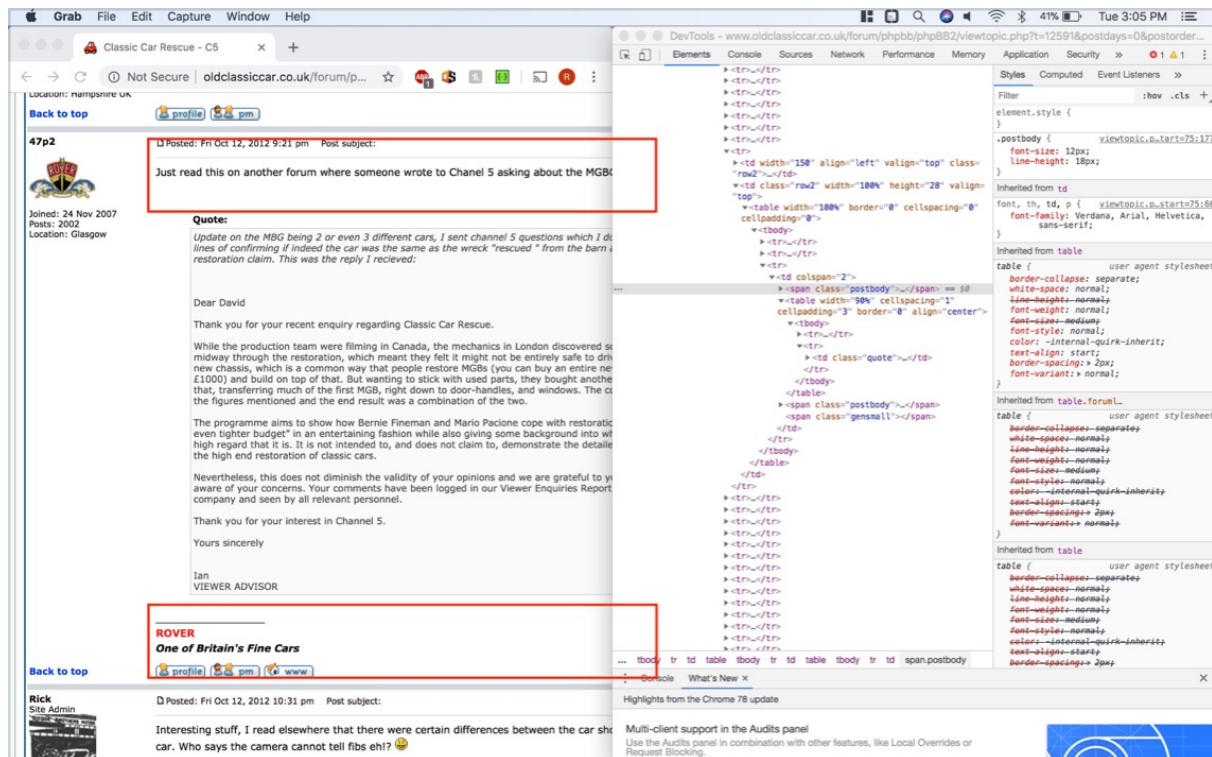# Design Decisions

## Performance

### Async vs Sync

In a synchronous web scraper, each website is visited in serial. This means while the page data is being fetched, the program is waiting, essentially doing nothing. This is lower the program's throughput and therefore performance. This may be fine for this use case, since the thread is only nine pages long. However, when dealing with enourmous websites, excecution time increases dramatically.

In order to improve this, I used a async http request library `grequests` in order to be able to fetch data from multiple web pages in parallel. Which means that there is less detriments to the program's throughput.

---

## Modularity and Scalability

### Returning as Lists in Raw Data

Prior to any custom transformations, the data is stored in list. This might not make sense at first since ideally each needed datapoint maps to one element per post. However, this is not the case for the example website since that when quotes are present there are multiple post bodies.



Since this pattern is present on this website, it could be present on others. With modularity in mind, it's better to keep the data raw at first so that this scraper can be reused on

multiple websites. The engineer can then modify the data as needed using custom transformations in the Scraper class' `data_config` paramter.

## Config Schema

```python
data_config = {
    'post_id': { # field the data will be saved in result record object
        'path': [
            { 'selector': 'span', 'options': { 'class': 'name' } }
        ],
        'selector': 'a',
        'options': {},
        'attribute': 'name',
        'transform': lambda res: res[0] # optional custom function to
manipulate response
    },
}
```

The scraper is designed to follow the path of the target's ancestors first before getting all of the targeted elements. Each object feeds into a typical BeautifulSoup `find` function. Written this way, not only does the scraper remain DRY, but it allows engineers to use the scraper without much knowledge of BeautifulSoup.

## Transform

Raw data is collected initially and can be transformed using custom functions specified by the engineer.

```
[u'Posted: Mon Sep 24, 2012 4:53 pm\xa0\xa0 \xa0Post subject: Classic Car Rescue —
C5']
```

```
=> u'Mon Sep 24, 2012 4:53'
```

Rather than hardcode these transformations into the webscraper, which would reduce modularity, the engineer has the ability to pass in custom transformations for their own needs. This means that the scraper can be used on multiple websites.

## Seperation of Nav, Body, and Data

The locations of the navigation, the body (where the data is located) and what data the engineer wants to scrape is different for every website. To maintain modularity, I decided to seperate their logic.