



gMonitor: The itsy bitsy spider

By: Gerhard van Andel

Web crawler

■ Rest Endpoint,

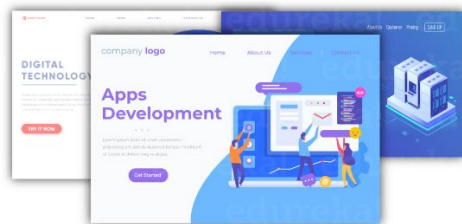
- POST -- "/api/url"

- body:

- ```
{ "url": "https://asynchronousgillz.github.io/" }
```

- response:

- ```
{"uuid": "1234-1234-1234-1234"}
```



Webpages



Web Scraping



Structured Data

edureka!



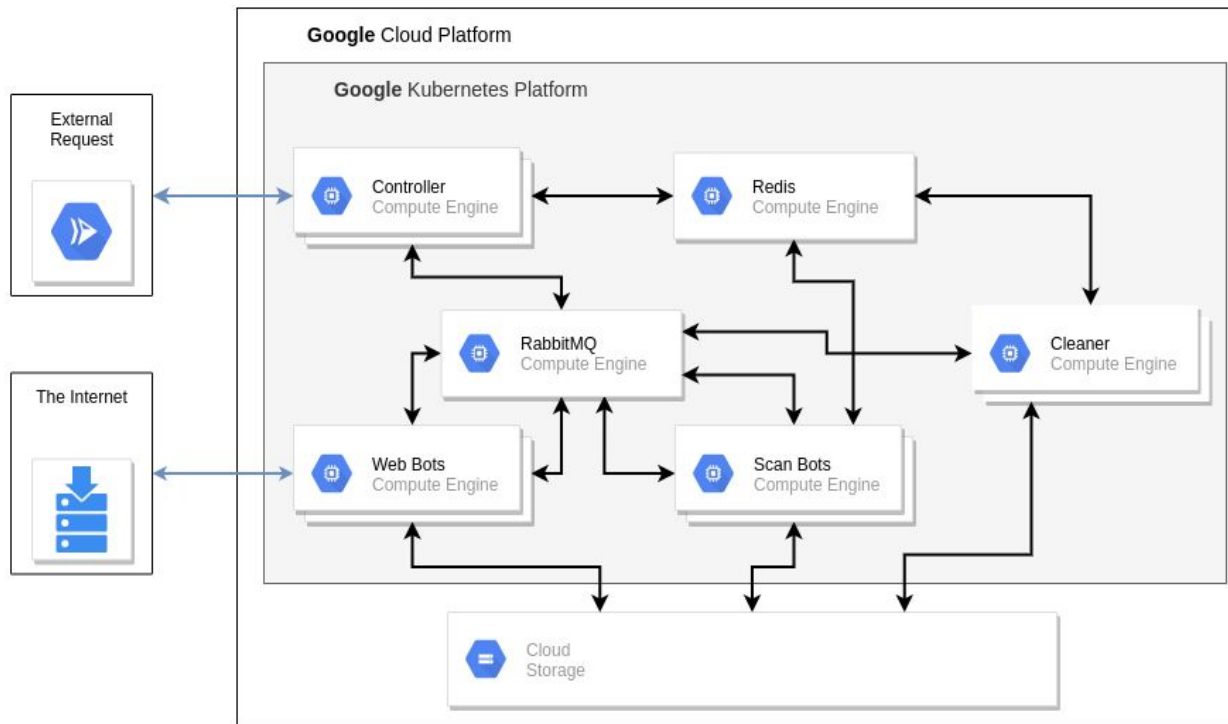
Components

- **Google Kubernetes Engine**
- **Google Cloud Storage**
- **RabbitMQ**
- **Redis**
- **Docker / python3 services**
 - **python packages:**
 - **Flask, Pika, Redis, Requests, BeautifulSoup, Google Cloud Libraries**

Architectural Diagram



Architecture: term-project web crawler



Debugging

■ Logging with containers

- Using docker-compose locally to test messaging / cache interactions

■ Cluster communication

- RabbitMQ messaging

■ Rubber Duck on my desk

- Doesn't give up any secrets of the universe but can helped me work through the problems I encountered.

