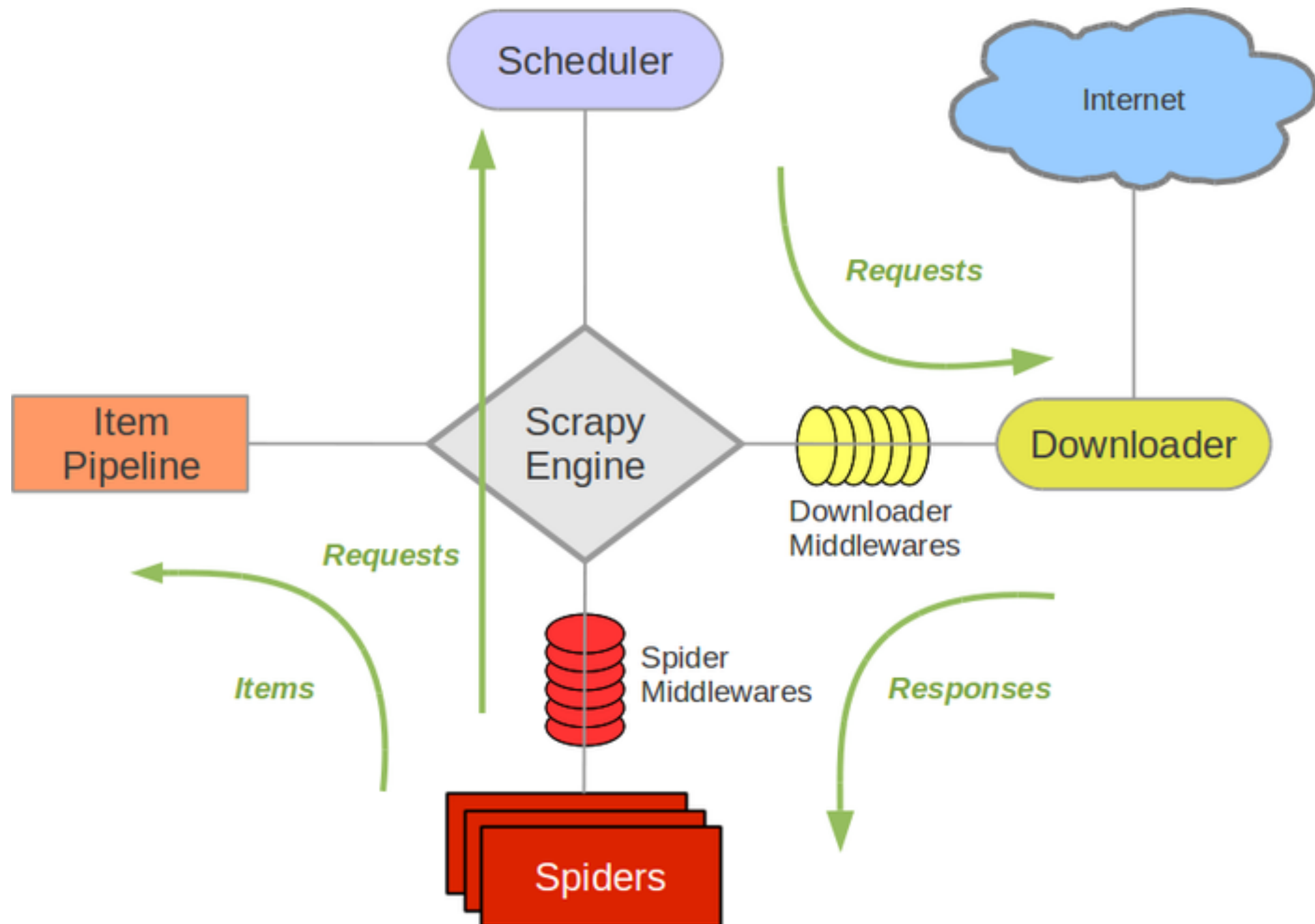
- ▶ Aggregates NYC listing data
- ▶ August 2013 - Purchased by Zillow for \$50mm
- ▶ Standardized listing data
- ▶ Over 19k active rental listings



PROJECT OVERVIEW

- ▶ Used Scrapy to crawl and parse 19k active rental listings
- ▶ Scrapy shell & Chrome
- ▶ The dataLayer
- ▶ Pipe data to PostgreSQL with SQLAlchemy
- ▶ Limitations / Future Improvements
 - ▶ Only rental data / expand to sale data
 - ▶ Current snapshot / scrape historical listing data
 - ▶ Machine learning to predict NYC real estate market

SCRAPY ARCHITECTURE



SPIDER CODE SNIPPET

```
2 import scrapy
3 import re
4 import ast
5 from nycrex.items import listingItem
6
7 import logging
8 from scrapy.utils.log import configure_logging
9
10 configure_logging(install_root_handler=False)
11 logging.basicConfig(
12     filename='nycrex.log',
13     format='%(asctime)s: %(name)s: %(levelname)s: %(message)s',
14     level=logging.INFO
15 )
16 logging.getLogger('sqlalchemy.engine').setLevel(logging.DEBUG)
17
18 class StreeteasySpider(scrapy.Spider):
19     name = "streeteasy"
20     allowed_domains = ["streeteasy.com"]
21     start_urls = [
22         'http://streeteasy.com/for-rent/manhattan',
23         'http://streeteasy.com/for-rent/brooklyn',
24         'http://streeteasy.com/for-rent/queens',
25         'http://streeteasy.com/for-rent/bronx',
26         'http://streeteasy.com/for-rent/staten-island'
27     ]
28
29     def parse(self, response):
30         for sel in response.xpath('//div[@class="details-title"]'):
31
32             url = 'http://streeteasy.com%s' % sel.xpath('a/@href').extract()[0]
33             yield scrapy.Request(url, callback=self.parseListing)
34
35         for href in response.xpath('//div[@class="pagination center_aligned bot'):
36             url = response.urljoin(href.extract())
37             yield scrapy.Request(url, callback=self.parse)
38
39     def parseListing(self, response):
40
41         # Parse <script> dataLayer object
42         data = response.xpath('//script/text()[contains(., "dataLayer")]')[0]
43         data = data.extract().strip()
```

SHINY DEMO