# CS 467 Graphical Web Crawler

**Kerensa Crump - crumpke@oregonstate.edu**
**Brent Freeman - freeembre@oregonstate.edu**
**Long Le - lelon@oregonstate.edu**

Capstone Project to create a web application that graphically demonstrates two methods of traversing the internet.

## Description

The Gudja online web crawler was Created by OSU's team Gudja for the CS 467 Capstone class. In less than 10 weeks, our team successfully built a single page application for a web crawler that allows the user to traverse the internet using two different algorithms: breadth first search (BFS) and depth first search (DFS).

BFS examines all links at each level starting from the user supplied URL. The algorithm expands the search at every level, like an expanding bubble. DFS searches one branch, or unbroken sequence of links, at a time. It begins with the first link and continues until it reaches the desired depth. As an option, the user can tell the crawler to end the search if it encounters a keyword on a web page.

The graphical result is a high level map of all the external pages that were visited by the web crawler. You'll find that BFS is much slower than DFS and each algorithm will give you different information about the connectivity of the internet.

Our graphical web crawler allows users to visually see a web crawling program in action while exploring the internet. The simulation will provide basic insight on how a real life web crawler, such as Googlebot, might visit pages on the web.

Kerensa Crump
https://github.com/cadelx

Brent Freeman
https://github.com/freeman-bw

Long Le
https://github.com/lelon32

## Oregon State University

## Nodejs server

## Cloud_dfs

## Client-Server Architecture

Internet

User1 — Client Browser
User2 — Client Browser
User3 — Client Browser

Google Cloud Server

Web Crawler

## Hosted App Architecture

Google Cloud Platform

App Engine — Node.js Server, Angular Frontend, D3 Design

Cloud Functions — Python Function, Beautiful Soup Scraper

## D3 Graphical Features

- Web page title - tooltip
- Web page URL - tooltip
- Active link - node onClick
- Domain name - node color
- Links between web pages
- Directionality of the link --->
- Node connectivity - diameter
- Keyword found - highlight

## Front-End Architecture

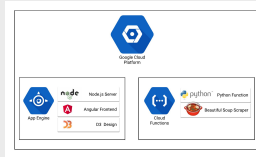**Angular Data Models**
Post, CookieData, CrawlerData

**Angular Components**
1. Header
2. Home
3. Footer
4. Post-Create
5. Activity
6. Crawler
7. Post-Display
8. Cookie

**Angular Service Calls**
1. None
2. None
3. None
4. PostsService
5. None
6. CrawlerService, PostsService
7. PostsService
8. CookieDataService

## GUI & Breadth First Search



## Depth First Search Result



## Technologies Used

**Angular** - Angular provided the framework to create a responsive single page application. Angular's MVC structure utilized models, components, and services along with observables to implement two way data binding and reducing coupling. **See Front-End Architecture, center right**

**D3**- D3 is a JavaScript library that uses native HTML, CSS, and Javascript with data binding to turn create a wide array of visualizations. The core of the network graph visualization used the forceSimulation module. Features like the tooltip, dynamic svg sizing animation, and pulsing highlight were managed separately using custom functions tied to the events like mouseover and simulation tick. **See D3 Graphical Features, center**

**Node.js** - The backbone of or app, Node is both familiar to the whole team and also simple to use. Node Packages: Express, Cors, Request-Promise, Body-Parser, Cookie-Parser. **See Node.js server, left**

**Python 3.x** - Python was chosen to develop the web crawler programs because it was an opportunity to learn something new. There are many positives to using Python including: clean code, higher levels of abstraction, fast development speed, great documentation, and a large number of useful libraries. Python Libraries: Beautiful Soup 4, lxml, NumPy.
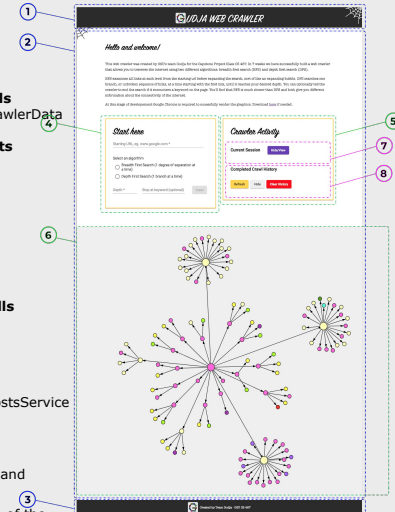
**Google Cloud Functions** - Google's answer to serverless technology. We created our individual functions written out as Python methods, Google hosts these functions, and we can access them with a http request. **See Cloud_dfs, left.**

**Google App Engine** - Google App Engine serves as the host for our application. This platform allows for a relatively easy method to share our application with the public with less worry about maintenance. It supports different languages giving us freedom to develop in any languages.

**Slack** - The go-to communication platform. Slack was instrumental in allowing us to share our progress, ideas and frustrations and fixes.

**Git and GitHub** - GitHub was used for version control during the development of the Gudja Web Crawler. It was particularly useful for managing four different compiled front ends needed for our different testing/deployment environments