Learnings

■ Resource limiting to not overload web servers

- distributed lock manager
 - one worker at a time per domain

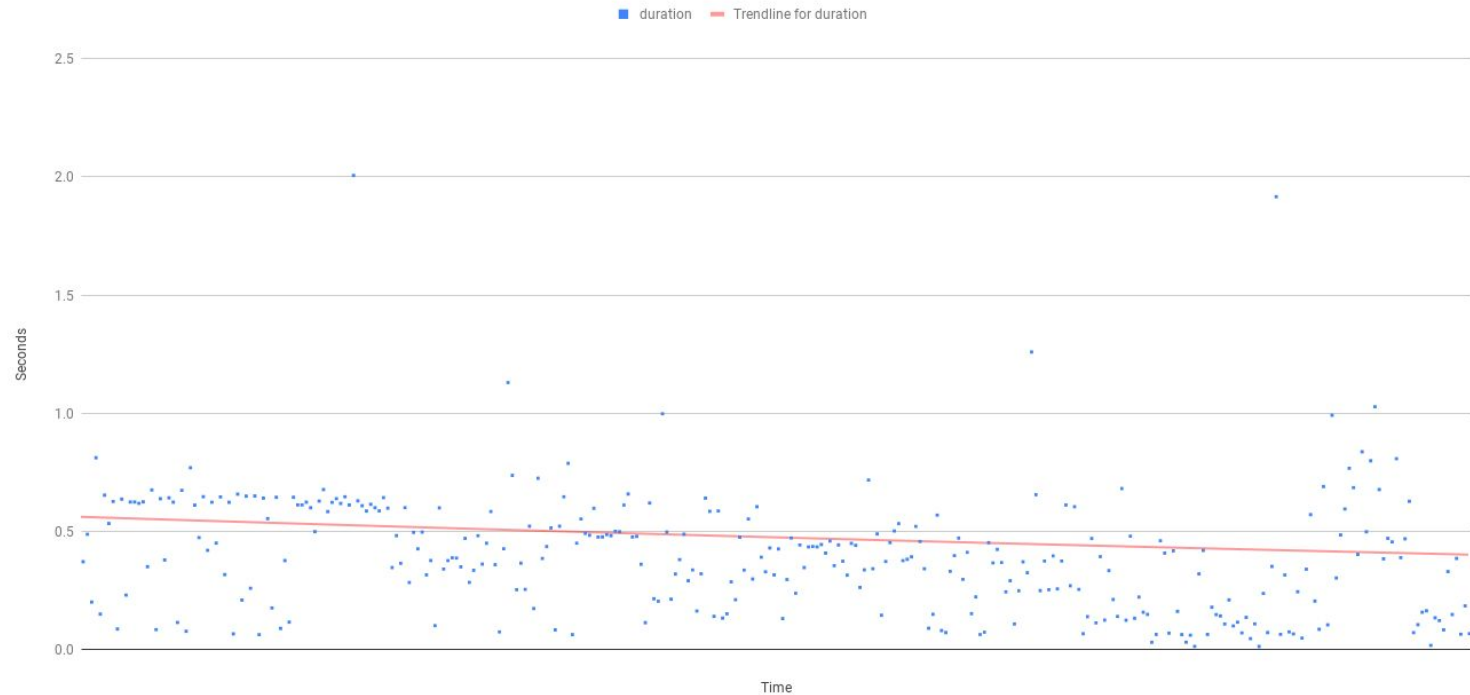
■ Following the rules of robots.txt

- crawl delay
 - time between requests
- sitemap
 - use a sitemap to provide information about specific types of content on your pages

