# Expert Search - Automated Crawler

Runs to discover Faculty Bio-pages by crawling University Websites.

Mohana Venkata Kalyan Cheerla
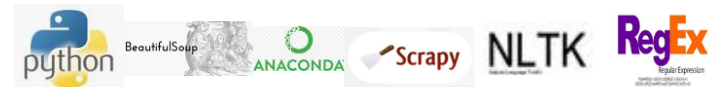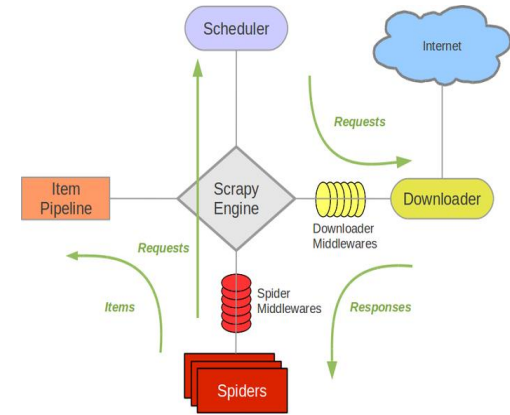
# Why an Automated Crawler?

- ➜ Manually searching the web and extracting faculty bio pages are time consuming

  - ◆ It's tedious

  - ◆ Doesn't scale

  - ◆ Navigation is difficult

- ➜ Our Automated Crawler crawls university websites to automatically fetches bio pages

- ➜ It converts the Markup text into normalized text data to feed its downstream components

# Technology Stack

➔ The crawler runs on a cron schedule to to passively crawl universities websites

➔ The crawler is developed in Python. It utilizes the Scrapy library and follows Breadth First Search (BFS) up to the designated depth specified by the User/Administrator

➔ Crawler has Spiders that recursively follow all the web links of the University Website and extracts the web page data along with the Website's "URL Tree" to understand, validate, troubleshoot and trace

# Crawler in Action

# Crawler in Action (contd...)



Faculty Bio Page URLs

# Crawler in Action (contd...)

# Challenges

➜ Crawler runs a long time and drains computational resources on the developer's machine when we run it for all Universities through all the web pages of the university websites without any limitation on its depth.

   ◆ Parallelization over multiple machines works best

➜ Optimizing the crawler was challenging

   ◆ We had to eliminate different branches of crawl tree that were uninteresting, e.g. publications, news, sports

# Questions?

Please reach out to cheerla3@illinois.edu if you have any question.

# Thank you!!