■ <https://support.google.com/webmasters/answer/7451184>

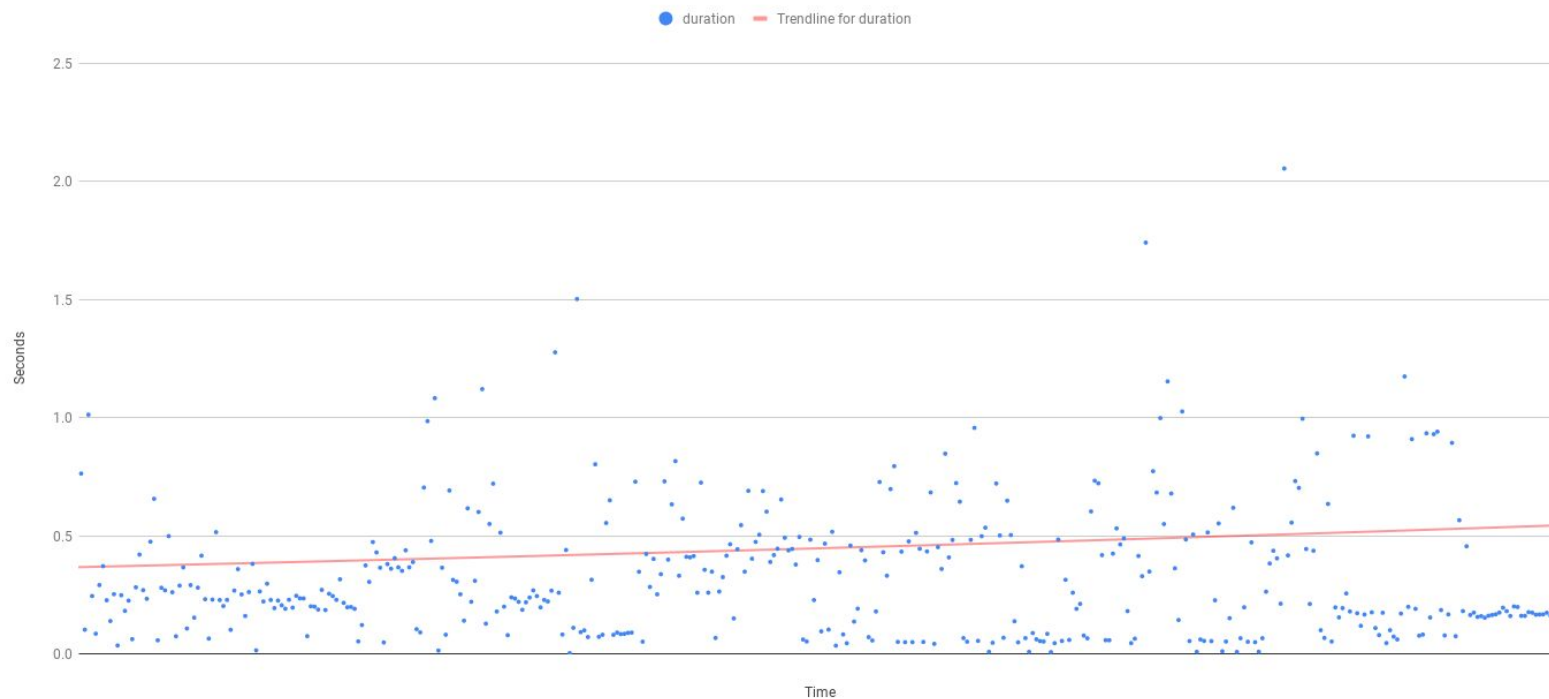
Results United-States

US Central Load Times



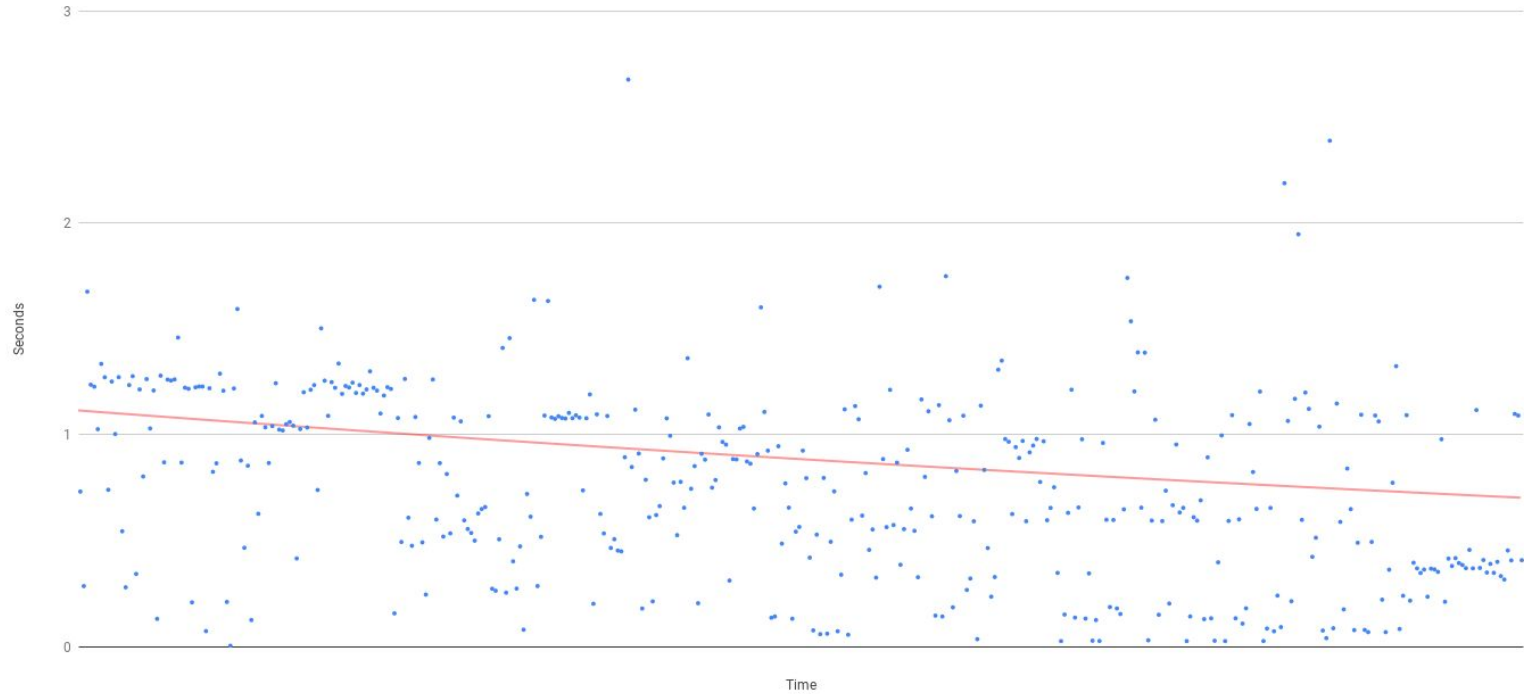
Results Europe

EU Central Load Time



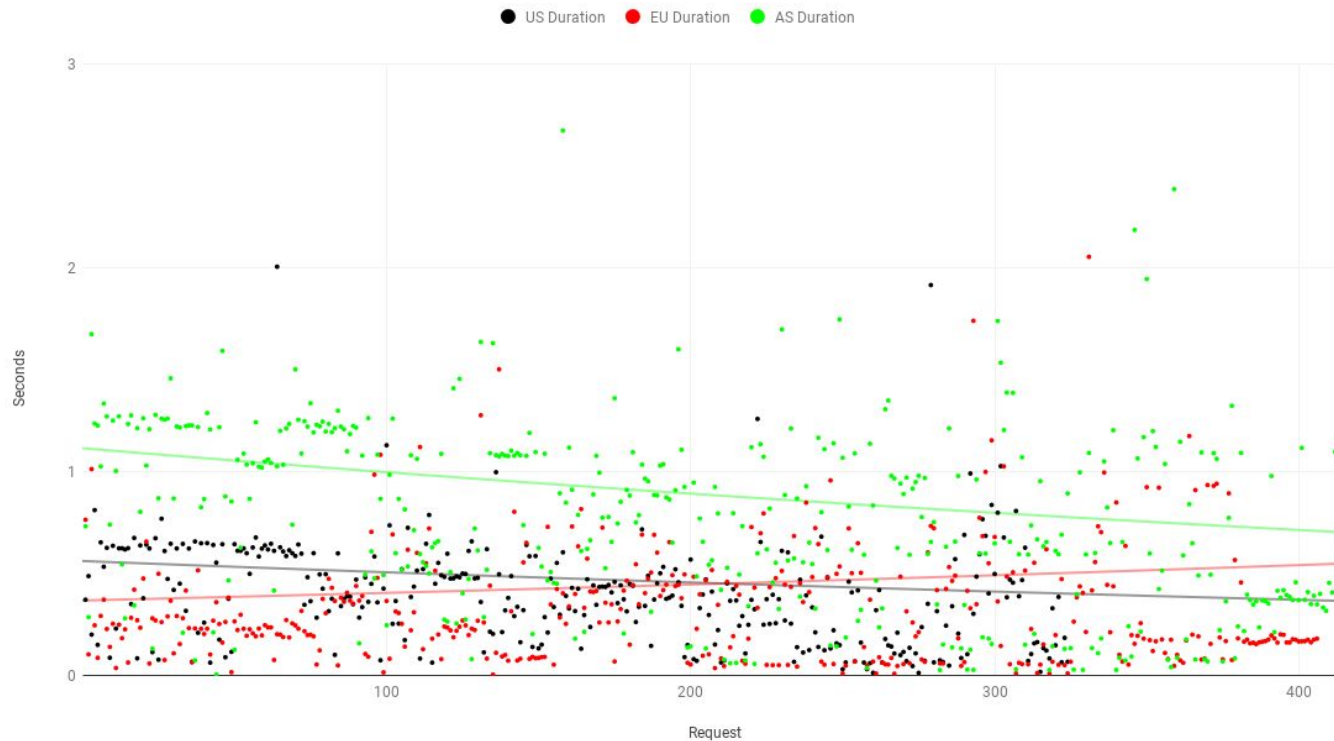
Results East Asia

AS East Load Times



Results Compared

Compared Duration





Thank You!

contact information

email: geva2368@colorado.edu